Introduction to Computational Linguistics

Gemini Pro

October 2025 Version 0.1 Engineered by Andrei Dubovik.

With the sole exception of the title page and this imprint page, this entire book, including all the diagrams and tables, has been automatically generated using a large language model. No information from this book has been independently verified. A reader discretion is advised.

This book is provided solely as an illustration of the capabilities of large language models. This book is released with the understanding that the author of this project, Andrei Dubovik, is not rendering technical or other professional advice. In particular, the project's author disclaims any liability that is incurred from the use or application of the contents of this book.

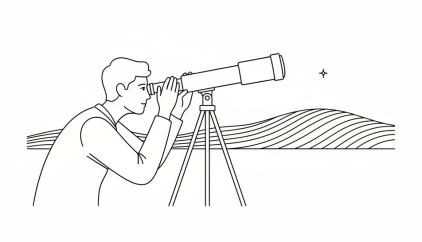
Personal names and company names used throughout this book should be viewed as fictitious. Any resemblance to actual persons or companies should be viewed as coincidental.

Contents

| T | Introduction to Computational Linguistics | 3 |
|----------|--|-----|
| 2 | Words, Regular Expressions, and Automata | 16 |
| 3 | Corpus Linguistics and Text Normalization | 33 |
| 4 | Language Modeling with N-grams | 49 |
| 5 | Part-of-Speech Tagging | 67 |
| 6 | Syntactic Parsing | 81 |
| 7 | Lexical and Compositional Semantics | 100 |
| 8 | Discourse, Coreference, and Dialogue | 119 |
| 9 | Machine Translation | 134 |
| 10 | Information Retrieval and Information Extraction | 151 |
| 11 | Sentiment Analysis and Opinion Mining | 166 |
| 12 | The Future: Large Language Models and Ethics | 180 |

Chapter 1

Introduction to Computational Linguistics



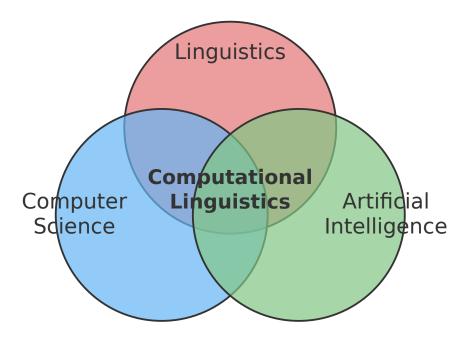


Figure 1.1: A Venn diagram illustrating computational linguistics as the intersection of three overlapping fields: linguistics, computer science, and artificial intelligence.

Computational linguistics is the scientific study of language from a computational perspective. At its core, the discipline is concerned with modeling human language in all its forms—from the sounds of speech to the structure of sentences and the flow of dialogue—using the formalisms and algorithms of computer science.

This pursuit has a dual nature. First, it is a *scientific* discipline. By creating computational models of linguistic phenomena, we can test the validity and predictive power of linguistic theories, leading to a deeper, more rigorous understanding of how language itself works. Second, it is an *engineering* discipline, focused on building useful technologies that can process, understand, interpret, and generate human language.

This interplay between scientific inquiry and technological application is a defining characteristic of the field. Theoretical insights about linguistic structure often pave the way for new applications, while the challenges of building practical systems frequently reveal the limitations of our current theories, driving the science forward. It asks not only how language is structured, but also how we can build machines that process it.

Computational linguistics is not a monolithic discipline but a vibrant synthesis of three distinct fields: linguistics, computer science, and artificial intelligence. As illustrated in the Venn diagram in Fig. 1.1, it resides at the intersection where the scientific study of language, the theory of computation, and the quest for intelligent behavior converge. Each parent field provides indispensable components, and understanding their unique contributions is key to grasping the nature of our subject.

From *linguistics*, the field inherits its core object of study and theoretical grounding. Linguistics provides the formal models for describing language structure, from the smallest sounds (phonology) and word forms (morphology) to sentence structure (syntax) and meaning (semantics and pragmatics). It frames the fundamental questions: What constitutes a grammatical sentence? How do words acquire meaning from context? How do we resolve ambiguity? Without the descriptive and theoretical frameworks developed over centuries of linguistic study, our computational efforts would be unguided, lacking a deep understanding of the phenomena we seek to model.

Computer science provides the 'computational' in the field's name, contributing the algorithms, data structures, and engineering principles necessary to make linguistic theo-

ries operational. It is the discipline that allows us to process language at a massive scale. Formal language theory, which we will encounter in Chapter 2, provides the mathematical foundation for modeling language structures, while efficient algorithms for tasks like parsing and searching enable the creation of practical applications. Computer science provides the rigor and scalability to transform abstract linguistic models into working systems that can handle the complexity and sheer volume of real-world text and speech.

Finally, artificial intelligence (AI) provides the overarching ambition and many of the most powerful modern techniques. The goal of building machines that can understand and generate language is a quintessential AI problem, often seen as a hallmark of true intelligence. Historically, computational linguistics has been considered a core subfield of AI. In the contemporary era, this link is stronger than ever, as AI's machine learning paradigm—and deep learning in particular—has become the dominant approach for nearly every language task. This synergy is what propels the field forward: linguistic insight informs the architecture of AI models, which are then implemented and evaluated using the powerful tools of computer science.

To appreciate the complexity that computational linguistics addresses, consider a seemingly simple sentence:

I saw a man on a hill with a telescope.

For a human, understanding this sentence feels effortless. We might briefly register a slight ambiguity, but context or common sense usually guides us to a single, intended meaning. For a computer, however, this sentence presents a significant challenge. The string of words is grammatically correct, yet it allows for at least two distinct interpretations based on its underlying structure:

- 1. There is a man on a hill, and I used a telescope to see him. In this case, the telescope is the instrument of the action 'saw'.
- 2. There is a man on a hill, and that hill also has a telescope on it. In this case, the telescope is part of the description of the hill.

This is a classic example of *structural ambiguity*, also known as syntactic ambiguity. The uncertainty doesn't come from the words themselves but from how they can be grammatically combined. The source of the confusion is the prepositional phrase *with a telescope*. What does it modify? Does it attach to the verb phrase (*saw...with a telescope*) or the noun phrase (*a hill with a telescope*)?

This ambiguity is visually represented in the two syntactic parse trees shown in Fig. 1.2. In one tree, the phrase with a telescope attaches to the verb phrase headed by saw, modifying the act of seeing. In the other, it attaches to the noun phrase a hill, modifying the location. Both are perfectly valid grammatical structures. A computational system, however, has no innate understanding of the world to help it decide. It cannot know if the speaker is an astronomer or is describing a landmark. Without a mechanism to resolve this ambiguity, a system cannot reliably answer questions like 'What did you use to see the man?' or 'What was on the hill?'

Resolving this kind of ambiguity is a central task in computational linguistics. It requires a powerful combination of disciplines. From *linguistics*, we borrow the formalisms of syntax that allow us to identify and represent these different possible structures. From *computer science* and *artificial intelligence*, we leverage algorithms to parse the sentence and statistical models to determine which structure is more probable. For instance, a model trained on a vast amount of text might learn that the verb 'saw' is far more frequently modified by the instrument 'with a telescope' than the noun 'hill' is. This single, illustrative sentence thus encapsulates the core challenge of the field: to make the implicit, probabilistic reasoning that humans perform so easily explicit and computable.

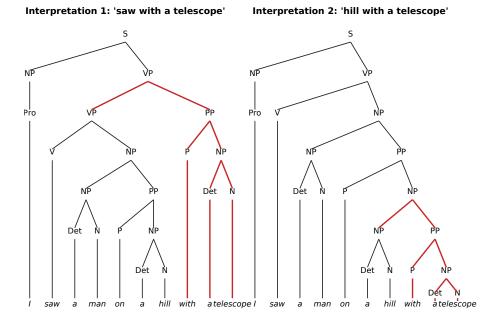


Figure 1.2: Two syntactic parse trees for the ambiguous sentence 'I saw a man on a hill with a telescope.' On the left, 'with a telescope' modifies the verb phrase headed by 'saw' (instrumental reading). On the right, it modifies the noun phrase 'a hill' (locative reading), visually representing structural ambiguity.

As you navigate the field, you will frequently encounter two terms: Computational Linguistics (CL) and Natural Language Processing (NLP). While often used interchangeably, especially in industry, they represent a subtle but important distinction in focus, best understood through the lens of science versus engineering.

Computational Linguistics is fundamentally a scientific discipline. Its primary goal is to understand the nature of human language itself, using computational models as its primary theoretical and experimental tool. A computational linguist asks questions like: What formal mechanisms can explain how humans resolve ambiguity? Can we create a model of syntax that aligns with psycholinguistic evidence of human sentence processing? The ultimate aim is to gain deeper insight into language and the human mind.

Natural Language Processing, on the other hand, is an engineering discipline. It is a subfield of artificial intelligence and computer science focused on building systems that can perform useful tasks involving language. An NLP engineer asks questions like: What is the most accurate and efficient algorithm for translating Spanish to Japanese? How can we build a chatbot that reliably answers customer questions? The success of an NLP system is measured by its performance, robustness, and utility in a real-world application.

Think of the relationship between physics and mechanical engineering. A physicist studies the laws of motion and thermodynamics to understand the universe. An engineer uses those laws to build a more efficient engine. The two are deeply intertwined and mutually beneficial. In the same way, CL and NLP exist in a symbiotic relationship:

- Scientific insights from CL about the structure of language often lead to new techniques and better-performing models in NLP.
- The practical challenges and empirical results from building **NLP** systems often challenge existing theories and drive new research questions in **CL**.

In this book, we will explore both facets. We will delve into the linguistic theories and formal models that form the scientific core of the field, while also grounding our discussion in the practical algorithms and applications that define modern NLP.

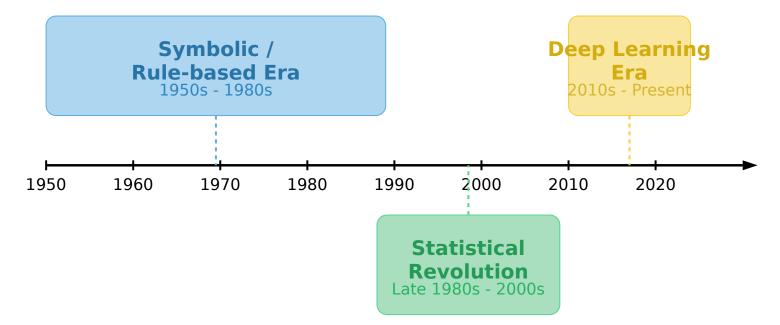


Figure 1.3: A timeline of the major eras in the history of computational linguistics, from the early symbolic approaches to the current deep learning paradigm.

The history of computational linguistics is not one of steady, linear progress but is better understood as a series of major paradigm shifts, each driven by new theoretical insights and advancements in computing power. As summarized in the timeline in Fig. 1.3, we can broadly identify three distinct eras that have shaped the field.

The journey began in the 1950s, alongside the dawn of artificial intelligence itself. This first phase, often called the *symbolic* or *rule-based era*, was characterized by attempts to explicitly encode linguistic knowledge. Researchers believed that by writing down a sufficient number of grammatical rules and dictionary entries, they could enable machines to understand and translate language. This approach dominated the field for several decades but often struggled with the vast complexity and ambiguity inherent in human language.

The late 1980s and 1990s witnessed the *statistical revolution*. This marked a profound shift away from handcrafted rules and towards learning probabilistic patterns directly from large digital text collections known as corpora. Finally, the 2010s ushered in the current *deep learning era*, where neural networks have redefined the state of the art, enabling unprecedented performance on a wide array of language tasks.

The early history of computational linguistics, from the 1950s through the 1980s, was defined by a *symbolic* or *rule-based* paradigm. The guiding philosophy was one of rationalism: human language was seen as a formal system that could be fully described by a set of explicit rules and structures. The task, then, was for linguists to meticulously codify this knowledge—the rules of grammar, syntax, and morphology—and for computer scientists to implement these rules in programs. This approach envisioned the computer not as a learner, but as a logician, mechanically applying human-created linguistic rules to process text.

The most prominent and publicized example of this era was the Georgetown-IBM machine translation experiment in 1954. This demonstration, designed to showcase the potential of automated translation, successfully translated over sixty carefully selected Russian sentences into English. The system was rudimentary by modern standards, relying on a small vocabulary of 250 words and just six handcrafted grammatical rules. A Russian word would be looked up in a bilingual dictionary, and its English equivalent

would be output, with the rules handling basic word reordering and inflectional changes. For instance, a rule might specify how to handle a Russian noun's case ending to produce a correct English prepositional phrase.

The demonstration was a resounding success in the public eye, creating a wave of optimism and securing substantial government funding for the nascent field. The press confidently predicted that automated interpreters would be commonplace within a few years. However, this initial success on a small, controlled problem masked a fundamental weakness. The rule-based approach failed to scale. The number of rules required to cover the vast complexity, ambiguity, and exceptions of a real language grew astronomically. As more rules were added, they would often conflict, leading to a cascade of errors. This 'combinatorial explosion' proved that hand-crafting a comprehensive grammar for a language was an intractable task, setting the stage for a major paradigm shift.

By the late 1980s, the limitations of purely symbolic, rule-based systems were becoming clear. Hand-crafting comprehensive grammars and lexicons proved to be a monumental task; for every rule a linguist could devise, real-world language presented countless exceptions. These systems were notoriously *brittle*: an unfamiliar word or a slightly unconventional grammatical construction could cause them to fail completely. Furthermore, they were difficult to scale or adapt to new domains, as each new application required another intensive, manual effort of rule creation.

This led to a major paradigm shift, often called the *statistical revolution*. The central idea was to move away from hand-crafted rules and towards models that learn patterns automatically from data. Instead of relying on a linguist's intuition, the field embraced probability theory and machine learning. The fundamental question changed from 'Is this sentence grammatically correct according to my rules?' to 'What is the most *probable* interpretation of this sentence given the evidence we have seen?'

This new empirical approach was fueled by two parallel developments: the exponential growth in computing power, which made it feasible to process huge datasets, and the increasing availability of large-scale digital text collections, known as *corpora*. Researchers could now train statistical models on millions or even billions of words of authentic text, allowing the models to learn the likelihood of different linguistic phenomena. Tasks were reframed as problems of statistical inference. For example, machine translation was no longer about applying transfer rules, but about finding the target sentence T that maximizes the probability P(T|S) given a source sentence S. Using Bayes' theorem, this can be expressed as finding the T that maximizes the product of two probabilities:

$$P(S|T) \times P(T)$$

This elegantly separates the problem into a *translation model* (the probability of the source given the target) and a *language model* (the a priori probability of the target sentence). This probabilistic mindset, pioneered in speech recognition and machine translation, made systems more robust and adaptable, laying the essential groundwork for the data-driven methods that define computational linguistics today.

Beginning in the early 2010s, the field underwent its most dramatic transformation to date, ushering in the current era of deep learning. This paradigm shift was fueled by the convergence of three key factors: the availability of massive datasets, significant advancements in parallel computing hardware like Graphics Processing Units (GPUs), and breakthroughs in neural network architectures. Unlike the statistical methods that relied on carefully engineered features and explicit probabilistic models, deep learning models learn relevant features automatically from raw data through a hierarchy of interconnected layers of 'neurons.'

The first major success of this era was the development of dense word representations, known as word embeddings. Models like Word2Vec and GloVe learned to map words to low-dimensional vectors, where semantically similar words were located close to each other in the vector space. This was a revolutionary departure from the sparse, high-dimensional vectors of the previous era. It allowed models to capture nuanced semantic relationships—for instance, learning a vector relationship such that $vector(`king') - vector(`man') + vector(`woman') \approx vector(`queen')$. This ability to represent meaning numerically became a cornerstone of modern NLP.

Following this, architectures specifically designed for sequential data, such as *Recurrent Neural Networks* (RNNs) and their more robust variants like Long Short-Term Memory (LSTM) networks, became dominant. These models process text one word at a time, maintaining an internal state or 'memory' that allows them to capture sequential dependencies, making them highly effective for tasks like language modeling and machine translation.

However, the most significant breakthrough of the deep learning era has been the Transformer architecture, introduced in 2017. Its core innovation, the self-attention mechanism, allows a model to weigh the influence of all words in an input sentence when processing any single word. This enables it to capture complex, long-range dependencies far more effectively than RNNs and, crucially, allows for massive parallelization during training. This architecture is the foundation for the massive Large Language Models (LLMs), such as BERT and GPT, that now define the state of the art. These models are pretrained on immense quantities of text, acquiring a broad understanding of language that can then be fine-tuned for exceptional performance on a wide range of specific tasks. This paradigm of pre-training and fine-tuning has largely replaced the need to train a new model from scratch for every problem, fundamentally changing how computational linguistics is practiced today.

At its core, computational linguistics pursues a deceptively simple ambition: to enable computers to process human language with a facility approaching our own. This overarching ambition can be distilled into two primary, complementary goals, often framed as understanding and generation.

- 1. Natural Language Understanding (NLU): This is the task of analysis—of mapping raw language input to a structured, unambiguous representation of its meaning. It goes far beyond simple keyword matching. True understanding requires a machine to dissect grammatical structure (a process we will explore in Chapter 6), resolve the meaning of words in context (Chapter 7), and infer a speaker's underlying intent. The goal is to *interpret* language, successfully navigating the kinds of ambiguities illustrated by the 'telescope' sentence to arrive at a coherent meaning.
- 2. Natural Language Generation (NLG): This is the inverse task of synthesis. NLG systems start with a non-linguistic, formal representation of information—such as data from a database or a knowledge graph—and produce fluent, grammatically correct, and stylistically appropriate human language. Applications range from automatically generating weather forecasts from meteorological data to summarizing long documents or crafting responses for a chatbot.

These two goals are not independent but form the two halves of a complete communication cycle. A machine translation system must first understand a sentence in a source language before it can generate its equivalent in a target language. A virtual assistant must interpret a user's command before it can generate a helpful response. The ultimate aim is to create systems that can both listen and speak, read and write, thereby bridging the gap between human communication and machine computation.

To see how the different branches of computational linguistics come together, consider one of its most ubiquitous applications: a virtual assistant like Amazon's Alexa or Google Assistant. These systems appear to have a seamless, singular ability to converse, but this

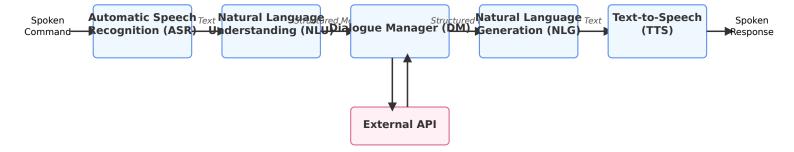


Figure 1.4: A flowchart showing the processing pipeline for a spoken command given to a virtual assistant. The process begins with Automatic Speech Recognition (ASR) converting speech to text. The text is then interpreted by Natural Language Understanding (NLU) to extract intent and entities. A Dialogue Manager (DM) uses this structured data to decide on an action, potentially querying an external API. The result is formulated into a sentence by the Natural Language Generation (NLG) module, and finally converted back into audible speech by a Text-to-Speech (TTS) engine.

illusion is built upon a cascade of distinct computational tasks. This entire sequence is visualized in the processing pipeline shown in Fig. 1.4. Let's trace the journey of a simple spoken command: 'Hey assistant, what's the weather like in Paris tomorrow?'

The first challenge is to convert the physical sound waves of your voice into digital text. This is the domain of **Automatic Speech Recognition (ASR)**. The assistant's microphone captures the utterance, and an ASR model, trained on thousands of hours of speech data, transcribes it into the text string: 'what's the weather like in paris tomorrow'. This process must be robust enough to handle a vast range of accents, speaking speeds, and background noise. A single transcription error at this stage can cause the entire interaction to fail.

Once the command exists as text, the system must decipher its meaning, a task known as **Natural Language Understanding (NLU)**. NLU's goal is to transform the unstructured string of words into a structured, machine-readable format. It does this by identifying two key things:

- Intent: The user's core goal. In this case, the intent is to GetWeatherForecast.
- Entities (or Slots): The specific pieces of information needed to fulfill the intent. Here, the entities are Location: Paris and Date: tomorrow.

The output of the NLU component is no longer just a sequence of words, but a structured representation like: { intent: GetWeatherForecast, location: "Paris", date: "tomorrow" }.

This structured data is passed to the **Dialogue Manager (DM)**, the system's brain. The DM maintains the state of the conversation and decides what action to take. Seeing the **GetWeatherForecast** intent, it knows it must query an external service, such as a weather API. It uses the extracted entities to populate this query, requesting the forecast for 'Paris' on the relevant date. If the NLU module had failed to identify a location (e.g., if the user had just asked, 'What's the weather like tomorrow?'), the DM's logic would decide to ask a clarifying question: 'For which city?'.

After the DM retrieves the information from the weather service—say, { forecast: "sunny", high: "15°C" }—it must deliver this answer to the user. This is where Natural Language Generation (NLG) comes in. The NLG module takes the structured data from the API and converts it into a well-formed, natural-sounding sentence, such as: 'The weather in Paris tomorrow will be sunny, with a high of 15 degrees Celsius.'

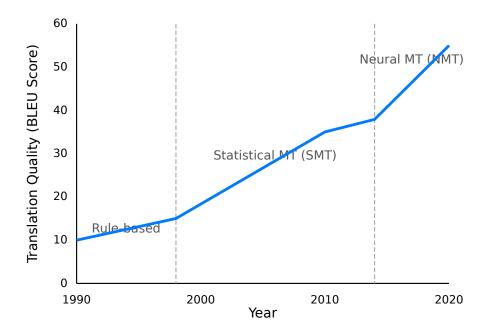


Figure 1.5: The evolution of Machine Translation quality over time, measured by a metric like the BLEU score. The chart highlights the significant performance improvements following the major paradigm shifts from rule-based systems to Statistical Machine Translation (SMT) in the 1990s and the subsequent, even more dramatic, leap with the advent of Neural Machine Translation (NMT) in the mid-2010s.

Finally, this generated text string is sent to a **Text-to-Speech (TTS)** or speech synthesis engine. The TTS system converts the text back into audible speech, producing the synthesized voice that you hear as the final response. This intricate, multi-stage process—from speech to text, text to meaning, meaning to action, action to text, and text back to speech—all happens in just a few seconds, perfectly illustrating how multiple subfields of computational linguistics are integrated to create a single, powerful application.

Perhaps no application demonstrates the transformative power of computational linguistics more vividly than automated Machine Translation (MT). The ambitious goal of instantly translating text from a source language to a target language has been a driving force in the field since its inception. While early systems produced translations that were often more comical than useful, today's platforms, like Google Translate or DeepL, have become indispensable tools for communication, commerce, and information access on a global scale.

The societal impact is immense. These systems empower travelers to navigate foreign countries, enable immigrants to access essential services, and allow researchers to read scientific papers published in other languages. They have fundamentally lowered the barrier to cross-cultural communication, making the world's vast repository of knowledge and culture more accessible than ever before. For businesses, MT facilitates international e-commerce by localizing product descriptions and customer support, opening up global markets that were once prohibitively expensive to enter. A small online retailer can now communicate with customers across dozens of countries, a feat unimaginable just a few decades ago.

This revolution did not happen overnight. It is the direct result of the major paradigm shifts in computational linguistics. As illustrated in **Fig. 1.5**, the quality of machine translation, often measured by automated metrics like the BLEU score, has seen dra-

¹BLEU (Bilingual Evaluation Understudy) is a metric for evaluating a generated sentence to a set of high-quality reference translations. It measures the correspondence between a machine's output and that

matic improvements that coincide with these shifts. The move from rule-based systems to Statistical Machine Translation (SMT) in the 1990s and 2000s led to a significant leap in performance by learning translation patterns from massive amounts of parallel text. The subsequent adoption of Neural Machine Translation (NMT) in the mid-2010s, powered by the deep learning techniques we will explore later in this book, caused an even more spectacular jump in quality. Modern NMT systems can produce stunningly fluent and accurate translations, particularly for high-resource language pairs like French and English.

The economic implications are just as significant. The global language services industry is valued in the tens of billions of dollars, and MT is a disruptive force within it. Rather than replacing human translators, it has created a new collaborative model where machines produce a high-quality first draft that a human expert then refines—a process known as post-editing. This synergy makes professional translation faster and more cost-effective, fueling the engine of global enterprise. Machine translation stands as a testament to the field's ability to tackle a profoundly complex and human problem with computational rigor, yielding a technology that has reshaped our interconnected world.

Search engines like Google or Bing represent one of the most profound and widely-used applications of computational linguistics. At first glance, a search engine might seem like a simple keyword-matching tool, retrieving documents that contain the exact words typed into the search bar. The reality, however, is far more complex. Modern search engines do not just *match* strings of text; they strive to *understand* a user's intent and the content of billions of web pages to deliver truly relevant results. This monumental leap from matching to understanding is powered by core techniques from computational linguistics.

The process begins the moment a user types a query. The engine must first interpret what the user is actually looking for. This involves several tasks. It performs spelling correction for typos and uses knowledge of lexical semantics to perform query expansion, recognizing that a search for 'car repair shops' is also related to 'auto maintenance services.' Furthermore, it must resolve ambiguity. A query for 'jaguar' could refer to the animal, the car brand, or an operating system. By analyzing a user's prior searches or geographic location, the engine can disambiguate the query to infer the most likely intent.

In parallel with understanding queries, the engine must also have a deep understanding of the documents it might return as results. Before a page is ever shown to a user, it is crawled and analyzed using a pipeline of language processing techniques. The text is tokenized, and words are normalized using stemming and lemmatization so that searches for 'run' can match documents containing 'ran' or 'running.' More advanced methods like *Named Entity Recognition* are used to identify and tag key entities like people, organizations, and locations, creating a rich, structured index of the web's unstructured content.

With an understood query and a well-structured index, the final challenge is to rank documents by relevance. This is the central task of $Information\ Retrieval$. Instead of a simple yes/no decision on whether a document contains a keyword, documents are assigned a relevance score. A foundational technique for this is weighting terms by their importance, a classic example being TF-IDF (Term Frequency-Inverse Document Frequency). The score for a term t in a document d is a product of its frequency in that document and its rarity across all documents:

$$TF - IDF(t, d) = TF(t, d) \times IDF(t)$$

The core idea is that a term is significant if it appears frequently in a specific document (TF) but is rare in the overall collection (IDF). This simple but powerful heuristic helps

the engine prioritize documents where the query terms are not just present, but a central part of the topic. By combining these computational linguistics techniques, a search engine orchestrates a complex process of understanding, analysis, and ranking to transform a simple query into a useful list of results in a fraction of a second.

Beyond the prominent examples of machine translation, search, and virtual assistants, the techniques of computational linguistics power a vast ecosystem of applications. These tools are often seamlessly integrated into the software we use daily, working behind the scenes to structure data, understand opinions, and facilitate communication. The following examples represent just a sample of this diverse and impactful landscape.

- Sentiment Analysis and Opinion Mining: This field focuses on automatically identifying and categorizing opinions expressed in text. The primary goal is often to determine the author's attitude—positive, negative, or neutral—towards a particular topic, product, or service. Businesses use this technology extensively to analyze customer feedback, monitor brand perception on social media, and gauge public response to marketing campaigns.
- Information Extraction (IE): While search engines retrieve relevant documents, information extraction systems go a step further by pulling structured data from unstructured text. A key subtask is Named Entity Recognition (NER), which identifies and classifies entities like people (e.g., 'Ada Lovelace'), organizations ('Google'), and locations ('Paris'). Another is Relation Extraction, which seeks to discover the relationships between these entities, such as identifying who works for which company from a collection of news articles.
- Text Summarization: With the overwhelming volume of online text, automatic summarization is an invaluable tool. The objective is to produce a concise and fluent summary of a longer document while retaining its most important information. This technology is used to generate news headlines on aggregator sites, create abstracts for scientific papers, and condense long reports into manageable digests for busy executives.
- Dialogue Systems and Chatbots: While related to the virtual assistants discussed earlier, this broad category also includes more focused conversational agents. Task-oriented chatbots guide users through specific processes like booking a flight or providing automated customer support, freeing up human agents for more complex issues.

These applications, along with others like grammar correction and plagiarism detection, demonstrate the broad utility of the field. Each one relies on a combination of the core linguistic and computational concepts we will explore throughout this book. A more detailed summary of these and other modern applications, outlining their primary goals and providing concrete examples, is presented in Fig. 1.6.

This book is structured to guide you systematically from the foundational principles of computational linguistics to its most advanced applications. We have designed a path that builds knowledge incrementally, ensuring each new concept rests on a solid prior understanding.

Our journey begins with the essential building blocks of language processing. The next three chapters focus on the fundamentals:

• Chapter 2 introduces the formal tools for pattern matching: regular expressions and finite automata.

| Application | Goal | Example |
|--------------------------------|---|--|
| Sentiment Analysis | Determine the attitude (positive, negative, neutral) expressed in a text. | Classify product reviews as positive or negative. |
| Named Entity Recognition (NER) | Identify and classify entities like people, organizations, and locations. | Identify 'Ada Lovelace' (Person) and 'Google' (Organization) in a news article. |
| Relation Extraction | Discover semantic relationships between identified entities. | Determine that a specific person works for a specific company from a text. |
| Text Summarization | Create a concise and fluent summary of a longer document. | Generate headlines for news aggregator sites or abstracts for scientific papers. |
| Dialogue Systems | Guide a user through a task-oriented conversation. | An automated chatbot to help a customer book a flight or check an account balance. |

Figure 1.6: A summary of modern computational linguistics applications, outlining their primary goals and providing concrete examples.

- Chapter 3 grounds us in the empirical reality of language data, covering *corpus linguistics* and essential *text normalization* techniques.
- Chapter 4 provides our first taste of statistical modeling by teaching you how to predict word sequences with *N-gram language models*.

With these core skills established, we will ascend the traditional linguistic hierarchy. We will tackle the core analytical tasks of assigning grammatical roles with *Part-of-Speech tagging* (Chapter 5), uncovering sentence structure through *syntactic parsing* (Chapter 6), and finally, representing meaning with *lexical and compositional semantics* (Chapter 7).

From there, we expand our view beyond single sentences to the broader context of discourse and dialogue (Chapter 8). The subsequent chapters showcase major real-world applications where these techniques converge: Machine Translation (Chapter 9), Information Retrieval and Extraction (Chapter 10), and Sentiment Analysis (Chapter 11). Our exploration culminates in Chapter 12 with a look at the current state-of-the-art—Large Language Models—and a crucial discussion of the ethical considerations that shape the future of our field. This roadmap will equip you with a comprehensive understanding of how we empower computers to process human language.

In this chapter, we have laid the groundwork for our study of computational linguistics. We began by defining the field as the scientific endeavor to understand and model human language using computation, highlighting its unique position at the intersection of computer science, linguistics, and artificial intelligence. We distinguished this scientific goal from the engineering applications of *Natural Language Processing* (NLP), while acknowledging that the two are deeply intertwined.

Our historical overview traced the field's evolution through three major paradigms:

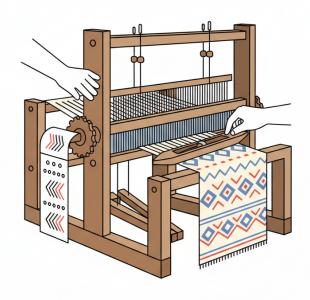
- The Symbolic Era (1950s–1980s): Dominated by hand-crafted grammatical rules and symbolic logic, exemplified by early machine translation efforts.
- The Statistical Revolution (late 1980s–2000s): A pivotal shift towards learning probabilistic models from large-scale text data, or corpora.

• The Neural Era (2010s—Present): Characterized by the remarkable success of deep learning and large-scale neural networks.

Finally, we saw how these approaches empower a vast array of real-world applications, from machine translation and search engines to virtual assistants and sentiment analysis. These examples underscore the core goal of enabling computers to understand, interpret, and generate language, a theme we will explore in technical detail throughout this book.

Chapter 2

Words, Regular Expressions, and Automata



To us, human language appears fluid and intuitive. We process words and sentences effortlessly. For a computer, however, language is merely a sequence of characters, devoid of inherent meaning or structure. To bridge this gap and enable computational processing, we cannot rely on intuition. We need precise, unambiguous methods for describing the patterns that govern language. This is where formal mathematical models become indispensable. They provide a rigorous framework for defining what constitutes a valid word, phrase, or sentence in a given language.

This chapter introduces the foundational tools for this task, drawn from a field known as formal language theory. We will begin by modeling the simplest components of language—the patterns that make up individual words. By defining simple 'languages' as specific sets of strings and designing abstract machines to recognize them, we create our first computational building blocks. This formal approach is not just a theoretical exercise; it provides the robust and efficient foundation upon which almost all complex language processing technologies are built.

To formally model language, we begin with the simplest yet surprisingly powerful class of formal languages: the *regular languages*. In formal language theory, a 'language' is simply a set of strings, and a regular language is a specific type of this set—one whose structure can be described using a constrained but effective set of rules.

Why start here? Because a vast number of word-level phenomena and simple text patterns can be modeled effectively as regular languages. For instance, the set of all English words ending in -ing, the structure of dates like MM-DD-YYYY, or the morphological variants of a noun are all patterns that fall into this class. The key advantage of regular languages is their computational elegance; they can be described by the practical tool of regular expressions and recognized by an efficient class of abstract machines called finite automata. This makes them both theoretically fundamental and highly practical for real-world text processing. They represent the first, essential step on a ladder of formalisms used to describe linguistic structure.

To formalize our study of language, we begin with three foundational concepts: the *alphabet*, the *string*, and the *language*. These concepts provide the mathematical bedrock for defining patterns and structures.

An **alphabet**, denoted by Σ (Sigma), is a non-empty, finite set of symbols. For example, the binary alphabet is $\Sigma = \{0, 1\}$, while the alphabet for lowercase English text is $\Sigma = \{a, b, \ldots, z\}$.

A **string** (or *word*) is a finite sequence of symbols drawn from an alphabet. The string **book** is formed from the English alphabet. The *length* of a string w, denoted |w|, is the number of symbols it contains; for example, |book| = 4. A special case is the **empty string**, which has zero length and is denoted by ϵ (epsilon).

Finally, a **formal language** is simply a set of strings over a given alphabet. Importantly, this set can be either finite or infinite. Given the alphabet $\Sigma = \{a, b\}$, consider the following languages:

- A finite language: $L_1 = \{\epsilon, a, b, ab, ba\}$
- An infinite language: $L_2 = \{a, aa, aaa, \dots\}$, the set of all strings consisting of one or more 'a's.

The set of all possible strings that can be formed from an alphabet Σ is denoted Σ^* (pronounced 'Sigma star'). Therefore, any language L over Σ is a subset of Σ^* , which can be formally stated as $L \subseteq \Sigma^*$.

To define interesting languages, we need formal operations for combining and repeating strings. The most fundamental of these is *concatenation*, which simply joins two strings

end-to-end. If we have a string s_1 = 'compu' and another string s_2 = 'ter', their concatenation, written as s_1s_2 , is the string computer. Concatenating any string s with the empty string ϵ results in the original string: $s\epsilon = \epsilon s = s$.

While concatenation builds longer strings, we also need a mechanism for repetition. This is provided by the *Kleene star* or *Kleene closure*. For a given language (a set of strings) L, its Kleene closure, denoted L^* , is the set of all strings formed by concatenating zero or more strings from L. The 'zero' case is important; it means the empty string ϵ is always an element of L^* . For example, if $L = \{`a'\}$, then $L^* = \{\epsilon, `a', `aa', `aaa', \ldots \}$.

A closely related operator is the Kleene plus or positive closure, denoted L^+ . This operation is similar to the Kleene star but requires one or more concatenations of strings from L. Therefore, L^+ includes all strings in L^* except for the empty string ϵ . These three operations—concatenation, Kleene star, and Kleene plus—are the essential building blocks for defining the patterns found in regular languages.

Having formally defined the concept of a regular language, we now turn to the standard notation used to describe them: **regular expressions**, often abbreviated as *regex*. A regular expression is a compact algebraic formula for specifying a set of strings. Think of it as a specialized 'pattern language,' purpose-built for finding and manipulating text. Every regular expression defines a regular language, and conversely, every regular language can be described by a regular expression.

The construction of regular expressions is recursive. We start with simple, atomic expressions and combine them using a small set of powerful operators to form more complex patterns. The most basic regular expressions are the individual characters of our alphabet Σ . For example, the regular expression a describes the language $L(a) = \{'a'\}$. From this foundation, we use three fundamental operations to build more elaborate expressions:

- Concatenation: This is the default operation, specified by placing one expression after another. The regex book is the concatenation of the expressions b, o, o, and k. It describes the language whose only member is the string 'book'. Formally, if R_1 and R_2 are regular expressions for languages $L(R_1)$ and $L(R_2)$, their concatenation R_1R_2 defines the language $L(R_1R_2) = \{xy \mid x \in L(R_1) \text{ and } y \in L(R_2)\}$. This means any string from the first language followed by any string from the second.
- Union (Disjunction): The vertical bar |, often called a pipe, represents choice. It functions like an 'or' operator, matching any of the expressions it separates. For instance, the regex gray | grey is useful for handling spelling variations, as it matches either the string 'gray' or 'grey'. The language it defines is the set-theoretic union of the languages of its components: $L(R_1|R_2) = L(R_1) \cup L(R_2)$.
- Kleene Star (Closure): The asterisk *, named after mathematician Stephen Kleene, signifies 'zero or more' repetitions of the immediately preceding expression. For example, a* matches the empty string, 'a', 'aa', 'aaa', and so on. The language it defines is the set of all finite strings of 'a's, including the empty string ϵ . To apply the Kleene star to a sequence of characters, we must group them with parentheses. This distinction is crucial: ab* matches 'a' followed by zero or more 'b's (e.g., 'a', 'ab', 'abb'), while (ab)* matches zero or more instances of the entire string 'ab' (e.g., '', 'ab', 'abab').

To avoid ambiguity, these operators have a defined order of precedence. The Kleene star has the highest precedence, followed by concatenation, and finally union has the lowest. Therefore, cat|dog* is interpreted as (cat)|(dog*), not as (cat|dog)*. Parentheses can be used to override this default order, just as in arithmetic. These three core operators, summarized for convenience in Fig. 2.1, are all that is needed to describe any regular language.

| Operator | Name | Example |
|---------------------|---|---|
| (juxtaposition) * | Concatenation Union (Disjunction) Kleene Star (Closure) | book matches 'book' gray grey matches 'gray' or 'grey' a* matches ϵ , 'a', 'aa', |

Figure 2.1: A summary of the three fundamental regular expression operators.

Let's put the basic operators of concatenation, union, and closure into practice. Imagine we are given a small block of raw text and our goal is to extract all dates and email addresses.

Consider the following text:

The project kickoff is 03/15/2024. Please contact either jane.doe@university.edu or the lead, smith-j@dept.org, for more details. The final deadline is set for 04-20-2024.

First, we will devise a pattern to find dates in the MM/DD/YYYY or MM-DD-YYYY format. We can observe that a date consists of two digits, a separator, two more digits, another separator, and finally four digits.

- Two digits can be represented by [0-9] [0-9]. This is a concatenation of two character sets, each allowing any digit from 0 to 9.
- The separator is either a hyphen or a forward slash. We can represent this choice using a character set as a form of union: [-/].
- Four digits are [0-9] [0-9] [0-9].

By concatenating these components, we arrive at our first regular expression for finding dates: [0-9][0-9][-/.][0-9][0-9][0-9][0-9]

When applied to the text, this pattern would successfully match both 03/15/2024 and 04-20-2024.

Next, let's build a pattern for email addresses. A simplified structure is a sequence of characters for the username, an @ symbol, and a sequence of characters for the domain name. We can use the positive closure operator, +, to mean 'one or more' of the preceding character type.

- A username can contain letters, numbers, periods, and hyphens. A pattern for this is [a-zA-Z0-9.-]+.
- This is followed by the literal @ symbol.
- The domain name is similar, often containing letters, numbers, and hyphens, followed by a period and the top-level domain. We can model this as [a-zA-Z0-9-]+ followed by a literal . and finally [a-zA-Z.]+.

Combining these gives us a basic expression for emails: [a-zA-Z0-9.-]+0[a-zA-Z0-9-] +\. [a-zA-Z.]+

This pattern will match jane.doe@university.edu and smith-j@dept.org. These examples show how a few simple operators can build powerful tools for locating structured information within unstructured text.

While the basic operators of union, concatenation, and Kleene star provide a complete system for defining regular languages, most modern regular expression engines offer a much richer and more convenient syntax. These extensions don't add theoretical power—they can all be simulated with the basic operators—but they make writing complex patterns significantly more concise and readable. We will now expand our toolkit to include four key

| Syntax | Description | Example | |
|-----------|---|--|--|
| Characte | r Classes | | |
| [abc] | Match any character from the set. | [aeiou] matches 'a' in 'cat'. | |
| [^abc] | Match any character not in the set. | [^0-9] matches 'a' in 'a1'. | |
| [a-z] | Match any character in the range. | [a-z] matches 'c' in 'cat'. | |
| \d | Match any digit character. | \d matches '5' in 'file_5'. | |
| \w | Match any word character (alphanumeric or _). | \w matches '_' in 'a_b'. | |
| \s | Match any whitespace character. | \s matches the space in 'a b'. | |
| Quantifie | rs | | |
| * | Match preceding element 0 or more times. | ab*c matches 'ac' and 'abc'. | |
| + | Match preceding element 1 or more times. | ab+c matches 'abc' and 'abbc'. | |
| ? | Match preceding element 0 or 1 time. | colou?r matches 'color'. | |
| {n} | Match preceding element exactly n times. | \d{3} matches '123'. | |
| $\{n,m\}$ | Match preceding element from n to m times. | \d{2,4} matches '99' and '1234'. | |
| Anchors | | | |
| ^ | Anchor match to the start of the string/line. | ^cat matches 'cat' in 'catdog'. | |
| \$ | Anchor match to the end of the string/line. | cat\$ matches 'cat' in 'tomcat'. | |
| \b | Match a word boundary position. | \bcat\b matches 'cat', not in 'caterpillar'. | |
| Capturin | Capturing Groups | | |
| () | Group pattern and capture the matched string. | (US EU) captures 'US' from 'Order: US 543'. | |

Figure 2.2: A summary of common extended regular expression syntax.

concepts: character classes, quantifiers, anchors, and capturing groups. A comprehensive list of these syntactic elements is provided in Fig. 2.2.

A character class, denoted by square brackets [], allows you to match any single character from a specified set. For instance, [aeiou] will match any lowercase vowel. Ranges are also permitted, so [a-z] matches any lowercase letter and [0-9] matches any digit. A ^ as the first character inside the brackets negates the class, so [^0-9] matches any character that is *not* a digit. For convenience, several predefined character classes exist, such as \d for digits ([0-9]), \w for 'word' characters (letters, numbers, and underscore), and \s for whitespace characters (spaces, tabs, newlines).

Quantifiers control how many times a preceding part of an expression can occur. The three most common are:

- The asterisk * matches the preceding item zero or more times. ab*c matches ac, abc, abc, etc.
- The plus + matches one or more times. ab+c matches abc and abbc, but not ac.
- The question mark? matches zero or one time, making the preceding item optional. The pattern colou?r matches both 'color' and 'colour'. For more precise control, braces can specify a range of repetitions. For example, \d{3} matches exactly three digits, while \d{2,4} matches between two and four digits.

Anchors are special characters that don't match any character in the string but instead match a *position*. The caret ^ anchors the match to the beginning of a line or string, while the dollar sign \$ anchors it to the end. The pattern ^cat\$ will only match the exact string 'cat' and nothing else. Another crucial anchor in computational linguistics is \b, which matches a word boundary. This is the zero-width position between a word character (\w) and a non-word character (\w). Using \bcat\b ensures you match the word 'cat' but not the 'cat' in 'caterpillar'.

Finally, **capturing groups**, created with parentheses (), serve two purposes. First, they group part of a pattern so that operators like * or + apply to the whole group, as in (ab)+, which matches ab, abab, and so on. Second, and more importantly, they *capture* the portion of the string that matched the enclosed pattern. This allows you to extract specific pieces of information from a larger text. For example, in the string 'Order: US 543', the pattern (US|EU) (\d+) would not only match 'US 543' but also allow you to separately extract the captured groups 'US' and '543'.

By combining these elements, we can build powerful and precise patterns. Consider a pattern for matching simple email addresses: $([\w.-]+)@([\w.-]+)\.([a-zA-Z]{2,})$. This expression uses character classes, quantifiers, and capturing groups to identify and extract the username, domain name, and top-level domain from a string. With this expanded syntax, regular expressions become an indispensable tool for information extraction, as we will see in the following case study.

With our expanded toolkit of advanced operators, regular expressions become more than just a tool for finding patterns; they are a mechanism for simple *Information Extraction* (IE). The goal of IE is to identify and pull structured data, like names, dates, or relationships, from unstructured text. While complex IE requires more powerful machine learning models, regular expressions are highly effective for tasks where the input data has a consistent, predictable format.

Let's consider a practical case study: extracting author and title information from a bibliography file. Imagine our file contains entries formatted in a consistent style, such as: Chomsky, N. (1957). *Syntactic Structures*. Mouton.

Our goal is to write a single regular expression that can read a line like this and extract 'Chomsky, N.' as the author and 'Syntactic Structures' as the title. We can achieve this by designing a pattern that matches the entire line while using capturing groups (...) to isolate the specific pieces of information we want.

Consider the following regular expression: $([\w\s,.-]+)\s(\d{4}\)\.\s*([^*]+)*.*$ Let's break this down:

- ^ asserts our pattern must start at the beginning of the line.
- ([\w\s,.-]+) is our first capturing group for the author. It matches one or more characters that are either a word character (\w), whitespace (\s), a comma, a period, or a hyphen.
- \s\(\d{4}\)\.\s matches the non-captured parts in between: a space, the parenthesized four-digit year, a period, and another space.
- *([^*]+)* is our second capturing group. It first matches a literal asterisk (*). The group itself, ([^*]+), is a crucial trick: it matches one or more characters that are *not* a closing asterisk. This non-greedy approach ensures we only capture the text of the title. It concludes by matching the final literal asterisk.
- .*\$ matches the rest of the line until the end.

When this expression is applied to each line of a bibliography, a program can directly access the content of the capturing groups for every successful match. The text matched by the first group is the author, and the text from the second is the title. As **Fig. 2.3** demonstrates, this technique effectively transforms a simple text file into a list of structured records. This approach is powerful but brittle; it relies on the strict consistency of the input format. A missing year or a title not enclosed in asterisks would cause the pattern to fail. Nonetheless, for well-formatted data, regular expressions provide a direct and efficient method for rudimentary information extraction.

| Input String | Extracted Data |
|--|---|
| Chomsky, N. (1957). *Syntactic Structures*. Mouton. | Author: Chomsky, N. Title: Syntactic Structures |
| Bird, S., Klein, E., and Loper, E. (2009). *Natural Language Processing with Python*. O'Reilly Media. | Author: Bird, S., Klein, E., and Loper, E. Title: Natural Language Processing with Python |
| Manning, C. D., Schutze, H. (1999). *Foundations of Statistical Natural Language Processing*. MIT Press. | Author: Manning, C. D., Schutze, H. Title: Foundations of Statistical Natural Language Processing |

Figure 2.3: Using a regular expression with capturing groups to extract structured data (author, title) from formatted bibliography entries.

While regular expressions provide a compact notation for *describing* patterns in text, they don't specify *how* a computer should perform the matching. For that, we turn to a more formal computational model. The abstract machine that recognizes the languages described by regular expressions is known as a **finite-state automaton**, often abbreviated as FSA or simply a finite automaton.

An FSA is an idealized computational device, one of the simplest models of computation. The core idea is that the machine can only be in one of a finite number of *states* at any given time. It processes an input string one symbol at a time from left to right. As it reads each symbol, it *transitions* from its current state to another state. The specific transition is determined by the current state and the input symbol being read.

Formally, an FSA can be visualized as a directed graph:

- The **states** are represented by nodes (circles).
- The **transitions** are represented by labeled edges between the nodes.
- There is one special **start state**, where the machine begins processing.
- One or more states are designated as **final states** or **accept states**, often drawn with a double circle.

The process of recognizing a string begins at the start state. The automaton consumes the input string symbol by symbol, following the corresponding transition edge from state to state. After the final symbol is read, the machine stops. If it has landed in a final (accept) state, the string is considered **accepted**. If it ends in any other state, the string is **rejected**. The set of all strings accepted by an FSA defines the formal language it recognizes. The crucial insight is that finite automata recognize precisely the class of regular languages. This means that for any regular expression, an equivalent FSA can be built to recognize the same set of strings, and vice versa. This equivalence forms the theoretical foundation for implementing regular expression engines.

While a state diagram provides an intuitive picture of a finite automaton, a formal definition allows us to be precise. We will begin with the **Deterministic Finite Automaton (DFA)**, a type of FSA where every move is uniquely determined by the current state and the input symbol. There is no ambiguity; from any state, a given character leads to exactly one other state.

Formally, a DFA is a 5-tuple $M = (Q, \Sigma, \delta, q_0, F)$ where:

- 1. Q is a finite set of states.
- 2. Σ is a finite set of input symbols, called the alphabet.

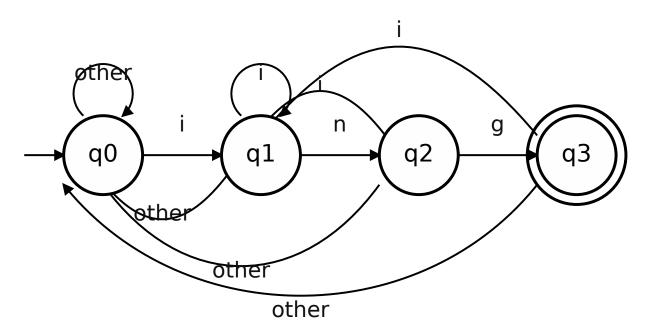


Figure 2.4: A state diagram of a simple Deterministic Finite Automaton (DFA) that recognizes strings ending in 'ing', as described in the text. The states q0, q1, q2, and q3 represent having seen '', 'i', 'in', and 'ing' respectively. q0 is the start state and q3 (double circle) is the final (accepting) state.

- 3. δ is the transition function, which maps a state and an input symbol to a state. Its signature is $\delta: Q \times \Sigma \to Q$. This function is the 'program' of the machine, dictating every possible move.
- 4. $q_0 \in Q$ is the designated start state.
- 5. $F \subseteq Q$ is the set of final or accepting states.

Let's ground this formal definition using the state diagram shown in **Fig. 2.4**, which depicts a DFA that recognizes the language of all strings ending in **ing**.

The components of this specific automaton can be formally defined as follows:

- States (Q): The circles in the diagram represent the states. We can name them $\{q_0, q_1, q_2, q_3\}$. Each state represents a stage in the recognition process, essentially remembering how much of the target suffix (ing) we have just seen. For instance, q_0 is the initial state, q_1 is the state after seeing an i, q_2 after seeing in, and q_3 after seeing ing.
- **Alphabet** (Σ): This is the set of all characters the machine can process. For simplicity, let's assume it's the set of all lowercase English letters, $\{a, b, c, \ldots, z\}$.
- Start State (q_0) : In the diagram, q_0 is marked as the start state by an incoming arrow that originates from no other state. This is where processing for any input string begins.
- Final States (F): The set of final states is $\{q_3\}$. In diagrams, final states are indicated by a double circle. If the machine's final state after reading an entire input string is a member of F, the string is accepted.
- Transition Function (δ): The labeled arrows represent the transition function. For example, the arrow from q_0 to q_1 labeled i means that $\delta(q_0, i) = q_1$. The

deterministic nature is clear: from q_1 , an n must lead to q_2 , while any other character (like a or i) might lead back to a different, specific state. For example, if the machine is in state q_1 (it has just seen an i) and the next character is another i, the machine would stay in state q_1 , since this new i could be the start of a new ing sequence. This corresponds to a transition $\delta(q_1,i)=q_1$. In contrast, if it sees a z, the sequence is broken, so the transition might be $\delta(q_1,z)=q_0$. A complete DFA specifies a transition for every state and every symbol in the alphabet.

To understand how a DFA operates, let's trace the execution of a machine designed to recognize a simplified 'sheep language.' This machine accepts strings that match the regular expression ba+!, such as ba!, baa!, and so on.

Let's formally define this DFA, M, with the 5-tuple $(Q, \Sigma, \delta, q_0, F)$:

- $Q = \{q_0, q_1, q_2, q_3, q_{dead}\}$ are the states.
- $\Sigma = \{a, b, !\}$ is the alphabet.
- q_0 is the start state.
- $F = \{q_3\}$ is the set of final (or accepting) states.
- δ is the transition function, defined as follows:
 - $-\delta(q_0,b) \to q_1$
 - $-\delta(q_1,a)\to q_2$
 - $-\delta(q_2,a)\to q_2$
 - $-\delta(q_2,!) \to q_3$
 - Any other transition leads to q_{dead} . For example, $\delta(q_0, a) \to q_{dead}$ or $\delta(q_1, b) \to q_{dead}$. Once in q_{dead} , the machine stays there for any input.

The machine works by reading an input string one character at a time from left to right, updating its current state according to the transition function δ . A string is accepted if the machine is in a final state after reading the *entire* string. Otherwise, it is rejected.

Example 1: Tracing an accepted string baa!

- 1. The machine begins in the start state, q_0 .
- 2. It reads the first character, b. Following the transition $\delta(q_0, b)$, it moves to state q_1 .
- 3. It reads the next character, a. Following $\delta(q_1, a)$, it moves to state q_2 .
- 4. It reads the third character, a. Following $\delta(q_2, a)$, it remains in state q_2 .
- 5. It reads the final character, !. Following $\delta(q_2,!)$, it moves to state q_3 .

The machine has now consumed the entire string and is in state q_3 . Since q_3 is in our set of final states F, the string baa! is accepted.

Example 2: Tracing a rejected string ba

- 1. The machine starts in q_0 .
- 2. It reads b and moves to state q_1 .
- 3. It reads a and moves to state q_2 .

The input string is now exhausted. The machine's final state is q_2 . Since q_2 is not in the set of final states F, the string ba is **rejected**. This example shows that simply passing through an accepting state is not enough; the machine must end in one.

Example 3: Tracing a rejected string aba!

- 1. The machine starts in q_0 .
- 2. It reads the first character, a. Our transition function does not have a defined path for a from q_0 , except to the 'dead' state. So, the machine moves to q_{dead} .
- 3. For every subsequent character (b, !), the machine will remain in q_{dead} .

At the end of the string, the machine is in state q_{dead} . Since $q_{dead} \notin F$, the string aba! is **rejected**. This demonstrates the deterministic nature of the DFA: for any state and any input symbol, there is exactly one move to make.

While deterministic finite automata provide a clear, unambiguous model of computation, they are not always the most intuitive tool for designing a language recognizer. For any given state and input symbol, a DFA has exactly one choice. This rigidity can lead to complex automata for relatively simple languages. To overcome this, we introduce a more flexible model: the *non-deterministic finite automaton* (NFA).

Non-determinism can be thought of as the power of 'choice' or 'parallel exploration.' An NFA can have multiple possible next states for a given input symbol. Furthermore, it can change state even without consuming any input, a feature known as an *epsilon transition* (ϵ -transition). These two features are the core of what makes NFAs non-deterministic.

Formally, an NFA is also a 5-tuple $(Q, \Sigma, \delta, q_0, F)$, but its transition function, δ , is defined differently.

- Q: a finite set of states.
- Σ : a finite input alphabet.
- $q_0 \in Q$: the start state.
- $F \subseteq Q$: the set of final states.
- $\delta: Q \times (\Sigma \cup \{\epsilon\}) \to 2^Q$: the transition function.

The crucial difference lies in δ . The function takes a state and a symbol from the alphabet (or the special symbol ϵ) and returns a *set* of possible next states. The notation 2^Q represents the power set of Q—the set of all possible subsets of Q. If, from state q_i on input 'a', the machine can transition to either q_j or q_k , we would write $\delta(q_i, a) = \{q_j, q_k\}$. If there is no transition from q_i on input 'a', then $\delta(q_i, a) = \emptyset$ (the empty set).

How does an NFA process an input string? Since it can be in multiple states at once, we can imagine it exploring all possible paths simultaneously. When the NFA reads an input symbol, it follows all available transitions for that symbol from all of its current states. The machine *accepts* a string if, after the entire string has been consumed, *at least one* of these parallel paths ends in a final state. It's as if the machine can magically guess the correct path to take to an accepting state if one exists.

Consider the NFA in **Fig. 2.5**, which recognizes the language of the regular expression a*b*. The automaton has a start state q_0 and a final state q_1 . From q_0 , there is a loop back to itself on the input 'a', allowing it to process any number of 'a's. The key non-deterministic feature is an ϵ -transition from q_0 to q_1 . This transition allows the machine to move from the 'a'-processing state to the 'b'-processing state without reading any input. Finally, state q_1 has a loop back to itself on input 'b', allowing it to process any number of 'b's.

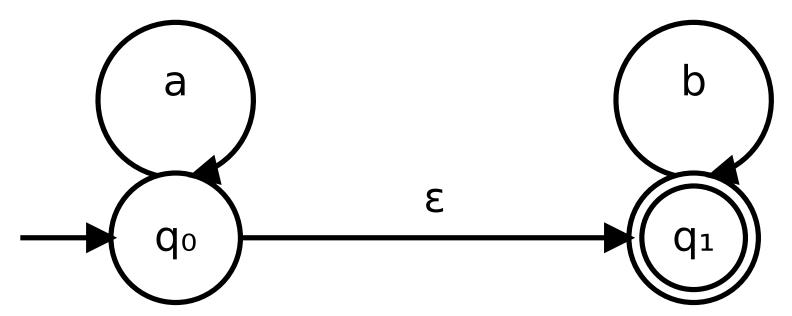


Figure 2.5: A Non-deterministic Finite Automaton (NFA) for the language described by the regular expression a^*b^* . State q_0 is the start state, and q_1 is the final state. The ϵ -transition allows the machine to move from the 'a'-processing state to the 'b'-processing state without consuming input.

