

# The Book of RAG – and more

Seth J. Chandler, with help from AI

August 21, 2025

# Contents

<b>1 Retrieval Augmented Generation: RAG</b>	<b>1</b>
1.1 Introduction: Beyond the Black Box – Grounding AI in the Rule of Law . . . . .	1
1.2 Deconstructing the RAG Engine: From Text to Meaningful Vectors . . . . .	2
1.2.1 The Art of Representation: Understanding Embedding Models . . . . .	2
1.2.2 The Digital Library: The Role of Vector Databases . . . . .	3
1.2.3 Measuring Relevance: Distance Metrics in Plain English . . . . .	3
1.3 RAG in Practice: A Paradigm Shift for Legal Services . . . . .	4
1.3.1 The Core Value Proposition: Grounding, Citing, and Trusting . . . . .	5
1.3.2 Transforming the Workflow: Key Legal Use Cases . . . . .	5
1.4 The Devil in the Details: Nuances of a High-Performance Legal RAG System . . . . .	6
1.4.1 The Art of the Chunk: Why Document Segmentation Matters . . . . .	6
1.4.2 The Retriever’s Dilemma: How Many Sources are Enough? (The ”Top-K” Problem) .	7
1.4.3 The Embedding Bottleneck: Not All Embeddings Are Created Equal . . . . .	7
1.5 The Frontier of RAG: Advanced and Hybrid Techniques . . . . .	8
1.5.1 Beyond Naive Retrieval: Hybrid Search and Re-ranking . . . . .	8
1.5.2 Semantic Enrichment Through Metadata and Pre-processing . . . . .	9
1.5.3 Fine-Grained Relevance with Multidimensional Embeddings . . . . .	10
1.5.4 Preserving Integrity with Hierarchical Structures . . . . .	10
1.6 The Law Student as Builder: Creating Your First RAG System . . . . .	10
1.6.1 Using OpenAI’s Custom GPTs . . . . .	10
1.6.2 Leveraging Anthropic’s Claude . . . . .	11
1.6.3 Exploring Google’s Gemini . . . . .	11
1.7 RAG in the Broader AI Ecosystem: A Comparative Analysis . . . . .	12
1.7.1 RAG vs. Model-Centric Prompting (MCP): Knowledge vs. Action . . . . .	12
1.7.2 RAG vs. Ultra-Large Context Windows: The ”RAG is Dead” Debate . . . . .	13
1.8 Conclusion: The Future-Ready Lawyer in the Age of Augmented Intelligence . . . . .	15

<b>2</b>	<b>Custom AI for the Modern Lawyer</b>	<b>21</b>
2.1	Defining the Tools: Gems, CustomGPTs, and Projects . . . . .	21
2.2	Strategic Choices: When to Use Each Platform . . . . .	22
2.3	Use Cases: Taming the Cumbersome in Law . . . . .	22
2.3.1	Example 1: The Case Brief Generator . . . . .	22
2.3.2	Example 2: The Discovery Drafter . . . . .	23
2.3.3	Example 3: The Bluebook Citation Checker . . . . .	23
2.4	Creating Well-Aligned Custom Assistants . . . . .	24
2.5	Additional High-Value Use Cases . . . . .	24
2.6	The Next Frontier: Agents vs. Custom AI . . . . .	25
<b>3</b>	<b>What is Prompt Caching?</b>	<b>26</b>
3.1	Primary Benefits . . . . .	26
3.2	Applicability: API vs. Web Interface . . . . .	26
3.3	Use Case: Legal Document Analysis . . . . .	27
3.3.1	Scenario . . . . .	27
3.3.2	Implementation with Caching . . . . .	27
3.3.3	Conclusion . . . . .	27



# Chapter 1

## Retrieval Augmented Generation: RAG

### 1.1 Introduction: Beyond the Black Box – Grounding AI in the Rule of Law

The legal profession stands at the precipice of a technological revolution, driven by the advent of Large Language Models (LLMs)—powerful artificial intelligence systems capable of generating sophisticated, human-like text. Yet, this promise is shadowed by a fundamental tension. The very nature of these generative models, which operate on probabilistic principles to predict the next most likely word, is often at odds with the stringent demands of legal practice: absolute accuracy, unwavering verifiability, and strict fidelity to authoritative sources. This tension manifests in the now-infamous phenomenon of “hallucinations,” where an LLM confidently fabricates information, from subtle misinterpretations to the outright invention of case law [1]. For a lawyer, citing a non-existent case is not a mere technical glitch; it is a sanctionable offense, a breach of professional ethics, and a critical risk that can undermine a case and a career. Ungrounded AI, in its raw form, is a professional minefield.

The solution to this critical challenge is not to abandon the technology but to re-architect it for the unique demands of the legal domain. This is the purpose of Retrieval-Augmented Generation, or RAG. RAG is not a single piece of software but an architectural pattern designed to impose discipline on a creative, probabilistic system [3]. It fundamentally transforms an LLM’s operational paradigm. Instead of functioning as a closed-book exam-taker, relying solely on the vast but static and sometimes outdated information it was trained on, a RAG system operates like an open-book expert. Before generating any response, it first retrieves relevant, timely information from an authoritative, pre-determined knowledge base—be it a curated legal database like Westlaw, a firm’s internal document management system, or a specific set of statutes and regulations [4]. RAG can be thought of as a form of “context engineering”: determining what the Large Language Model sees and relies on in addition to the instructions from the user in generating a response.

This process of “grounding” the LLM’s output in verifiable facts is the key to unlocking its potential as a trustworthy tool for lawyers [3]. The probabilistic nature of an LLM means it excels at synthesis, summary, and drafting, but it lacks an inherent concept of truth. The legal profession, conversely, is built upon a foundation of determinism and verifiability; every argument must be traceable to a specific source of authority. RAG creates a crucial bridge between these two worlds. It compels the probabilistic generator to condition its creative output on a deterministic, retrieved set of facts [4]. In doing so, RAG imposes a legalistic requirement for evidence onto the AI, forcing it to, in essence, “show its work.” This architectural shift is what makes it possible to harness the power of generative AI while upholding the rigorous standards of the rule of law.

## 1.2 Deconstructing the RAG Engine: From Text to Meaningful Vectors

To appreciate the transformative power of RAG, it is essential to look under the hood and understand its core mechanical components. At first glance, the process by which a machine "reads" a legal brief and "understands" a complex query may seem like magic. In reality, it is a logical, multi-stage process of mathematical transformation and comparison. By demystifying these components—embedding models, vector databases, and distance metrics—we can move from being passive users of AI to informed architects of its application in law.

### 1.2.1 The Art of Representation: Understanding Embedding Models

The first and most fundamental step in any RAG system is translating human language into a format a computer can process. Computers do not understand words; they understand numbers. Embedding models are the sophisticated translation engines that perform this critical task [7]. They take a piece of text—a word, a sentence, a paragraph from a judicial opinion—and convert it into a high-dimensional numerical vector, which is simply an ordered list of numbers [10].

Early methods for this were quite literal, often relying on counting keyword occurrences, a technique known as "sparse" representation [12]. While simple, these methods failed to capture the rich semantic nuance of language. Modern "dense" embedding models, often built on complex neural network architectures, are far more powerful. They are trained on vast quantities of text to learn the subtle relationships and contexts in which words appear [7]. The core principle that makes these models so effective is that texts with similar meanings will be mapped to vectors that are mathematically close to one another in a high-dimensional "semantic space" [13]. The word "contract" will have a vector representation closer to "agreement" than to "tort." It will also be the case that the embedding for King minus the embedding for Queen will be similar to the embedding for man minus the embedding for woman.

To make this abstract concept concrete, imagine a simplified, two-dimensional "legal concept space." Here are three short, related, but distinct legal sentences:

1. The plaintiff filed a motion for summary judgment.
2. The defendant's attorney submitted a responsive pleading to the motion for a new trial.
3. The court will hear arguments on the motion to dismiss for lack of personal jurisdiction.

In the context of a language model, these sentences would be converted into high-dimensional vectors, or embeddings, that capture their meaning and semantic relationships [1]. To visualize this in two dimensions, one might imagine the sentences positioned on a 2D plane based on their similarities and differences.

- Axis 1 (Phase): The "early" phase of a case could be represented by a value closer to -1, while a "later" or more advanced phase is closer to 1.
- Axis 2 (Actor): A plaintiff or their action could be represented by a positive value, a defendant by a negative value, and a neutral party like the court by a value close to zero.

Based on this, a 2D embedding space might look something like this:

- Sentence 1: Positioned at (-0.6, 0.9) because it is at a fairly early stage but done by the plaintiff.
- Sentence 2: Positioned at (0.9, -0.8) because it is at a late stage but done by the defendant.
- Sentence 3: Positioned at (-0.9, 0) because it is early but done by a neutral court.

This visualization illustrates how a language model uses numerical vectors to capture the complex relationships—in this case, semantic and procedural—between pieces of text, enabling it to understand context and meaning. Just keep in mind that in the real world (a) the embeddings may have thousands of dimensions; 2 is definitely not enough and (b) the embeddings are learned, not predetermined, through a complex process involving neural networks, backpropagation and a variety of clever algorithms. Embedding models that you may encounter are likely to include the old-but-still-used BERT (and its many variants such as Modern-BERT and Sentence-BERT), Voyage-3-large, NV-Embed-v2, OpenAI models, and gte-Qwen2-7B-instruct. The **voyage-law-2 model** is touted as being useful for embedding legal documents.

### 1.2.2 The Digital Library: The Role of Vector Databases

Once every document, or chunk of a document, in a legal corpus has been converted into a vector, these numerical representations must be stored in a specialized repository designed for efficient retrieval. This is the function of a vector database [16]. Companies like Pinecone, Zilliz (which commercializes Milvus), Weaviate, and Qdrant make their living in this business. Other tech giants such as MongoDB may also offer vector databases as one of their products.

