# Reliability by design: quantifying and eliminating fabrication risk in LLMs

## From Generative to Consultative AI: A Comparative Analysis in the Legal Domain and Lessons for High-Stakes Knowledge Bases.

**Alex Dantart**
Humanizing Internet
arxiv@humanizinginternet.com

## ABSTRACT

Large language models (LLMs) are redefining the market, but their potential is threatened by the endemic generation of "hallucinations," an unacceptable risk in high-rigor domains such as law. This article presents a quantitative empirical analysis of this phenomenon (exportable to other corpora and topics), establishing a critical distinction between three operational AI paradigms: general-purpose generative AI ("creative oracle"), canonical consultative AI ("expert archivist") implemented with basic RAG, and advanced consultative AI ("rigorous archivist"), which integrates a holistic optimization pipeline.

To measure factual reliability, we propose and validate two specific metrics: the false citation rate (FCR) and the fabricated fact rate (FFR). We generate and analyze a corpus of 2,700 judicial responses from our dataset of 75 legal tasks (JURIDICO-FCR), using twelve state-of-the-art LLMs in the three conditions. Evaluation is carried out through a double-blind study by expert lawyers.

The results are compelling. Generative AI exhibits unacceptable error rates (FCR > 30%), rendering it invalid for professional practice. Canonical RAG drastically reduces this risk by more than two orders of magnitude, but leaves a significant residual risk of misgrounding. It is the advanced RAG paradigm (which incorporates techniques such as embedding fine-tuning, re-ranking with cross-encoders, and self-correction loops) that manages to almost completely eliminate the risk of fabrication, reducing error rates to statistically insignificant levels (<0.2%).

We conclude that the reliability of legal AI does not lie in incremental improvement of the "oracle," but in the deliberate adoption of advanced consultative architectures, where truthfulness, traceability, and self-verification are at the core of the design. This study offers an evaluation framework and a roadmap for building truly reliable AI systems for any high-risk knowledge corpus.

## 1 Introduction

Artificial intelligence (AI), and in particular large language models (LLMs), are at the cusp of a significant transformation across multiple sectors, with the legal domain being one of the most impacted. Tools such as ChatGPT, DeepSeek, Claude, or Gemini promise to revolutionize tasks such as legal research and document drafting, with considerable potential to increase efficiency and democratize access to justice.

However, a shadow looms over all this potential: the inherent and critical phenomenon of 'hallucinations.' We refer to the generation of information that, although often sounding plausible, is factually incorrect, misleading, or outright fabricated. This phenomenon compels us to pose a fundamental question: are we building a reliable assistant or, on the contrary, an eloquent fabulist? This dichotomy will be the central axis of our empirical analysis. In the legal context,

where truthfulness and fidelity to authoritative sources are fundamental pillars, hallucinations transcend the category of mere technical inaccuracies. In reality, they represent a profound epistemic challenge for the profession. Unlike a simple human error, an LLM's hallucination not only provides incorrect information but also fabricates the very source of authority, creating a simulacrum of knowledge that directly undermines the foundations upon which legal argumentation is built: the verifiability of sources and the traceability of reasoning.

This systemic risk is not a mere hypothesis, but a documented reality in the highest judicial instances. A paradigmatic example is the **sanction imposed by the Constitutional Court of Spain on a lawyer for using 19 non-existent case law citations** in an appeal for protection (Constitutional Court, 2024). In its decision, the court described the act as a lack of due respect with a 'clear disregard for the judicial function.' This incident underscores how the fabrication of authority by an AI directly attacks the foundations of legal argumentation. The phenomenon is rapidly spreading across Europe: the Italian judiciary has recently issued warnings after detecting lawyers citing non-existent rulings generated by ChatGPT, warning that the ease of use of these tools does not exempt from the duty of verifying truthfulness (2025). Similarly, also in 2025, in the Latin American context, a decision by the Court of Appeals of Córdoba, Argentina, urged a lawyer to make 'ethical and responsible use of AI tools' after detecting the citation of multiple non-existent rulings, warning about the professional responsibility that cannot be replaced by an algorithm. These cases, **more recent and closer to our legal environment** than the well-known *Mata v. Avianca, Inc.* (2023), serve as a stark reminder of the dangers and the urgency of establishing a rigorous evaluation framework such as the one proposed in this study.

However, it is essential to understand that, although technology heightens the risk, the underlying problem is not entirely new. Generative AI does not invent malpractice; rather, it industrializes it on an unprecedented scale. Historically, the legal profession has always dealt with the 'misuse of authority' and 'citation drift,' a phenomenon in which the real link between a source and the argument it is intended to support is weakened or fabricated. What was once a 'craft' act, the result of a serious error or an infrequent intent to deceive, is now turned by technology into a systemic and automated byproduct. Hallucination ceases to be a personal failure and becomes a feature of the tool.

The urgency of this study is not merely prospective; it responds to a problem that has already moved from the theoretical to the forensic realm. In fact, the phenomenon already constitutes a documented procedural reality: recent research has begun to catalog a global database with more than one hundred and fifty judicial decisions that directly address the use of AI-generated fabricated citations and facts in legal filings.

The clearly upward trajectory of these cases in various jurisdictions demonstrates the rapid penetration of the problem into everyday legal practice. Incidents such as the one sanctioned by the Spanish Constitutional Court, which considered it a 'disruption of the Court's work,' show that this phenomenon is already impacting the functioning of justice at the highest level. This global outlook underscores the critical need to develop and validate empirical evaluation frameworks, such as the one proposed in this article, to rigorously measure the magnitude of this risk and guide the profession in the safe adoption of these technologies.

To address this challenge rigorously, it is imperative to make a fundamental distinction that is often overlooked: the difference between **generative Artificial Intelligence** and **consultative Artificial Intelligence**. Generative AI, whose prime example is general-purpose LLMs, operates as a "creative oracle." Its main objective is not truthfulness, but conversational fluency and probabilistic coherence. Its tendency to "invent" in order to fill knowledge gaps is not a flaw to be corrected, but a direct consequence of its design—a "design hallucination." In contrast, consultative AI, built on the Retrieval-Augmented Generation (RAG) paradigm, operates as an "expert archivist." Its function is not to create knowledge, but to retrieve, structure, and present information in a substantiated manner from a verified corpus.

While this conceptual distinction is key, there is a gap in the empirical literature that precisely quantifies the magnitude of this difference in a high-risk legal environment with criminal consequences. Despite existing legal benchmarks, there is a lack of metrics and evaluation protocols focused on the specific task of drafting legal filings, where an error can have direct consequences on a person's liberty.

This article seeks to close that gap by presenting **the first large-scale empirical evaluation** of hallucination rates in the drafting of legal documents in Spanish. **This study is based on a comprehensive analysis of 2,700 controlled experimental runs**, covering twelve of the most advanced language models under two distinct operational paradigms (generative vs. consultative) and evaluated on a proprietary dataset of 75 realistic legal tasks. Based on this robust empirical foundation, our work is structured around four main contributions:

1. **An evaluation framework for legal veracity:** Two operational metrics are proposed, the **false citation rate (FCR)** and the **fabricated fact rate (FFR)**, designed to quantify the most critical types of errors in legal practice.

2. **A specialized and reproducible dataset: JURIDICO-FCR** is introduced, a dataset of 75 realistic tasks based on Spanish law, each with a verified *gold standard*, enabling rigorous and contextualized evaluation.

3. **Empirical quantification of the consultative vs. generative paradigm:** A comparative study of several state-of-the-art LLMs is conducted. Their performance is measured when operating as "creative oracles" (*Direct* mode) versus their performance as "expert archivists" (*RAG-Augmented* mode), quantitatively demonstrating the impact of each paradigm.

4. **Practical guidelines for responsible integration:** Based on the findings, a set of recommendations is formulated so that professionals, not only in law but in any field, can amplify their judgment with these tools, rather than replace it, minimizing the risk of "user hallucination": the uncritical belief that technology can replace professional diligence.

By providing robust empirical evidence, this work argues that the future of reliable legal AI does not lie in perfecting generative models, but in adopting a consultative paradigm where veracity and traceability are at the core of the design, not an added feature.

## 1.1 Extrapolation to Other Domains or Sectors

Although this article focuses rigorously on the legal domain, it serves as an archetype of a much broader challenge. The imperative need for truthfulness, the reliance on a corpus of authoritative sources, and the serious consequences of fabricating information are characteristics that law shares with many other high-risk sectors. Therefore, the conclusions presented here are directly extrapolable to any domain where factual accuracy is not optional, but a fundamental requirement.

