
INTELIGENCIA ARTIFICIAL JURÍDICA Y EL DESAFÍO DE LA VERACIDAD: ANÁLISIS DE ALUCINACIONES, OPTIMIZACIÓN DE RAG Y PRINCIPIOS PARA UNA INTEGRACIÓN RESPONSABLE

INFORME TÉCNICO

Alex Dantart

CIO LittleJohn

Paseo de la Castellana 194

28046, Madrid, España

arxiv@littlejohn.ai

ABSTRACT

Los grandes modelos de lenguaje (LLMs) están redefiniendo aceleradamente la práctica, la educación y la investigación jurídicas. Sin embargo, su vasto potencial se ve significativamente amenazado por la generación endémica de "alucinaciones" – resultados textuales que, aunque a menudo plausibles, son fácticamente incorrectos, engañosos o inconsistentes con las fuentes legales autorizadas. Este ensayo presenta una revisión exhaustiva y un análisis crítico multidimensional del fenómeno de las alucinaciones en LLMs aplicados al derecho. Se documentan las tendencias y manifestaciones de las alucinaciones a través de jurisdicciones, tipos de tribunales y clases de tareas legales, fundamentándonos en la creciente evidencia empírica de estudios recientes que evalúan tanto LLMs públicos como herramientas comerciales especializadas de Inteligencia Artificial (IA) legal. Se analizan en profundidad las causas subyacentes de estas alucinaciones, desde las deficiencias en los datos de entrenamiento y las limitaciones inherentes a la arquitectura probabilística de los modelos, hasta las complejidades del lenguaje jurídico y la tensión fundamental entre fluidez generativa y factualidad estricta.

Se examina con detalle la Generación Aumentada por Recuperación (RAG) como la principal estrategia de mitigación propuesta, evaluando críticamente su efectividad teórica, sus implementaciones prácticas y sus limitaciones persistentes en el singular contexto legal, incluyendo los puntos de fallo en sus fases de recuperación y generación. Más allá del RAG canónico, se discuten y proponen estrategias holísticas y avanzadas para la optimización y mitigación, abarcando desde la curación estratégica de datos y la ingeniería de prompts sofisticada, hasta la consideración de agentes de IA conscientes de la jerarquía normativa (como la pirámide de Kelsen), el fine-tuning enfocado en la fidelidad, y la implementación de robustos mecanismos de verificación post-hoc y calibración de confianza. Se ilustra la gravedad de estos fenómenos mediante el análisis de estudios de caso detallados de incidentes judiciales reales, extrayendo lecciones tangibles sobre las consecuencias de la confianza acrítica en la IA.

Con una mirada prospectiva, se explora el camino hacia una IA legal más fiable, delineando los desarrollos necesarios en modelos inherentemente más explicables (XAI, del inglés *Explainable Artificial Intelligence*), sistemas técnicamente auditables y la adopción de un paradigma de IA responsable por diseño. Finalmente, se exploran las profundas implicaciones éticas y regulatorias, con especial atención al marco normativo europeo y español, enfatizando el rol irreductible e insustituible de la supervisión humana y el juicio profesional del abogado en la era de la inteligencia artificial. Se concluye subrayando la imperiosa necesidad de una integración cautelosa, crítica y supervisada de los LLMs en la práctica legal, se proponen una tipología refinada de alucinaciones legales con el fin de guiar y estructurar la investigación futura en este campo crucial, y también se propone un nuevo marco de trabajo que distingue entre la IA Generativa de propósito general y la IA Consultiva especializada, ofreciendo una tipología refinada de alucinaciones legales que guiará la investigación futura hacia una integración verdaderamente responsable.

Sin embargo, es imperativo realizar también una distinción fundamental que a menudo se pasa por alto en el debate actual: la diferencia entre la Inteligencia Artificial de propósito general (como los LLMs públicos) y la Inteligencia Artificial especializada y consultiva diseñada específicamente para el dominio legal. Mientras que la primera, por su naturaleza generativa, es inherentemente propensa a las alucinaciones al "inventar" respuestas para mantener la fluidez conversacional, la segunda opera bajo un principio radicalmente distinto. Una IA consultiva no crea conocimiento, sino que lo recupera, estructura y presenta de forma fundamentada, actuando como un asistente experto que cita sus fuentes en lugar de un oráculo creativo. Este informe argumentará que la mitigación efectiva de las alucinaciones en el sector legal no reside en mejorar incrementalmente los modelos generativos, sino en adoptar un paradigma consultivo donde la veracidad y la trazabilidad son el núcleo del diseño, no una característica añadida. La tecnología, en este contexto, no es un sustituto del juicio humano, sino una herramienta para amplificarlo, cumpliendo la máxima de humanizar la tecnología en lugar de simplemente automatizar procesos.

Keywords Alucinaciones IA · Large Language Models (LLM) · Retrieval-Augmented Generation (RAG) · Derecho · ética legal · evaluación IA · mitigación de alucinaciones · Inteligencia Artificial Jurídica

1 Introduction

La inteligencia artificial (IA), y en particular los grandes modelos de lenguaje (LLMs), se encuentran en la cúspide de una transformación significativa en múltiples sectores, siendo el dominio legal uno de los más impactados y debatidos (Choi et al., 2022; Katz et al., 2023; Rodgers, Armour, and Sako 2023). Herramientas como ChatGPT de OpenAI, Gemini de Google, DeepSeek, y Llama de Meta, junto con plataformas especializadas de IA legal, prometen revolucionar tareas fundamentales como la investigación jurídica, la redacción de documentos, el análisis de contratos y la asistencia en litigios (Guha et al. 2023; Livermore, Herron, and Rockmore 2024). El potencial para aumentar la eficiencia, reducir costos y **democratizar el acceso a la justicia es considerable** (Perlman 2023; Tan, Westermann, and Benyekhlef 2023).