Let's trace the input 'ab'. We start in state q_0 . On reading 'a', we stay in q_0 . Now, the machine is in state q_0 with 'b' left to read. It must first decide whether to take the ϵ -transition. In essence, we track the set of all possible current states. Before reading 'b', the machine could be in q_0 or, by following the ϵ -transition, it could also be in q_1 . From this set of states $\{q_0, q_1\}$, we read 'b'. There is no 'b' transition from q_0 , so that path dies. There is a 'b' transition from q_1 to itself, so that path continues. After reading 'b', the machine is in the set of states $\{q_1\}$. Since the input is consumed and the set of active states contains a final state (q_1) , the string 'ab' is accepted.

The primary advantage of NFAs is that they are often far simpler and more compact than their deterministic counterparts. Designing an NFA to recognize a language described by a regular expression is frequently a more direct and intuitive process. This simplicity is a significant benefit, even if the underlying model of computation seems more abstract. As we will see, despite their apparent power of 'guessing,' NFAs are not more powerful than DFAs; they recognize the exact same class of languages—the regular languages. The next sections will prove this fundamental equivalence.

We have now introduced two formalisms: regular expressions, a compact notation for describing string patterns, and finite automata, an abstract machine for recognizing them. The connection between these two is not a coincidence; it is a deep and fundamental result in computer science, formalized by Kleene's theorem.¹

Kleene's theorem states that any language that can be described by a regular expression can be recognized by a finite automaton, and vice-versa. This establishes a perfect equivalence in expressive power. The class of languages that can be defined by either of these formalisms is known as the **regular languages**. This equivalence is a cornerstone of formal language theory and has immense practical implications for computational linguistics.

The theorem can be understood as a two-part statement:

¹Kleene, Stephen C. (1956). 'Representation of events in nerve nets and finite automata'. In C. Shannon and J. McCarthy (eds.), *Automata Studies*. Princeton University Press.

- Part 1: From Regular Expression to Automaton. For any regular expression r, there exists an equivalent non-deterministic finite automaton (NFA) that accepts the language L(r). This is the constructive direction that powers practical applications. It gives us a blueprint for taking a human-readable pattern and compiling it into an executable machine.
- Part 2: From Automaton to Regular Expression. For any finite automaton M (either deterministic or non-deterministic), there exists a regular expression r that describes the exact same language, such that L(r) = L(M). This completes the proof of equivalence, showing the mapping works in both directions.

This powerful connection means that we can use the concise, declarative syntax of regular expressions to specify complex patterns, while relying on the efficient, well-understood mechanics of automata to perform the actual string processing. When a programmer writes a regular expression, the system's underlying engine can compile it into a highly optimized automaton to match text. This provides a guaranteed bridge between a formal description and its practical implementation. In the following sections, we will demystify this process by examining the standard algorithms—Thompson's construction and the subset construction—that make this conversion possible.

The equivalence between regular expressions and finite automata, established by Kleene's theorem, is not merely a theoretical curiosity. It provides a practical pathway for implementation: we can convert a regular expression, which is easy for humans to write, into an NFA, which is a straightforward computational model. Thompson's construction is a classic and elegant algorithm that accomplishes exactly this. The algorithm is recursive, meaning it defines how to build an NFA for a complex expression by first building NFAs for its simpler sub-expressions.

The process starts with two base cases for the simplest possible regular expressions:

- A single symbol: For a regular expression consisting of a single symbol a from the alphabet Σ , we construct an NFA with a new start state and a new final state, linked by a single transition labeled a.
- The empty string: For the empty string, ϵ , we construct an NFA with a new start and final state, linked by a single ϵ -transition.

From these atomic components, we build larger NFAs using rules that directly mirror the operators in the regular expression. Let's assume we have already constructed NFAs for the sub-expressions s and t, which we will call N(s) and N(t).

- 1. Union (s|t): To construct the NFA for the union s|t, we create a new start state and a new final state. We then add ε-transitions from the new start state to the original start states of both N(s) and N(t). We also add ε-transitions from the original final states of N(s) and N(t) to our new final state. The original final states are no longer marked as final. This construction creates a machine that can non-deterministically choose to traverse either the path for s or the path for t.
- 2. Concatenation (st): For the concatenation st, the start state of N(s) becomes the overall start state, and the final state of N(t) becomes the overall final state. We connect the two machines by adding an ε-transition from the final state of N(s) to the start state of N(t). The final state of N(s) is then unmarked. This forces the machine to complete the path for s before beginning the path for t.
- 3. **Kleene Star** (s^*): To construct the NFA for s^* , we again create a new start state and a new final state. We add an ϵ -transition from the new start state directly to the

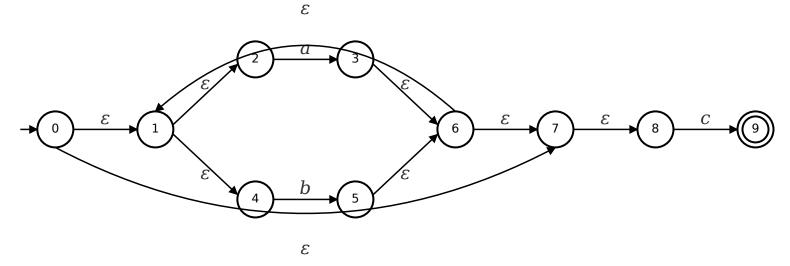


Figure 2.6: A diagram illustrating Thompson's construction for the regular expression (a|b)*c. The NFA is built systematically by first creating a machine for the union 'a|b', then applying the Kleene star rule to create a machine for '(a|b)*', and finally applying the concatenation rule to join it with the machine for 'c'. Each step in the construction creates a new NFA with a single start and final state from the components of its sub-expressions, connected by ϵ -transitions.

new final state to handle the case of zero occurrences of $\mathfrak s$. We also add ϵ -transitions from the new start state to the start state of $\mathbb N(\mathfrak s)$ and from the final state of $\mathbb N(\mathfrak s)$ to the new final state. Crucially, to allow for one or more repetitions, we add a 'loop-back' ϵ -transition from the final state of $\mathbb N(\mathfrak s)$ to its own start state.

The power of this method lies in its compositionality. As illustrated in **Fig. 2.6**, to build an NFA for a complex expression like (a|b)*c, we simply apply these rules recursively. We first construct the NFAs for a and b, combine them with the union rule to get N(a|b), apply the Kleene star rule to this result to get N((a|b)*), and finally use the concatenation rule to join it with N(c). The resulting NFA has a predictable structure—always with a single start state and a single final state—making it a systematic and reliable method for converting any regular expression into a functionally equivalent NFA.

While Non-deterministic Finite Automata (NFAs) offer great flexibility for design, their operational mechanics can be inefficient to simulate directly. A computer must track multiple possible paths simultaneously for a given input string. In contrast, Deterministic Finite Automata (DFAs) are computationally straightforward, as there is only one possible path for any input. Fortunately, a fundamental algorithm known as **subset construction** provides a systematic way to convert any NFA into an equivalent DFA that recognizes the same language. This process is crucial because it proves that NFAs, despite their non-determinism, do not possess any more computational power than DFAs and guarantees that we can always build an efficient deterministic machine from a simpler non-deterministic design.

The core idea behind subset construction is that each state in the new DFA corresponds to a *set of states* from the original NFA. A DFA state represents the set of all possible states the NFA could be in after processing a particular sequence of input symbols. The algorithm works as follows:

1. Create the DFA start state: The start state of the DFA is the set containing the NFA's start state and any other states reachable from it via ϵ -transitions (the ϵ -closure).

| DFA State | NFA State Set | Resu | lting DFA State on Input |
|-----------------|----------------|----------------|--------------------------|
| D | ${\mathcal S}$ | 'a' | 'b' |
| $\rightarrow A$ | $\{q_0, q_1\}$ | \overline{A} | В |
| B^* | $\{q_2\}$ | B | A |

Figure 2.7: A trace of the subset construction algorithm converting a sample NFA to a DFA. The start state is indicated by \rightarrow , and final states by an asterisk (*). Each new DFA state corresponds to a set of states from the original NFA.

- 2. Iteratively create new states and transitions: For each newly created DFA state D_{new} and for each symbol c in the alphabet Σ :
 - Find the set of all states in the NFA that can be reached from any state in D_{new} by following a transition on symbol c.
 - Compute the ϵ -closure of this resulting set. This new set of NFA states becomes a state in our DFA.
 - Add a transition in the DFA from D_{new} to this new state on symbol c.
- 3. **Repeat until convergence:** Continue this process until no new DFA states are generated.
- 4. **Define final states:** Any state in the DFA that contains at least one of the original NFA's final states becomes a final state in the new DFA.

This procedure can seem abstract, but it becomes clear when traced. The table in Fig. 2.7 provides a step-by-step trace of the subset construction algorithm applied to a sample NFA. Each row corresponds to a newly discovered DFA state, showing the set of NFA states it represents and how its transitions are computed for each symbol in the alphabet. By following the table, you can see how the algorithm systematically explores all reachable combinations of NFA states, creating a finite and deterministic map of the NFA's behavior. The subset construction algorithm is a cornerstone result, proving the equivalence of NFAs and DFAs and providing a practical method for turning conceptual models into efficient implementations.

Finite-state automata are excellent tools for *recognition*. They provide a formal mechanism to answer a simple yes/no question: 'Is this string a member of the language we have defined?' While this is a fundamental capability, many tasks in computational linguistics require more. We often need to *transform* one string into another, not just accept or reject it.

This process of mapping an input string to a corresponding output string is called transduction. Consider the task of morphological analysis, where we might want to generate a plural noun from its singular form (e.g., mapping goose to geese) or translate a word from one language to another. An FSA, by its nature, cannot perform this task. To model such relationships, we extend the finite automaton to a more powerful machine: the Finite-State Transducer (FST).

An FST is conceptually very similar to a finite automaton, with one crucial enhancement: each transition is augmented with an output symbol. Where an FSA transition from one state to another is labeled with an input symbol, an FST transition is labeled with an *input:output* pair. As the FST consumes an input string by moving from state to state, it simultaneously generates an output string by appending the output symbol from each traversed arc. The final output is the sequence of symbols generated along a successful path from the start state to a final state.

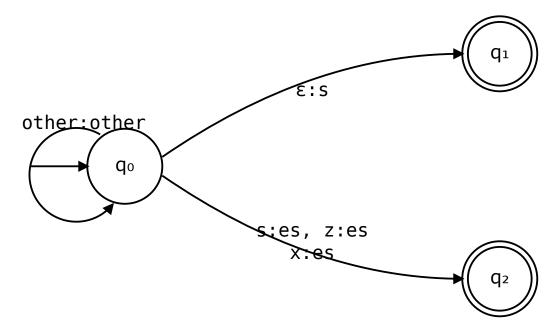


Figure 2.8: A state diagram of a Finite-State Transducer (FST) for English noun pluralization, corresponding to Fig. 2.17. The diagram shows a start state (q_0) with a self-loop for most characters (labeled 'other:other'). From this state, a transition on no input (ϵ) outputs an 's' to form regular plurals like 'cats', leading to the final state q_1 . A separate path handles nouns ending in 's', 'z', or 'x'; it transduces these characters to 'es' and leads to a second final state, q_2 , correctly forming plurals like 'foxes'.

Formally, an FST is like an NFA, but its transition function δ also produces output. For an input alphabet Σ and an output alphabet Γ , the transition function maps a state and an input symbol to a set of pairs, where each pair contains a new state and an output string:

$$\delta: Q \times (\Sigma \cup \{\epsilon\}) \to 2^{Q \times \Gamma^*}$$

This definition allows a single input to map to multiple possible outputs, or for output to be produced even on an ϵ -input, making FSTs a flexible tool for modeling complex string-to-string relationships. Their utility in computational linguistics is immense, especially for tasks involving morphology. They can elegantly capture rules for inflectional morphology, such as pluralization and verb conjugation, and derivational morphology, like adding prefixes and suffixes. We will now explore a practical case study of this application.

Let's consider a practical application of Finite-State Transducers: modeling English noun morphology. While FSAs can recognize valid words, FSTs can be used to transform them, making FSTs a natural fit for tasks like inflectional morphology, which involves changing a word's form to reflect grammatical features like number. The formation of noun plurals in English is a classic example that, while seemingly simple, is governed by a set of well-defined rules.

Most nouns form their plural by simply adding an -s (e.g., $cat \rightarrow cats$). However, nouns ending in sounds like /s/, /z/, or $/\int/$ (spelled with s, z, x, sh, ch) require an -es suffix (e.g., $fox \rightarrow foxes$). We can capture this logic elegantly in an FST, as shown in Fig. 2.8. An FST for pluralization maps an input string (the singular noun) to an output string (the plural form).

Consider the input word cat. The transducer would process this by following a path where each input character is mapped to itself as output. Upon reaching a final state after consuming the entire word, a final transition is taken that reads no input symbol (ϵ) but

outputs an s. The full transduction is thus c:c, a:a, t:t, ϵ :s, mapping the input cat to the output cats.

Now, consider the input fox. The machine again maps f:f and o:o. However, upon reading the final x, the machine is designed to follow a different path. Instead of a simple x:x transition, it takes a specific transition labeled x:es. This single transition handles the morphological rule for this class of nouns, correctly producing the output foxes.

This simple FST demonstrates the power of this formalism. It encodes multiple linguistic rules into a single, efficient computational model. Different paths through the machine correspond to different morphological processes. While a comprehensive FST for English morphology would be far more complex, incorporating rules for words ending in -y ($fty \rightarrow flies$) and handling irregular forms ($mouse \rightarrow mice$), this case study illustrates the fundamental utility of FSTs in computational morphology.

Despite their utility for modeling word-level phenomena, regular expressions and finite automata have a crucial limitation: their finite memory. An FSA can only remember which of its finite set of states it is currently in. This is sufficient for recognizing local patterns, but it proves inadequate for capturing certain long-distance dependencies found in the syntax of many natural languages.

A classic example of a structure that regular languages cannot handle is *center embedding*, where a phrase is placed in the middle of another phrase of the same type. Consider the following sentences:

- The cat meowed.
- The cat the dog chased meowed.
- The cat the dog the rat bit chased meowed.

To verify the grammatical correctness of these sentences, a model must match each noun phrase subject ('the cat', 'the dog', 'the rat') with its corresponding verb ('bit', 'chased', 'meowed'). As the depth of embedding increases, the model must 'remember' an increasing number of subjects before it encounters their verbs. An FSA, with its fixed number of states, cannot store an unbounded number of pending subjects.

Formally, this problem is equivalent to recognizing the language $L = \{a^nb^n | n \geq 0\}$. This language consists of some number of a's followed by the exact same number of b's. No finite automaton can recognize this language. To do so, it would need to count the number of a's, which could be infinite. This would require an infinite number of states, contradicting the definition of an FSA. This fundamental inability to count and store unbounded information means we require more powerful formalisms, like context-free grammars, to model the hierarchical structure of sentences.

In this chapter, we have laid the formal and computational groundwork for processing language at its most fundamental level: the word. We began with regular expressions, a practical and powerful tool for describing and matching patterns in text. We then delved into their theoretical counterpart, the finite-state automaton, exploring both deterministic (DFA) and non-deterministic (NFA) variants. The essential takeaway, formalized in Kleene's theorem, is that these two perspectives are equivalent in their descriptive power; anything a regular expression can describe, a finite automaton can recognize, and vice versa.

We extended this model from simple recognition to transformation using Finite-State Transducers (FSTs), demonstrating their utility in tasks like morphological analysis. However, the power of this machinery has clear boundaries. The finite-state model is fundamentally limited by its lack of memory. It cannot, for instance, recognize a language like $L = \{a^nb^n \mid n \geq 0\}$, which requires counting the number of a's to ensure an equal number of b's.

This limitation is not merely a theoretical curiosity. It means that regular languages cannot model many syntactic phenomena in human language, such as the nested structures found in center-embedded clauses. To capture the hierarchical nature of sentences, we need a more powerful class of grammars. Having mastered the tools for modeling words, we are now prepared to move up the linguistic hierarchy. The subsequent chapters will build upon this foundation, introducing the formalisms necessary to parse phrases, clauses, and full sentences, and thereby unlock a deeper computational understanding of syntactic structure.

Chapter 3

Corpus Linguistics and Text Normalization



The evolution of computational linguistics is marked by a profound paradigm shift, away from the purely theoretical and towards the empirically grounded. Early efforts in the field often followed a rationalist tradition, where linguists and computer scientists attempted to hand-craft intricate sets of rules to capture a language's grammar. The goal was to build a logical model based on linguistic theory and introspection. This approach, however, frequently stumbled when faced with the immense complexity and variability of real-world language. The resulting systems were often brittle, unable to handle exceptions, idiomatic expressions, or the constant evolution of language use. They lacked the flexibility to learn from experience.

The modern, data-driven era is built on a different philosophy. The central tenet is simple yet powerful: to create computational models that understand human language, we must learn its properties directly from vast quantities of it as it is actually written and spoken. This represents a turn towards an *empirical* methodology, where direct observation of language data is the primary source of knowledge. Rather than starting with abstract rules, we start with the data itself and use statistical methods and machine learning to discover the patterns, frequencies, and relationships that govern communication. The guiding question is not 'What *should* the rule be?' but 'What does the evidence in the data suggest?'

This empirical foundation is made tangible through the **linguistic corpus**. A corpus (Latin for 'body'; plural *corpora*) serves as the bedrock for nearly all contemporary NLP. It is not just any random collection of text, but a large, structured, and principled sample of authentic language use, curated for analysis. The rise of the internet and the explosion in digital text, combined with the availability of powerful computing resources, made it possible to compile and process corpora on a scale previously unimaginable. This synergy of data and computation is the engine behind the field's most significant breakthroughs. Instead of a few thousand hand-analyzed sentences, models can now learn from billions of words, enabling them to capture subtle statistical nuances of language that would be impossible to codify by hand. From machine translation to sentiment analysis, the algorithms we will explore in this book all derive their power from learning from a corpus. This chapter, therefore, is dedicated to this cornerstone of our field, exploring what a corpus is and how we prepare its raw text for analysis.

At its core, a linguistic corpus is far more than just a large collection of text. To be useful for computational analysis, a corpus must be a *principled* collection. We formally define a corpus as a large, structured, and machine-readable collection of authentic text or speech samples, which has been assembled according to explicit design criteria to be representative of a specific language or language variety. Each component of this definition is critical and distinguishes a corpus from a simple accumulation of digital text.

The requirement for the collection to be *large* is a direct consequence of the statistical nature of modern computational linguistics. Language is rife with phenomena that are relatively rare; a small sample of text might not contain enough examples of a particular word, phrase, or grammatical construction to allow for meaningful statistical analysis or the robust training of a machine learning model. The samples must also be *authentic*, meaning they consist of naturally-occurring language produced by real speakers or writers for a genuine communicative purpose. This is a crucial distinction from collections of sentences invented by linguists to illustrate a specific point. We study language as it is actually used, so we draw from real-world sources like news articles, fiction, academic papers, and transcribed conversations.

Perhaps the most important characteristics are that a corpus is *principled* and *structured*. A principled collection is one built with a clear purpose and methodology, unlike a random assortment of documents downloaded from the internet. The design principles often aim for the corpus to be *balanced* and *representative* of the language variety under

study. For instance, a corpus of 'general American English' would not just contain news articles, but would be carefully sampled to include a proportional mix of different genres like fiction, conversation, and academic prose. The term *structured* refers to the fact that the data is organized in a consistent, machine-readable format. At a minimum, this includes consistent character encoding and document boundaries. More advanced corpora contain rich metadata, such as the author, publication date, and genre for each document, which is essential for many types of analysis.

Ultimately, a corpus serves as the empirical bedrock for computational linguistics. It is the digital laboratory where we can observe linguistic phenomena, test hypotheses, and gather frequency data. For the data-driven models that dominate the field today, from N-gram models to large neural networks, the corpus is the source of the training and evaluation data that allows them to learn the patterns of human language. Without well-designed corpora, the quantitative, evidence-based study of language at scale would be impossible.

Not all corpora are created equal. The design, content, and structure of a corpus are dictated by the research questions it aims to answer or the NLP application it is intended to support. This leads to several distinct types of corpora, each with its own strengths and typical use cases. The most common distinctions are based on the breadth of coverage (balanced vs. specialized), the number of languages included (monolingual vs. parallel), and the presence of linguistic metadata (raw vs. annotated).

A balanced corpus aims to be a representative snapshot of a particular language or language variety. To achieve this balance, its creators carefully sample texts from a wide array of genres and domains—such as news reports, fiction, academic articles, and transcripts of spoken conversations. The goal is to mirror the proportions of these genres in the wider language, preventing any single text type from dominating. This makes balanced corpora invaluable for general linguistic inquiry or for training general-purpose language models that are not tailored to a specific task. Classic examples include the pioneering Brown Corpus for American English and the larger British National Corpus (BNC).

In contrast, a specialized corpus (or domain-specific corpus) deliberately focuses on a narrow slice of language. Instead of broad coverage, it offers deep coverage of a particular subject matter, genre, or communicative setting. For example, a corpus might consist solely of legal contracts, biomedical research abstracts, or social media posts. This allows a model to learn the specific vocabulary, jargon, and stylistic conventions of its domain, which is essential for building high-performance applications like a medical information extraction system or a financial sentiment analyzer.

Another crucial type is the parallel corpus. This is a multilingual resource that contains the same texts in two or more languages. The key feature of a parallel corpus is alignment, where sentences or segments in one language are explicitly linked to their translations in the other(s). They are the cornerstone of statistical and neural machine translation, as they provide the raw data from which a system can learn translational equivalences between words and phrases. Prominent examples include the Europarl Corpus, which contains proceedings from the European Parliament, and the Canadian Hansards, which document parliamentary debates in English and French.

The aforementioned types primarily describe the *content* of a corpus. An orthogonal and critically important dimension is whether a corpus is *annotated*. An annotated corpus contains not just the raw text, but also explicit linguistic information added by human experts. This annotation can vary in complexity:

- Part-of-speech (POS) tags: Each word is labeled with its grammatical category (e.g., noun, verb).
- Syntactic annotation: Sentences are parsed into tree structures representing their

| Corpus Type | Description | Primary Use Case | Examples |
|-------------|--|--|--|
| Balanced | Aims to be a representative snapshot of a language by sampling texts from a wide array of genres and domains. | General linguistic inquiry; training general-purpose lan- guage models. | Brown Corpus, British National Corpus (BNC) |
| Specialized | Deliberately focuses on a narrow subject matter or genre to provide deep, domain-specific coverage. | Building high-performance, domain-specific NLP applications (e.g., medical information extraction). | A corpus of legal contracts, biomed- ical research ab- stracts, or social media posts. |
| Parallel | Contains the same texts in two or more languages, with sentences or segments explic- itly aligned to their transla- tions. | Training and evaluating statistical and neural machine translation systems. | Europarl Corpus, Canadian Hansards. |
| Annotated | Contains raw text enriched with explicit linguistic metadata (e.g., POS tags, syntactic trees) added by human experts. | Providing 'ground truth' data for training and evaluating supervised machine learning models. | Penn Treebank |

Figure 3.1: A concise summary of major corpus types, their primary use cases, and well-known examples.

grammatical constituency or dependency relations. A corpus with such annotation is often called a *treebank*.

• **Semantic annotation:** Words or phrases are tagged with information about their meaning, such as their role in an event (semantic role labeling).

Annotated corpora are fundamental for supervised machine learning, as they provide the 'ground truth' data needed to train and evaluate models for specific tasks like POS tagging or parsing. The influential Penn Treebank is a prime example, containing detailed POS and syntactic annotation.

These categories are not mutually exclusive; a corpus can be, for instance, an annotated, specialized, parallel corpus. The choice of which type to use is one of the most important decisions in any NLP project, as the data ultimately determines the capabilities and limitations of the resulting model. Fig. 3.1 provides a concise summary of these major corpus types, their primary use cases, and well-known examples.

To make the concept of a balanced corpus concrete, we can examine one of the most influential early examples: the *Brown University Standard Corpus of Present-Day American English*, universally known as the Brown Corpus. Compiled in the 1960s by Henry Kučera and W. Nelson Francis at Brown University, it was a landmark achievement in linguistics. Its primary objective was to create a one-million-word sample of American English text that was broadly representative of the language as it was used in a single year, 1961. This made it the first truly balanced, computer-readable corpus of its kind, setting a standard for decades of future work in the field.

The corpus consists of 500 text samples, or *documents*, each approximately 2,000 words long. These samples were meticulously drawn from 15 distinct genre categories to ensure a wide-ranging representation of written English. As shown in **Fig. 3.2**, these categories span a diverse spectrum, including *Press: Reportage*, *Religion*, *Skills and Hobbies*, various genres of *Fiction*, and academic prose from the *Learned* and *Scientific* categories. The

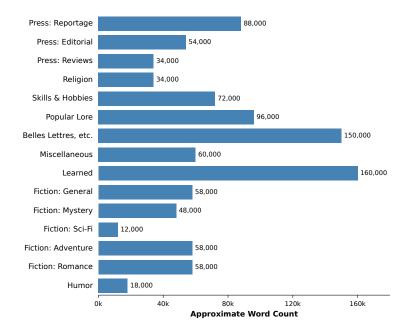


Figure 3.2: Composition of the Brown Corpus, showing the approximate number of words from each of the 15 genre categories. This distribution illustrates the corpus's design as a balanced, representative sample of 1961 American English.

creators ensured that the number of samples from each category was proportional to its prevalence in the publishing world at the time, providing a snapshot of the literary landscape.

This careful stratification is the essence of a balanced corpus. The goal is not merely to amass text, but to create a microcosm of the language as it is actually produced and consumed. By sampling from romance fiction as well as from government documents, the Brown Corpus allows researchers to make more generalizable claims about word frequencies, syntactic structures, and lexical patterns across the language as a whole, rather than just within a single, narrow domain. For instance, one could reliably compare the average sentence length in journalistic writing versus scientific prose using this data.

Although a million words is considered small by today's standards—where corpora often contain billions of words—the methodological rigor of the Brown Corpus was revolutionary. It provided the empirical data for some of the first large-scale computational analyses of English vocabulary and grammar. Its design principles have been replicated and adapted for creating similar corpora in many other languages, a testament to its enduring influence. It demonstrated that a principled, data-driven approach was essential for the scientific study of language.

Building a corpus is a more deliberate and structured process than simply downloading a large amount of text. The goal is to create a resource that is balanced, representative, and suitable for the research questions it is intended to answer. This construction process generally follows a series of essential steps, as illustrated in the workflow in **Fig. 3.3**.

The first stage is data acquisition. The source of the text depends entirely on the corpus's purpose. For a corpus of contemporary news, one might use web scraping to gather articles from various online newspapers. For a historical corpus, the process might involve digitizing physical books and manuscripts using Optical Character Recognition (OCR). Other sources include parliamentary records, social media feeds, or transcribed speech. A crucial part of acquisition is sampling. Rather than collecting all available data, which may be impractical or introduce bias, we select a representative sample. For a balanced corpus like the Brown Corpus, this meant carefully selecting texts from 15

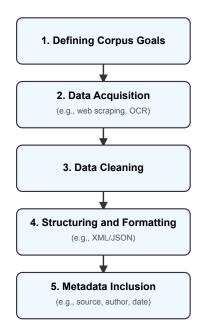


Figure 3.3: A flowchart showing the key stages of corpus construction, from defining goals and acquiring data to cleaning, formatting, and including metadata.

different genres to ensure the sample reflected the broad range of published American English at the time.

Once the raw data is collected, it must undergo extensive data cleaning. Text from the web is often riddled with non-linguistic noise, such as HTML tags, advertisements, navigation menus, and JavaScript code, all of which must be stripped away. Scanned documents may contain OCR errors that need correction. This stage is often labor-intensive but is critical for the quality of the final corpus; the principle of 'garbage in, garbage out' applies with full force.

After cleaning, the text is structured and formatted, often using standards like XML (eXtensible Markup Language) or JSON (JavaScript Object Notation). This provides a consistent format and, more importantly, allows for the inclusion of *metadata*—literally, 'data about the data.' Metadata provides essential contextual information for each document. For a newspaper article, metadata might include:

- The publication source (e.g., The New York Times)
- The publication date
- The author's name
- The article's section (e.g., 'Business,' 'Sports')

This structured information is invaluable. It allows researchers to not only analyze the language itself but also to investigate how language use varies across different sources, time periods, or genres. By following these steps, a simple collection of raw text is transformed into a powerful and scientifically valuable linguistic corpus.

Raw text, while useful, is merely a sequence of characters. For a computer to perform sophisticated linguistic analysis, this raw data must be enriched with explicit, structured linguistic information. The process of adding this expert-level information to a corpus is known as *annotation* or *markup*. Annotation transforms a simple collection of texts into a powerful, machine-readable resource that can be used to train and evaluate computational models. It provides the 'ground truth' data that data-driven algorithms learn from. For

example, to build a system that automatically identifies nouns and verbs, we first need a corpus where human experts have already labeled thousands of examples of nouns and verbs.

The annotation process is often a meticulous and costly undertaking, typically performed by trained linguists following a detailed set of instructions known as an annotation guideline. To ensure the quality and consistency of the labels, it is common practice to have multiple annotators label a portion of the same data independently. The level of consensus between them is then measured using a metric called *inter-annotator agreement (IAA)*. A high IAA score indicates that the annotation guidelines are clear and the resulting labels are reliable, which is crucial for building robust NLP systems.

Annotation can be performed at various linguistic levels, with the specific choice of what to mark up—the *annotation schema*—depending on the intended use of the corpus. Some of the most common types of annotation include:

- Part-of-Speech (POS) Tagging: This is one of the most fundamental layers of annotation. Each word or token in the text is assigned a tag indicating its grammatical category. For instance, the sentence 'The cat sat on the mat' might be annotated as: The/DT cat/NN sat/VBD on/IN the/DT mat/NN ./. Here, DT stands for determiner, NN for noun, VBD for past-tense verb, and IN for preposition.
- Syntactic Annotation (Parsing): This layer reveals the grammatical structure of sentences. One popular approach is *constituency parsing*, which groups words into nested phrases. The same sentence could be represented with brackets to show its phrasal structure: (S (NP (DT The) (NN cat)) (VP (VBD sat) (PP (IN on) (NP (DT the) (NN mat))))). This annotation shows that 'The cat' is a Noun Phrase (NP) and 'sat on the mat' is a Verb Phrase (VP), which together form the sentence (S).
- Semantic Annotation: This involves marking up meaning-related information. This can range from word sense disambiguation, where a word like 'bank' is tagged with its specific meaning (e.g., river bank vs. financial institution), to named entity recognition, where mentions of people, organizations, and locations are identified. Another crucial type is coreference resolution, which links expressions that refer to the same entity, such as connecting 'Dr. Evans,' 'she,' and 'the professor' throughout a text.

Ultimately, an annotated corpus is an embodiment of linguistic data tailored for a specific computational task. The richness and reliability of the annotations directly determine the potential and limitations of the models trained on it. This deep connection between linguistic theory, annotation practice, and computational modeling is exemplified by seminal projects like the Penn Treebank.

While the concept of annotation can be applied simply, such as labeling documents for sentiment, some corpora feature incredibly rich, multi-layered annotation schemas. Perhaps the most influential example of such a resource in the history of computational linguistics is the **Penn Treebank (PTB)**. Developed at the University of Pennsylvania, the PTB project annotated a large corpus, primarily consisting of over a million words from *Wall Street Journal* articles, with detailed linguistic information. Its release provided a large, consistent, and publicly available dataset that became the standard benchmark for developing and evaluating a generation of NLP models, particularly syntactic parsers.

The annotation in the Penn Treebank is applied in layers. The first and most fundamental layer is **Part-of-Speech (POS) tagging**. Every single word in the corpus is assigned a grammatical tag based on a detailed tagset. The PTB tagset uses 36 distinct

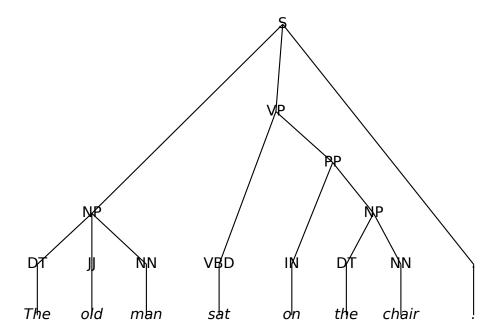


Figure 3.4: A Penn Treebank syntactic parse tree for the sentence 'The old man sat on the chair.' The diagram shows the hierarchical constituent structure, with nodes for phrasal categories (S, NP, VP, PP), part-of-speech tags (DT, JJ, NN, VBD, IN), and the individual words.

POS tags and 12 other tags for punctuation and currency symbols. This level of granularity is crucial for computational analysis; for instance, the tagset distinguishes between a singular noun (NN), a plural noun (NNS), a proper noun (NNP), a base form verb (VB), a verb in the past tense (VBD), and a verb in the present participle form (VBG). This fine-grained information is far more useful for a computational model than simply labeling all of these as 'noun' or 'verb'.

The second and most famous layer of annotation is the **syntactic treebank**, from which the project gets its name. On top of the POS tags, the sentences were parsed into their constituent structures, creating a 'tree' that represents the sentence's grammatical organization. This type of annotation is known as a *constituency parse* or *phrase-structure parse*. It hierarchically groups words into phrases, which are in turn grouped into larger phrases, until the entire sentence is accounted for. Common phrasal labels include NP (Noun Phrase), VP (Verb Phrase), PP (Prepositional Phrase), and ADJP (Adjective Phrase), with S representing the entire sentence clause.

A visual representation of this structure for the sentence 'The old man sat on the chair' is shown in **Fig. 3.4**. At the lowest level of the tree, each word is connected to its POS tag (e.g., 'man' is a NN, 'sat' is a VBD). These tagged words form the *terminal nodes*. Above them, the tree groups words into phrasal constituents. 'The', 'old', and 'man' are grouped into a Noun Phrase (NP). Similarly, 'on the chair' forms a Prepositional Phrase (PP). These phrases are then nested within larger ones; the PP is part of the Verb Phrase (VP), 'sat on the chair'. Finally, the initial NP and the VP combine to form the complete sentence, labeled with the root symbol S. In the corpus files, this hierarchical structure is typically stored using a bracketed notation, like this:

(S (NP (DT The) (JJ old) (NN man)) (VP (VBD sat) (PP (IN on) (NP (DT the) (NN chair)))) (. .))

The influence of the Penn Treebank extends beyond these two layers. The same underlying text has since been enriched with further layers of annotation in subsequent projects. For example, **PropBank** added a layer of semantic role labels, identifying who did what

to whom. Later, projects like **OntoNotes** added annotations for coreference (linking pronouns like 'he' and 'it' to the entities they refer to), word senses, and named entities. This demonstrates how a well-designed base corpus can serve as a scaffold for building even richer linguistic resources over time. The legacy of the Penn Treebank is a testament to the power of high-quality, deeply annotated data in driving scientific progress and enabling the development of sophisticated language understanding technologies.

Having established the importance of corpora as the source of our language data, we now confront a practical challenge. Raw text, as it exists in the wild, is inherently 'noisy' and inconsistent from a computational standpoint. To a computer, which sees text as a simple sequence of characters, the strings run, Run, and running are entirely distinct entities. Similarly, abbreviations like U.S.A. and USA, or words connected by punctuation like state-of-the-art, are treated as unique items, different from their constituent parts. This surface-level variation, while trivial for a human reader to interpret, creates significant problems for our algorithms. If left unaddressed, it leads to a needlessly large vocabulary and exacerbates the problem of data sparsity, where many meaningful word forms appear too infrequently in the corpus for a statistical model to learn from them effectively.

This brings us to the indispensable process of **text normalization**, often referred to as text pre-processing. Text normalization is the task of cleaning raw text and transforming it into a more canonical, standardized format. The primary goal is to group together different surface forms of what is essentially the same conceptual word, ensuring they are treated consistently by downstream models. This process is best conceptualized as a *pipeline*, a sequence of operations where raw text is passed from one stage to the next. In the following sections, we will dissect the core components of this pipeline, starting with the fundamental task of **tokenization**—breaking text into words—before moving on to techniques like **stemming** and **lemmatization** for reducing words to their root forms. These foundational steps are the bedrock upon which nearly all subsequent NLP analysis is built.

After acquiring a raw text, the first and most fundamental step in preparing it for analysis is *tokenization*. This is the process of segmenting a continuous stream of characters into a sequence of discrete, meaningful units called *tokens*. For English and many other languages, these tokens often correspond to words, numbers, and punctuation marks. To a computer, a text file is simply a long sequence of bytes; it has no inherent understanding of where one word ends and another begins. Tokenization imposes the first layer of linguistic structure onto this raw data, transforming it from an undifferentiated string into a list of items that our algorithms can count, compare, and analyze.

Consider the simple sentence: The quick brown fox jumps. A standard tokenizer would break this string into the following list of six tokens:

```
['The', 'quick', 'brown', 'fox', 'jumps', '.']
```

Notice that the final period is treated as a separate token. This is a common and important choice, as it separates the word 'jumps' from the sentence-boundary marker. By isolating punctuation, we can analyze it as a distinct feature or choose to ignore it later, but the decision is preserved. Without tokenization, we would be left to work with individual characters or complex string-matching rules, making even simple tasks like counting word occurrences unnecessarily difficult. The sequence of tokens becomes the primary input for most subsequent NLP tasks, from part-of-speech tagging to parsing and language modeling.

This process also allows us to introduce two foundational concepts in corpus linguistics: tokens and types. A token is an instance of a word or symbol in the text. In the sentence, 'The cat sat on the mat.', there are seven tokens in total. A type, on the other hand, is a unique word or symbol in the vocabulary. The set of types for that sentence would be {The, cat, sat, on, the, mat, .}. If we are being case-sensitive, this set contains

seven unique types. If we ignore case, then 'The' and 'the' would be instances of the same type, reducing our type count to six. The collection of all unique types in a corpus is known as its **vocabulary**, often denoted as V. The distinction is crucial: the total number of tokens tells us the length of the corpus, while the number of types tells us the size of its vocabulary.

At first glance, tokenization might seem as simple as splitting a string by its whitespace characters. This naive approach, however, quickly breaks down when faced with the messiness of real-world text. Language is filled with ambiguities that challenge simple rules. For instance, how should we handle a date like 03/25/2024? Is it one token or five (03, /, 25, /, 2024)? What about a monetary value like \$45.50? Or a single contraction like 'don't', which could be treated as one token (don't) or two (do, n't) to separate the verb from the negation. As we will see next, successfully navigating these edge cases is crucial for building a reliable NLP pipeline.

While tokenization might seem as trivial as splitting a string by its whitespace characters, this simplistic approach quickly fails on real-world text. A simple programmatic split on the sentence 'She can't go.' would likely yield ['She', "can't", 'go.'], a result that is suboptimal for most analyses. The token go. incorrectly bundles the word with its sentence-terminating punctuation, and can't merges two distinct morphemes: the modal verb can and the negation particle not. A robust tokenizer must navigate a landscape of such linguistic ambiguities.

The primary challenges revolve around punctuation and multi-word expressions. Consider the following common issues:

- Contractions and Possessives: English contractions like can't, they're, or it's are a frequent hurdle. For analytical consistency, these are often split into their constituent parts. A standard approach, used in corpora like the Penn Treebank, is to separate the clitic. For instance, can't becomes two tokens, ca and n't, while it's becomes It and 's.¹ This separation preserves morphological information that is crucial for downstream tasks like part-of-speech tagging. Similarly, sentence-final punctuation, such as periods, question marks, and exclamation points, should almost always be treated as separate tokens.
- Ambiguous Punctuation: The period is particularly problematic. It can signal the end of a sentence, but it also appears within abbreviations (e.g., U.S., Dr.), numbers (e.g., 3.14), and URLs. A naive tokenizer that splits on every period would incorrectly break apart these meaningful units. A sophisticated tokenizer must use contextual clues to determine if a period is a sentence boundary or an internal character of a token.
- **Hyphenation:** Hyphens present another form of ambiguity. In a term like *state-of-the-art*, the hyphenated phrase acts as a single adjectival unit and is best kept as one token. However, hyphens can also be used for line-breaking (soft hyphens) or to connect words that do not form a permanent compound, such as in *a well-written but ultimately flawed argument*. The tokenizer's behavior for hyphenated words is often guided by a combination of pre-defined rules and lexical resources.

These examples illustrate that tokenization is not a solved problem but a series of language-specific engineering decisions. A well-designed tokenizer must implement a cascade of rules to correctly segment a text stream. The table in Fig. 3.5 provides several examples of challenging input strings and the desired token sequences, highlighting the nuanced logic required for this foundational step.

 $^{^1}$ The split of can't into ca and n't may seem unusual, but it reflects a convention designed to regularize patterns for parsing algorithms.

| Input String | Correct Token Sequence |
|--|--|
| She can't go. It's a state-of-the-art device. He lives in the U.S. | ['She', 'ca', 'n't', 'go', '.'] ['It', ''s', 'a', 'state-of-the-art', 'device', '.'] ['He', 'lives', 'in', 'the', 'U.S.', '.'] |

Figure 3.5: A table illustrating common tokenization challenges with example input strings and their correctly tokenized sequences.

Once text has been tokenized, a common next step is to reduce the different inflectional forms of a word to a common base. For instance, in a search application, we want a query for 'study computational linguistics' to match documents containing the words 'studies', 'studying', or 'studied'. The simplest and most computationally efficient way to achieve this is through **stemming**. Stemming is a heuristic-based process for crudely chopping off word endings—primarily suffixes—to obtain a common base form, known as a stem.

It is crucial to understand that stemming is not a formal linguistic analysis. It operates on word strings using a pre-defined set of rules and has no knowledge of context or part of speech. The goal is not to produce a linguistically correct root word, but rather to ensure that a group of related words maps to the same stem. For example, the words computation, computational, and compute might all be reduced to the stem comput. Notice that the stem itself is not always a valid English word, which is perfectly acceptable for the purposes of stemming. The process is a form of conflation—collapsing distinctions between different word forms.

The heuristic nature of stemming, while fast, inevitably leads to two types of errors:

- Over-stemming: This occurs when too much of a word is removed, causing words with different meanings to be conflated into the same stem. For example, a stemmer might incorrectly reduce both *universal* and *university* to the stem univers, treating them as the same concept. This is a false positive.
- Under-stemming: This is the opposite problem, where the stemmer fails to reduce words that are actually related to the same stem. A simple stemmer might fail to relate data and datum, or adhesion and adhesive, because its rules are not sophisticated enough to handle their specific morphological patterns. This is a false negative.

Despite these limitations, stemming is a powerful and widely used technique, especially in applications like information retrieval where speed is critical and the occasional error has a minimal impact on overall performance. The simplicity of the approach makes it an attractive first step in the normalization pipeline. Numerous algorithms for stemming exist, each with its own set of rules and trade-offs between aggression and accuracy. The most influential of these is the Porter Stemmer, which we will examine next.

One of the most influential and widely used stemming algorithms is the Porter Stemmer, developed by Martin Porter in 1980. It provides a simple yet effective heuristic for removing common morphological and inflectional endings from English words. The algorithm does not rely on a dictionary but instead applies a sequence of suffix-stripping rules in a series of steps.

The Porter Stemmer operates through five distinct phases, applied in a fixed order. Each phase consists of a set of rules, such as "if the word ends in -ational, replace it with -ate'. For a rule to fire, the word must not only match the suffix but the remaining stem must also satisfy a certain length condition. This condition, known as the 'measure," roughly counts the number of vowel-consonant sequences in the stem. This prevents the algorithm from being overly aggressive on short words; for instance, it stops relate from

| Step | Rule Applied | Word Form |
|---------------|---------------------------------------|---------------|
| Initial State | _ | normalization |
| Step 2 | $\texttt{-ational} \to \texttt{-ate}$ | normalize |
| Step 4 | -ize $ ightarrow \epsilon$ | normal |
| Final Output | _ | normal |

Figure 3.6: Step-by-step reduction of the word 'normalization' using the Porter Stemmer algorithm.

being incorrectly reduced to *rel*. The key principle is that if multiple rules within a step could apply, only one is chosen—typically the one with the longest matching suffix.

To make this concrete, let's consider how the algorithm would process the word normalization. The process, which is detailed step-by-step in Fig. 3.6, demonstrates the sequential nature of the rules. The word is transformed iteratively as it passes through the phases. In Step 2, the rule (-ational \rightarrow -ate) reduces normalization to normalize. In a subsequent step, the rule (-ize \rightarrow) further reduces it to normal. Finally, after passing through all five steps, the algorithm terminates, yielding the final stemmed form normal.

It is crucial to remember that the Porter Stemmer is a **heuristic**. Its output is not guaranteed to be a linguistically correct root (e.g., university becomes univers). However, its strength lies in its speed and its ability to consistently map variants like normalize, normalized, and normalization to the same token. This consistency is often exactly what is needed for downstream tasks like information retrieval, where the goal is to match query terms to documents regardless of their specific inflectional form.

While stemming provides a quick and often effective method for normalization, its heuristic nature can be a blunt instrument, sometimes producing non-words. A more sophisticated and linguistically principled approach is **lemmatization**. The goal of lemmatization is to reduce a word not to a 'stem,' but to its *lemma*—the canonical or dictionary headword form. For instance, the lemma for cars is car, the lemma for feet is foot, and the lemma for running is run. The lemma represents the abstract concept of the word, under which all its inflected variants (like plural nouns or conjugated verbs) are grouped.

Unlike stemming's crude suffix-stripping rules, a lemmatizer relies on a deeper understanding of a language's morphology. The process involves more than just algorithmic transformation; it is a form of analysis. To correctly identify the lemma, a system typically requires two components:

- 1. A **dictionary** (or a more complex lexical resource) containing valid words and their corresponding lemmas.
- 2. Morphological analysis rules to handle inflections. For regular forms like books → book, this might seem simple, but the real power of lemmatization lies in its ability to handle irregular forms that do not follow simple patterns.

Crucially, effective lemmatization often requires knowing the word's part-of-speech (POS) in the context of the sentence. A word's lemma can change completely depending on its grammatical role. Because of this dependency, lemmatization is often performed as a step immediately following POS tagging (see Chapter 5). Consider the word saw in two different contexts:

• In 'I saw a bird,' saw is a verb. A POS-aware lemmatizer would correctly identify its lemma as the verb see.

| Original Word | Stemmed Form (Porter) | Lemmatized Form |
|---------------|-----------------------|-----------------|
| studies | studi | study |
| studying | studi | study |
| study | studi | study |
| better | better | good |
| mice | mice | mouse |

Figure 3.7: A comparative table contrasting the outputs of stemming and lemmatization. The table shows how stemming provides a crude, heuristic reduction (e.g., 'studi') while lemmatization returns the correct dictionary form (e.g., 'study', 'good', 'mouse').

• In 'He used a saw,' saw is a noun. Its lemma is simply saw.

Without the part-of-speech context, a lemmatizer would be forced to guess or default to the most common form, potentially introducing errors. This ability to disambiguate based on context is a significant advantage over stemming, which would likely leave the word saw unchanged in both cases.

This reliance on linguistic knowledge allows lemmatization to handle a wide range of complex and irregular cases that are impossible for a stemmer to resolve. For example, it correctly maps irregular verbs like went to go and was to be. It can also handle non-standard plurals, such as geese to goose and corpora to corpus. Furthermore, it can consolidate comparative and superlative adjectives, reducing both better and best to the lemma good. This ability to group words based on their shared meaning, rather than just their orthographic similarity, is invaluable for many semantic analysis tasks. Many lemmatization algorithms, such as the one included in the popular lexical database WordNet², are built upon these large, hand-crafted resources that map inflected forms to their base lemmas.

The choice between stemming and lemmatization represents a fundamental trade-off between computational efficiency and linguistic accuracy. While both techniques aim to reduce words to a base form, their approaches and the quality of their outputs differ significantly. Stemming is a heuristic, rule-based process that works by crudely chopping off word endings. Lemmatization, in contrast, uses morphological analysis and a dictionary to return the proper base form, or *lemma*, of a word.

This distinction is most apparent when observing their outputs. A stemmer like the Porter algorithm will reduce *studies*, *studying*, and *study* to the common stem 'studi'—a character sequence that is not a valid English word. A lemmatizer, on the other hand, will correctly map all three forms to the dictionary headword 'study'. The gap in sophistication widens with irregular forms. Stemmers are typically unable to connect words like *better* to *good*, or *mice* to *mouse*, as there is no simple affix-stripping rule that governs these relationships. Lemmatization, with its reliance on a comprehensive lexicon, handles these cases correctly. The comparative examples in **Fig. 3.7** clearly illustrate these differences.

So, when should one be preferred over the other?

• Stemming is chosen when speed and performance are the primary concerns. For applications like search engines, the goal is often just to ensure a query for 'studying' also retrieves documents containing 'studies'. The linguistic inelegance of the stem 'studi' is irrelevant as long as it successfully collapses the variants together. Stemming can also process unknown words that are not in a dictionary.

²WordNet is a large lexical database of English. Nouns, verbs, adjectives and adverbs are grouped into sets of cognitive synonyms (synsets), each expressing a distinct concept. It is often used for its lemmatization capabilities.

• Lemmatization is essential for tasks requiring a deeper semantic understanding. In machine translation, question answering, or sentiment analysis, knowing that 'better' relates to 'good' is crucial for correct interpretation. The cost of this accuracy is a slower process that depends on a lexicon and often requires part-of-speech information to resolve ambiguity (e.g., to lemmatize 'saw' to 'see' or 'saw').

Ultimately, the decision rests on the specific needs of your NLP pipeline. If you require a fast and simple method for collapsing word forms and can tolerate some imprecision, stemming is a reasonable choice. If your application demands linguistic correctness, lemmatization is the superior, more principled approach.

In addition to morphological normalization, two other simple yet powerful techniques are often included in a standard pre-processing pipeline: case folding and stop word removal. Case folding is the process of converting all characters in the text to a single, uniform case—typically lowercase. This ensures that words like Book, book, and BOOK are all treated as a single token. The primary benefit is a reduction in the size of the vocabulary, which helps consolidate word counts and is advantageous for many statistical models. However, this is a lossy transformation. Casing can carry important information, such as distinguishing the proper noun US (United States) from the pronoun us, or the name Brown from the color brown. For tasks like Named Entity Recognition, preserving the original case is often crucial.

Stop word removal is the process of filtering out common words that are considered to have little semantic value for a given task. These *stop words* are typically the most frequent words in a language and include:

- Articles (a, an, the)
- Prepositions (in, on, of, for)
- Conjunctions (and, but, or)
- Pronouns (he, she, it)

Because these words appear in nearly all documents, they can dominate frequency-based analyses without contributing much to the overall topic or meaning. The removal process involves comparing tokens against a pre-compiled, language-specific list of stop words and discarding any matches. Like case folding, this technique must be applied with care. For sentiment analysis, removing the word not could completely invert a phrase's meaning. For machine translation, these function words are grammatically indispensable. Thus, the decision to use these normalization steps depends entirely on the goals of the downstream application.

To solidify our understanding, let's walk through the complete normalization of a single sentence from raw string to a final, processed list of tokens. This multi-stage pipeline, which combines the techniques we have just discussed, is a standard prerequisite for most language analysis tasks. The entire process is visually summarized in the flowchart shown in Fig. 3.8.

Our starting point is the simple declarative sentence:

```
'The quickest foxes are jumping.'
```

The first step is **tokenization**. The raw string is segmented into a sequence of individual tokens, carefully separating punctuation from the final word. This yields the list:

```
['The', 'quickest', 'foxes', 'are', 'jumping', '.']
```

Next, we apply **case folding** to standardize the text by converting all alphabetic tokens to lowercase. This step ensures that 'The' and 'the' are not treated as two distinct words, resulting in:

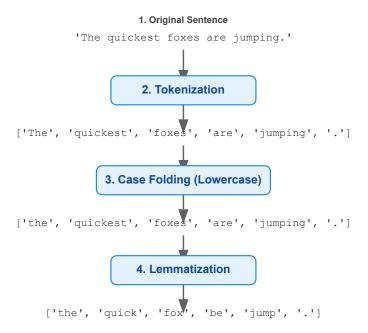


Figure 3.8: A diagram of the text normalization pipeline. The sentence 'The quickest foxes are jumping.' undergoes tokenization, case folding, and lemmatization. The output is shown at each stage, culminating in the final list of normalized tokens.

Finally, we perform **lemmatization**. Each token is reduced to its morphological root, or lemma. A morphologically-aware lemmatizer will convert the superlative adjective 'quickest' to 'quick', the plural noun 'foxes' to 'fox', and the inflected verb forms 'are' and 'jumping' to their base forms, 'be' and 'jump', respectively. The determiner and punctuation remain unchanged. The final, normalized output is a list of standardized tokens:

This sequence, stripped of its surface-level grammatical variations, is now in a canonical form suitable for subsequent quantitative analysis, such as frequency counting or feature extraction for a machine learning model.

Fortunately, implementing the text normalization pipeline from scratch is rarely necessary. A rich ecosystem of open-source software libraries provides robust, well-tested, and often multilingual implementations of these fundamental tasks. These tools allow practitioners to focus on higher-level modeling rather than re-implementing core components.

Some of the most prominent libraries include:

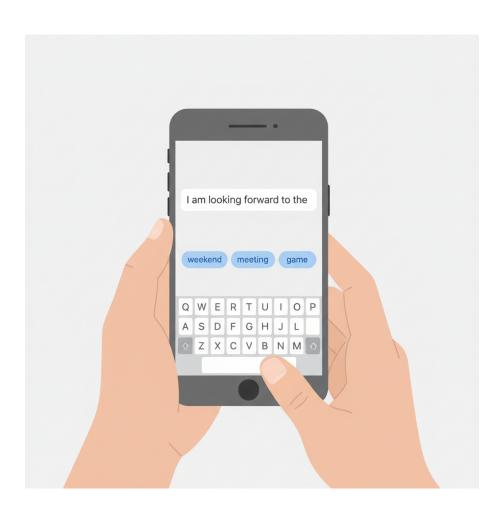
- NLTK (Natural Language Toolkit): A comprehensive library that is excellent for educational purposes. It provides a wide array of tokenizers, multiple stemmer implementations including the Porter and Lancaster algorithms, and a lemmatizer that leverages the WordNet lexical database.
- **spaCy:** Designed for production use, spaCy is known for its speed and efficiency. Its processing pipeline handles tokenization, lemmatization, and other tasks in a highly optimized manner, offering pre-trained models for dozens of languages.
- Scikit-learn: While primarily a machine learning library, its text feature extraction modules (e.g., CountVectorizer) seamlessly integrate normalization steps like tokenization, lowercasing, and stop word removal into the model training workflow.

• Stanza: Developed by the Stanford NLP Group, Stanza offers high-accuracy, pretrained models for many human languages, providing a full pipeline from tokenization and lemmatization to more advanced syntactic analysis.

This chapter has established the two foundational pillars of empirical language processing: the linguistic corpus and text normalization. These are not merely administrative preliminaries; they are the indispensable first steps that enable every subsequent task we will explore. A principled, well-constructed corpus acts as the ground truth—the source of authentic language data from which statistical patterns and linguistic structures are learned. However, this raw data is rarely suitable for direct computation.

The text normalization pipeline, with its crucial stages of tokenization, case folding, lemmatization, and stop word removal, transforms chaotic, unstructured text into a clean, consistent format. This standardization is what allows algorithms to treat 'run,' 'ran,' and 'running' as instances of the same underlying concept. Without this careful preparation, the statistical models and parsers discussed in the following chapters would be confounded by superficial variations, leading to poor performance. In short, the quality of corpus construction and normalization directly determines the potential success of any downstream NLP application, from part-of-speech tagging to machine translation.

Chapter 4 Language Modeling with N-grams



How can a machine tell the difference between a fluent, natural-sounding sentence and a jumble of words? Consider the following two English sentences:

- 1. The students opened their books.
- 2. Books their opened students the.

For a human speaker, the first sentence is immediately recognizable as valid, while the second is nonsensical. But what if the choice is more subtle?

- 1. The cat sat on the mat.
- 2. The cat sat on the sky.

Both are grammatically well-formed, yet we intuitively know that the first is far more plausible. The core task of **language modeling** is to quantify this intuition by assigning a probability to a sequence of words. A language model is a statistical tool that learns the patterns of a language from a large body of text and can then compute the likelihood of a given word sequence. Its fundamental goal is to calculate the probability P(W), where W is a sequence of words (w_1, w_2, \ldots, w_n) . A good language model will assign a higher probability to sentence (1) than to (2), and a higher probability to sentence (3) than to (4).

This ability to score text for its likelihood is one of the most essential components in computational linguistics, powering a vast array of applications:

- Predictive Text: When your smartphone keyboard suggests the next word, it is using a language model to predict the most probable word to follow what you have already typed.
- Speech Recognition: A speech recognizer might be unsure whether a person said 'recognize speech' or 'wreck a nice beach.' A language model can determine that the former phrase is far more probable in common usage, helping to disambiguate the audio signal.
- Machine Translation: When translating a sentence, a system may generate multiple possible outputs. A language model helps select the most fluent and natural-sounding translation.

It is crucial to understand that a language model computes probability, not grammaticality. The famous sentence 'Colorless green ideas sleep furiously,' coined by linguist Noam Chomsky, is perfectly grammatical. However, because the sequence of words is semantically bizarre and has likely never appeared in any text corpus, a language model would assign it an extremely low probability. Conversely, a common but ungrammatical utterance like 'Gonna go to the store' might be assigned a relatively high probability if it appears frequently in the training data. Language models are a reflection of what is likely based on past observation, not what is correct based on formal rules. In this chapter, we will learn how to build our first language models using a simple yet powerful technique: N-grams.