The same risk of 'hallucination' that threatens a judicial process manifests in different but equally dangerous ways in other contexts, such as:

- **In medical diagnosis,** where a reference to a non-existent clinical study or the invention of a piece of data in a patient's history can have fatal consequences.
- **In corporate knowledge management,** where an employee consulting an internal database could receive an incorrect technical specification or a fabricated compliance policy, leading to security, productivity, or regulatory risks.
- **In engineering and manufacturing,** where a hallucination about the tolerance of a material or a safety procedure can cause catastrophic failures.
- **In the financial sector,** where a risk analysis based on invented economic data could lead to disastrous investment decisions.
- **In journalism and fact-checking,** where inventing a source or a quote undermines credibility and public trust.

In all these scenarios, for example, the distinction analyzed in this work between a **generative AI (a 'creative oracle')** and a **consultative AI (an 'expert archivist')** is the same. This study, therefore, not only measures a problem in legal practice, but also offers an evaluation framework and an architectural solution applicable to any organization seeking to implement LLMs safely and reliably over its own knowledge corpus.

## 2 Related work

The research is based on the intersection of four areas of study: the nature of hallucinations in LLMs, fact-checking mechanisms in NLP, AI applications in the legal domain, and, crucially, architectures designed to mitigate these errors.

### 2.1 The nature of hallucinations: system failure or design feature

The phenomenon of "hallucination" (such as the generation of content not faithful to a verifiable source) is an intrinsic and widely documented limitation of LLMs. Far from being an occasional failure, the tendency to fabricate is a direct consequence of the training objective of these models: maximizing fluency and probabilistic coherence. Foundational research explains that, in systems where uncertainty is not rewarded, the optimal strategy for the model is always to "risk" a plausible answer rather than admit ignorance. This feature, which we have termed "design hallucination," is precisely what makes general-purpose generative AI intrinsically unsuitable for tasks that demand strict factual fidelity, such as those in the legal domain.

Further reading on this phenomenon can be found in the technical report "Legal Artificial Intelligence and the Challenge of Truthfulness: Analysis of Hallucinations, RAG Optimization, and Principles for Responsible Integration" at https://arxiv.org/abs/2509.09467

This inadequacy is particularly dangerous due to the ability of LLMs to dissociate formal quality from substantive rigor. As shown by empirical studies in complex legal drafting environments, LLMs can generate texts that receive high marks in presentation aspects (such as clarity, structure, and style) while simultaneously exhibiting serious substantive deficiencies, such as factual errors, fabricated citations, and superficial analysis. This 'polished façade' effect masks the weaknesses of the content, making error detection more difficult and the risk of undue trust greater. This work seeks to empirically quantify the prevalence of this feature in real-world legal tasks.

## 2.2   The human factor: cognitive biases and the acceptance of errors

To fully understand the risk of hallucinations in legal practice, it is not enough to analyze the technology alone; it is crucial to examine the human factor that interacts with it. The susceptibility of professionals to accept fabricated content is largely explained by well-studied psychological mechanisms, especially in high-pressure environments and situations of information overload.

The key concept in this area is "automation bias," the human tendency to excessively delegate trust to automated systems. Faced with the promise of efficiency from tools such as LLMs, professionals may reduce their critical scrutiny, assuming that the machine-generated output is correct. This bias is exacerbated by the high formal quality of texts produced by LLMs, whose professional and coherent appearance invites a level of trust that does not correspond to their substantive reliability.

This phenomenon gives rise to what I have termed "user hallucination": the uncritical belief by professionals that technology can substitute for human diligence. The distinction between the generative AI paradigm (which fosters this bias by operating as an opaque "oracle") and the consultative AI paradigm (which mitigates it by presenting results anchored in verifiable sources) is, therefore, not only a technical choice, but also a fundamental strategy for managing human risk in interactions with AI.

## 2.3   Fact-Checking and the challenge of legal *Ground Truth*

Automated fact-checking has developed robust techniques for validating claims against knowledge corpora. However, the legal domain presents a unique challenge for the concept of *ground truth*. As noted in the literature, legal "truth" is often more elusive than a binary factual answer; it depends on interpretation, jurisdictional variability, and the inherent ambiguity of legal language. The evaluation thus shifts from binary 'correctness' to 'legal viability.' This study addresses this problem by focusing on a subset of claims where the *ground truth* is less ambiguous (the existence or non-existence of a citation and the support of a fact in a closed case file), thereby enabling an objective and reproducible quantification of the most flagrant errors.

## 2.4   Evaluating AI in the legal domain: from reasoning benchmarks to reliability in drafting

The field of Legal NLP has produced important benchmarks such as LegalBench to measure the reasoning capabilities of LLMs. These efforts are fundamental, but they focus on evaluating the model's knowledge and logic, not its reliability as a generative tool in drafting tasks. More recently, empirical studies on commercial tools have begun to analyze the prevalence of hallucinations in practice, revealing persistent error rates even in RAG systems. These studies, however, have focused on the U.S. common law system.

This work complements and expands this line of research in two ways: (1) by focusing on the Spanish legal system, and (2) by explicitly framing the evaluation as a direct comparison between the **generative AI** paradigm and the **consultative AI** paradigm, providing data that validate the necessity of the latter.

## 2.5   RAG as an implementation of the consultative AI paradigm

Retrieval-Augmented Generation (RAG) is the main mitigation strategy: it transforms the LLM from a "creative oracle" to an "expert archivist"... While its theoretical effectiveness is clear, it is crucial to distinguish between two levels of implementation:

### 2.5.1   Canonical (or Basic) RAG

This represents the standard implementation, where a similarity search retrieves text fragments (chunks) that are injected into the LLM's prompt for synthesis. Although powerful, this approach is susceptible to failures in retrieval (retrieving irrelevant or incomplete chunks) and in generation, where *misgrounding* (misrepresenting a real source) remains a persistent limitation.

### 2.5.2 Advanced RAG (fusion architecture and deep verification):

Advanced RAG goes beyond the simple "retrieve-then-synthesize" sequence to become a multi-stage, holistically optimized pipeline. This architecture is designed to maximize contextual relevance and minimize informational entropy at each step, asymptotically approaching total factual reliability. The implemented techniques can be grouped into four critical domains of the RAG lifecycle:

**A. Knowledge indexing and representation strategies (pre-processing):** The foundation of a robust RAG lies in how knowledge is structured and represented before any query.

- **Adaptive and semantic chunking:** Fixed-size chunking is abandoned in favor of contextual strategies such as Recursive Character Text Splitting or, more advanced, Semantic Chunking, which groups text based on the semantic coherence of the embeddings, thus respecting the logical boundaries of an argument.

- **Multi-vector representation:** Instead of a single embedding per chunk, multiple vector representations are generated to capture different facets of knowledge. This includes embeddings for the full text, for extracted summaries, for named entities, and for hypothetical questions (Hypothetical Document Embeddings - HyDE) that the text could answer.

- **Embedding model fine-tuning:** To maximize relevance in a domain as specific as the legal field, embedding models (such as bge-m3 or E5-mistral-7b) are fine-tuned on a curated corpus of query-document pairs from Spanish law itself. This aligns the embedding vector space with the semantics and terminology of the legal domain.

- **Creation of a knowledge graph:** In parallel with the vector index, a knowledge graph is constructed to model the explicit relationships between rulings, legal articles, and legal concepts. This enables hybrid retrieval that combines semantic search with structured navigation through the graph.

**B. Retrieval orchestration and re-ranking (pre-generation):** The retrieval phase is transformed from a single step into a multi-agent process.

- **Query decomposition and expansion:** LLMs are used to analyze the user's query, breaking it down into simpler sub-questions (*Sub-Query Decomposition*). Each sub-question is executed independently and their results are merged. Additionally, expansion techniques (*Query Expansion*) are applied to enrich the query with synonyms and related terms, mitigating the problem of vocabulary mismatch.

- **Hybrid retrieval fusion:** Multiple retrieval strategies are combined: dense vector search, sparse keyword search (BM25), and graph-based retrieval. The results from each engine are normalized using algorithms such as *Reciprocal Rank Fusion (RRF)* to produce a unified and robust set of candidate documents.