Sin embargo, este potencial transformador se ve obstaculizado por un desafío inherente y crítico: el fenómeno de las "alucinaciones" (Ji, Lee, et al. 2023). Las alucinaciones en LLMs se refieren a la **generación de información que, aunque a menudo plausible y lingüísticamente coherente, es fácticamente incorrecta, engañosa, inconsistente** con las fuentes proporcionadas o completamente fabricada (Dahl et al., 2024; Magesh et al., 2024).

En el contexto legal, donde la precisión, la fidelidad a las fuentes autorizadas (precedentes, estatutos) y la argumentación basada en hechos son primordiales, las alucinaciones no son meras inexactitudes técnicas, sino que representan un **riesgo sustancial** que puede llevar a errores estratégicos, consejos legales perjudiciales, sanciones profesionales e incluso la **erosión de la confianza pública en el sistema legal** (Roberts 2023; Weiser 2023). El infame caso *Mata v. Avianca, Inc.* (2023), donde abogados fueron sancionados por presentar un escrito judicial citando casos inexistentes generados por ChatGPT, sirve como un claro recordatorio de los peligros (Lantyer, 2024).

Es crucial matizar, sin embargo, que el desafío de la veracidad en el derecho trasciende la mera corrección factual. A diferencia de otros dominios, en el ámbito jurídico una afirmación no es solo "verdadera" o "falsa"; su validez a menudo reside en la solidez de su interpretación y argumentación, que es precisamente el terreno del juicio profesional experto. Por tanto, el peligro de la IA no es solo que genere falsedades verificables, sino también que construya argumentos legalmente inviables o interpretaciones superficiales que, sin el filtro crítico de un abogado, pueden conducir a estrategias erróneas. El análisis de la veracidad debe, por consiguiente, abarcar tanto la fidelidad a la fuente como la viabilidad interpretativa.

Este ensayo se embarca en una exploración exhaustiva de las alucinaciones de los Grandes Modelos de Lenguaje (LLMs) en el dominio legal, con el objetivo de ir más allá de los informes anecdóticos y proporcionar un análisis sistemático basado en la creciente evidencia empírica y la literatura académica. Para ello, primero definiremos y categorizaremos las alucinaciones específicas del contexto legal, explorando sus causas fundamentales y su impacto particular en la práctica jurídica (Sección 2). Seguidamente, examinaremos en detalle los métodos y desafíos inherentes a la evaluación de la prevalencia y naturaleza de estas alucinaciones, revisando críticamente los estudios recientes sobre LLMs generales y las herramientas comerciales de IA legal (Sección 3). A continuación, analizaremos en profundidad la Generación Aumentada por Recuperación (RAG) como la principal estrategia de mitigación propuesta, evaluando tanto sus promesas conceptuales como sus limitaciones inherentes y su efectividad empírica en el contexto legal (Sección 4). Posteriormente, discutiremos un abanico de estrategias complementarias y avanzadas para la optimización y mitigación de alucinaciones, abarcando desde la curación de datos y la ingeniería de prompts hasta la consideración de agentes de IA conscientes de la jerarquía normativa y los mecanismos de verificación post-hoc (Sección 5). Para ilustrar la gravedad y las consecuencias tangibles de estos fenómenos, presentaremos y analizaremos estudios de caso

detallados de incidentes reales donde las alucinaciones de la IA han impactado procedimientos judiciales (Sección 6). Con una mirada prospectiva, exploraremos el camino hacia una IA legal más fiable, discutiendo el desarrollo de modelos explicables, auditables y responsables por diseño (Sección 7). Finalmente, reflexionaremos sobre las cruciales consideraciones éticas y regulatorias que surgen, con especial atención al marco normativo europeo y español (Sección 8), para concluir sintetizando los hallazgos y enfatizando el camino a seguir hacia una integración responsable y efectiva de la IA en la práctica legal (Sección 9).

2 El fenómeno de las alucinaciones en LLMs legales: naturaleza, causas e impacto

La integración de los Grandes Modelos de Lenguaje (LLMs) en el ecosistema legal representa una de las transformaciones tecnológicas más profundas y potencialmente disruptivas de la era moderna. Estas arquitecturas de inteligencia artificial (IA), capaces de procesar y generar lenguaje natural con una fluidez sin precedentes, prometen optimizar radicalmente tareas intensivas en conocimiento como la investigación jurídica, la redacción de contratos, el análisis de pruebas (discovery) y la generación de escritos procesales (Choi et al., 2022; Livermore, Herron, and Rockmore 2024). Sin embargo, esta promesa se ve ensombrecida por un desafío inherente y omnipresente: el fenómeno de las "alucinaciones" (Ji, Lee, et al. 2023; Marcus & Davis, 2022). Lejos de ser una anomalía ocasional, las alucinaciones constituyen una característica intrínseca del funcionamiento actual de los LLMs, manifestándose como la generación de contenido que, aunque a menudo sintácticamente y semánticamente plausible, carece de fundamento fáctico, es lógicamente inconsistente o contradice directamente las fuentes de autoridad establecidas. En el dominio legal, donde la precisión factual, la fidelidad a la autoridad (leyes, precedentes, doctrina...) y la integridad argumentativa son pilares fundamentales, la propensión de los LLMs a alucinar no es un mero inconveniente técnico, sino un riesgo sistémico con profundas implicaciones éticas, profesionales y sociales (Roberts 2023). El fenómeno de las alucinaciones no es un mero inconveniente técnico, sino que ha sido identificado como uno de los desafíos críticos que definen la frontera actual de la investigación en IA legal. Revisiones exhaustivas del campo señalan que, a pesar de los avances transformadores de los LLMs, la "**alucinación en reclamaciones legales, manifestada como citas espurias o fabricaciones normativas**", junto con los déficits de explicabilidad y la adaptación jurisdiccional, constituyen las principales barreras para su adopción generalizada y fiable (Shao et al., 2025).

2.1 Un paradigma fundamental: IA generativa vs. IA consultiva

Antes de diseccionar el fenómeno de las alucinaciones, es imperativo establecer una distinción conceptual que el debate actual a menudo ignora, generando una peligrosa confusión: la diferencia fundamental entre la **Inteligencia Artificial generativa** y la **Inteligencia Artificial consultiva**. El término "IA Legal" se utiliza de forma monolítica, cuando en realidad describe dos arquitecturas con propósitos, mecanismos y perfiles de riesgo radicalmente distintos. Entender esta dicotomía no es un mero ejercicio académico; es la clave para una integración responsable y efectiva de la IA en la práctica jurídica.