A traditional relational database, which stores data in rows and columns, is optimized for finding exact matches based on keywords. This is insufficient for the conceptual search required by RAG. A vector database, by contrast, is purpose-built to store and index high-dimensional vectors and to perform one task exceptionally well: finding the “nearest neighbors” to a given query vector based on their mathematical proximity [16].

Searching through millions or even billions of high-dimensional vectors to find the closest matches for every query would be computationally intractable if done naively. To solve this, vector databases employ sophisticated indexing algorithms. One of the most common and effective is the Hierarchical Navigable Small World (HNSW) algorithm [18]. An analogy helps to understand how HNSW works. Imagine trying to find a specific book in a massive, unorganized library. A brute-force search would require you to check every single book on every shelf—an impossibly slow process. HNSW is like creating a multi-layered map of that library.

1. *The Top Layer (The Neighborhood Map)*: The algorithm first creates a sparse, high-level graph connecting distant but representative “landmark” documents. When a query comes in, the system starts here, quickly navigating across the entire library to find the general “neighborhood” of conceptually similar documents.
2. *The Middle Layers (The Street Map)*: As the search zooms in, it moves to denser layers of the graph, which are like street maps connecting different sections or aisles within that neighborhood.
3. *The Bottom Layer (The Shelf Map)*: Finally, at the lowest and densest layer, the algorithm performs a detailed search among a small, highly relevant set of documents on a specific “shelf” to find the exact nearest neighbors.

By using this hierarchical approach, HNSW and similar algorithms can perform incredibly fast and accurate similarity searches over vast datasets without having to compare the query vector to every single vector in the database [18].

### 1.2.3 Measuring Relevance: Distance Metrics in Plain English

The final piece of the retrieval puzzle is the specific mathematical formula used to calculate the “distance” or “similarity” between two vectors. This is known as a distance metric [21]. While many exist, two are particularly common in text analysis:

- *Manhattan Distance*: Also known as “city block” or “taxicab” distance, this is one of the most intuitive metrics. It calculates the distance a taxi would have to travel on a grid-like city street plan to get

from point A to point B. Mathematically, it is simply the sum of the absolute differences of the vector coordinates along each dimension [21]. The formula for two vectors,  $x$  and  $y$ , in an  $n$ -dimensional space is:

$$d_1(x, y) = \sum_{i=1}^n |x_i - y_i|$$

- *Cosine Similarity*: This metric is often preferred for text analysis because it focuses on the *orientation* of the vectors rather than their magnitude (or length). In natural language, the length of a document doesn't necessarily correlate with its relevance. A concise, on-point paragraph might be more relevant than a long, rambling chapter. Cosine similarity measures the cosine of the angle between two high-dimensional vectors. If the vectors point in the exact same direction, the angle is 0 degrees, and the cosine similarity is 1 (maximum similarity). If they are orthogonal (90 degrees apart), the similarity is 0. If they point in opposite directions (180 degrees), the similarity is -1 [22].

Let's return to our simplified 2D legal concept space to see this in action with a tangible calculation. A law student poses a query: *"Find cases in which the plaintiff's attorney was sanctioned for filing a frivolous motion to dismiss."* The system embeds this query into a vector using the same system employed to embed the sentences in the main text. Perhaps our query would land at the coordinates  $(-0.85, 0.6)$  legal space. The vector database then measures the distance from the query to the vectors in the database. Here I use a simple so-called "Manhattan Distance" metric. In the real world, the vector database would likely use a more sophisticated metric such as the "cosine distance" which addresses the angle between multidimensional vectors.

$$\begin{aligned} &|-0.85 - (-0.6)| + |0.6 - 0.9| = |-0.25| + |-0.3| = 0.25 + 0.3 = 0.55 \\ &|-0.85 - 0.9| + |0.6 - (-0.8)| = |-1.75| + |1.4| = 1.75 + 1.4 = 3.15 \\ &|-0.85 - (-0.9)| + |0.6 - 0| = |0.05| + |0.6| = 0.05 + 0.6 = 0.65 \end{aligned}$$

Based on these computations, the system would be more likely to retrieve the first sentence (*The plaintiff filed a motion for summary judgment.*) than the third (*The court will hear arguments on the motion to dismiss for lack of personal jurisdiction.*) And use it to augment the context by which the large language model responds to the query.

The ultimate utility of a RAG system in a legal context is not a matter of chance; it is the end result of a precise causal chain. The quality of a final generated legal memo depends directly on the relevance of the context provided to the LLM [6]. The relevance of that context is determined entirely by the accuracy of the nearest neighbor search performed by the vector database [16]. That search accuracy, in turn, is a function of how well the vector representations capture the true semantic meaning of the underlying legal documents [7]. Finally, this semantic fidelity is dictated by the choice and quality of the embedding model used at the very beginning of the process [26]. Therefore, a seemingly arcane technical decision made by a system's architect—which embedding model to use—has a profound and direct downstream effect on whether the final output is a reliable piece of legal work product or a fatally flawed and professionally dangerous document. This establishes an unbreakable link between a technical choice and a professional-legal outcome.

### 1.3 RAG in Practice: A Paradigm Shift for Legal Services

Understanding the technical architecture of RAG is the foundation for appreciating its practical impact. The transition from theory to practice reveals RAG not merely as a technical improvement but as a catalyst for a paradigm shift in how legal services are researched, analyzed, and delivered. By addressing the core requirements of the legal profession—accuracy, verifiability, and confidentiality—RAG enables a new class of applications that can transform the daily workflows of legal professionals.



### 1.3.1 The Core Value Proposition: Grounding, Citing, and Trusting

At its heart, RAG's value proposition for the legal field rests on three pillars that directly counteract the inherent weaknesses of standalone LLMs.

- *Mitigating Hallucinations:* The single greatest barrier to the adoption of generative AI in law is its propensity to hallucinate. RAG's primary function is to provide a factual "ground" for the LLM's generation process. By forcing the model to synthesize its answer from a specific, provided set of retrieved documents, the system's operation shifts from an act of pure recall from a static memory to an act of comprehension and summary based on authoritative evidence [3]. This is analogous to the difference between a closed-book and an open-book exam. In the open-book scenario, the risk of fabricating case citations or misstating legal standards is dramatically reduced, as the model is constrained to the provided texts [5].
- *Creating an Audit Trail:* For a legal argument to have any weight, its premises must be verifiable. A critical feature of well-designed RAG systems is their ability to provide citations for the information they generate [1]. When a RAG-powered tool summarizes a point of law, it can and should link back to the specific paragraph in the judicial opinion or the specific section of the statute from which that information was derived. This creates a transparent and auditable trail from the final work product back to the primary source material. For a lawyer, this is not a "nice-to-have" feature; it is a requirement for professional diligence and ethical practice. It allows for immediate human verification, building the trust necessary to integrate these tools into high-stakes workflows.
- *Confidentiality and Security:* Law firms are custodians of highly sensitive and privileged client information. Using public, third-party AI tools for substantive legal work raises profound confidentiality concerns, as queries and documents may be logged or used for future model training, potentially violating attorney-client privilege [1]. RAG architecture provides a powerful solution to this problem. A firm can deploy a RAG system that operates entirely within its own secure, private infrastructure. The knowledge base can consist of the firm's proprietary data—internal research memos, client files, deposition transcripts, and previous work product—allowing lawyers to leverage this internal knowledge without ever sending sensitive information to an external server [1]. The LLM can be accessed via a secure API or even run locally, ensuring that all data remains within the firm's control.

### 1.3.2 Transforming the Workflow: Key Legal Use Cases

The practical applications of a grounded, verifiable, and secure AI system are vast and poised to reshape core legal tasks.

- *Legal Research:* Traditionally, legal research has been a keyword-driven process, requiring lawyers to guess the precise terms a judge might have used in a relevant opinion. RAG transforms this into a conceptual, semantic search. A lawyer can now ask a complex, natural-language question such as, "What are the key precedents in the Ninth Circuit regarding the 'inevitable disclosure' doctrine as applied to software engineers leaving a company to join a direct competitor?" A RAG system does not just look for these keywords. It understands the *meaning* behind the query and retrieves conceptually related documents, including case law that discusses the doctrine, relevant statutes on trade secrets, and, crucially, any internal firm memos or briefs that have previously analyzed this exact issue [30]. This leads to faster, more comprehensive, and more on-point research results.
- *Document Review & eDiscovery:* The process of eDiscovery in litigation often involves manually reviewing hundreds of thousands or even millions of documents to find relevant evidence—a time-consuming and expensive task. RAG can automate and accelerate this process with remarkable accuracy. An attorney can issue a natural language instruction like, "Find all communications from Jane Doe to John Smith discussing 'Project Titan' between May 1, 2023, and July 31, 2023, that express a negative sentiment or concern about project delays." The system can semantically search across a vast

repository of emails, chat logs, and internal documents, retrieving not just exact keyword matches but communications that are conceptually related to the query, and even analyze their sentiment [30].

- *Contract Analysis and Drafting:* RAG systems can be configured with a law firm’s or a corporate legal department’s entire repository of contracts and clause playbooks as their knowledge base. When presented with a third-party contract, the system can perform a rapid analysis, comparing its clauses against the firm’s preferred language and risk tolerance. It can flag non-standard indemnity clauses, identify missing provisions required by a specific regulation, or summarize key obligations and deadlines [2]. Furthermore, it can assist in drafting by suggesting alternative clauses drawn directly from the firm’s approved templates, ensuring consistency and adherence to best practices.