To formally assign a probability to a sequence of words $W = (w_1, w_2, \dots, w_n)$, we turn to a fundamental tool from probability theory: the *chain rule*. The chain rule allows us to decompose the joint probability of a sequence of events into a product of conditional probabilities.

Let's start with a simple two-word sequence, (w_1, w_2) . The joint probability $P(w_1, w_2)$ is defined as:

$$P(w_1, w_2) = P(w_1)P(w_2|w_1)$$

In plain language, the probability of seeing the word w_1 followed by w_2 is the probability of seeing w_1 at the start of a sequence, multiplied by the probability of seeing w_2 given that we have just seen w_1 .

We can extend this logic to a sequence of any length. For a three-word sequence, the rule expands as follows:

$$P(w_1, w_2, w_3) = P(w_1)P(w_2|w_1)P(w_3|w_1, w_2)$$

The general form for a sequence of n words is a product of these conditional probabilities. The probability of the entire sequence W is given by:

$$P(w_1, \dots, w_n) = P(w_1) \times P(w_2|w_1) \times P(w_3|w_1, w_2) \times \dots \times P(w_n|w_1, \dots, w_{n-1})$$

This can be expressed more compactly using product notation:

$$P(w_1, \dots, w_n) = \prod_{k=1}^{n} P(w_k | w_1, \dots, w_{k-1})$$

This equation is the exact, unsimplified probability of a word sequence. For example, the probability of the sentence 'the cat sat' would be calculated as $P(\text{the}) \times P(\text{cat}|\text{the}) \times P(\text{sat}|\text{the}, \text{cat})$. While theoretically sound, this formulation presents a major practical problem. To calculate the probability of the last word, $P(w_n|w_1,\ldots,w_{n-1})$, we would need to estimate the probability of that word occurring given its *entire* history. For any reasonably long sentence, this specific sequence of preceding words will be unique and will almost certainly never have appeared in our training data. This makes it impossible to reliably estimate these probabilities. To build a workable model, we will need to introduce a simplifying assumption.

Calculating the true probability of a word given its entire history is practically impossible. The conditional probability term from the chain rule, $P(w_i|w_1,\ldots,w_{i-1})$, requires us to estimate the likelihood of a word appearing after a unique, and often very long, sequence of preceding words. For any reasonably sized corpus, most long sequences will have never appeared before. For example, the probability of the word exam following the sequence The students who had been studying all night for the final computational linguistics... would be impossible to compute directly, as this specific history is unlikely to exist in our training data. This extreme data sparsity makes the full chain rule model intractable.

To make the problem manageable, we introduce a crucial simplifying assumption. The **Markov assumption**, named after the mathematician Andrey Markov, states that the probability of the next word depends only on a fixed window of the k most recent words, rather than the entire preceding sequence. We effectively assume that the distant past is irrelevant for predicting the immediate future. While this is not strictly true for human language—which contains long-distance dependencies—it provides a powerful and computationally feasible approximation.

Formally, we approximate the true conditional probability as follows:

$$P(w_i|w_1, w_2, \dots, w_{i-1}) \approx P(w_i|w_{i-1}, \dots, w_{i-1})$$

This assumption dramatically reduces the complexity of our model. Instead of needing to see every possible sentence prefix to make a prediction, we only need to have seen the most recent few words. **Fig. 4.1** provides a visual representation of this simplification,

Full Chain Rule

Markov Assumption (Trigram)

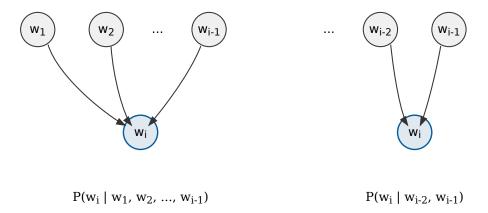


Figure 4.1: A diagram contrasting the full chain rule model with a simplified trigram model based on the Markov assumption. The left side shows the probability of a word being dependent on all preceding words, while the right side shows the probability being dependent on only the two preceding words, illustrating a significant reduction in complexity.

contrasting the complex web of dependencies in the full chain rule with the streamlined, local dependencies used in a model based on the Markov assumption. This simplification is the foundation of the N-gram model.

The size of the conditioning window, k, determines the order of the Markov model. This leads directly to the different types of N-gram models we will explore:

- Unigram Model (k=0): We assume a word's probability is independent of any prior words. This is the simplest, but weakest, assumption. $P(w_i|w_1,\ldots,w_{i-1}) \approx P(w_i)$
- Bigram Model (k=1): The probability of a word depends only on the single immediately preceding word. $P(w_i|w_1,\ldots,w_{i-1}) \approx P(w_i|w_{i-1})$
- Trigram Model (k=2): The probability of a word depends on the previous two words. $P(w_i|w_1,\ldots,w_{i-1})\approx P(w_i|w_{i-2},w_{i-1})$

By making this trade-off—sacrificing perfect linguistic representation for statistical feasibility—we can build models that are surprisingly effective at capturing local word patterns and predicting likely sequences.

The Markov assumption provides a practical way to approximate the chain rule of probability. The class of models built on this simplifying assumption is known as **N-gram models**, where an N-gram is simply a contiguous sequence of N words. These models form a foundational technique in statistical language modeling. The value of 'N' determines the size of the context window we consider when predicting the next word. While theoretically N can be any positive integer, in practice we most commonly encounter N=1, 2, or 3.

The simplest case is the **unigram model** (N=1). This model makes the radical assumption that each word is generated independently of all others. The probability of a word depends only on its own frequency in the training corpus, not on any preceding context. The approximation of the conditional probability thus becomes:

| Unigrams (N=1) | Bigrams (N=2) | Trigrams (N=3) |
|----------------|---------------|-----------------------------|
| The | <s> The</s> | <s> <s> The</s></s> |
| cat | The cat | $\langle s \rangle$ The cat |
| sat | cat sat | The cat sat |
| on | sat on | cat sat on |
| the | on the | sat on the |
| mat | the mat | on the mat |
| | mat. | the mat. |
| | . | mat . |

Figure 4.2: A breakdown of the sentence 'The cat sat on the mat.' into unigrams, bigrams, and trigrams. The special tokens <s> and </s> are used to model the start and end of the sentence for the bigram and trigram models.

$$P(w_i|w_1,\ldots,w_{i-1})\approx P(w_i)$$

Calculating the probability of an entire sentence $W = (w_1, \dots, w_k)$ under a unigram model is therefore just the product of the individual probabilities of its words:

$$P(W) = \prod_{i=1}^{k} P(w_i)$$

While this model discards all syntactic and semantic context, its simplicity makes it a useful baseline for comparison.

A more powerful and common model is the **bigram model** (N=2). Here, we apply the first-order Markov assumption: the probability of a word depends only on the single word that immediately precedes it.

$$P(w_i|w_1,...,w_{i-1}) \approx P(w_i|w_{i-1})$$

To properly handle the beginning and end of a sentence, we introduce special tokens. We prepend a start-of-sentence token, $\langle s \rangle$, and append an end-of-sentence token, $\langle s \rangle$. The probability of the first word is now conditioned on the start token, $P(w_1|< s \rangle)$, and the model must also predict the end token, $P(\langle s \rangle | w_k)$. The full probability of the sentence W is:

$$P(W) = P(w_1|< s>) \times \prod_{i=2}^{k} P(w_i|w_{i-1}) \times P(|w_k)$$

Extending this logic, a **trigram model** (N=3) uses a second-order Markov assumption, where the probability of a word is conditioned on the two preceding words.

$$P(w_i|w_1,...,w_{i-1}) \approx P(w_i|w_{i-2},w_{i-1})$$

For trigram models, we need two start-of-sentence tokens, <s> <s>, to provide the necessary two-word context for the first word of the sentence.

The general pattern holds for any N, but bigram and trigram models represent a sweet spot. They capture useful local context without becoming computationally unwieldy. The table in **Fig. 4.2** provides a concrete example, breaking down a simple sentence to illustrate how unigrams, bigrams, and trigrams are extracted. As N increases, the model can capture longer-distance dependencies and potentially become more accurate. However, this comes at a cost. The number of possible N-grams grows exponentially with

N, meaning we are far less likely to have seen any specific long N-gram in our training data. This issue, known as data sparsity, is a central challenge we will address next.

Once we have defined the N-gram model, the central task becomes learning its parameters from data. The parameters of an N-gram model are the conditional probabilities themselves, such as $P(w_i|w_{i-1})$ for a bigram model. The standard method for estimating these probabilities from a text corpus is called **Maximum Likelihood Estimation**, or **MLE**. The intuition behind MLE is simple yet powerful: the best estimate for a parameter is the one that maximizes the probability of the training corpus we observed. In the context of N-grams, this principle simplifies to a very direct calculation: the relative frequency. We assume that the probability of a future event is best represented by its frequency in the past, as recorded in our corpus.

The general formula for the Maximum Likelihood Estimate of an N-gram probability is a ratio of counts. The probability of a word w_i given its preceding history of n-1 words is the count of the full N-gram sequence divided by the count of the prefix (the history).

Let $C(w_1, ..., w_k)$ be the function that counts the number of times the sequence of words $(w_1, ..., w_k)$ appears in the corpus. The MLE probability is then:

$$P_{MLE}(w_i|w_{i-n+1},\ldots,w_{i-1}) = \frac{C(w_{i-n+1},\ldots,w_{i-1},w_i)}{C(w_{i-n+1},\ldots,w_{i-1})}$$

This single equation is the foundation for training simple N-gram models. Let's see how it applies to our specific cases:

- Trigram Model: For a trigram model (n = 3), we estimate the probability of a word w_i given the two preceding words, w_{i-2} and w_{i-1} . The MLE is the count of the trigram $C(w_{i-2}, w_{i-1}, w_i)$ divided by the count of the bigram prefix $C(w_{i-2}, w_{i-1})$. $P(w_i|w_{i-2}, w_{i-1}) = \frac{C(w_{i-2}, w_{i-1}, w_i)}{C(w_{i-2}, w_{i-1})}$
- **Bigram Model:** For a bigram model (n = 2), we estimate the probability of w_i given the single preceding word w_{i-1} . This is calculated as the count of the word pair $C(w_{i-1}, w_i)$ divided by the count of the single-word prefix $C(w_{i-1})$. $P(w_i|w_{i-1}) = \frac{C(w_{i-1}, w_i)}{C(w_{i-1})}$
- Unigram Model: The unigram model (n = 1) is the simplest case. It does not depend on any prior context, so the probability of a word w_i is just its relative frequency in the entire corpus. Here, the denominator is the total number of word tokens in the corpus, which we can denote as N. $P(w_i) = \frac{C(w_i)}{N}$

This process of 'training' an N-gram model, then, is not a complex optimization procedure but rather a straightforward task of data collection. It involves scanning a large corpus and counting all occurrences of relevant N-grams. For instance, to build a bigram model, we would need two main data structures: one to store the counts of all individual words (for the denominators) and another to store the counts of all adjacent word pairs (for the numerators). Once these counts are tabulated, calculating any specific conditional probability is a simple matter of division. This direct, count-based approach is both computationally efficient and easy to interpret, as the resulting probabilities are directly tied to the evidence seen in the training data. The next section will walk through a concrete example of this calculation. However, as we will soon see, this direct reliance on observed counts leads to a significant problem when we encounter N-grams that never appeared in our training text.

To make the process of training an N-gram model concrete, let's walk through a detailed example with a small corpus. A hands-on calculation will clarify how word counts from a text are transformed into a probabilistic model of language.

Imagine our entire corpus consists of just three sentences. We have already preprocessed the text by lowercasing all words and adding special start-of-sentence <s> and end-of-sentence </s> tokens.

- <s> the cat sat on the mat </s>
- <s> the dog sat on the rug </s>
- <s> the cat sat on the cat </s>

Our goal is to build a bigram model. The core task is to compute the Maximum Likelihood Estimate (MLE) for every conditional probability $P(w_i|w_{i-1})$ that can be formed from the vocabulary in this corpus. The formula, as a reminder, is:

$$P(w_i|w_{i-1}) = \frac{C(w_{i-1}, w_i)}{C(w_{i-1})}$$

This requires two sets of counts: unigram counts for the denominators and bigram counts for the numerators.

First, we calculate the unigram counts, which represent the number of times each word appears. These counts will serve as the denominators in our probability formula. For example, the word *the* appears as the first word in three bigrams—(the, cat), (the, dog), and (the, mat) in the first sentence alone. Counting across the entire corpus, we find:

- C(<s>) = 3
- C(the) = 6
- C(cat) = 3
- C(sat) = 3
- C(on) = 3
- C(dog) = 1
- C(mat) = 1
- C(rug) = 1

Next, we systematically count every sequence of two words to get our bigram counts. For instance, the sequence (<s>, the) occurs at the start of all three sentences, so its count is 3. The sequence (the, cat) occurs three times in total (once in the first sentence, twice in the third). The first part of **Fig. 4.3** shows the complete matrix of these bigram counts derived from our corpus.

With both sets of counts, we can now calculate the MLE probabilities. Let's compute a few examples:

- 1. What is the probability of the word cat following the word the? $P(\text{cat}|\text{the}) = \frac{C(\text{the, cat})}{C(\text{the})} = \frac{3}{6} = 0.5$
- 2. What is the probability of the word sat following the word cat? $P(\text{sat}|\text{cat}) = \frac{C(\text{cat, sat})}{C(\text{cat})} = \frac{2}{3} \approx 0.67$
- 3. What is the probability that a sentence starts with the word the? $P(\text{the}|<\text{s}>) = \frac{C(<\text{s}>, \text{the})}{C(<\text{s}>)} = \frac{3}{3} = 1.0$

Part 1: Bigram Counts $C(w_{i-1}, w_i)$

| | w_i | | | | | | | |
|----------------|------------------------|-----|------|----|----------------------|-----|-----|---|
| w_{i-1} | $\overline{	ext{the}}$ | cat | sat | on | dog | mat | rug | |
| <s></s> | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| \mathbf{the} | 0 | 3 | 0 | 0 | 1 | 1 | 1 | 0 |
| \mathbf{cat} | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 1 |
| \mathbf{sat} | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 0 |
| on | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| \mathbf{dog} | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| \mathbf{mat} | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| \mathbf{rug} | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |

Part 2: Bigram Probabilities $P(w_i|w_{i-1})$

| | w_i | | | | | | | |
|----------------------|-------|-----|----------------|-----|----------------|------|------|------|
| w_{i-1} | the | cat | \mathbf{sat} | on | \mathbf{dog} | mat | rug | |
| <s></s> | 1.0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| ${f the}$ | 0 | 0.5 | 0 | 0 | 0.17 | 0.17 | 0.17 | 0 |
| cat | 0 | 0 | 0.67 | 0 | 0 | 0 | 0 | 0.33 |
| \mathbf{sat} | 0 | 0 | 0 | 1.0 | 0 | 0 | 0 | 0 |
| on | 1.0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| \mathbf{dog} | 0 | 0 | 1.0 | 0 | 0 | 0 | 0 | 0 |
| \mathbf{mat} | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1.0 |
| rug | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1.0 |

Figure 4.3: Bigram counts (top) and corresponding Maximum Likelihood Estimation (MLE) probabilities (bottom) for a small corpus. Each probability in the bottom table is calculated by dividing the corresponding count in the top table by the sum of all counts in that row (the unigram count of the conditioning word).

Following this procedure, we can fill out a complete probability matrix, as shown in the second part of **Fig. 4.3**. This matrix *is* our trained bigram model.

We can now use this model to assign a probability to a full sentence. Let's take a sentence from our training corpus, <s> the dog sat on the rug </s>. Its probability is the product of the probabilities of its constituent bigrams:

 $P(\langle s \rangle \text{ the dog sat on the rug } \langle s \rangle) = P(\text{the}|\langle s \rangle) \times P(\text{dog}|\text{the}) \times P(\text{sat}|\text{dog}) \times P(\text{on}|\text{sat}) \times P(\text{the}|\text{on}) \times P(\text{rug}|\text{sat}) \times P(\text{the}|\text{on}) \times P(\text{the}|\text{on})$

Plugging in the values from our calculated model:

$$P = 1.0 \times \frac{1}{6} \times \frac{1}{1} \times \frac{3}{3} \times \frac{3}{3} \times \frac{1}{6} \times \frac{1}{1} = \frac{1}{36} \approx 0.0278$$

The model correctly assigns a non-zero probability to this seen sentence. But what happens when we encounter a new sentence that is perfectly valid English, but contains word sequences not present in our training corpus? Consider the sentence: <s> the cat sat on the floor </s>.

To calculate its probability, we would need to compute, among others, the conditional probability P(floor|the). Following our MLE procedure:

$$P(\text{floor}|\text{the}) = \frac{C(\text{the, floor})}{C(\text{the})}$$

Because the bigram (the, floor) never appeared in our tiny corpus, its count is 0. This means:

$$P(\text{floor}|\text{the}) = \frac{0}{6} = 0$$

The consequence is dire. Since the probability of the entire sentence is a product of individual probabilities, a single zero-probability event makes the probability of the whole sequence zero. The model predicts that this perfectly reasonable sentence is impossible. This is the fundamental problem of *data sparsity*, and it reveals the critical flaw of the simple Maximum Likelihood Estimation approach. In the next sections, we will explore techniques designed specifically to solve this problem.

The Maximum Likelihood Estimation (MLE) approach is simple and intuitive, but it harbors a significant flaw: it only works for N-grams that have appeared in the training corpus. What happens when our model encounters a new sentence containing a perfectly valid bigram, such as 'sentient robots', that simply wasn't present in our training text? The count for this bigram, C(sentient robots), will be zero. Consequently, the MLE calculation will yield a probability of zero. This isn't a minor edge case; it is a fundamental and pervasive challenge in statistical language modeling known as **data sparsity**.

Human language is vast and creative. We constantly generate new, sensible phrases and sentences. Any finite training corpus, even one containing billions of words, represents only a minuscule fraction of all possible valid word sequences. The number of potential N-grams explodes combinatorially as we increase our vocabulary size. For a modest vocabulary of 20,000 words, there are $20,000^2=400$ million possible bigrams and $20,000^3=8$ trillion possible trigrams. The vast majority of these will not appear in any given corpus, no matter how large. The bigram count matrix derived from our small sample corpus, shown in **Fig. 4.4**, provides a stark visual illustration of this problem. The matrix is overwhelmingly populated with zeros, each one representing a word combination that our model was never taught.

This leads directly to the **zero-probability problem**. The probability of a word sequence is calculated by multiplying the probabilities of its constituent N-grams:

| | the | chef | prepared | fish | tofu | sentient | robots |
|------------------------|-----|------|----------|------|------|----------|--------|
| $\overline{	ext{the}}$ | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| chef | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| prepared | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| fish | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| tofu | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| sentient | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| \mathbf{robots} | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

Figure 4.4: A sparse bigram count matrix derived from a small sample corpus. The matrix is overwhelmingly populated with zeros, visually representing the data sparsity problem. The model has observed bigrams from 'the chef prepared the fish', but has zero counts for plausible but unseen bigrams like ('prepared', 'tofu') or ('sentient', 'robots').

$$P(w_1, w_2, \dots, w_k) = \prod_{i=1}^k P(w_i | w_{i-1}, \dots, w_{i-N+1})$$

If the probability of just one of the N-grams in this sequence is zero, the probability of the entire sequence becomes zero. This is a catastrophic outcome for a language model. It means the model has concluded that a potentially grammatical and meaningful sentence is absolutely impossible. For instance, if our training data included 'the chef prepared the fish' but not 'the chef prepared the tofu', our bigram model would assign P(tofu|prepared) = 0. As a result, the sentence 'the chef prepared the tofu' would be deemed impossible.

In practical applications, this is debilitating. A speech recognition system might refuse to output a correct transcription because it contains an unseen bigram. A machine translation system could discard a perfectly good translation for the same reason. A model that assigns zero probability to legitimate events is not just inaccurate; it's brittle and fails to generalize beyond the specific examples it has seen.

Clearly, we cannot trust a model that is so easily broken by novelty. To build robust language models, we must find a way to assign some small, non-zero probability to unseen N-grams. This process involves taking a small amount of probability mass from the N-grams we *have* seen and redistributing it to those we *haven't*. This family of techniques is known as **smoothing**, and it is the essential next topic we must address.

The Maximum Likelihood Estimation approach, while intuitive, is brittle. It assigns a probability of zero to any n-gram that did not appear in the training corpus. As a result, a single unseen bigram in a test sentence would cause the probability of the entire sentence to become zero, rendering the model useless for evaluating novel text. To overcome this critical flaw, we must use a set of techniques known collectively as **smoothing**, or sometimes **discounting**.

The core idea behind all smoothing algorithms is to reallocate probability mass. We take a small amount of probability from the events we *have* seen and distribute it among the events we *have not* seen. Since the total probability must always sum to 1, we cannot simply invent new probability for unseen n-grams; we must 'steal' it from the more frequent, seen n-grams. This ensures that every possible n-gram is assigned a small, non-zero probability.

This redistribution is conceptually illustrated in **Fig. 4.5**. In a standard MLE model, the entire probability mass is divided exclusively among the n-grams observed in the corpus. In a smoothed model, a fraction of the probability assigned to each seen event is shaved off. This collected probability mass is then redistributed among the vast number of unseen events, lifting their probability from zero to a small positive value.

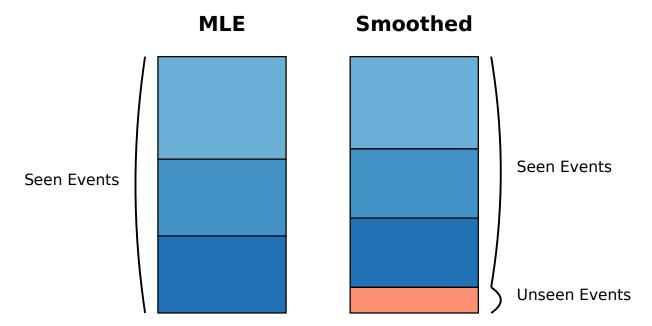


Figure 4.5: Conceptual diagram illustrating probability redistribution in smoothing. The Maximum Likelihood Estimation (MLE) bar shows probability mass divided exclusively among seen events. The Smoothed bar shows a fraction of that probability shaved off from the seen events and reallocated to a new block for unseen events, ensuring they have a non-zero probability.

In practice, this is achieved by adjusting the counts used in the probability calculation. Instead of using the raw frequency $C(w_{i-1}w_i)$, we use a discounted or adjusted count, c^* . For a seen n-gram, its adjusted count c^* will be slightly less than its real count C. For an unseen n-gram, its adjusted count will be a small positive value instead of zero. This adjustment acknowledges a fundamental truth: our training corpus is just one sample of the language, and the fact that an event is unseen does not make it impossible. The next sections introduce several specific algorithms for calculating these adjusted counts, from the simple to the more sophisticated.

The most direct and intuitive solution to the zero-probability problem is *Laplace smoothing*, also known as *add-one smoothing*. The core idea is simple: we pretend we have seen every possible event one more time than we actually have. This ensures that no event, seen or unseen, has a count of zero. By artificially inflating the counts, we can reallocate a small portion of the total probability mass to cover the events we didn't encounter in our training corpus.

Let's formalize this for our bigram model. The standard Maximum Likelihood Estimation (MLE) for a bigram probability is:

$$P_{MLE}(w_i|w_{i-1}) = \frac{C(w_{i-1}w_i)}{C(w_{i-1})}$$

To apply Laplace smoothing, we simply add one to every bigram count in the numerator. This means a bigram that occurred k times is now treated as if it occurred k+1 times, and an unseen bigram that occurred 0 times is now treated as having a count of 1. The adjusted numerator becomes $C(w_{i-1}w_i) + 1$.

However, we must also adjust the denominator to ensure the probabilities still sum to one. Since we added one to the count of *every* possible bigram that starts with the word w_{i-1} , we have effectively increased the total count of words following w_{i-1} by the total number of words in our vocabulary, |V|. For every word $v \in V$, we added one to the

| | (| Counts | Probabilities | | |
|----------|-----------------|----------------|---------------|----------------------|--|
| Bigram | $C(w_{i-1}w_i)$ | Adjusted Count | P_{MLE} | $P_{Laplace}$ | |
| I am | 8 | 9 | 8/10 = 0.800 | $9/19 \approx 0.474$ | |
| I am not | 0 | 1 | 0/10 = 0.000 | $1/19 \approx 0.053$ | |
| a fish | 1 | 2 | 1/5 = 0.200 | $2/14 \approx 0.143$ | |
| a tree | 0 | 1 | 0/5 = 0.000 | $1/14 \approx 0.071$ | |

Figure 4.6: A comparison of Maximum Likelihood Estimation (MLE) probabilities and Laplace (add-one) smoothed probabilities for a sample of bigrams. This table illustrates how Laplace smoothing redistributes probability mass from seen events (like 'I am') to unseen events (like 'I fish'), eliminating zero probabilities. The calculations assume the unigram counts are C(I)=10 and C(a)=5, and the total vocabulary size is |V|=19.

count of the bigram $C(w_{i-1}v)$. Therefore, the denominator, which is the total count for the prefix w_{i-1} , must be adjusted to $C(w_{i-1}) + |V|$.

This gives us the final formula for Laplace-smoothed bigram probability:

$$P_{Laplace}(w_i|w_{i-1}) = \frac{C(w_{i-1}w_i) + 1}{C(w_{i-1}) + |V|}$$

The practical effect of this redistribution is stark, as illustrated in **Fig. 4.6**. This table compares the original MLE probabilities from our earlier example with the new probabilities calculated after applying Laplace smoothing. Observe how bigrams that previously had a count and probability of zero, such as 'a fish', now have a small, non-zero probability. Conversely, the probabilities of frequent, observed bigrams like 'I am' are slightly reduced. This is the central trade-off of smoothing: we 'steal' a bit of probability mass from the events we have seen and give it to the events we have not.

While elegant in its simplicity, Laplace smoothing has a significant drawback: it is a blunt instrument. In a typical language model, the vocabulary |V| can be very large, often tens of thousands of words. By adding one, we reallocate a substantial amount of probability mass to a massive number of unseen N-grams. The majority of these unseen N-grams are not just unattested in our corpus—they are nonsensical combinations that are extremely unlikely to ever occur in real language (e.g., 'the of a'). Giving all of them equal probability mass is wasteful and systematically underestimates the probabilities of the events we actually did see.

For a vocabulary of 50,000 words, a bigram model would have $50,000^2 = 2.5$ billion possible bigrams. Laplace smoothing treats all unseen bigrams as equally plausible, which is an unrealistic assumption. This problem is exacerbated for trigram and higher-order N-gram models, where the number of potential N-grams ($|V|^N$) grows exponentially. Because it moves too much probability mass to the 'unseen' column, add-one smoothing often results in poor performance and is rarely used in modern systems. However, its simplicity makes it an essential pedagogical tool for understanding the fundamental goal of smoothing, paving the way for more sophisticated techniques. A common variant, add-k smoothing (or Lidstone smoothing), replaces the 1 with a small fractional value k (e.g., 0.01), which can mitigate the issue but doesn't fully solve it.

While Laplace smoothing offers a simple fix for the zero-probability problem, it is a blunt instrument. It reallocates probability mass equally among all unseen n-grams, regardless of how likely they might be, and it does so by taking too much probability away from high-frequency counts. A more sophisticated and empirically grounded approach is **Good-Turing smoothing**, named after its creators Alan Turing and I.J. Good.

The core intuition behind Good-Turing is both clever and powerful: to estimate the

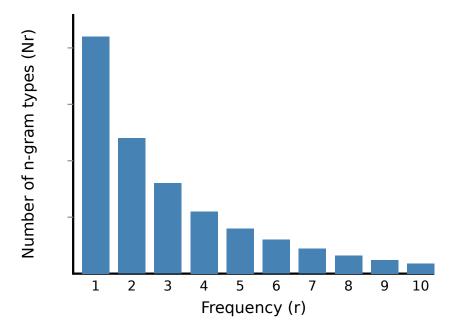


Figure 4.7: A bar chart showing the 'Frequency of Frequencies' for n-grams, illustrating the long-tail distribution used in Good-Turing smoothing.

probability of things you haven't seen, you should look at the frequency of things you've only seen once. Why? An n-gram that has occurred just once (a singleton) is on the verge of novelty. It was, until recently, an unseen event. The population of these singletons, therefore, provides our best evidence for the rate at which new, previously unseen events will appear in the future.

To formalize this, we first need to count the frequency of frequencies. Let N_c be the number of n-gram types that occur exactly c times in the corpus. For example, N_1 is the count of all singleton n-grams (those that appeared only once), N_2 is the count of all n-gram types that appeared twice, and so on. When we plot these values for a large corpus, we typically see a long-tailed distribution, as shown in Fig. 4.7. The number of singletons (N_1) is very large, followed by a smaller N_2 , an even smaller N_3 , and so on.

![Fig. 4.7](description: A bar chart titled 'Frequency of Frequencies'. The x-axis represents a frequency 'r' (n-grams that appeared 1 time, 2 times, etc.), and the y-axis represents the number of n-gram types (Nr) that had that frequency. This visualizes the long-tail distribution that Good-Turing smoothing leverages to estimate the probability of unseen events.)

Good-Turing uses this distribution to estimate the total probability mass that should be reserved for all zero-count n-grams combined. The estimate is remarkably simple:

$$P_{\text{unseen}} = \frac{N_1}{N}$$

Here, N is the total number of n-gram tokens in the corpus. We are effectively reassigning the probability mass of the singleton bucket to the zero-count bucket. The logic is that the N_1 singletons we observed are a good proxy for the number of new n-grams we would expect to find if we doubled the size of our corpus.

Of course, if we allocate this probability to unseen events, we must reduce the probability of the events we have seen to ensure all probabilities still sum to 1. Good-Turing does this by adjusting the count c for any n-gram that appeared c > 0 times. This adjusted count, denoted c^* , is used in place of the original MLE count. The formula is:

$$c^* = (c+1)\frac{N_{c+1}}{N_c}$$

The new probability for an n-gram that was seen c times is then $P_{GT} = \frac{c^*}{N}$. The formula for c^* may seem complex, but its effect is intuitive. Since for most corpora $N_{c+1} < N_c$ (there are fewer bigrams that appear 10 times than 9 times), the ratio $\frac{N_{c+1}}{N_c}$ will be less than 1. This means the adjusted count c^* will almost always be slightly less than the original count c ($c^* < c$), effectively discounting it. The amount of the discount is not arbitrary, as in Laplace smoothing, but is determined by the actual distribution of word counts in the corpus.

A significant practical problem arises with this formula: what if we encounter an N_c for which there is no observed N_{c+1} ? For example, we might have n-grams that occurred 9 times $(N_9 > 0)$ but none that occurred 10 times $(N_{10} = 0)$. This would make c^* for n-grams with a count of 9 equal to zero, which is clearly not right. To solve this and other issues related to sparse frequency-of-frequency counts, a variation called **Simple Good-Turing** is typically used in practice. It involves fitting a smooth curve to the plot of N_c vs. c, providing a reliable estimate for N_c even when the observed value is zero.

In summary, Good-Turing smoothing provides a far more principled method for handling unseen events than Laplace smoothing. By leveraging the distribution of the corpus itself—specifically using the count of singletons as an estimate for the unknown—it creates a more accurate and reliable language model.

After building and smoothing an N-gram model, a crucial question arises: how good is it? Is a trigram model with Good-Turing smoothing superior to a bigram model with Laplace smoothing for a given task? To answer such questions, we need a way to evaluate and compare language models. While the ultimate test is often *extrinsic*—how much the model improves a real-world application like machine translation—it is incredibly useful to have an *intrinsic* evaluation metric that measures the quality of the model independent of any specific application. The most common intrinsic metric for language models is **perplexity**.

Intuitively, perplexity measures how 'surprised' or 'confused' a model is by a sequence of words it has not seen during training. A good language model should assign a high probability to sentences that are well-formed and likely to occur, and a low probability to sentences that are nonsensical or ungrammatical. If a model assigns a high probability to a test sentence, it means it was not very surprised by it, which is the desired behavior. Perplexity is a formalization of this concept of surprise. A lower perplexity score indicates a better language model, one that is less surprised by the test data. You can think of perplexity as the average 'branching factor' of the language according to the model. If a model has a perplexity of 50, it suggests that at each word, the model is as confused as if it had to choose uniformly and independently from 50 possible words.

Formally, the perplexity of a language model on a test set $W = w_1, w_2, ..., w_N$ is the inverse probability of the test set, normalized by the number of words N. It is calculated as the N-th root of the inverse of the test set's probability:

$$PP(W) = \sqrt[N]{\frac{1}{P(w_1, w_2, ..., w_N)}} = P(w_1, w_2, ..., w_N)^{-\frac{1}{N}}$$

Using the chain rule of probability, which our N-gram models approximate, we can expand the joint probability $P(w_1, w_2, ..., w_N)$ into a product of conditional probabilities:

$$P(w_1, w_2, ..., w_N) = \prod_{i=1}^{N} P(w_i | w_1, ..., w_{i-1})$$

Substituting this into the perplexity formula gives us:

$$PP(W) = \left(\prod_{i=1}^{N} P(w_i|w_1, ..., w_{i-1})\right)^{-\frac{1}{N}}$$

The normalization by N is crucial because it ensures that the length of the test set does not unfairly penalize a model. Without it, a longer sentence would almost always have a lower probability than a shorter one, simply because it is a product of more fractions.

Perplexity has a close and important relationship with another information-theoretic concept: **cross-entropy**. The cross-entropy of a model M on a sequence W is defined as:

$$H(W) = -\frac{1}{N}\log_2 P(w_1, w_2, ..., w_N)$$

By rewriting the probability term as a product of conditional probabilities, we get:

$$H(W) = -\frac{1}{N} \sum_{i=1}^{N} \log_2 P(w_i | w_1, ..., w_{i-1})$$

Cross-entropy can be interpreted as the average number of bits needed to encode each word in the test set, given our language model. A better model (one that is less surprised) provides a more efficient compression scheme, resulting in a lower cross-entropy. Looking at the formulas for perplexity and cross-entropy, we can see their direct relationship:

$$PP(W) = 2^{H(W)}$$

This identity shows that minimizing perplexity is equivalent to minimizing cross-entropy. A lower cross-entropy—meaning the model's probability distribution is a better match for the empirical distribution of the test data—will always correspond to a lower perplexity.

When using perplexity, several practical considerations are vital:

- Unseen Test Data: The test set used for evaluation must be entirely separate from the training corpus. Evaluating a model on the data it was trained on will result in an artificially low perplexity, as the model has already memorized the sequences. The primary goal is to measure how well the model *generalizes* to new, unseen data.
- Tokenization Consistency: Perplexity scores are only comparable if the models being evaluated use the exact same vocabulary and tokenization scheme. For example, if one model treats 'don't' as a single token and another splits it into 'do' and 'n't', their vocabularies and the total word count N of the test set will differ, making their perplexity scores incomparable.
- Out-of-Vocabulary Words: The handling of unknown words (words in the test set but not in the training vocabulary) also significantly impacts the final score. A common practice is to replace all such words with a special <UNK> token in both training and testing, allowing the model to learn a probability for encountering an unknown word.

Perplexity is a powerful tool for quickly comparing different language models (e.g., bigram vs. trigram, or different smoothing techniques) on a standardized test set. A lower perplexity score generally correlates with better performance on downstream tasks, making it an essential first step in the model development and evaluation pipeline. However, it remains an intrinsic measure. It tells us how well the model predicts text, not necessarily



Figure 4.8: A clean illustration of a smartphone's virtual keyboard. A partial sentence is typed in the text field (e.g., 'I am heading to the...'), and the predictive text bar above the keyboard shows three probable next words (e.g., 'park', 'store', 'gym') that an N-gram model would suggest.

how useful it is for a specific application. An improvement in perplexity should always be confirmed with an extrinsic evaluation on the final task if possible.

While the theory of N-grams and perplexity can seem abstract, these models power a feature millions of people use every day: predictive text on smartphone keyboards. This application serves as a perfect, concrete illustration of language modeling in action. The core task of a predictive keyboard is to guess the next word you are likely to type, given the words you have already entered. This is precisely the problem that N-gram models are designed to solve.

Let's imagine you are typing a message. The system has access to the sequence of words you have written so far, $w_1, ..., w_{n-1}$. Its goal is to present a small set of highly probable next words, w_n . A system based on a bigram model would look only at the immediately preceding word, w_{n-1} , and calculate the probability $P(w_n|w_{n-1})$ for all words w_n in its vocabulary. The words with the highest conditional probability are the ones suggested to the user. The system seeks to find:

$$\underset{w \in V}{\operatorname{argmax}} P(w|w_{n-1})$$

Consider the scenario depicted in Fig. 4.8, where a user has typed 'I am heading to the'. A bigram model would use 'the' as the context. It would then query its stored probability distributions, which were learned from a massive training corpus. Based on billions of sentences from the web, the model might find that the probabilities of the words following 'the' are, for instance:

- P(store|the) = 0.008
- P(park|the) = 0.006
- P(gym|the) = 0.005
- P(airport|the) = 0.004

The system would then display 'store', 'park', and 'gym' as the three most likely completions. The model is simply reflecting the statistical patterns it observed during training; in its 'experience,' the phrase 'the store' was more frequent than 'the park.'

A more sophisticated system might use a trigram model, which considers the last two words as context. In our example, the context is 'to the'. The model now calculates $P(w_n|\text{to the})$. The suggested words might change. For example, while 'store' is a common word to follow 'the', the trigram 'to the store' might be even more common. Conversely, a word like 'airport' might have a much higher probability given the context 'to the' than it would given only 'the'. This demonstrates the power of using a larger N-gram window to capture more specific contextual information.

Several practical considerations are crucial for making such a system work effectively.

- 1. **Smoothing:** The importance of smoothing cannot be overstated. Without it, any word that never appeared after 'to the' in the training corpus would be assigned a probability of zero and could *never* be suggested. A user trying to type a perfectly valid but unseen phrase like 'to the conclave' would receive no useful predictions. Smoothing ensures that the model can handle the creativity and variability of human language.
- 2. Backoff and Interpolation: What happens if the trigram 'heading to the' has never been seen? A robust system doesn't just give up. It 'backs off' to a smaller context. If it has no data for the trigram, it uses the bigram model $(P(w_n|\text{the}))$. If that also fails, it could back off to the unigram model, suggesting words based on their general frequency.¹ This ensures the system can always provide some prediction.
- 3. **Efficiency:** The model must be incredibly fast. Suggestions need to appear instantaneously with each keystroke. To achieve this, the N-gram counts and probabilities are pre-calculated and stored in efficient data structures, like hash maps or tries, allowing for near-instant lookups.

This simple keyboard feature elegantly demonstrates the core concepts of this chapter: approximating the probability of a word sequence, the trade-off between context size and data sparsity, and the absolute necessity of smoothing to create a useful, generalizable model.

In this chapter, we have established the fundamental principles of statistical language modeling. We began with the core task of assigning a probability to a sequence of words, $P(w_1, w_2, ..., w_n)$. By formalizing this with the chain rule of probability and then simplifying it using the practical Markov assumption, we arrived at the N-gram model. You have learned how to train these models by calculating Maximum Likelihood Estimation (MLE) probabilities directly from corpus counts and seen how simple unigram, bigram, and trigram models capture local word patterns.

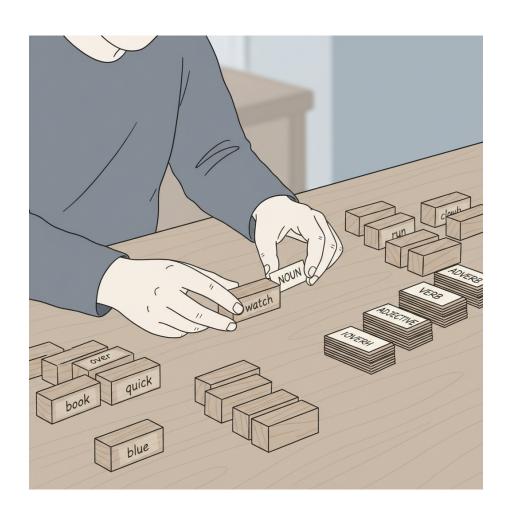
A central theme of the chapter was the critical problem of data sparsity. The fact that any finite training corpus will be missing countless valid word sequences forces naive MLE models to assign a probability of zero to unseen events. We addressed this by introducing smoothing, a class of techniques designed to reallocate probability mass from seen to unseen N-grams. We examined the intuitive but flawed Laplace (add-one) smoothing before moving to the more sophisticated Good-Turing smoothing, which estimates the probability of unseen events from the frequency of events seen only once. Finally, we introduced perplexity as a standard intrinsic metric for evaluating and comparing the predictive power of different language models.

¹A related technique is *interpolation*, where the final probability is a weighted average of the unigram, bigram, and trigram probabilities, combining evidence from all N-gram levels.

While foundational, N-gram models are limited by their fixed-context window. They cannot capture long-range dependencies or complex syntactic structures that span many words. This inability to model deeper linguistic phenomena is a key motivation for the more advanced models we will encounter later in this book. In the era of deep learning, neural networks have provided powerful new architectures for language modeling that can learn from much broader contexts, a topic we will return to in our final chapter.

Chapter 5

Part-of-Speech Tagging



Part-of-Speech (POS) tagging is the process of assigning a grammatical category to each word in a text. It is one of the most fundamental tasks in computational linguistics, acting as a foundational layer of analysis for many higher-level applications. The categories, or tags, correspond to the familiar parts of speech: nouns, verbs, adjectives, adverbs, prepositions, conjunctions, and more. For example, in the sentence 'The quick brown fox jumps over the lazy dog,' a POS tagger would analyze the sequence of words and label 'quick' as an adjective, 'jumps' as a verb, and 'dog' as a noun. This process adds a layer of shallow syntactic annotation to raw text, making it more amenable to further computational analysis.

Formally, the task is to map an input sentence, represented as a sequence of tokens $W = w_1, w_2, \ldots, w_n$, to a corresponding sequence of tags $T = t_1, t_2, \ldots, t_n$. Each tag t_i is drawn from a predefined inventory of tags, known as a tagset. This transforms a simple string of text into a sequence of annotated tokens, enriching the data with valuable structural information. This is a classic example of a sequence labeling problem, where the goal is to assign a categorical label to each member of a sequence, often taking into account the properties of both the individual item and its neighbors to make a decision.

The primary challenge that makes POS tagging a non-trivial computational problem is *lexical ambiguity*. A single word form can belong to multiple POS categories depending on its context. The task is not a simple dictionary lookup; rather, it requires an algorithm that can use the surrounding words to correctly disambiguate the word's role. Resolving this ambiguity is the central problem that modern tagging methods are designed to solve.

At first glance, POS tagging might seem as simple as looking up each word in a dictionary and assigning its category. This approach, however, quickly fails due to a fundamental property of language: *lexical ambiguity*. Many words can belong to more than one lexical category, and only the surrounding context can resolve the uncertainty. This is the central challenge that a POS tagger must overcome.

A classic example is the word *book*. It can function as a verb, meaning 'to make a reservation,' or as a noun, referring to a written work. The words in the immediate vicinity provide the crucial evidence for disambiguation. Consider the following sentences:

- 1. Book that flight for me.
- 2. I need to read that book.

In the first sentence, the imperative structure and the object 'flight' signal that book is an action—a verb. In the second, the preceding verb 'read' makes it clear that book is the object being acted upon—a noun. This contextual dependency, illustrated in Fig. 5.1, is intuitive for human speakers but must be computationally modeled for a machine.

This problem is far from isolated. Common English words like watch, fly, saw, and duck are also highly ambiguous. A successful POS tagger cannot simply consider a word in isolation; it must develop a systematic method for leveraging contextual patterns to make the correct choice. This is precisely why simple dictionary lookups are insufficient and more sophisticated probabilistic models are required.

Part-of-speech tagging is rarely an end in itself. Rather, it functions as a crucial preprocessing step for a vast range of more complex Natural Language Processing tasks. By resolving the kind of lexical ambiguity we saw with the word 'book', tagging provides a foundational layer of syntactic information that is essential for higher-level analysis. Many NLP systems are structured as a *pipeline*, where raw text is passed through a series of modules, each one enriching the data for the next. As shown in Fig. 5.2, POS tagging is a fundamental stage in this process, transforming a simple sequence of words into an annotated stream ready for deeper interpretation.

The most direct consumer of POS information is *syntactic parsing* (covered in Chapter 6). The goal of a parser is to uncover the grammatical structure of a sentence, identifying

Bookthat flight

Read that **book**



Figure 5.1: A diagram illustrating lexical ambiguity. The surrounding context determines the part-of-speech tag for the word 'book'. On the left, the imperative structure suggests 'Book' is a verb. On the right, the preceding verb 'Read' indicates that 'book' is a noun.

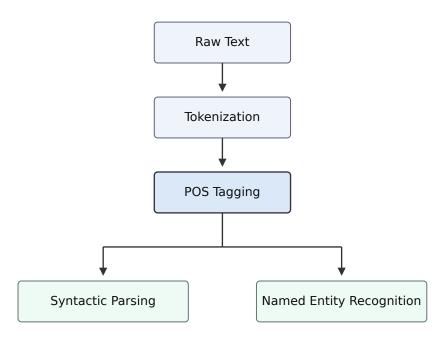


Figure 5.2: A standard Natural Language Processing (NLP) pipeline. Raw text is processed through sequential stages like Tokenization and Part-of-Speech (POS) Tagging. The resulting annotated data serves as a crucial input for more complex downstream tasks such as Syntactic Parsing and Named Entity Recognition.

| Tag | Description | Example |
|-----|-----------------------------------|---------|
| NN | Noun, singular or mass | book |
| NNS | Noun, plural | books |
| NNP | Proper noun, singular | London |
| JJ | Adjective | big |
| RB | Adverb | quickly |
| VB | Verb, base form | walk |
| VBD | Verb, past tense | walked |
| VBG | Verb, gerund/present participle | walking |
| VBZ | Verb, 3rd person singular present | walks |
| IN | Preposition/conjunction | in, of |
| DT | Determiner | the, a |
| PRP | Personal pronoun | he, it |

Figure 5.3: A selection of the most common tags from the Penn Treebank (PTB) tagset, with their descriptions and examples.

phrasal constituents like Noun Phrases (NP) or Verb Phrases (VP). This task is immensely more tractable when the parser already knows the grammatical category of each word. For instance, a parser can use a simple rule like $DET + ADJ + NOUN \rightarrow NP$ only if the words in the input have already been tagged as a determiner, an adjective, and a noun. Without this prior information, the parser would face a combinatorial explosion of possibilities.

POS tags are also a vital feature for *Information Extraction* (IE) systems. A key IE task is Named Entity Recognition (NER), which aims to find and classify entities like persons, organizations, and locations. A powerful heuristic for identifying named entities is to look for sequences of capitalized words that have been tagged as proper nouns (NNP). In a sentence like 'Dr. Evelyn Reed joined Globex Corporation,' a tagger would label 'Evelyn,' 'Reed,' 'Globex,' and 'Corporation' as proper nouns, giving the NER system a strong signal that 'Dr. Evelyn Reed' and 'Globex Corporation' are entities of interest.

Finally, accurate tagging benefits many other applications:

- Machine Translation: The grammatical role of a word is critical for correct translation. The German word Weg can be a noun ('the way') or an adverb ('away'), and a translator must know its POS tag to choose the right equivalent.
- Sentiment Analysis: Adjectives and adverbs often carry the strongest sentiment, so identifying them helps a system focus on the most opinionated parts of a text.

To perform Part-of-Speech tagging consistently, we need a standardized set of labels. Such a collection of tags is known as a *tagset* or an *annotation schema*. While dozens of tagsets have been developed over the years, each with varying levels of granularity, a few have become de facto standards within the computational linguistics community, ensuring that models and corpora from different research groups are comparable. Having a shared tagset is crucial for evaluating and replicating work in the field.

For English, the most influential and widely used schema is the **Penn Treebank** (**PTB**) tagset. It was developed in the early 1990s as part of the Penn Treebank Project at the University of Pennsylvania, which involved annotating a large corpus of Wall Street Journal articles with both part-of-speech and syntactic structure information. Its widespread adoption has made it the baseline for countless NLP systems and research papers, and familiarity with it is essential for any practitioner.

The PTB tagset uses short, mnemonic labels to represent grammatical categories. For instance, NN represents a singular common noun, JJ an adjective, and VB a verb in its base form. The tags also capture important morphological distinctions. A plural noun is

tagged NNS (the 'S' signifies plural), while a third-person singular present tense verb like 'walks' is tagged VBZ (the 'Z' is mnemonic for the '-s' ending). A selection of the most common tags from this schema, along with their descriptions and examples, is provided in **Fig. 5.3**. This table serves as a crucial reference for understanding the output of many standard POS taggers.

The core PTB tagset consists of 36 distinct part-of-speech tags, plus additional tags for punctuation and symbols. This number represents a deliberate trade-off. A smaller set of tags might be too coarse, grouping words with different grammatical behaviors (e.g., lumping all verb forms together). Conversely, a much larger tagset, while linguistically richer, would increase the difficulty of manual annotation and create data sparsity issues for statistical models, as many tags would appear too infrequently in the training corpus. The Penn Treebank tagset strikes a practical balance that has proven effective for a wide range of computational tasks.

Before the dominance of statistical methods, the first successful attempts at Part-of-Speech tagging were *rule-based systems*. These approaches sought to directly encode linguistic knowledge into a computer program through a combination of two primary components: a comprehensive dictionary and a set of handcrafted disambiguation rules. The core logic was intuitive and mirrored how a human linguist might approach the task.

The process typically occurred in two stages. First, the system would consult a large lexicon that mapped each known word to its possible set of POS tags. For a word like 'book,' the dictionary would return both NN (noun) and VB (verb). For an unambiguous word like 'and,' it would simply return CC (coordinating conjunction). This initial step would leave many words in a sentence with multiple potential tags.

The second, more critical stage involved applying a set of disambiguation rules to prune the incorrect tags. These rules were manually written by linguists to capture patterns of English grammar. For instance, a simple but effective rule might be:

If an ambiguous word follows a determiner (like the, a), it is most likely a noun.

This rule would correctly tag 'book' as a noun in the phrase 'the book.' Another rule could state:

If an ambiguous word follows a modal verb (like will, can, must), it is most likely a verb.

This would correctly tag 'book' as a verb in 'we will book a flight.' More sophisticated rules would handle morphology to guess the tags of unknown words. For example, a word ending in -ing is likely a gerund or present participle (VBG), a word ending in -ly is likely an adverb (RB), and a capitalized word not at the beginning of a sentence is likely a proper noun (NNP). One of the most famous early systems, the Brill tagger, even learned these kinds of contextual rules automatically from a corpus in a process called transformation-based learning.

While rule-based systems could achieve high accuracy, they suffered from significant drawbacks. Creating and maintaining the large set of rules required immense effort from expert linguists and was extremely time-consuming. These systems were also *brittle*; rules developed for one domain, like news articles, would often perform poorly on another, like poetry. Most importantly, the entire rule set had to be re-developed from scratch for each new language. These challenges of cost, scalability, and portability created a strong incentive for the field to develop more robust, data-driven methods, which led directly to the stochastic approaches we will explore next.

While rule-based taggers can achieve reasonable accuracy, they are brittle and laborious to create. A single new rule can have unintended consequences, and the entire system must be handcrafted by linguistic experts. To overcome these limitations, the field shifted towards a more robust and scalable paradigm: **stochastic tagging**. Instead of relying on hand-written rules, stochastic approaches use probability and statistics to determine the

most likely tag for a word in its given context.

The core idea is to leverage a large, pre-annotated corpus as a source of statistical evidence. By analyzing the frequencies of words and tags in this training data, a stochastic model can *learn* the likelihood of different tag assignments. This data-driven approach is a hallmark of modern computational linguistics, allowing models to automatically capture complex linguistic patterns without explicit instruction. The task is no longer to define what is grammatically *possible*, but rather to calculate what is statistically *probable*.

Stochastic taggers typically rely on two fundamental pieces of information derived from the corpus:

- 1. Lexical Probability (Emission Probability): The likelihood of a word appearing with a particular tag. For example, we can calculate the probability of seeing the word book given the tag is a noun, NN. This is represented as P(book|NN).
- 2. Contextual Probability (Transition Probability): The likelihood of a tag following another tag. For instance, we can calculate the probability that a noun (NN) follows a determiner (DT), such as in the phrase 'the book'. This is represented as P(NN|DT).

The central challenge of POS tagging is to combine these probabilities to find the best possible tag sequence for an entire sentence. We don't just want the most likely tag for each word in isolation; we want the sequence of tags $T = t_1, t_2, \ldots, t_n$ that is most probable for the sequence of words $W = w_1, w_2, \ldots, w_n$. Formally, the goal is to find the sequence \hat{T} that maximizes the probability P(T|W). To accomplish this, we need a formal model that can elegantly integrate lexical and contextual probabilities to search for this optimal path. The Hidden Markov Model, which we turn to next, provides exactly such a framework.

To formalize the stochastic approach to tagging, we introduce one of the most important statistical tools for modeling sequential data: the **Hidden Markov Model (HMM)**. An HMM is a probabilistic model designed to explain or generate a sequence of observable events that depend on a sequence of underlying, unobservable (or *hidden*) states. This structure makes it perfectly suited for tasks where we can see the output (like words in a sentence) but need to infer the process that generated it (like the sequence of part-of-speech tags).

Imagine you are in a windowless room and want to guess the weather outside. You cannot see the weather directly, making it a hidden state (e.g., Rainy or Sunny). However, each day a colleague enters the room, and you can observe whether or not they are carrying an umbrella. This is your observable event. You notice that your colleague is more likely to carry an umbrella on a rainy day than on a sunny day. Furthermore, you know that a rainy day is more likely to be followed by another rainy day than by a sunny day. An HMM provides a mathematical framework to combine these pieces of information to infer the most likely sequence of weather patterns (the hidden states) given the sequence of umbrella observations.

Formally, a Hidden Markov Model is defined by the following five components:

- A set of N hidden states, $Q = \{q_1, q_2, ..., q_N\}$ In our analogy, this would be the set of possible weather conditions, such as $Q = \{\text{Rainy, Sunny}\}.$
- A set of M possible observations, $O = \{o_1, o_2, \dots, o_M\}$ These are the observable outputs. In the analogy, this would be $O = \{\text{Umbrella}\}$.
- A state transition probability distribution, A This is a matrix $A = \{a_{ij}\}$ representing the probability of moving from state q_i to state q_j . The core 'Markov'

property is embedded here: the probability of the next state depends only on the current state.

$$a_{ij} = P(q_t = j \mid q_{t-1} = i)$$

• An observation emission probability distribution, B This is a matrix $B = \{b_j(k)\}$ representing the probability of seeing observation o_k when the model is in hidden state q_j . This connects the hidden states to the observable data.

$$b_i(k) = P(o_t = k \mid q_t = j)$$

• An initial state probability distribution, π This is a vector $\pi = \{\pi_i\}$ that specifies the probability of the model starting in state q_i at the beginning of a sequence.

$$\pi_i = P(q_1 = i)$$

For an HMM to be computationally tractable, it relies on two fundamental simplifying assumptions. First is the **Markov Assumption**: the probability of a particular state depends only on the previous state. This means $P(q_t \mid q_{t-1}, q_{t-2}, \ldots, q_1) = P(q_t \mid q_{t-1})$. While this ignores longer-range dependencies, it dramatically simplifies the model. Second is the **Output Independence Assumption**: the probability of an observation depends only on the state that produced it, not on any other states or previous observations. Thus, $P(o_t \mid q_t, q_{t-1}, \ldots, q_1, o_{t-1}, \ldots, o_1) = P(o_t \mid q_t)$. Together, these assumptions allow us to efficiently model complex sequential phenomena, making HMMs a cornerstone of classic computational linguistics.

To apply a Hidden Markov Model to Part-of-Speech tagging, we must frame the task in terms of the model's core components: hidden states, observations, and the probabilities that connect them. The mapping is direct and intuitive. Given a sentence, the words are what we can *see*, while the grammatical tags are the *hidden* structure we wish to uncover.

Formally, we define the components of the HMM as follows:

- Observations (O): The sequence of words in the sentence, w_1, w_2, \ldots, w_n . For the sentence 'Book that flight,' the observation sequence is ('Book', 'that', 'flight'). These are the evident, observable outputs of our model.
- Hidden States (Q): The sequence of Part-of-Speech tags, t_1, t_2, \ldots, t_n , corresponding to each word. The goal of the tagging process is to find the most probable sequence of these states. For our example, a plausible hidden state sequence is (Verb, Determiner, Noun). The set of all possible states is the complete tagset we are using (e.g., the 36 main tags of the Penn Treebank).

With the states and observations defined, the HMM is characterized by two key sets of probabilities, which are learned from an annotated corpus during a training phase.

First, we have the **transition probabilities**. These answer the question: given that we are in a certain state (tag), what is the probability of moving to another state (tag)? This is the probability of a tag t_i given the previous tag t_{i-1} , denoted as $P(t_i|t_{i-1})$. For instance, English grammar dictates that a determiner is very likely to be followed by a noun. Therefore, the transition probability P(Noun|Determiner) will be high. Conversely, the probability of a determiner following another determiner, P(Determiner|Determiner), will be extremely low. These probabilities effectively encode syntactic knowledge of the language.

Second, we have the **emission probabilities**. These answer the question: given that we are in a particular hidden state (tag), what is the probability of observing a specific word? This is the probability of a word w_i given its tag t_i , denoted as $P(w_i|t_i)$.

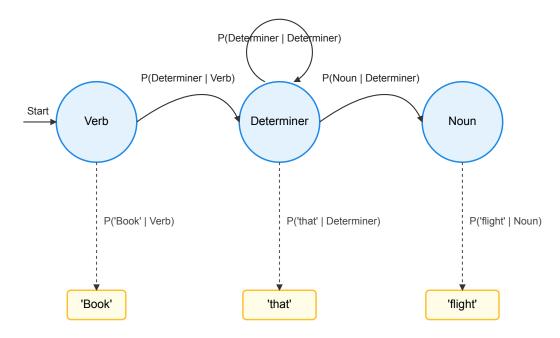


Figure 5.4: A diagram of a Hidden Markov Model (HMM) for Part-of-Speech (POS) tagging. The hidden states (Verb, Determiner, Noun) are shown in circles, and they generate observable words ('Book', 'that', 'flight'), shown in rectangles. The solid arrows between states represent transition probabilities (the likelihood of one tag following another), while the dashed arrows from states to words represent emission probabilities (the likelihood of a tag generating a specific word).