2.1.1 Inteligencia Artificial generativa: el oráculo creativo

La IA Generativa, cuyo máximo exponente son los LLMs de propósito general como GPT, Gemini o Claude, opera como un "imitador avanzado" o un "sabelotodo creativo". Su objetivo principal no es la veracidad, sino la **fluidez conversacional y la coherencia probabilística**.

- **Definición y mecanismo:** Estos modelos funcionan prediciendo la siguiente palabra más probable en una secuencia, basándose en los patrones estadísticos aprendidos de un vasto y heterogéneo corpus de datos de internet. No "comprenden" el contenido ni "razonan" a partir de principios lógicos, sino que ensamblan texto que *suenan* plausible. Su conocimiento es paramétrico y está "congelado" en el momento de su entrenamiento. Investigaciones fundamentales sobre las causas de las alucinaciones explican que este comportamiento no es un fallo a corregir, sino una consecuencia directa de su diseño. Los modelos son optimizados para ser buenos "examinandos": en un sistema donde no se premia la incertidumbre, la estrategia más efectiva para obtener una "buena nota" (una respuesta plausible) es siempre arriesgar una respuesta en lugar de admitir desconocimiento. Por tanto, su tendencia a "inventar" es el resultado esperado de su entrenamiento (Kalai et al., 2025).
- **Ventajas:** Su fortaleza reside en tareas creativas: redacción de borradores, lluvia de ideas, resumen de textos no críticos y la generación de contenido donde la originalidad es más importante que la precisión factual.
- **Desventajas y riesgos inherentes:** Para el sector legal, su diseño es una receta para el desastre.
- **Alucinaciones "de diseño":** La propensión a alucinar no es un fallo, es una característica intrínseca de su arquitectura. Para evitar silencios y mantener la coherencia, el modelo "rellenará los huecos" o "inventará" hechos, sentencias o estatutos.

- **Opacidad total ("Caja Negra"):** Es imposible trazar el origen de una afirmación específica. La respuesta es un producto final opaco, sin referencias verificables.
- **Riesgo de "incesto de IAs":** Al ser entrenadas con el internet público, corren el riesgo de retroalimentarse con contenido de baja calidad generado por otras IAs, degradando su fiabilidad en un ciclo vicioso.

2.1.2 Inteligencia Artificial consultiva: el archivero experto

La IA Consultiva representa un cambio de paradigma. Su objetivo no es crear, sino **recuperar, estructurar y presentar conocimiento verificado**. Su arquitectura fundamental se basa en la Generación Aumentada por Recuperación (RAG), operando como un "archivero experto" o un "detective" que investiga antes de hablar.

- **Definición y mecanismo:** Este modelo no confía en su conocimiento paramétrico interno. Ante una consulta, su primer paso es buscar en un corpus de datos externo, curado y autorizado (ej. bases de datos de legislación, jurisprudencia, documentos internos de un despacho). Solo después de recuperar los fragmentos de información más relevantes, genera una respuesta que debe estar *estrictamente fundamentada* en dichos fragmentos.
- **Ventajas:** Diseñada para la fiabilidad en dominios críticos.
- **Mitigación de alucinaciones:** Reduce drásticamente la fabricación de hechos, ya que las respuestas están ancladas a fuentes explícitas.
- **Transparencia y trazabilidad:** La respuesta no es una "caja negra". Un sistema consultivo bien diseñado debe citar sus fuentes, permitiendo al profesional legal verificar la información y asumir la responsabilidad final con conocimiento de causa. Es la materialización del principio de **"no sustituir, sino amplificar"** el juicio humano.
- **Conocimiento actualizado:** Su fiabilidad depende de la actualidad de su base de datos, que es mucho más fácil y barata de actualizar que reentrenar un LLM masivo.
- **Desventajas y limitaciones:** No es una panacea. Su efectividad depende críticamente de la calidad de su corpus documental y de la sofisticación de su módulo de recuperación. Aún puede producir alucinaciones sutiles, como el *misgrounding* (tergiversar una fuente real), pero el riesgo de invención flagrante se minimiza.

2.1.3 Tabla comparativa de paradigmas

Característica	IA generativa (propósito general)	IA consultiva (especializada)
Objetivo principal	Fluidez y coherencia conversacional.	Precisión, fiabilidad y fundamentación.
Fuente de conocimiento	Paramétrico, interno, estático ("libro cerrado").	Externo, curado, dinámico ("libro abierto").
Riesgo de alucinación	Alto, especialmente fabricación de hechos ("de diseño").	Bajo en fabricación, riesgo de <i>misgrounding</i> .
Transparencia	Baja ("caja negra").	Alta (debe citar fuentes y razonamiento).
Caso de uso ideal	Brainstorming, borradores creativos, tareas no críticas.	Investigación jurídica, <i>due diligence</i> , respuestas factuales.
Analogía	Un "sabelotodo" elocuente pero a veces poco fiable.	Un "archivero" meticoloso que siempre muestra sus fichas.

La adopción de este marco dual es esencial para navegar la complejidad de la IA Legal. Confundir ambos paradigmas lleva a expectativas irreales y a una aplicación irresponsable de la tecnología. Las secciones subsiguientes de este informe analizarán en profundidad los desafíos inherentes al modelo generativo y cómo las arquitecturas consultivas, principalmente a través de RAG, intentan construir un camino hacia una IA legal verdaderamente fiable.