The ability of RAG to connect to a firm’s private, internal document repositories unlocks one of its most profound long-term benefits: the preservation and scaling of institutional knowledge [1]. A law firm’s most valuable asset is the cumulative experience of its attorneys, yet this knowledge is often tacit, stored in the minds of senior partners or scattered across decades of disparate memos, emails, and briefs. When a senior partner retires, their unique expertise often walks out the door, forcing junior associates to “reinvent the wheel” on complex legal issues. A RAG system, by ingesting and indexing the firm’s entire historical work product, effectively digitizes this institutional memory. A junior associate facing a novel issue can then query the system, asking, “*How has our firm previously argued against a motion to dismiss for lack of personal jurisdiction in California federal court for a European client?*” In seconds, the system can retrieve relevant past briefs, internal strategy memos, and email discussions on the topic, instantly providing the associate with access to decades of the firm’s collective wisdom. This not only creates massive efficiencies but also functions as a powerful, on-demand training tool, democratizing and scaling senior-level expertise across the entire organization.

## 1.4 The Devil in the Details: Nuances of a High-Performance Legal RAG System

While the conceptual framework of RAG is powerful, its practical implementation is fraught with nuance. A naive or poorly configured RAG system can produce results that are no better, and sometimes worse, than a standalone LLM. The difference between a transformative legal tool and a frustratingly inaccurate one lies in a series of critical design choices that must be carefully considered, especially given the precision required for legal work. Designing a better RAG system is a big and valuable business.

### 1.4.1 The Art of the Chunk: Why Document Segmentation Matters

Legal documents—from lengthy judicial opinions to complex multi-part statutes—are often too large to be converted into a single vector embedding. They must first be broken down into smaller, manageable “chunks” of text before being processed by the embedding model [26]. The strategy used for this segmentation, or “chunking,” is a critical and often overlooked factor that profoundly impacts the quality of retrieval.

Several chunking strategies exist, each with significant trade-offs for legal documents:

- *Fixed-Size Chunking:* This is the simplest method, where a document is split into chunks of a fixed number of characters or tokens (e.g., every 500 tokens). While easy to implement, this approach is highly problematic for legal text. It pays no attention to the logical structure of the document, often slicing sentences, paragraphs, or even critical clauses in half. This can completely destroy the context and meaning of the text, leading to the retrieval of incoherent or misleading information [34].
- *Sentence-Based or Recursive Chunking:* A more sophisticated approach involves splitting the text along natural boundaries, such as sentences or paragraphs. Recursive chunking attempts to split by paragraphs first, then by sentences, and so on, respecting the grammatical and logical structure of the

text. This is a significant improvement over fixed-size chunking as it keeps coherent thoughts together, which is essential for preserving the meaning of legal arguments and contractual provisions [34].

- *Semantic Chunking*: This is the most advanced and often most effective strategy for complex documents. Instead of relying on fixed rules or punctuation, semantic chunking uses an embedding model to analyze the relationships between sentences. It groups sentences that are semantically related into a single, coherent chunk, even if they are not perfectly contiguous. For example, it could group all sentences related to the "indemnification" clause in a contract together, regardless of their exact position. This method, though computationally expensive, is ideal for legal documents where a single concept may be discussed across multiple paragraphs, ensuring that the retrieved context is conceptually complete [34].

The choice of chunking strategy is not a minor technical detail. A poorly chunked contract that separates a key definition from the clause that uses it can lead to a catastrophic misinterpretation by the LLM. For legal applications, strategies that preserve semantic context are paramount.

### 1.4.2 The Retriever's Dilemma: How Many Sources are Enough? (The "Top-K" Problem)

When a query is made, the RAG system's retriever searches the vector database and returns the 'k' most similar chunks to be passed to the LLM. This parameter, 'k', which determines the number of documents to retrieve, presents a critical dilemma with no single correct answer.

The trade-off is clear:

- If 'k' is too low (e.g., retrieving only the top 1 or 2 chunks), the system might miss crucial context. The single most similar chunk may not contain the complete answer, and important nuance or a critical exception might be located in the third or fourth most relevant chunk. This can lead to answers that are correct but incomplete.
- If 'k' is too high (e.g., retrieving 20 or 50 chunks), the system risks introducing noise and irrelevant information into the context provided to the LLM. This can confuse the model, dilute the relevance of the most important chunks, and lead to less focused and potentially inaccurate answers. Furthermore, sending a large number of chunks increases the number of tokens processed by the LLM, which directly increases both the monetary cost and the time (latency) required to generate a response [37].

The optimal value for 'k' is highly dependent on the specific task and the nature of the documents. For answering a simple factual question, a small 'k' might suffice. For generating a complex legal memorandum that requires synthesizing multiple sources, a larger 'k' may be necessary. Determining the right number often requires empirical testing and tuning to find the balance that maximizes relevance without introducing excessive noise [37].

### 1.4.3 The Embedding Bottleneck: Not All Embeddings Are Created Equal

As established earlier, the entire RAG pipeline is causally dependent on the quality of the vector embeddings. A critical nuance here is that the choice of the embedding model itself can be a major point of failure. Most general-purpose embedding models are trained on a vast corpus of general internet text, like Wikipedia and web crawls. While this makes them broadly knowledgeable, they may lack a deep understanding of the highly specialized and nuanced terminology of the law [40].

For example, a generic model might see the terms "liability," "indemnity," and "guarantee" as broadly similar concepts related to financial responsibility. However, a lawyer knows that these terms have precise and distinct legal meanings with vastly different implications. An embedding model that fails to capture

these distinctions will create a flawed semantic map of the legal domain, leading to inaccurate retrieval. This has given rise to the development of domain-specific embedding models that are pre-trained or fine-tuned on massive corpora of legal, financial, or medical texts. For high-stakes legal work, using a generic embedding model represents a significant and often unacceptable risk. The most advanced “Legal-RAG” systems place a strong emphasis on this domain-specific tuning to ensure that the foundational vector representations accurately reflect the intricate semantics of the law [42].

There are neural networks designed specifically for embedding legal information such as Legal-BERT, but users of these embedders need to examine precisely the data on which it was trained. It may, for example, have emphasized European law where the resulting embeddings may not map well onto American understandings. The well-known computer science principle of “Garbage In, Garbage Out” (GIGO) is not just applicable but is dangerously amplified in RAG systems. A standard LLM’s failure mode is a hallucination, which a discerning user might identify as suspicious. A RAG system, however, introduces new and more insidious potential points of failure earlier in the pipeline. If the source documents in the knowledge base are themselves outdated, inaccurate, or incomplete, the RAG system will confidently and authoritatively generate answers based on this flawed information. Worse, it will provide citations to these incorrect sources, lending a false veneer of credibility to the output [6]. Similarly, if the chunking strategy is poor or the embedding model is ill-suited for legal text, the system will fail to retrieve the correct information even if it exists in the knowledge base. This shifts the burden of diligence for the legal professional. It is no longer enough to simply evaluate the final generated text; one must also scrutinize the quality and currency of the underlying knowledge base and validate the integrity of the entire retrieval pipeline.

## 1.5 The Frontier of RAG: Advanced and Hybrid Techniques

The basic RAG architecture—retrieve then generate—is a powerful foundation, but the field of AI is evolving at a breakneck pace. Researchers and engineers are continuously developing more sophisticated techniques to enhance the precision, relevance, and reasoning capabilities of these systems. For law students entering the profession, understanding these advanced concepts is key to appreciating the state-of-the-art in legal technology and recognizing the limitations of simpler implementations.

### 1.5.1 Beyond Naive Retrieval: Hybrid Search and Re-ranking

The initial retrieval step in a basic RAG system is often a compromise between speed and accuracy. Advanced techniques introduce additional layers to refine this process, ensuring the context provided to the LLM is of the highest possible quality.

- *Hybrid Search:* While pure semantic (vector) search is excellent for understanding the conceptual meaning of a query, it can sometimes struggle with queries that depend on specific, literal keywords. For instance, a query for “*cases citing Marbury v. Madison*” requires an exact match for the case name, which a purely semantic search might miss if it prioritizes conceptual similarity over lexical precision. Hybrid search solves this by combining the best of both worlds. It simultaneously performs a traditional keyword-based search (using algorithms like BM25, which are excellent at finding exact terms) and a semantic vector search. The results from both searches are then intelligently merged to produce a final ranked list that benefits from both keyword precision and conceptual relevance [3]. This is particularly powerful in the legal domain, where research often involves a mix of specific entities (case names, statute numbers, party names) and broad legal concepts [46].
- *Re-ranking:* The first retrieval stage in a RAG system is designed to be fast, casting a wide net to recall a set of potentially relevant documents (e.g., the top 50 chunks). However, this initial ranking may not be perfect. A re-ranker introduces a second, more computationally intensive quality control step. This component takes the initial list of retrieved chunks and uses a more powerful model (often a “cross-encoder” that jointly analyzes the query and each document in detail) to re-evaluate and

re-order them based on a more nuanced understanding of relevance. The re-ranker's job is to push the absolute best, most on-point chunks to the very top of the list before they are sent to the LLM for generation [3]. This two-pass system—a fast, broad recall step followed by a slower, high-precision re-ranking step—significantly improves the quality of the context and, consequently, the final generated answer [48].

The evolution from a simple retrieve-and-generate model to an advanced pipeline incorporating hybrid search and re-ranking is not merely an arbitrary technical progression. It is an algorithmic reflection of the established and effective cognitive workflow of a human lawyer. A junior associate conducting research often begins with a broad, conceptual search to identify a pool of potentially relevant cases (analogous to vector retrieval). They then refine this search using specific keywords, party names, or citations to narrow the results (hybrid search). From this collected set of documents, they do not simply read them in the order the database presented them; they critically evaluate each one, prioritizing the most authoritative, on-point, and jurisdictionally relevant cases (re-ranking). Only after this rigorous filtering and prioritization process do they begin to synthesize the information into a legal argument (generation). Thus, the architecture of an advanced RAG system mirrors the methodical and proven process of human legal research and analysis.

### Frontier methods

While the architectural components described thus far form the bedrock of a functional RAG system, the frontier of development is focused on elevating the quality and precision of the retrieval stage itself. The ultimate efficacy of the generative model is wholly dependent on the relevance and completeness of the context it receives. Consequently, advanced strategies have emerged that move beyond simple vector similarity to incorporate deeper layers of semantic understanding and structural awareness into the retrieval process. For the legal field, where nuance and context are paramount, these techniques represent a critical evolution toward building truly reliable AI systems.