Emission probabilities capture lexical knowledge. For example, the probability of seeing the word 'book' given the tag Noun, P(`book'|Noun), will be significantly higher than the probability of seeing 'book' given the tag Adverb, P(`book'|Adverb). This is how the model handles lexical ambiguity; while 'book' can be both a Noun and a Verb, one is typically more probable than the other, and the model learns these likelihoods from the corpus.

The conceptual model, illustrated in Fig. 5.4, shows this generative process. The HMM starts in an initial state, transitions to a hidden tag (e.g., Verb), and then emits an observable word (e.g., 'Book'). It then transitions to the next tag (Determiner), emits the next word ('that'), and so on. Our task, however, is the inverse of generation. We are given the sequence of observations—the words—and our goal is to find the single most likely sequence of hidden states—the tags—that could have produced them. This decoding problem is precisely what the Viterbi algorithm, which we will discuss next, is designed to solve.

Once we have formally defined a Hidden Markov Model, the next logical step is to train it. Training an HMM simply means learning its parameters from data. For a POS tagger, this data is a large, pre-tagged corpus—a collection of text where each word has already been assigned its correct tag by a human annotator. We use this 'gold-standard' data to estimate the two sets of probabilities that define our model: transition probabilities and emission probabilities. This section focuses on the former.

Transition probabilities answer the question: given that we have just seen a particular tag, what is the likelihood of seeing another specific tag next? Formally, this is the conditional probability of tag t_i occurring given that the previous tag was t_{i-1} , which we denote as $P(t_i|t_{i-1})$. This captures the grammatical structure of a language. For instance, in English, it is highly probable that a determiner (like 'the') will be followed by a noun (like 'cat'), so we would expect the probability P(Noun|Determiner) to be high. Conversely,

| | Current Tag (t_i) | | | |
|--------------------------|---------------------|------|------|------|
| Previous Tag (t_{i-1}) | Determiner | Noun | Verb | |
| <s></s> | 0.60 | 0.30 | 0.10 | 0.00 |
| Determiner | 0.05 | 0.90 | 0.05 | 0.00 |
| Noun | 0.05 | 0.15 | 0.50 | 0.30 |
| Verb | 0.40 | 0.30 | 0.10 | 0.20 |

Figure 5.5: A simplified example of a transition probability matrix. Each cell shows the probability of the current tag (t_i) occurring, given the previous tag (t_{i-1}) . The special tag texttt<s> marks the start of a sentence and

texttt</s> marks the end. Note that each row's probabilities sum to 1.

the probability of a determiner following another determiner, P(Determiner|Determiner), would be extremely low.

We calculate these probabilities directly from the frequencies in our annotated corpus using a method called *Maximum Likelihood Estimation* (MLE). The formula is straightforward and intuitive:

$$P(t_i|t_{i-1}) = \frac{C(t_{i-1}, t_i)}{C(t_{i-1})}$$

Here, $C(t_{i-1}, t_i)$ is the count of how many times the tag sequence t_{i-1} followed by t_i appears in the corpus. The denominator, $C(t_{i-1})$, is the total count of the tag t_{i-1} appearing anywhere in the corpus. For example, to calculate P(Noun|Determiner), we would count every instance of a determiner followed by a noun and divide it by the total count of all determiners.

A special case arises at the beginning of a sentence, as the first word has no preceding tag. To handle this, we introduce a special *start-of-sentence* tag, often denoted as <s>. The transition probability from this special state, $P(t_i|<s>)$, represents the likelihood of a sentence beginning with tag t_i . This is calculated by counting how many sentences in the corpus start with tag t_i and dividing by the total number of sentences.

Once we compute these probabilities for every possible pair of tags in our tagset, we can organize them into a transition probability matrix. As shown in Fig. 5.5, the rows of this matrix represent the preceding tag (t_{i-1}) and the columns represent the current tag (t_i) . Each cell contains the corresponding conditional probability. A crucial property of this matrix is that the probabilities in each row must sum to 1, since some tag must always follow the preceding one.

$$\sum_{t_i \in T} P(t_i | t_{i-1}) = 1$$

This matrix forms one of the core components of our HMM, encoding the learned grammatical patterns of the language. The second component, the emission probabilities, will tell us which words are likely for each tag.

While transition probabilities model the relationship between tags, emission probabilities model the link between a tag and a word. They answer the question: given that we are in a particular hidden state (a POS tag), what is the likelihood of observing a specific word? Formally, this is the conditional probability of a word w given a tag t, denoted as P(w|t). These probabilities are also known as lexical likelihoods.

Like transition probabilities, emission probabilities are estimated directly from a large, tag-annotated corpus using Maximum Likelihood Estimation (MLE). The calculation is a straightforward counting exercise. To find the probability of a word w being generated by

| | Word | | | | |
|------------------------------|-----------------|--------|--------|--------------------|------------------|
| Tag | book | flight | walk | the | <unk></unk> |
| NN (Noun) | | | 0.0008 | | 0.0001 |
| VB (Verb) DT (Determiner) | 0.0004 0.0000 | | 0.000 | $0.0000 \\ 0.4500$ | 0.00005 0.0000 |

Figure 5.6: A sample emission probability matrix. Cells contain the conditional probability P(word | tag). For instance, the word 'book' is more likely to be a Noun (0.0015) than a Verb (0.0004). The special token <UNK> handles out-of-vocabulary words, with non-zero probabilities for open-class tags like Noun.

a tag t, we count how many times the word w appeared with the tag t in the corpus and divide it by the total count of the tag t.

The formula is as follows:

$$P(w_i|t_i) = \frac{\text{Count}(t_i, w_i)}{\text{Count}(t_i)}$$

Here:

- Count (t_i, w_i) is the number of times the word w_i is tagged with t_i in the corpus.
- Count (t_i) is the total number of times the tag t_i appears in the corpus.

For example, let's calculate the emission probability for the word 'book' given the tag Noun (NN). Suppose our training corpus contains 50,000 instances of the NN tag. We find that the word 'book' is tagged as a noun 75 times. The emission probability would be:

$$P(\text{'book'}|\text{NN}) = \frac{\text{Count}(\text{NN}, \text{'book'})}{\text{Count}(\text{NN})} = \frac{75}{50,000} = 0.0015$$

Now, suppose the tag Verb (VB) appears 30,000 times, and the word 'book' is tagged as a verb ('to book a flight') only 12 times. The probability is:

$$P(\text{'book'}|\text{VB}) = \frac{\text{Count}(\text{VB}, \text{'book'})}{\text{Count}(\text{VB})} = \frac{12}{30,000} = 0.0004$$

These probabilities quantify the lexical ambiguity we discussed earlier. The word 'book' can be a verb, but it is far more likely to be a noun. This information is crucial for the HMM when it needs to decide between competing tag sequences. We can pre-calculate these probabilities for every word-tag pair in our training vocabulary and store them in a large emission probability matrix. A small sample of such a matrix is shown in **Fig. 5.6**.

A significant challenge arises with words that appear in our test data but were never seen during training. These are often called out-of-vocabulary (OOV) words. Using the formula above, any unseen word would have a count of zero, resulting in an emission probability of $P(w_{OOV}|t)=0$ for all tags t. This is problematic because it makes any sentence containing that word impossible to parse. A standard technique to handle this is to replace all rare words (e.g., those appearing only once) in the training data with a special 'unknown word' token, <UNK>. We then calculate emission probabilities for <UNK> just like any other word. When an OOV word is encountered during tagging, we assign it the pre-computed probabilities for <UNK>, providing a non-zero chance for the model to proceed.

Once our Hidden Markov Model is trained—meaning we have calculated all the necessary transition and emission probabilities from our corpus—we face the central task of

decoding. Given a new sentence, which is a sequence of observations (words), how do we find the single most likely sequence of hidden states (POS tags) that produced it?

A simple but flawed approach would be to assign each word the tag that is most probable for it in isolation. This greedy method, however, fails to consider the context. A word might be slightly more likely to be a noun than a verb in general, but if it follows a modal verb like 'will,' the probability of it being a verb is much higher. The optimal tag sequence must account for the *entire* sequence, balancing both the emission probabilities (P(word|tag)) and the transition probabilities $(P(tag_i|tag_{i-1}))$. Exhaustively calculating the probability of every possible tag sequence is computationally intractable, as the number of sequences grows exponentially with the length of the sentence.

The solution to this problem is the **Viterbi algorithm**, an efficient dynamic programming approach that guarantees finding the optimal tag sequence without this exponential overhead. Dynamic programming works by breaking a complex problem into a series of simpler, overlapping subproblems and storing the solutions to these subproblems to avoid re-computation. In the context of POS tagging, the Viterbi algorithm moves from the first word to the last, and at each step, it calculates the most probable path to that point.

Imagine building a grid, or *trellis*, where the columns represent the words in the sentence and the rows represent all possible POS tags. The Viterbi algorithm fills this trellis one column at a time. For each cell in the grid—representing a specific tag for a specific word—the algorithm calculates the probability of the most likely tag sequence that ends in that very cell. It does this by considering all possible tags for the *previous* word, extending each of those paths to the current cell, and keeping only the path with the highest total probability.

Formally, we define $v_t(j)$ as the probability of the most probable tag sequence for the first t words of the sentence that ends with tag j. This value can be calculated recursively. To find the value for the cell at word t and tag j, we look at all the values calculated for the previous word, t-1. For each previous tag i, we calculate a candidate probability by multiplying three values: the stored Viterbi probability of the previous cell $(v_{t-1}(i))$, the transition probability from tag i to tag j, and the emission probability of the current word $(word_t)$ given tag j. The new Viterbi probability, $v_t(j)$, is the maximum of these candidate probabilities.

The core recursive step is defined as: $v_t(j) = \max_{i=1}^{N} [v_{t-1}(i) \cdot P(tag_j|tag_i) \cdot P(word_t|tag_j)]$ While calculating this maximum probability, the algorithm must also store a *back-pointer* for each cell, indicating which of the previous tags (i) led to this maximum value. This is crucial for reconstructing the final path.

The entire process unfolds in three stages:

- 1. **Initialization:** For the first word in the sentence, the Viterbi probability for each tag is simply the initial or start-state probability of that tag multiplied by the emission probability of the first word given that tag.
- 2. **Recursion:** For each subsequent word from t = 2 to the end of the sentence, and for each possible tag j, compute $v_t(j)$ using the formula above and store a backpointer to the previous tag that yielded the maximum value.
- 3. **Termination and Path Retrieval:** After processing the final word, the algorithm identifies the tag with the highest overall Viterbi probability. This tag is the final state in the most likely sequence. From there, it follows the chain of backpointers backward through the trellis to the beginning of the sentence, thereby recovering the complete, optimal sequence of POS tags.

To make the Viterbi algorithm concrete, let's walk through the process with a simple sentence: 'Book that flight'. Our goal is to find the most probable sequence of Part-of-

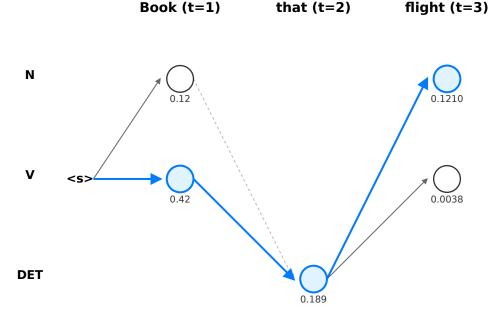


Figure 5.7: A Viterbi trellis diagram for a short sample sentence. The horizontal axis represents the sequence of words, and the vertical axis represents the possible states (tags). Each cell in the grid shows the calculated probability for that tag at that position, with backpointers indicating the most likely path, visually demonstrating the decoding process.

Speech tags. For this example, we will use a highly simplified tagset consisting of only three tags: Noun (N), Verb (V), and Determiner (DET). We will also assume we have already trained an HMM on a corpus, yielding the following transition and emission probabilities.

Transition Probabilities $P(tag_i|tag_{i-1})$

This table shows the probability of a tag given the preceding tag. <s> represents the start of a sentence.

| N | V | DET |
|-----|-------------------|-------------------------------|
| 0.3 | 0.7 | 0.0 |
| 0.2 | 0.7 | 0.1 |
| 0.3 | 0.2 | 0.5 |
| 0.8 | 0.2 | 0.0 |
| | 0.3 0.2 0.3 | 0.3 0.7 0.2 0.7 0.3 0.2 |

Emission Probabilities P(word|tag)

This table shows the probability of observing a word, given a specific tag.

| | Book | that | flight |
|----------------|------|------|--------|
| N | 0.4 | 0.0 | 0.8 |
| \mathbf{V} | 0.6 | 0.1 | 0.1 |
| \mathbf{DET} | 0.0 | 0.9 | 0.0 |

The Viterbi algorithm proceeds by constructing a probability matrix, often visualized as a trellis (see Fig. 5.7). Let $v_t(j)$ be the highest probability of a tag sequence ending at time step t with tag j. We will also store backpointers, $b_t(j)$, which tell us which previous tag led to this highest probability.

Initialization Step (t=1, word='Book')For the first word, the Viterbi probability is the product of the start transition probability and the emission probability.

• For tag N:
$$v_1(N) = P(\text{`Book'}|N) \times P(N|< s>) v_1(N) = 0.4 \times 0.3 = 0.12$$

• For tag V:
$$v_1(V) = P(\text{`Book'}|V) \times P(V|< s>) v_1(V) = 0.6 \times 0.7 = 0.42$$

At this point, the most likely tag for 'Book' in isolation is Verb (0.42 > 0.12), but we must consider the entire sequence. These values form the first column of our Viterbi trellis, as depicted in Fig. 5.7.

Recursive Step (t=2, word='that')Now we move to the second word, 'that'. For each possible tag of 'that', we calculate the probability by considering all possible paths from the previous step. The general formula is: $v_t(j) = P(w_t|tag_j) \times \max_{i=1..N} [v_{t-1}(i) \times P(tag_j|tag_i)]$

In our simplified example, 'that' can only be a Determiner (DET) as its emission probability for other tags is zero.

- For tag DET: We need to find the most probable path leading to DET from the previous tags (N or V).
 - Path from N: $v_1(N) \times P(DET|N) = 0.12 \times 0.1 = 0.012$
 - Path from V: $v_1(V) \times P(DET|V) = 0.42 \times 0.5 = 0.210$

The maximum of these two values is 0.210, which came from the previous tag being V. We now multiply this by the emission probability of 'that' given DET.

$$v_2(DET) = P(\text{`that'}|DET) \times 0.210 \ v_2(DET) = 0.9 \times 0.210 = 0.189$$

We also store the backpointer: $b_2(DET) = V$. This means the most likely path to 'that' being a DET comes from 'Book' being a V. In Fig. 5.7, this is shown as an arrow pointing from the V cell at t=1 to the DET cell at t=2.

Recursive Step (t=3, word='flight')We repeat the process for the final word, 'flight'. In our tables, 'flight' can be a Noun or a Verb.

- For tag N: The only possible preceding tag is DET (from step t=2).
 - Path from DET: $v_2(DET) \times P(N|DET) = 0.189 \times 0.8 = 0.1512$

The maximum is trivially 0.1512. Now we include the emission probability.

$$v_3(N) = P(\text{`flight'}|N) \times 0.1512 \ v_3(N) = 0.8 \times 0.1512 = 0.12096$$

The backpointer is $b_3(N) = DET$.

• For tag V:

- Path from DET:
$$v_2(DET) \times P(V|DET) = 0.189 \times 0.2 = 0.0378$$

The maximum is 0.0378.

$$v_3(V) = P(\text{`flight'}|V) \times 0.0378 \ v_3(V) = 0.1 \times 0.0378 = 0.00378$$

The backpointer is $b_3(V) = DET$.

The final column of our trellis now contains the probabilities $v_3(N) = 0.12096$ and $v_3(V) = 0.00378$.

Termination and Path RetrievalThe final step is to identify the highest probability in the last column and trace the backpointers to the beginning of the sentence.

1. **Identify the best final tag:** The highest probability in the final column is $v_3(N) = 0.12096$. Therefore, the most likely tag for the last word, 'flight', is **N**.

- 2. Trace back: We now follow the backpointers from this final state.
 - The backpointer for 'flight'/N is $b_3(N) = DET$. So, the tag for 'that' is **DET**.
 - The backpointer for 'that'/DET is $b_2(DET) = V$. So, the tag for 'Book' is V.

By reversing the sequence we traced, we get our final answer: V - DET - N.

The probability of this entire sequence is the final value we calculated: 0.12096. This step-by-step process of filling the trellis from left to right, storing only the maximum probability and a backpointer at each cell, is the essence of the Viterbi algorithm. It uses dynamic programming to efficiently find the optimal path without needing to score all possible tag sequences, which would be computationally intractable for any real sentence. The trellis in Fig. 5.7 provides a clear visual map of this computation, with the final traced path representing the single most likely grammatical structure for the sentence according to our model.

After training a Part-of-Speech tagger, we must measure its performance to understand its effectiveness. The standard procedure involves evaluating the tagger's output against a *gold-standard* test set—a corpus of text that has been manually annotated by expert linguists and was held out from the training data. These human-assigned tags are considered the correct answers, or ground truth.

The primary metric used for this evaluation is *accuracy*, which is the percentage of words that the tagger assigns the correct tag to when compared against the gold standard. The calculation is simple and intuitive:

$$Accuracy = \frac{\text{Number of correctly tagged words}}{\text{Total number of words}}$$

For instance, if a test set contains 10,000 words and the tagger correctly labels 9,750 of them, its accuracy is 97.5%. For a well-resourced language like English, state-of-the-art POS taggers routinely achieve accuracy scores in the 97–98% range.

While impressive, this single accuracy figure can be slightly misleading. Many words in a language have only one common POS tag, making them easy to classify. A more insightful evaluation, therefore, also measures performance on specific subsets of the data. One of the most important secondary metrics is the accuracy on *unknown words* (also called out-of-vocabulary words)—those that were not present in the training corpus. A tagger's ability to correctly label these words is a critical test of its generalization power, as it cannot rely on learned emission probabilities and must infer the correct tag purely from the sequence context provided by the transition probabilities.

This chapter has introduced Part-of-Speech tagging, a cornerstone task in computational linguistics that resolves lexical ambiguity by assigning grammatical categories to words. We traced the evolution from early rule-based systems to the dominant stochastic methods, which leverage annotated corpora to learn probabilistic models.

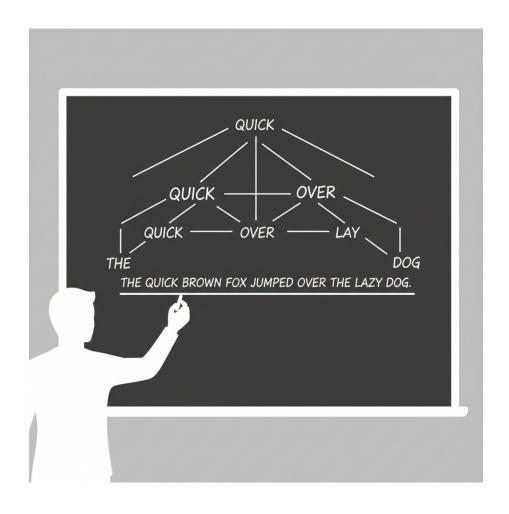
At the heart of our discussion was the *Hidden Markov Model (HMM)*, a powerful framework for sequence labeling. The HMM models POS tagging by combining two essential probabilities learned from data:

- Transition Probabilities: The likelihood of a tag sequence, $P(t_i|t_{i-1})$.
- Emission Probabilities: The likelihood of a word given a tag, $P(w_i|t_i)$.

We demonstrated that finding the optimal tag sequence is not a trivial task. The solution lies in the *Viterbi algorithm*, an elegant and efficient dynamic programming method for decoding the most probable sequence of hidden states. This combination of HMMs for modeling and Viterbi for decoding provides a robust template for solving a wide range of sequence-based problems you will encounter later in this book.

Chapter 6

Syntactic Parsing



Having learned to assign a part-of-speech tag to each word, we now move from a linear sequence of tags to a deeper, more structured understanding of a sentence. While knowing that a sentence contains a noun followed by a verb and another noun is useful, it doesn't fully capture the relationships between them. It doesn't tell us who is doing what to whom. This is the domain of **syntactic parsing**, the process of analyzing a string of words to reveal its underlying grammatical structure according to a formal grammar. The goal is not merely to verify if a sentence is grammatically correct, but to produce a formal representation—typically a tree—that makes its syntactic relationships explicit.

This structural analysis is crucial for resolving ambiguity and enabling true language understanding. Consider the classic sentence: 'I saw the man with the telescope.' This sentence has two distinct interpretations. Did I use a telescope to see the man? Or did I see a man who was holding a telescope? A simple sequence of POS tags is identical for both meanings. Syntactic parsing, however, forces us to commit to a specific structure that resolves this ambiguity. In the first interpretation, the phrase 'with the telescope' modifies the verb 'saw,' describing the instrument of seeing. In the second, it modifies 'the man,' describing an attribute of the man. The ability to computationally produce the distinct structural representations for these two meanings is a core goal of parsing.

Uncovering this structure is a foundational step for nearly all complex downstream NLP tasks. For a machine translation system to correctly translate a sentence, it must first understand the relationships between its components. For an information extraction system to identify that a company acquired another, it must parse the sentence to identify the subject, verb, and object. Similarly, a question-answering system must parse both the question and the potential answer text to find a structural match. By transforming a flat string of text into a hierarchical or relational structure, syntactic parsing provides the scaffolding upon which meaning can be built. In this chapter, we will explore the dominant computational formalisms and algorithms developed to achieve this.

Once we accept the goal of revealing a sentence's grammatical structure, a fundamental question arises: what should that structure look like? In computational linguistics, two dominant paradigms have emerged to answer this question, each offering a different philosophical perspective on what constitutes syntax. These are *constituency parsing* and *dependency parsing*.

Constituency parsing, also known as phrase-structure parsing, is built on the idea that sentences are composed of nested components called *constituents* or *phrases*. According to this view, words group together to form larger units, which in turn group together to form even larger ones, culminating in the complete sentence. The primary output of a constituency parser is a hierarchical tree that explicitly shows this grouping. For instance, the words 'the', 'old', and 'dog' combine to form a Noun Phrase (NP), and 'chased the cat' forms a Verb Phrase (VP). The NP and VP then combine to form the full sentence (S). This approach focuses on identifying the boundaries of these phrasal units and arranging them in a hierarchical structure.

Dependency parsing, in contrast, models syntax as a set of relationships between individual words. The structure is not a hierarchy of phrases but a graph of dependencies. In this framework, every word in the sentence, except for one (the root, typically the main verb), is a *dependent* of another word, which is its *head*. These directed links, or *arcs*, are labeled with the specific grammatical function they represent, such as nsubj for a nominal subject or dobj for a direct object. This approach forgoes phrasal nodes entirely and instead focuses on representing how words modify or depend on each other, revealing the functional architecture of the sentence.

The fundamental difference between these two formalisms is starkly illustrated in Fig. 6.1.

[DIAGRAM: Fig. 6.1]

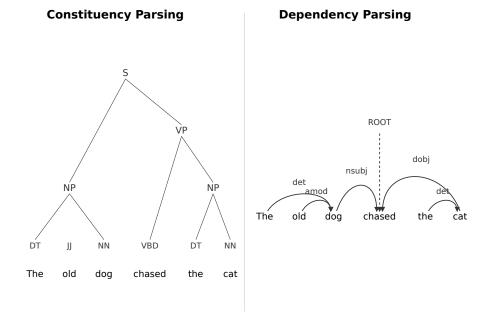


Figure 6.1: A side-by-side comparison of Constituency Parsing and Dependency Parsing for the sentence 'The old dog chased the cat'. The left side shows the hierarchical phrase structure, grouping words into constituents like Noun Phrases (NP) and Verb Phrases (VP). The right side shows the dependency graph, where directed, labeled arcs represent the grammatical relationships (e.g., subject, object) between individual words.

A side-by-side comparison diagram showing the same sentence, e.g., 'The old dog chased the cat', analyzed using both constituency and dependency parsing. The left side will show the hierarchical phrase structure tree (constituency), and the right side will show the dependency graph with directed, labeled arcs between words.

As the figure shows, the constituency tree for 'The old dog chased the cat' contains non-terminal nodes like NP and VP that are not words from the original sentence. It tells us that 'The old dog' is a single, cohesive unit. The dependency graph, on the other hand, consists only of the words themselves, connected by labeled arcs. It explicitly shows that 'dog' is the subject of 'chased' and 'cat' is its object, while 'old' directly modifies 'dog'. In essence, constituency parsing asks, 'What components make up this sentence?', while dependency parsing asks, 'Which words depend on which other words, and how?'. The remainder of this chapter will explore the formalisms and algorithms behind each of these powerful paradigms.

We first explore constituency parsing, an approach rooted in the idea that sentences are not flat strings of words but have an underlying hierarchical structure. The fundamental goal of this paradigm is to identify and organize groups of words that function as a single grammatical unit. These units are called **constituents** or **phrases**.

Consider the sentence, 'The little dog barked loudly'. Intuitively, we understand that 'the little dog' acts as a single block; it is the subject of the sentence. This group of words is a constituent. Similarly, 'barked loudly' functions as another block describing the action. This is also a constituent. Constituency parsing is the task of systematically breaking the sentence down into these nested blocks.

How do we formally identify a constituent? Linguists use several diagnostic tests to provide evidence for these groupings.

• Substitution: A sequence of words is likely a constituent if it can be replaced by a single word, often a pronoun or pro-form, without making the sentence ungrammatical. In our example, 'The little dog barked loudly', we can replace 'The little

dog' with 'It' to get 'It barked loudly'. This successful substitution strongly suggests 'The little dog' is a valid constituent.

• Movement: Constituents can often be moved to different positions within a sentence. For example, using a construction known as a cleft sentence, we can transform 'The cat sat on the mat' into 'It was on the mat that the cat sat.' The fact that the entire phrase 'on the mat' can be moved as a single, coherent unit is evidence of its status as a constituent.

These constituents are not just arbitrary clusters; they belong to specific grammatical categories. For example, 'the little dog' is a **Noun Phrase (NP)** because its most important word, or *head*, is a noun ('dog'). Likewise, 'barked loudly' is a **Verb Phrase (VP)**, and 'on the mat' is a **Prepositional Phrase (PP)**. The entire sentence itself is considered the top-level constituent, typically labeled **S**.

The central insight of this approach is that these phrases are nested within each other to form the complete sentence. An NP and a VP might combine to form a sentence (S). A VP might contain within it another NP (as in 'saw the big cat') or a PP (as in 'sat on the mat'). This creates a recursive, hierarchical structure. The output of a constituency parser is a **phrase-structure tree** (or simply a parse tree) that visually represents this hierarchy. In this tree, the words of the sentence appear as the leaves (or terminal nodes) at the bottom. The labels for the constituent phrases (NP, VP, PP, etc.) are the internal non-terminal nodes. The very top node, the root of the tree, is the S node representing the entire sentence.

This tree-based representation makes a sentence's internal grammar explicit. But how do we know which combinations of words and phrases are valid? How can a computer systematically build such a tree from a string of words? To answer these questions, we need a formal mechanism for defining the rules of sentence structure. This mechanism is the **Context-Free Grammar**, which we will explore next.

To computationally model the hierarchical structure of constituency, we need a formal mechanism that can define the set of all valid structures in a language. The most common tool for this is the **Context-Free Grammar (CFG)**, a formalism with deep roots in both linguistics, through the work of Noam Chomsky, and in computer science, where it is used to define the syntax of programming languages. A CFG provides a finite set of rules that can generate an infinite set of grammatical sentences, capturing the productive and recursive nature of human language.

Formally, a Context-Free Grammar G is a 4-tuple:

$$G = (N, \Sigma, R, S)$$

Let's break down each of these four components.

- 1. N: A finite set of **non-terminal symbols**. These are the variables of our grammar, representing syntactic categories or phrases. They are 'non-terminal' because they can be broken down further into other symbols. Think of them as the labels for the internal nodes in a parse tree—symbols like *Noun Phrase* (NP), *Verb Phrase* (VP), or *Prepositional Phrase* (PP). By convention, non-terminals are represented with uppercase letters or abbreviations.
- 2. Σ : A finite set of **terminal symbols**, sometimes called the *alphabet* or the *lexicon*. These are the actual words of the language that will appear in our sentences, such as 'cat', 'sees', or 'with'. They are 'terminal' because they cannot be rewritten into other symbols; they represent the end points of the generation process and correspond to the leaves of a parse tree. The set of non-terminals and terminals must be disjoint, meaning no symbol can be both: $N \cap \Sigma = \emptyset$.

| Component | Definition |
|---|--|
| Non-terminals (N) Terminals (Σ) Start Symbol | $ \begin{cases} S, NP, VP, PP, Det, N, V, P \\ \text{`the', 'cat', 'sat', 'on', 'mat'} \\ S \end{cases} $ |
| Production Rules (R) | $S \rightarrow NP \ VP$ $VP \rightarrow V \ PP$ $PP \rightarrow P \ NP$ $NP \rightarrow Det \ N$ $Det \rightarrow \text{'the'}$ $N \rightarrow \text{'cat'} \mid \text{'mat'}$ $V \rightarrow \text{'sat'}$ $P \rightarrow \text{'on'}$ |

Figure 6.2: The components of a simple Context-Free Grammar (CFG) based on the text.

- 3. R: A finite set of **production rules**. These rules are the heart of the grammar, specifying how non-terminals can be rewritten. Each rule is of the form $A \to \beta$, where A is a single non-terminal from the set N, and β is a sequence of one or more symbols from the combined set of non-terminals and terminals $(N \cup \Sigma)^*$. The arrow \to can be read as 'can be rewritten as,' 'consists of,' or 'can be expanded to.' For instance, the rule $S \to NP$ VP states that a sentence S can consist of a Noun Phrase NP followed by a Verb Phrase VP. It is common to use the pipe symbol | as a convenient shorthand to group rules with the same left-hand side. For example, NP -> Det N | N is equivalent to the two separate rules:
 - NP -> Det N
 - NP -> N

Production rules can be broadly categorized into two types. Syntactic rules specify how phrase types are composed of other phrases (e.g., VP -> V NP), while lexical rules (or lexicon rules) connect the syntactic categories to the terminal words (e.g., N -> 'cat').

4. S: A designated **start symbol**, which must be a member of the non-terminal set N. The start symbol is the entry point for the grammar, representing the highest-level category from which all valid sentences in the language defined by the grammar can be generated. For natural language grammars, this is almost always the symbol for sentence, typically denoted as S.

Let's make this concrete with the grammar we will use to parse our upcoming example sentence. The components of this simple grammar are laid out in **Fig. 6.2**.

According to **Fig. 6.2**, our set of non-terminals is $N = \{S, NP, VP, PP, Det, N, V, P\}$. Our terminals are the words themselves: $\Sigma = \{\text{'the'}, \text{'cat'}, \text{'sat'}, \text{'on'}, \text{'mat'}\}$. The start symbol is S. The production rules in the table define the valid sentence structures. For example:

- S -> NP VP is the fundamental rule for an English sentence.
- VP -> V PP defines a Verb Phrase that consists of a Verb followed by a Prepositional Phrase.

¹The asterisk notation $(N \cup \Sigma)^*$ means 'a sequence of zero or more symbols' from the given set. For most practical parsing grammars, rule right-hand sides have at least one symbol.

- NP -> Det N defines a Noun Phrase as a Determiner followed by a Noun.
- PP -> P NP defines a Prepositional Phrase as a Preposition followed by a Noun Phrase.

Notice the recursive potential here. A VP can contain a PP, which in turn contains an NP. This ability for rules to call upon other rules, including themselves (either directly or indirectly), is what gives CFGs the power to generate the complex, nested structures characteristic of human language. The process of using these rules to generate a string of terminals from the start symbol is called a **derivation**. A derivation starts with S and repeatedly applies rules from R, replacing a non-terminal with the right-hand side of one of its rules, until only terminal symbols remain. The set of all terminal strings that can be derived from the start symbol constitutes the *language* generated by the grammar.

To see how a CFG generates a sentence structure, let's parse the classic example: 'The cat sat on the mat'. We first need a grammar—a set of production rules—that can account for the structure of this specific sentence. Consider the following simple grammar, which we will call G1:

```
• S → NP VP (A sentence is a Noun Phrase followed by a Verb Phrase)
```

```
• NP → Det N (A Noun Phrase is a Determiner followed by a Noun)
```

```
• VP -> V PP (A Verb Phrase is a Verb followed by a Prepositional Phrase)
```

```
• PP → P NP (A Prepositional Phrase is a Preposition followed by a Noun Phrase)
```

```
• Det → 'the'
```

```
• N → 'cat' | 'mat'
```

V → 'sat'

• P → 'on'

The first four are *phrasal rules* that define how non-terminals can be rewritten into other non-terminals. The last four are *lexical rules* which ground the grammar by connecting non-terminal symbols (parts of speech) to the actual words, or *terminals*, in our language.

Using these rules, we can generate the sentence from the start symbol S through a sequence of substitutions called a **derivation**. In a *left-most derivation*, we always expand the left-most non-terminal at each step:

```
1. S
2. ⇒ NP VP (by applying the rule S → NP VP)
3. ⇒ Det N VP (by applying NP → Det N to the left-most non-terminal, NP)
4. ⇒ 'The' N VP
5. ⇒ 'The' 'cat' VP
6. ⇒ 'The' 'cat' V PP
7. ⇒ 'The' 'cat' 'sat' PP
8. ⇒ 'The' 'cat' 'sat' P NP
```

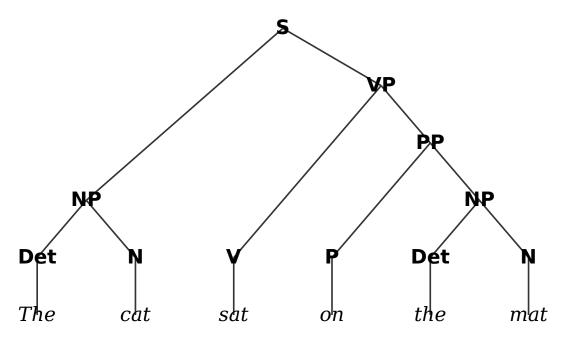


Figure 6.3: A phrase structure tree for the sentence 'The cat sat on the mat', derived from the grammar G1.

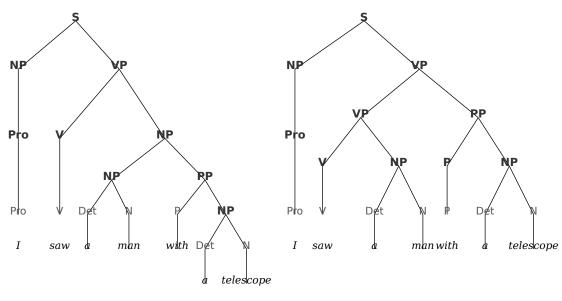
```
9. \Rightarrow 'The' 'cat' 'sat' 'on' NP 
 10. \Rightarrow 'The' 'cat' 'sat' 'on' Det N 
 11. \Rightarrow 'The' 'cat' 'sat' 'on' 'the' N 
 12. \Rightarrow 'The' 'cat' 'sat' 'on' 'the' 'mat'
```

This derivation confirms that the sentence is grammatical according to G1. More importantly, this process directly maps to a hierarchical structure. Each time a rule like $A \rightarrow B$ C is applied, we establish that B and C are the immediate constituents of A. The entire derivation can be visualized as a **phrase structure tree**, or parse tree, as shown in **Fig. 6.3**.

The tree in **Fig. 6.3** provides an explicit, visual representation of the sentence's syntax. The root of the tree is the start symbol S. The internal nodes are the non-terminals (NP, VP, PP), and the leaves are the terminals—the words of the sentence. The diagram clearly shows that 'The cat' forms a Noun Phrase (the subject), while 'sat on the mat' forms a complete Verb Phrase (the predicate). Within that VP, the tree reveals further nested structure: 'on the mat' is its own constituent, a Prepositional Phrase. This ability to capture the nested, hierarchical grouping of words is the fundamental power of constituency grammars.

While a Context-Free Grammar provides a formal and powerful mechanism for describing sentence structure, it also immediately brings a crucial challenge of natural language to the forefront: ambiguity. In many cases, a single sentence can be assigned more than one valid parse tree according to the rules of a given grammar. This phenomenon is known as structural ambiguity or syntactic ambiguity, and it is one of the most significant problems that a computational parser must address. Unlike lexical ambiguity, where a word has multiple meanings (e.g., 'bank' can be a financial institution or a river's edge), structural ambiguity arises from different ways of combining constituents.

A classic illustration of this problem is the sentence 'I saw a man with a telescope.' This sentence has at least two distinct interpretations, and each corresponds to a different syntactic structure.



Low Attachment (PP modifies NP)

High Attachment (PP modifies VP)

Figure 6.4: Two parse trees for the sentence 'I saw a man with a telescope' illustrating PP-attachment ambiguity. The left tree shows low attachment, where the prepositional phrase (PP) 'with a telescope' modifies the noun phrase (NP) 'a man'. The right tree shows high attachment, where the PP modifies the verb phrase (VP) 'saw a man'.

- Interpretation 1: The man I saw was holding a telescope.
- Interpretation 2: I used a telescope as an instrument to see the man.

These two meanings are not a result of any word having multiple definitions; they are a direct consequence of how the prepositional phrase (PP) 'with a telescope' attaches to the rest of the sentence. In the first interpretation, the PP modifies the noun phrase (NP) 'a man'. The phrase 'a man with a telescope' acts as a single, large NP, which is the direct object of the verb 'saw'. This is often called low attachment because the PP attaches to a constituent low in the parse tree. In the second interpretation, the PP modifies the verb phrase (VP) 'saw a man'. Here, the PP describes how the seeing action was performed. This is known as high attachment, as the PP attaches to the main VP, higher up in the tree. As illustrated in Fig. 6.4, these two interpretations correspond to two distinct and equally valid parse trees that can be generated by a plausible English grammar.

This problem, known as *PP-attachment ambiguity*, is just one of several types of structural ambiguity that are pervasive in language. Other common forms include:

- Coordination Ambiguity: This occurs with conjunctions like 'and'. For example, in the phrase 'young men and women', does 'young' modify only 'men', resulting in the structure [young men] and [women], or does it modify the entire coordinated phrase, giving young [men and women]? The correct interpretation depends entirely on context.
- Noun Compound Ambiguity: The structure of noun compounds can be ambiguous. The phrase 'student feedback system' could mean a system for providing feedback to students student [feedback system] or a system that processes feedback from students [student feedback] system.

The existence of structural ambiguity has profound implications. First, and most importantly, it means that syntax and semantics are deeply intertwined; a different parse

tree leads directly to a different meaning. Any system aiming for genuine language understanding must have a strategy for disambiguation—choosing the most plausible structure among the many possibilities. Second, the number of potential parse trees for a sentence can grow exponentially with its length, a problem known as combinatorial explosion. A single, moderately complex sentence can have hundreds or even thousands of valid parses. This makes parsing a significant computational challenge. Consequently, a parser must not only be able to find valid structures according to a grammar but must also do so efficiently and have a principled way to select the most likely parse, a topic that motivates the algorithms and probabilistic models we will explore next.

While a Context-Free Grammar provides the formal rules for generating a language, it doesn't specify how to find a valid parse for a given sentence. A naive approach might try to recursively expand every possible rule from the start symbol, a method that would quickly fail. Given the high degree of ambiguity in natural language, a brute-force search through all possible derivations would face a combinatorial explosion, rendering it computationally intractable for all but the most trivial sentences. To overcome this, we need a more systematic and efficient method. The solution lies in dynamic programming, an algorithmic paradigm that solves a complex problem by breaking it down into a collection of simpler subproblems, solving each subproblem just once, and storing their solutions to avoid redundant computation.

The Cocke-Younger-Kasami (CYK) algorithm is a classic example of dynamic programming applied to parsing. It systematically fills a table, often called a *chart*, with the set of non-terminals that can generate every possible substring of the input sentence. By building up solutions for progressively longer substrings from the solutions for shorter substrings it has already found, CYK can efficiently determine if a sentence can be generated by a given grammar.

However, the standard CYK algorithm comes with a crucial prerequisite: the Context-Free Grammar must be converted into **Chomsky Normal Form (CNF)**. A grammar is in CNF if all of its production rules take one of two simple forms:

- 1. $A \to BC$ (A non-terminal rewrites to exactly two non-terminals)
- 2. $A \rightarrow a$ (A non-terminal rewrites to a single terminal)

This strict binary-branching structure is the key that makes the CYK algorithm's dynamic programming approach work. The first rule type ensures that any constituent spanning multiple words is always formed by combining exactly two smaller, adjacent constituents. The second rule type, known as a lexical rule, grounds the parse by connecting the non-terminal symbols of the grammar to the actual words (terminals) in the sentence. While this may seem like a significant constraint, it has been proven that any context-free grammar can be converted into an equivalent one in CNF, meaning the algorithm loses no generative power.

The CYK algorithm operates on an $n \times n$ upper-triangular chart, let's call it P, where n is the number of words in the sentence. A cell P[i,j] in this chart (where $1 \le i \le j \le n$) will store the set of all non-terminals that can derive the substring of words from position i to position j. The algorithm proceeds bottom-up in two main stages:

• Initialization (Spans of length 1): First, we fill the diagonal of the chart. For each word w_i in the sentence (from i = 1 to n), we look at all rules of the form $A \to w_i$ in our CNF grammar. For every such rule found, we add the non-terminal A to the chart cell P[i,i]. After this step, the diagonal contains all possible phrasal categories for each individual word.

• Main Loop (Spans of length > 1): Next, the algorithm iteratively fills the rest of the chart, moving from shorter spans to longer ones. It iterates over the span length, 1en, from 2 up to n. For each length, it iterates over all possible starting positions, i, for a span of that length. The ending position j is simply i + 1en - 1. To fill the cell P[i,j], the algorithm must consider every possible way to split the substring $w_i...w_j$ into two smaller, contiguous substrings. Let's call the split point k, where k ranges from i to j-1. For each split, we have two sub-spans: $w_i...w_k$ and $w_{k+1}...w_j$. The algorithm then consults the chart cells that have already been filled for these smaller spans, namely P[i,k] and P[k+1,j]. If it finds any rule $A \to BC$ in the grammar such that non-terminal $B \in P[i,k]$ and non-terminal $C \in P[k+1,j]$, it adds non-terminal A to the cell P[i,j].

This process continues until the entire upper-triangular chart is filled. The final step is to check the top-rightmost cell, P[1,n], which corresponds to the span of the entire sentence. If the grammar's start symbol (e.g., S) is present in P[1,n], the sentence is recognized as being grammatically valid according to the grammar. If S is not in this cell, the sentence cannot be generated by the grammar and the parse fails.

The time complexity of the CYK algorithm is $O(n^3 \cdot |G|)$, where n is the length of the sentence and |G| is the size of the grammar. The n^3 factor comes from the three nested loops: one for the span length, one for the start position, and one for the split point. This cubic complexity is a vast improvement over the exponential complexity of naive search methods, making parsing feasible for reasonably long sentences. While the basic algorithm is a recognizer—it only tells you if a sentence is valid—it can be easily extended to become a full parser. To do this, whenever a non-terminal is added to a cell, we store back-pointers indicating which rule and which sub-cells were used to derive it. After the chart is filled, we can trace these pointers back from the start symbol in the top cell to reconstruct all possible valid parse trees.

To make the CYK algorithm concrete, let's walk through the process of parsing the sentence 'she eats fish'. We will use a simple Context-Free Grammar already converted to Chomsky Normal Form. Assume our grammar contains the following rules:

• Lexical Rules:

- $-PRP \rightarrow \text{`she'}$
- $-V \rightarrow \text{`eats'}$
- $-N \rightarrow \text{`fish'}$

• Structural Rules:

- $-S \rightarrow NP VP$
- $-NP \rightarrow PRP$
- $-NP \rightarrow N$
- $-VP \rightarrow VNP$

Our goal is to populate a parsing chart, often called a CYK chart or well-formed substring table, to determine if the start symbol S can generate the entire sequence of words. This chart, shown in **Fig. 6.5**, is an upper-triangular matrix where each cell P_{\downarrow} [i, j] will store the set of non-terminals that can generate the substring of words from position i to j. Let our words be indexed as w_1, w_2, w_3 .

Step 1: Initialization (Spans of Length 1)

We begin by filling the diagonal of the chart, which corresponds to spans of a single word. We look for all lexical rules $A \to w_k$ in our grammar.

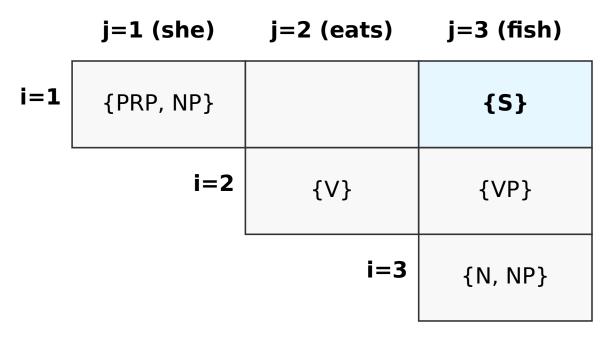


Figure 6.5: A completed CYK parsing chart for the sentence 'she eats fish'. The chart is an upper-triangular matrix where cell (i, j) contains the set of non-terminals that can generate the substring from word i to word j. The presence of the start symbol 'S' in the top-right cell (1, 3) confirms the sentence is grammatically valid according to the provided grammar.

- Cell P[1, 1] for 'she' (w_1): The rule $PRP \rightarrow$ 'she' applies. We add PRP to the cell. We must also consider any unit productions that can be reached from PRP. The rule $NP \rightarrow PRP$ allows us to add NP as well. So, P[1, 1] = {PRP, NP}.
- Cell P[2, 2] for 'eats' (w_2) : The rule $V \to$ 'eats' applies. There are no further unit productions from V. So, P[2, 2] = $\{V\}$.
- Cell P[3, 3] for 'fish' (w_3) : The rule $N \to$ 'fish' applies. From N, the rule $NP \to N$ allows us to add NP. So, P[3, 3] = {N, NP}.

After this step, the diagonal of our chart is filled, representing the possible grammatical categories for each individual word.

Step 2: Filling the Chart (Spans of Length 2)

Now we move to spans of length two. For each cell P[i, j] where the span length j - i + 1 is two, we look for all possible split points k. We check if any rule $A \to B$ C exists where $B \in P[i,k]$ and $C \in P[k+1,j]$.

- Cell P[1, 2] for 'she eats' (w_1w_2) : The only possible split is at k=1. We need to find a rule $A \to B$ C where $B \in P[1,1]$ and $C \in P[2,2]$.
 - We check combinations from $P[1,1] = \{PRP, NP\}$ and $P[2,2] = \{V\}$.
 - Is there a rule $A \to PRP \ V$? No.
 - Is there a rule $A \to NP V$? No.
 - Since no rules match, cell P[1, 2] remains empty.
- Cell P[2, 3] for 'eats fish' (w_2w_3) : The only split is at k=2. We need a rule $A \to B$ C where $B \in P[2,2] = \{V\}$ and $C \in P[3,3] = \{N,NP\}$.

- Is there a rule $A \to V N$? No.
- Is there a rule $A \to V NP$? Yes, $VP \to V NP$.
- Therefore, we add VP to this cell. $P[2, 3] = {VP}$.

Step 3: Filling the Final Cell (Spans of Length 3)

Finally, we compute the entry for the top-right cell, P[1, 3], which covers the entire sentence. The span length is three, so there are two possible split points: k = 1 and k = 2.

- Split 1 (k = 1): Check for rules $A \to B$ C where $B \in P[1, 1]$ and $C \in P[2, 3]$.
 - $-P[1,1] = \{PRP, NP\} \text{ and } P[2,3] = \{VP\}.$
 - We check for rules like $A \to PRP\ VP$ (no) and $A \to NP\ VP$ (yes, $S \to NP\ VP$).
 - This split successfully finds the start symbol S. We add S to P[1, 3].
- Split 2 (k=2): Check for rules $A \to B$ C where $B \in P[1,2]$ and $C \in P[3,3]$.
 - Since P[1, 2] is empty, this split cannot produce any non-terminals.

The final content of the top cell is P[1, 3] = {S}. The completed chart is shown in Fig. 6.5. Because the start symbol S is present in the cell corresponding to the entire sentence (P[1, 3]), we conclude that 'she eats fish' is a valid sentence according to our grammar. If S were not found in this cell, the sentence would be rejected as ungrammatical. To reconstruct the actual parse tree, a practical implementation of CYK also stores backpointers in each cell to remember which rule and which sub-constituents were used to add each non-terminal.

Having established how constituency grammars model sentence structure through hierarchical phrases, we now shift our focus to an alternative and equally influential paradigm: dependency parsing. Rather than focusing on nested constituents, dependency parsing represents grammatical structure by identifying the relationships between individual words. The core idea is that for any sentence, the syntactic structure consists of a set of directed, binary, and asymmetrical links between words.

In this framework, every relationship connects a head (also called a governor) to a dependent (or modifier). The head is the word that is grammatically dominant, while the dependent is the word that modifies or is an argument of the head. For example, in the phrase 'big dogs', the noun dogs is the head, and the adjective big is its dependent, as big modifies dogs. These directed links, called dependencies, are typically labeled to specify the exact nature of their grammatical function. These labels come from a standardized set, such as the Universal Dependencies (UD) tagset, which includes common functions like nsubj (nominal subject), obj (direct object), amod (adjectival modifier), and det (determiner).

The complete dependency structure for a sentence is represented as a dependency tree (or, more formally, a directed acyclic graph). To ensure that every word is connected and the structure forms a single tree, a special pseudo-word, often labeled ROOT, is introduced. The ROOT node acts as the head of the main predicate of the sentence, which is usually the main verb. A well-formed dependency tree for a sentence with n words must satisfy three key constraints:

- 1. The ROOT node has no incoming arcs (it has no head).
- 2. Every other word has exactly one incoming arc (it has exactly one head).

3. The graph is acyclic.

Together, these constraints ensure that the resulting structure is a tree where every word in the sentence is ultimately a dependent of the ROOT node. This approach provides a lexicalized representation of grammar, as the structure is built directly between the words themselves, without the need for intermediate phrasal nodes like NP or VP.

Let's return to our familiar example, 'The cat sat on the mat,' to see how a dependency parser would analyze it. The resulting dependency tree would consist of the following set of labeled arcs, represented here as (head, label, dependent) triples:

- (ROOT, root, sat): sat is the root of the sentence.
- (sat, nsubj, cat): cat is the nominal subject of sat.
- (cat, det, The): The is the determiner for cat.
- (sat, prep, on): on is a prepositional modifier of sat.
- (on, pobj, mat): mat is the object of the preposition on.
- (mat, det, the): the is the determiner for mat.

Notice how this structure differs from a constituency parse. There are no nodes for 'the cat' (NP) or 'on the mat' (PP). Instead, the relationships are direct: *The* modifies *cat*, and *cat* is the subject of *sat*. This direct, word-to-word representation can be particularly advantageous for languages with more flexible word order, where constituents may not be contiguous. Furthermore, the dependency labels often map more directly to the semantic roles that words play, making the output of a dependency parser a valuable input for downstream tasks like information extraction and question answering.

The output of a dependency parser is a dependency graph (or dependency tree), a structure that represents the grammatical relationships within a sentence. Unlike a constituency tree that groups words into phrases, a dependency graph consists of directed, labeled arcs connecting individual words. In this graph, words are the nodes, and the grammatical relationships are the arcs. Each arc connects a head (the word that governs the relationship) to a dependent (the word that modifies or is an argument of the head). By convention, the arrow points from the head to the dependent.

To make this concrete, let's analyze the sentence 'She ate a green apple.' The dependency parse for this sentence is illustrated in Fig. 6.6. The structure is rooted in the main verb of the sentence, ate. A special pseudo-node, ROOT, is typically used to mark the head of the entire sentence, with an arc pointing to this main verb. From ate, arcs extend to its arguments and modifiers.

Fig. 6.6 A dependency graph for the sentence 'She at a green apple'. In the graph, we can observe the following key relationships:

- An arc from ate to She is labeled nsubj, indicating that 'She' is the *nominal subject* of the verb 'ate'. This tells us who performed the action.
- An arc from ate to apple is labeled dobj, identifying 'apple' as the *direct object*. This tells us what was acted upon.
- The noun apple is itself a head. An arc from apple to green is labeled amod (adjectival modifier), showing that 'green' describes the apple.
- Similarly, an arc from apple to a is labeled det (determiner), specifying the article associated with the noun.

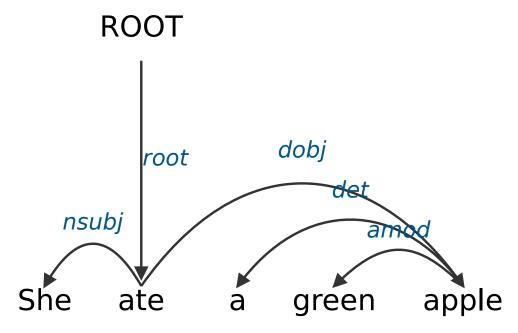


Figure 6.6: A dependency graph for the sentence 'She at a green apple', showing words connected by labeled, directed arcs that represent grammatical relationships.

This representation explicitly captures the functional relationships between words. The triplet nsubj(ate, She) and dobj(ate, apple) directly encodes the core predicate-argument structure: who did what. This is a significant advantage of dependency grammar, as it often makes extracting semantic information more straightforward than navigating a complex phrase structure tree. The labels for these grammatical relations typically come from a standardized set, such as the one defined by the Universal Dependencies (UD) project, which aims to create consistent annotation schemes across different languages. This standardization is crucial for building robust, multilingual parsing tools. The entire structure forms a directed, acyclic graph where each word (except the root) has exactly one head, ensuring a coherent grammatical analysis.

While graph-based methods attempt to find the highest-scoring dependency tree from all possible trees for a sentence, an alternative and highly efficient approach is **transition-based parsing**. Instead of scoring entire trees, a transition-based parser processes a sentence word by word, making a sequence of locally optimal decisions to build the final dependency graph. This method is often conceptualized as a state machine. It starts in an initial configuration, and at each step, it applies a *transition* (or *action*) to move to a new configuration, gradually constructing the dependency arcs. The process terminates when the sentence is fully processed and a complete tree has been formed.

At the heart of a transition-based parser is its configuration, which typically consists of three components:

- A stack (σ) , which holds words that are currently being processed.
- A **buffer** (β) , which holds the sequence of words from the sentence that have not yet been processed.
- A set of dependency arcs(A), which is initially empty and is incrementally built by the parser.

The parser is initialized in a standard starting configuration. The stack contains a special ROOT token, representing the root of the dependency tree. The buffer contains all

the words of the sentence in their original order. The set of arcs is empty. The goal of the parser is to reach a terminal configuration where the buffer is empty and the stack contains only the ROOT token.

The parser moves from one configuration to another by applying one of a small set of defined actions. The most common system, known as an *arc-standard* parser, uses three primary actions:

- 1. **SHIFT**: This action removes the first word from the buffer and pushes it onto the top of the stack. This is the primary way words are brought into consideration for forming dependency relationships. It is used when no dependency can yet be formed involving the word at the front of the buffer and the word on top of the stack.
- 2. **LEFT-ARC** (L_y) : This action asserts a head-dependent relationship. Let the word on top of the stack be w_j and the second word from the top be w_i . The LEFT-ARC action creates a dependency arc $w_j \to w_i$, with the label y. This means w_j is the head and w_i is the dependent. After adding the arc to the set A, the dependent $(w_i,$ the second word) is popped from the stack. This action is typically chosen when w_i has found its head and will not be the head for any other words remaining on the stack or in the buffer.
- 3. **RIGHT-ARC** (R_y) : This action also asserts a head-dependent relationship. As before, let the word on top of the stack be w_j and the second word from the top be w_i . The RIGHT-ARC action creates a dependency arc $w_i \to w_j$, with the label y. Here, w_i is the head and w_j is the dependent. After the arc is added to A, the dependent $(w_j$, the top word) is popped from the stack. This action is used when the word on top of the stack has found its head and has already been assigned all of its own dependents.

The critical question, of course, is how the parser decides which action to take at each step. This decision is made by a component called an **oracle**. In modern transition-based parsers, this oracle is a supervised machine learning classifier. At each step in the parsing process, a feature vector is extracted from the parser's current configuration. These features might include the part-of-speech tags and word forms of the top few words on the stack, the first few words in the buffer, and information about previously constructed dependency arcs.

The classifier, which has been trained on a large *treebank* (a corpus of sentences annotated with correct dependency trees), takes this feature vector as input and predicts the most likely valid action (SHIFT, LEFT-ARC, or RIGHT-ARC) for that specific configuration. The parser then executes this predicted action and moves to the next state. This process repeats until a terminal state is reached.

The primary advantage of this approach is its efficiency. Since the parser makes a single pass through the sentence, making a constant number of decisions for each word, its computational complexity is linear in the length of the sentence, or O(n). This makes transition-based parsers very fast and suitable for large-scale applications. However, their greedy nature—making locally optimal decisions without the ability to backtrack—means that an early mistake can lead to a cascade of subsequent errors, a problem known as error propagation. Despite this, the combination of speed and high accuracy from powerful machine learning oracles has made transition-based parsing a dominant and influential paradigm in the field. To make these abstract components and actions concrete, the next section will walk through a step-by-step example, illustrating how these pieces work together to parse a simple sentence.