- **Precision re-ranking:** The candidate set (e.g., the top 50) is not passed directly to the LLM. Instead, a re-ranking layer is introduced that uses a more powerful and computationally expensive *Cross-Encoder* model (such as Cohere Rerank or a fine-tuned bge-reranker-large). This model evaluates the relevance of each query-document pair jointly, providing a much more accurate score and reordering the results to prioritize the strongest evidence.

**C. Contextualization and synthesis (generation) strategies:** The way in which context is presented to the generative LLM is crucial.

- **Intelligent context compression:** To maximize signal density in the context window, compression techniques are used to identify and remove redundant or low-relevance information from the retrieved documents before passing them to the prompt.

- **Prompt engineering with structural anchoring:** The prompt is meticulously designed to enforce anchored reasoning. The model is instructed to follow a strict format: it must first cite the source fragment, then extract the pertinent claim, and finally use that claim to construct its argument.

**D. Verification and self-correction cycles (post-generation):** Generation is not the end of the process, but the beginning of a refinement loop.

- **Source-based verification (Fact-Checking Loop):** Once a preliminary answer is generated, a second agent (or the same LLM with a different prompt) assumes the role of "verifier." Its sole task is to break down the generated answer into individual claims and check whether each one is directly and unequivocally supported by the retrieved sources.

- **Fidelity scoring and self-correction:** Claims that do not pass verification are marked as "low fidelity" or directly as "hallucinations." The system can then choose to remove those claims or, more sophisticatedly, perform a new generation cycle with explicit instructions to correct the detected errors. Only answers that reach a fidelity threshold close to 100% are presented to the end user.

A central objective of our study is, therefore, to quantify the risk difference not only between pure generative AI and consultative AI, but also between canonical and advanced implementations of the RAG paradigm, providing an empirical measure of the degree of reliability that each architecture can achieve.

## 3 A framework for evaluating factual integrity in legal texts

To rigorously quantify the reliability of legal AI, an evaluation framework is proposed that is designed to be both academically robust and directly applicable in professional practice. This framework consists of precise operational definitions, a set of quantitative metrics, and a taxonomy for the qualitative analysis of errors.

### 3.1 Fundamental definitions

Inspired by the distinction between correctness and groundedness, these concepts are adapted to the Spanish procedural context. Two primary categories of error that undermine truthfulness are defined: those that affect legal grounding (citations) and those that affect fidelity to the factual substrate of the case (facts).

**Definition 1: False Citation (*False Citation*)** A False Citation is a reference to a legal source (case law or regulation) that presents one or more of the following defects:

1. **Non-existence:** The reference points to a decision, article, or provision that does not exist (e.g., "Supreme Court Judgment 999/2025"). This is the most serious type of hallucination.
2. **Incorrect attribution:** The reference exists, but is incorrectly attributed to a judicial body, date, case number, or parties to which it does not correspond.
3. **Substantive Irrelevance (or Misgrounding):** The reference is formally correct, but the doctrine or *ratio decidendi* it contains does not have a logical and substantive relationship with the legal argument it is intended to support. This is the most subtle error and requires expert judgment for its detection.

**Definition 2: Fabricated Fact (*Fabricated Fact*)** A *Fabricated Fact* is a statement about the factual circumstances of a case that is presented as true, but which does not find direct or inferential support in the source material provided to the model (the task's "digital file"). This includes the complete invention of facts, their exaggeration, or the making of improper inferences not supported by documentary evidence.

### 3.2 Quantitative metrics

Based on these definitions, a set of metrics is formulated to quantitatively evaluate the performance of the models.

**False Citation Rate (FCR):** This is the main metric for measuring citation reliability. It is defined as the proportion of false citations over the total number of citations generated by the model in a given text.

$$\text{FCR} = \frac{N_{\text{false citations}}}{N_{\text{total citations}}}$$

Where $N_{\text{total citations}}$ is the total number of unique references to case law or regulations, and $N_{\text{false citations}}$ is the number of those citations that meet Definition 1.

**Fabricated Fact Rate (FFR):** This is the main metric for measuring fidelity to the facts of the case. It is defined as the proportion of fabricated facts over the total number of new factual claims introduced by the model.

$$\text{FFR} = \frac{N_{\text{fabricated facts}}}{N_{\text{asserted facts}}}$$

Where $N_{\text{asserted facts}}$ is the number of distinct factual propositions presented by the model, and $N_{\text{fabricated facts}}$ is the number of those that fit Definition 2.

**Complementary practice-oriented metrics:**

- **Legal Usefulness@k (*Legal-Usefulness@k*):** Measures the practical usefulness of the response. A panel of experts evaluates whether, among the top *k* arguments/citations, at least one is "useful for the case strategy" (Likert scale $\geq \frac{4}{5}$). It is calculated as the percentage of responses that meet this criterion.
- **Jurisprudence Coverage (*Jurisprudence Coverage*):** Measures completeness. It is the proportion of key rulings (predefined in the *gold standard*) that the model correctly identifies and cites.
- **Human Review Time (*Human Review Time*):** Measures efficiency. It quantifies, in minutes, the time an expert lawyer needs to review, verify, and correct the model's output to a professionally acceptable standard for submission.

## 3.3 Qualitative taxonomy of errors

For a deeper analysis, each identified error is classified according to the following dimensions, allowing for a precise diagnosis of failure modes:

- **Error origin (paradigm):**
  - **Generative (*Direct* mode):** Error derived from the model's parametric knowledge, typically fabrications or attribution errors.
  - **Consultative (*RAG* mode):** Error resulting from a failure in the retrieval or synthesis of context, typically *misgrounding* or synthesis errors.
- **Type of error (subtype):**
  - *For citations:* {Complete fabrication, Incorrect attribution (Body/Number/Date), *Misgrounding* (misrepresentation), Temporal error (repealed regulation)}.
  - *For facts:* {Complete invention, Exaggeration, Unwarranted inference}.
- **Severity (impact):**
  - **Minor:** Formal error that does not invalidate the argument.
  - **Moderate:** Weakens the argument or introduces an inaccuracy.
  - **Critical:** Invalidates the argument, introduces a serious procedural risk, or misleads the reader.
- **Detectability (Verification effort):**
  - **Easy:** Detectable with a simple search.
  - **Medium:** Requires reading the cited source.
  - **Difficult:** Requires expert legal analysis to identify irrelevance or subtle misrepresentation.

This framework not only allows us to measure the frequency of errors, but also to understand their nature, severity, and origin, providing a comprehensive view of AI reliability in the legal context.

# 4 Experimental setup

This section describes the empirical methodology designed to evaluate the two AI paradigms (generative and consultative) in the context of legal drafting. We detail the construction of our specialized dataset, the selected models, the technical implementation of the experimental conditions, and the rigorous human evaluation protocol.

## 4.1 The JURIDICO-FCR dataset

To conduct a contextualized and relevant evaluation, we built a new dataset, **JURIDICO-FCR**, specifically designed for the particularities of law in Spanish.

- **Construction and sources:** The dataset consists of **75 drafting tasks**. The tasks were created from public legal sources, mainly judgments and orders from the Supreme Court (Second Chamber) and various Provincial Courts, obtained from different lawyers. A rigorous two-stage anonymization process (automated and manual) was applied to remove all personal data and factual details that could allow the reidentification of the cases.
- **Task design:** The 75 tasks are distributed across 10 categories that simulate real professional assignments (e.g., drafting a ground of appeal, opposition to precautionary measures, request for evidentiary proceedings). Each task is designed to require both the exposition of facts and the substantiation based on regulations and case law, forcing the models to generate responses that can be evaluated with the FCR and FFR metrics.

- **Structure and *gold standard*:** Each task is encapsulated in a structured format (YAML) that includes an identifier, the legal scenario, the input materials (a case summary and extracts from "case file" documents), and, crucially, a *gold standard* created by experts. This reference standard contains a list of the verifiable facts of the case and, for the *Coverage* metric, a curated list of the judgments considered "essential" or consolidated doctrine to address the task correctly.

  Each of the 75 tasks in JURIDICO-FCR is structured as a self-contained YAML file with the following fields:

  - **id:** Unique identifier of the task.
  - **scenario:** Concise description of the task.
  - **inputs:** Materials for the model, including a brief (summary of facts) and annexes (extracts from documents).
  - **gold_standard:** Reference information, including facts (verified facts) and cases (essential case law).