#### 1.5.2 Semantic Enrichment Through Metadata and Pre-processing

A foundational limitation of naive retrieval is that it operates on the raw text of a document corpus. This approach treats every word as equally important and can be misled by superficial keyword matches. Advanced RAG systems address this through a pre-processing step of semantic enrichment, where structured metadata is generated and associated with each document chunk. This is data about the data, designed to provide a higher-level, more abstract understanding of the underlying text. Two powerful forms of this metadata include:

- **Generated Summaries:** Before the embedding process, an auxiliary language model can be used to generate a concise summary for each document, section, or even paragraph. The RAG system can then be configured to perform an initial search against this corpus of summaries. This allows the retriever to first identify chunks that are thematically and conceptually aligned with the query's intent, rather than just lexically similar. It is analogous to a lawyer first reviewing the headnotes of a case to grasp its core holdings before diving into the full text.
- **Hypothetical Question-Answer Pairs:** Another sophisticated technique involves pre-generating a set of potential questions that each text chunk could answer. This creates a list of question-answer pairs that are tightly bound to the source text. When a user's query closely matches one of these pre-generated questions, it serves as a powerful signal of that chunk's relevance. This method effectively anticipates user intent and maps it directly to the most pertinent segments of the knowledge base.

By searching over these enriched metadata layers using queries that have been similarly enriched, the retrieval process becomes less about matching strings of text and more about aligning with the core meaning and purpose of the source documents, leading to a significant increase in the precision of the context provided to the LLM.

### 1.5.3 Fine-Grained Relevance with Multidimensional Embeddings

The standard practice of representing a text chunk as a single vector—a single point in semantic space—is a powerful abstraction, but it is also a lossy one. A dense paragraph of legal analysis may contain multiple distinct concepts, arguments, and citations. Compressing this multifaceted information into a single vector can obscure important details and fail to capture the full richness of the text. To overcome this, advanced models like ColBERT (Contextualized Late Interaction over BERT) employ a multidimensional or token-level embedding strategy. Instead of creating one vector for an entire chunk, these models generate a vector for every single token (or word) within the chunk. This transforms the representation of a document from a single point into a matrix of vectors—a fine-grained semantic map of the text. (The mathematicians among you will be fascinated to learn that dimensionality reduction techniques are used to reduce the number of columns in the matrix). The retrieval process is likewise more nuanced. The user’s query is also converted into a matrix of token-level vectors. The system then performs a “late interaction” comparison, meticulously scoring the similarity between each query vector and all of the document vectors. This method allows the system to identify multiple points of conceptual overlap and understand how different parts of a query relate to different parts of a document. This concept is reminiscent of why transformers (the backbone of large language models) use multihead attention heads instead of a single attention head. For a complex legal query with multiple sub-issues, this granular approach is vastly superior to single-vector comparison, as it can identify documents that address the full constellation of concepts in the query, not just the dominant one.

### 1.5.4 Preserving Integrity with Hierarchical Structures

A significant risk in RAG is the retrieval of isolated facts that are misleading or incorrect when stripped of their surrounding context. A single sentence from a judicial opinion might appear to support an argument, but the preceding and succeeding sentences could entirely reverse its meaning or limit its applicability. Hierarchical retrieval addresses this by explicitly indexing the logical structure of the source documents. Text chunks are not treated as a flat, undifferentiated list but are organized into a parent-child hierarchy. A specific sentence chunk is linked to its parent paragraph, which is in turn linked to its parent section or chapter. This structure enables a multi-stage retrieval process. The system can first perform a precise search to identify the most relevant sentence-level chunks (the needles”). However, because of the indexed hierarchy, the system can then automatically retrieve the parent chunks (the haystack context”) that provide the necessary framing. This ensures that the LLM receives not only the specific, targeted fact but also the broader context required for a complete and accurate synthesis. For legal applications, where a single clause can be rendered meaningless without its defining section, this ability to preserve contextual integrity is not merely an enhancement; it is a requirement for responsible and ethical use.

## 1.6 The Law Student as Builder: Creating Your First RAG System

The concepts behind RAG can feel abstract, but the technology has become remarkably accessible. For law students, building a simple, personal RAG system is not only possible but is also an invaluable exercise. Using no-code or low-code platforms provided by major AI companies, you can transform your own course materials into a queryable knowledge base, making the technology tangible and immediately useful.

### 1.6.1 Using OpenAI’s Custom GPTs

One of the most straightforward ways to build a RAG system is through OpenAI’s Custom GPTs feature, available to ChatGPT Plus subscribers. The “Knowledge” feature within the GPT builder is, in effect, a fully managed, user-friendly RAG system [51].



A simple walkthrough to create a personal study assistant is as follows:

1. *Access the GPT Builder:* Navigate to the "Explore GPTs" section in ChatGPT and select "Create a GPT."
2. *Configure Your GPT:* In the "Configure" tab, give your GPT a name (e.g., "Contracts Law Assistant") and a brief description.
3. *Upload Your Knowledge Base:* This is the crucial step. Under the "Knowledge" section, click "Upload files." Select the documents you want your assistant to be an expert on. This could include PDFs of your class notes, case briefs you've written, course outlines, and relevant articles.
4. *Enable Retrieval:* Ensure the "Knowledge" capability is checked. This tells the GPT that it is allowed to retrieve information from the files you've uploaded.
5. *Set the Instructions:* In the "Instructions" box, it is critical to provide a clear directive to ground the model. Include a sentence like: *"You are an expert legal assistant. When answering questions, you must base your answers exclusively on the information contained in the uploaded documents. If the answer cannot be found in the documents, state that clearly."*
6. *Test and Refine:* You can now interact with your custom GPT in the preview pane. Ask it questions about your course material and see how it performs.

Behind the scenes, OpenAI is automatically performing the entire RAG pipeline: chunking your documents, creating embeddings, storing them in a vector store, and performing a semantic search to retrieve relevant chunks whenever you ask a question [51]. You have successfully built a RAG-powered study assistant without writing a single line of code.

### 1.6.2 Leveraging Anthropic's Claude

Anthropic's Claude platform offers a similar, powerful feature through its "Projects" functionality. A project can serve as a massive, RAG-enabled knowledge base that allows you to converse with a large collection of documents simultaneously [53].

A walkthrough for Claude is as follows:

1. *Create a Project:* Within the Claude interface (available on paid plans), create a new "Project."
2. *Upload Files:* Add all relevant files to the project. Claude's projects can handle a significantly larger amount of content than a standard context window, effectively creating a substantial knowledge base.
3. *Interact with Your Knowledge:* When you start a conversation within that project, Claude can access all the uploaded documents. As the amount of information in the project grows and exceeds the standard context window, Claude automatically enables a RAG-like "project knowledge search tool" to intelligently find and retrieve the most relevant information from your documents in response to your queries [53].

This allows you to have a continuous conversation with your entire semester's worth of readings, asking complex questions that require synthesizing information across multiple sources.

### 1.6.3 Exploring Google's Gemini

Google's Gemini models also offer capabilities that can be used for RAG-like workflows, often by combining their ultra-large context windows with powerful search functionalities [3]. The simplest way to access this capability is to use Gemini "Gems." These not only permit a custom query template that can be used to

structure future responses but allows for incorporation of a large knowledge base. From the user’s perspective it does not differ greatly from a CustomGPT.

In addition to Gems, Google offers two other vehicles for using RAG technology (or technology that is certainly “RAG-like.”) The first is its famous NotebookLM. The whole point of this platform is to let an LLM gather sources from across the web or from your own files and use it to ground responses, including multimodal responses. It even offers an integrated “Discovery mode” that lets NotebookLM itself find the sources it thinks will be most valuable. The second is experimental and not yet well known: Opal. This platform lets users construct AI agents that perform a variety of sequenced or parallel tasks. The agents can be given access to “Assets” (like a bunch of documents) that can ground the agent’s performance.

The act of building one of these personal RAG systems is more than just a technical exercise or a way to build a productivity tool. It is a powerful exercise in knowledge management and metacognition—the process of thinking about one’s own thinking. To construct an effective RAG system, a student must first carefully curate, organize, and structure their knowledge base [54]. This forces a critical engagement with the course material, requiring the student to identify key concepts, understand the relationships between different topics, and structure their notes in a logical, coherent manner. When a query to their personal system yields a poor or incomplete response, it provides immediate feedback not only on the AI’s performance but, more importantly, on the quality, clarity, and organization of their own notes and understanding. This feedback loop transforms the RAG system into a tool for learning how to learn, compelling a structured approach to knowledge acquisition and synthesis that is an essential skill for successful legal practice.

## 1.7 RAG in the Broader AI Ecosystem: A Comparative Analysis

Retrieval-Augmented Generation is a powerful and increasingly essential architecture, but it is not the only method for enhancing the capabilities of Large Language Models. To make informed decisions about legal technology, it is crucial to understand where RAG fits within the broader ecosystem of AI tools and techniques. A comparative analysis of RAG against other prominent approaches, such as Model-Centric Prompting (MCP) and the use of ultra-large context windows, reveals a landscape of specialized tools, each with distinct strengths, weaknesses, and ideal use cases.

### 1.7.1 RAG vs. Model-Centric Prompting (MCP): Knowledge vs. Action

A common point of confusion is the distinction between RAG and another emerging concept, Model Context Protocol (MCP). While both serve to connect an LLM to external resources, they are designed for fundamentally different purposes.