To understand how a transition-based parser operates, let's walk through a concrete example. We will parse the simple sentence 'He ate fish.' Our parser uses three data

| Step | Stack | Buffer | Action Taken | Generated Arcs |
|------|-------------------|-----------------|--------------------|--|
| 0 | [ROOT] | [He, ate, fish] | _ | {} |
| 1 | [ROOT, He] | [ate, fish] | SHIFT | {} |
| 2 | [ROOT, He, ate] | [fish] | SHIFT | {} |
| 3 | [ROOT, ate] | [fish] | LEFT-ARC $(nsubj)$ | $\{(ate, nsubj, He)\}$ |
| 4 | [ROOT, ate, fish] | | SHIFT | $\{(ate, nsubj, He)\}$ |
| 5 | [ROOT, ate] | Ö | RIGHT-ARC $(dobj)$ | $\{(ate, nsubj, He), (ate, dobj, fish)\}$ |
| 6 | [ROOT] | | RIGHT-ARC $(root)$ | $\{(ate, nsubj, He), (ate, dobj, fish), (ROOT, root, ate)\}$ |

Figure 6.7: A step-by-step trace of a transition-based parser operating on the sentence 'He ate fish.' The table shows the state of the stack, buffer, and generated dependency arcs after each action.

structures: a *stack*, which holds words that are currently being processed; a *buffer*, which holds the remaining words of the sentence; and a list of *arcs*, which will store the final dependency relations we construct.

The parser begins in an initial configuration. The stack contains a special [ROOT] symbol, which acts as a placeholder for the ultimate root of the sentence. The buffer contains all the words of the sentence in order. The set of generated arcs is empty.

Initial State:

• Stack: [ROOT]

• Buffer: [He, ate, fish]

• Arcs: {}

The parser's job is to sequentially apply one of three actions—SHIFT, LEFT-ARC, or RIGHT-ARC—until it reaches a terminal state, which is defined as having an empty buffer and a stack containing only the [ROOT] element. At each step, a model known as an *oracle* decides which action is correct based on the current configuration. For this walkthrough, we will act as the oracle, making decisions that lead to a grammatically correct parse. The entire process is traced step-by-step in Fig. 6.7.

Step 1: The stack contains [ROOT] and the buffer starts with He. We cannot yet establish a relationship between ROOT and He, as we haven't seen the main verb of the sentence. Therefore, we must move a word from the buffer to the stack for further consideration. The only valid action is SHIFT.

• Stack: [ROOT, He]

• Buffer: [ate, fish]

• Action: SHIFT

• Arcs: {}

Step 2: The top two items on the stack are ROOT and He. We still cannot relate them. The main verb, ate, is still in the buffer. So, we SHIFT again.

• Stack: [ROOT, He, ate]

Buffer: [fish]Action: SHIFT

• Arcs: {}

Step 3: Now, the stack is [ROOT, He, ate]. We can examine the top two words: He and ate. A dependency exists between them: ate is the head of He (the verb is the head of its subject). The arc is thus ate -> He. Since the head (ate) is to the right of the dependent (He) on the stack, this calls for a LEFT-ARC action. We will label this arc with the relation nsubj (nominal subject). This action adds the arc (ate, nsubj, He) to our set and removes the dependent, He, from the stack.

```
Stack: [ROOT, ate]
Buffer: [fish]
Action: LEFT-ARC(nsubj)
Arcs: {(ate, nsubj, He)}
```

Step 4: The stack is now [ROOT, ate] and the buffer still contains fish. We cannot create a dependency between ate and fish because fish is not yet on the stack. The only possible action is to SHIFT.

```
Stack: [ROOT, ate, fish]Buffer: []Action: SHIFTArcs: {(ate, nsubj, He)}
```

• Stack: [ROOT]

Step 5: The buffer is now empty. We must process the remaining items on the stack. The top two words are ate and fish. Here, ate is the head of fish (the verb is the head of its direct object). The arc is ate -> fish. Because the head (ate) is to the left of the dependent (fish), this is a RIGHT-ARC action. We label it with the relation dobj (direct object). This action adds the arc (ate, dobj, fish) to our set and removes the dependent, fish, from the stack.

```
Stack: [ROOT, ate]
Buffer: []
Action: RIGHT-ARC(dobj)
Arcs: {(ate, nsubj, He), (ate, dobj, fish)}
```

Step 6: We are nearly finished. The buffer is empty and the stack contains [ROOT, ate]. The final dependency connects the special ROOT symbol to the main predicate of the sentence, which is ate. The arc is ROOT -> ate. Since the head (ROOT) is to the left of the dependent (ate), this is another RIGHT-ARC action. This final arc establishes ate as the head of the entire sentence.

```
Buffer: []
Action: RIGHT-ARC(root)
Arcs: {(ate, nsubj, He), (ate, dobj, fish), (ROOT, root, ate)}
```

The parser has now reached its terminal state: the buffer is empty and the stack contains only [ROOT]. The set of arcs generated represents the final dependency graph for the sentence 'He ate fish.'

A crucial question remains: how does the parser *decide* which action to take at each step? In our example, we used our own linguistic knowledge. A real transition-based parser uses a machine learning classifier as its oracle. This classifier is trained on a large treebank—a corpus of sentences that have been manually annotated with their correct dependency structures. For each configuration in the training data, the oracle learns the correct historical action.

During parsing, the oracle makes a prediction based on features of the current configuration. These features might include the top words on the stack, the first words in the buffer, their part-of-speech tags, and any existing dependency arcs. Based on these features, the classifier predicts the most likely action (SHIFT, LEFT-ARC, or RIGHT-ARC) to take next. The power of this approach lies in its ability to learn complex grammatical patterns from data, rather than relying on hand-crafted rules. This data-driven method makes transition-based parsers both highly accurate and computationally efficient, as they build the syntactic structure in a single linear pass over the sentence.

In this chapter, we have journeyed through the foundational principles of syntactic parsing, the task of uncovering the grammatical structure of a sentence. We explored two primary paradigms for representing this structure. The first, constituency parsing, views sentences as a hierarchy of nested phrases, or constituents, formally described by Context-Free Grammars. This approach excels at identifying the phrasal units that group together, such as noun phrases and verb phrases. The second, dependency parsing, models grammatical relationships directly as a directed graph of dependencies between words, linking heads to their dependents with specific functional labels like nsubj (nominal subject). This representation is often more useful for downstream tasks that need to know 'who did what to whom.'

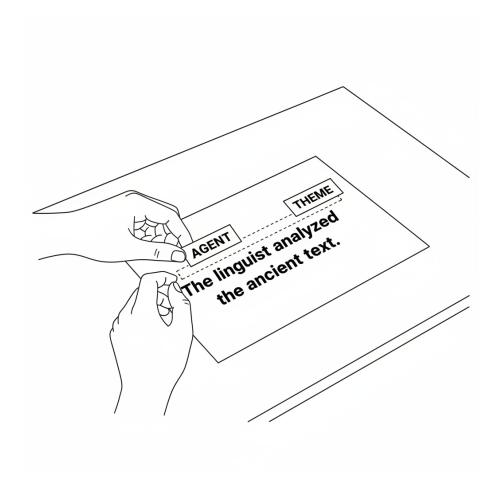
To operationalize these ideas, we examined classic algorithms that form the bedrock of parsing theory. For constituency, we detailed the Cocke-Younger-Kasami (CYK) algorithm, a dynamic programming method that systematically and exhaustively finds all possible parses for a sentence according to a Chomsky Normal Form grammar. For dependency parsing, we investigated transition-based systems, which offer a highly efficient, linear-time approach by processing a sentence from left to right, making a series of local decisions to build a dependency graph using a stack and a buffer. These methods provide a crucial conceptual framework for thinking about syntactic structure computationally.

While these classic techniques are fundamental, the state of the art in syntactic parsing has decisively shifted towards models based on neural networks. The primary advantage of neural parsers lies in their ability to overcome the feature engineering and data sparsity problems that plagued earlier statistical models. Instead of relying on hand-crafted, sparse, symbolic features (e.g., 'the current word is cat' and 'the previous tag was DET'), neural parsers learn dense vector representations—or embeddings—for words, part-of-speech tags, and other contextual cues. These dense vectors capture rich semantic and syntactic similarities, allowing the model to generalize far more effectively from the training data.

In a modern neural transition-based parser, for instance, the decision of which action to take (SHIFT, LEFT-ARC, etc.) is made by a small neural network. This network takes as input the vector embeddings of the words on top of the stack and in the front of the buffer, and then outputs a probability distribution over the possible actions. This approach replaces brittle feature templates with a powerful, learned decision-making function. A different family of neural models, known as graph-based parsers, calculates a score for every possible dependency arc between all words in the sentence and then finds the highest-scoring tree structure.

More recently, large-scale architectures like the Transformer, which we will explore in detail in Chapter 12, have further revolutionized the field. By using mechanisms like self-attention, these models can capture complex, long-distance syntactic dependencies with unprecedented accuracy. The fundamental concepts of constituency and dependency remain as vital as ever, but the computational engine used to uncover these structures is now overwhelmingly powered by the deep learning techniques that define modern computational linguistics.

Chapter 7 Lexical and Compositional Semantics



Our journey so far has equipped us with tools to analyze the structure of language. We can identify parts of speech and map the grammatical architecture of a sentence through parsing. Yet, structure is only half the story. A computer that can parse 'The rover is exploring Martian geology' but cannot infer that a machine is on another planet has not truly *understood* the text. This deeper level of comprehension is the domain of **semantics**, the study of meaning in language. Semantics moves beyond the *form* of language to investigate its *content* and communicative intent. It is the bridge between linguistic expressions and the concepts, objects, and situations they represent in the world.

The distinction between syntax and semantics is fundamental. In Chapter 6, we focused on *syntax*, the set of rules that tells us *how* words can be legally combined to form grammatically valid sentences. Semantics, in contrast, is concerned with what those combinations *mean*. The classic sentence 'Colorless green ideas sleep furiously' serves as a perfect illustration. It adheres flawlessly to the rules of English grammar—it has a subject, verb, and adverb in the right places—making it syntactically impeccable. However, it is semantically incoherent because the meanings of its constituent words clash. 'Ideas' cannot be 'green,' and they certainly cannot 'sleep furiously.' This stark contrast demonstrates that a computational model must go beyond parsing grammatical trees to grapple with the logic, truth, and real-world reference embedded in language.

In this chapter, we will dissect the problem of meaning into two core components. We begin with **lexical semantics**, which focuses on the meaning of individual words, or *lexemes*. How can we represent the meaning of words like 'rover', 'planet', or 'explore' in a way a computer can process? We will explore both classic, knowledge-based approaches and modern, data-driven methods that learn word meanings from vast text corpora. Subsequently, we will tackle **compositional semantics**, which examines how these individual word meanings combine systematically to create the meaning of phrases and sentences. Understanding the meaning of 'rover' and 'explore' is one thing; understanding 'the rover explores' as a complete proposition is another. This part of our exploration will introduce foundational principles and computational techniques for building meaning from the ground up, laying the groundwork for advanced applications like question answering and information extraction.

To computationally model meaning, we must divide the problem into two interconnected challenges, which form the central theme of this chapter. The first is *lexical semantics*, which concerns the meaning of individual words or *lexemes*. This is the task of representing what words like *algorithm*, *analyze*, or *data* signify in isolation. It involves understanding relationships between words (e.g., that a *poodle* is a type of *dog*) and resolving ambiguity (e.g., distinguishing the river *bank* from the financial *bank*). Our initial focus will be on methods that capture these individual word meanings, treating them as the fundamental building blocks for any larger semantic analysis.

The second, more complex challenge is *compositional semantics*. This area explores how the meanings of these building blocks combine, guided by syntax, to form the meanings of phrases and sentences. Knowing the definitions of 'the,' 'model,' 'predicts,' and 'outcome' is essential but insufficient. We must also have a computational mechanism to assemble them to understand the proposition conveyed by 'The model predicts the outcome,' and to recognize that it means something different from 'The outcome predicts the model.' The principle that the meaning of a complex expression is a function of the meanings of its parts and the rules used to combine them is the cornerstone of compositional semantics.

This fundamental division is visualized in **Fig. 7.1**. On one side, we have lexical semantics, mapping individual words to their core concepts. On the other, compositional semantics takes these concepts and, guided by the sentence's grammatical structure, assembles them into a coherent representation of the full sentence's meaning. This chapter

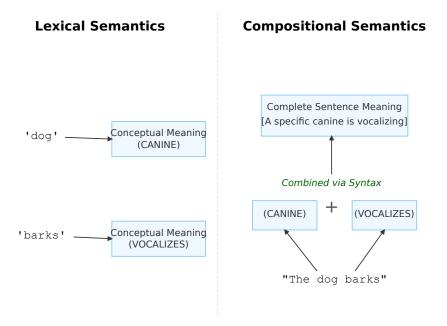


Figure 7.1: A two-panel diagram contrasting lexical and compositional semantics. The left panel, labeled 'Lexical Semantics,' shows individual words ('dog,' 'barks') pointing to their conceptual meanings. The right panel, 'Compositional Semantics,' shows how these individual meanings combine within a sentence structure ('The dog barks') to produce a complete sentence meaning.

follows this natural progression: we will first master techniques for representing word meaning before turning to the methods used to combine them.

Our exploration of lexical semantics begins with the traditional approach: if we want a computer to know what a word means, we can try to explicitly provide that knowledge. This strategy relies on creating large, structured knowledge bases, often called *lexical resources*, which are meticulously handcrafted by linguists and lexicographers. These resources aim to capture the complex network of meanings that constitute a human lexicon, encoding definitions and, more importantly, the relationships between words.

The most influential and widely used resource of this kind is **WordNet**¹, a large lexical database of English developed at Princeton University. WordNet is far more than a simple digital dictionary. While a dictionary is typically organized alphabetically, WordNet is organized conceptually. It attempts to model how humans store and relate lexical knowledge. Instead of treating words as isolated entries, it groups them into sets of synonyms and connects them through a rich network of semantic relationships.

The fundamental idea is that a word's meaning is not defined in a vacuum but by its connections to other words. For example, knowing that a poodle is a type of dog, which is a type of canine, which is a type of mammal, provides a deep, structured understanding that a simple definition alone cannot. WordNet represents this knowledge as a large graph, where concepts are the nodes and semantic relations are the edges. This allows a computational system to traverse the graph, inferring relationships and measuring semantic similarity between concepts.

The great strength of such knowledge-based approaches is the high quality and explicit nature of the encoded information. However, they are enormously expensive and time-consuming to create and maintain. They can also struggle to cover niche domains or keep pace with the rapid evolution of language. This challenge sets the stage for the

¹WordNet was created by a team led by psychologist George A. Miller and its development began in 1985.

data-driven, distributional methods we will explore next.

To understand how a resource like WordNet represents meaning, we must look beyond individual words and focus on its fundamental building blocks and the relationships between them. At its core, WordNet is not a dictionary of words but a network of concepts. Its structure is built upon two key elements: synsets and the semantic relations that link them.

The primary unit in WordNet is the **synonym set**, or *synset*. A synset is a group of words or short phrases that are treated as semantically equivalent in a particular context; they all represent the same underlying concept. For example, the words car, auto, automobile, machine, motorcar can be grouped into a single synset that represents the concept of a four-wheeled passenger vehicle. Each synset is more than just a list of words; it also contains:

- A specific **part-of-speech** (noun, verb, adjective, or adverb). The word *back*, for instance, belongs to different synsets depending on whether it's used as a noun (the rear part of the body), an adverb (in the reverse direction), or a verb (to support).
- A gloss, which is a brief, dictionary-style definition of the concept. For our car synset, the gloss might be 'a motor vehicle with four wheels; usually propelled by an internal combustion engine.'
- Example sentences that illustrate the usage of the synset's members.

The real power of WordNet, however, emerges from the rich network of connections between these synsets. These connections are not arbitrary; they represent well-defined semantic relations that allow a machine to navigate the landscape of meaning. The most important of these relations are:

- Hypernymy and Hyponymy: This is the foundational 'is-a' relationship that organizes most nouns and verbs into a large, hierarchical structure. A concept Y is a hypernym of a concept X if X is a type of Y. Conversely, X is a hyponym of Y. For example, the synset motor vehicle is a hypernym of car, auto, ..., because a car is a type of motor vehicle. This relationship is transitive: since a sedan is a hyponym of car, auto, ..., it is also, by extension, a hyponym of motor vehicle. Traversing these links allows a program to generalize (by moving up to hypernyms) or specialize (by moving down to hyponyms).
- Meronymy and Holonymy: This relationship captures the 'part-of' or 'has-a' connection between concepts. A concept X is a meronym of Y if X is a part of Y. Conversely, Y is the holonym of X. For instance, the synset wheel is a meronym of car, auto, ..., because a wheel is a part of a car. Unlike the 'is-a' hierarchy, this relation describes composition rather than classification. WordNet further divides this into different kinds of part-whole relations, such as part_of (engine is a part of a car), member_of (player is a member of a team), and substance_of (wood is the substance of a table).

As illustrated in Fig. 7.2, these relationships form a complex, directed graph. The synset car, auto, automobile appears as a node in this network. A directed edge labeled is-a points from this node to its hypernym, motor vehicle, signifying its place in the classification hierarchy. Simultaneously, other edges, such as one labeled has-part, connect it to its meronyms, like the synset wheel. By following these labeled edges, a computational system can infer that a car is a vehicle and that it has wheels, all without understanding the raw text of the definitions.

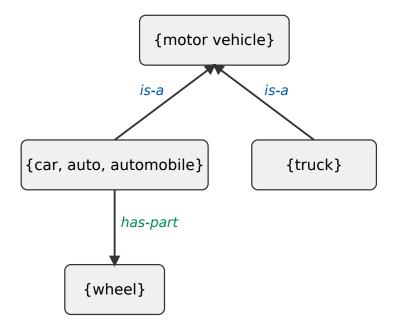


Figure 7.2: A network diagram illustrating semantic relations in WordNet, as described in Fig. 7.4. Nodes represent synsets, such as car, auto, automobile and its hypernym motor vehicle. Directed edges show the 'is-a' (hypernymy) relation from a concept to its more general type and the 'has-part' (meronymy) relation from a whole to its component, such as wheel.

WordNet also encodes other relations tailored to specific parts of speech. For adjectives, a key relation is **antonymy**, which links synsets with opposite meanings (e.g., good is an antonym of bad). For verbs, a crucial relation is **entailment**, where one verb logically implies another (e.g., the act of snoring entails the state of sleeping). This carefully engineered structure, developed over many years at Princeton University, provides a formal, machine-readable model of the human lexicon that has been an invaluable resource for countless computational linguistics tasks.²

The true power of WordNet's structure is realized when we traverse its graph computationally. By programmatically following the pointers that define relations like hyponymy and meronymy, we can explore a word's semantic neighborhood, quantify relationships between concepts, and leverage this knowledge for downstream tasks. A common and intuitive example is to trace the hypernym ('is-a') hierarchy for a given word.

Let's illustrate this by finding the generalization path for the word *car*. We begin by looking up the word in WordNet to find its corresponding synonym sets. The most frequent sense of 'car' is represented by the synset <code>car.n.01</code>, which has the gloss: 'a motor vehicle with four wheels; usually propelled by an internal combustion engine'. This synset serves as our starting node in the graph.

From car.n.01, we can iteratively follow the hypernym pointers to ascend the 'is-a' hierarchy, moving to a more general concept at each step. The direct hypernym of car.n.01 is motor_vehicle.n.01. Following this chain upward reveals a path of increasing abstraction, as depicted in Fig. 7.3. The full traversal from the specific concept of a car to the most general physical concept looks like this:

- $car.n.01 \longrightarrow motor_vehicle.n.01$
- --- wheeled_vehicle.n.01

 $^{^2}$ Miller, G. A. (1995). WordNet: A Lexical Database for English. Communications of the ACM, 38(11), 39-41.

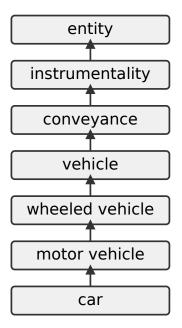


Figure 7.3: A hierarchical graph demonstrating the traversal of WordNet's hypernym hierarchy for 'car'. Arrows point upwards, showing the 'is-a' relationship through a chain of increasingly general concepts, from 'car' at the bottom to 'entity' at the top.

- ullet \longrightarrow vehicle.n.01
- --- conveyance.n.03
- --> instrumentality.n.03
- → artifact.n.01
- \longrightarrow whole.n.02
- ullet \longrightarrow object.n.01
- ullet \longrightarrow physical_entity.n.01
- ullet \longrightarrow entity.n.01

This ability to navigate the semantic hierarchy is invaluable. For instance, we can compute a measure of semantic similarity between words by finding their lowest common ancestor in the hypernym tree. A 'car' and a 'truck' are semantically close because their common ancestor, motor_vehicle, is very specific. This type of reasoning is fundamental for tasks like textual entailment and query expansion, where understanding these hierarchical relationships allows a system to infer meaning beyond literal word matching.

While powerful, knowledge-based resources like WordNet have inherent limitations. They are expensive and time-consuming to create, may not cover all words or specialized domains, and struggle to adapt to the ever-evolving nature of language. To overcome these challenges, modern computational linguistics has largely shifted to data-driven methods for capturing meaning, all of which are built upon a foundational principle known as the **distributional hypothesis**.

The core idea is captured perfectly by the linguist J.R. Firth's famous dictum: 'You shall know a word by the company it keeps.' This hypothesis posits that the meaning

³Firth, J. R. (1957). A synopsis of linguistic theory, 1930-1955.

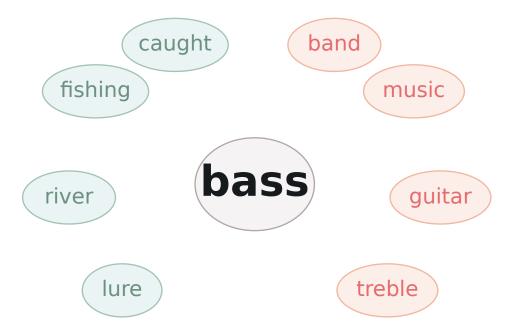


Figure 7.4: An illustration of the distributional hypothesis using the polysemous word 'bass'. The central word is surrounded by two distinct clouds of contextual words. One cloud, containing words like 'fishing', 'river', 'lure', and 'caught', indicates its meaning as a type of fish. The other cloud, with words like 'music', 'guitar', 'treble', and 'band', points to its meaning as a musical instrument.

of a word is not an intrinsic property but is defined by the contexts in which it typically appears. Words that occur in similar contexts are assumed to have similar meanings. For example, consider the words cat and dog. Both frequently appear alongside words like pet, food, walk, play, and vet. This shared contextual environment suggests they occupy a similar semantic space—in this case, the category of common household pets. We don't need a manually curated dictionary to tell us they are related; the pattern of their usage in a large body of text reveals this relationship automatically.

This principle is particularly powerful for handling polysemy, where a single word has multiple meanings. The surrounding words, or the distribution, provide the necessary clues to disambiguate the intended sense. As illustrated in Fig. 7.4, the word bass can refer to a type of fish or a musical instrument. If bass appears in a sentence with words like fishing, river, lure, and caught, its context clearly points to the aquatic meaning. Conversely, if it co-occurs with music, guitar, treble, and band, its meaning as a low-frequency instrument is activated. The context doesn't just suggest the meaning; in the distributional view, it constitutes the meaning.

The profound implication of this hypothesis is that we can transform the abstract concept of meaning into something computationally tractable. If a word's meaning is the set of all its contexts, we can represent that meaning by aggregating those contexts statistically. This insight paves the way for representing words as numerical vectors, where each dimension corresponds to some aspect of their contextual environment. Instead of relying on a predefined graph of semantic relations, we can *learn* these relationships directly from raw, unstructured text, forming the basis for the vector space models we will explore next.

This principle is formalized through **vector space models (VSM)**, a paradigm where words are represented as points or vectors in a high-dimensional geometric space. The fundamental idea is that a word's meaning can be captured by a vector that summarizes its distribution across different contexts. The relationships between words, such as similarity

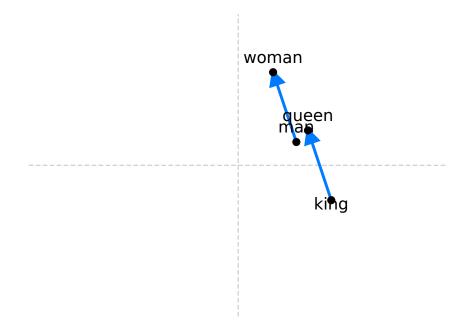


Figure 7.5: A simplified 2D vector space model illustrating semantic relationships. Words like 'king', 'queen', 'man', and 'woman' are plotted as points. The diagram shows that the vector relationship (direction and distance) between 'man' and 'woman' is nearly identical to the one between 'king' and 'queen', capturing the abstract concept of gender as a spatial offset.

or analogy, can then be modeled as geometric relationships between their corresponding vectors.

The most straightforward way to build such a model is by constructing a **term-context co-occurrence matrix**. In this matrix, each row represents a target word from our vocabulary, and each column represents a context word. A 'context' is typically defined as a fixed-size window of words surrounding the target word. A cell at the intersection of a target word w and a context word c contains a count of how many times c appeared within the context window of w across a large corpus.

Each row of this matrix is, therefore, a vector for a target word. For a vocabulary of 50,000 words, the vector for the word 'cat' would have 50,000 dimensions, with each dimension corresponding to a potential context word like 'meow', 'pet', or 'kitchen'. These raw count vectors are typically very sparse, as most words do not co-occur with most other words.

Once we have these vectors, the core insight of the distributional hypothesis—that words appearing in similar contexts have similar meanings—can be measured computationally. We can quantify the similarity between two word vectors, **A** and **B**, using a distance metric. The most common is **cosine similarity**, which measures the cosine of the angle between two vectors:

$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^{n} A_i B_i}{\sqrt{\sum_{i=1}^{n} A_i^2} \sqrt{\sum_{i=1}^{n} B_i^2}}$$

A cosine value near 1 indicates the vectors point in a very similar direction, implying high semantic similarity, while a value near 0 indicates orthogonality, or dissimilarity. This allows us to find the closest words to 'dog' by simply calculating its cosine similarity to all other vectors in the space. As shown in the simplified two-dimensional projection in Fig. 7.5, this spatial arrangement can capture remarkably nuanced relationships. Notice how the vector from 'man' to 'woman' is very similar to the vector from 'king' to 'queen',

preserving the concept of gender as a direction in the space.

While co-occurrence vectors formalize the distributional hypothesis, they have significant practical limitations. A typical vocabulary can contain tens of thousands of words, which means each vector is also tens of thousands of dimensions long. Furthermore, most entries in these vectors are zero, a property known as *sparsity*. These sparse, high-dimensional vectors are computationally inefficient and often fail to capture subtle semantic relationships, as they only model direct, hard-count co-occurrences. They can tell you that car appears near road, but they struggle to generalize that car is also similar to automobile, which might appear in slightly different contexts.

To overcome these issues, we use **word embeddings**. An embedding is a *learned*, *dense*, and relatively *low-dimensional* vector representation of a word. Instead of a vector with 50,000 dimensions where each dimension corresponds to a specific word in the vocabulary, an embedding might have only 50 to 300 dimensions. These dimensions do not have an obvious interpretation; instead, they represent latent semantic features that are learned automatically from the data. The term *dense* means that most values in the vector are non-zero, with each dimension contributing in some small way to defining the word's meaning.

The crucial insight behind modern embeddings is that we can learn these dense vectors by training a simple neural network on an auxiliary prediction task. The goal isn't to build a state-of-the-art prediction model, but to use the task to force the model to learn good word representations. For instance, the model might be tasked with predicting a word given its surrounding context words. In the process of learning to make these predictions, the network's internal weights are adjusted. These learned weights, which represent each word, become the word embeddings. This 'learning' process forces the model to distill the essence of a word's typical contexts into its fixed-size vector.

The power of this approach is that the resulting vector space encodes semantic similarity as geometric proximity. Words that appear in similar contexts, such as boat and ship, will be pushed towards having similar vectors, making them close to each other in the space. Conversely, semantically distinct words like boat and apple will end up with distant vectors. We can measure this distance formally using metrics like cosine similarity, which calculates the cosine of the angle between two vectors:

$$\text{similarity}(\vec{u}, \vec{v}) = \cos(\theta) = \frac{\vec{u} \cdot \vec{v}}{\|\vec{u}\| \|\vec{v}\|}$$

This remarkable property—encoding complex meaning in vector geometry—makes embeddings a foundational component of nearly all modern NLP systems. The following sections will detail two of the most influential algorithms for learning these representations: Word2Vec and GloVe.

The Word2Vec model, introduced by Mikolov et al. at Google in 2013, provided a highly efficient and scalable framework for learning such embeddings. Its key innovation was to reframe the learning problem. Instead of directly counting co-occurrences, Word2Vec trains a simple neural network on a proxy task: predicting a word from its neighbors. The goal is not the task itself, but the learned weights of the network's hidden layer. These weights, once training is complete, become the dense vector representations for each word in the vocabulary. Word2Vec proposes two distinct architectures for this proxy task: the Continuous Bag-of-Words (CBOW) and Skip-gram.

The first architecture, Continuous Bag-of-Words (CBOW), works by predicting a target (center) word from its surrounding context words. For the sentence 'the quick brown fox jumps', if 'fox' is our target word and we use a context window of two, the context words would be 'quick', 'brown', 'jumps', and 'over' (if the next word was 'over'). As shown on the left side of Fig. 7.6, the CBOW model takes the embedding vectors for

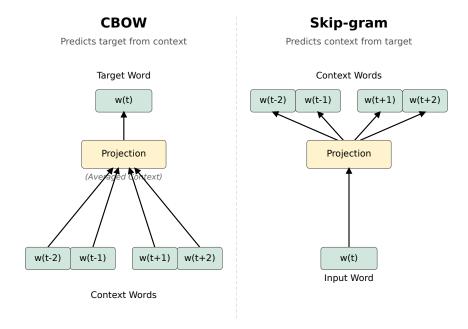


Figure 7.6: A side-by-side comparison of the Word2Vec architectures. On the left, the Continuous Bag-of-Words (CBOW) model predicts a central target word from its surrounding context words. On the right, the Skip-gram model does the reverse, using a single input word to predict its context.

these context words, averages them to create a single context vector, and then uses this vector to predict the target word, 'fox'. The term 'bag-of-words' signifies that the order of the context words is disregarded; their influence is combined through averaging. During training, the model's prediction is compared to the actual target word. The resulting error is then backpropagated to adjust the embedding vectors of the context words, nudging them to be better predictors of that specific target.

The Skip-gram model, depicted on the right of Fig. 7.6, reverses this logic. Instead of predicting the center from the context, it uses a single input word to predict its surrounding context words. Given the input word 'fox', the model attempts to predict 'quick', 'brown', 'jumps', and 'over'. For each word in the training corpus, the Skip-gram model creates multiple training examples: the input word is paired with each of its context words. This means that for a single instance of 'fox', the model learns to predict its neighbors, effectively generating more training data from the same amount of text compared to CBOW. This makes Skip-gram slower to train but allows it to capture the semantics of the input word more precisely, often resulting in higher-quality embeddings, especially for infrequent words.

A naive implementation of either architecture would be computationally prohibitive. The final layer of the network must produce a probability distribution over the entire vocabulary, which can contain hundreds of thousands of words. Calculating the softmax function over this large vector is a major bottleneck. To make training feasible, Word2Vec introduces two key optimizations:

• **Hierarchical Softmax:** This technique replaces the flat softmax layer with a binary tree structure, typically a Huffman tree built from word frequencies. Predicting a word becomes a sequence of binary decisions navigating the tree from the root to the target word's leaf node. This reduces the computational complexity from being proportional to the vocabulary size, \$0(|V|)\$, to being proportional to the logarithm of the vocabulary size, \$0(\log_2|V|).

• Negative Sampling: This is the more common and often more effective approach. Instead of updating weights for the entire vocabulary, the model updates only a small sample. For a given training instance (e.g., target: fox, context: jumps), the model is trained to increase the probability of this 'positive' pair. Concurrently, it randomly selects a few 'negative' samples—words that do not appear in the context (e.g., car, ocean, music). The model is then trained to decrease the probability of these negative pairs. This transforms a complex multi-class classification problem into a much simpler set of binary classification tasks, dramatically speeding up training.⁴

In practice, a choice must be made between the two architectures. CBOW is faster and tends to perform better for frequent words, as the averaging of context vectors acts as a form of regularization. Skip-gram, while slower, excels at learning representations for rare words and is generally considered to produce more robust and higher-quality embeddings overall, making it a popular choice for many applications. Through these architectures and optimizations, Word2Vec provided the first truly practical method for training high-quality, dense word embeddings on massive text corpora.

While the local context window approach of Word2Vec is powerful, another influential model called **GloVe**, short for *Global Vectors for Word Representation*, argues that such predictive models do not efficiently leverage the vast statistical information available in the corpus. Instead of learning from individual, local context windows as they appear, GloVe is a *count-based* model that first aggregates global co-occurrence statistics from the entire corpus and then learns vectors that best explain these statistics.

The process begins by constructing a large co-occurrence matrix, X, from the corpus. An entry in this matrix, X_{ij} , represents the number of times word j appears in the context of word i. This matrix is a direct, global statistical summary of the corpus. The core insight behind GloVe is that the ratios of co-occurrence probabilities can encode meaning more effectively than the raw probabilities themselves. For instance, consider the words ice and steam.

- A probe word like solid is semantically related to *ice* but not steam. We would expect the ratio P(solid|ice)/P(solid|steam) to be very large.
- A probe word like gas is related to steam but not ice. We would expect the ratio P(gas|ice)/P(gas|steam) to be very small.
- A word like water, which is related to both, should have a ratio close to 1.0.
- An unrelated word like fashion should also have a ratio close to 1.0.

GloVe is designed to learn word vectors that capture these ratios as linear relationships in the vector space. The model aims to learn two vectors for each word: a word vector w_i and a separate context vector \tilde{w}_j . The training objective is to learn these vectors such that their dot product models the logarithm of their co-occurrence count. More formally, the model is trained on the following objective:

$$w_i^T \tilde{w}_j + b_i + \tilde{b}_j = \log(X_{ij})$$

Here, b_i and \tilde{b}_j are bias terms for each word and context word, respectively. The model is trained to minimize a weighted least-squares cost function that measures the difference between the two sides of this equation for all word pairs in the vocabulary.

⁴The technique is a simplified variant of Noise Contrastive Estimation (NCE), which aims to train a model to differentiate true data from noise.

$$J = \sum_{i,j=1}^{V} f(X_{ij}) (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log(X_{ij}))^2$$

The weighting function, $f(X_{ij})$, is crucial. It serves two purposes: it prevents the model from being dominated by extremely frequent co-occurrences (e.g., 'the' appearing with 'is'), which provide little semantic information, and it gives zero weight to pairs that never co-occur $(X_{ij} = 0)$, making the computation tractable.

In essence, the fundamental distinction between the two models is this:

- Word2Vec is a *predictive* model. It iterates through the corpus and learns by making local predictions (e.g., predicting a center word from its context). It captures the global statistics *implicitly*.
- GloVe is a *count-based* model. It first pre-computes a global co-occurrence matrix and then directly learns vectors that best explain this explicit statistical information.

This makes GloVe a hybrid method, combining the strengths of traditional count-based matrix factorization methods (like Latent Semantic Analysis) with the local context window methods of Word2Vec. In practice, both GloVe and Word2Vec produce high-quality embeddings that excel at capturing semantic regularities, and the final vector for each word is typically the sum of its word and context vectors $(w_i + \tilde{w}_i)$.

One of the most remarkable and illustrative properties of dense word embeddings is their ability to capture semantic relationships as consistent vector offsets. This goes beyond simply placing similar words near each other; it means the geometric arrangement of words in the vector space encodes analogies. This allows us to use simple vector arithmetic to perform a kind of conceptual algebra.

The classic case study that demonstrates this power is the analogy: 'man is to king as woman is to queen.' A well-trained word embedding model learns this relationship not as a rule, but as a geometric configuration. We can express this analogy through vector arithmetic:

$$\vec{v}_{\rm king} - \vec{v}_{\rm man} + \vec{v}_{\rm woman} \approx \vec{v}_{\rm queen}$$

When we perform this calculation, the word whose embedding is closest (typically measured by cosine similarity) to the resulting vector is, indeed, queen. The intuition behind this is that the difference vector, $\vec{v}_{\rm king} - \vec{v}_{\rm man}$, can be thought of as capturing the abstract concept of royalty or monarchy. By adding this 'royalty vector' to the vector for woman, we effectively navigate the semantic space to the location representing a female monarch. The visual representation in **Fig. 7.7** illustrates this process, showing how starting at the 'king' vector, subtracting the 'man' vector, and then adding the 'woman' vector lands us in the immediate vicinity of the 'queen' vector.

This property is not a special case; it generalizes to many other types of relationships, showcasing the rich structure learned by these models. Examples include:

- Country-Capital: $\vec{v}_{\text{Paris}} \vec{v}_{\text{France}} + \vec{v}_{\text{Germany}}$ results in a vector very close to \vec{v}_{Berlin} .
- Verb Tense: $\vec{v}_{\text{walking}} \vec{v}_{\text{walk}} + \vec{v}_{\text{swam}}$ results in a vector close to $\vec{v}_{\text{swimming}}$.
- Company-Product: $\vec{v}_{\text{Microsoft}} \vec{v}_{\text{Windows}} + \vec{v}_{\text{Google}}$ points towards vectors for products like Android or Search.

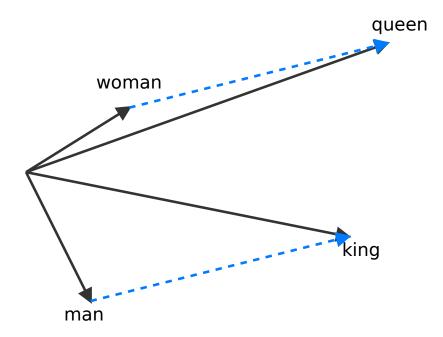


Figure 7.7: A 2D visualization of word vector arithmetic for the analogy 'man is to king as woman is to queen'. The parallel dashed arrows represent the constant vector offset corresponding to the concept of 'royalty'. This geometric arrangement demonstrates how the relationship $v(king) - v(man) \approx v(queen) - v(woman)$ is captured in the vector space.

These analogy-solving capabilities provided early, compelling evidence that distributional models do more than just count co-occurrences. They learn a high-dimensional semantic space where the directions and distances between words are meaningful. This latent structure is what gives word embeddings their power and makes them such a successful foundational component for more complex natural language understanding tasks. The ability to capture relational meaning was a significant leap beyond the capabilities of earlier sparse, knowledge-based representations.

Having explored how to represent the meaning of individual words, we now turn to a fundamental question: how do these individual meanings combine to create the meaning of phrases and sentences? We rarely communicate in isolated words; we combine them into complex structures to express complex ideas. This is the central concern of *compositional semantics*.

The guiding principle in this field is the **Principle of Compositionality**, a concept most often attributed to the German logician and philosopher Gottlob Frege.⁵ The principle states that the meaning of a complex expression is a function of the meanings of its parts and of the syntactic rules by which they are combined. This idea is powerfully intuitive. We understand the sentence 'The student read the book' because we know the meanings of *student*, *read*, and *book*, and we understand how the English grammatical structure assigns roles to each of them. If we change the structure to 'The book read the student,' the sentence becomes nonsensical, even though the component words and their lexical meanings remain identical. The rules of combination are just as important as the meanings of the parts.

We can formalize this principle abstractly. If an expression E is composed of two parts, α and β , its meaning M(E) can be expressed as a function of the meanings of its parts:

⁵While the principle is widely known as 'Frege's Principle,' its explicit formulation in his work is a subject of scholarly debate. Regardless of attribution, it remains the foundational assumption of most formal and computational semantics.

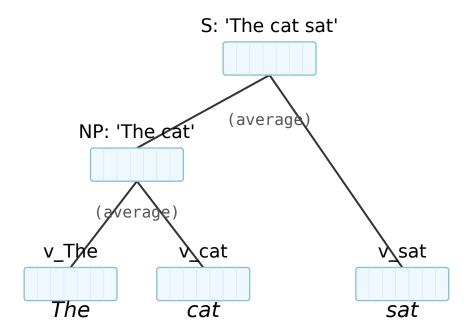


Figure 7.8: A diagram of compositional semantics using vector averaging. The vectors for 'The' and 'cat' are averaged to form a noun phrase (NP) vector. This is subsequently averaged with the vector for 'sat' to create a single vector for the entire sentence (S).

$$M(E) = f_{\text{compose}}(M(\alpha), M(\beta))$$

Here, f_{compose} is the composition function that is guided by the syntactic structure combining α and β . The central challenge for computational semantics is defining this function. Now that we have powerful vector representations for lexical meaning (i.e., word embeddings), the question becomes: what is the right f_{compose} for combining word vectors into phrase or sentence vectors? How do we computationally model the way meaning is constructed? The following sections will explore several approaches to this very problem.

Having established a method for representing individual word meanings as vectors, we now turn to the challenge of compositional semantics. How do we combine these word vectors to create meaningful representations for multi-word expressions like phrases and sentences? Frege's principle suggests we need rules for combining the meanings of the parts. The most direct computational approach is to apply simple vector arithmetic.

Two common baseline methods for composing a sequence of word vectors $\vec{v}_{w_1}, \vec{v}_{w_2}, \dots, \vec{v}_{w_n}$ into a single phrase vector \vec{v}_{phrase} are element-wise addition and averaging.

• Vector Sum: The phrase vector is the sum of its constituent word vectors.

$$\vec{v}_{\text{phrase}} = \sum_{i=1}^{n} \vec{v}_{w_i}$$

• Vector Average: The phrase vector is the average of its constituent word vectors.

$$\vec{v}_{\text{phrase}} = \frac{1}{n} \sum_{i=1}^{n} \vec{v}_{w_i}$$

These methods are intuitive and computationally inexpensive. For example, as illustrated in **Fig. 7.8**, we could compute a vector for the noun phrase 'The cat' by averaging the vectors for *The* and *cat*. This new phrase vector could then be combined, perhaps

again through averaging, with the vector for *sat* to produce a single vector representing the entire sentence. This technique effectively creates a 'bag of vectors,' where the final representation is an amalgamation of all word meanings in the phrase, disregarding their structural relationships.

Despite their simplicity, these additive and averaging methods have severe limitations. Most critically, they are *insensitive to word order*. Because vector addition is commutative, the phrase 'the student taught the professor' would yield the exact same vector representation as 'the professor taught the student,' even though their meanings are entirely different. The compositional operation fails to capture the crucial syntactic structure that assigns the roles of teacher and learner. Similarly, the sentence 'The cat is not on the mat' would likely produce a vector very similar to 'The cat is on the mat,' as the vector for *not* is simply averaged in with the others, failing to properly negate the proposition. These simple operations cannot adequately model the complex semantic impact of syntax, function words, and logical operators. While useful for tasks where word order is less critical (like document classification), they are insufficient for capturing the nuanced meaning required for true sentence understanding, motivating the need for more structured approaches.

While simple compositional techniques like vector addition can capture the general topic of a phrase, they are fundamentally limited. The meaning of a sentence like 'The cat chased the mouse' is not merely the sum of its parts; it is a structured event with distinct participants playing specific roles. Averaging the vectors for cat, chased, and mouse would give us a representation in the semantic neighborhood of felines, rodents, and pursuit, but it critically fails to preserve the information of who did the chasing and who was chased. To capture this vital relational information, we need a more sophisticated model of sentence-level semantics.

This brings us to the task of **Semantic Role Labeling (SRL)**, a process that aims to uncover the underlying predicate-argument structure of a sentence. Often described as determining 'who did what to whom, where, when, and how,' SRL provides a shallow semantic representation of events. The core idea is to identify the main predicate of a sentence—typically a verb—and then identify and label the phrases (arguments) that fill the semantic roles associated with that predicate. This moves beyond simple co-occurrence to model the explicit relationships between participants in an event.

Consider the sentence: 'The startup will launch its innovative new app in Singapore next quarter.' An SRL system would first identify the predicate, which is the event-denoting word launch. It would then identify and label the arguments associated with this predicate's frame of meaning:

- Agent (The doer): The startup
- Theme (The thing being acted upon): its innovative new app
- Location (Where the event happens): in Singapore
- Time (When the event happens): next quarter

By identifying these roles, SRL extracts a structured representation of the sentence's meaning: LAUNCH(Agent: The startup, Theme: its innovative new app, Location: in Singapore, Time: next quarter). This structured output is far more useful for downstream tasks like question answering or information extraction than a single sentence vector.

It is crucial to distinguish semantic roles from syntactic roles. Syntactic parsing identifies grammatical functions like *subject* and *direct object*, which are tied to the sentence's structure. Semantic roles, in contrast, are tied to the underlying meaning. For example, consider these two sentences:

| Role | Description | Annotated Example |
|-----------------|--|--|
| Agent | The entity that willfully instigates or performs the action. | [The chef] _{AGENT} cooked the meal. |
| Patient / Theme | The entity affected by the action or that undergoes a change of state. | The chef cooked [the meal] $_{\mathrm{PATIENT}}$. |
| Instrument | The means or tool used to perform the action. | He sliced bread [with a knife] Instrument. |
| Recipient | The entity that receives something as a result of the action. | She gave the book to [her friend] RECIPIENT. |
| Location | The spatial setting where the action occurs. | They met [in the park]_LOCATION. |
| Source | The origin or starting point of a motion or transfer. | He walked [from the library] $_{\rm SOURCE}$. |
| Goal | The endpoint or destination of a motion or transfer. | He walked to [the station] $_{\mathrm{GOAL}}$. |

Figure 7.9: A summary of fundamental semantic roles common to frameworks like Prop-Bank and FrameNet, with descriptions and examples.

- 1. *[The researcher] Subject published [the paper] Object.
- 2. *[The paper] Subject was published by _[the researcher] Object of Preposition.

Syntactically, the subject of sentence (1) is 'The researcher,' while the subject of sentence (2) is 'The paper.' However, in terms of meaning, 'The researcher' is the *Agent* (the one doing the publishing) in both sentences, and 'the paper' is the *Patient* or *Theme* (the thing being published). SRL abstracts away from these syntactic variations to produce a consistent semantic representation, capturing the fact that both sentences describe the very same event. This ability to normalize meaning across different surface forms makes SRL a powerful tool for computational language understanding. The next section will explore the lexical resources that formally define these roles.

To formalize the predicate-argument structure identified by Semantic Role Labeling, we require a standardized inventory of roles. A system cannot simply invent labels like 'doer' or 'thing being done' on the fly; it needs a consistent lexicon to describe the part each constituent plays in the event. Two major resources, PropBank and FrameNet, provide comprehensive frameworks for defining and annotating these roles, albeit with different philosophical approaches. **Fig. 7.9** provides a summary of several fundamental roles common to these frameworks, offering a quick reference for the concepts discussed here.

The Proposition Bank, or PropBank, takes a verb-centric approach. It develops a specific set of roles for each individual verb, meaning the definition of a role is tied to the verb that introduces it. However, PropBank maintains consistency by defining general, numbered 'arguments' whose semantic meaning is broadly similar across verbs. The core roles are:

- **Arg0:** Typically the Agent or 'doer' of the action. This corresponds to the *Proto-Agent*, the entity that instigates or controls the event.
- Arg1: Typically the Patient or 'theme' of the action. This corresponds to the *Proto-Patient*, the entity that is affected by the event or undergoes a change of state.
- Arg2, Arg3, ...: These higher-numbered arguments represent other essential participants whose meanings are highly specific to the verb. For the verb give, Arg2 would be the recipient. For move, it might be the destination.

Consider the verb to eat:

[The students] Arg0 ate [the pizza] Arg1 [in the cafeteria] ArgM-LOC.

Here, the students are the eaters (Arg0) and the pizza is the thing eaten (Arg1). PropBank also includes a set of general-purpose modifier arguments, prefixed with **ArgM**, that capture circumstantial information not essential to the core verb meaning. Common modifiers include **ArgM-LOC** for location, **ArgM-TMP** for time, **ArgM-MNR** for manner, and **ArgM-DIR** for direction. These modifiers can be attached to nearly any verb's argument structure.

An alternative, more semantically-driven resource is FrameNet. Instead of focusing on individual verbs, FrameNet is organized around **semantic frames**, which are schematic representations of situations or events (e.g., a commercial transaction, a journey, a conflict). Each frame defines a set of **Frame Elements (FEs)**, which are the semantic roles specific to that conceptual scenario. A key insight of FrameNet is that multiple verbs, nouns, or adjectives can evoke the same frame.

For example, the COMMERCE_BUY frame describes a situation involving a commercial transaction. Its core Frame Elements include the Buyer, the Seller, the Goods, and the Money. This single frame can be evoked by various words:

The customer Buyer bought [a new laptop]Goods from [the store]Seller.

The store **Seller sold** [a new laptop]**Goods** to [the customer]**Buyer**.

1. The price of [the new laptop]Goods was [\$1200]Money.

Notice how FrameNet captures the same underlying semantic event even when the syntactic structure and triggering words change dramatically. In the first two sentences, the Buyer and Seller switch between being the subject and an object of a prepositional phrase, but their semantic roles within the transaction frame remain constant. This abstraction from syntax to a deeper semantic level is FrameNet's primary strength.

In practice, PropBank's verb-centric, generalized roles provide broader coverage and are often simpler to annotate, making it a highly practical resource for building robust SRL systems. FrameNet offers a more fine-grained and semantically nuanced representation, which can be more powerful for deep understanding tasks, but its complexity means it may not cover as wide a range of linguistic phenomena. By converting a sentence from a linear string of words into a structured representation like this—a predicate with its labeled arguments—we equip a machine to begin reasoning about the events described in the text, a crucial step towards true language understanding.

To see the power of Semantic Role Labeling (SRL) in action, let's consider a practical case study: building a system to automatically extract information about corporate acquisitions from financial news. The goal is to move beyond simple keyword searches and create a structured database of these events that can answer sophisticated questions.

Imagine our system processes the following sentence from a news wire:

Innovate Corp, the cloud-computing giant, finalized its acquisition of the AI startup Synapse Dynamics for \$2.5 billion on Tuesday.

A traditional information retrieval system might identify keywords like 'Innovate Corp', 'acquisition', and 'Synapse Dynamics'. However, it would struggle to understand the precise relationship between them. Was Innovate Corp acquired? Did the deal fail? How much was it worth?

This is where an SRL system provides a much deeper level of understanding. By focusing on the predicate—the action word finalize (or its underlying verb acquire)—the system identifies and labels the semantic roles of the surrounding phrases. The output would look something like this:

- Predicate: finalize (within the event frame of Acquisition)
- Agent (the entity performing the action): Innovate Corp, the cloud-computing giant
- Theme (the entity being acted upon): its acquisition of the AI startup Synapse Dynamics
- ARG-Value (the price or value): for \$2.5 billion
- ARG-Temporal (the time of the action): on Tuesday

This structured output is far more valuable than a list of keywords. The linguistic analysis can be directly mapped into a structured format, like a database record or a JSON object, that represents the core event:

```
{
  "event_type": "ACQUISITION",
  "acquirer": "Innovate Corp",
  "acquired_entity": "Synapse Dynamics",
  "value": 2500000000,
  "currency": "USD",
  "date": "Tuesday"
}
```

By processing thousands of such articles, we transform a chaotic stream of unstructured text into a structured knowledge base. This enables an advanced question-answering capability that would be impossible with keyword matching alone. We can now ask precise questions and get reliable answers:

- 'Who acquired Synapse Dynamics?'
 - The system queries for the acquirer field where acquired_entity is 'Synapse Dynamics'.
- 'List all acquisitions with a value greater than \$1 billion.'
 - The system filters all records where the value field is greater than 1,000,000,000.
- 'What companies did Innovate Corp acquire?'
 - The system retrieves all acquired_entity values from records where acquirer is 'Innovate Corp'.

Crucially, SRL distinguishes roles. If another article stated, 'Synapse Dynamics rejected an acquisition offer from Innovate Corp,' the system would correctly identify 'Innovate Corp' as the Agent of the 'offer' but would not log a completed acquisition. This case study demonstrates how SRL serves as a critical bridge from raw text to actionable knowledge, enabling applications to understand not just what a text is about, but who did what to whom, when, and for how much.

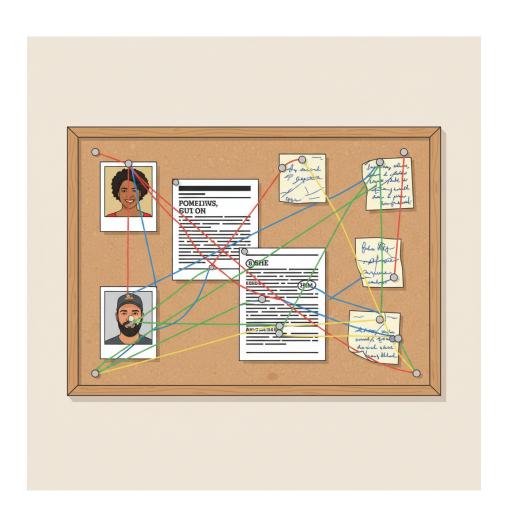
This chapter has navigated the fundamental challenge of representing meaning computationally, charting a course from individual words to the structure of simple sentences. We contrasted two major paradigms: the symbolic, knowledge-based approach of lexical resources like *WordNet*, and the statistical, data-driven approach of *distributional semantics*. The shift from meticulously handcrafted semantic networks to learning dense vector representations like Word2Vec and GloVe directly from text marks a pivotal development

in computational linguistics, allowing models to capture subtle semantic relationships from raw data.

However, these representations are not an end goal. Whether we are using word embeddings to measure semantic similarity or Semantic Role Labeling to identify 'who did what to whom,' we are building foundational components. A robust model of semantics is the bedrock upon which more sophisticated language understanding is built. The ability to analyze discourse, translate between languages, or extract structured information—all topics we will soon explore—depends critically on the techniques for representing meaning that we have established in this chapter.

Chapter 8

Discourse, Coreference, and Dialogue



Thus far, our journey has focused primarily on the sentence as the fundamental unit of analysis. We have learned to model its internal structure through parsing and represent its meaning through compositional semantics. These are powerful tools, but they operate under a significant limitation: human language rarely consists of isolated, self-contained sentences. We communicate through extended narratives, interactive dialogues, and structured documents where each sentence is part of a larger, interconnected whole. To build systems that truly comprehend language, we must move beyond the sentence boundary and analyze the web of connections that link utterances together.

Consider the challenges posed by a simple text snippet: 'Maria trained for months. Her dedication was inspiring. In the end, she won it.' Analyzing the final sentence alone, a syntactic parser can identify its structure, and a semantic model might understand the concept of 'winning.' Yet, fundamental questions remain unanswered:

- Who is she? (Maria)
- What is it that she won? (Presumably, a competition mentioned earlier)
- How does this sentence logically follow from the previous ones? (It is the culmination of her training and dedication)

Answering these questions is impossible without considering the broader context, or discourse. The meaning is not just in the sentences, but between them. This chapter is dedicated to the computational modeling of this larger context. We will explore the principles that make a text a coherent whole, tackle the task of coreference resolution to link pronouns like 'she' to their antecedents, and investigate the structure of dialogue. By understanding these mechanisms, we can build applications that don't just process sentences, but comprehend stories and conversations.

To truly understand language, we must look beyond the boundaries of individual sentences. The meaning of an utterance is rarely self-contained; it is deeply shaped by the surrounding text and the context in which it is used. Two linguistic fields are dedicated to studying this phenomenon: *pragmatics* and *discourse analysis*.

Pragmatics is the study of how context influences the interpretation of meaning. It focuses on the gap between what is literally said (locution) and what is actually meant (illocution). Consider the statement, 'It's cold in here.' While the literal meaning is a simple observation about the temperature, the pragmatic meaning might be an indirect request to close a window or turn up the heat. Pragmatics deals with this kind of implicature, where a speaker's intention must be inferred. For a computer, this is a profound challenge. It requires modeling the speaker's goals, the listener's beliefs, and a vast repository of shared world knowledge—information that is almost never present in the text itself.

Discourse analysis, on the other hand, is concerned with the structure of language beyond the sentence. It examines how sentences are linked together to form coherent and meaningful texts, whether they are conversations, articles, or stories. While syntax, which we explored in Chapter 6, provides the rules for well-formed sentences, discourse analysis provides the principles for well-formed texts. It seeks to answer questions such as:

- How do sentences connect to one another to create a logical flow?
- What makes a paragraph a unified whole, rather than a random list of sentences?
- How are topics introduced, developed, and shifted over the course of a text?

For instance, in the sequence 'Maria arrived late. She had missed the bus,' we effortlessly infer a causal relationship. Discourse analysis aims to formalize these implicit connections that bind a text together.

While distinct, the two fields are closely related and often overlap. Pragmatics often examines the inferential work needed in specific conversational turns, whereas discourse analysis tends to focus on the larger structural properties of a complete text. Both are indispensable for building systems that can comprehend narratives, summarize documents, or engage in meaningful dialogue, as they provide the theoretical tools to model meaning in context.

Consider the following two sequences of sentences:

Sequence A: The cat sat on the mat. It was a fluffy Persian. It began to purr loudly. Sequence B: The cat sat on the mat. The sky is blue. Jupiter is the largest planet in our solar system.

Both sequences are composed of grammatically correct English sentences. Yet, you intuitively recognize Sequence A as a proper text and Sequence B as a random list. The crucial property that Sequence A possesses and Sequence B lacks is **discourse coherence**. Coherence is the quality that makes a sequence of sentences a unified, meaningful whole. It is the underlying logical and semantic connectedness that allows us to interpret a text as a single unit of communication, whether it's a story, an argument, or a set of instructions. Without coherence, we have only a jumble of disconnected propositions.