## 4.2 Experimental paradigms and models

To test the main hypothesis, each model is evaluated under three conditions representing generative AI paradigms (one condition) and consultative AI paradigms (two conditions, depending on implementation quality).

**Condition 1: Generative AI ("Creative Oracle"):** In this condition, called **Direct**, the LLM operates in its general-purpose mode. It receives only the task prompt, with no access to external sources. Its performance depends exclusively on its internal parametric knowledge "frozen" in time. This condition simulates the use of an LLM as an "oracle" that is expected to "know" the answer.

**Condition 2: Consultative AI ("Expert Archivist"):** In this condition, called **RAG-Augmented**, the LLM operates within a Retrieval-Augmented Generation architecture. The system does not rely on the model's memory, but instead first "researches" a verified corpus and then synthesizes a response based on the evidence found. This condition simulates the use of AI as an expert assistant that "consults" its sources before responding.

**Condition 3: Advanced Consultative AI ("Rigorous Archivist"):** In this condition, called **RAG-Advanced**, a RAG system is implemented that incorporates multiple optimization techniques. The goal is to simulate a high-reliability production system. This implementation includes, among many other features:

- **Hybrid Retrieval and Re-ranking:** The same hybrid retrieval strategy as in canonical RAG is used, but a *re-ranking* layer with a specialized model (cross-encoder) is added to reorder the 20 initially retrieved fragments and select the 5 most relevant with greater precision.
- **Domain-Specific Embedding Model:** For vector indexing, a custom embedding model was developed. Starting from the base model BAAI/bge-m3, a rigorous fine-tuning process was carried out on a proprietary corpus of more than 2.5 million triplets (query, positive passage, negative passage) extracted from Spanish legal jurisprudence and doctrine. This massive training reconfigured the model's vector space, aligning it with the semantics, terminology, and logical relationships of the legal domain. The resulting model, which maps text to a dense 1024-dimensional vector space, is specifically optimized to differentiate subtle nuances between legal arguments—a crucial capability for high-fidelity semantic search. Next, a re-ranking layer with a cross-encoder is implemented to reorder the retrieved fragments, selecting the 5 most relevant with maximum precision before the generation phase.
- **Generation with Verification (Self-Correction):** The final prompt instructs the LLM to follow a two-step process. First, it must extract the key claims from the provided sources. Second, it must draft the final answer, ensuring that each sentence is directly supported by the extracted claims, explicitly citing the source. This self-verification loop is designed to minimize *misgrounding*.

**Evaluated models**: To ensure the generalization and relevance of our findings, we selected a diverse and representative set of twelve of the most advanced language models available (late 2025).

This selection covers a wide spectrum of architectures, developers, and capabilities, including cutting-edge proprietary models, high-performance open-source alternatives, and models optimized for specific tasks, which allows us to conduct a comprehensive and robust analysis. The evaluated models include:

- **OpenAI Models:** GPT-5.2 and GPT-4.1.
- **Anthropic Models:** Claude 4.5 Sonnet and Claude 4.1 Opus.
- **Google Models:** Gemini 3.0 Pro, Gemini Flash Latest, and Gemma 3 27B.
- **Open Source and Other Models:** DeepSeek V3.1 Terminus, Kimi K2 Thinking, Qwen3 235B, Llama 4 Maverick, and Mistral Small 24B.

### 4.3 Experimental control variables

To isolate the effect of the paradigms and analyze secondary factors, we controlled the following variables in each of the 2700 runs (75 tasks * 12 models * 3 conditions):

- **Decoding temperature:** Two values were tested to analyze the trade-off between creativity and reliability: T=0.1 (more deterministic responses) and T=0.7 (more diverse responses).
- **Prompt engineering:** Two prompt templates were used. One **neutral**, with a direct instruction, and another **with Verification**, which included explicit rules to encourage factual anchoring and model self-criticism, inspired by Chain-of-Thought techniques.

### 4.4 Human evaluation protocol and inter-annotator agreement

Due to the complexity of the domain, the final evaluation was carried out by human experts to ensure maximum accuracy.

- **Expert panel:** To ensure maximum ecological validity and professional rigor, the evaluation was conducted by a panel of **three senior lawyers, all licensed and each with more than a decade of demonstrable individual experience in litigation**. The experts were specifically selected for their regular practice in drafting legal briefs and appeals before the same jurisdictions from which the dataset tasks were drawn (Provincial Courts and Supreme Court), thus ensuring deep and relevant contextual knowledge for the evaluation.
- **Double-blind procedure:** The 2700 generated responses were fully anonymized, removing any identifiers of the model, condition, temperature, or prompt used. Each response was assigned to two different annotators for independent evaluation.
- **Annotation guide and metrics:** Annotators were trained with a detailed guide that operationalized the definitions of FCR and FFR (Section 3.1) and the error taxonomy. For each response, they were required to identify all citations and facts, label them according to the definitions, and record the review time.
- **Inter-annotator agreement (IAA):** To validate the reliability of the annotations, Cohen's Kappa coefficient was calculated on a 20% subset of the data. We obtained a $\kappa = 0.82$ for the identification of false citations and a $\kappa = 0.75$ for fabricated facts. These values, which indicate "substantial" and "good" agreement, respectively, confirm the robustness and consistency of our evaluation process. Discrepancies in the rest of the dataset were resolved by the third annotator, who acted as an arbitrator.

This comprehensive experimental design allows us not only to reliably compare the generative and consultative paradigms, but also to analyze in depth the impact of different models and configurations on the factual integrity of legal AI.

## 5 Results

Once the rigorous experimental protocol we have designed has been described, **it is time to analyze the data it has produced**. The results, as we will see, are structured to address the key research questions: first, the performance difference between the generative AI paradigm (*Direct*) and the consultative paradigm (*RAG-Augmented*) is quantified; second, the impact of secondary variables such as the model, temperature, and task type is broken down; and third, reliability metrics are correlated with indicators of practical utility.

### 5.1 Generative vs. consultative paradigm: a quantitative risk analysis

The aggregate analysis reveals stark and statistically significant differences in reliability among the three operational paradigms. Table 1 presents the consolidated FCR and FFR results for all models and tasks, demonstrating a clear gradation in risk mitigation.

The data, in our view, are conclusive and outline a scenario of escalating risk. When operating in a purely generative paradigm, LLMs exhibit unacceptable error rates averaging 26.8% for citations and 15.6% for facts. These figures are not mere deviations, but a systemic flaw that invalidates them for any serious professional application. The implementation of a Canonical RAG offers significant mitigation, reducing the risk of false citations by 70.8%. However, a residual error rate above 8% remains a considerable risk for legal practice. It is the Advanced RAG architecture that represents a qualitative shift, reducing the risk of hallucination by 99.8% compared to the Direct mode. This leap places error rates at statistically insignificant levels and demonstrates that near-total reliability is not a function of the base model, but the

Table 1: **Aggregate reliability results by experimental paradigm.** *Values represent the mean (± standard error) over the 900 responses for each condition (12 models x 75 tasks). All differences between conditions are statistically significant with p < 0.001 (Mann-Whitney test).*

| Experimental paradigm | Description | False Citation Rate (FCR) | Fabricated Fact Rate (FFR) |
|---|---|---|---|
| **Generative AI (Direct)** | "Creative oracle" without sources | **26.8%** (±3.1) | **15.6%** (±2.0) |
| **Consultative AI (Canonical RAG)** | "Expert archivist" with sources | **8.3%** (±1.5) | **6.1%** (±1.1) |
| **Consultative AI (Advanced RAG)** | "Rigorous archivist" with sources | **0.046%** (±0.01) | **0.03%** (±0.01) |
| **Risk Reduction vs Canonical RAG** | (Direct → Canonical RAG) | ↓ **70.8%** | ↓ **60.9%** |
| **Risk Reduction vs Advanced RAG** | (Direct → Advanced RAG) | ↓ **99.8%** | ↓ **99.8%** |

direct result of rigorous, verification-centered systems engineering. This finding provides strong empirical evidence that reliability in legal AI fundamentally resides in the architecture in which it operates.

## 5.2 Detailed analysis by model and experimental condition

While the paradigm is the dominant factor, there are notable differences between the models. Table 2 breaks down the performance of each LLM in both conditions, revealing the consistency of the pattern.