- *Defining the Difference:* The simplest way to frame the distinction is *Knowledge vs. Action*.
  - RAG is for *knowing*. Its purpose is to retrieve static or semi-static information from a knowledge base to provide factual grounding for the LLM’s response. It answers the question, “What does the model need to *know* to answer this query accurately?” [55].
  - MCP is for *doing*. It is a standardized protocol that allows an LLM to interact with external “tools,” which are typically live APIs. It enables the model to take actions in other systems, such as querying a real-time database, sending an email, booking a meeting, or creating a ticket in a project management system [57].
- *Complementary Nature:* RAG and MCP are not competing technologies; they are highly complementary and synergistic [55]. The most powerful AI systems, often referred to as “agents,” will use both in a coordinated workflow. Consider an advanced legal assistant tasked with filing a motion.

1. The user asks, “Please draft and file a motion to compel discovery in the *Smith v. Jones* case.”



2. The agent first uses *RAG* to query the firm’s internal knowledge base and a legal database to retrieve the relevant federal rules of civil procedure, local court rules for that specific jurisdiction, and the firm’s own templates for motions to compel [61].
3. Using this retrieved knowledge, the LLM drafts the motion.
4. The agent then uses an *MCP tool* that connects to the court’s live e-filing API to submit the drafted document [62].

In this scenario, RAG provides the necessary knowledge, and MCP provides the ability to act on that knowledge.

### 1.7.2 RAG vs. Ultra-Large Context Windows: The “RAG is Dead” Debate

The recent emergence of LLMs with “ultra-large” context windows—the ability to process one million or even two million tokens (the equivalent of thousands of pages) in a single prompt—has sparked a debate about the future of RAG [63]. The argument is straightforward: if you can simply “stuff” all relevant documents directly into the prompt, is there still a need for a separate retrieval system? While the capability of these models is impressive, the declaration that “RAG is dead” is premature. A closer look reveals significant trade-offs that ensure RAG’s continued relevance.

- *Cost and Latency:* The business models for most commercial LLMs are based on token usage. Sending millions of tokens with every single query is, for most applications, prohibitively expensive and slow. RAG is far more efficient, as it retrieves only a small, targeted set of the most relevant text chunks, minimizing both cost and the time it takes to get a response [38].
- *The “Lost in the Middle” Problem:* Research from institutions like Stanford has demonstrated that LLMs, even with large context windows, do not pay uniform attention to the information provided. They exhibit strong primacy and recency bias, meaning they are much better at recalling information from the very beginning and very end of a long context. Performance degrades significantly when the model needs to find and use a crucial piece of information—a “needle in a haystack”—that is buried in the middle of a massive amount of text [65].
- *Precision vs. Brute Force:* RAG is a precision instrument. It is designed to surgically extract the most relevant snippets of information for the task at hand. Stuffing the entire context is a brute-force approach that inevitably includes a large amount of irrelevant noise. This noise can distract the model and degrade the quality of its reasoning and final output [38]. A smaller amount of high-relevance context will almost always produce a better result than a massive amount of low-relevance context.
- *Scalability and Currency:* A context window, no matter how large, is finite. A law firm’s or a corporation’s total knowledge base is effectively infinite and grows daily. It is not feasible to fit an entire enterprise’s document management system into a single prompt. RAG is designed to scale to these massive, dynamic knowledge bases, ensuring that the system can always access the most current information by simply updating the indexed documents [66].

The consensus among experts is that RAG is not dead; in fact, large context windows make RAG even more powerful. The future is likely a hybrid approach: using RAG’s efficient and precise retrieval mechanism to find the most relevant set of documents, and then placing that high-quality, curated context into a large context window to allow the LLM more space for complex synthesis and multi-document reasoning [67].

For a law student who will soon be a professional evaluating legal tech, a clear comparative framework is essential. The choice of underlying architecture has direct implications for a tool’s cost, accuracy, and reliability. The following table provides a synthesized overview of these different LLM augmentation techniques, mapping their technical characteristics to core legal and business concerns.

Table 1.1: A synthesized overview of LLM augmentation techniques. Data synthesized from [28, 71].

Technique	Primary Use Case	Data Currency	Cost & Complexity	Verifiability	Hallucination Risk
<b>Prompt Engineering</b>	Guiding the behavior of the LLM on a per-query basis.	Static (relies on model's training data).	Low. Can be done manually.	Low. No external sources are used.	High.
<b>Retrieval-Augmented Generation (RAG)</b>	Grounding the LLM's response in a specific, external knowledge base.	Real-Time Capable. Can query updated databases.	Medium. Requires setting up a data pipeline and vector store.	High. Can provide citations to retrieved sources.	Low.
<b>Fine-Tuning</b>	Specializing the model's skills and knowledge for a specific domain (e.g., legal writing).	Static (at time of tuning). Becomes part of the model's weights.	High. Requires large, curated datasets and significant computation.	Low. Knowledge is baked in; no direct source attribution.	Medium.
<b>Model Context Protocol (MCP)</b>	Enabling the LLM to perform actions by calling external tools and APIs.	Real-Time. Can access live data sources.	High. Requires building or integrating with secure API servers.	High. Actions and data queries are explicit and auditable.	Low (for data retrieval).
<b>Large Context Windows</b>	In-prompt analysis and synthesis of large but finite amounts of text.	Static (at time of query). Limited to what is provided in the prompt.	Very High (per query). Cost scales with input size.	Medium. Sources are in the prompt but can be hard to trace in the output.	Medium ("Lost in the Middle").

## 1.8 Conclusion: The Future-Ready Lawyer in the Age of Augmented Intelligence

Retrieval-Augmented Generation is not merely another tool in the expanding legal tech landscape; it represents a foundational architectural approach that directly addresses the legal profession's core, non-negotiable requirements for evidence-based reasoning, verifiability, and factual accuracy. By tethering the creative, probabilistic power of Large Language Models to the deterministic ground of authoritative data, RAG forges a necessary bridge between the potential of artificial intelligence and the principles of legal practice. It transforms the LLM from an unreliable savant into a diligent, open-book research assistant capable of augmenting, rather than replacing, human legal expertise.

For the law students of today who will be the legal professionals of tomorrow, understanding this architecture is no longer a niche or optional skill. As AI becomes increasingly integrated into the fabric of legal services—from research and discovery to contract analysis and compliance—the ability to discern how these tools work, to understand their strengths, and, most critically, to recognize their limitations, will become a fundamental component of professional competence. The future-ready lawyer will not be the one who can code an AI from scratch, nor will they be the one who is replaced by it. The successful lawyer will be the one who understands how to strategically *augment* their own judgment and expertise with these powerful new capabilities.

They will be the ones who can critically evaluate a vendor's claims by asking not just *what* a tool does, but *how* it does it. Is it using a generic or a domain-specific embedding model? What is its chunking strategy for complex contracts? How does it ensure the currency of its knowledge base? These are the questions that separate a savvy adopter of technology from a potential victim of its pitfalls. RAG is not a panacea, and it does not eliminate the need for human oversight and critical judgment. Rather, it is a framework that, when implemented thoughtfully, provides that human judgment with a more powerful, reliable, and comprehensive evidentiary foundation. The ultimate goal is not to automate legal reasoning, but to empower it, enabling lawyers to deliver faster, more accurate, and more efficient services while upholding the highest standards of their professional and ethical obligations.

# Bibliography

- [1] Why Every Law Firm Will Need RAG: The Strategic Imperative for Legal AI - Medium, accessed August 19, 2025, [https://medium.com/@chris\\_hobbick/why-every-law-firm-will-need-rag-the-strategic-imperative-for-legal-ai-c9c2df4e3601](https://medium.com/@chris_hobbick/why-every-law-firm-will-need-rag-the-strategic-imperative-for-legal-ai-c9c2df4e3601)
- [2] AI and You: The Critical Role of Retrieval Augmented Generation (RAG) in Legal Practice, accessed August 19, 2025, <https://www.americanbar.org/groups/gpsolo/resources/ereport/2024-july/ai-you-critical-role-retrieval-augmented-generation-rag-legal-practice/>
- [3] What is Retrieval-Augmented Generation (RAG)? - Google Cloud, accessed August 19, 2025, <https://cloud.google.com/use-cases/retrieval-augmented-generation>
- [4] What is RAG? - Retrieval-Augmented Generation AI Explained - AWS, accessed August 19, 2025, <https://aws.amazon.com/what-is/retrieval-augmented-generation/>
- [5] Retrieval Augmented Generation - IBM, accessed August 19, 2025, <https://www.ibm.com/architectures/patterns/genai-rag>
- [6] Intro to retrieval-augmented generation (RAG) in legal tech, accessed August 19, 2025, <https://legal.thomsonreuters.com/blog/retrieval-augmented-generation-in-legal-tech/>
- [7] Word embedding - Wikipedia, accessed August 19, 2025, [https://en.wikipedia.org/wiki/Word\\_embedding](https://en.wikipedia.org/wiki/Word_embedding)
- [8] What is Embedding? - Embeddings in Machine Learning Explained - AWS, accessed August 19, 2025, <https://aws.amazon.com/what-is/embeddings-in-machine-learning/>
- [9] Word Embeddings in NLP - GeeksforGeeks, accessed August 19, 2025, <https://www.geeksforgeeks.org/nlp/word-embeddings-in-nlp/>
- [10] www.ibm.com, accessed August 19, 2025, <https://www.ibm.com/think/topics/word-embeddings#:~:text=Word%20embeddings%20are%20a%20way,relationships%20among%20the%20corresponding%20words.>
- [11] The Ultimate Guide to Word Embeddings - neptune.ai, accessed August 19, 2025, <https://neptune.ai/blog/word-embeddings-guide>
- [12] The Beginner's Guide to Text Embeddings — deepset Blog, accessed August 19, 2025, <https://www.deepset.ai/blog/the-beginners-guide-to-text-embeddings>
- [13] What Are Word Embeddings? — IBM, accessed August 19, 2025, <https://www.ibm.com/think/topics/word-embeddings>
- [14] Word embeddings in NLP: A Complete Guide - Turing, accessed August 19, 2025, <https://www.turing.com/kb/guide-on-word-embeddings-in-nlp>
- [15] Word Embedding Demo: Tutorial, accessed August 19, 2025, <https://www.cs.cmu.edu/~dst/WordEmbeddingDemo/tutorial.html>