2.2 Definición y taxonomía de las alucinaciones legales

Definir la "alucinación" en el contexto de la IA legal requiere ir más allá de la simple dicotomía correcto/incorrecto. Una alucinación legal se materializa cuando un LLM genera una afirmación, cita, argumento o conclusión que se desvía de la realidad jurídica verificable o de la información contextual proporcionada, a menudo presentándola con una confianza injustificada (Khmaïess Al Jannadi, 2023). Es crucial entender que, si bien el término 'alucinación' es comúnmente usado, su aplicación en el derecho presenta desafíos únicos. A diferencia de dominios con verdades fácticas singulares, en el ámbito legal la 'corrección' de una afirmación interpretativa o un argumento puede ser objeto de debate entre

expertos. Por ello, más allá de la simple desviación factual, una alucinación legal también puede entenderse como la generación de una propuesta que, **aunque plausible, resulta legalmente inviable o indefendible bajo un escrutinio experto**, incluso si no contradice directamente una fuente explícita. La aparente coherencia de estas salidas puede enmascarar su falta de solidez jurídica, haciendo su detección particularmente compleja. Esta desviación puede adoptar múltiples formas, cada una con implicaciones distintas para la práctica legal.

Para enriquecer esta taxonomía, proponemos una dimensión adicional de clasificación basada en el origen arquitectónico de la IA. Las alucinaciones manifestadas por una IA de propósito general (ej. ChatGPT) suelen ser más graves (como la invención completa de jurisprudencia), ya que su objetivo es la coherencia conversacional a toda costa. Por el contrario, los errores en una IA consultiva especializada (basada en RAG) tienden a ser más sutiles, como el *misgrounding* o errores de síntesis, derivados de fallos en la recuperación o interpretación de un corpus documental controlado.

Esta distinción es crucial, pues mientras el primer tipo de alucinación representa un fallo sistémico de diseño para el uso legal, el segundo es un problema de implementación que puede ser mitigado con técnicas de optimización, como se discutirá más adelante. Ignorar esta diferencia es como confundir la opinión de un aficionado elocuente con el análisis documentado de un archivero experto.

Para definir y clasificar las alucinaciones legales con rigor, es útil adoptar un marco analítico que distinga las dos dimensiones clave del error. El estudio seminal de Magesh et al. (2025) sobre herramientas de IA legal propone una distinción fundamental entre:

- **Corrección (Correctness):** Si la afirmación es fácticamente verdadera en el mundo real.
- **Fundamentación (Groundedness):** Si la afirmación está correctamente respaldada por la fuente citada.

A partir de este marco, una "alucinación" se define como una respuesta que es incorrecta (contiene información falsa) o mal fundamentada (*misgrounded*, es decir, cita una fuente que no respalda la afirmación). Esta desviación puede adoptar múltiples formas, cada una con implicaciones distintas.

Las siguientes categorías detallan las manifestaciones específicas de estos fallos en la práctica jurídica:

- **Alucinaciones factuales/extrínsecas** (inconsistencia con los hechos del mundo legal): este es quizás el tipo más peligroso en la investigación y el asesoramiento legal directo. Se refiere a la generación de contenido que contradice el cuerpo establecido y verificable del derecho y los hechos relacionados.
 - *Misstatement* (declaración errónea) de la Ley o Precedente: El LLM describe incorrectamente el contenido o el holding de una ley o decisión judicial existente. Esto puede ir desde sutiles tergiversaciones hasta contradicciones directas con la autoridad citada o conocida.
 - Fabricación de autoridad: el modelo inventa por completo casos, estatutos, regulaciones o incluso jueces y académicos inexistentes. El caso *Mata v. Avianca, Inc.* (2023) es el ejemplo paradigmático, donde ChatGPT generó múltiples citas judiciales ficticias que fueron incorporadas a un escrito judicial.
 - Error de aplicación jurisdiccional o temporal: el LLM aplica incorrectamente principios legales de una jurisdicción a otra, o presenta como vigente una ley o precedente derogado u obsoleto, fallando en reconocer la dinámica temporal y espacial del derecho.
- **Alucinaciones basadas en fuentes** (errores de *groundedness* en Sistemas RAG): particularmente relevantes para los sistemas de Retrieval-Augmented Generation (RAG), que se discuten en la Sección 4. Estas ocurren cuando la respuesta generada es inconsistente con los documentos específicos recuperados por el sistema para fundamentar dicha respuesta.
 - *Misgrounding* (fundamentación errónea): el LLM cita correctamente una fuente existente (recuperada por el sistema RAG), pero hace una afirmación sobre su contenido que la fuente no respalda o incluso contradice (Magesh et al., 2024). Esto crea una falsa apariencia de soporte documental.
 - *Ungrounding* (falta de fundamentación): el LLM realiza afirmaciones factuales específicas que deberían estar respaldadas por el material recuperado, pero no proporciona citas o las fuentes recuperadas no contienen la información afirmada.
- **Alucinaciones de inferencia y razonamiento:** implican fallos en la estructura lógica del argumento legal o en la caracterización de las relaciones entre conceptos o autoridades.
 - Argumentación ilógica o inválida: el modelo construye una línea de razonamiento que viola principios lógicos básicos o que no se sostiene bajo el escrutinio legal, aunque pueda parecer superficialmente persuasiva.

- Miscaracterización de argumentos, partes o posturas procesales: el LLM confunde los argumentos de una parte con el holding del tribunal, o describe incorrectamente la postura procesal o las relaciones entre las partes en un litigio (Dahl et al., 2024).
- **Alucinaciones intrínsecas** (inconsistencia con el prompt o corpus de entrenamiento): aunque potencialmente menos frecuentes en respuestas directas a consultas legales factuales, pueden surgir en tareas de dominio cerrado como la summarización de textos legales extensos o la redacción de documentos basada en instrucciones detalladas, donde el resultado final se desvía sustancialmente o contradice el contenido o las directrices del input proporcionado.

Es crucial reconocer que estas categorías no son mutuamente excluyentes; una única respuesta alucinada puede exhibir múltiples tipos de errores simultáneamente. La característica unificadora es la **desconexión entre la salida generada y una base de verdad relevante** (sea esta los hechos del mundo legal, las fuentes recuperadas o el prompt inicial), a menudo enmascarada por la fluidez lingüística del modelo (Ji, Lee, et al. 2023).