It is useful to distinguish coherence from a related concept, *cohesion*. Cohesion refers to the explicit linguistic devices—the 'glue'—that create surface-level links between sentences. These include:

- Pronouns: Using 'it' to refer back to 'the cat.'
- Conjunctions: Words like 'and,' 'but,' and 'so' that signal relationships.
- Lexical repetition: Repeating key words or using synonyms.

Sequence A uses cohesion (the pronoun 'it'). However, a text can be coherent with minimal cohesive markers. For example: 'The power went out. The room was plunged into darkness.' No explicit conjunction links these sentences, but the causal relationship is immediately clear to a human reader, making the text coherent. Cohesion, therefore, is one tool for achieving coherence, but coherence itself is a deeper property of meaning and logical flow. It relies on our real-world knowledge, our expectations about how events unfold, and our ability to infer relationships that are not explicitly stated.

From a computational perspective, understanding and generating coherent text is a fundamental challenge. For language understanding tasks like question answering or text summarization, a model must grasp the implicit relationships between sentences to build a complete representation of the text's meaning. For language generation tasks, such as creating summaries or chatbot responses, the system must produce not just grammatically correct sentences, but sentences that follow one another in a logical and natural way. Simply stringing together high-probability sentences often leads to the kind of nonsensical output seen in Sequence B. To build systems that can truly process language, we need ways to model the structure of coherent discourse, moving beyond individual sentences to the relationships that bind them together.

To explain how texts achieve coherence, researchers have developed formal frameworks to model their internal structure. One of the most influential is **Rhetorical Structure Theory (RST)**, developed by William Mann and Sandra Thompson. RST is not just a theory of discourse; it provides a detailed analytical method for describing the relations between different parts of a text. The core premise of RST is that a coherent text can be described by a hierarchical tree structure, where every part of the text serves a specific rhetorical purpose in relation to another part.

RST analysis begins by breaking a text down into its minimal, non-overlapping spans, called *Elementary Discourse Units* (EDUs). An EDU typically corresponds to a single

clause. These EDUs form the leaves of the RST tree. The magic happens in how these leaves are connected. Adjacent spans of text—either individual EDUs or larger, previously combined spans—are linked by **rhetorical relations**. These relations describe the functional connection between them. There are dozens of defined relations, each capturing a specific type of connection. Some common examples include:

- **Elaboration**: One span provides more detail, a definition, or an example related to the other.
- Cause / Result: One span describes a situation that caused the situation in the other.
- Evidence: One span provides information intended to increase the reader's belief in the other.
- Contrast: The two spans present information that is seen as being in contrast.
- **Justification**: One span provides a reason for an action or belief presented in the other.
- Background: One span provides context or background information for the other.

A crucial aspect of RST is the asymmetry of most relations. Each relation connects a nucleus and a satellite. The **nucleus** is the more central and essential part of the text, while the **satellite** is supplementary, providing supporting information. A simple test is to imagine deleting one of the spans; if the text can still function without the satellite, but not without the nucleus, the distinction is clear. For instance, in an Evidence relation, the claim being made is the nucleus, and the information supporting it is the satellite. While most relations are monosatellite, some, like Contrast or Sequence, are multi-nuclear, connecting two or more spans of equal importance.

These components combine recursively to form a complete tree structure that spans the entire text. This process is illustrated in **Fig. 8.1**. At the bottom of the diagram, we see the individual EDUs extracted from the paragraph. Moving up the tree, pairs of spans are joined by a rhetorical relation. For example, one EDU might act as the satellite in an *Elaboration* relation, providing more detail for its neighboring nucleus. This newly formed, larger span then acts as a single unit—perhaps as a nucleus in a *Cause* relation with another span. This continues until the entire text is unified under a single root node, representing the primary rhetorical purpose of the whole passage.

A text is considered coherent under RST if such a comprehensive tree can be constructed. The tree itself is an explicit model of the text's coherence, demonstrating precisely how each clause contributes to the overall message. The automatic creation of these trees, known as RST parsing, is a complex but valuable computational task. Models that can identify these deep rhetorical structures can perform more sophisticated text understanding, enabling applications like automated text summarization (by prioritizing the extraction of key nuclei), more coherent text generation, and advanced argumentation mining.

To make the abstract concept of coherence concrete, let's contrast two short texts. Consider the following paragraph:

Maria woke up feeling anxious. The final exam was only a few hours away. She had studied for weeks, but the material on syntactic parsing was particularly difficult. Taking a deep breath, she decided to review her notes on the CYK algorithm one last time.

This text is *coherent*. Each sentence logically follows from the one before it, creating a unified narrative. The first sentence establishes the main entity (*Maria*) and her emotional state (*anxious*). The following sentences provide the *cause* for this state (the exam), an

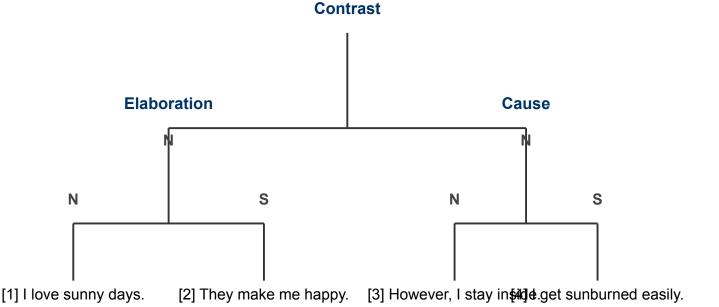


Figure 8.1: A diagram illustrating the hierarchical structure of a text according to Rhetorical Structure Theory (RST), referred to as Fig. 8.4. The analysis breaks the text into Elementary Discourse Units (EDUs) at the leaf nodes. These units are recursively combined into larger spans through rhetorical relations like 'Elaboration', 'Cause', and 'Contrast'. Each relation connects a central Nucleus (N) with a supplementary Satellite (S), or in the case of 'Contrast', multiple Nuclei.

elaboration on the cause (the difficult material), and a resulting action (reviewing her notes). The sentences are linked by a consistent topic, causal relationships, and the flow of time. They work together to build a single, understandable situation.

Now, consider this second text, composed of individually sound sentences:

Maria woke up feeling anxious. The capital of Peru is Lima. Most modern cars use internal combustion engines. Taking a deep breath, she decided to review her notes on the CYK algorithm one last time.

This text is *incoherent*. While each sentence is grammatically correct, there is no logical connection or shared context linking them. The mention of Peru's capital and car engines is entirely unrelated to Maria's anxiety or her study plans. The final sentence abruptly returns to the topic of the first, but the intervening sentences have broken any narrative or logical thread. There are no clear discourse relations—no cause, elaboration, or logical progression—between consecutive sentences. This example highlights that coherence is not a property of individual sentences but an emergent quality of the entire text, arising from the meaningful relationships that bind its parts into a whole. Without these relationships, we are left with a random collection of facts, not a text.

While coherence frameworks like RST explain how sentences relate logically, another crucial element of a connected text is its ability to consistently refer to the same entities—people, places, and things—as the narrative unfolds. This brings us to the task of **coreference resolution**: the process of identifying all linguistic expressions in a text that refer to the same real-world entity. These expressions can take many forms, from proper names and definite descriptions to pronouns.

Consider the following short passage:

Dr. Anya Sharma, a renowned physicist, published her latest findings. The lead researcher from the institute explained that the discovery could change the field. Friends noted that she had worked on the project for years.

In this text, a human reader effortlessly understands that multiple phrases refer to the same person. The goal of coreference resolution is to make this same connection computationally explicit. The expressions involved are:

- Dr. Anya Sharma (a proper name)
- a renowned physicist (a description)
- her (a possessive pronoun)
- The lead researcher (another description)
- she (a personal pronoun)

All these phrases are *coreferent*—they refer to the same entity. A set of coreferent expressions is called a **coreference chain** or an *entity cluster*. The computational challenge is to design an algorithm that can take a text and automatically group these referring expressions into their correct chains.

Solving this task is fundamental for any application that requires deep language understanding. For a question-answering system to correctly answer 'What did the lead researcher publish?', it must first establish that 'the lead researcher' is the same entity as 'Dr. Anya Sharma'. Similarly, for information extraction, grouping these mentions allows a system to aggregate all known facts about a single, unique individual. Without coreference resolution, a text is just a collection of disconnected statements; with it, a rich and interconnected model of the entities and events within the text begins to emerge.

To systematically solve the coreference problem, we must first establish a precise vocabulary. The fundamental unit is the *referring expression*, which is any noun phrase that points to an entity in the world. These expressions can be proper nouns (e.g., *Ada Lovelace*), pronouns (e.g., *she*, *it*), or descriptive noun phrases (e.g., *the first computer programmer*). Coreference resolution is the task of clustering these referring expressions into sets, where each set corresponds to a single real-world entity.

The most common relationship between referring expressions is *anaphora*, where an expression, called the *anaphor*, refers back to a previously introduced entity, known as the *antecedent*. This backward-looking reference is the backbone of textual coherence. Consider the following:

The system failed the initial test. It was not robust enough.

Here, it is the anaphor, and its antecedent is the system. The relationship links the two sentences together. A less frequent but structurally important pattern is cataphora, which reverses this ordering. In cataphora, a referring expression points forward to an entity that has not yet been mentioned. This is often used for stylistic or dramatic effect:

Although he didn't expect it, **David** won the competition.

In this case, the pronoun he is a cataphor that refers forward to David.

Referring expressions themselves can be categorized, and understanding these categories is crucial for building computational models. The most obvious type of anaphor is a pronoun, leading to what is called *pronominal anaphora*. This includes personal pronouns (he, she, it, they), possessive pronouns (his, her, its, their), and reflexive pronouns (himself, herself).

However, many references are not made with pronouns but with other noun phrases. This is known as *nominal anaphora*. This category is diverse and includes several common patterns:

• Repetition: A noun phrase is simply repeated for clarity.

We need to analyze the algorithm. The algorithm has a high time complexity.

| Term | Definition | Example |
|----------------------|--|---|
| Referring Expression | A noun phrase that points to an entity in the world. | Ada Lovelace, she, the programmer |
| Anaphora | An expression that refers back to a previously introduced entity (the antecedent). | The system failed. It was not robust. |
| Cataphora | An expression that refers forward to an entity that has not yet been in- troduced. | Although he didn't expect it, David won. |
| Antecedent | The entity to which an anaphor refers. | The system failed. It was not robust. |
| Pronominal Anaphora | Anaphora where the referring expression is a pronoun. | We found the book . It was on the shelf. |
| Nominal Anaphora | Anaphora where the referring expression is a noun phrase, not a pronoun. | The automaker issued a recall. The company |

Figure 8.2: Key terminology for coreference resolution, with definitions and examples.

• Synonyms/Related Nouns: An entity is referred to using a synonym, a hypernym (a more general term), or a related word.

The automaker issued a recall. The company cited a manufacturing defect.

• **Definite Descriptions:** An indefinite noun phrase (*a company*) introduces an entity, which is later referred to by a definite noun phrase (*the company*).

We hired a new engineer. The engineer starts on Monday.

Each of these patterns presents a unique challenge for a computational system. Linking a pronoun to its antecedent often relies on grammatical cues like gender and number agreement, whereas resolving nominal anaphora may require world knowledge or access to a knowledge base like WordNet. For a structured summary of these key terms and their relationships, please refer to the table in Fig. 8.2. This vocabulary provides the formal framework needed to discuss and develop algorithms for coreference resolution.

To make the abstract task of coreference resolution concrete, let's walk through a brief narrative case study. The goal is to identify all atextual 'mentions'—the referring expressions—and group them into sets, or *chains*, where each set corresponds to a single real-world entity. This process is fundamental for any system that needs to track participants and objects through a text.

Consider the following short paragraph:

Dr. Elara Aris, a renowned astrophysicist, announced a groundbreaking discovery. She confirmed the existence of an exoplanet with a breathable atmosphere, a finding made possible by the James Webb Space Telescope. Her work has electrified the scientific community. The powerful instrument had been observing the star system for months before it yielded the crucial data. The discovery itself is so significant that it will likely redefine the search for extraterrestrial life.

Even in this simple text, there are multiple entities and numerous expressions referring to them. A successful coreference resolution system must untangle these references to build a coherent model of the text's meaning. As visualized in **Fig. 8.3**, we can trace three distinct coreference chains within this narrative.

Dr. Elara Aris,a renowned astrophysicist, announceda groundbreaking discovery. She confirmed the existence of an exoplanet with a breathable atmosphere, a finding made possible by the James Webb Space Telescope. Herwork has electrified the scientific community. The powerful instrument had been observing the star system for months before it yielded the crucial data. The discovery itself is so significant that it will likely redefine the search for extraterrestrial life.

Entity 1: The ScientistEntity 2: The TelescopeEntity 3: The Discovery

Figure 8.3: An illustration of coreference resolution within a narrative paragraph. Mentions belonging to the same entity are highlighted in a consistent color, revealing three distinct coreference chains: the scientist (blue), the telescope (green), and the discovery (orange).

- Entity 1: The Scientist. The first entity is introduced with the proper name *Dr. Elara Aris*. This mention establishes the entity. Subsequent expressions that co-refer with this entity include:
 - a renowned astrophysicist: An appositive phrase that provides additional information about the same person.
 - She: An anaphoric pronoun whose antecedent is clearly Dr. Elara Aris.
 - Her: A possessive pronoun that also refers back to Dr. Aris. All four of these expressions form a single chain referring to one individual.
- Entity 2: The Telescope. The second entity is introduced as the James Webb Space Telescope. Later mentions include:
 - The powerful instrument: A definite noun phrase that describes the telescope.
 Resolving this requires world knowledge or contextual cues to know that a space telescope is a type of instrument.
 - it: The first instance of the pronoun 'it' in the text. Its antecedent is *The powerful instrument*, and by extension, the telescope itself.
- Entity 3: The Discovery. The third entity is introduced with the indefinite noun phrase a groundbreaking discovery. Other mentions are:
 - a finding: Another noun phrase that acts as a synonym for the discovery.
 - The discovery itself: A repetition that reinforces the reference.
 - it: The second instance of 'it' in the text, referring to the significance of the discovery.

This example highlights the complexity of the task. Pronoun resolution, especially for ambiguous pronouns like it, is a significant challenge. A system must correctly determine whether it refers to the telescope, the discovery, or even the star system. Furthermore, connecting Dr. Elara Aris to a renowned astrophysicist requires recognizing that the latter is a description of the former, a task that often relies on syntactic cues and semantic knowledge. Having established this goal of identifying chains, we now turn to the computational models designed to solve this problem automatically.

With a grasp of the coreference task, we turn to its computational implementation. Before the rise of deep learning, the dominant approach was to frame coreference resolution as a supervised machine learning problem, relying heavily on carefully engineered linguistic features. The most common framework for this was the mention-pair model. This model works by considering every possible pair of mentions (m_i, m_j) in a text, where the candidate antecedent m_i precedes the anaphor m_j . For each pair, a binary classifier is trained to decide: are these two mentions coreferent?

To make this decision, the classifier cannot see the raw text directly. Instead, we must extract a set of descriptive properties—features—that capture the relationship between the two mentions. This process, known as feature engineering, was the art of classic coreference resolution. The system's performance depended entirely on the quality and comprehensiveness of these features, which typically included:

- Lexical Features: These capture surface-level string properties.
 - String Match: A binary feature indicating if the strings of m_i and m_j are identical (e.g., 'Apple' and 'Apple').
 - Head Noun Match: Do the head nouns of the two mentions match? (e.g., 'the tall man' and 'the man in the hat').
 - Substring Match: Is one mention a substring of the other? (e.g., 'Dr. Evans' and 'Evans').
- **Grammatical Features:** These leverage syntactic information, often derived from a parse tree.
 - Mention Type: Are the mentions proper nouns, common nouns, or pronouns?
 A pronoun is very likely to have an antecedent.
 - Number/Gender Agreement: Do the mentions agree in number (singular/plural) and gender (masculine/feminine/neuter)? 'The company' and 'she' are an unlikely match.
 - Syntactic Constraints: Does the syntactic relationship between the mentions obey grammatical rules like Binding Theory? For instance, in 'Dr. Smith saw him,' the pronoun him cannot refer to Dr. Smith, whereas in 'Dr. Smith saw himself,' it must.
- Semantic Features: These features aim to capture meaning compatibility.
 - Semantic Class: Do both mentions belong to the same named entity class (e.g., PERSON, ORGANIZATION)? 'Google' and 'he' would be a poor semantic match.
 - WordNet Similarity: How closely related are the head words in a semantic lexicon like WordNet? This helps link synonyms or related concepts.

• Positional Features:

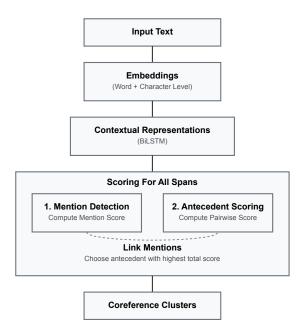


Figure 8.4: A simplified architectural diagram of a modern end-to-end neural coreference resolution model. The diagram shows the flow of information from input text to final coreference clusters. Key components illustrated include: 1) Word and character embeddings, 2) A BiLSTM to create contextual representations, 3) A mention detection module that scores all possible text spans, and 4) An antecedent scoring module that links each detected mention to its most likely antecedent.

- Sentence Distance: How many sentences separate the two mentions? Coreference is much more likely between mentions that are close to each other.

A classifier, such as a decision tree or logistic regression model, is trained on an annotated corpus where true coreference links are marked. It learns to weigh these diverse features to predict the probability that any given pair of mentions is coreferent.

Finally, since the model only makes pairwise decisions, a separate clustering step is needed to form the final coreference chains. A common method is to apply transitive closure: if the model predicts that (A, B) is a coreferent pair and (B, C) is also a coreferent pair, then A, B, and C are all grouped into a single entity cluster. This feature-driven approach was powerful but brittle; it depended on a complex pipeline of other NLP tools (parsers, NER taggers) and an immense, time-consuming effort in feature design.

While classic machine learning models for coreference resolution were a significant step forward, they relied heavily on *feature engineering*. This process involved manually designing hundreds of linguistic features—such as number agreement, gender, semantic class, and syntactic relationships derived from a parser—to help the model make its decisions. This approach was not only labor-intensive and required deep linguistic expertise, but it also resulted in brittle systems that performed poorly on text from new domains where the hand-crafted features were less effective. A more integrated and data-driven approach was needed to overcome these limitations.

The breakthrough came with the development of end-to-end neural network models. The term 'end-to-end' signifies a paradigm shift: instead of relying on a pipeline of separate, pre-processing tools (like POS taggers or parsers), a single neural network learns to perform the entire task, from raw text input to final coreference cluster output. This approach allows the model to learn its own internal representations and features directly from the data, making it more robust and adaptable.

A typical architecture for a modern neural coreference model is shown in **Fig. 8.4**. The process begins by converting the input text into a sequence of rich vector representations. This is done by combining pre-trained *word embeddings* (as discussed in Chapter 7) with *character-level embeddings*, which help the model handle unknown words. These initial vectors are then fed into a bidirectional Long Short-Term Memory network (BiLSTM). The BiLSTM processes the sequence both forwards and backwards, producing a context-aware representation for every word that incorporates information from the entire sentence.

From these contextualized word vectors, the model performs two interconnected tasks:

- 1. **Mention Detection:** The model considers every possible contiguous span of text (e.g., 'Dr. Evans,' 'the professor,' 'she') as a potential mention. For each span, it computes a *mention score*, which represents the model's confidence that the span is a valid referring expression. This is a crucial step, as the model does not know beforehand which words constitute mentions.
- 2. **Antecedent Scoring:** For every span i that is identified as a potential mention, the model considers every preceding span j as a potential antecedent. It then computes a coreference score, s(i,j), for the pair. This score indicates the likelihood that mention i and mention j refer to the same entity.

The scoring function s(i, j) is the heart of the model. It is a neural network component that takes the vector representations of the two spans as input and outputs a single scalar value. The span representations themselves are sophisticated, often created using an attention mechanism over the words within the span. This allows the model to focus on the most important words (e.g., the head noun) when representing a mention. The full scoring function for a mention i can be expressed as:

$$s(i, j) = s_{\text{mention}}(i) + s_{\text{mention}}(j) + s_{\text{antecedent}}(i, j)$$

Here, $s_{\rm mention}$ is the score for an individual span being a mention, and $s_{\rm antecedent}$ is the score for the pairwise link. The model also considers the possibility that a mention is new and has no antecedent by linking it to a special null antecedent, ϵ . During training, the model learns to assign high scores to correct antecedent links found in an annotated corpus. At inference time, for each mention, the model simply chooses the preceding mention (or the null antecedent) that yields the highest score. This elegant integration of mention detection and linking within a single, trainable system has largely replaced feature-based methods, setting the current standard for coreference resolution.

Having explored how meaning is built across sentences in static text, we now shift our focus to a more dynamic and interactive form of discourse: *dialogue*. The principles of coherence and coreference are arguably even more critical here, as a computer must understand and participate in a back-and-forth conversation with a human user. The systems designed for this purpose are known as **dialogue systems** or **conversational agents**. You interact with them daily in the form of virtual assistants like Siri and Alexa, customer service chatbots, and voice-controlled smart home devices.

Broadly, we can distinguish between two major types of conversational agents:

- Task-Oriented Systems: These agents are engineered to help a user complete a specific goal within a well-defined domain. Examples include booking a flight, ordering a pizza, or finding a local restaurant. Their dialogue is structured, purposeful, and aimed at efficiently collecting the information needed to perform an action.
- Open-Domain Chatbots: In contrast, these systems, often called 'chit-chat' bots, aim to engage in general, unconstrained conversation. Their primary goal is not

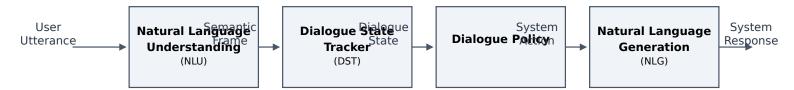


Figure 8.5: A diagram of the pipeline architecture for a traditional task-oriented dialogue system. It illustrates the sequential flow of information from the user's utterance through four core modules: Natural Language Understanding (NLU), Dialogue State Tracker (DST), Dialogue Policy, and Natural Language Generation (NLG), culminating in a system response back to the user.

task completion but rather social engagement, providing companionship, or simply mimicking human-like interaction for entertainment.

While open-domain systems pose fascinating and complex research challenges, the foundational principles of dialogue modeling are best illustrated through the modular architecture of task-oriented systems. They provide a clear framework for applying core computational linguistics concepts to a dynamic setting. In the following sections, we will dissect the typical components of such a system to understand precisely how it interprets user requests, manages the conversation's state, decides what to do next, and generates a coherent response.

Having explored how sentences connect to form a coherent text, we now shift our focus to modeling an entire conversation. The principles of discourse are central to building dialogue systems, also known as conversational agents or chatbots. While open-domain chatbots like those in entertainment aim for general conversation, many of the most practical systems are task-oriented. Their goal is to help a user complete a specific task, such as booking a flight, ordering a pizza, or finding a restaurant. Historically, these systems are built using a modular pipeline architecture, where the problem is broken down into a series of distinct components. As illustrated in Fig. 8.5, this classic architecture consists of four main modules: Natural Language Understanding (NLU), a Dialogue State Tracker (DST), a Dialogue Policy, and Natural Language Generation (NLG).

The first component, Natural Language Understanding (NLU), serves as the system's 'ears.' Its responsibility is to interpret the user's raw text utterance and convert it into a structured, machine-readable format. This process typically involves two key sub-tasks:

- Intent Classification: This is the task of identifying the user's primary goal or purpose. For example, in the utterance 'I need to book a flight to Berlin,' the intent is book_flight. Other possible intents in a travel system might include check_booking_status or ask_for_weather.
- Slot Filling: This task, which is a form of information extraction, identifies and extracts the key pieces of information, or *slots*, from the utterance that are necessary to fulfill the intent. In our example, Berlin would be identified as the value for the destination slot.

The output of the NLU module is a *semantic frame*, a simple data structure that combines the intent and the filled slots. For the utterance 'I need a flight to Berlin for next Tuesday,' the NLU would produce a frame like: {intent: "book_flight", destination: "Berlin", departure_date: "next Tuesday"}. This structured representation is far easier for the subsequent components of the system to process than the original, ambiguous natural language text.

The semantic frame is then passed to the **Dialogue State Tracker (DST)**, which acts as the system's memory. The DST's job is to maintain the *dialogue state* throughout the conversation. The state is a summary of everything that has happened so far, aggregating information across multiple turns. It tracks all the constraints the user has specified (e.g., destination, time, price range) and what the system's last action was. For each turn, the DST takes the NLU output from the current user utterance and the dialogue state from the previous turn, and produces an updated state. If the user first says 'I need a flight to Berlin' and then, in response to a system query, says 'I want to leave tomorrow,' the DST would ensure the final state contains both pieces of information: {destination: "Berlin", departure_date: "2023-10-27"}.

With an updated state from the DST, the **Dialogue Policy** (π) must decide what the system should do next. The policy is the 'brain' of the operation, mapping the current dialogue state to a specific system action. It functions as a decision-making component. For a state s, the policy selects an action a, such that $\pi(s) = a$. These actions are abstract, formal representations, not yet natural language. Examples of system actions include:

- request(origin_city): Ask the user for a required piece of information.
- confirm(destination: "Berlin"): Explicitly confirm a piece of information with the user.
- inform(flight_options): Provide the user with the information they requested.
- api_call(search_flights): Query an external database or API.

Policies can range from simple, hand-crafted rules (e.g., "IF destination is known but origin is not, THEN request(origin)") to complex statistical models learned using Reinforcement Learning, where the system learns the optimal action to take in any state by maximizing a reward signal, such as successful task completion.

Finally, the abstract action selected by the policy is sent to the **Natural Language Generation (NLG)** module, which serves as the system's 'mouth.' The NLG's sole task is to translate the system action into a fluent, human-readable text response. For example, the action request(origin_city) might be converted into 'Where will you be departing from?' or 'What city are you flying out of?'. The simplest NLG systems use templates with slots to fill (e.g., 'OK, a flight to [destination]. Where from?'). More sophisticated neural models can generate more varied and natural-sounding responses, making the interaction feel less robotic. This generated text is then presented to the user, completing one full turn of the dialogue and preparing the system for the user's next utterance.

To make these components concrete, let's trace a simplified interaction with a task-oriented dialogue system designed to book flights. The system's goal is to fill a set of required *slots* (e.g., origin, destination, date) to complete the task. The entire flow of this conversation, showing how information is processed at each turn, is detailed in **Fig. 8.6**.

Imagine the user begins with the following utterance:

• User: 'I want to book a flight to Berlin next Tuesday.'

The system's pipeline springs into action. First, the **NLU** module processes this utterance to extract the user's *intent* and any relevant information. It identifies the intent as book_flight and fills two slots: destination: Berlin and date: 2024-10-29 (assuming the system can resolve 'next Tuesday' to a specific date).

Next, this structured information is passed to the **Dialogue State Tracker (DST)**. The DST maintains the system's belief about the conversation's state. It receives the NLU output and updates its internal record, which now indicates that the destination and

| Speaker | Utterance | $egin{array}{ll} { m NLU} & { m (Intent} & \& \ { m Slots}) \end{array}$ | Dialogue State (Updated Slots) | Policy (System Action) | NLG (System Re sponse) |
|---------|--|--|--|------------------------|---|
| User | 'I want to book a flight to Berlin next Tuesday.' | _ | _ | | _ |
| System | _ | <pre>intent: book_flight slots: { destina- tion: Berlin, date: 2024-10-29}</pre> | {destination: Berlin, date: 2024-10-29, origin: null} | request(origin) | 'Certainly. Wher will you be departin from?' |
| User | 'I'm flying from San Francisco.' | _ | _ | _ | _ |
| System | | $slots$: $\{origin$: San Francisco $\}$ | {destination: Berlin, date: 2024-10-29, origin: San Francisco} | confirm_and search | 'Got it. A fligh from San Francisc to Berlin for Octobe 29th. Shall I search?' |

Figure 8.6: Trace of a simplified flight-booking dialogue. The table shows how user input is processed through the NLU, DST, and Policy components to generate a system response over two turns.

date are known. Crucially, the DST is also aware of the required slots for a book_flight intent and recognizes that a key piece of information, the origin city, is still missing.

The updated dialogue state is then fed to the **Dialogue Policy**. This component acts as the system's brain, deciding what to do next. Seeing that the origin slot is empty, the policy's strategy is to request this missing information. It selects an abstract system action, such as request(origin).

Finally, this abstract action is sent to the **NLG** module. The NLG's job is to translate the system's intent into a cooperative, human-readable sentence. For the **request** (origin) action, it might generate:

• System: 'Certainly. Where will you be departing from?'

The user's subsequent response ('I'm flying from San Francisco') would initiate the same cycle again. The NLU would extract origin: San Francisco, the DST would update the state to reflect that all necessary slots are now filled, and the Policy would decide on a new action, like confirm_and_search. This turn-by-turn interaction, as visualized in Fig. 8.6, demonstrates the pipeline's deterministic flow. Each component has a specialized role: NLU interprets, the DST remembers, the Policy decides, and NLG communicates, working in concert to achieve a specific goal.

To conclude, the three pillars of this chapter—discourse analysis, coreference resolution, and dialogue modeling—are not isolated fields but are fundamentally interconnected. Understanding language requires moving beyond the analysis of single sentences. Discourse analysis provides the theoretical framework for understanding how multiple sentences combine to form a *coherent* and *structured* text, explaining the logical and rhetorical links between utterances.

Within this broader structure, coreference resolution serves as a critical enabling task. It provides the mechanism for tracking entities—people, places, and concepts—as they are mentioned and re-mentioned throughout a text or conversation. Without it, a system cannot connect 'the scientist' to a later mention of 'she.'

Ultimately, both of these concepts are prerequisites for building effective dialogue systems. A conversational agent relies on an implicit understanding of discourse structure to generate relevant responses and uses coreference resolution to maintain context across

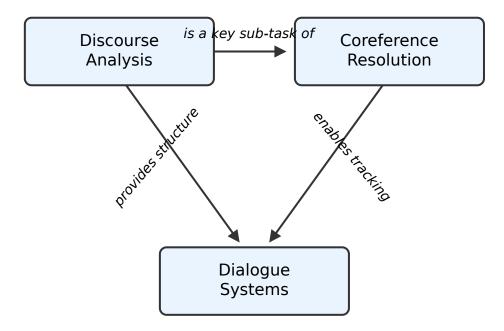
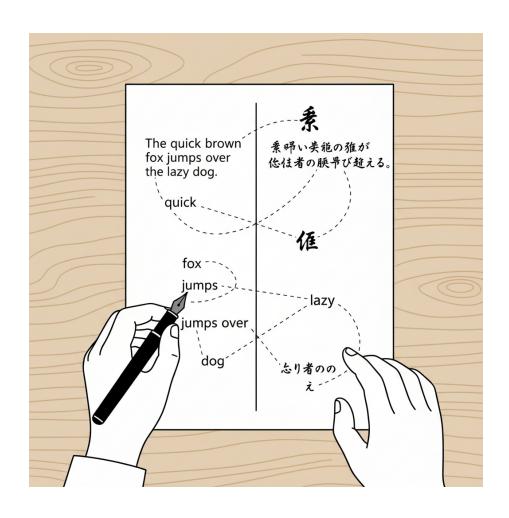


Figure 8.7: A conceptual diagram illustrating the foundational role of Discourse Analysis and Coreference Resolution in the development of Dialogue Systems.

turns. As visually summarized in Fig. 8.7, these foundational components are what elevate a system from a simple command-processor to a truly conversational partner.

Chapter 9

Machine Translation



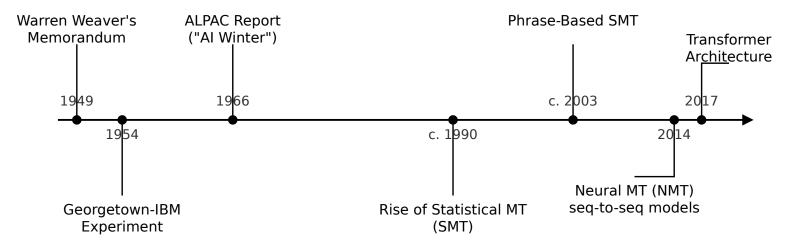


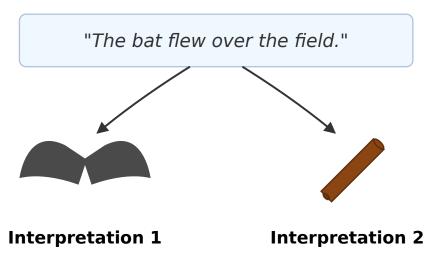
Figure 9.1: A timeline showing key milestones in the history of Machine Translation (MT), from early concepts to modern neural architectures.

Machine Translation (MT) stands as one of the oldest and most ambitious goals in the history of computer science. The dream of automatic translation is nearly as old as the electronic computer itself. Long before the term 'artificial intelligence' was coined, researchers envisioned machines that could break down language barriers, a quest that began in earnest after World War II. The field's origins are often traced to a 1949 memorandum by scientist Warren Weaver, who famously proposed treating translation as a problem of cryptography—essentially 'decoding' a foreign language.

This early optimism fueled efforts like the 1954 Georgetown-IBM experiment, a public demonstration that successfully translated a small set of Russian sentences into English. The initial hype, however, soon met the profound complexity of human language. The influential 1966 ALPAC report concluded that progress was slow and expensive, ushering in a long 'AI winter' for MT research. It was only with the rise of new data-driven statistical and neural methods that the field experienced its modern renaissance. This dramatic history, from early rule-based systems to today's powerful models, is outlined in the timeline in Fig. 9.1.

At its core, machine translation (MT) is the task of automatically converting text from a source language to a target language. The goal is to produce a translation that is both fluent in the target language and faithful to the source text's meaning. Over the decades, the computational approach to this challenge has evolved through three distinct paradigms, each with a fundamentally different philosophy. This chapter will focus on the latter two, but understanding all three provides essential context.

- 1. Rule-Based Machine Translation (RBMT): The earliest systems relied on direct human expertise. RBMT uses vast sets of hand-crafted rules created by linguists, including large bilingual dictionaries for word-level translation and explicit rules for grammatical transformations (e.g., word reordering to match target syntax). This approach is highly interpretable but brittle; it struggles with linguistic ambiguity and exceptions, and requires immense manual effort to build and maintain.
- 2. Statistical Machine Translation (SMT): The dominant paradigm from the 1990s to the mid-2010s, SMT reframes translation as a probabilistic problem. Instead of rules, SMT systems learn how to translate by analyzing massive parallel corpora. The guiding principle is to find the most probable translation \mathbf{t} for a given source sentence \mathbf{s} , a task often modeled using Bayes' theorem: $P(\mathbf{t}|\mathbf{s})$. These systems learn word and phrase correspondences directly from data.



Spanish: El murciélago voló... Spanish: El bate voló...

(bat = flying mammal) (bat = sporting equipment)

Figure 9.2: An illustration of lexical ambiguity using the English sentence 'The bat flew over the field.' The word 'bat' can be interpreted as either a flying mammal (murciélago in Spanish) or as sporting equipment (bate in Spanish), leading to two different translations.

3. Neural Machine Translation (NMT): The current state-of-the-art, NMT uses deep neural networks for translation. Architectures like sequence-to-sequence models encode the source sentence into a dense vector representation and then decode this representation into the target language. This holistic, end-to-end approach allows the model to capture complex context and long-range dependencies, overcoming many limitations of the more fragmented SMT pipeline.

At its heart, translation is a process of mapping meaning from a source language to a target language. If language were a simple code where each word had a single, unambiguous counterpart, this would be a trivial task. However, the reality is far more complex. The fundamental difficulty of machine translation stems from the inherent ambiguity of human language, which manifests in several critical ways.

The most straightforward challenge is *lexical ambiguity*, where a single word can have multiple meanings. Consider the English word 'bat.' As shown in Fig. 9.2, the sentence 'The bat flew over the field' presents two distinct interpretations. In Spanish, these meanings require entirely different words: *murciélago* for the flying mammal and *bate* for the sporting equipment. Without sufficient context, an MT system cannot know which translation is correct. Similarly, words like 'bank' (a financial institution or a river's edge) or 'light' (not heavy or a form of energy) present the same problem. The system must learn to disambiguate based on the surrounding words.

Equally challenging is *structural ambiguity*, where the grammatical structure of a sentence allows for multiple interpretations. The classic example is 'I saw the man with the telescope.' This sentence could mean:

- 1. I used a telescope to see the man.
- 2. I saw a man who was holding a telescope.

A human speaker often resolves this ambiguity unconsciously, but an MT system must make an explicit choice. This is crucial because the two interpretations might require vastly different grammatical constructions in the target language. Finally, languages fundamentally differ in their syntactic structure, such as word order. English follows a Subject-Verb-Object (SVO) order ('She reads the book'), while languages like Japanese use Subject-Object-Verb (SOV) and Irish uses Verb-Subject-Object (VSO). These differences require more than a simple reordering of words; they demand a deep analysis of the sentence's grammatical structure. These challenges demonstrate that translation is not merely a process of substitution but a complex task of interpretation and generation.

Modern machine translation systems, particularly the statistical and neural paradigms we will explore, learn to translate not from hand-crafted linguistic rules but from data. The essential resource for this data-driven approach is the *parallel corpus*, also known as a bitext. A parallel corpus is a collection of texts in a source language presented alongside its translation in a target language.

Crucially, these texts are aligned, most commonly at the sentence level. This means each sentence in the source text is explicitly linked to its corresponding translated sentence. For example:

- EN: The book is on the table.
- DE: Das Buch liegt auf dem Tisch.

By processing millions or even billions of these aligned sentence pairs, MT models discover powerful statistical patterns. They learn the probability that a source word like 'book' corresponds to a target word like 'Buch,' and that the phrase 'on the table' often translates to 'auf dem Tisch.'

Sources for large-scale parallel corpora are often official documents that require translation by law, such as the proceedings of the Canadian Parliament (the Hansard corpus) or the European Parliament (Europarl). Religious texts and translated literary works also serve as valuable sources. The size and quality of this training data are paramount; they are arguably the single most important factors influencing the performance of a modern MT system.

How do we automatically measure the quality of a translation? While human evaluation is the gold standard, it is slow, expensive, and difficult to replicate consistently. For the rapid development and comparison of systems, we need automated metrics. The most influential of these is **BLEU** (Bilingual Evaluation Understudy). The core intuition behind BLEU is straightforward: a good machine translation will share a significant number of words and phrases with a professional human translation.

BLEU measures the correspondence between a machine-generated translation (the candidate) and one or more high-quality human translations (the references). Its score is based on a modified form of n-gram precision. For different n-gram sizes (typically 1-grams to 4-grams), we calculate the proportion of n-grams in the candidate that also appear in any of the reference translations. To prevent systems from gaining an unfair advantage by over-generating common words, this precision is 'clipped.' For example, if the word 'the' appears once in a reference but three times in the candidate, its count for the precision calculation is clipped to one.

Furthermore, a system could achieve high precision by producing a very short but accurate output. To counteract this, BLEU incorporates a *brevity penalty* (BP) that penalizes candidate translations that are shorter than their corresponding references. The final score combines these elements as a geometric mean of the modified n-gram precisions, multiplied by the brevity penalty:

$$BLEU = BP \cdot \exp\left(\sum_{n=1}^{N} w_n \log p_n\right)$$

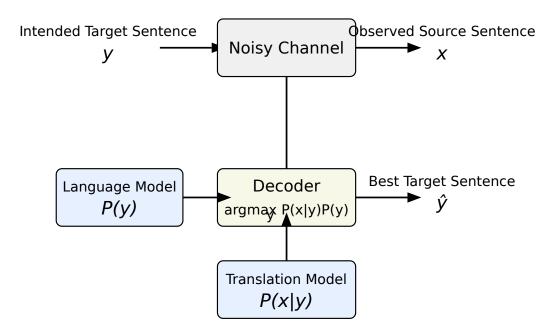


Figure 9.3: A diagram of the noisy channel model for Statistical Machine Translation. The top half shows the conceptual model: an intended target sentence 'y' is transformed by a 'noisy channel' into an observed source sentence 'x'. The bottom half illustrates the decoding task: a decoder takes 'x' as input and uses a Language Model P(y) and a Translation Model P(x|y) to find the most probable target sentence 'ŷ'.

Here, p_n is the modified n-gram precision for n-grams of size n, and N is typically 4. While indispensable for research, it's crucial to remember that BLEU is an imperfect proxy for quality and does not always correlate perfectly with human judgments of fluency or adequacy.

The breakthrough of Statistical Machine Translation (SMT) was to reframe translation not as a linguistic puzzle of rule application, but as a problem of probability. Given a source sentence x (e.g., in French), the goal is to find the target sentence y (e.g., in English) that is most probable. We want to find the y that maximizes the conditional probability P(y|x):

$$\hat{y} = \operatorname*{arg\,max}_{y} P(y|x)$$

Directly modeling P(y|x) is difficult. The key insight of SMT was to use Bayes' theorem to decompose this problem into two more manageable components:

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

Since we are searching for the best y for a fixed input x, the denominator P(x) is a constant and can be ignored. Our maximization problem thus becomes:

$$\hat{y} = \arg\max_{y} P(x|y)P(y)$$

This formulation is famously known as the *noisy channel model*, a concept borrowed from information theory (see Fig. 9.3). We imagine that the 'true' message is the target sentence y, which was corrupted by a 'noisy channel' to produce the source sentence x that we observe. The decoder's task is to recover the original message by combining two separate models:

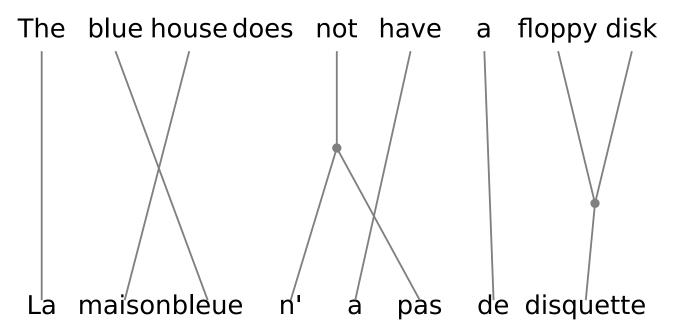


Figure 9.4: An illustration of word alignment between an English source sentence and its French translation, as described in Fig. 9.7. The lines show different alignment patterns: one-to-one (e.g., *house* \rightarrow *maison*), one-to-many (e.g., *not* \rightarrow *n' pas*), and many-to-one (e.g., *floppy disk* \rightarrow *disquette*). The English word *does* is unaligned, demonstrating a NULL alignment.

- The Translation Model, P(x|y): This component answers the question: how likely are we to see the source sentence x if the original sentence was y? It models the *faithfulness* of the translation. For instance, it would assign a high probability to P(le chat|the cat).
- The Language Model, P(y): This component, which we covered in Chapter 4, answers the question: how likely is y to be a sentence in the target language? It models the *fluency* of the translation, ensuring the output is grammatical and natural. For example, P(the cat sat) would be much higher than P(cat the sat).

SMT's power lies in this separation. It searches for a translation that is not only a plausible rendering of the source text (a high translation model probability) but is also a well-formed sentence in the target language (a high language model probability). The system must balance these two competing objectives to find the optimal output.

Before an SMT system can learn that the French phrase $ne \dots pas$ is a good translation for the English phrase $do \ not$, it must first learn a more fundamental correspondence: that the word pas often appears in French sentences when the word not appears in their English counterparts. This foundational task of identifying word-level correspondences between sentences in a parallel corpus is known as **word alignment**. It is the essential first step in learning a statistical translation model. The goal is to produce a mapping that indicates which word or words in the target sentence are translations of each word in the source sentence.

This mapping is rarely a simple one-to-one correspondence between words in sequence. Languages have different grammatical structures, word orders, and ways of expressing concepts. As illustrated in **Fig. 9.4**, word alignments can exhibit several patterns:

• One-to-one: A single source word maps to a single target word (e.g., $blue \rightarrow bleu$). This is the most straightforward type of alignment.

- One-to-many: A single source word maps to multiple target words. This often occurs when a concept encapsulated in one word in the source language requires a multi-word phrase in the target language (e.g., a German verb like *mitgehen* might align to *qo with* in English).
- Many-to-one: Multiple source words map to a single target word (e.g., the English phrase floppy disk aligns to the single French word disquette).
- Many-to-many: Multiple source words map to multiple target words, a common pattern for idioms or complex phrases.

Furthermore, some words in one language may have no direct equivalent in the other. For instance, English requires an auxiliary verb in 'I do not agree,' where the word do has a purely grammatical function. In its French translation, 'Je ne suis pas d'accord,' the word do has no corresponding word. We model this by allowing source words to align to a special NULL token in the target sentence, effectively marking them as unaligned. The number of target words a source word aligns to is sometimes called its fertility.

Word alignments are not present in our parallel corpora; we only have the source and target sentences. The alignments are therefore *latent variables* that must be inferred. The core challenge is a circular one: if we had correct word alignments, we could easily estimate the probability of a word being the translation of another. Conversely, if we knew these translation probabilities, we could determine the most likely alignments. Breaking this cycle is the primary goal of the first SMT models. They are designed to learn these alignments automatically from nothing more than thousands of sentence pairs, forming the statistical bedrock upon which the entire translation model is built.

The foundational work on learning word alignments automatically was conducted at IBM in the early 1990s. The result was a series of five increasingly complex statistical models, known collectively as the *IBM Models*, which provided a principled, data-driven method for this crucial sub-task. These models treat the alignment as a hidden variable; while we observe the source and target sentences in a parallel corpus, the true word-level correspondences are unknown and must be inferred.

The core of the IBM Models is a generative story that probabilistically describes how a target sentence T and its alignment A could be generated from a source sentence S. The objective is to learn model parameters that maximize the probability P(T,A|S). The simplest of these, IBM Model 1, imagines a straightforward, if linguistically naive, process: for each position in the target sentence, it first chooses a source word to align with and then generates the target word based on a learned translation probability, p(t|s).

The key challenge is that we know neither the alignments nor the translation probabilities. The IBM Models solve this chicken-and-egg problem using the Expectation-Maximization (EM) algorithm. The process iterates between two steps:

- E-Step (Expectation): Using the current estimates of translation probabilities, calculate the expected alignment counts for all word pairs in the corpus. Essentially, we 'softly' guess the alignments.
- M-Step (Maximization): Use these expected counts to re-estimate and improve the translation probabilities. We update our model based on the guesses from the E-step.

This iterative process, which might start with uniform probabilities, converges towards a robust set of lexical translation probabilities. The subsequent models (Model 2 through 5) build upon this foundation, adding more sophisticated parameters to account for linguistic phenomena like word reordering (distortion) and the number of words a source

word translates to (fertility). These models were revolutionary, providing the statistical bedrock upon which phrase-based SMT was built.

While the IBM Models provided a powerful mathematical foundation for learning word-to-word correspondences, translating text one word at a time is fundamentally flawed. The limitations of this approach quickly become apparent. Languages have different word orders, and meaning is often conveyed by groups of words, not single lexical items. A word-for-word translation of the English phrase the green house into Spanish might incorrectly yield el verde casa, when the correct translation is la casa verde, demonstrating a change in both noun-adjective order and gender agreement. Furthermore, idiomatic expressions are impossible to translate literally; the English idiom he kicked the bucket has no meaningful word-for-word equivalent in most other languages.

The solution to these challenges was the leap from word-based to *phrase-based* statistical machine translation. This paradigm shift is centered on a simple but powerful idea: the fundamental unit of translation should not be a single word, but a contiguous sequence of words, or a 'phrase.' Instead of translating *green* and then *house*, the system learns to translate the block *green house* directly into *casa verde*. This seemingly simple shift allows the model to inherently learn local reordering, collocations, and short idiomatic expressions directly from the data.

The core of a phrase-based SMT system is a **phrase table**, which stores a massive collection of source phrases, their corresponding target phrase translations, and the probabilities associated with each pairing. These phrase pairs are not defined by linguistic rules but are extracted automatically from the parallel corpus using the word alignments generated by the IBM Models. The guiding heuristic is simple: any group of source words that is consistently aligned *only* with a specific group of target words across the corpus is a candidate phrase pair.

For instance, from the alignment between the green house and la casa verde, the system could extract the following valid phrase pairs:

- (the, la)
- (house, casa)
- (green house, casa verde)
- (the green house, la casa verde)

Each of these pairs is assigned a translation probability, such as p(casa verde|green house), estimated from its frequency of co-occurrence in the aligned corpus. By learning to translate multi-word chunks, phrase-based SMT achieved a dramatic improvement in translation fluency and adequacy over its word-based predecessors. The model no longer had to piece together translations from individual words; it could now use larger, more meaningful building blocks. With a massive phrase table constructed, the task of translation transforms into a search problem: finding the optimal way to segment the source sentence into phrases and combine their most probable translations to form a coherent and fluent target sentence.

Once we have a phrase table and a language model, the central task is to combine them to translate a new source sentence, \mathbf{f} . The challenge is not just to find a translation, but the *best* one. Formally, we are searching for the target sentence, $\hat{\mathbf{e}}$, that maximizes the probability of the translation. This is expressed as finding the sentence that maximizes

¹It is crucial to note that in this context, the term 'phrase' does not refer to a formal linguistic constituent like a noun phrase or a verb phrase. It simply means any contiguous sequence of one or more words from a sentence.

a score combining the translation model (TM), which ensures faithfulness to the source, and the language model (LM), which ensures fluency in the target language:

$$\hat{\mathbf{e}} = \arg\max_{\mathbf{e}} \operatorname{score}_{TM}(\mathbf{f}, \mathbf{e}) \cdot \operatorname{score}_{LM}(\mathbf{e})$$

Finding this optimal translation is a monumental search problem. Consider the possibilities: a single source sentence can be segmented into phrases in many different ways. Each source phrase might have dozens of possible target phrase translations in our phrase table. Finally, the chosen target phrases can often be reordered to fit the grammatical structure of the target language. This combination of segmentation, translation, and reordering choices leads to a combinatorial explosion, creating an astronomically large space of candidate translations. Generating and scoring every single one to find the absolute best is computationally impossible for all but the shortest of sentences.

This search for the highest-scoring translation is known as *decoding*. To make this problem tractable, SMT systems rely on heuristic algorithms. The dominant approach is a form of beam search. Instead of exploring all paths, a beam search decoder builds the translation incrementally, typically from left to right. At each step, it considers all possible ways to extend its current set of partial translations, or *hypotheses*. It then scores these newly expanded hypotheses and prunes the list, keeping only a fixed number of the most promising ones—a collection known as the 'beam.' All other, lower-scoring hypotheses are permanently discarded. This process of extending and pruning continues until complete sentences are formed.

Beam search is a greedy approach; by discarding hypotheses at each step, it is not guaranteed to find the single best translation. However, it provides an essential trade-off between computational cost and translation quality. It effectively navigates the vast search space to find a very high-quality translation in a fraction of the time an exhaustive search would require.

To make the decoding process concrete, let's walk through a simplified case study of a phrase-based SMT system translating the English sentence, 'The small house is on the hill', into Spanish. The goal of the system is to find the Spanish sentence \hat{s} that maximizes a probabilistic score combining three components: a phrase translation model, a reordering (or distortion) model, and a target language model.

Step 1: SegmentationFirst, the system considers all possible ways to segment the source sentence into phrases that exist in its pre-learned phrase table. For instance, a few of the many possible segmentations are:

- [The small house] [is on] [the hill]
- [The] [small house] [is] [on the hill]
- [The small] [house] [is on the hill]

Step 2: Translation Options For each segmented phrase, the system looks up potential translations and their associated probabilities, $\phi(\text{target phrase}|\text{source phrase})$, from the phrase table. A sample of these options might include:

- The small house \rightarrow La casa pequeña (prob. 0.7), La pequeña casa (prob. 0.2)
- is on \rightarrow está en (prob: 0.9)
- the hill \rightarrow la colina (prob: 0.8), el cerro (prob: 0.1)
- small house \rightarrow casa pequeña (prob: 0.6)

Step 3: Decoding and ScoringThe decoder's job is to navigate the immense search space of these options to find the highest-scoring complete translation. It does this by building hypotheses incrementally. Let's compare two competing hypotheses for translating the first part of the sentence, 'The small house':

- Hypothesis A: This hypothesis is generated by translating the single phrase [The small house] to La pequeña casa. This translation is monotonic—it directly preserves the English word order. It would incur a very low reordering penalty. Its score would be based on the phrase translation probability (0.2) and the language model's probability for 'La pequeña casa'.
- **Hypothesis B:** This hypothesis is generated from the same source phrase [The small house] but uses the higher-probability translation La casa pequeña. This phrase pair implicitly captures the necessary reordering of the adjective.

The crucial step is scoring these alternatives with the Spanish language model. A strong language model, trained on vast amounts of Spanish text, would know that casa pequeña (house small) is far more common and natural than pequeña casa (small house). Therefore, the language model score for Hypothesis B, P_{LM} (La casa pequeña), would be significantly higher than for Hypothesis A.

Even if Hypothesis A had a lower reordering cost, the powerful signal from the language model makes Hypothesis B the clear winner for this part of the sentence. The decoder prunes away less likely hypotheses like A and continues building upon promising ones like B. By extending this process, it combines La casa pequeña with the best translations for the rest of the sentence (e.g., está en and la colina), ultimately yielding the final, fluent output: La casa pequeña está en la colina.

For over a decade, Statistical Machine Translation, particularly phrase-based systems, represented the state-of-the-art. These models were a significant leap over their rule-based predecessors and produced useful translations for many language pairs. However, their core architecture contained inherent weaknesses that ultimately limited their potential and motivated a paradigm shift toward a new, more integrated approach.

The primary limitations of SMT systems stemmed from their complexity and fragmented design. An SMT system was not a single, unified model but a complex pipeline of several independently optimized components, including:

- A translation model to map source phrases to target phrases.
- A language model to ensure the fluency of the target output.
- A reordering model to handle different word orders between languages.

Tuning these disparate parts to work together was a formidable task, and errors from one stage would inevitably propagate to the next. Furthermore, SMT systems required extensive and painstaking *feature engineering*. Human experts had to manually design hundreds or even thousands of features to guide the translation process, making the models brittle and difficult to adapt to new language pairs or domains.

Perhaps most critically, phrase-based SMT had a fundamentally local view of the translation task. By processing sentences chunk by chunk, it struggled to capture long-range dependencies and broader contextual nuances, often resulting in disfluent or grammatically awkward output. The desire for a single, powerful model that could learn to translate in an *end-to-end* fashion—reading an entire source sentence and generating a translation without separate components—led directly to the rise of Neural Machine Translation.

The paradigm shift to Neural Machine Translation (NMT) introduced a fundamentally different and more elegant architecture for tackling the translation task. Instead of a

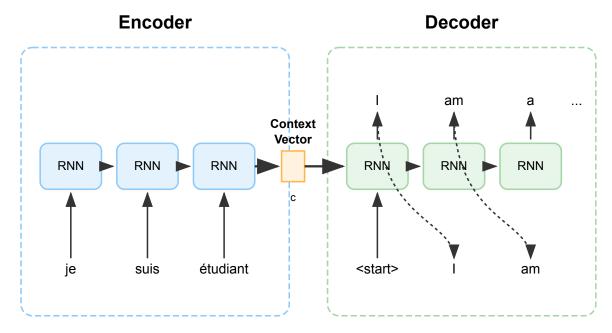


Figure 9.5: A high-level architectural diagram of a sequence-to-sequence (seq2seq) model, as described in Fig. 9.13. The Encoder RNN processes the input sentence step-by-step, compressing it into a fixed-size context vector. The Decoder RNN then uses this vector to generate the output sentence token-by-token in an auto-regressive manner.

complex pipeline of separately engineered components, as seen in SMT, NMT uses a single, large neural network that is trained *end-to-end*. The foundational architecture for this approach is the **sequence-to-sequence** (**seq2seq**) model, which is designed to map an input sequence of one length to an output sequence of a potentially different length. This model consists of two primary components: an *encoder* and a *decoder*.

The **encoder** is tasked with processing the source sentence and compressing its meaning into a dense, fixed-size numerical representation. As depicted in **Fig. 9.5**, the encoder is typically a Recurrent Neural Network (RNN), or one of its more powerful variants like a Long Short-Term Memory (LSTM) or Gated Recurrent Unit (GRU) network.² It reads the source sentence one word (or token) at a time, from left to right. At each step, the RNN updates its internal *hidden state*, which incorporates information about the current word and the context from all previous words. After processing the final word of the input sentence, the encoder's last hidden state serves as a summary of the entire sentence. This final vector is often called the **context vector** (or sometimes a 'thought vector'). The central idea is that this single vector encapsulates the complete meaning—the semantic essence—of the source sentence.

The **decoder** is the second major component, and its job is to take the context vector and generate the target sentence word by word. The decoder is also an RNN, and its initial hidden state is initialized with the context vector produced by the encoder. This 'primes' the decoder with the meaning of the sentence it needs to generate. The generation process then unfolds sequentially:

- 1. The decoder is first given a special start-of-sentence token (e.g., <start>).
- 2. Using the context vector and this initial token, it predicts the most probable first word of the target sentence.

²LSTMs and GRUs are advanced types of RNNs specifically designed to better handle long-range dependencies in sequences, mitigating issues like the vanishing gradient problem that can affect simpler RNNs.

- 3. The word it just generated is then fed back as the input for the next time step.
- 4. The decoder then predicts the second word, conditioned on both the context vector and the first word it generated.
- 5. This auto-regressive process continues, with each newly generated word becoming the input for the subsequent step, until the model predicts a special end-of-sentence token (e.g., <end>).

The entire seq2seq model is trained as a single system on a parallel corpus. The network's goal is to learn to maximize the conditional probability of producing the correct target sentence $Y = (y_1, ..., y_m)$ given a source sentence $X = (x_1, ..., x_n)$. Formally, it optimizes its parameters to maximize P(Y|X). During training, the model's output at each step is compared to the ground-truth word from the training data, and the error is calculated. This error is then propagated back through the entire network—from the decoder all the way back through the encoder—to adjust the model's millions of parameters. This end-to-end training allows the model to learn a complex, non-linear mapping from source text to target text directly, without any need for explicit word alignments or phrase tables.