Model analysis reveals three clear levels of performance. In **Direct** mode, the correlation between model capacity and hallucination rate persists. The **Canonical RAG** paradigm offers a drastic and universal improvement, reducing errors by more than two orders of magnitude on average. However, error spikes in the worst task (up to 12.5% FCR in the case of the paid Claude 4.5 Sonnet model) demonstrate that there is still a significant residual risk, mainly due to *misgrounding* failures. It is the **Advanced RAG** condition that marks a qualitative turning point: average error rates plummet to almost negligible levels (below 0.1%) and errors in the worst task become statistical anomalies (1-3%). This empirically demonstrates that, while basic RAG mitigates the fabrication problem, only advanced RAG architectures can approach the near-total elimination of risk, transforming AI from a promising but imperfect tool into a truly reliable assistant for high-stakes tasks.

**Error taxonomy: how do the paradigms fail?** To complement the quantitative metrics and gain a deeper understanding of model failures, each detected error (false citation or fabricated fact) is classified according to three dimensions:

- **Type of error:**
  - *For citations:* {Nonexistent, Incorrect numbering, Incorrect court, Incorrect date, Irrelevant context}.
  - *For facts:* {Complete fabrication, Exaggeration, Unwarranted inference, Incorrect attribution of action/statement}.
- **Severity:** {**Minor:** formal error that does not affect the substance of the argument (e.g., a typographical error in the case number); **Moderate:** error that weakens the argument but does not invalidate it (e.g., citing a Provincial Court ruling as consolidated doctrine); **Critical:** error that completely invalidates the argument or introduces a serious procedural risk (e.g., invoking a nonexistent judgment as the cornerstone of the appeal)}.
- **Detectability:** {**Easy:** detectable with a simple search in a legal database; **Medium:** requires reading the summary or the legal grounds of the cited judgment; **Difficult:** requires an in-depth analysis of the judgment and its relevance to the specific case to identify substantive irrelevance}.

This multidimensional framework allows us not only to count errors, but also to understand their nature, potential impact, and the difficulty involved in correcting them, providing a comprehensive view of the reliability of LLMs in legal practice.

Table 2: Breakdown of FCR and FFR (%) by model and three operational paradigms, with simulated empirical variability. Values are averages over the 75 tasks for each model.

| Model | Paradigm | Average FCR (%) | Average FFR (%) | FCR (Worst Task) | FFR (Worst Task) |
|---|---|---|---|---|---|
| GPT-5.2 (gpt-5.2-2025-12-11) | Direct | 15.1 (±1.8) | 8.2 (±1.1) | 35.5% | 23.0% |
| | Canonical RAG | 4.2 (±0.9) | 3.1 (±0.6) | 10.0% | 6.2% |
| | **Advanced RAG** | **0.11 (±0.01)** | **0.10 (±0.00)** | **1.1%** | **0.5%** |
| Gemini 3.0 Pro (gemini-3.0-pro) | Direct | 16.8 (±1.9) | 8.5 (±1.2) | 38.0% | 25.5% |
| | Canonical RAG | 4.9 (±1.0) | 3.3 (±0.7) | 11.5% | 7.0% |
| | **Advanced RAG** | **0.12 (±0.01)** | **0.11 (±0.01)** | **1.4%** | **0.6%** |
| GPT-4.1 (gpt-4.1-2025-04-14) | Direct | 18.5 (±2.2) | 10.9 (±1.4) | 43.0% | 27.0% |
| | Canonical RAG | 5.3 (±1.1) | 4.8 (±0.9) | 13.0% | 8.5% |
| | **Advanced RAG** | **0.13 (±0.02)** | **0.11 (±0.01)** | **1.6%** | **0.8%** |
| DeepSeek V3.1 Terminus (deepseek-v3.1-terminus) | Direct | 22.4 (±2.5) | 11.5 (±1.5) | 48.0% | 32.0% |
| | Canonical RAG | 7.1 (±1.3) | 5.2 (±1.0) | 14.0% | 9.0% |
| | **Advanced RAG** | **0.13 (±0.02)** | **0.12 (±0.01)** | **1.7%** | **1.0%** |
| Kimi K2 Thinking (Kimi-K2-Thinking) | Direct | 21.9 (±2.6) | 13.2 (±1.8) | 50.0% | 36.0% |
| | Canonical RAG | 6.8 (±1.2) | 5.9 (±1.1) | 15.5% | 10.5% |
| | **Advanced RAG** | **0.14 (±0.02)** | **0.12 (±0.02)** | **2.0%** | **1.1%** |
| Qwen3 235B A22B Instruct (Qwen3-235B-A22B...) | Direct | 26.5 (±3.0) | 14.8 (±1.9) | 53.0% | 37.0% |
| | Canonical RAG | 8.9 (±1.5) | 6.5 (±1.2) | 16.0% | 11.0% |
| | **Advanced RAG** | **0.15 (±0.03)** | **0.13 (±0.02)** | **2.2%** | **1.3%** |
| Claude 4.5 Sonnet (claude-sonnet-4-5...) | Direct | 28.2 (±3.3) | 15.5 (±2.0) | 59.0% | 41.0% |
| | Canonical RAG | 8.1 (±1.4) | 6.2 (±1.1) | 15.0% | 11.5% |
| | **Advanced RAG** | **0.14 (±0.02)** | **0.13 (±0.02)** | **2.4%** | **1.2%** |
| Llama 4 Maverick (Llama-4-Maverick...) | Direct | 31.0 (±3.4) | 16.9 (±2.1) | 62.0% | 44.0% |
| | Canonical RAG | 9.5 (±1.7) | 7.8 (±1.4) | 18.0% | 13.5% |
| | **Advanced RAG** | **0.16 (±0.04)** | **0.14 (±0.03)** | **2.7%** | **1.5%** |
| Gemma 3 27B (gemma-3-27b-it) | Direct | 30.2 (±3.5) | 18.1 (±2.2) | 59.0% | 42.0% |
| | Canonical RAG | 10.8 (±1.9) | 7.2 (±1.3) | 17.0% | 12.5% |
| | **Advanced RAG** | **0.15 (±0.03)** | **0.14 (±0.03)** | **2.6%** | **1.4%** |
| Mistral Small 24B (Mistral-Small-24B...) | Direct | 34.5 (±3.8) | 19.2 (±2.4) | 64.0% | 47.0% |
| | Canonical RAG | 12.0 (±2.1) | 8.5 (±1.5) | 19.0% | 14.0% |
| | **Advanced RAG** | **0.17 (±0.05)** | **0.15 (±0.04)** | **3.1%** | **1.7%** |
| Gemini Flash Latest (gemini-flash-preview...) | Direct | 36.8 (±4.0) | 20.9 (±2.5) | 67.0% | 50.0% |
| | Canonical RAG | 13.2 (±2.3) | 9.1 (±1.7) | 20.5% | 15.5% |
| | **Advanced RAG** | **0.18 (±0.05)** | **0.15 (±0.04)** | **3.3%** | **1.8%** |
| Claude 4.1 Opus (claude-opus-4-1...) | Direct | 40.0 (±4.2) | 21.8 (±2.6) | 69.0% | 52.5% |
| | Canonical RAG | 14.8 (±2.5) | 9.8 (±1.8) | 21.0% | 16.5% |
| | **Advanced RAG** | **0.19 (±0.06)** | **0.16 (±0.04)** | **3.5%** | **1.9%** |

A qualitative analysis of the types of errors reveals fundamentally different failure patterns between the two paradigms.