Más allá de la fabricación de información, una forma más insidiosa de desviación se produce a través de la alteración del contenido existente, que puede inducir sesgos cognitivos en el profesional. La investigación ha cuantificado cómo los LLMs, en tareas de resumen, alteran el encuadre del texto original, por ejemplo, cambiando el sentimiento de neutro a positivo o negativo. En un estudio, se observó que esto ocurre en un 21.86% de los casos (Alessa et al., 2025). En el contexto legal, esto podría manifestarse como un resumen de una sentencia que enfatiza los argumentos de una parte sobre la otra, o que presenta un análisis doctrinal de manera más favorable o crítica de lo que realmente es, influyendo sutilmente en la evaluación inicial del abogado.

Adicionalmente, se ha identificado el sesgo de primacía, donde el resumen generado por el LLM se enfoca desproporcionadamente en la información presentada al inicio de un documento, ocurriendo en un 5.94% de las ocasiones (Alessa et al., 2025). Esto representa un riesgo significativo en la revisión de largos expedientes judiciales o contratos, donde los detalles críticos pueden encontrarse en secciones posteriores que el LLM podría minimizar u omitir.

Esta taxonomía general se complementa con esfuerzos de la comunidad para categorizar errores más granulares específicos de los sistemas RAG. Por ejemplo, el benchmark LibreEval (Arize AI) identifica tipos de fallo como 'Overclaim', donde el modelo excede lo soportado por las fuentes, 'Incompleteness', cuando la respuesta omite información crucial presente en el contexto, o 'Relational-error', que denota fallos al sintetizar correctamente la información de múltiples fragmentos recuperados. Estos errores específicos de RAG pueden considerarse manifestaciones detalladas de nuestras categorías más amplias, subrayando la complejidad de asegurar la fidelidad en estos sistemas.

2.3 Causas raíz de las alucinaciones en LLMs legales

Comprender por qué los Grandes Modelos de Lenguaje (LLMs) generan alucinaciones, especialmente cuando se aplican al riguroso dominio legal, es un paso indispensable para desarrollar estrategias efectivas de mitigación y evaluación. Las causas son multifactoriales, arraigadas tanto en las propiedades fundamentales de la tecnología actual de LLMs como en las complejidades específicas del conocimiento y el lenguaje jurídico. Estos factores interactúan de maneras complejas, dando lugar a las diversas manifestaciones de errores que hemos categorizado previamente.

Una causa fundamental reside en las **limitaciones inherentes a los datos de entrenamiento**. La vasta escala de los corpus utilizados para entrenar LLMs, a menudo extraídos de la web, implica una inevitable variabilidad en la calidad, veracidad y actualidad (Bender et al., 2021). En el ámbito legal, esto es particularmente problemático. Los textos legales disponibles públicamente pueden ser incompletos o representar solo una fracción del panorama jurídico total. Más críticamente, el derecho es un sistema dinámico; leyes y precedentes cambian constantemente, haciendo que cualquier LLM entrenado en un conjunto de datos estático contenga inevitablemente información obsoleta (Khmaïess Al Jannadi, 2023). A esto se suma la presencia de sesgos históricos –sociales, económicos, raciales– codificados en los textos legales y judiciales del pasado. Al aprender patrones estadísticos de estos datos, **los LLMs corren el riesgo no solo de reproducir, sino de amplificar estas desigualdades**, generando respuestas que pueden perpetuar injusticias sistémicas (Gebu et al., 2018; O'Neil 2016; Barocas and Selbst 2016). La escasez relativa de datos legales verificados, de alta calidad y representativos de todas las jurisdicciones y áreas del derecho sigue siendo un cuello de botella significativo.

Además, se puede argumentar que el problema se origina en la propia cultura de evaluación de la inteligencia artificial. Los modelos son entrenados y evaluados predominantemente con métricas que penalizan severamente las respuestas que expresan incertidumbre (como "no lo sé"). Como resultado, los LLMs aprenden que una conjetura plausible es preferible a una abstención honesta, perpetuando un comportamiento de "adivinar siempre" (Kalai et al., 2025). Este modo de operar, análogo al de un estudiante que nunca deja preguntas en blanco en un examen, es fundamentalmente incompatible con la prudencia que exige la práctica legal.

Table 1: Taxonomía de alucinaciones en LLMs legales

Categoría Principal	Subtipo / Descripción Breve	Ejemplo Ficticio Legal
Alucinaciones factuales / extrínsecas	Inconsistencia con los hechos del mundo legal.	
	<i>Misstatement de Ley o precedente:</i> Declaración errónea del contenido o holding de una autoridad.	"El LLM afirma que la Ley de Arrendamientos Urbanos de 2022 permite el desahucio inmediato sin notificación." (Cuando la ley exige 30 días)
	<i>Fabricación de autoridad:</i> Invención de casos, estatutos, o académicos inexistentes.	"Según el caso Martínez c. Constructora Sol (2025), la responsabilidad objetiva es inaplicable." (El caso no existe)
	<i>Error de aplicación jurisdiccional/temporal:</i> Aplicación incorrecta de normas a otra jurisdicción o presentación de normas derogadas como vigentes.	"El LLM cita un artículo del Código Civil de 1950 para resolver una disputa contractual actual, ignorando reformas posteriores."
Alucinaciones basadas en fuentes (errores de RAG)	Inconsistencia con los documentos recuperados por el sistema RAG.	
	<i>Misgrounding (Fundamentación errónea):</i> Cita una fuente real, pero afirma algo que la fuente no respalda o contradice.	"El documento X dice 'el contrato es válido', pero el LLM reporta: 'Según el documento X, el contrato es nulo'."
	<i>Ungrounding (Falta de fundamentación):</i> Realiza afirmaciones que deberían estar respaldadas por el material recuperado, pero no proporciona citas o las fuentes no lo contienen.	"El demandado actuó con negligencia. (Sin citar ninguna prueba o documento recuperado que lo sustente)."
Alucinaciones de inferencia y razonamiento	Fallos en la estructura lógica del argumento legal.	
	<i>Argumentación ilógica o inválida:</i> Construye una línea de razonamiento que viola principios lógicos.	"Si todos los contratos requieren oferta y aceptación, y este documento es un contrato, entonces el cielo es azul." (Conclusión no sigue)
	<i>Miscaracterización de argumentos/partes:</i> Confunde los argumentos de una parte con el holding, o describe incorrectamente posturas procesales.	"El LLM presenta la petición del demandante como si fuera la sentencia final del juez."
Alucinaciones intrínsecas	Inconsistencia con el prompt o corpus de entrenamiento (en tareas de dominio cerrado).	
	Desviación sustancial del contenido o directrices del input en tareas como sumariación o redacción basada en instrucciones.	"Prompt: 'Resume el siguiente contrato en 100 palabras enfocándote en las cláusulas de penalización.' Respuesta del LLM: Un resumen de 500 palabras sobre la historia de la empresa contratante."