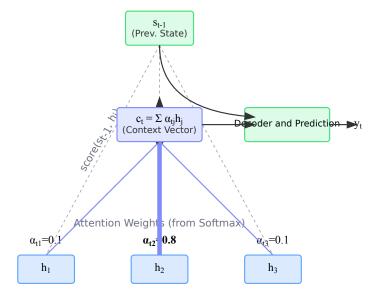
While the sequence-to-sequence architecture was a major breakthrough, its initial design harbored a critical weakness. The encoder is tasked with processing an input of arbitrary length and compressing its entire semantic content into a single, fixed-length vector. This vector, sometimes called the *context vector* or 'thought vector,' must encapsulate everything the decoder needs to know to generate a correct translation—the words, their meanings, their relationships, and their order.

This design imposes a severe **information bottleneck**. Consider the challenge of summarizing a long, complex sentence versus a short, simple one. A single vector of a predefined size has a finite capacity to store information. As the length and complexity of the source sentence increase, the model's ability to cram all the necessary nuance into this single representation diminishes. Inevitably, information is lost.

This architectural flaw had a direct and measurable impact: the performance of these early NMT models degraded significantly as sentence length grew. The model might successfully translate the first part of a long sentence but 'forget' key details or context from the beginning by the time it starts generating the end of the target sentence. The decoder is working from a single, static summary and has no way to refer back to the original source words. This limitation made it clear that a more dynamic way of accessing the source information was needed, leading directly to the development of the attention mechanism.

The fixed-length context vector used in early encoder-decoder models, as we have seen, creates a significant *information bottleneck*. It forces the model to compress the entire meaning of a source sentence—regardless of its length or complexity—into a single vector of a few hundred dimensions. This is an immense burden. A long, nuanced sentence about political philosophy and a short, simple sentence like 'The cat sat' must both be squeezed into the same fixed-size representation. Unsurprisingly, the performance of these models degrades sharply as sentence length increases, as information from early in the sentence is often lost by the time the vector is fully constructed.

To overcome this critical limitation, a powerful and elegant solution was introduced: the **attention mechanism**. The central idea is intuitive yet revolutionary: instead of forcing the decoder to rely on a single, static context vector, we allow it to *attend* to different parts of the source sentence at each step of the decoding process. At the moment it generates the word 'la', the French translation of 'the', it might need to focus on the English word 'The'. When it later generates 'souris', it should focus on 'mouse'. The



Encoder Hidden States

Figure 9.6: The attention mechanism at a single decoder time step. The decoder's previous state, s(t-1), 'attends' to all encoder hidden states by calculating alignment scores. These scores are converted via softmax into attention weights (α) , which determine each hidden state's contribution to the dynamic context vector c(t). The varied thickness of the purple arrows represents these weights, showing the model focusing heavily on the second hidden state. This context vector is then used with the decoder state to predict the next word, y(t).

attention mechanism gives the model a way to learn and implement this dynamic focus automatically.

At its core, the attention mechanism provides the decoder with access to all of the encoder's hidden states $(h_1, h_2, ..., h_N)$ at every decoding time step, t. Instead of using a pre-computed summary of the source, the decoder actively decides which encoder states are most relevant for predicting the current target word, y_t . It then computes a context-specific vector based on this decision. This process, visualized in **Fig. 9.6**, can be broken down into a few key steps.

First, we need to score how well each encoder hidden state h_j aligns with the decoder's hidden state from the *previous* time step, s_{t-1} . This score quantifies the relevance of the j-th source word to the t-th target word we are about to generate. This is accomplished using an *alignment score function*, which can be as simple as a dot product or a small feed-forward neural network that takes both states as input. We can represent this as:

$$e_{tj} = \operatorname{score}(s_{t-1}, h_j)$$

This gives us a vector of scores, one for each source word position.

Second, these raw alignment scores are normalized into a probability distribution using the softmax function. The resulting values, α_{tj} , are called the *attention weights*. Each weight represents the amount of attention the decoder should place on the *j*-th source word when generating the *t*-th target word. Because they are the output of a softmax function, these weights conveniently sum to 1 over all source words.

$$\alpha_{tj} = \frac{\exp(e_{tj})}{\sum_{k=1}^{N} \exp(e_{tk})}$$

Third, a dynamic context vector, c_t , is computed as a weighted sum of all the encoder

hidden states. The weights used in this sum are precisely the attention weights, α_{tj} , we just calculated.

$$c_t = \sum_{j=1}^{N} \alpha_{tj} h_j$$

This context vector is a powerful construct. If a particular source word at position j is deemed highly relevant for the current decoding step, its attention weight α_{tj} will be high, and its corresponding hidden state h_j will dominate the context vector c_t . If another source word is irrelevant, its weight will be near zero, effectively silencing its contribution.

Finally, this tailored context vector c_t is combined with the decoder's current hidden state, typically by concatenation, and then used to predict the target word y_t . The entire mechanism, from scoring to prediction, is implemented with differentiable operations, meaning it can be trained end-to-end with backpropagation, just like the rest of the network.

The introduction of attention was a watershed moment for NMT. By relieving the encoder of the impossible task of perfect compression, it dramatically improved translation quality, especially for long sentences. Furthermore, it offered a rare and valuable glimpse into the model's inner workings. By visualizing the attention weights α_{tj} as a matrix, we can see which source words the model 'looked at' when generating each target word. This often reveals plausible, soft alignments between languages, making the model more interpretable and easier to debug. This ability to dynamically weigh input components proved so effective that it became the foundational concept for the next generation of neural architectures.

The abstract concept of the attention mechanism becomes much clearer when visualized. At each step of generating a target word, the decoder calculates a set of attention weights, one for each word in the source sentence. These weights, which are positive and sum to one, determine the influence each source word has on the generation of the current target word. We can arrange these weights into a matrix, where rows correspond to the generated target words and columns correspond to the source words.

Fig. 9.7 presents a heatmap of such an attention matrix for a sample translation from French to English. The source sentence is laid out along the x-axis, and the generated target sentence is along the y-axis. The intensity of the color in each cell (i, j) represents the attention weight α_{ij} , which indicates how much the model focused on the j-th source word when producing the i-th target word. A brighter cell signifies a higher weight, meaning a stronger focus.

Notice the strong diagonal pattern in the visualization. This is common for languages with similar word order, indicating that the model often aligns the *n*-th target word with the *n*-th source word. However, the true power of attention is revealed in the non-diagonal alignments, which show the model handling differences in grammar and word order. For instance, in translating the French phrase 'la voiture verte' to the English 'the green car', the model must reorder the adjective and noun. When generating the word 'green', the attention mechanism would allow it to place a high weight on the source word 'verte', regardless of their different positions in their respective sentences. This would be visible as a bright, off-diagonal cell in the matrix.

This ability to dynamically link target words to relevant source words is precisely what overcomes the fixed-context-vector bottleneck of earlier models. Instead of relying on a single, compressed summary of the entire input, the decoder has direct, selective access to the source representation at every step. This not only dramatically improves translation quality but also offers a valuable form of *interpretability*. By visualizing attention, we gain a direct insight into the model's internal decision-making process, observing how it aligns words across languages to construct a coherent translation.

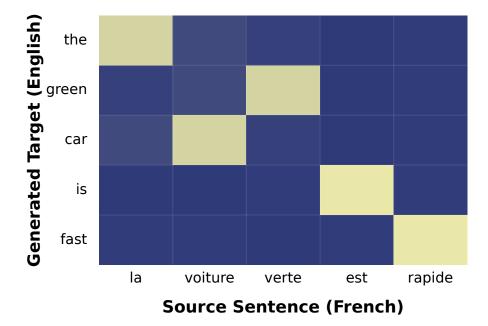


Figure 9.7: A heatmap of an attention matrix for a sample translation from French to English. The source sentence is on the x-axis and the generated target sentence is on the y-axis. The brightness of a cell (i, j) indicates the attention weight, showing how much the model focused on the j-th source word when generating the i-th target word. The strong off-diagonal alignments for 'green'/'verte' and 'car'/'voiture' demonstrate the model's ability to handle word reordering.

While the attention mechanism solved the information bottleneck of simple encoder-decoder models, the underlying recurrent neural networks (RNNs) still possessed a fundamental limitation: their sequential nature. Because the computation for each word depended on the hidden state of the previous word, the process was inherently difficult to parallelize, making it slow to train on very large datasets and long sentences. In 2017, a landmark paper titled 'Attention Is All You Need' introduced the *Transformer* architecture, a new model for NMT that dispensed with recurrence entirely.

The Transformer's core innovation is a mechanism called *self-attention*. Instead of passing information sequentially, self-attention allows every word in the source sentence to directly look at and weigh the importance of all other words in that same sentence. This process is repeated in multiple layers, allowing the model to build a deeply contextualized representation for each word by aggregating information from across the entire input. The Transformer maintains the familiar encoder-decoder structure, but each component is now a stack of these self-attention layers rather than a recurrent network.

This architectural shift had two profound consequences. First, by allowing direct connections between any two words, it became exceptionally good at modeling long-range dependencies. Second, and perhaps more importantly, the removal of recurrence made the computations massively parallelizable. The representation for every word could be calculated simultaneously, leading to dramatic reductions in training time on modern hardware like GPUs. The Transformer's superior performance and efficiency quickly established it as the new state-of-the-art for machine translation, and its architecture became the foundational blueprint for the large language models that would come to dominate the field, as we will see in the final chapter.

The theoretical advantages of Neural Machine Translation, especially the attention mechanism's capacity to model long-range dependencies, translate directly into a dramatic and observable leap in translation quality. While Statistical Machine Translation

| Source Sentence | SMT Output | NMT Output | |
|---|------------|--|--|
| The policy, which had been under discussion for months, was eventually shot down by the opposition. | , | The policy, which had been under discussion for months, was ultimately rejected by the opposition. | |

Figure 9.8: A comparative table showcasing the qualitative difference between Statistical Machine Translation (SMT) and Neural Machine Translation (NMT) outputs. The NMT model correctly interprets the idiomatic phrase 'shot down', whereas the SMT model produces a literal and disfluent translation.

often produces literal and sometimes disjointed translations by stitching together memorized phrases, NMT systems generate text that is significantly more fluent, coherent, and accurate.

This qualitative gap is starkly illustrated in the comparative example shown in **Fig. 9.8**. The English source sentence, 'The policy, which had been under discussion for months, was eventually shot down by the opposition,' presents challenges in its complex syntax and idiomatic phrasing ('shot down').

The SMT output demonstrates the typical weaknesses of the phrase-based approach. It struggles with the long-distance dependency between 'The policy' and 'was...shot down,' and it is likely to produce a clumsy, literal translation of the idiom—perhaps as 'was fired down' or 'was gunned down.' The result is a grammatically awkward sentence that fails to convey the correct political meaning and would require significant human post-editing to be usable.

In contrast, the NMT output reveals a much deeper, more holistic understanding of the source text. By processing the entire sentence contextually, the model correctly interprets 'shot down' as a metaphor for 'rejected' or 'defeated' and finds an appropriate, non-literal equivalent in the target language. It smoothly handles the embedded clause, maintaining the sentence's logical flow. This ability to move beyond local word and phrase substitutions to re-express sentence-level meaning is the hallmark of the NMT revolution. It marks the transition from a mechanical process of fragment replacement to a more nuanced act of linguistic regeneration.

The evolution from Statistical Machine Translation (SMT) to Neural Machine Translation (NMT) marks a profound shift in methodology and core philosophy. The SMT paradigm is one of decomposition. It dissects the translation process into distinct, manageable sub-problems: word alignment, phrase extraction, and decoding, each optimized separately. Its power comes from combining a translation model, which ensures faithfulness to the source, with a language model, which ensures fluency in the target. The system relies heavily on explicit, count-based statistical tables and carefully engineered features to guide the search for the best translation.

NMT, in contrast, champions a *holistic*, end-to-end approach. It replaces the complex pipeline of SMT with a single, large neural network trained to perform the entire translation task. This model learns to directly map a sequence of source words to a sequence of target words without relying on explicit phrase tables or intermediate steps. By representing language in continuous vector spaces and using mechanisms like attention, NMT models can capture long-range dependencies and subtle contextual nuances that were difficult for SMT to manage, leading to more fluent and human-like translations.

Despite the remarkable advances brought by Neural Machine Translation, the field is far from solved. Significant challenges remain, and exciting new research directions are constantly emerging. The data-hungry nature of NMT models creates a major bottleneck, leading to a stark digital divide between high-resource languages with abundant parallel corpora (like English and French) and the vast majority of the world's languages, which are considered *low-resource*. Current research actively tackles this problem through techniques like transfer learning, where a model trained on a high-resource pair is fine-tuned on a smaller low-resource dataset, and unsupervised methods that learn translations using only monolingual data.

Beyond data availability, several other key challenges persist:

- Domain Mismatch: A model trained on news articles will struggle to translate legal documents or social media posts. Adapting MT systems to specific domains and handling noisy, informal text remains a crucial area of work.
- Evaluation: As we have seen, automated metrics like BLEU are an imperfect proxy for human judgment. They can penalize valid translations that use different phrasing and often fail to capture subtle nuances of fluency and adequacy. The development of better evaluation metrics is an ongoing pursuit.
- Bias and Fairness: Like other large models trained on web data, MT systems can learn and amplify societal biases related to gender, race, and culture. Ensuring fair and unbiased translations is a critical ethical imperative.

Looking forward, the future of machine translation is likely to be more contextual and interactive. One of the most promising frontiers is multimodal machine translation, which aims to incorporate information from other modalities, such as images or audio, into the translation process. For example, an accompanying image could help a model correctly disambiguate the word 'bat' (animal vs. sports equipment) when translating a sentence. Another exciting direction is developing controllable translation systems, allowing users to specify constraints such as formality, tone, or the inclusion of specific terminology. This would transform MT from a static tool into a dynamic, collaborative assistant, particularly for professional translators. The quest for perfect translation continues, driving innovation at the heart of computational linguistics.

Chapter 10

Information Retrieval and Information Extraction



This chapter focuses on two critical tasks for managing and extracting value from the vast amounts of text available today: *Information Retrieval* (IR) and *Information Extraction* (IE). While both disciplines aim to satisfy a user's need for information, their goals and outputs are fundamentally different.

Information Retrieval is the task of finding a subset of documents from a large collection that are relevant to a user's query. This is the classic problem solved by search engines. When you search for a topic, an IR system sifts through millions of documents to return a ranked list of the ones most likely to contain the information you seek. The core output of an IR system is a set of *relevant documents*.

In contrast, Information Extraction operates at a more granular level. Its goal is not to retrieve entire documents, but to automatically identify and pull out specific, pre-defined types of information from within the text. Instead of returning a list of articles, an IE system outputs *structured data*, such as filling a database table with names, locations, or events mentioned in the text. This chapter will delve into the foundational models for both of these essential fields.

To grasp the fundamental difference between these two fields, consider a practical scenario. Imagine you are a financial analyst tasked with tracking the performance of a fictional company, 'InnovateCorp,' during the fourth quarter of 2023.

Your first step might involve Information Retrieval. You would use a search engine with a query like InnovateCorp Q4 2023 earnings. The IR system's goal is to locate documents relevant to this query from a vast collection. It would return a ranked list of results: perhaps the company's official press release, a news article from a financial publication, and a detailed regulatory filing. The system has successfully retrieved potentially relevant sources, but your work is not done. You must now manually read through these documents to locate the specific figures you need. In short, IR helps you find the right haystack.

This is where **Information Extraction** provides a more direct solution. An IE system is not concerned with returning entire documents. Instead, it is designed to parse the text *within* those documents to identify and pull out specific, predefined types of information. Its goal is to turn unstructured prose into structured data. Given the same set of documents, an IE system would aim to produce a structured output, like a table or a database entry:

• Company: InnovateCorp

• Fiscal Period: Q4 2023

• Revenue: \$15.2 Billion

• Net Income: \$2.1 Billion

• Earnings Per Share: \$1.35

In this contrast, the core distinction becomes clear. IR systems return *documents* that are relevant to a query. IE systems return *facts* extracted from text. IR helps you find *where* the information is, while IE aims to deliver the information itself, structured and ready for analysis.

At its core, the Information Retrieval (IR) problem is about satisfying a user's information need. We can formally define the classic IR task using three main components. First, there is a static **document collection**, which is a set of texts $D = \{d_1, d_2, ..., d_N\}$ through which we want to search. Second, a user provides a **query**, q, which is a representation of their information need, typically expressed as a short string of keywords.

The objective of an IR system is to process the query q against the collection D and return a **ranked list of documents**. This list is ordered by relevance, with the documents

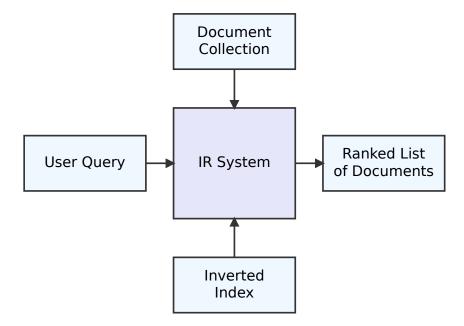


Figure 10.1: High-level architecture of a classic Information Retrieval (IR) system. A user query is processed by the system, which utilizes a document collection and an inverted index to produce a ranked list of relevant documents.

that the system predicts are most likely to satisfy the user's need placed at the top. The high-level architecture of this process, from query input to ranked output, is visualized in Fig. 10.1. The central challenge, therefore, is to accurately and efficiently model this notion of relevance. The upcoming sections will explore the fundamental data structures and algorithms used to accomplish this.

How can a system efficiently find documents containing a query term from a collection of millions or even billions of texts? A brute-force approach, which involves scanning every document for the term at query time, is computationally infeasible. The solution to this problem is a data structure at the heart of nearly all modern search engines: the **inverted index**. The core idea is simple and elegant, analogous to the index at the back of a physical book. Instead of mapping from a document to the terms it contains, an inverted index maps from a term to the documents that contain it.

An inverted index consists of two main components:

- Dictionary (or Vocabulary): This is a list of all unique terms that appear in the entire document collection. For fast lookups, the dictionary is typically sorted or stored in a complex data structure like a B-tree.
- Postings Lists: For every term in the dictionary, there is an associated list of pointers to the documents that contain that term. Each entry in this list, called a *posting*, is typically a document identifier (DocID).

When a single-word query, such as 'computational', arrives, the system performs a fast lookup in the dictionary to find the term. It then retrieves the corresponding postings list, which immediately provides the set of all relevant documents.

The real power of the inverted index becomes apparent with multi-word queries, such as "computational" AND "linguistics". The system retrieves the postings list for 'computational' and the postings list for 'linguistics'. It then computes the *intersection* of these two lists—the set of DocIDs that appear in both. This operation is vastly more efficient than scanning every document for both terms. The inverted index trades a slower,

| I. Document Collection | | | | | |
|------------------------------|---|--|--|--|--|
| Document ID | Content | | | | |
| 1 | Company growth report. | | | | |
| 2 | Earnings report shows revenue growth. | | | | |
| 3 | New company strategy. | | | | |
| 4 | Quarterly revenue exceeds expectations. | | | | |
| II. Resulting Inverted Index | | | | | |
| Token | Postings List | | | | |
| company | [1, 3] | | | | |
| earnings | [2] | | | | |
| exceeds | [4] | | | | |
| expectations | [4] | | | | |
| growth | [1, 2] | | | | |
| new | [3] | | | | |
| quarterly | [4] | | | | |
| report | [1, 2] | | | | |
| revenue | [2, 4] | | | | |
| shows | [2] | | | | |
| strategy | [3] | | | | |

Figure 10.2: Building an inverted index. The top part shows the source document collection. The bottom part shows the resulting index, mapping each unique term to its postings list (the documents in which it appears). The example for the term 'revenue' with postings list [2, 4] matches the text's description.

one-time indexing process for extremely fast query resolution, making it the foundational component for modern Information Retrieval.

To illustrate how an inverted index is built, let's consider the small collection of four documents shown in the top half of **Fig. 10.2**. The process begins by collecting all unique, normalized tokens from the entire collection to form the dictionary. For our example, this dictionary would include terms like company, earnings, growth, report, revenue, and so on.

Next, we iterate through each term in the dictionary and compile its *postings list*—the list of all document IDs containing that term. To find the postings for the term 'revenue', we scan the collection and note that it appears in Document 2 and Document 4. Therefore, its postings list is [2, 4]. This process is repeated for every unique term.

The complete inverted index is shown in the bottom half of **Fig. 10.2**. It consists of the dictionary of terms, where each term points to its corresponding postings list. Once built, this structure allows for extremely fast retrieval. Finding documents that mention 'revenue' is now a simple lookup operation, not an expensive scan of every document.

While an inverted index efficiently finds documents containing a query term, it doesn't inherently rank them by relevance. To do this, we need a way to quantify the thematic content of a document. The **vector space model (VSM)** provides an elegant algebraic framework for this task by representing documents and queries as numerical vectors in a high-dimensional space.

In this model, each unique term in the entire document collection corresponds to a dimension. A document is then represented as a vector, \vec{d} , where each component of the vector reflects the importance of the corresponding term within that document. For example, a document about corporate acquisitions would have large vector components in the dimensions for terms like 'acquisition', 'merger', and 'shares'.

The fundamental insight of the VSM is that proximity in this vector space equates to

thematic similarity. Documents whose vectors point in a similar direction are considered semantically related. By also casting the user's query as a vector, \vec{q} , we can transform relevance ranking into a geometric problem: finding the document vectors that are closest to the query vector.

To represent documents and queries as numerical vectors, we use a foundational data structure called the **term-document matrix**. Conceptually, this is a large table where each row corresponds to a unique term in the entire collection's vocabulary, and each column corresponds to a single document. The value in a cell at the intersection of a term t and a document d quantifies the importance of that term within that document.

In its most basic form, this value is simply the raw frequency of the term. For a vocabulary of size |V| and a collection of |D| documents, we have a matrix M of size $|V| \times |D|$. The entry M_{td} would be the count of term t in document d. Consequently, each column of this matrix is the vector representation for a specific document, existing in a |V|-dimensional space. For instance, if 'algorithm' is the 50th term in our vocabulary and it appears 4 times in Document 12, the value at $M_{50,12}$ would be 4. While intuitive, these raw counts are just a starting point; a more sophisticated weighting is needed for an effective system.

Not all terms in the term-document matrix carry equal weight. A word like 'the' might appear frequently in every document but provides little information about a document's specific topic. Conversely, a technical term that appears many times in one document but is absent from most others is likely a strong indicator of its content. The **TF-IDF** (Term Frequency-Inverse Document Frequency) weighting scheme is designed to formalize this intuition, assigning a numerical score to each term that reflects its importance to a document within a collection.

TF-IDF is composed of two distinct parts:

• Term Frequency (TF): This measures how often a term t appears in a document d. A higher frequency implies greater importance within that specific document. To prevent very frequent terms from overly dominating the score, we often use a logarithmically scaled frequency. The formula is:

$$tf_{t,d} = \log(1 + f_{t,d})$$

where $f_{t,d}$ is the raw count of term t in document d.

• Inverse Document Frequency (IDF): This measures how common or rare a term is across the entire document collection. The intuition is that a term appearing in many documents is less informative than one appearing in only a few. The IDF score is high for rare terms and low for common ones. It is calculated as:

$$idf_t = \log \frac{N}{df_t}$$

where N^* is the total number of documents in the collection, and df_t is the number of documents that contain the term t (the document frequency).

The final TF-IDF weight for a term t in a document d is the product of these two components:

$$w_{t,d} = t f_{t,d} \times i d f_t$$

The effect of this calculation is powerful. A term receives a high TF-IDF score if it appears frequently in a document (high TF) but is rare in the overall collection (high IDF). These are precisely the terms that best characterize a document's unique topic.

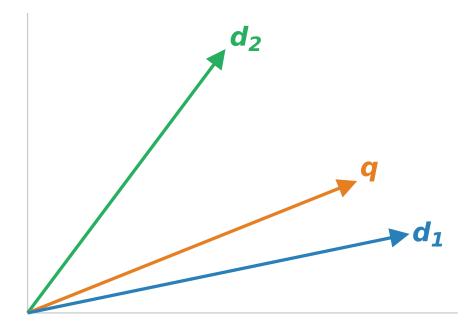


Figure 10.3: A 2D vector space diagram to intuitively explain cosine similarity. The diagram shows a query vector ('q') and two document vectors ('d1' and 'd2') originating from the same point. The angle between 'q' and 'd1' is visibly smaller than the angle between 'q' and 'd2', illustrating that 'd1' is considered more relevant to the query.

By replacing the raw counts in the term-document matrix with these TF-IDF scores, we create much more meaningful vectors for use in the vector space model, leading to more accurate relevance rankings.

With both documents and the query represented as TF-IDF vectors, the task of ranking becomes a geometric problem: measuring the proximity between the query vector \mathbf{q} and each document vector \mathbf{d} . In the vector space model, the standard measure is *cosine similarity*. This metric is effective because it measures the angle between two vectors, making it sensitive to the orientation (the topics) rather than the magnitude (the length of the documents).

The similarity is calculated as the cosine of the angle θ between the two vectors:

$$\mathrm{similarity}(\mathbf{q},\mathbf{d}) = \cos(\theta) = \frac{\mathbf{q} \cdot \mathbf{d}}{||\mathbf{q}|| \; ||\mathbf{d}||}$$

The resulting score ranges from 0 to 1 for non-negative TF-IDF vectors. A score closer to 1 signifies a smaller angle and thus higher conceptual similarity. As illustrated in Fig. 10.3, the document vector **d1** is deemed more relevant to the query **q** than **d2** because the angle between them is smaller. An IR system computes this similarity score for every document in the collection and then sorts the documents in descending order to generate the final ranked list of results.

Once an IR system returns a list of documents, how do we measure its success? The quality of a search result is not a simple binary judgment. To formally evaluate performance, we rely on metrics that capture different aspects of relevance. The two most fundamental of these are *precision* and *recall*.

Precision answers the question: Of the documents the system retrieved, what fraction were actually relevant? It is a measure of exactness or fidelity. If every document you see is on-topic, the system has high precision. Recall, on the other hand, answers: Of all the relevant documents that exist in the collection, what fraction did the system find? It is a measure of completeness. If the system found every possible relevant document, it has high recall.

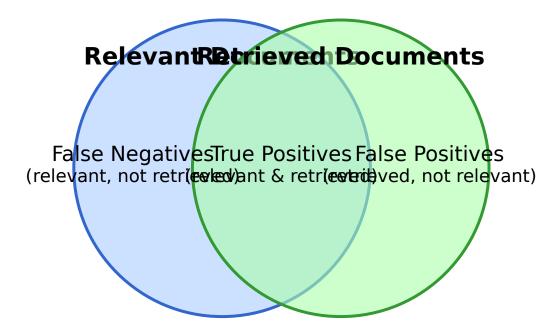


Figure 10.4: A Venn diagram illustrating the concepts of precision and recall by showing the overlap between the set of relevant documents and the set of retrieved documents.

These concepts are often visualized using the sets of retrieved and relevant documents, as shown in **Fig. 10.4**. To formalize this, we can categorize any document as a *True Positive* (relevant and retrieved), a *False Positive* (retrieved but not relevant), or a *False Negative* (relevant but not retrieved). The formulas are then:

- Precision = $\frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$
- Recall = $\frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$

There is an inherent trade-off between these two metrics. A system designed for maximum recall could simply return every document in the collection, guaranteeing it finds all relevant ones (100% recall) but with abysmal precision. Conversely, a cautious system that returns only a single document it is extremely confident about might achieve perfect precision but have terrible recall. Balancing this trade-off is a central challenge in designing IR systems.

Often, there is an inverse relationship between precision and recall; improving one can hurt the other. This trade-off makes it difficult to compare systems using two separate numbers. To address this, the F1-score provides a single metric that combines them. The F1-score is the harmonic mean of precision and recall:

$$F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

Unlike a simple average, the harmonic mean penalizes extreme values more heavily. This means a system must achieve both reasonably high precision *and* high recall to get a high F1-score. For instance, a system with perfect precision (1.0) but very low recall (0.1) would have a low F1-score of only about 0.18. This property makes it one of the most common and robust evaluation metrics in information retrieval.

While precision and recall evaluate an unordered set of documents, most IR systems produce a ranked list. A good system must place relevant documents at the top. Metrics designed for this purpose are crucial, with Mean Average Precision (MAP) being a standard measure.

MAP is calculated in two stages. First, for a single query, we compute the *Average Precision (AP)*. This is the average of the precision values at each point in the ranked list where a relevant document is found. It rewards systems for ranking relevant documents higher. The formula is:

$$AP = \frac{\sum_{k=1}^{n} (P(k) \times \text{rel}(k))}{\text{number of relevant documents}}$$

Here, P(k) is the precision at rank k, and rel(k) is 1 if the document at rank k is relevant. Second, the MAP is simply the mean of these AP scores calculated over a large set of queries.

Let's consolidate the theoretical concepts of Information Retrieval by building a miniature search engine. Our goal is to rank a small collection of news headlines based on their relevance to a user query.

Imagine our document collection consists of four headlines:

- **D1**: 'The Federal Reserve hinted at another interest rate hike to combat inflation.'
- **D2:** 'Global stock markets reacted nervously to the Federal Reserve's policy statement.'
- D3: 'A new tech stock surged today after a successful product launch.'
- **D4:** 'Market analysts predict a stock hike following the interest rate decision.'

Our first step is pre-processing. We tokenize the text, convert it to lowercase, and remove common stop words ('the', 'at', 'to', 'a', etc.). This process yields a vocabulary of unique terms. From this, we construct our inverted index, which maps each term to the documents containing it:

• federal: D1, D2

• reserve: D1, D2

• interest: D1, D4

• rate: D1, D4

• hike: D1, D4

• stock: D2, D3, D4

• market: D2, D4

• ...and so on for every term in the vocabulary.

Next, we transform these documents into numerical vectors using the TF-IDF weighting scheme. We create a term-document matrix where each cell will hold the TF-IDF score for a term in a specific document. For example, let's calculate the weight for the term federal in D1. The Term Frequency (TF) is high, as it appears once in a short document. The Inverse Document Frequency (IDF) is calculated as $log(\frac{N}{n_t})$, where N is the total number of documents (4) and n_t is the number of documents containing the term federal (2). Thus, the IDF is $log(\frac{4}{2}) \approx 0.301$. The final TF-IDF score combines these two values. In contrast, the term stock has a lower IDF of $log(\frac{4}{3}) \approx 0.125$ because it is more common in our collection. This process is repeated for every term in every document, resulting in four vectors, $v_{D1}, v_{D2}, v_{D3}, v_{D4}$, in a high-dimensional space where each dimension corresponds to a term in our vocabulary.

Now, a user submits the query: 'Federal Reserve stock'. We apply the same preprocessing and TF-IDF vectorization process to the query itself, creating a query vector, v_q . The final step is to measure the relevance between the query and each document. We do this by calculating the cosine similarity, which measures the angle between the query vector and each document vector. The formula is:

$$\text{similarity}(q,d) = \cos(\theta) = \frac{v_q \cdot v_d}{\|v_q\| \|v_d\|}$$

A higher score (closer to 1) indicates a smaller angle and thus greater similarity. After performing these calculations, we might get the following hypothetical scores:

• similarity(q, D2): 0.91

• similarity(q, D1): 0.78

• similarity(q, D4): 0.45

• similarity(q, D3): 0.31

Based on these scores, our search engine would return the ranked list of documents: [D2, D1, D4, D3]. This simple case study demonstrates the complete IR pipeline, from processing raw text with an inverted index to representing documents in a vector space and ranking them by relevance to a query.

While Information Retrieval systems excel at finding relevant documents from a large collection, they typically stop there, leaving the user to read the document to find the specific facts they need. We now shift our focus to a different, though related, task: Information Extraction (IE).

The goal of Information Extraction is to automatically identify and pull structured information from unstructured or semi-structured text. Instead of returning an entire document, an IE system returns specific pieces of data formatted for direct use in a database or spreadsheet. Where IR might return a dozen news articles about corporate mergers, an IE system would process those articles to extract the specific names of the acquiring company, the company being acquired, and the monetary value of the deal. This process transforms a sea of free-form text into a structured, queryable knowledge base, turning raw prose into actionable data. The subsequent sections will detail the core subtasks required to build such a system.

To make the goal of Information Extraction more concrete, consider the difference between finding a document and understanding its contents. While an IR system would successfully retrieve a news article about a corporate acquisition, an IE system would process that same article to pull out the key facts. Imagine the system encounters the following text in a press release: 'Global Tech Inc. announced today its definitive agreement to acquire Innovate Solutions for a staggering \$2.5 billion.'

The objective of IE is to automatically populate a structured record, like a database table, from this unstructured sentence. As shown in **Fig. 10.5**, the system must identify specific text fragments and map them to predefined slots. It needs to recognize that 'Global Tech Inc.' is the *Acquirer*, 'Innovate Solutions' is the *Target Company*, and '\$2.5 billion' is the *Acquisition Price*. This process effectively transforms a block of prose into a row of data that can be easily queried, aggregated, and analyzed. The core challenge lies in building models that can perform this mapping accurately across countless variations in phrasing and sentence structure.

With the overall goal of Information Extraction established, we begin with its most fundamental subtask: **Named Entity Recognition (NER)**. The objective of NER is to locate and classify named entities—mentions of real-world objects—into pre-defined

Unstructured Text

Structured Output

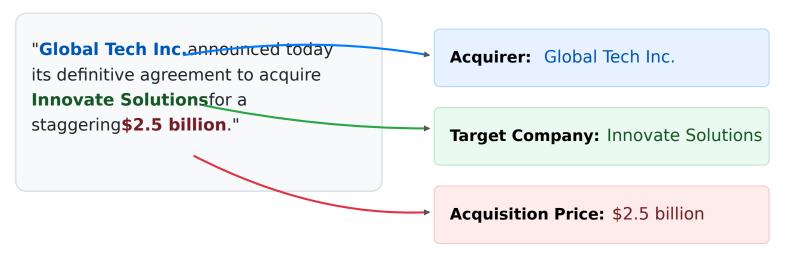


Figure 10.5: Illustration of Information Extraction (IE). An IE system processes an unstructured sentence about a corporate acquisition (left) to identify and extract key pieces of information, populating them into a structured data record (right). Arrows indicate the mapping from text fragments to their corresponding data slots.

categories. These entities are typically rigid designators, meaning they refer to a specific, unique thing. While the set of categories can be tailored to any domain, a common starting point includes:

- **PER**: Person (e.g., Ada Lovelace, the CEO)
- ORG: Organization (e.g., Google, The United Nations)
- LOC: Location (e.g., Silicon Valley, Mount Everest)
- **GPE**: Geopolitical Entity (e.g., France, Tokyo)

Essentially, NER is a classification task performed on spans of text. For instance, given the sentence, 'In 2014, Amazon acquired Twitch,' an NER system would identify *Amazon* as an ORG and *Twitch* as another ORG. By identifying these key players, NER provides the foundational building blocks—the nouns of an information database—upon which more complex tasks, such as Relation Extraction, can be built. This process is the first crucial step in transforming a sea of unstructured text into structured, queryable knowledge.

Before the widespread adoption of statistical methods, early approaches to NER relied on human expertise and manually curated resources. Two of the most prominent techniques from this era were the use of *gazetteers* and *hand-crafted rules*.

A gazetteer is essentially a large dictionary or pre-defined list of specific entity names. For instance, an NER system could be equipped with separate gazetteers containing:

- All countries and major cities in the world (for LOCATIONS).
- A list of well-known corporations (for ORGANIZATIONS).
- Common first names and surnames (for PERSONS).

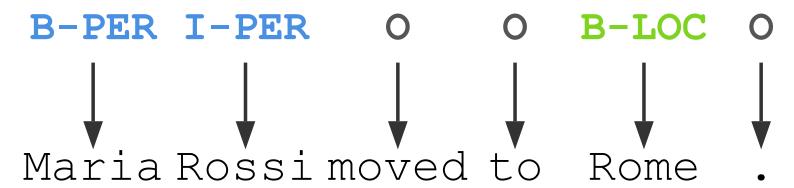


Figure 10.6: Named Entity Recognition (NER) framed as a sequence labeling task using the BIO encoding scheme. Each token in the input sentence is assigned a tag indicating whether it is at the beginning (B-TYPE), inside (I-TYPE), or outside (O) of a named entity.

The system performs a simple lookup, matching sequences of words in the text against entries in these lists. If a match is found, the corresponding entity type is assigned. The primary limitation of this approach is coverage; a gazetteer is only as good as its contents and will fail to identify any new or unlisted entities.

To capture more general patterns, developers wrote hand-crafted rules, often using the power of regular expressions (as discussed in Chapter 2). These rules encode orthographic and contextual clues. A simple rule might state that any capitalized word followed immediately by 'Inc.' or 'Corp.' should be tagged as an ORGANIZATION. Another rule could identify a title like 'Mr.' or 'Dr.' followed by a capitalized name as a PERSON. While effective in specific, well-understood domains, rule-based systems are brittle, labor-intensive to create, and difficult to maintain or adapt to new types of text.

While rule-based systems are intuitive, they are often brittle and require significant manual effort to create and maintain. Modern approaches have largely shifted towards statistical models by framing NER as a *sequence labeling* task. This paradigm should be familiar, as it is the same fundamental approach we used for Part-of-Speech tagging in Chapter 5. The goal is to assign a label from a predefined set of tags to each token in an input sequence.

However, NER presents a challenge not found in POS tagging: entities can span multiple tokens. For instance, in 'Maria Rossi,' we need to know that both words refer to a single entity. A simple tag set like {PER, LOC, ORG} is insufficient. To address this, we use a more expressive tagging scheme, the most common of which is the **BIO** encoding:

- B-TYPE: Marks the *Beginning* of an entity of a certain type.
- I-TYPE: Marks the *Inside* of an entity, used for the second and subsequent tokens.
- O: Marks a token as being Outside any named entity.

This scheme allows a model to explicitly mark the boundaries of multi-word entities. As shown in Fig. 10.6, the sentence 'Maria Rossi moved to Rome.' is transformed into a sequence of token-tag pairs. The token Maria is labeled B-PER as it begins a person entity, while Rossi is labeled I-PER, indicating it is inside the same person entity. The single-word entity Rome is tagged with B-LOC¹, and the remaining tokens are tagged with D

¹Note that for single-token entities, the B- tag is used by convention. There is no corresponding I-tag.

By converting raw text into this labeled sequence format, we can leverage powerful machine learning algorithms. Models like Hidden Markov Models (HMMs), Conditional Random Fields (CRFs), and modern neural networks can be trained on annotated corpora to learn how to predict the most probable sequence of BIO tags for a new, unseen sentence. These models make their predictions by learning from the features of the words themselves and the context in which they appear.

For a statistical model to perform Named Entity Recognition, it cannot rely on the word itself; it needs descriptive clues, or *features*, about the word and its environment. These features are the evidence the model uses to make its classification decisions. A rich set of features is crucial for building an accurate NER system.

Common features can be grouped into several categories:

- Word-Level Features: These capture information internal to the word token.
 - Word Shape: Abstracts away from the specific word to its orthographic pattern.
 For instance, Apple might have the shape Xxxxx, IBM would be XXXX, and 12/03/2024 could be dd/dd/dddd. This helps the model generalize to unseen words that follow common capitalization or formatting patterns for entities.
 - Affixes: Prefixes and suffixes of a word (e.g., -berg, -ton, -shire) can be strong indicators of names or locations.
 - Part-of-Speech Tag: The word's grammatical category is a powerful feature. A
 proper noun tag (NNP in the Penn Treebank tagset) is a very strong signal for
 an entity.
- Contextual Features: These look at the words surrounding the target token.
 - Surrounding Words: The identity of the previous and next few words (e.g., the two words before and two words after) provides immediate context. The word Washington is more likely a person if preceded by Mr. and a location if preceded by in.
 - Surrounding POS Tags: Similarly, the POS tags of neighboring words are also highly informative.
- Gazetteer Features: This involves using pre-existing lists of entities, known as gazetteers. A simple binary feature can be created: 'Is this word present in a list of known U.S. cities?' This allows the model to incorporate external world knowledge directly.

Identifying named entities is a crucial first step, but it often leaves us with a collection of disconnected labels. To truly understand a text, we must uncover how these entities relate to one another. This is the goal of **Relation Extraction (RE)**: the task of identifying and categorizing the specific semantic relationships that exist between the entities found by an NER system. The objective is to convert unstructured text into structured, relational data.

Typically, these relationships are represented as tuples, often of the form (entity1, relation_type, entity2). For instance, from the sentence 'Tesla, a company based in Austin, announced a new factory,' an RE system should be able to identify the LocatedIn (Tesla, Austin) relation. Other examples of target relations could include WorksFor (Person, Organization) or Acquired(Company, Company). By transforming unstructured sentences into a set of such relational facts, RE systems build a scaffold of knowledge that can be used to populate databases, construct knowledge graphs, or answer complex questions about the information contained within a text corpus.

Once entities have been identified, the simplest way to find relations between them is to search for recurring patterns in the text. These pattern-based approaches, while straightforward, can be surprisingly effective, especially in domains with regular, conventional language like news reports.

The most direct method involves applying regular expressions over the raw text string. After an NER system has tagged entities, we can search for specific lexical patterns between them. For instance, a pattern like (ORGANIZATION) acquired (ORGANIZATION) is a strong indicator of an acquisition relationship. Similarly, (PERSON), CEO of (ORGANIZATION) clearly signals an is-CEO-of relation. The main weakness of this approach is its brittleness; a slight variation in phrasing, such as 'the acquisition of Y by X,' would fail to match the initial pattern.

A more robust technique uses patterns over syntactic structures rather than linear text. By first parsing a sentence to get its dependency graph (as discussed in Chapter 6), we can identify relations based on the dependency path between two entities. For the sentence 'Steve Jobs founded Apple in Cupertino,' the dependency path between Jobs and Apple might be represented as:

Jobs <-nsubj- founded -dobj-> Apple

This pattern, which captures the grammatical subject-verb-object relationship, is far more general. It can successfully identify the *founded_by* relation in various syntactic constructions like 'Apple was founded by Steve Jobs' or 'In 1976, Jobs, along with Wozniak, founded Apple,' because the core dependency structure remains consistent. These methods provide a strong baseline for relation extraction, often used to create high-precision rules or to generate training data for more advanced models.

While pattern-based systems are effective, they can be brittle and require manual effort. A more robust and generalizable approach is to use supervised machine learning, which frames relation extraction as a classification problem. The goal is to train a model that can predict the correct relation type for any given pair of entities, (e_1, e_2) , within a sentence.

This process requires a corpus annotated with entities and the relations between them. For each sentence containing a marked entity pair, we extract a feature vector that represents the pair and its context. A classifier, such as a Support Vector Machine (SVM) or Logistic Regression model, is then trained on these vectors and their corresponding relation labels. The set of possible labels includes all pre-defined relation types (e.g., Located_In, Founded_By) plus a crucial No_Relation label for pairs that are co-located in a sentence but are not semantically related.

Common features used for this task include:

- Entity-based features: The text of the entities, their types (e.g., *PERSON*, *OR-GANIZATION*), and their order in the sentence.
- Contextual features: The sequence of words and their Part-of-Speech tags between the two entities. The words immediately preceding e_1 and following e_2 are also often included.
- Syntactic features: The dependency path between the two entities in the sentence's parse tree. This is a particularly powerful feature, as it captures the grammatical relationship directly, abstracting away from the specific surface words.

Once trained, the model can be applied to new, unseen text to identify and classify relationships between entities.

To make this concrete, let's illustrate how relations are extracted from a biographical text. Consider the sentence:

After leaving Atari, Steve Jobs co-founded Apple with Steve Wozniak in 1976.

First, a Named Entity Recognition (NER) system would identify the entities: [PER Steve Jobs], [ORG Apple], and [PER Steve Wozniak]. The relation extraction task then processes pairs of these entities.

For the pair (Steve Jobs, Apple), a pattern-based system could use a simple rule that looks for the verb 'found' (or its variants) between a PER and an ORG entity, triggering the extraction. A supervised machine learning classifier, in contrast, would convert the context—the words 'co-founded' and 'with', their POS tags, and the syntactic path between the entities—into a feature vector. This vector would then be fed into a model trained to predict relation types.

In either case, the system would ideally extract the structured fact founded_by (Apple, Steve Jobs). Notice that the same sentence also yields founded_by(Apple, Steve Wozniak), demonstrating how multiple relations can be harvested from a single piece of text.

To see Information Extraction in action, let's consider a practical case study: automatically populating a knowledge base of corporate executives and their roles from financial news reports. The goal is to create a structured database that can answer queries like 'Who is the CEO of InnovateCorp?' or 'List all executives at TechSolutions.'

Our input is a stream of unstructured text from news articles:

- 'Following the quarterly report, Jane Doe, the newly appointed Chief Financial Officer of InnovateCorp, presented the financial outlook.'
- 'Meanwhile, John Smith will be stepping down as Chief Operating Officer at rival firm TechSolutions.'
- 'InnovateCorp's founder, Emily Chen, remains on the board.'

An IE pipeline would process this text in two main stages. First, *Named Entity Recognition* (NER) identifies the key entities. The system would tag the text as follows:

- \bullet Jane Doe (PERSON), Chief Financial Officer (ROLE), InnovateCorp (ORGANIZATION)
- \bullet John Smith (PERSON), Chief Operating Officer (ROLE), TechSolutions (ORGANIZATION)
- InnovateCorp (ORGANIZATION), Emily Chen (PERSON)

Next, the *Relation Extraction* stage analyzes the context surrounding these entities to identify relationships between them. It looks for linguistic patterns like '...[PERSON], the [ROLE] of [ORGANIZATION]...' or '[ORGANIZATION]'s founder, [PERSON]...'. From our examples, the system would extract a set of structured relational triples:

- (Jane Doe, is_cfo_of, InnovateCorp)
- (John Smith, was_coo_of, TechSolutions)
- (Emily Chen, is_founder_of, InnovateCorp)

These triples are then added to our knowledge base. As the system processes thousands of articles, it builds a rich, interconnected graph of people, their roles, and their affiliations. This task is challenging; the system must handle ambiguity (many people named John Smith), track changes in roles over time, and resolve coreferences (e.g., understanding that 'she later joined the board' refers to Jane Doe). Nonetheless, this application demonstrates

| Aspect | Information Retrieval | Information Extraction |
|--------|--------------------------|--------------------------|
| Goal | Find relevant documents | Extract structured facts |
| Output | Ranked list of documents | Filled database records |
| Scope | Document-level | Sub-document level |

Figure 10.7: A summary of the core differences between Information Retrieval and Information Extraction.

the core power of IE: transforming vast quantities of messy, unstructured text into a valuable, queryable asset.

Having explored both Information Retrieval and Information Extraction, we can now synthesize their relationship. The fundamental differences between them, summarized concisely in Fig. 10.7, are critical. At its core, IR answers the question, 'Which documents are about my topic?' Its goal is discovery, and its output is a ranked list of documents. In contrast, IE answers the question, 'What specific facts are in this text?' Its goal is structuring, and its output is structured data, such as filled database records or a knowledge graph.

Despite these distinct objectives, IR and IE are highly complementary and are often combined in a powerful pipeline. A typical workflow begins with an IR system retrieving a set of relevant documents from a vast corpus—for example, finding all articles about company acquisitions. This step acts as a coarse filter, dramatically reducing the search space. Subsequently, an IE system processes this smaller, topic-focused collection to extract precise, structured details like the acquirer, the target, and the acquisition price for each event. This synergistic process efficiently transforms a broad information need into a structured, queryable knowledge base, highlighting the practical interplay between these two essential fields.

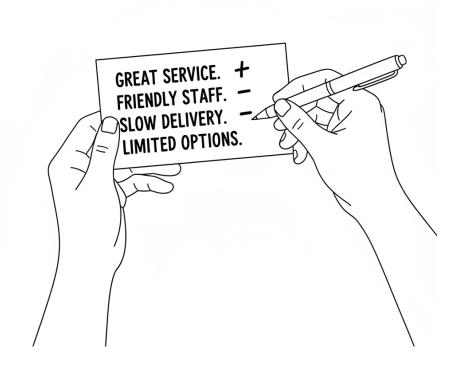
While the statistical and rule-based methods described in this chapter represent the foundational principles of Information Retrieval and Information Extraction, the field has been profoundly reshaped by modern deep learning. These techniques have moved beyond the limitations of sparse, keyword-based representations and manual feature engineering, setting a new standard for performance.

In Information Retrieval, the paradigm has shifted from sparse TF-IDF vectors to dense vector representations, or embeddings, generated by large neural models. This approach, often called dense retrieval, allows search systems to move beyond simple keyword matching and capture the semantic intent behind a query. A search for 'films about the US Civil War' can thus intelligently match a document titled 'A historical drama on the conflict between North and South' even without significant term overlap.

For Information Extraction, transformer-based models like BERT (which we will explore in Chapter 12) have become the state-of-the-art. By pre-training on vast amounts of text, these models learn rich, contextual representations of words. When fine-tuned for tasks like Named Entity Recognition or Relation Extraction, they consistently outperform older systems that rely on hand-crafted features like word shape or gazetteers. This powerful contextual understanding enables them to accurately identify entities and their relationships with unprecedented accuracy, pushing the boundaries of what can be automatically structured from raw text.

Chapter 11

Sentiment Analysis and Opinion Mining



This chapter introduces *sentiment analysis*, the computational task of identifying and extracting opinions, sentiments, and emotions from text data.¹ While previous chapters often focused on the objective, factual content of language, our focus now shifts to its *subjective* dimension. The core goal is to automatically determine the author's attitude, evaluation, or emotional state concerning a particular topic or entity.

The most common task in sentiment analysis is determining a text's *polarity*: is the expressed opinion positive, negative, or neutral? This seemingly simple classification powers countless applications. As we will see, the immense volume of user-generated content on the web—from product reviews and social media posts to political commentary—has made sentiment analysis an indispensable tool for understanding public opinion at scale.

The practical impact of sentiment analysis is vast, touching nearly every industry that deals with human-generated text. Its core value lies in transforming massive volumes of unstructured opinions into structured, actionable data, enabling automated, large-scale analysis that would be impossible for humans to perform manually. This capability provides powerful insights for decision-making across numerous domains.

Key application areas include:

- Business Intelligence: Companies systematically process customer feedback from product reviews, support tickets, and surveys. This helps them to monitor brand perception, track customer satisfaction over time, and quickly identify product or service flaws.
- Social and Political Analysis: Political campaigns and social scientists gauge public opinion on candidates and policies by analyzing millions of posts on social media platforms. This provides a real-time pulse on public mood and reaction to events.
- Financial Markets: Analysts sift through news articles and investor forums to predict market trends, a practice known as *financial sentiment analysis*, by capturing the mood surrounding specific stocks or the economy as a whole.
- **Healthcare:** Researchers can analyze patient narratives from online forums or reviews to understand real-world experiences with treatments and medical services.

This chapter will guide you through the core methodologies for sentiment analysis, charting a course from transparent, rule-based systems to more complex, data-driven models. Our exploration begins with the most intuitive approach: *lexicon-based* methods. These techniques operate on a simple principle, using pre-compiled dictionaries of words with associated sentiment scores to 'tally up' the overall feeling of a text.

We will then transition to more powerful and adaptable techniques based on *supervised machine learning*. This paradigm reframes sentiment analysis as a classification task. A model is trained to predict sentiment from a large corpus of pre-labeled examples (e.g., positive or negative reviews). This section will cover the classic machine learning pipeline, from feature extraction to training classifiers, and briefly introduce how modern neural networks are applied to the task.

Before classifying sentiment, we must first decide on the scope of the text to analyze. Sentiment analysis is not a monolithic task; it operates at different levels of granularity, and the appropriate level depends on the required detail. We can distinguish three main levels:

¹The terms *sentiment analysis* and *opinion mining* are often used interchangeably, and we will follow that convention in this chapter.



Figure 11.1: A diagram illustrating the three levels of sentiment analysis: document-level, sentence-level, and aspect-based, using a hotel review as an example.

- **Document-Level:** This is the broadest scope. The goal is to classify the overall sentiment of an entire document, such as a product review or a blog post. It assumes the text expresses a single, primary opinion. For instance, is a movie review, taken as a whole, positive, negative, or neutral?
- Sentence-Level: This level refines the analysis by determining the sentiment of each individual sentence. It can capture nuances that document-level analysis misses, as a single review might contain both positive and negative statements about different things.
- Aspect-Based (or Fine-Grained): This is the most detailed level. It moves beyond a single polarity to identify opinions about specific aspects or features of an entity. For a hotel review, it would identify the sentiment towards the 'service' separately from the 'cleanliness' or the 'location'.

The choice of level directly influences the complexity of the task and the specificity of the insights gained.

To make these distinctions concrete, consider a typical hotel review, which we can analyze at each of the three levels as illustrated in Fig. 11.1.

'The location was perfect, but the room was a bit noisy. Overall, a good stay.'
The insights derived depend entirely on the chosen granularity:

- 1. **Document-Level:** At the coarsest level, a model considers the entire text as a single document. Influenced by the phrase 'a good stay,' it would likely assign a single, overarching *Positive* polarity to the review. This gives a general summary but obscures important details.
- 2. **Sentence-Level:** Here, the review is broken down into its constituent sentences or clauses. A system would classify 'The location was perfect' as *Positive* and 'the room was a bit noisy' as *Negative*. This reveals the conflicting opinions present within the same review.

3. **Aspect-Based Level:** This most fine-grained analysis links sentiment directly to specific features, or *aspects*, of the hotel. The model identifies that the sentiment towards the aspect 'location' is Positive ('perfect'), while the sentiment towards the aspect 'room' is Negative ('noisy'). This level provides the most specific and actionable feedback.

The first major approach we will cover is the **lexicon-based method**. This technique is intuitive and powerful, operating without the need for model training on labeled data. At its core, this method relies on a pre-compiled sentiment *lexicon*—essentially a large dictionary where individual words and phrases are tagged with a prior sentiment polarity. This polarity can be a simple categorical label (e.g., positive, negative, neutral) or a numerical score indicating sentiment intensity.

For instance, a lexicon might assign numerical scores such as:

 \bullet "wonderful": +2

• "terrible": -2

ullet "adequate": +0.5

• "ordinary": 0

The sentiment of a given text is determined by first tokenizing it into words and then aggregating the polarity scores of the words found in the lexicon. A common aggregation function is a simple sum. Let S(T) be the sentiment score for a text T composed of words w_1, w_2, \ldots, w_n . A basic model would be:

$$S(T) = \sum_{i=1}^{n} \text{polarity}(w_i)$$

If the final score S(T) is above a positive threshold, the text is classified as positive; if below a negative threshold, it's deemed negative.

A variety of pre-compiled sentiment lexicons are available, each constructed differently and offering unique properties. Among the most well-known are the MPQA Subjectivity Lexicon and SentiWordNet.

The MPQA Subjectivity Lexicon was created through extensive manual annotation. It contains over 8,000 words, each tagged with its polarity (positive, negative, or neutral) and its subjectivity strength (strong or weak). For example, the word excellent is tagged as strong-positive, while adequate might be weak-positive. Its direct, human-curated labels make it a reliable and straightforward resource.

In contrast, **SentiWordNet** is a resource built on top of the WordNet lexical database. Instead of labeling individual words, it assigns three scores to each *synset* (a group of synonymous words representing a single concept): positivity, negativity, and objectivity. These three scores for any given synset always sum to 1.0. For instance, the primary sense of *gripping* might be assigned {Positive: 0.75, Negative: 0.0, Objective: 0.25}. This structure provides a more nuanced view, acknowledging that a word's sentiment can depend on its specific meaning in context.

The fundamental algorithm for lexicon-based sentiment analysis is both intuitive and computationally efficient. It operates by aggregating the polarity of individual words to determine the overall sentiment of a text, a process summarized visually in Fig. 11.2. The procedure can be broken down into four main steps:

1. **Tokenization:** First, the input text is segmented into a list of words, or tokens.

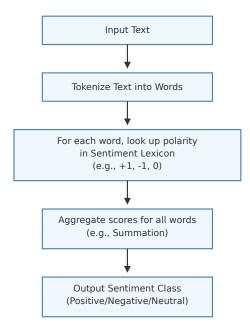


Figure 11.2: A flowchart diagram illustrating the steps of a simple lexicon-based sentiment analysis algorithm, from tokenization to classification.

- 2. **Polarity Scoring:** Each token is looked up in the sentiment lexicon. If the word is present, its pre-assigned polarity score is retrieved (e.g., +1 for positive, -1 for negative). Words not found in the lexicon are typically assigned a neutral score of 0.
- 3. **Aggregation:** The scores for all tokens in the text are combined, usually through simple summation, to produce a final sentiment score. This is calculated as: $S_{text} = \sum_{i=1}^{n} \text{polarity}(w_i)$ where w_i is the *i*-th word in the *n*-word text.
- 4. Classification: The final aggregated score, S_{text} , is then compared against a threshold to assign a category. For instance:
 - Positive if $S_{text} > 0$
 - Negative if $S_{text} < 0$
 - Neutral if $S_{text} = 0$

This method provides a transparent and easily implemented baseline for sentiment classification, serving as a valuable first approach before moving to more complex models.

To illustrate, let's calculate the sentiment of a short movie review using this aggregation method. Consider the sentence: 'The acting was brilliant and the story was truly amazing, but the ending felt a bit predictable.'

First, we consult a simple sentiment lexicon where words are assigned polarity scores. For this example, we'll use +1 for positive, -1 for negative, and 0 for neutral or out-of-lexicon words.² Our lexicon might contain:

- \bullet brilliant: +1
- amazing: +1
- predictable: -1

²In practice, stop words like 'the', 'was', and 'a' are often filtered out, and words not found in the lexicon are ignored. We treat them as having a score of 0 for simplicity.