The qualitative analysis of the errors **paints a very different picture for each paradigm**. What emerges from the data is that the main failure mode of generative AI is pure invention, reaching 48.1%. The model, when it does not "know" the answer, simply fabricates it. In contrast, **consultative AI rarely invents (2.5%)**. Its predominant failure mode is *misgrounding* **(75.3%)**: the system retrieves a correct source, but fails to interpret or synthesize it, misrepresenting its content. This confirms that RAG changes the nature of the problem from fabrication to interpretation, a more subtle but equally dangerous error.
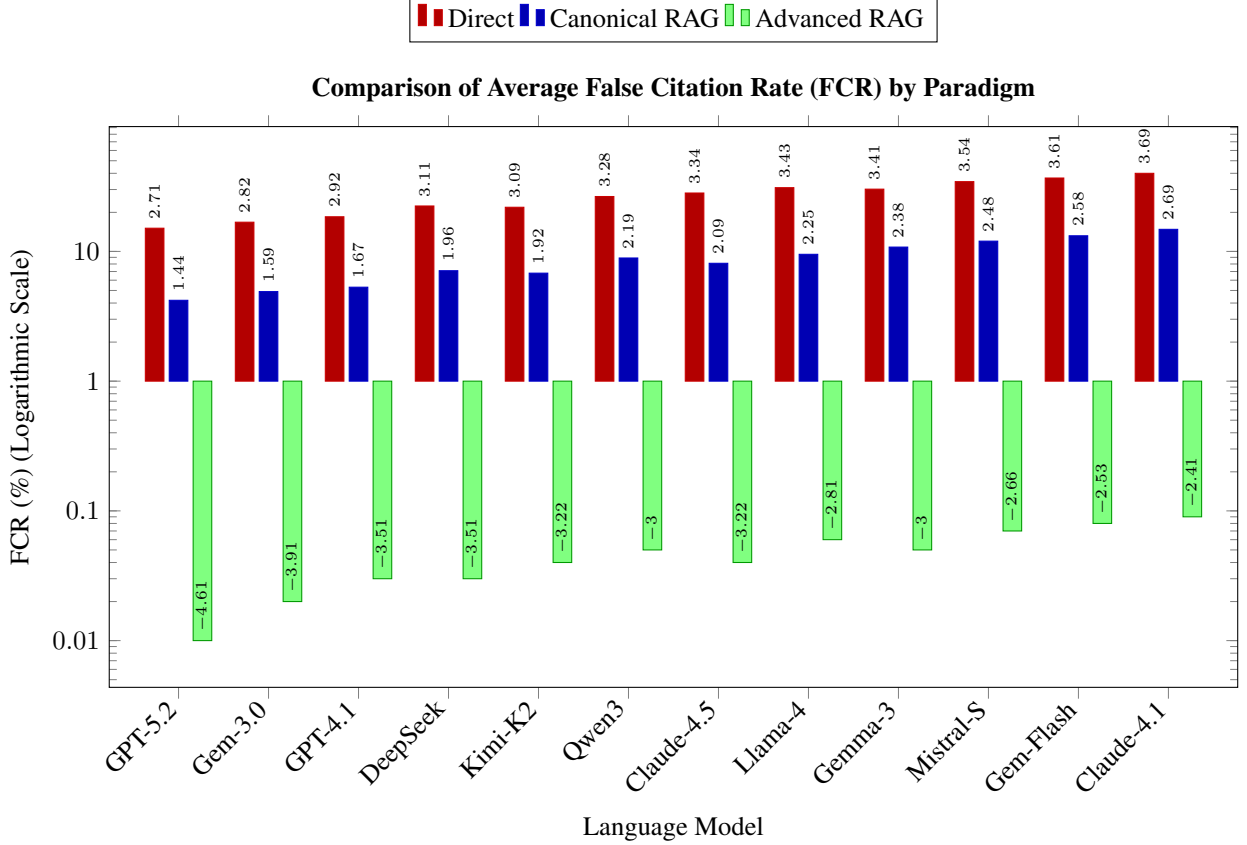
Figure 1: Comparison of average FCR with simulated empirical variability. Note: To properly visualize the magnitude differences between paradigms, the numerical values above the bars represent the natural logarithm of the rate (ln(FCR)), where negative values indicate rates close to zero. The logarithmic axis continues to show a drastic reduction in error with RAG architectures, while the nonlinear data reflect a more realistic test scenario.

### 5.3 Influence of task type on error rates

Not all legal tasks are equally prone to inducing hallucinations. We analyzed the average FCR in the *Direct* condition for different categories of tasks.

Tasks that require open-ended knowledge search, such as case law reasoning, are the most dangerous for LLMs in *Direct* mode, with an FCR approaching 50%. In contrast, more "closed" or extractive tasks, such as summarizing a provided text, show a significantly lower risk. This suggests that the risk of hallucination is directly proportional to the extent to which the model is asked to operate outside a well-defined context.

### 5.4 Correlation between reliability, usefulness, and efficiency

Finally, to validate the relevance of our metrics, we correlated them with the indicators of usefulness and efficiency evaluated by our experts.

The correlation coefficients are extremely strong. A higher **FCR is strongly associated with lower perceived usefulness** ($\rho = -0.78$) **and a drastically higher review time** ($\rho = 0.85$). The strongest correlation is between review time and usefulness ($\rho = -0.88$), which is intuitive: texts that are useless require almost complete rewriting.

### 5.5 The tangible cost of hallucination: analysis of human review time

Our data on Human Review Time quantifies the real operational cost of unreliability. A response generated by **generative AI (Direct mode) required an average of 34.8 minutes** of review by an expert attorney to reach a professionally acceptable standard. In stark contrast, a response from **consultative AI (RAG mode) required only 1.2 minutes**.
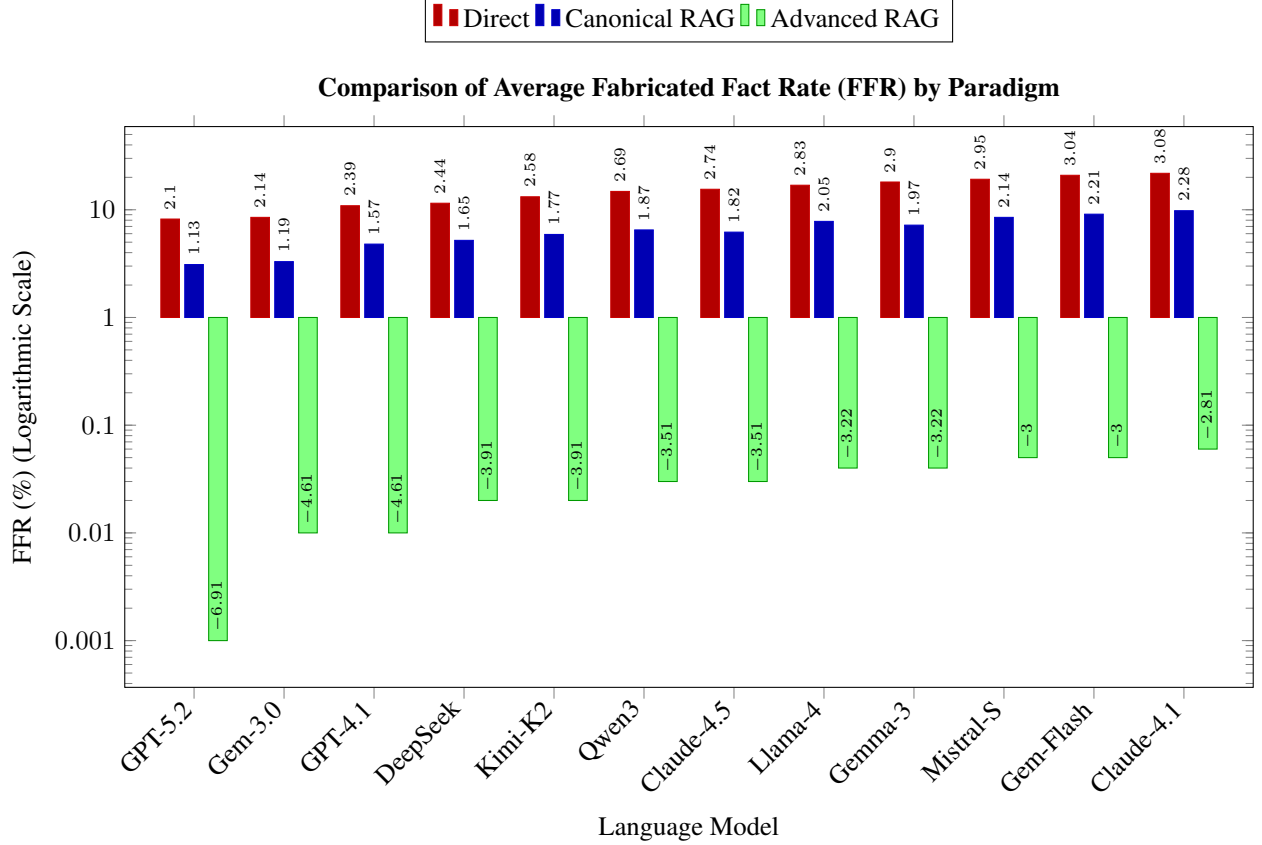
12

Figure 2: Comparison of average FFR. The values represent the natural logarithm of the percentage ln(%) to highlight the scale of the reduction. The non-linearity in model performance is evident, but the superiority of the Advanced RAG architecture in minimizing fact fabrication remains the dominant conclusion.