La escasez de datos verificados y representativos de jurisdicciones específicas es una causa directa y demostrable de las alucinaciones. Un estudio empírico sobre el rendimiento de los LLMs en una jurisdicción no anglosajona reveló que, si bien modelos como GPT y Claude destacaban en tareas de redacción, todos los modelos fallaban sistemáticamente en la investigación jurídica, generando de forma frecuente citas a casos inexistentes. El autor concluye que esta deficiencia se debe a que los LLMs están entrenados predominantemente con datos de los sistemas jurídicos dominantes (como el estadounidense), careciendo de una base de conocimiento suficiente sobre la jurisprudencia de otras regiones, lo que les obliga a "alucinar" para completar la tarea (Hemrajani, 2025).

Íntimamente ligada a los datos está la **naturaleza probabilística y la arquitectura misma de los LLMs**. Estos modelos, a pesar de su impresionante capacidad para generar texto coherente, no operan a través de una comprensión semántica profunda o un razonamiento lógico análogo al humano (Searle 1980; Marcus & Davis, 2022). Son fundamentalmente motores predictivos que calculan la secuencia de palabras más probable basándose en las correlaciones estadísticas aprendidas de sus vastos datos de entrenamiento. Esta orientación hacia la predicción estadística, optimizada a menudo para la fluidez lingüística por encima de la factualidad estricta, los hace intrínsecamente propensos a generar afirmaciones que *suenan* correctas pero que carecen de base real (Ji et al., 2023; Bowman 2015).

Esta optimización intrínseca para la fluidez puede llevar a lo que a veces se describe como 'confabulación', un proceso mediante el cual el modelo, enfrentado a una falta de información fáctica directa o a una ambigüedad en la consulta, **'inventa' detalles o narrativas coherentes para mantener la continuidad del discurso**, aunque estos elementos fabricados carezcan de una base real. La confabulación, en este sentido, es una manifestación directa de la arquitectura predictiva del LLM priorizando la apariencia de comprensión sobre la factualidad estricta, llevando a la generación de alucinaciones que, aunque erróneas, pueden ser engañosamente persuasivas por su coherencia superficial.

La dificultad de los LLMs para conectar los principios legales abstractos con los hechos concretos de un caso es una causa fundamental de las alucinaciones. Un estudio que condicionó a los LLMs con diferentes niveles de conocimiento del sistema legal alemán para detectar el discurso de odio lo demostró de manera concluyente. Cuando los modelos eran "condicionados" únicamente con el conocimiento más abstracto (como el título de una norma constitucional o estatutaria), **mostraban una falta de comprensión profunda de la tarea, llegando a contradecirse y a alucinar respuestas** cuando se les presentaban normas ficticias o irrelevantes (Ludwig et al., 2025). Esto sugiere que la arquitectura probabilística de los LLMs, en ausencia de un anclaje en definiciones y ejemplos concretos, lucha por aplicar correctamente el razonamiento jurídico, recurriendo a la invención.

Fenómenos como el sobreajuste (*overfitting*), donde el modelo memoriza patrones específicos del entrenamiento en lugar de aprender principios generales, pueden exacerbar este problema, limitando su capacidad para generalizar correctamente a situaciones nuevas o ligeramente diferentes (Khmaïess Al Jannadi, 2023). Además, su capacidad inherente para la extrapolación (aunque esencial para la generalización) puede desviarse fácilmente hacia la invención o la conexión espuria de conceptos cuando se enfrenta a consultas que bordean los límites de su conocimiento o requieren inferencias complejas (Shaip, 2022; Huang et al. 2021; Domingos 2015).

El **dominio legal en sí mismo presenta una complejidad intrínseca** que amplifica estos desafíos. El lenguaje jurídico es notoriamente técnico, denso en significado, altamente dependiente del contexto y plagado de ambigüedades y términos polisémicos (Khmaïess Al Jannadi, 2023). Interpretar correctamente un estatuto, un contrato o una sentencia requiere no solo comprender el significado literal de las palabras, sino también el contexto legislativo, la intención de las partes, la historia procesal y la red de precedentes relevantes – tareas que exigen un nivel de comprensión contextual y razonamiento que desafía a los LLMs actuales. El razonamiento jurídico per se, con sus métodos analógicos, deductivos basados en reglas y principios, y su constante ponderación de factores, representa una forma de cognición de orden superior que los LLMs, basados en patrones estadísticos, luchan por emular fielmente (Ashley 2017; Choi and Schwarcz, 2024).

Incluso las **estrategias diseñadas para mitigar las alucinaciones, como RAG, introducen sus propios puntos de vulnerabilidad**. Como se discutirá en detalle en la Sección 4, la efectividad de RAG depende críticamente de la calidad de su módulo de recuperación de información. Si la información recuperada es irrelevante, incorrecta o incompleta, el LLM generador, incluso si intenta ser fiel al contexto proporcionado, producirá una respuesta defectuosa. Además, el propio LLM generador puede fallar en integrar correctamente la información recuperada, priorizando su conocimiento paramétrico erróneo o sintetizando incorrectamente las fuentes (Addleshaw Goddard, 2024).