The overall sentiment score, S, for the review is the sum of the polarity scores, $P(w_i)$, for each relevant word w_i :

$$S = \sum_{w_i \in \text{review}} P(w_i)$$

For our sentence, the calculation is:

$$P(\text{brilliant}) + P(\text{amazing}) + P(\text{predictable}) = (+1) + (+1) + (-1) = +1$$

The final score of +1 is positive, so the model classifies the review as having a positive overall sentiment. This simple approach effectively captures the dominant opinion despite the presence of a negative term.

While elegant in their simplicity, lexicon-based methods face significant challenges that limit their accuracy in real-world scenarios. Their reliance on context-free word scores means they often fail to capture the nuances of human language. Key difficulties include:

- Handling Negation. A simple negation word like 'not' or 'never' can completely reverse the polarity of a subsequent phrase. For instance, in 'The movie was not good,' a basic algorithm might incorrectly register a positive sentiment from the word 'good' while failing to account for the critical negation that inverts its meaning.
- Accounting for Intensifiers. Language is filled with modifiers that strengthen or weaken sentiment. Words like 'very' or 'extremely' amplify polarity (e.g., 'very happy'), while words like 'slightly' or 'barely' diminish it. Simple lexicon lookups that assign a fixed score to each word often ignore the crucial impact of these adverbs.
- Adapting to Domain-Specific Context. The sentiment of a word can be highly dependent on its domain. The word 'unpredictable' is likely positive when describing a thriller novel's plot but negative for a car's braking system. A general-purpose lexicon will fail to capture these vital domain-specific meanings.

Case Study: E-commerce Product LaunchConsider an e-commerce company, *InnovateTech*, launching a new smart speaker. To gauge customer reaction in real-time, they need to analyze a high volume of tweets and on-site product reviews. Building a supervised machine learning model would require time and labeled data they don't yet have. Instead, they rapidly deploy a lexicon-based system.

The system ingests incoming feedback, tokenizes the text, and calculates a sentiment score for each comment by summing the polarities of its words from a general-purpose lexicon. A simple dashboard visualizes the aggregate sentiment. If the score trends sharply negative, the product team is immediately alerted. This enables them to pinpoint and address specific issues—like a software bug or unexpected shipping delays—long before these problems would impact sales figures, providing an invaluable and low-cost early warning system.

While lexicon-based methods provide a strong and interpretable baseline, they can be rigid and struggle with domain-specific language or subtle contexts. A more powerful and adaptable alternative is to frame sentiment analysis as a supervised machine learning problem. In this paradigm, instead of providing the machine with a dictionary of sentimental words, we provide it with a large corpus of texts that have already been labeled with the correct sentiment.

The task is thus transformed into a classic text classification problem. The model's objective is to learn a function that maps an input text (e.g., a movie review) to a predefined sentiment category such as *positive*, *negative*, or *neutral*. During a 'training' phase,

| Document | the | movie | was | good | boring |
|------------------------|-----|-------|-----|-----------------------|--------|
| 'The movie was good' | 1 | 1 | 1 | 1 | 0 |
| 'The movie was boring' | 1 | 1 | 1 | 0 | 1 |

Figure 11.3: An illustration of the Bag-of-Words (BoW) feature extraction process. Two sentences are converted into sparse numerical vectors based on the word counts from a combined vocabulary. Most values are zero for any given document, a characteristic of BoW representations.

a machine learning algorithm examines thousands or even millions of these labeled examples. It automatically learns to identify the textual features—words, phrases, and their combinations—that are predictive of each sentiment class. For instance, it might learn that words like 'brilliant' and 'amazing' strongly correlate with a positive label, while 'disappointing' and 'awful' suggest a negative one, without being explicitly told the polarity of these words. The key is that the model discovers these associations from the data. The ultimate goal is to produce a trained classifier that can generalize from these examples to accurately predict the sentiment of new, unseen text.

For supervised machine learning to work, the model needs a large corpus of examples where the correct answer—in this case, the sentiment label—is already known. This collection of labeled texts is called the training dataset. While one could manually pay human annotators to label thousands of documents as positive or negative, this process is incredibly slow and expensive.

A far more common and efficient strategy is to find data where labels already exist as a natural byproduct. Online reviews are a perfect source. A product review on an e-commerce site or a movie review on a platform like IMDb is typically accompanied by a star rating (e.g., 1 to 5 stars). We can use this explicit rating as a proxy for the sentiment label. A common heuristic for creating a binary classification dataset is:

- Reviews with 4 or 5 stars are automatically assigned a **positive** label.
- Reviews with 1 or 2 stars are automatically assigned a **negative** label.

Notice that 3-star reviews are often discarded in this process. They can represent neutral sentiment, mixed opinions, or other ambiguities, and removing them creates a cleaner training signal for the model. This method allows for the rapid creation of massive datasets, where each instance is a pair (d, c) consisting of a document d (the review text) and its derived class label c (positive or negative).

Machine learning algorithms operate on numerical data, not raw text. Therefore, a crucial step in building a supervised sentiment classifier is to convert text documents into fixed-length numerical feature vectors. This process is known as *feature extraction* or *vectorization*.

The most common and intuitive approach is the **bag-of-words** (**BoW**) model. This model represents text by disregarding grammar and word order, treating it as an unordered collection—or 'bag'—of its words. The process involves two steps: first, a vocabulary of all unique words in the entire training corpus is compiled. Second, each document is converted into a vector with the same length as the vocabulary. Each dimension of the vector corresponds to a unique word, and its value is typically the frequency (count) of that word in the document. This transformation is illustrated in Fig. 11.3, which shows how two sentences are converted into sparse vectors where most values are zero.

A limitation of BoW is that very common words (e.g., 'the,' 'is,' 'a') dominate the frequency counts but often carry little semantic weight. An effective alternative is the

Term Frequency-Inverse Document Frequency (TF-IDF) weighting scheme. TF-IDF evaluates how important a word is to a document within a collection. It is the product of two statistics:

- Term Frequency (TF): Measures how frequently a term t appears in a document d. tf(t,d)
- Inverse Document Frequency (IDF): Measures how informative a term is by down-weighting common terms. It is calculated as the logarithm of the total number of documents |D| divided by the number of documents containing the term t. $idf(t,D) = \log \frac{|D|}{|\{d \in D: t \in d\}|}$

The resulting TF-IDF score, $tfidf(t,d,D) = tf(t,d) \times idf(t,D)$, replaces the raw counts in the feature vector. This gives higher weight to terms that are frequent in a specific document but rare across the entire corpus, making them more discriminative for classification. These resulting vectors, whether from BoW or TF-IDF, serve as the input for training a classification model.

Once our text is represented as numerical feature vectors, the task of sentiment classification transforms into a standard supervised machine learning problem. We can feed these vectors into any number of classification algorithms to train a model that learns to distinguish between different sentiment polarities. We will briefly overview two classic and highly effective models for this task: Naive Bayes and Support Vector Machines.

The **Naive Bayes** classifier is a simple yet powerful probabilistic model based on Bayes' theorem. It calculates the probability of a document belonging to a class (e.g., positive or negative) given its features (the words it contains). Its 'naive' assumption is that every feature is conditionally independent of every other feature, given the class. While this is almost never true for natural language, the model often performs surprisingly well. To classify a new document, we choose the class \hat{c} that maximizes this probability:

$$\hat{c} = \underset{c \in \{\text{pos, neg}\}}{\arg \max} P(c) \prod_{i=1}^{n} P(w_i|c)$$

Here, P(c) is the prior probability of a class (how common it is), and $P(w_i|c)$ is the likelihood of word w_i appearing in a document of that class, both learned from the training data. Due to its simplicity and computational efficiency, Naive Bayes serves as an excellent baseline model.

In contrast, **Support Vector Machines (SVMs)** are discriminative classifiers. Instead of modeling probabilities, an SVM aims to find an optimal hyperplane that separates the data points of different classes in the high-dimensional feature space. For text, this means finding a decision boundary that best separates the vectors of positive reviews from negative ones. The optimal hyperplane is the one that maximizes the *margin*—the distance to the nearest data point from either class. This large-margin principle helps the model generalize better to unseen data. Furthermore, SVMs can employ the 'kernel trick' to learn complex, non-linear boundaries, often achieving high accuracy and making them a robust choice for sentiment classification.

To illustrate, let's walk through training a simple Naive Bayes classifier on a small, labeled dataset. Imagine our training corpus consists of just four sentences:

• Positive: 'a great film', 'I love this film'

• Negative: 'a bad plot', 'I hate the acting'

Positive Negative True Positive False Negative Positive (FN) (TP) **True Label** False Positive True Negative Negative (FP) (TN) Accuracy = (TP + TN) / (TP + TN + FP + FN)Precision = TP / (TP + FP) Recall = TP / (TP + FN)

Predicted Label

Figure 11.4: A confusion matrix for a binary classification task. It visualizes the four possible outcomes—True Positive (TP), True Negative (TN), False Positive (FP), and False Negative (FN)—which are used to calculate key evaluation metrics.

F1-score = 2TP / (2TP + FP + FN)

Our goal is to classify a new, unseen sentence: 'I love the film'. The classifier works by calculating a score for each potential class (positive and negative) and selecting the one with the higher score. The calculation is based on the Naive Bayes formula, where the score for a class is proportional to the class's overall probability multiplied by the probabilities of each word in the new sentence, given that class:

$$P(\text{class}|\text{document}) \propto P(\text{class}) \prod_{i=1}^{n} P(\text{word}_{i}|\text{class})$$

First, we compute the score for the *positive* class. The class probability, or *prior*, P(positive) is $\frac{2}{4} = 0.5$, as two of our four training documents are positive. We then calculate the *likelihood* of each word from our new sentence given the positive data. Words like 'love' and 'film' appear in our positive training examples, giving them relatively high likelihood values.

Next, we compute the score for the negative class. The prior P(negative) is also 0.5. However, the words 'love' and 'film' do not appear in any negative training sentences. A raw calculation would give them a probability of zero, which would incorrectly make the entire score for the negative class zero. To solve this, we use a technique called **Laplace smoothing** (or add-1 smoothing), which adds one to every word count. This ensures all words have a small, non-zero probability. Still, the likelihoods P(love|negative) and P(film|negative) will be much lower than their positive counterparts.

Finally, by multiplying the prior and the word likelihoods for each class, we find that the score for *positive* is significantly higher. The classifier therefore correctly labels 'I love the film' as positive.

Once a sentiment classifier is trained, we need a systematic way to measure its performance. A simple 'percentage correct' can be misleading, especially with imbalanced data. A more nuanced evaluation starts with a *confusion matrix*, which visualizes how the model's predictions align with the true labels, as shown in Fig. 11.4. This matrix categorizes predictions into four types:

- True Positives (TP): Positive instances correctly identified as positive.
- True Negatives (TN): Negative instances correctly identified as negative.
- False Positives (FP): Negative instances incorrectly labeled as positive (a Type I error).
- False Negatives (FN): Positive instances incorrectly labeled as negative (a Type II error).

From these counts, we derive several key metrics. *Accuracy* is the most intuitive, representing the overall fraction of correct predictions:

$$\label{eq:accuracy} \text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

To understand the trade-offs, we use *Precision* and *Recall*. Precision measures exactness: of all instances the model predicted as positive, how many were actually positive?

$$Precision = \frac{TP}{TP + FP}$$

Recall measures completeness: of all the instances that were truly positive, how many did the model find?

$$Recall = \frac{TP}{TP + FN}$$

The F1-score combines these two into a single number by calculating their harmonic mean, providing a balanced measure that is useful when both precision and recall are important.

$$F_1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2TP}{2TP + FP + FN}$$

Together, these metrics provide a robust framework for comparing different sentiment analysis models.

While powerful, the machine learning approaches discussed so far share a fundamental limitation: bag-of-words and TF-IDF representations discard word order. For these models, the sentences 'The service was not good' and 'The service was good, not!' appear very similar, despite having opposite meanings. To capture the sequential nature of language, we can turn to neural network architectures.

Recurrent Neural Networks (RNNs) and their variants, such as LSTMs (Long Short-Term Memory), are explicitly designed for sequential data. An RNN processes a sentence one word at a time, maintaining a hidden state that acts as a form of memory. At each step, the network updates its hidden state based on the current word and the state from the previous word.

This process allows the model to build a rich vector representation that encodes information from the entire sequence. The final hidden state, which represents a contextualized summary of the sentence, can then be fed into a classifier to predict the overall sentiment. By modeling dependencies between words, these neural approaches can better understand complex linguistic phenomena like negation and long-distance relationships, often leading to significantly improved classification performance.

While document-level and sentence-level classification provide a useful high-level summary, they often lack the necessary granularity for detailed analysis. A single product review, for example, can contain conflicting opinions about different features. A customer might praise a phone's *battery life* while simultaneously criticizing its *camera*. Assigning a single positive or negative label to the entire text would obscure this crucial detail, leading to a loss of valuable information.

To achieve a more nuanced understanding, we introduce **Aspect-Based Sentiment Analysis (ABSA)**. This more sophisticated task aims to identify opinions about specific entities or their attributes, which are formally referred to as *aspects*. Instead of asking 'Is this review positive?', ABSA asks 'What specific features are being discussed, and what is the sentiment towards each one?' This fine-grained approach provides far more actionable insights, moving beyond a simple polarity score to a detailed breakdown of what people like and dislike. The core challenge of ABSA is twofold: first, identifying the aspects themselves, and second, determining the sentiment expressed towards each one.

To achieve this fine-grained analysis, Aspect-Based Sentiment Analysis (ABSA) is typically broken down into two fundamental sub-tasks, often executed as a pipeline. Successfully solving the first task is a prerequisite for attempting the second.

- 1. **Aspect Term Extraction (ATE):** The initial goal is to identify the specific entities or their attributes that are the target of an opinion. These are the 'aspects' themselves. In a product review, these might be nouns or noun phrases like 'battery life,' 'screen resolution,' or 'customer service.' The output of this stage is a list of all explicit opinion targets mentioned in the text. For example, given the sentence, 'The phone's camera is amazing, but its speakers are tinny,' an ATE system should extract the terms *camera* and *speakers*.
- 2. Aspect Sentiment Classification (ASC): Once the aspect terms have been extracted, the next task is to determine the polarity of the sentiment expressed towards each specific aspect. This is a targeted classification problem. For each term identified by the ATE step, the ASC system assigns a sentiment label, such as positive, negative, or neutral. Continuing the previous example, the sentiment for the aspect camera would be classified as positive, while the sentiment for speakers would be classified as negative.

Document-level sentiment analysis can be a blunt instrument, often failing to capture the full picture when a text contains mixed opinions. To illustrate the power and necessity of Aspect-Based Sentiment Analysis (ABSA), consider this common type of review:

• 'The pizza was delicious, but the service was slow.'

A document-level classifier might label this sentence as neutral, as the positive and negative sentiments could effectively cancel each other out. This overlooks crucial information. ABSA, in contrast, dissects the sentence to extract more granular insights. The goal is to identify the specific *aspects* being discussed and then determine the *polarity* of the opinion expressed towards each one.

In this example, there are two distinct aspect-opinion pairs:

1. Aspect: pizza

Opinion Expression: delicious
Sentiment Polarity: Positive (+)

2. Aspect: service

• Opinion Expression: slow

• Sentiment Polarity: Negative (-)

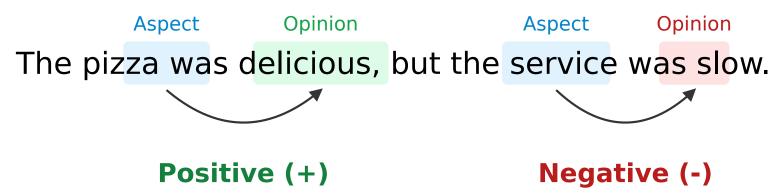


Figure 11.5: A diagram illustrating Aspect-Based Sentiment Analysis (ABSA). The sentence 'The pizza was delicious, but the service was slow' is deconstructed into its two core aspect-opinion pairs. The 'pizza' aspect is linked to the positive opinion 'delicious', while the 'service' aspect is linked to the negative opinion 'slow', visually representing the granular analysis provided by ABSA.

This process, visually deconstructed in Fig. 11.5, provides a much richer understanding of the reviewer's feedback. A restaurant owner reading this review learns precisely what customers liked (the food) and what needs improvement (the service). This level of detail is far more actionable than a simple 'neutral' rating and demonstrates the core value of the ABSA task.

The first sub-task in ABSA is aspect term extraction—identifying the specific features or entities being discussed. Techniques for this range from straightforward linguistic rules to sophisticated machine learning models.

Rule-based approaches leverage common grammatical patterns to find aspects. A simple heuristic might be to extract any noun or noun phrase that is modified by a sentiment-bearing adjective. For example, in 'The *pizza* was delicious,' a rule could identify pizza as an aspect because it is a noun connected to the positive adjective delicious. These methods often rely on outputs from earlier pipeline stages, such as Part-of-Speech tagging and dependency parsing, to create precise patterns. While easy to implement, they can be brittle and struggle with linguistic variation.

A more robust and dominant approach frames extraction as a supervised sequence labeling task. Here, the goal is to assign a tag to each word in a sentence, indicating whether it is part of an aspect term. A common tagging schema is BIO (Beginning, Inside, Outside). For the sentence 'The customer service was slow,' the desired output would be:

• The/O customer/B-ASP service/I-ASP was/O slow/O

In this schema, B-ASP marks the beginning of an aspect, I-ASP marks the continuation of an aspect, and O marks words that are outside any aspect. Models are trained on a manually annotated corpus to learn to predict these tag sequences. Historically, powerful models for this task include Conditional Random Fields (CRFs), which excel at sequence labeling by considering the context of the entire sentence when predicting each word's tag.

Once aspect terms like 'pizza' or 'service' have been identified, the subsequent task is to assign a specific sentiment polarity to each one. This is crucial because the overall sentence sentiment might be mixed or misleading. The sentiment associated with an aspect is almost always determined by its immediate local context.

A straightforward method is to use a *window-based* approach. For each aspect term, we define a context window (e.g., the 5 words before and after it) and apply a lexicon-based sentiment calculation only to the words within that window. The sentiment for the

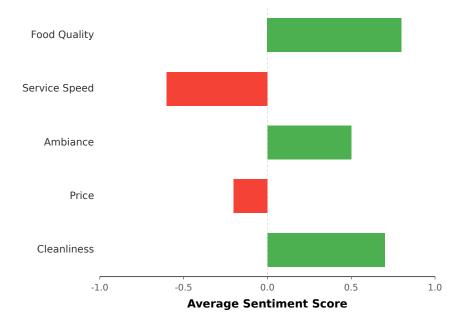


Figure 11.6: A horizontal bar chart visualizing the results of an Aspect-Based Sentiment Analysis on restaurant reviews. Key aspects are listed on the Y-axis, and their aggregated sentiment scores (from -1 to +1) are shown on the X-axis. The chart clearly distinguishes strengths like 'Food Quality' and 'Cleanliness' (positive bars) from weaknesses like 'Service Speed' (negative bar), turning unstructured feedback into actionable business intelligence.

aspect a can be a simple sum of polarity scores:

$$S(a) = \sum_{w \in \text{window}(a)} \text{polarity}(w)$$

A more powerful technique frames this as a supervised classification problem. A model is trained on examples where each instance consists of a sentence, a target aspect, and its corresponding polarity label (positive, negative, or neutral). Features for the classifier are engineered from the context, such as the bag-of-words within the window or syntactic dependencies linking the aspect to opinion words. For example, a dependency parse could explicitly link the adjective 'delicious' to the noun 'pizza.'

Modern neural network architectures can learn this context automatically, using attention mechanisms to focus on the most relevant opinion words for a given aspect, regardless of their distance from the aspect term itself.

To illustrate the practical power of ABSA, consider a large restaurant chain analyzing thousands of online customer reviews. A simple document-level analysis might reveal an average sentiment score of 'mildly positive,' but this provides no direction for improvement. By applying ABSA, the chain can dissect this feedback into more meaningful components.

The system first extracts common aspects discussed by customers—such as *food quality*, service speed, ambiance, price, and cleanliness. For each mention of an aspect, it then classifies the associated sentiment. After aggregating the data across all reviews, a clear picture emerges. The analysis might reveal that while 'food quality' consistently receives high praise, 'service speed' is a significant source of negative sentiment, a critical nuance that would otherwise be lost in an overall average score.

This granular breakdown, visualized in **Fig. 11.6**, transforms vague feedback into actionable intelligence. The chart clearly highlights the restaurant's operational strengths and critical weaknesses at a glance. Armed with this specific insight, management can bypass generic improvement plans and focus resources precisely where they are needed—

for instance, by retraining staff to improve service efficiency rather than needlessly altering a popular menu. ABSA thus provides a direct bridge from unstructured customer opinion to strategic business decisions.

In summary, we have explored the main paradigms of sentiment analysis, each presenting a unique set of trade-offs. The optimal choice depends on the specific task, available resources, and desired level of analytical depth.

- Lexicon-Based Methods: These are valued for their simplicity, speed, and independence from labeled training data. They provide an excellent baseline but are often brittle, struggling to handle contextual nuances like negation, sarcasm, or domain-specific jargon.
- Supervised Machine Learning: This approach, which frames sentiment analysis as a classification task, generally yields higher accuracy. Its power, however, is contingent upon the availability of a large, high-quality annotated corpus, making data acquisition a significant bottleneck.
- Aspect-Based Sentiment Analysis (ABSA): Offering the most granular insight, ABSA pinpoints opinions about specific features. This detailed output is highly actionable but comes at the cost of substantially increased complexity in both model design and annotation.

Despite the effectiveness of the methods covered, several significant challenges persist in sentiment analysis, often stemming from the complex nature of human expression. These problems require moving beyond simple word polarity and towards a deeper contextual understanding.

Key challenges include:

- Sarcasm and Irony: These figures of speech intentionally use positive words to convey a negative sentiment. For example, in the sentence 'I just love being stuck in traffic,' the literal positivity of 'love' is inverted by the real-world context. Simple models struggle to detect this contradiction.
- Context-Dependent Sentiment: The polarity of a word is not always fixed. The word *unpredictable* is desirable for a movie plot but highly undesirable for a car's brakes. This domain-specific nature challenges universal, one-size-fits-all sentiment lexicons.

Solving these issues often requires more sophisticated models that can incorporate world knowledge and nuanced pragmatic understanding.

The field of opinion mining continues to evolve rapidly. While the methods discussed provide a strong foundation, the future lies in leveraging the contextual understanding of Large Language Models (LLMs), which we will explore in the final chapter. These models excel at detecting subtle sentiment, sarcasm, and complex aspect-based opinions with minimal training. This capability, however, introduces critical ethical questions. We must consider the potential for manipulating public opinion, the amplification of biases present in training data, and the privacy implications of analyzing personal sentiments at an unprecedented scale.

Chapter 12

The Future: Large Language Models and Ethics



| Feature | RNN/LSTM Approach | Transformer Approach |
|-------------------------------|--|---|
| Handling of Sequential Data | Processes tokens sequentially (one by one). | Processes all tokens simultaneously. |
| Parallelization Capability | Inherently sequential; difficult to parallelize. | Highly parallelizable due to self-attention. |
| Long-Range Dependency Capture | Struggles with long distances (vanishing gradients). | Directly models all token-pair relationships. |

Figure 12.1: A comparison of architectural trade-offs between Recurrent Neural Networks (RNNs/LSTMs) and the Transformer architecture, highlighting the key limitations that motivated the shift in paradigm.

Our journey through computational linguistics has covered foundational techniques, from N-gram models and syntactic parsers to the first wave of word embeddings. These methods provided crucial building blocks for language understanding, each tackling a specific slice of the problem. However, recent years have witnessed a paradigm shift of unprecedented scale and speed, moving away from specialized, feature-engineered systems. This chapter bridges the gap from that established landscape to the revolutionary era of Large Language Models (LLMs), exploring the architecture and capabilities that have fundamentally redefined the field.

Before the models that now define the state-of-the-art, the field relied heavily on Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU) networks. These architectures process language sequentially, ingesting one word at a time and updating an internal 'memory' or hidden state. This approach was intuitive and effective for many tasks, from part-of-speech tagging to machine translation.

However, this sequential design presented two fundamental obstacles to further progress. The first was computational: since the processing of word w_t depends on the hidden state from word w_{t-1} , the architecture is inherently non-parallelizable. This created a major bottleneck for training on ever-larger datasets. The second was a context-length limitation. While LSTMs were a significant improvement over simple RNNs, they still struggled to effectively connect words separated by long distances. Information from the beginning of a paragraph could be 'forgotten' by the time the model reached the end. A new paradigm was needed—one that could process all input simultaneously and model long-range dependencies more directly. Fig. 12.1 provides a detailed comparison of these architectural trade-offs.

The architecture that powered this revolution is the **Transformer**, introduced by Vaswani et al. in their seminal 2017 paper, 'Attention Is All You Need.' Its core innovation was to dispense entirely with the sequential processing of Recurrent Neural Networks (RNNs) and instead rely exclusively on a mechanism called *self-attention*. This allows the model to process all input tokens simultaneously and directly weigh the influence of every word on every other word in a sequence, capturing complex, long-range dependencies more effectively than its predecessors.

As illustrated in Fig. 12.2, the original Transformer has a sophisticated encoder-decoder structure, designed for sequence-to-sequence tasks like machine translation. The *encoder* stack takes the input sequence and builds a rich, context-aware numerical representation of it. The *decoder* stack then uses this representation, along with the output it has generated so far, to produce the next token in the output sequence. Each stack is composed of multiple identical layers containing two primary sub-components: a multi-head self-attention mechanism and a simple, position-wise feed-forward network. Residual connections and layer normalization are used around each sub-component to facilitate training.

This flexible architecture became the blueprint for most modern LLMs. Models like BERT (*Bidirectional Encoder Representations from Transformers*) leverage the encoder stack to excel at tasks requiring deep language understanding. In contrast, models in

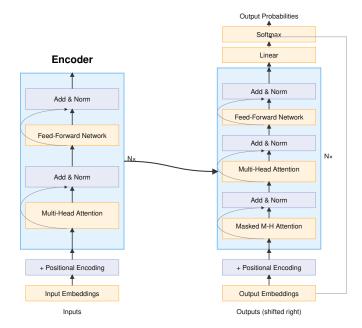


Figure 12.2: A high-level architectural diagram of the Transformer model, illustrating the encoder and decoder stacks. The diagram shows the flow of input embeddings through multi-head attention and feed-forward network layers, including residual connections and layer normalization, providing a complete visual map of the model's structure.

the GPT (*Generative Pre-trained Transformer*) family primarily use the decoder stack, making them exceptionally powerful for text generation. The key to their success lies in the self-attention mechanism, which we will now explore in detail.

At the heart of the transformer's success is the *self-attention mechanism*, a novel component that allows a model to weigh the influence of different words when processing a particular word in a sequence. Unlike recurrent architectures like LSTMs, which process sentences sequentially and can lose information over long distances, self-attention allows every word to directly interact with every other word in the sentence. This direct access is crucial for capturing complex dependencies, such as resolving the referent of a pronoun or understanding the scope of negation, regardless of how far apart the words are.

The mechanism operates on a simple but powerful principle, often described by an analogy to a library retrieval system. For each word in an input sequence, we create three distinct vectors: a **Query** (Q), a **Key** (K), and a **Value** (V).

- The **Query** vector represents a word's request for information. It's like asking, 'What other words in this sentence are relevant to my meaning?'
- The **Key** vector acts as a label or an index for a word. It advertises its own properties, saying, 'This is the kind of information I hold.'
- The Value vector contains the actual content or substance of a word. It's the information that will be passed on if a query finds its key relevant.

These three vectors are not fixed; they are generated by multiplying a word's input embedding by three separate weight matrices (W^Q, W^K, W^V) that are learned during the model's training process. This allows the model to learn the most effective way to query, label, and represent each word for the task at hand.

The self-attention process, depicted in Fig. 12.3, unfolds in three main steps. First, to determine how much attention the word at position i should pay to the word at position

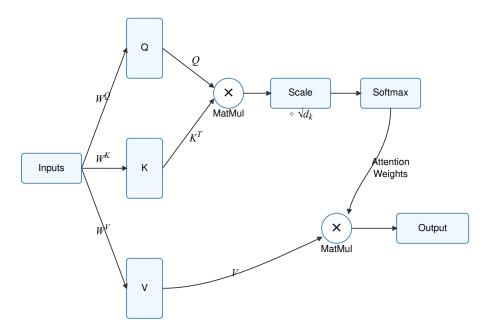


Figure 12.3: A diagram of the scaled dot-product self-attention mechanism, as described in the text.

j, we compute a compatibility score. This is done by taking the dot product of the query vector of word i (q_i) with the key vector of word j (k_j). A higher score implies greater relevance.

Second, these raw scores are scaled and normalized. The dot product scores are divided by the square root of the dimension of the key vectors, $\sqrt{d_k}$. This scaling factor prevents the dot products from growing too large, which helps stabilize the learning process. Afterwards, a softmax function is applied across all the scores for word i. This converts the scores into a set of positive weights that sum to one, effectively creating a probability distribution of attention over the entire sentence.

AttentionWeights_{ij} = softmax
$$\left(\frac{q_i \cdot k_j^T}{\sqrt{d_k}}\right)$$

Finally, the new representation for word i, which we can call its context vector z_i , is computed as a weighted sum of all the *Value* vectors in the sentence. The weights used in this sum are precisely the attention weights calculated in the previous step.

$$z_i = \sum_j \text{AttentionWeights}_{ij} \cdot v_j$$

In essence, the final representation of each word is a blend of all other words' Value vectors, where the amount contributed by each word is determined by its query-key compatibility. This entire process is performed in parallel for every word in the sequence using highly optimized matrix operations, a key advantage over sequential models. The result is a set of output vectors, each one enriched with contextual information from the entire input, forming the foundation for the transformer's deep understanding of language.

To make the concept of self-attention concrete, let's walk through an example using the sentence: 'The robot programmed the logic'. The goal is to compute a new, context-aware representation for each word that reflects its relationships with other words in the sentence. We will focus on the word 'programmed'.

First, for every word, we generate three distinct vectors from its initial embedding: a Query (Q), a Key (K), and a Value (V). These are produced by multiplying the word's

embedding by three unique weight matrices (W_Q, W_K, W_V) that are learned during training.

- The **Query** vector for 'programmed' ($q_{programmed}$) represents its search for context: 'Who did the programming, and what was programmed?'
- The **Key** vector for each word (e.g., k_{robot}) acts as a label describing the information it provides.
- The Value vector (e.g., v_{robot}) is the actual content or representation of that word.

The calculation to update the representation for 'programmed' unfolds in four steps:

- 1. **Score:** We measure the relevance of every other word to 'programmed' by taking the dot product of its query vector $(q_{programmed})$ with every key vector in the sentence. The score for the 'robot' connection is $q_{programmed} \cdot k_{robot}$. A high dot product indicates high relevance.
- 2. Scale: To aid in stable training, all scores are scaled by dividing them by the square root of the dimension of the key vectors, $\sqrt{d_k}$.
- 3. **Softmax:** A softmax function is applied to the scaled scores. This transforms them into a probability distribution—a set of *attention weights* that sum to 1. For the word 'programmed', we would expect the weights corresponding to 'robot' and 'logic' to be high, while the weights for 'the' would be very low.
- 4. Weighted Sum: The final, context-aware representation for 'programmed', which we'll call $z_{programmed}$, is a weighted sum of all the *Value* vectors in the sentence, where each value vector is multiplied by its corresponding attention weight.

This entire operation, performed for all words simultaneously using matrix operations, is elegantly captured by the formula:

$$\operatorname{Attention}(Q, K, V) = \operatorname{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

This process yields a new set of embeddings where each word's representation is now richly infused with information about its role and relationships within the entire sentence.

While a single self-attention mechanism is powerful, it forces the model to average all types of linguistic relationships—syntactic, semantic, and anaphoric—into a single representation. This can create an informational bottleneck. To overcome this, the transformer architecture employs *multi-head attention*, which allows the model to jointly attend to information from different perspectives simultaneously.

The core idea is to run the scaled dot-product attention mechanism multiple times in parallel. As illustrated in Fig. 12.4, rather than performing a single attention function on the model's full-dimensional inputs, multi-head attention first linearly projects the queries, keys, and values h times using different, learned projection matrices. These projected versions are then fed into h separate attention 'heads.' Each head can focus on a different representation subspace, enabling one head to capture, for example, subject-verb agreement while another tracks semantic similarity between distant words.

The outputs from each of the h heads are then concatenated and passed through one final linear projection to produce the layer's output. This process is formalized as:

$$MultiHead(Q, K, V) = Concat(head_1, ..., head_h)W^O$$

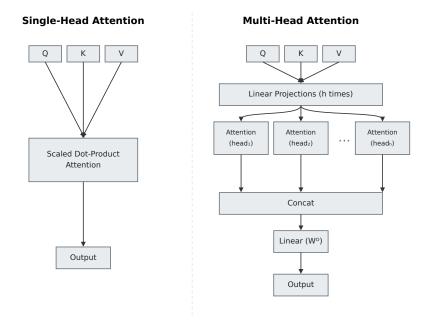


Figure 12.4: A diagram contrasting single-head attention (left) with multi-head attention (right). Multi-head attention projects the queries (Q), keys (K), and values (V) through parallel linear layers, feeding the results into multiple independent attention 'heads'. The outputs of these heads are concatenated and passed through a final linear transformation to produce the output, allowing the model to focus on different information subspaces simultaneously.

where $head_i = Attention(QW_i^Q, KW_i^K, VW_i^V)$.

The parameters to be learned are the projection matrices for each head (W_i^Q, W_i^K, W_i^V) and the final output projection matrix (W^O) . This parallel structure provides a richer, more nuanced way for the model to process relationships within a sequence, significantly enhancing its expressive power.

The self-attention mechanism, as powerful as it is, has a fundamental limitation: it is permutation-invariant. It processes the input as an unordered set of vectors, meaning it has no inherent sense of word order. To this mechanism, 'man bites dog' and 'dog bites man' would appear identical. To solve this, the Transformer architecture injects information about the position of each token directly into its input representation. This is achieved through positional encodings.

Instead of learning a separate embedding for each position, the original Transformer uses a clever, fixed formula based on sine and cosine functions of varying frequencies. For a token at position pos in the sequence and for each dimension i of the embedding vector, the positional encoding PE is calculated as follows:

 $PE_{(pos,2i)} = \sin(\frac{pos}{10000^{2i/d_{\text{model}}}}) \ PE_{(pos,2i+1)} = \cos(\frac{pos}{10000^{2i/d_{\text{model}}}})$ Here, d_{model} is the dimension of the embeddings. Each dimension of the positional encoding corresponds to a sinusoid of a different wavelength, from low to high frequency. This design is significant because it allows the model to easily learn to attend to relative positions, since the encoding for any position can be represented as a linear function of any other. The final input vector for each token, which is fed into the first Transformer block, is simply the sum of its word embedding and its corresponding positional encoding. This simple addition equips the model with the crucial sequence order information that self-attention alone lacks, enabling it to understand grammatical structure.

Prior to the current era, the standard approach involved training a model from scratch for each new natural language processing task. A model intended for sentiment analy-

Stage 1: Pre-training

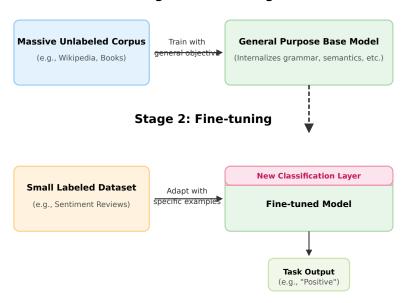


Figure 12.5: The pre-training and fine-tuning paradigm. In the first stage, a large model is pre-trained on a massive, unlabeled text corpus to develop a general understanding of language. In the second stage, this base model is fine-tuned using a much smaller, task-specific labeled dataset to adapt it for a particular application, such as sentiment analysis.

sis, for instance, would be trained solely on a dataset of labeled reviews, learning about language only through that narrow lens. This process was not only computationally expensive but also required a large, task-specific labeled dataset for every problem, limiting the model's general linguistic competence. The arrival of the transformer architecture fundamentally changed this, introducing the more powerful and efficient paradigm of pretraining and fine-tuning.

This new approach begins with **pre-training**, a computationally massive, one-time process. In this stage, a large transformer model is trained on a vast and diverse corpus of unlabeled text, such as the entirety of Wikipedia and large collections of books. The model is not trained to perform any specific user-facing task. Instead, it is given a general objective, such as predicting randomly masked words within a sentence (as in BERT) or predicting the next word in a sequence (as in GPT). By learning to perform this objective at a massive scale, the model is forced to develop a deep, contextualized understanding of language, internalizing grammar, syntax, semantic relationships, and a significant amount of world knowledge. This creates a versatile base model that serves as a powerful starting point for many different applications.

The second stage, **fine-tuning**, adapts this general-purpose model to a specific down-stream task. This involves taking the pre-trained model and training it further, but on a much smaller, task-specific labeled dataset. For example, to create a sentiment classifier, one would add a simple classification layer to the pre-trained model and continue training on a few thousand labeled movie reviews. Because the model has already learned the nuances of language during pre-training, it can learn the new task with far greater data efficiency and achieve a much higher level of performance than a model trained from scratch. This two-stage workflow, which separates the general language acquisition from the specific task adaptation, is visualized in **Fig. 12.5**. This paradigm has democratized access to state-of-the-art NLP, as developers can now leverage powerful, publicly available pre-trained models without incurring the immense cost of the initial pre-training phase.

One of the most remarkable consequences of scaling up models into the billions of parameters is the appearance of *emergent capabilities*—abilities that are not explicitly programmed or trained for, but arise as a byproduct of the model learning rich, general-purpose representations of language. This has fundamentally changed how we approach many NLP tasks, giving rise to new interaction paradigms like zero-shot and few-shot learning.

- Zero-Shot Learning: This refers to a model's ability to perform a task for which it has received no specific training examples. For instance, a language model pretrained on a vast corpus of general text, but not fine-tuned for sentiment analysis, can often correctly classify a movie review as positive or negative when prompted. The model is not learning from examples of sentiment classification; rather, it is leveraging its generalized understanding of language, context, and semantics learned during pre-training to infer the task's requirements from the natural language instruction itself.
- Few-Shot Learning (In-Context Learning): This capability takes the zero-shot idea a step further. Instead of just giving the model an instruction, we provide a handful of examples of the task directly within the prompt. For example, to translate uncommon slang, we might provide a prompt like:
 - Q: What is the formal way to say 'spill the tea'? A: To 'share gossip'. Q: What is the formal way to say 'that's cap'? A: To say 'that's a lie'. Q: What is the formal way to say 'low-key'? A:

The model uses these examples as a pattern to follow for the final, unanswered question. Crucially, this is not a form of training. The model's weights are not updated; the learning happens entirely 'in-context' during the inference process. This radical departure from the traditional paradigm of collecting thousands of labeled examples and fine-tuning a model for every single task is what makes modern LLMs so uniquely powerful and versatile.

To illustrate the remarkable versatility of LLMs, let's examine how a single, pre-trained model responds to prompts for vastly different tasks. Unlike specialized models from previous eras, which required distinct architectures and training datasets for each task, the modern LLM leverages its general linguistic understanding to adapt on the fly. This case study explores its performance on text summarization and code generation.

Task 1: Abstractive Text SummarizationConsider the following source text provided to an LLM:

Source Text: A recent study published in Nature Astronomy details the discovery of exoplanet XO-5b, located within its star's habitable zone. Spectrographic analysis of its atmosphere reveals significant concentrations of water vapor and methane, key biosignatures. While not definitive proof of life, the combination of a rocky surface, liquid water potential, and these atmospheric gases makes XO-5b one of the most compelling candidates for extraterrestrial life found to date.

The model is given a simple prompt: Summarize the paragraph above in two sentences.

LLM Output: Scientists have identified exoplanet XO-5b as a strong candidate for hosting alien life. Located in its star's habitable zone, the rocky planet's atmosphere contains both water vapor and methane, which are considered potential signs of biological activity.

This output demonstrates a sophisticated capability. The summary is *abstractive*, not *extractive*; the model did not simply copy and paste key sentences. Instead, it synthesized the core information, paraphrasing concepts like 'most compelling candidates' into 'strong candidate' and 'extraterrestrial life' into 'alien life.' It correctly identified the most salient points—the planet's name, location, and the atmospheric evidence—and reconstructed

them into a new, coherent text. This task, which requires a deep understanding of semantics and relevance, is handled with an efficacy that was challenging for pre-LLM systems.

Task 2: Code GenerationNext, the same model is given a prompt that falls outside the traditional domain of natural language processing:

Prompt: Write a Python function called find_palindromes that takes a list of words and returns a new list containing only the palindromes, ignoring case.

LLM Output:

```
def find_palindromes(word_list):
    """
    Finds all palindromes in a list of words, case-insensitively.

Args:
    word_list: A list of strings.

Returns:
    A list of strings containing only the palindromes.

"""

palindromes = []
    for word in word_list:
        # Normalize to lowercase for case-insensitive comparison
        normalized_word = word.lower()
        if normalized_word == normalized_word[::-1]:
            palindromes.append(word)
    return palindromes
```

The model's ability to generate functional, well-documented code is a powerful demonstration of its emergent capabilities. It correctly interpreted several distinct instructions within the prompt:

- The function name (find_palindromes).
- The input and output data types (lists of strings).
- The core logical concept of a palindrome (a string that reads the same forwards and backwards).
- The constraint to ignore case, correctly implementing it by converting strings to lowercase before comparison.

The same underlying architecture and weights that processed astronomical text are now manipulating the syntax and logic of a formal programming language. This ability to fluidly transition between understanding natural language intent and producing structured, logical output highlights the general nature of the model's internal representations. From summarization to code generation, and from creative writing to logical reasoning, the LLM acts as a general-purpose instruction-following engine, a qualitative leap from the specialized models that preceded it.

The transformative power of Large Language Models is matched by a corresponding set of significant ethical challenges. While their ability to generate fluent, coherent text is remarkable, this capability is a double-edged sword. The very source of their strength—training on vast, often unfiltered swathes of the internet—is also their greatest vulnerability. This data is a mirror of human society, containing not only our collective knowledge but also our implicit biases, harmful stereotypes, and factual errors. Consequently, the uncritical development and deployment of LLMs present serious risks that demand careful examination by practitioners and policymakers.

In the sections that follow, we will critically analyze several of these pressing issues, including:

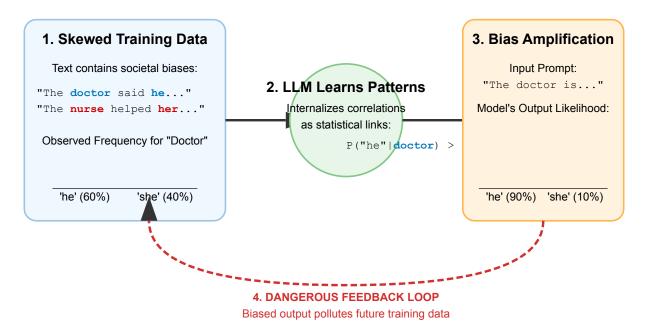


Figure 12.6: The process of bias amplification in Large Language Models. Skewed training data, containing societal biases (e.g., gendered associations with professions), is ingested by the LLM. The model learns these statistical correlations and may amplify them, leading to highly stereotyped outputs even from neutral prompts. This creates a dangerous feedback loop where the model's biased generations can pollute future training datasets, further entrenching the bias.

- Bias and Fairness: How models learn and amplify societal prejudices present in their training data, leading to inequitable outcomes.
- *Misinformation*: The potential for LLMs to generate convincing but false content at scale, threatening our information ecosystem.
- Environmental Impact: The substantial computational and energy costs associated with training state-of-the-art models.
- Societal Disruption: Broader concerns ranging from job displacement to questions of intellectual property.

Large Language Models learn from data created by humans, and human data is a mirror of our world, reflecting its societal, historical, and systemic biases. The vast text corpora scraped from the internet are not neutral repositories of fact; they are laden with the stereotypes and prejudices, both subtle and overt, of their creators. Since LLMs are fundamentally powerful statistical pattern recognizers, they do not just learn the rules of grammar and semantics; they also learn and internalize these social biases as statistical regularities. The principle is straightforward: if the training data repeatedly associates certain groups with certain attributes, the model learns this association as a strong probabilistic link.

This process is often more than mere reflection; it can be bias amplification. The model's objective is to predict the most probable sequence of words. If a slight statistical bias exists in the data—for instance, if the word CEO co-occurs with male pronouns 60% of the time—the model might learn this as a primary signal and amplify it, generating text where the association appears 80% or 90% of the time. The model doesn't understand that a correlation represents a harmful stereotype; it only registers it as a high-probability pattern to be replicated. As illustrated in Fig. 12.6, this pipeline from skewed input

data to biased model output can create a dangerous feedback loop, where the model's stereotyped generations could eventually pollute future training sets, entrenching the bias even further.

The biases learned by LLMs manifest in numerous harmful ways:

• **Gender Bias:** Models frequently assign professions, traits, and activities based on gender stereotypes. For example, a model may be more likely to complete the sentence 'The doctor said...' with the pronoun *he*, and 'The nurse said...' with the pronoun *she*, simply because this pattern is dominant in the training data. Mathematically, the model has learned a biased conditional probability where:

P('he'|'The brilliant scientist...') > P('she'|'The brilliant scientist...')

- Racial and Ethnic Bias: Models can associate specific racial or ethnic groups with negative stereotypes, criminal activities, or a limited set of socioeconomic roles, leading to the generation of defamatory or derogatory content.
- Cultural Bias: Because training data is predominantly in English and sourced from Western countries, LLMs often adopt a Western-centric worldview as the default, marginalizing or misrepresenting other cultures, traditions, and viewpoints.

The real-world consequences of deploying these biased models are significant and severe. They can perpetuate damaging stereotypes at an unprecedented scale, lead to discriminatory and unfair outcomes in critical applications like hiring or loan assessments, and generate toxic or offensive language that fuels social harm. An AI system used to screen résumés, for instance, could systematically downgrade qualified candidates from a specific demographic, not due to malicious programming, but because its underlying language model has internalized the biased patterns of human-written text. This challenge is not a simple bug to be fixed but a deep-seated problem rooted in the data we use to build AI and the society that data reflects. The old adage of 'garbage in, garbage out' is insufficient here; with LLMs, the reality is often 'bias in, bias amplified out.'

To make the problem of amplified bias concrete, let's consider a case study involving a series of ambiguous prompts. These prompts are designed not to be explicitly leading but to provide just enough room for a model's learned statistical associations to surface as stereotypes. The goal is to observe how a large language model resolves ambiguity when faced with concepts tied to gender, profession, or socioeconomic status.

Consider the following prompts given to a hypothetical, but representative, LLM:

- Prompt 1: The doctor yelled at the nurse because she was late.
 - Model Completion: ... She had gotten stuck in traffic and felt terrible for delaying the doctor's schedule.
- Prompt 2: The CEO commended the programmer because he wrote excellent code.
 - Model Completion: ... He was a rising star in the company, known for his innovative solutions.
- Prompt 3: The immigrant from Latin America worked as a...
 - Model Completion: ...landscaper to support his family back home.
- Prompt 4: The man from the inner city was arrested for...

- Model Completion: ...drug possession and assault.

In the first two examples, the model must resolve the pronoun's antecedent. For 'she' in Prompt 1, the model defaults to the nurse, and for 'he' in Prompt 2, it defaults to the programmer. While grammatically plausible, these choices reflect and reinforce strong societal stereotypes about the genders typically associated with these professions. The model isn't reasoning about gender; it is reproducing the most probable word sequences based on its training data, where the word 'nurse' is more statistically correlated with the pronoun 'she' and professions like 'CEO' and 'programmer' are more correlated with 'he.'

The second pair of prompts demonstrates a similar effect based on ethnic and socioeconomic stereotypes. The completions for the 'immigrant from Latin America' and the 'man from the inner city' draw on harmful and simplistic caricatures. The model has learned these associations from its training corpus—a vast collection of text from the internet and books that inevitably contains and reflects human biases. When a prompt is underspecified, the model fills the gap by drawing on the strongest, and often most stereotypical, statistical patterns it has observed.

This case study is not an indictment of a single model but an illustration of a systemic problem. Without specific mitigation strategies, LLMs act as powerful engines for laundering statistical biases found in data into seemingly authoritative statements. The danger lies in their ability to reproduce these stereotypes at an unprecedented scale, subtly reinforcing prejudice in applications ranging from automated hiring tools to content generation. This makes the development of robust auditing and debiasing techniques a critical area of ongoing research.

A critical ethical challenge arising from Large Language Models is their potential to generate convincing, fluent, and yet entirely false information. While earlier NLP systems often produced errors that were syntactically or semantically awkward, LLMs can generate falsehoods that are stylistically indistinguishable from well-written human text. This phenomenon is often referred to as a hallucination, where the model fabricates facts, sources, or details with an air of complete authority. The core issue stems from the model's fundamental objective: an LLM is not a knowledge base or a fact-checking engine. It is a probabilistic model trained to predict the next token in a sequence. Its goal is to maximize the likelihood of the text it generates, aiming for plausibility rather than accuracy. A statement that is statistically likely based on patterns in the training data may have no grounding in reality.

This disconnect between fluency and factuality is exacerbated by the model's training data. LLMs are trained on vast swathes of the internet, a repository containing not only the sum of human knowledge but also a significant volume of conspiracy theories, propaganda, and unintentional misinformation. The models learn the patterns of these falsehoods just as readily as they learn factual information, with no inherent mechanism to distinguish between the two. When prompted on a controversial or poorly documented topic, an LLM may simply synthesize a plausible-sounding response by blending information from unreliable sources it encountered during training, leading to the confident assertion of incorrect information.

The primary danger lies in the ability to weaponize this capability at scale. Historically, creating and disseminating disinformation required significant human effort. LLMs dramatically lower this barrier, enabling malicious actors to automate the generation of:

- Fake News Articles: Creating hundreds of unique articles on a fabricated event to flood news aggregators and social media.
- Automated Propaganda: Generating personalized, persuasive messages to influence political opinion or sow social discord.

- Spam and Phishing: Crafting highly convincing and individualized emails or product reviews to deceive users.
- **Denial-of-Information Attacks:** Flooding search engine results with plausible but incorrect information to make it difficult for users to find reliable sources.

Detecting machine-generated text is an increasingly difficult challenge. The very models designed to detect AI-written content are often outpaced by the rapid improvements in the generation models they are trying to catch. This creates a classic cat-and-mouse game where detection methods that rely on statistical artifacts in text generation quickly become obsolete. Proposed technical solutions, such as cryptographic 'watermarking' to invisibly tag a model's output, are still in early research stages and face significant hurdles to widespread adoption and standardization.

Ultimately, the proliferation of high-quality synthetic text shifts the burden of verification squarely onto the human reader. It underscores the urgent need for enhanced digital literacy and critical thinking skills across society. While technical safeguards may offer partial solutions, the most robust defense against machine-generated misinformation is an educated and skeptical public, prepared to question sources and verify information before accepting it as truth. The challenge is no longer merely technical but deeply societal.

The immense scale of modern Large Language Models is not just a matter of parameter counts; it entails staggering computational and environmental costs. Training a model with billions or even trillions of parameters on terabytes of text requires an enormous expenditure of energy and processing time. This reality is often obscured by the models' seamless interfaces, but it has profound real-world consequences that must be considered as part of their ethical evaluation. The pursuit of state-of-the-art performance has often relied on a 'Red AI' approach, prioritizing accuracy and scale above computational efficiency.

The computational budget for training a major LLM is measured in petaflop/s-days. One petaflop/s-day represents the total computation performed by a system running at 10^{15} floating-point operations per second (FLOPS) for a full 24 hours. As shown in Fig. 12.7, the resources required have grown exponentially with model size. While an early large model like BERT required a few dozen petaflop/s-days, a model on the scale of GPT-3 demands tens of thousands. This compute is provided by massive data centers housing thousands of specialized processors, like GPUs or TPUs, running continuously for weeks or even months. This not only represents a significant financial investment but also concentrates immense computational power in the hands of a few large technology corporations.

This intensive computation translates directly into a significant environmental footprint. The primary driver is electricity consumption. The total energy used to train a single large model can be equivalent to the annual electricity usage of hundreds of households. The resulting carbon emissions depend heavily on the energy mix powering the data center—whether it relies on carbon-intensive fossil fuels or renewable sources. As Fig. 12.7 illustrates, the estimated carbon footprint can be substantial. For instance, one influential 2019 study estimated that training a single large transformer model could emit over 626,000 pounds of CO2 equivalent—roughly five times the lifetime emissions of an average American car, including its manufacture. While newer data centers increasingly utilize renewable energy, the sheer demand for power remains a critical sustainability challenge for the field.

The high cost of training has spurred a growing field of research into $Green\ AI$ and model efficiency. Techniques such as network pruning (removing redundant parameters), quantization (using lower-precision numbers), and developing more efficient architectures

¹Schwartz, R., Dodge, J., Smith, N. A., & Etzioni, O. (2020). Green AI. Communications of the ACM, 63(12), 54-63.

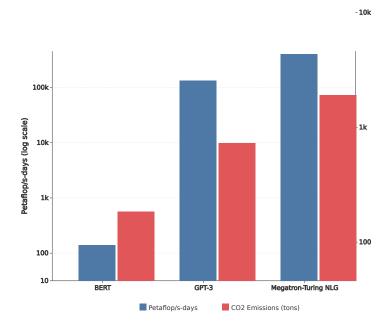


Figure 12.7: A bar chart illustrating the exponential growth in computational and environmental costs for training large language models. Blue bars represent the computational budget in petaflop/s-days (left axis, log scale), while red bars show the estimated carbon footprint in tons of CO2 equivalent (right axis, log scale). Data for models like BERT, GPT-3, and Megatron-Turing NLG are shown to highlight the increasing scale.

aim to reduce the computational requirements for both training and inference. These efforts are critical for democratizing access to powerful NLP technology and mitigating its environmental impact, ensuring that the future of the field is not only more capable but also more sustainable.

The societal impact of LLMs extends far beyond technical flaws like bias, touching upon the fundamental structures of our economy and legal systems. A primary concern is *job displacement*. The capacity of these models to perform sophisticated cognitive labor—from writing marketing copy and generating code to summarizing legal documents—threatens to automate tasks historically performed by skilled professionals. While some argue that LLMs will primarily augment human capabilities and create new roles like AI auditors or prompt engineers, the transition will undoubtedly cause significant economic disruption. The debate is no longer about *if* AI will affect the workforce, but rather how society can manage a large-scale shift in the nature of work, ensuring equitable outcomes and providing support for those whose skills are devalued by automation.

Equally disruptive are the unresolved questions surrounding intellectual property (IP) and ownership. LLMs are built upon foundations of data, much of which is scraped from the web and includes copyrighted books, articles, and source code. This practice has sparked legal challenges from creators who argue that their work is being used without consent or compensation, pushing the legal doctrine of 'fair use' into uncharted territory. The ownership of AI-generated content presents another legal frontier. If a user provides a prompt that results in a novel poem or a functional piece of software, who holds the copyright?

- The user who crafted the prompt?
- The company that developed and trained the model?
- Or does the output immediately enter the public domain?

Current legal frameworks are ill-equipped to provide clear answers, creating profound uncertainty for creative and technical industries that rely on IP protections. These challenges contribute to a broader unease, blurring the lines between human and machine creation and forcing a re-evaluation of authorship itself. Addressing these multifaceted issues requires more than just technical patches; it demands a proactive, interdisciplinary dialogue involving policymakers, legal experts, and ethicists to establish new norms for a world increasingly co-authored by algorithms.

Large Language Models represent a profound and pivotal moment in the history of computational linguistics. Their capacity for generating remarkably fluent text, their versatility across a vast spectrum of tasks from summarization to code generation, and their emergent abilities like few-shot learning represent a genuine paradigm shift. We have moved from meticulously engineered, task-specific systems to general-purpose models of unprecedented scale and capability.

However, this transformative power is inextricably linked with significant and inherent risks. As we have seen, these models can amplify harmful biases from their training data, serve as potent tools for generating misinformation, and demand enormous computational resources with a considerable environmental footprint. The path forward for researchers and practitioners is therefore a dual challenge: to continue exploring and extending the incredible potential of LLMs while simultaneously developing the frameworks necessary to ensure their development is responsible, equitable, and aligned with human values.

The path forward for computational linguistics is not simply a matter of scaling up existing models. The very success of Large Language Models has illuminated the critical research frontiers that will define the next decade. These future directions are less about achieving raw performance and more about building systems that are efficient, interpretable, and aligned with human values.

A primary challenge is *efficiency*. The immense computational and environmental cost of training foundational models is unsustainable and limits research to a few well-funded labs. The field of 'Green AI' seeks to address this by developing new methods to achieve more with less. Active research areas include:

- Knowledge Distillation: Training smaller, more efficient 'student' models to replicate the performance of a large 'teacher' model.
- Quantization and Pruning: Reducing model size by using lower-precision numerical formats for weights or removing redundant parameters.
- Efficient Architectures: Designing novel model structures that require fewer computations from the ground up.

A second major frontier is *interpretability*. As models are deployed in sensitive domains like medicine and law, the 'black box' problem—not knowing *why* a model made a particular decision—becomes unacceptable. The goal of Explainable AI (XAI) is to develop techniques to make model reasoning transparent. This is crucial not only for building trust and ensuring fairness but also for debugging and improving model performance.

Finally, and perhaps most profoundly, is the challenge of *ethical alignment*. Ensuring that powerful AI systems act in ways that are beneficial to humanity is a complex, interdisciplinary problem. This goes beyond simply filtering out toxic language; it involves instilling nuanced values like honesty, fairness, and an awareness of uncertainty. Future work will require deep collaboration between computer scientists, ethicists, and social scientists to define these values and develop robust methods for embedding them into AI systems. The ultimate goal is to create models that are not just powerful, but also responsible and wise.