- [16] An Introduction to Vector Databases for Beginners - Xomnia, accessed August 19, 2025, <https://xomnia.com/post/an-introduction-to-vector-databases-for-beginners/>
- [17] What Is A Vector Database? - IBM, accessed August 19, 2025, <https://www.ibm.com/think/topics/vector-database>
- [18] What is a Vector Database? - Elastic, accessed August 19, 2025, <https://www.elastic.co/what-is/vector-database>
- [19] Understanding Vector Databases — Microsoft Learn, accessed August 19, 2025, <https://learn.microsoft.com/en-us/data-engineering/playbook/solutions/vector-database/>
- [20] What is a Vector Database? - AWS, accessed August 19, 2025, <https://aws.amazon.com/what-is/vector-databases/>
- [21] Introduction to Data Mining Distances & Similarities, accessed August 19, 2025, [https://www.ccs.neu.edu/home/vip/teach/DMcourse/1\\_intro\\_sim\\_knn/notes\\_slides/DistancesSimilarities.pdf](https://www.ccs.neu.edu/home/vip/teach/DMcourse/1_intro_sim_knn/notes_slides/DistancesSimilarities.pdf)
- [22] 9 Distance Measures in Data Science - Maarten Grootendorst, accessed August 19, 2025, <https://www.maartengrootendorst.com/blog/distances/>
- [23] Understanding and Using Common Similarity Measures for Text Analysis, accessed August 19, 2025, <https://programminghistorian.org/en/lessons/common-similarity-measures>
- [24] Introduction to Similarity Measures (Cosine, DotProduct, etc) : r/learnmachinelearning, accessed August 19, 2025, [https://www.reddit.com/r/learnmachinelearning/comments/172ilqs/introduction\\_to\\_similarity\\_measures\\_cosine/](https://www.reddit.com/r/learnmachinelearning/comments/172ilqs/introduction_to_similarity_measures_cosine/)
- [25] What is Cosine Distance? A Deep Dive - DataCamp, accessed August 19, 2025, <https://www.datacamp.com/tutorial/cosine-distance>
- [26] What is Retrieval Augmented Generation (RAG)? - DataCamp, accessed August 19, 2025, <https://www.datacamp.com/blog/what-is-retrieval-augmented-generation-rag>
- [27] Embedding Models Explained: A Guide to NLP's Core Technology - Medium, accessed August 19, 2025, <https://medium.com/@nay1228/embedding-models-a-comprehensive-guide-for-beginners-to-experts-0cfc11d449f1>
- [28] What is Retrieval Augmented Generation (RAG)? - Databricks, accessed August 19, 2025, <https://www.databricks.com/glossary/retrieval-augmented-generation-rag>
- [29] RAG-as-a-Service for Legal Tech — Nuclia AI, accessed August 19, 2025, <https://nuclia.com/rag-for-legal/>
- [30] How Law Firms Use RAG to Boost Legal Research - Datategy, accessed August 19, 2025, <https://www.datategy.net/2025/04/14/how-law-firms-use-rag-to-boost-legal-research/>
- [31] RAG in Legal Research: A Guide — e! - Lexemo's, accessed August 19, 2025, <https://e.lexemo.com/uncategorized-en/retrieval-augmented-generation-rag-in-legal-research/>
- [32] Retrieval-Augmented Generation with Vector Stores, Knowledge Graphs, and Hierarchical Non-negative Matrix Factorization - arXiv, accessed August 19, 2025, <https://arxiv.org/html/2502.20364v1>
- [33] 3 practical applications of Retrieval-Augmented Generation for legal teams - Onna, accessed August 19, 2025, <https://onna.com/blog/practical-applications-of-retrieval-augmented-generation-for-legal-teams>
- [34] 5 RAG Chunking Strategies for Better Retrieval-Augmented Generation - Lettria, accessed August 19, 2025, <https://www.lettria.com/blogpost/5-rag-chunking-strategies-for-better-retrieval-augmented-generation>

- [35] 7 Chunking Strategies in RAG You Need To Know - F22 Labs, accessed August 19, 2025, <https://www.f22labs.com/blogs/7-chunking-strategies-in-rag-you-need-to-know/>
- [36] How to Choose the Right Chunking Strategy for Your LLM Application — MongoDB, accessed August 19, 2025, <https://www.mongodb.com/developer/products/atlas/choosing-chunking-strategy-rag/>
- [37] RAG is failing when the number of documents increase - API - OpenAI Community Forum, accessed August 19, 2025, <https://community.openai.com/t/rag-is-failing-when-the-number-of-documents-increase/578498>
- [38] RAG vs Long Context Models [Discussion] : r/MachineLearning - Reddit, accessed August 19, 2025, [https://www.reddit.com/r/MachineLearning/comments/1ax6j73/rag\\_vs\\_long\\_context\\_models\\_discussion/](https://www.reddit.com/r/MachineLearning/comments/1ax6j73/rag_vs_long_context_models_discussion/)
- [39] Optimizing Document-Level Retrieval in RAG: Alternative Approaches? - Reddit, accessed August 19, 2025, [https://www.reddit.com/r/Rag/comments/1ii0b8b/optimizing\\_documentlevel\\_retrieval\\_in\\_rag/](https://www.reddit.com/r/Rag/comments/1ii0b8b/optimizing_documentlevel_retrieval_in_rag/)
- [40] What are the limitations of embeddings? - Zilliz Vector Database, accessed August 19, 2025, <https://zilliz.com/ai-faq/what-are-the-limitations-of-embeddings>
- [41] Natural language processing in the legal domain - InK@SMU.edu.sg, accessed August 19, 2025, [https://ink.library.smu.edu.sg/cgi/viewcontent.cgi?article=6485&context=sol\\_research](https://ink.library.smu.edu.sg/cgi/viewcontent.cgi?article=6485&context=sol_research)
- [42] Legal-RAG vs. RAG: A Technical Exploration of Retrieval Systems ..., accessed August 19, 2025, <https://www.truelaw.ai/blog/legal-rag-vs-rag-a-technical-exploration-of-retrieval-systems>
- [43] Beyond Keywords: Optimizing Legal Information Retrieval through Embeddings, Cross-Encoders, and Large Language Models - DiVA, accessed August 19, 2025, <https://kth.diva-portal.org/smash/get/diva2:1963310/FULLTEXT01.pdf>
- [44] Hybrid Search RAG: Revolutionizing Information Retrieval — by Alex ..., accessed August 19, 2025, <https://medium.com/@alexrodriguesj/hybrid-search-rag-revolutionizing-information-retrieval-9905d3437cdd>
- [45] Hybrid Search a method to Optimize RAG implementation — by Akash Chandrasekar, accessed August 19, 2025, <https://medium.com/@csakash03/hybrid-search-is-a-method-to-optimize-rag-implementation-98d9d0911341>
- [46] Optimizing RAG with Hybrid Search & Reranking — VectorHub by Superlinked, accessed August 19, 2025, <https://superlinked.com/vectorhub/articles/optimizing-rag-with-hybrid-search-reranking>
- [47] Re-Ranking Mechanisms in Retrieval-Augmented Generation ..., accessed August 19, 2025, <https://medium.com/@adnanmasood/re-ranking-mechanisms-in-retrieval-augmented-generation-pipelines-an-overview-8e24303ee789>
- [48] RankRAG: Unifying Context Ranking with Retrieval-Augmented Generation in LLMs - arXiv, accessed August 19, 2025, <https://arxiv.org/html/2407.02485v1>
- [49] RankRAG: Unifying Context Ranking with Retrieval-Augmented Generation in LLMs, accessed August 19, 2025, [https://openreview.net/forum?id=S1fc92uemC&referrer=%5Bthe%20profile%20of%20Jiaxuan%20You%5D\(%2Fprofile%3Fid%3D~Jiaxuan\\_You2\)](https://openreview.net/forum?id=S1fc92uemC&referrer=%5Bthe%20profile%20of%20Jiaxuan%20You%5D(%2Fprofile%3Fid%3D~Jiaxuan_You2))
- [50] Introducing RAG 2.0 - Contextual AI, accessed August 19, 2025, <https://contextual.ai/introducing-rag2/>