En última instancia, el análisis de las causas raíz de las alucinaciones sería incompleto si se limitara al modelo. **La causa fundamental más peligrosa no reside en la máquina, sino en el humano que la utiliza sin criterio**. La IA tiene el potencial de amplificar las capacidades de los profesionales diligentes, mientras que puede inducir a errores a aquellos que la utilizan sin un criterio crítico y una supervisión adecuada. Un profesional con pensamiento crítico y experiencia utilizará el LLM como un catalizador para acelerar su investigación, validando cada resultado. Sin embargo, un usuario sin estas bases, seducido por la aparente facilidad, delegará su razonamiento y caerá en la trampa de la

complacencia. Por tanto, la mayor causa de riesgo no es la alucinación del modelo, sino la "**alucinación del usuario**": la creencia de que una herramienta puede sustituir la responsabilidad, el esfuerzo y el juicio profesional. Este fenómeno, alimentado por una cultura de la inmediatez y los atajos, es el verdadero desafío a mitigar en la integración de la IA en el sector legal.

2.4 Impacto específico y riesgos asociados en la práctica legal

La manifestación de estas causas en forma de alucinaciones tiene un impacto tangible y multifacético en el ecosistema legal:

1. **Socavamiento de la investigación y el análisis jurídico:** la base de cualquier trabajo legal riguroso es la investigación precisa. Las alucinaciones, al introducir información falsa o fabricada, contaminan este proceso fundamental, haciendo perder tiempo en la verificación, llevando a análisis erróneos y, en última instancia, a estrategias legales defectuosas.
2. **Riesgos profesionales y éticos:** para los abogados, confiar en información alucinada puede tener consecuencias devastadoras. Puede llevar a la presentación de escritos judiciales defectuosos (resultando en sanciones como la popular del caso *Mata v. Avianca*), al incumplimiento del deber de competencia y diligencia, a la violación del deber de franqueza ante el tribunal, y a posibles reclamaciones por negligencia profesional (Yamane, 2020; Schwarcz et al., 2024). La necesidad de verificar exhaustivamente cada resultado de la IA puede, irónicamente, anular los beneficios de eficiencia prometidos (Gottlieb 2024).
3. **Erosión de la confianza:** la prevalencia de alucinaciones, especialmente si no se aborda con transparencia, puede minar la confianza de los profesionales legales, los clientes y el público en general hacia las herramientas de IA y, por extensión, hacia quienes las utilizan (Khmaïess Al Jannadi, 2023). Esta erosión de la confianza puede obstaculizar la adopción de tecnologías potencialmente beneficiosas.
4. **Impacto en el acceso a la justicia:** existe una paradoja preocupante: los LLMs se promocionan como una herramienta para democratizar el acceso a la información legal para litigantes pro se o personas de bajos recursos. Sin embargo, estos mismos usuarios son los menos equipados para detectar y verificar alucinaciones sofisticadas, lo que los hace particularmente vulnerables a recibir información legal incorrecta y perjudicial (Draper and Gillibrand 2023; Dahl et al., 2024). En lugar de cerrar la brecha, la IA no fiable podría ampliarla.
5. **Integridad del sistema judicial:** a nivel sistémico, la introducción de información falsa o fabricada en los procedimientos judiciales, ya sea inadvertidamente por abogados o potencialmente de forma maliciosa, amenaza la integridad fundamental del proceso contradictorio y la búsqueda de la verdad.
6. **Riesgos Cognitivos y de Juicio Sutil:** Quizás el riesgo más subestimado no es que la IA proporcione información falsa, sino que presente información verídica de una manera que explote los sesgos cognitivos humanos. Los LLMs pueden actuar como "vectores de sesgo", induciendo sesgos de encuadre (framing bias) que alteran la percepción de un problema sin cambiar los hechos. Por ejemplo, al resumir los argumentos de la parte contraria, un LLM podría seleccionar un lenguaje que los haga parecer más débiles de lo que son. De igual forma, el sesgo de autoridad puede llevar a un abogado a aceptar una conclusión generada por la IA con menos escrutinio del que aplicaría a un colega humano, simplemente por la presentación fluida y aparentemente lógica del modelo (Alessa et al., 2025). Este efecto erosiona la objetividad del juicio profesional desde dentro, de una forma mucho más difícil de detectar que una simple cita falsa.

Abordar las causas raíz de las alucinaciones legales no es, por consiguiente, una mera optimización técnica, sino un imperativo ético y funcional para el futuro de la IA en el derecho.

La gravedad de este impacto no ha pasado desapercibida para los legisladores, y el riesgo inherente a la diseminación de información legal incorrecta o fabricada es una de las preocupaciones centrales que animan los esfuerzos regulatorios a nivel global. En este sentido, el Reglamento (UE) 2024/1689 del Parlamento Europeo y del Consejo, conocido como Ley de Inteligencia Artificial de la Unión Europea (en adelante, la Ley de IA de la UE, el Reglamento o EU-AIAct), un marco legislativo pionero y ambicioso, establece un precedente significativo. Al adoptar un enfoque basado en el riesgo, la Ley de IA de la UE busca imponer requisitos más estrictos a aquellos sistemas de IA cuyas fallas podrían tener consecuencias severas para los derechos fundamentales, la seguridad o el correcto funcionamiento de instituciones clave. Aunque la categorización específica de todas las herramientas de IA legal bajo este marco aún está por definirse en su aplicación práctica, es plausible anticipar que aquellos sistemas destinados a influir en la administración de justicia o a proporcionar asesoramiento en áreas críticas podrían ser objeto de un escrutinio regulatorio intensificado precisamente por el potencial disruptivo de fenómenos como las alucinaciones.

En conclusión, las alucinaciones no son un fallo técnico menor, sino una manifestación de las limitaciones fundamentales en la forma en que los LLMs actuales procesan la información y modelan el mundo, con ramificaciones particularmente

críticas en el sensible y normativo dominio del derecho. Abordar este desafío es una condición previa para cualquier integración responsable y beneficiosa de la IA en la profesión legal.