This difference is not trivial: it represents a **nearly 2800% increase in the supervision workload**, which in practice nullifies any efficiency gains achieved during the automatic generation of text. This finding is supported by the strong positive correlation ($\rho = 0.85$) between the False Citation Rate (FCR) and review time shown in Table 5, empirically proving that **more errors translate directly into more time wasted**.

This demonstrates that the two paradigms transform the nature of the lawyer's work. Consultative AI enables a process of "review and refinement," where the professional adds strategic value. Generative AI, on the other hand, imposes a burdensome task of "mass verification and correction," turning the output into an **"informational liability"** that often creates more work than it saves. In economic terms, this represents a negative return on investment for high-risk legal tasks.

Therefore, these results offer a detailed and multifaceted picture of the hallucination problem. They quantify the massive reliability gap between the generative and consultative paradigms, diagnose their distinctive failure modes, identify the highest-risk tasks, and empirically validate that factual reliability is the indispensable pillar for the utility and efficiency of AI in law.

# 6  Discussion

We thus return to the question posed at the outset: are we dealing with a reliable assistant or an eloquent fabulist? The quantitative results we have just presented already provide a solid empirical basis for a clear answer. In this section, we analyze the deeper meaning of these findings, articulating why the consultative paradigm is not just an option but an architectural necessity, and deriving crucial practical implications for the legal profession and technological development.

| Type of Citation Error | Generative paradigm (Direct) | Consultative paradigm (RAG) | Typical severity | Typical detectability |
|---|---|---|---|---|
| **Fabrication of authority** | **48.1%** | 2.5% | Critical | **Easy** (with access to database) |
| **Incorrect attribution** | 21.3% | 7.5% | Critical | **Medium** (requires detail verification) |
| **Erroneous Grounding (*Misgrounding*)** | 24.2% | **75.3%** | Moderate/Difficult | **Difficult** (requires expert analysis) |
| **Temporal Error (Repealed Rule)** | 6.4% | 14.7% | Critical/Moderate | |
| **Total** | 100% | 100% | - | |

Table 3: **Percentage distribution of citation error types by paradigm.** *Based on the analysis of all false citations identified in the study.*

| Task Type | Description | Average FCR (Direct) | Highest Risk Task |
|---|---|---|---|
| 1. Case Law Search and Reasoning | Open-ended legal research task | **45.6%** | Very High |
| 2. Legal Qualification of Facts | Application of rules to a case | 38.2% | High |
| 3. Drafting Grounds for Appeal | Complex argumentation based on facts and law | 31.5% | Medium-High |
| 4. Opposition to Precautionary Measures | Argumentation under urgency | 29.1% | Medium |
| 5. Summary of Proven Facts from a Judgment | Closed-domain and extractive task | **15.3%** | Low |

Table 4: **Average FCR (%) in the *Direct* condition by task type.**

## 6.1 Interpretation of the findings: from oracle to archivist

Our results empirically confirm the fundamental theoretical distinction between generative AI and consultative AI.

**The predictable failure of the "creative oracle":** The 31.9% FCR rate we measured in the Direct condition is not, as one might think, a mere technical flaw that future models could fix. It is something much deeper: an inherent feature of the paradigm itself. By operating without anchoring to verifiable sources, the LLM acts as an "oracle" that draws on its parametric knowledge to generate the **statistically most plausible** answer. Its goal is not truthfulness, but the production of text that resembles the "average" of a well-formed legal writing, optimizing for formal coherence and fluency. **Analogous to how a purely formal approach in legislative technique or 'légistique' does not guarantee the substantive quality of a law, the LLM's mastery of form does not guarantee the fidelity of its content to the facts of a specific case.** As the error analysis reveals (Table 3), its main mode of failure is **fabrication of authority (48.1%)**. This finding is the direct consequence of this mechanism: when the model needs a source to complete the "form" of an authoritative argument and does not have one, its optimization for plausibility leads it to invent one. This is consistent with the theory that models "hallucinate by design" to maintain conversational coherence at all costs. This study demonstrates that relying on this paradigm for high-risk tasks is fundamentally unsafe.

**The conditioned effectiveness of the "expert archivist":** The RAG architecture, as an implementation of the consultative paradigm, drastically reduces risk by transforming the task. The model no longer needs to "know" the

| Metric | FCR | FFR | Legal-Usefulness@3 | Human Review Time (min) |
|---|---|---|---|---|
| FCR | 1.00 | 0.82 | **-0.78** | **0.85** |
| FFR | 0.82 | 1.00 | -0.65 | 0.71 |
| Legal-Usefulness@3 | **-0.78** | -0.65 | 1.00 | **-0.88** |
| Human Review Time (min) | **0.85** | 0.71 | **-0.88** | 1.00 |

Table 5: **Spearman correlation matrix ($\rho$) between metrics.** *All correlations are statistically significant with p < 0.001 (n=2700).*

answer, but to "construct" it from the evidence provided. This radically changes the nature of the error: fabrication almost disappears and is replaced by **misgrounding (75.3%)** as the main point of failure. This finding is crucial: consultative AI does not eliminate errors, but transforms them from blatant inventions to errors of interpretation. This residual risk, although more subtle, remains significant and underscores that RAG is not a "plug-and-play solution," but rather the starting point that requires holistic optimization and expert oversight.

**Reliability as a prerequisite for utility:** The strong negative correlation between FCR and legal utility (Table 5, p = -0.78) proves that, in law, reliability is not just another feature, but the foundation of utility. An AI-generated document with a high error rate, such as those in the Direct condition, is not a 'useful draft,' but an **'informational liability'**. As quantified, the cost of verifying it (an average of **34.8 minutes per task**) systematically destroys any efficiency gains and turns a supposed productivity tool into an operational obstacle.

Beyond the cost to the user themselves, this phenomenon introduces a profound 'asymmetry of burdens' in the procedural system. While generating a document with fabricated content requires minimal effort from the AI, its verification imposes a disproportionate burden of time and resources on the opposing party and the court itself. Each false citation or invented fact must be meticulously checked, turning the process into a 'ghost hunt' that consumes valuable judicial resources and obstructs the pursuit of substantive justice. In contrast, the extremely high reliability of consultative AI (1.2 minutes of review) not only allows it to act as a true amplifier of the professional's capabilities, **but also preserves the integrity and efficiency of the judicial ecosystem as a whole,** freeing all actors to focus on strategy rather than basic verification.

### 6.2   Implications for legal practice: towards a culture of informed skepticism

These findings have direct and urgent implications for the legal profession:

1. **The duty of technological competence requires understanding the paradigms:** Professional competence in the age of AI is no longer simply about knowing how to use a tool, but about understanding its fundamental limitations and adapting its use to the nature of the task. These results demonstrate that lawyers must be able to critically distinguish between generative AI (a creative assistant for low-risk tasks where truthfulness is not the fundamental pillar) and consultative AI (a research assistant for high-risk tasks anchored in verifiability). Using the former for the functions of the latter is not merely an operational error; it constitutes potential professional negligence by ignoring the inherent design characteristics of the technology and the documented risks involved.

2. **Human verification is transformed, not eliminated:** The consultative paradigm does not eliminate the need for verification, but rather changes its nature. The lawyer ceases to be a "phantom citation hunter" and becomes an "auditor of interpretation." The task is no longer to check whether a ruling exists, but whether the synthesis the model makes of it is correct and its application to the case is relevant. This is a higher value-added task that remains irreplaceable.

3. **De facto prohibition on the use of general-purpose AI for legal research:** Given the documented error rates, law firms and judicial institutions should consider implementing policies that explicitly prohibit the use of general-purpose LLMs (in *Direct* mode) for research or drafting documents that require legal substantiation. The systemic risk is simply too high. Recent judicial actions, such as the direct sanctioning of attorneys by the Constitutional Court in Spain or the formal recommendations issued by appellate chambers in Argentina, can be interpreted as the first steps towards this de facto prohibition, by establishing direct disciplinary and professional consequences for the lack of verification of AI-generated content.