- [51] Retrieval Augmented Generation (RAG) and Semantic Search for ..., accessed August 19, 2025, <https://help.openai.com/en/articles/8868588-retrieval-augmented-generation-rag-and-semantic-search-for-gpts>
- [52] Building GPTs - OpenAI Help Center, accessed August 19, 2025, <https://help.openai.com/en/collectons/8475422-building-gpts>
- [53] Retrieval Augmented Generation (RAG) for Projects — Anthropic Help ..., accessed August 19, 2025, <https://support.anthropic.com/en/articles/11473015-retrieval-augmented-generation-rag-for-projects>
- [54] How to Create a RAG System: A Complete Guide to Retrieval ..., accessed August 19, 2025, <https://www.mindee.com/blog/build-rag-system-guide>
- [55] MCP vs RAG : Know The Key Differences - TrueFoundry, accessed August 19, 2025, <https://www.truefoundry.com/blog/mcp-vs-rag>
- [56] MCP vs RAG: how they overlap and differ - Merge.dev, accessed August 19, 2025, <https://www.merge.dev/blog/rag-vs-mcp>
- [57] Is MCP a better alternative to RAG for Observability? — Parseable Blog, accessed August 19, 2025, <https://www.parseable.com/blog/mcp-better-alternative-to-rag-for-observability>
- [58] How RAG & MCP solve model limitations differently - DEV Community, accessed August 19, 2025, <https://dev.to/aws/how-rag-mcp-solve-model-limitations-differently-pjm>
- [59] MCP Implementation using RAG: A Step-by-step Guide - ProjectPro, accessed August 19, 2025, <https://www.projectpro.io/article/mcp-with-rag/1144>
- [60] MCP vs RAG: Making Sense of Agentic AI's Alphabet Soup - Segun Akinyemi, accessed August 19, 2025, <https://segunakinyemi.com/blog/mcp-vs-rag/>
- [61] RAG vs MCP vs Agents — What's the right fit for my use case? : r/LLMDevs - Reddit, accessed August 19, 2025, [https://www.reddit.com/r/LLMDevs/comments/1l2j6s4/rag\\_vs\\_mcp\\_vs\\_agents\\_whats\\_the\\_right\\_fit\\_for\\_my/](https://www.reddit.com/r/LLMDevs/comments/1l2j6s4/rag_vs_mcp_vs_agents_whats_the_right_fit_for_my/)
- [62] MCP and RAG: A Powerful Partnership for Advanced AI Applications — by Plaban Nayak — The AI Forum — Medium, accessed August 19, 2025, <https://medium.com/the-ai-forum/mcp-and-rag-a-powerful-partnership-for-advanced-ai-applications-858c074fc5db>
- [63] Long Context Models Explained: Do We Still Need RAG? - Louis Bouchard, accessed August 19, 2025, <https://www.louisbouchard.ai/long-context-vs-rag/>
- [64] Understanding the Impact of Increasing LLM Context Windows - Meibel, accessed August 19, 2025, <https://www.meibel.ai/post/understanding-the-impact-of-increasing-llm-context-windows>
- [65] Exploring The Dynamics Between Long Context Windows And RAG Systems In AI, accessed August 19, 2025, <https://customgpt.ai/long-context-windows-vs-rag/>
- [66] Navigating the Challenges of Context Window Limitations — by Awarity.ai - Medium, accessed August 19, 2025, <https://medium.com/awarity-ai-blog/navigating-the-challenges-of-context-window-limitations-eef2dcfc02e1>
- [67] What are your thoughts on the 'RAG is dead' debate as context windows get longer? - Reddit, accessed August 19, 2025, [https://www.reddit.com/r/LLMDevs/comments/1mt51h5/what\\_are\\_your\\_thoughts\\_on\\_the\\_rag\\_is\\_dead\\_debate/](https://www.reddit.com/r/LLMDevs/comments/1mt51h5/what_are_your_thoughts_on_the_rag_is_dead_debate/)
- [68] Can Large Context Windows Replace RAG? - Infogain, accessed August 19, 2025, <https://www.infogain.com/blog/can-large-context-windows-replace-rag/>

- [69] RAG vs fine-tuning vs. prompt engineering - IBM, accessed August 19, 2025, <https://www.ibm.com/think/topics/rag-vs-fine-tuning-vs-prompt-engineering>
- [70] RAG vs prompt engineering: Getting the best of both worlds - K2view, accessed August 19, 2025, <https://www.k2view.com/blog/rag-vs-prompt-engineering/>
- [71] Understanding the Key Differences: Prompting, Fine-Tuning, and Retrieval-Augmented Generation (RAG) — by Swaroop Piduguralla — Medium, accessed August 19, 2025, <https://medium.com/@tejaswaroop2310/understanding-the-key-differences-prompting-fine-tuning-and-retrieval-augmented-generation-rag-d685aefcc1a7>



## Chapter 2

# Custom AI for the Modern Lawyer

*Customized large language models are reusable sets of instructions, optionally combined with source documents, that direct an AI to perform a specific, recurring task. For law students and practitioners, they offer a powerful way to automate cognitive heavy lifting, from summarizing depositions to drafting discovery requests. What if you could build a specialist assistant for every tedious part of your legal workflow, and do it in under five minutes?*

These tools—sold as Gemini Gems, CustomGPTs, and Claude Projects, among others—represent a significant step beyond generic, one-shot prompting. They are, in essence, an exceptionally user-friendly implementation of a technique called Retrieval-Augmented Generation (RAG), which allows a language model to ground its responses in a specific set of external documents. We define these tools, show when to choose Gems versus CustomGPTs versus Projects, explain the “saved prompting plus knowledge” mechanism, provide concrete legal examples with alignment patterns, and contrast them with agentic AI.

### 2.1 Defining the Tools: Gems, CustomGPTs, and Projects

At their core, Gemini Gems (Google), CustomGPTs (OpenAI), and Claude Projects (Anthropic) are functionally similar. They allow a user to create a specialized version of a base model by providing it with a set of persistent instructions and, critically, a corpus of “knowledge” files.

**Mechanism:** The process is straightforward. First, you give the custom AI a name and a core instruction set. This is not a single prompt, but a constitution—a set of rules and a persona it should adopt for every interaction. For example, a “Deposition Summarizer” Gem might be instructed: “You are a paralegal specializing in tort law. Your task is to receive deposition transcripts, identify key admissions and contradictions related to the elements of negligence, and present them in a bulleted list with page and line citations.”

Second, you upload source documents. This is the “Retrieval-Augmented” part of RAG. When you upload a “knowledge” file, the platform automatically converts it into a format the AI can search through embeddings and vector databases. When you issue a prompt, the system first searches your uploaded documents for relevant snippets and then feeds those snippets to the model as part of a much more complex, hidden prompt. For our Deposition Summarizer, this could be the case complaint, which lists the formal elements of the negligence claim. The AI will now analyze transcripts *in light of* the specific legal claims outlined in the complaint.

**Ease of Creation:** The defining feature of these tools is their accessibility. Creating a basic version requires no coding. The interface is conversational. You essentially tell the AI what you want it to be, and it builds itself. You can then refine its behavior through a simple chat interface and test it in a preview

window. The entire process of creating a useful, task-specific AI can take less time than drafting a single memorandum.

## 2.2 Strategic Choices: When to Use Each Platform

While functionally similar, the platforms have key differences that inform when to use each.

**Gemini Gems:** As a Google product, Gems are deeply integrated into the Google ecosystem. Recent updates place Gems directly in the Docs/Sheets/Slides/Gmail side panel, allowing one-click invocation while drafting. They're best suited for personal productivity workflows where your source material resides within Google's walled garden. Currently, true Gem sharing is limited—you can share a chat (with optional instructions) but not the Gem object itself. Mechanism: Google stores the configuration and surfaces the Gem in Workspace apps; you can attach files from your device or Drive with documented limits on file size and number.

**CustomGPTs:** OpenAI has taken a different approach, creating a GPT Store that makes CustomGPTs searchable and shareable. A law firm could create a suite of internal CustomGPTs for its associates, or a legal tech company could publish a CustomGPT for public use. OpenAI specifies clear file/size/token limits for Knowledge (up to 20 files, 512 MB each), and documents how files are chunked and embedded for retrieval. This transparency and shareability make them powerful tools for knowledge management within an organization. If the goal is to create a tool for a team, CustomGPTs are the logical choice.

**Claude Projects:** Anthropic's Claude is known for its large context window, meaning it can handle and analyze extremely long documents in a single prompt. A "Project" in Claude is akin to a dedicated workspace with its own instructions, chats, and Project Knowledge. For Team/Enterprise plans, Projects can be shared with colleagues with "can use" or "can edit" permissions. Projects automatically switch to a built-in RAG mode as the knowledge base grows, scaling the amount of usable reference material beyond the raw context window. This makes them excellent for deep analysis of specific document sets, such as contract review or preparing a witness based on their entire deposition history.

The choice of platform is therefore tactical: for personal, Google-integrated tasks, use a Gem. For shareable, team-oriented tools, use a CustomGPT. For deep, document-intensive analysis with team collaboration, use a Claude Project.

## 2.3 Use Cases: Taming the Cumbersome in Law

The value of these tools is unlocked by identifying high-frequency, low-creativity tasks that are nonetheless time-consuming. The practice and study of law are replete with such activities.

### 2.3.1 Example 1: The Case Brief Generator

A foundational chore for any first-year law student is briefing cases. The process is mechanical: identify the facts, procedure, issue, holding, and reasoning. While essential for learning, the fiftieth time is less about pedagogy and more about endurance.

A student could create a CustomGPT named "1L Case Briefer."

- **Instructions:** "You are a law professor's research assistant. When I upload a judicial opinion, generate a concise case brief. The brief must follow the FIRAC (Facts, Issue, Rule, Application, Conclusion) structure. Be precise and extract direct quotes for the specific legal rule announced by the court. The target audience is a 1L law student preparing for a cold call."

- **Knowledge Files:** A PDF explaining the FIRAC structure and a dozen examples of well-written case briefs the student has previously drafted. This helps the AI mimic the desired style and level of detail.
- **Mechanism:** The assistant enforces the schema every time; attached casebook PDFs ground quotes through retrieval rather than memory.

Now, instead of spending 45 minutes manually dissecting each case, the student uploads the opinion and receives a structured brief in seconds. The task shifts from laborious extraction to critical review and refinement—a higher-value activity.

### 2.3.2 Example 2: The Discovery Drafter

In litigation practice, drafting initial discovery documents—interrogatories, requests for production, and requests for admission—is often a matter of adapting standard templates to the facts of a new case. This is ripe for automation.

A junior associate could build a Gemini Gem called “Standard Discovery Set - Premises Liability.”

- **Instructions:** “You are a litigation paralegal. When I provide you with a short summary of a new slip-and-fall case, you will generate a standard set of initial discovery requests based on the uploaded templates. You will replace placeholders like [PLAINTIFF’S NAME] or [DATE OF INCIDENT] with the information I provide in my prompt. The output should be three separate documents: Interrogatories, Requests for Production of Documents, and Requests for Admissions.”
- **Knowledge Files:** The firm’s approved template discovery documents for premises liability cases (as .docx files).
- **Mechanism:** Side-panel presence in Google Docs prevents context-loss; templates ensure jurisdictional compliance.

The associate’s workflow transforms. Instead of opening three templates, performing a tedious find-and-replace, and risking clerical errors, they write a single prompt with case specifics. The Gem produces a clean, ready-to-review draft.

### 2.3.3 Example 3: The Bluebook Citation Checker

Legal citation is a universe of baroque rules. Mastering *The Bluebook* is a rite of passage, but its application is a repetitive, error-prone task that distracts from the substantive legal argument.

A law review editor could create a Claude Project named “Bluebook Style Assistant.”

- **Instructions:** “You are a meticulous law review editor with deep expertise in *The Bluebook* (21st ed.). When I paste a legal citation, you will check it for correctness and provide a corrected version if necessary, explaining the specific rule that was violated. Your analysis should be limited to citation format and not the substance of the source.”
- **Knowledge Files:** Style guides created by the journal, cheat sheets, and examples of correctly formatted citations for various source types (cases, statutes, articles).
- **Mechanism:** Project Knowledge keeps terminology consistent across chats; automatic RAG scaling handles large style manuals.

This tool allows students to get instant feedback on their citations, improving accuracy and freeing up cognitive bandwidth for the core task of legal writing and analysis.

## 2.4 Creating Well-Aligned Custom Assistants

To create a generic, well-aligned custom AI, follow these principles:

**Role + Audience:** Begin instructions with a clear persona. “You are a tax attorney’s assistant for a 2L student” is better than “Help me with tax.” Be explicit about what it *should not* do. “Do not provide investment advice or speculate on market movements.”

**Scope:** “Use only the uploaded record and controlling authorities; if a request asks for factual assertions outside the record, say ‘Out of record—needs source.’” Mechanism: this steers the RAG system toward retrieval before free-form generation.

**Output Contract:** Don’t just ask for a summary; ask for a “summary in three bullet points, followed by a list of key actors and a one-sentence conclusion.” This structures the output and reduces variability. For legal documents: “Always return: (1) short answer; (2) doctrine; (3) application; (4) risks/alternatives; (5) citations with pincites.”

**Source-First Rule:** “If the answer can be grounded in uploaded knowledge, quote and cite it; otherwise ask for a source or state uncertainty.” This constraint makes the assistant usable by anyone while maintaining accuracy.

**Truthfulness Guardrails:** “Never fabricate a citation, docket number, or page cite. If unsure, output a ‘verification needed’ list.” Add: “After answering, create a ‘Verify’ checklist listing which quotes and pincites to confirm, with page numbers.”

**Academic Integrity Line:** “Do not write graded work product verbatim; provide outlines, structures, and feedback; require me to apply the rules to facts.”

**Style Hooks:** “Prefer headings over boldface within text; keep bullet lists to checklists; show one example, then switch to analysis.”

**Iterate:** Your first version will not be perfect. Use the preview window to test it. If it makes a mistake, don’t just correct the output; go back and amend the core instructions to prevent that class of error in the future.

## 2.5 Additional High-Value Use Cases

Beyond the three detailed examples, these tools excel at numerous other legal tasks:

1. **Reading-to-Outline Pipeline:** “Turn today’s reading notes into a doctrinal outline section: rule-statement, elements, defenses, policy.” Mechanism: side-panel use in Google Docs keeps the outline in the same file.
2. **Socratic Drill Coach:** “Ask one question at a time on Commerce Clause doctrine; escalate difficulty if I nail it; if I miss, supply the controlling case and a short corrective explanation.”
3. **Closed-Universe Research Assistant:** Upload the record, key cases, and bench memos; instruct, “Answer only from these files; cite to page/line; if outside scope, say ‘Not in record.’” Mechanism: RAG confines outputs to uploaded sources.
4. **Motion Skeletons and Checklists:** “Given a fact memo, emit a motion outline with headings, the standard of review, the key elements to prove, and a list of required attachments.”
5. **Deposition Prep Outlines:** “Given a witness description and case theory, draft a tiered outline: lock-ins, exhibits, impeachment points, and ‘funnel’ sequences.”
6. **Discovery Request Narrower:** “Draft 10 RFPs that meet proportionality; for each, add a sentence explaining why it’s not overbroad.”

7. **Journal Production Pipeline:** Shared with the editorial board: “Proof to journal style; check citations for required signals; generate table of authorities; label any ambiguous sources for human review.”
8. **Exam-Mode Practice:** “Emit a 60-minute issue-spotter on Admin Law with 5 latent issues; timebox follow-ups; grade with a rubric I supply; give only one hint per issue.”

## 2.6 The Next Frontier: Agents vs. Custom AI

It is crucial to distinguish these custom AIs from “agentic AI.” A custom AI, like our Deposition Summarizer, is a passive tool. It waits for you to upload a transcript and give it a command. It executes a well-defined, if complex, task.

An AI agent, by contrast, is autonomous. It can be given a high-level goal and can then independently formulate and execute a sequence of tasks to achieve it. An agent might be tasked with “monitoring the dockets for all new patent infringement lawsuits filed in the Eastern District of Texas and preparing a summary of each complaint.” This agent would not wait for instruction. It would (hypothetically) access a court docket database, identify new filings, download the relevant documents, feed them to a specialized analysis tool, and deliver a summary without human intervention at each step. Mechanism: agents decide when to call tools and iterate through function calling, code execution, and API orchestration—capabilities documented in OpenAI’s Assistants platform and Google’s agent development frameworks.

Custom AIs are not agents. They are, however, a critical component of a future agentic workflow. One can easily imagine an AI agent for litigation management that, upon receiving a new case file, automatically invokes a series of specialized CustomGPTs: one to brief the key precedents, another to draft initial discovery, and a third to summarize the opposing party’s corporate filings. The custom AIs are the skilled specialists; the agent is the project manager.

Safe patterns for combining them include:

1. **Custom AI as the “policy module,” agent as the “operator.”** Put your standards (style, refusal rules, output schema) into a CustomGPT/Project. The agent then calls that assistant to generate drafts, while the agent itself handles file I/O, scraping, calendaring, or API actions under human supervision.
2. **Agent performs collection; saved assistant performs analysis.** Let an agent gather cases via APIs; hand the curated corpus to your CustomGPT/Project to summarize and check cites—keeping generation anchored in a vetted set.
3. **Saved assistant as a safe default; escalate to agent only on demand.** Teach your assistant: “If this requires browsing, calendars, or external data pulls, ask me to escalate to the agent.”

In this model, law students and junior lawyers should focus on building a library of custom AIs that automate the discrete, repetitive components of their work. Start tiny: create one assistant for one course or one journal task with only the minimum sources needed. Fewer files mean cleaner retrieval, as both OpenAI and Anthropic documentation implicitly rewards small, high-signal corpora.

The future of legal work is not a single, all-knowing AI, but a dynamic collaboration between human lawyers, specialized AI tools they have built and trained, and autonomous agents that manage the workflow. These custom AI tools—Gems, CustomGPTs, and Projects—are your entry point. They require no coding, take minutes to create, and immediately reduce the tedious parts of legal work. The time to start building your toolkit is now.

## Chapter 3

# What is Prompt Caching?

Prompt caching is a technique used to optimize interactions with large language models (LLMs). It works by storing the processed state of a large, static portion of a prompt that is used repeatedly across multiple API calls.

Conceptually, it is analogous to a web browser's cache. When a browser loads a webpage for the first time, it downloads and stores assets like images and scripts. On subsequent visits, it loads these assets from the local cache, significantly speeding up the process. Similarly, an LLM expends considerable computational resources to process the input text (tokens) it receives. By caching the processed version of a recurring part of a prompt (such as a detailed system instruction), the model can bypass the redundant step of reprocessing this same information for every new request. For subsequent calls, only the new, variable part of the prompt is sent, and the model begins its generation process from the already-processed cached state.

### 3.1 Primary Benefits

The main advantages of implementing prompt caching are speed and cost-efficiency.

- **Reduced Latency:** By eliminating the need to re-process a potentially lengthy system prompt or a set of few-shot examples for every call, the model can generate a response much more quickly. The time-to-first-token is significantly reduced, leading to a more responsive application.
- **Cost Savings:** Most LLM APIs operate on a pay-per-token basis, counting both input and output tokens. If an application uses a 2,000-token system prompt for every 50-token user query, it is charged for 2,050 input tokens each time. With caching, after the initial call, the application would only be charged for the 50 new tokens, leading to substantial cost savings at scale.

### 3.2 Applicability: API vs. Web Interface

Prompt caching is fundamentally a developer-oriented feature designed for programmatic use via an API. It is not typically a feature that end-users can control through a standard web-based chat interface. While the service providers of such interfaces may employ similar optimization techniques on their backend, direct control over caching is exposed at the API level for developers building applications on top of the model.

## 3.3 Use Case: Legal Document Analysis

A compelling application for prompt caching is in the legal field, where practitioners often need to analyze voluminous texts.

### 3.3.1 Scenario

Suppose a legal team needs to run a series of queries against the complete texts of several lengthy court cases or contracts. These documents can easily contain tens of thousands of tokens. The goal is to extract specific pieces of information, summarize arguments, or find precedents within these texts.

### 3.3.2 Implementation with Caching

1. **The Cached Content:** The full, immutable text of the legal cases or contracts would form the static part of the prompt. An initial API call would be made to process these documents and store them in a cache, associated with a unique identifier.
2. **The Variable Content:** Each subsequent query from the legal team would be a short, specific question. For example:
  - “Summarize the plaintiff’s primary argument in Case A.”
  - “Identify all clauses related to intellectual property rights in this contract.”
  - “In which of these documents is ‘breach of fiduciary duty’ discussed?”
3. **The API Call:** For each question, the API call would simply reference the cached content’s identifier and provide the new question. The LLM would not need to re-read and re-process the entire set of legal documents each time.

### 3.3.3 Conclusion

In this scenario, prompt caching makes perfect sense. It transforms an otherwise slow and expensive task into a rapid and cost-effective one. The legal team can conduct an interactive and iterative analysis, asking dozens of follow-up questions without incurring the high latency and cost of reprocessing the source documents for every single query.