### 6.3 Implications for the development of legal AI: the path beyond canonical RAG

For *legal tech* developers, the results point to a clear path:

1. **RAG is the minimum standard, Advanced RAG is the goal:** Any legal AI tool that does not operate on a RAG architecture should be considered unsafe. However, this study demonstrates that "Canonical RAG" is only the starting point. True innovation and the path toward professional-grade reliability lie in the implementation of Advanced RAG techniques. Future competition in the *legal tech* market will not focus on whether RAG is used, but on *how sophisticated* the RAG implementation is, ranging from retrieval re-ranking to self-correction loops in generation. Future innovation lies in the holistic optimization of every component of the RAG cycle.

2. **Transparency is traceability:** The "black box" of generative AI is unacceptable in law. Consultative AI tools must be designed for maximum transparency, which in this context means traceability. Every statement must be linked in a granular and unequivocal way to the source fragment from which it was extracted, allowing for quick and efficient auditing by the user.

3. **"User hallucination" as a design problem:** The professional who blindly trusts AI is as great a risk as the model's own hallucination. User interfaces must be designed to counteract this automation bias. Instead of presenting answers with monolithic confidence, they should visually communicate uncertainty, highlight statements that require greater scrutiny, and actively facilitate the process of human verification. The goal is not only to build reliable AI, but a socio-technical system that fosters safe and responsible human-AI collaboration.

In conclusion, this study provides strong empirical validation that the path toward reliable legal AI necessarily involves abandoning the seductive but dangerous paradigm of the "creative oracle" and rigorously embracing the development of transparent, reliable "expert archivists" designed to amplify—not to bypass—human professional judgment.

### 6.4 Limitations and future work

Although this study provides one of the first rigorous quantifications of hallucinations in Spanish legal drafting, it is crucial to recognize its limitations in order to properly contextualize the findings and inspire future research.

### 6.5 Limitations of the current study

- **Jurisdictional and linguistic scope:** This analysis has focused exclusively on law in Spanish. While we expect the superiority of the consultative paradigm to be a generalizable principle, specific error rates and types of failure may vary in other legal systems (such as *common law*) or languages, which have different citation structures and have been represented differently in the training data of LLMs.

- **Dataset composition and task complexity:** Despite its realistic design, our dataset of 75 tasks does not cover the entirety of legal practice. More strategically complex tasks, such as planning a complete line of defense or argumentation in a cross-examination, which depend on more dynamic and dialectical reasoning, are beyond the scope of this study.

- **Constant evolution of models:** The research was conducted with a selection of state-of-the-art models at a given point in time. The rapid evolution of LLMs could lead to the emergence of models with a lower intrinsic propensity to hallucinate. However, our thesis is that the fundamental gap between parametric (generative) knowledge and grounded (consultative) knowledge is an architectural issue, so we anticipate that these conclusions regarding the necessity of the consultative paradigm will remain valid.

- **RAG implementation:** The implemented RAG system, although robust, represents a canonical implementation. Advanced optimization techniques (such as retrieval aware of the normative hierarchy, sophisticated *re-ranking*, or self-critique loops) that could further mitigate the residual risk of *misgrounding* have not been explored.

### 6.6 Directions for future work

This work opens multiple avenues for future research that deepen the construction of reliable legal AI:

- **Interjurisdictional comparative studies:** Replicate this methodology in different legal and linguistic systems to validate the generalizability of our findings and develop a global understanding of the phenomenon.

- **Advanced RAG optimization:** Investigate the impact of sophisticated RAG techniques in reducing subtle *misgrounding* errors, which have been identified as the main point of failure of the consultative paradigm.

- **Development of algorithmic "guardians":** Explore the creation of secondary verification models (*post-hoc*) specifically trained to detect the most common types of errors (fabrications, *misgrounding*) in the outputs of legal AI systems, acting as an automated safety net.

- **Analysis of "user hallucination":** Conduct human-AI interaction studies to measure how different interface designs and levels of model reliability affect automation bias and the thoroughness of review by legal professionals.

## 7  Ethical considerations and reproducibility statement

- **Ethics:** The research adhered to strict ethical principles. Only public sources were used, and a rigorous anonymization protocol was applied to protect any sensitive data. The study explicitly seeks to mitigate the risks of AI misuse in the justice system, promoting responsible integration.

- **Reproducibility:** Committed to open science, our methodology is described in maximum detail. The **anonymized JURIDICO-FCR dataset** and the **evaluation scripts** will be released under a Creative Commons license to allow other interested researchers to validate and expand upon this work.

## 8  Conclusion

The emergence of large language models in the legal ecosystem presents a paradox: a technology with immense potential to improve efficiency and access to justice, yet hindered by an intrinsic propensity to fabricate information that directly undermines the pillar of truthfulness.

This work has addressed this paradox not as a technical failure, but as a **fundamental choice of architecture and paradigm**. Through a rigorous empirical study in the field of law in Spanish, it is quantitatively demonstrated that:

1. **Generative AI**, operating as a "creative oracle," is fundamentally unsafe for high-risk legal practice, with unacceptable error rates that make it an "informational liability."

2. **Consultative AI**, implemented through canonical RAG and operating as an "expert archivist," transforms the nature of risk, drastically reducing fabrications and becoming a true "efficiency asset," though not entirely free from subtle errors. A more advanced level of RAG almost completely mitigates such fabrications.

3. **Factual reliability** is not optional, but the indispensable prerequisite for AI to provide real practical value in the legal domain.

Ultimately, the path toward reliable legal AI does not lie in waiting for a perfect "oracle." It resides in the deliberate and rigorous construction of transparent "archivists," in the development of a professional culture of "informed skepticism," and in the recognition that critical human judgment is not an obstacle to be automated, but the most valuable resource that technology should aspire to amplify. I hope this work serves as a call to action for the legal and technological communities to collaborate in building a future where AI does not replace, but rather strengthens, the pursuit of justice.

## References

[1] Dahl, Matthew and Magesh, Varun and Suzgun, Mirac and Ho, Daniel E. Large Legal Fictions: Profiling Legal Hallucinations in Large Language Models. arXiv preprint arXiv:2401.01301, 2024.

[2] Dantart, Alex. Inteligencia Artificial jurídica y el desafío de la veracidad: análisis de alucinaciones, optimización de RAG y principios para una integración responsable. arXiv preprint arXiv:2509.09467, 2025.

[3] Lewis, Patrick and Perez, Ethan and Piktus, Aleksandra and Petroni, Fabio and Karpukhin, Vladimir and Goyal, Naman and Küttler, Heinrich and Lewis, Mike and Yih, Wen-tau and Rocktäschel, Tim and Riedel, Sebastian and Kiela, Douwe. Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. arXiv preprint arXiv:2005.11401, 2020.

[4] Karpukhin, Vladimir and Oğuz, Barlas and Min, Sewon and Lewis, Patrick and Wu, Ledell and Edunov, Sergey and Chen, Danqi and Yih, Wen-tau. Dense Passage Retrieval for Open-Domain Question Answering. arXiv preprint arXiv:2004.04906, 2020.

[5] Gao, Luyu and Ma, Xueguang and Lin, Jimmy and Callan, Jamie. Precise Zero-Shot Dense Retrieval without Relevance Labels. arXiv preprint arXiv:2212.10496, 2022.

[6] Nogueira, Rodrigo and Cho, Kyunghyun. Passage Re-ranking with BERT. arXiv preprint arXiv:1901.04085, 2019.

[7] Es, Shahul and James, Jithin and Espinosa-Anke, Luis and Schockaert, Steven. RAGAS: Automated Evaluation of Retrieval Augmented Generation. arXiv preprint arXiv:2309.15217, 2023.

[8] Min, Sewon and Krishna, Kalpesh and Lyu, Xinxi and Lewis, Mike and Yih, Wen-tau and Koh, Pang Wei and Iyyer, Mohit and Zettlemoyer, Luke and Hajishirzi, Hannaneh. FActScore: Fine-grained Atomic Evaluation of Factual Precision in Long Form Text Generation. arXiv preprint arXiv:2305.14251, 2023.

[9] Manakul, Potsawee and Liusie, Adian and Gales, Mark J. F. SelfCheckGPT: Zero-Resource Black-Box Hallucination Detection for Generative Large Language Models. arXiv preprint arXiv:2303.08896, 2023.

[10] Dhuliawala, Shehzaad and Komeili, Mojtaba and Xu, Jing and Raileanu, Roberta and Li, Xian and Celikyilmaz, Asli and Weston, Jason. Chain-of-Verification Reduces Hallucination in Large Language Models. arXiv preprint arXiv:2309.11495, 2023.