3 Evaluación de alucinaciones en aplicaciones legales de IA: metodologías, desafíos y estado actual

La mera existencia del fenómeno de las alucinaciones en LLMs aplicados al derecho, detallada en la sección anterior, impone una necesidad crítica e ineludible: el desarrollo y la aplicación de metodologías rigurosas para su evaluación, detección y cuantificación. Dada la naturaleza de alto riesgo del dominio legal, donde las decisiones basadas en información incorrecta pueden tener consecuencias jurídicas, financieras y sociales devastadoras, la simple confianza en las afirmaciones de los desarrolladores o en la aparente plausibilidad de las respuestas generadas es insostenible. La evaluación empírica sistemática se convierte, por tanto, no solo en un ejercicio académico deseable, sino en un prerrequisito fundamental para la integración responsable de estas tecnologías en la práctica profesional, la educación jurídica y el sistema de justicia en general. Sin embargo, como exploraremos en esta sección, la evaluación de las alucinaciones legales es una tarea intrínsecamente compleja, plagada de desafíos metodológicos y conceptuales únicos que requieren enfoques matizados y un escrutinio constante.

3.1 Desafíos fundamentales en la evaluación de la IA legal

Evaluar la factualidad y detectar alucinaciones en los LLMs cuando operan sobre conocimiento legal presenta un conjunto de desafíos particulares que van más allá de los encontrados en dominios más generales o con hechos más objetivos. Estos desafíos limitan la aplicabilidad directa de muchas métricas de evaluación estándar y exigen una consideración cuidadosa del contexto específico.

3.1.1 El problema del *Ground Truth* legal

A diferencia de preguntas con respuestas factuales únicas y objetivas (p. ej., "*¿Quién ganó el Mundial de 2022?*"), la "verdad" legal es a menudo más elusiva. El *ground truth* en derecho está intrínsecamente ligado a:

- *Interpretación*: las leyes y los precedentes requieren interpretación, y los expertos legales pueden discrepar razonablemente sobre el significado o la aplicación de una norma a un conjunto específico de hechos.
- *Variabilidad jurisdiccional y temporal*: el derecho aplicable varía enormemente entre jurisdicciones (locales, estatales, internacionales) y evoluciona constantemente con nuevas leyes y decisiones judiciales. Lo que es "correcto" en una jurisdicción o momento puede ser incorrecto en otro.
- *Ambigüedad lingüística*: como se mencionó, el lenguaje legal está repleto de términos técnicos, estándares vagos ("razonable", "debido proceso") y ambigüedades inherentes que desafían una verificación binaria simple. Esta complejidad inherente significa que, para muchas tareas legales que trascienden la mera recuperación de información (como el análisis de problemas jurídicos complejos o la formulación de estrategias), el concepto de un único *ground truth* contra el cual medir una respuesta de IA se vuelve inaplicable. En tales escenarios, la evaluación se desplaza de la 'corrección' binaria hacia la 'viabilidad legal': la capacidad de una respuesta para ser argumentativamente sostenible y coherente dentro del marco normativo y doctrinal, aun cuando puedan existir múltiples enfoques válidos. Establecer esta viabilidad requiere, por tanto, una profunda experiencia legal y, a menudo, implica juicios interpretativos que pueden ser objeto de debate. Establecer un *ground truth* fiable para evaluar las respuestas de un LLM requiere, por tanto, una profunda experiencia legal y, a menudo, implica juicios interpretativos que pueden ser objeto de debate.

3.1.2 Opacidad de los sistemas comerciales (el problema de la "Caja Negra")

Una barrera significativa para la evaluación independiente y rigurosa es la naturaleza propietaria y cerrada de muchas de las herramientas de IA legal más avanzadas disponibles comercialmente (Magesh et al., 2024). Los proveedores rara vez divulgan detalles cruciales sobre:

- *Datos de entrenamiento*: la composición exacta, fuentes, actualidad y posibles sesgos de los datos masivos utilizados para entrenar sus modelos base o especializados.
- *Arquitectura del modelo y algoritmos*: las especificidades de la arquitectura del LLM subyacente, los algoritmos de RAG empleados, o los métodos de fine-tuning aplicados.

- *Procesos internos*: los mecanismos específicos de recuperación de información, los prompts internos utilizados, o los filtros post-generación aplicados.

Esta opacidad impide a los investigadores y usuarios comprender plenamente por qué un sistema produce una respuesta particular (alucinada o no), aislar las fuentes de error, replicar los resultados de forma independiente o comparar de manera justa el rendimiento entre diferentes plataformas. La evaluación a menudo debe basarse únicamente en el análisis de la salida final, tratando al sistema como una "caja negra".

La opacidad inherente a los sistemas comerciales de IA legal trasciende la problemática puramente técnica para entrar en el ámbito de la estrategia de mercado y la gestión del riesgo. Se pueden identificar varias dinámicas clave:

- Señalización de mercado vs. transparencia técnica: El término "IA" funciona como una potente señal de mercado para atraer capital y clientes. Sin embargo, esta estrategia de marketing no siempre se corresponde con una divulgación transparente de las arquitecturas, los datos de entrenamiento o las tasas de error de los sistemas. Esto crea una asimetría de información que dificulta la evaluación objetiva por parte de los usuarios.
- Riesgo reputacional sistémico: La estrategia de "caja negra", si bien puede ofrecer ventajas comerciales a corto plazo, genera un riesgo sistémico. Un fallo notorio en un sistema opaco (p. ej., una alucinación con consecuencias judiciales) no solo daña la reputación del proveedor, sino que puede mermar la confianza en toda la categoría de productos de IA legal, ralentizando su adopción generalizada.
- El valor de la auditabilidad: En consecuencia, un factor diferenciador clave para la madurez del sector será la transición desde modelos que priorizan la percepción de la innovación hacia aquellos que demuestran su valor a través de la transparencia y la auditabilidad. Un sistema cuyo rendimiento puede ser verificado y comprendido por terceros ofrece una base más sólida para la confianza y la integración responsable en flujos de trabajo críticos.

3.1.3 Complejidad inherente de las tareas y habilidades legales

La práctica legal implica una gama diversa de tareas cognitivas que van mucho más allá de la simple recuperación de información o respuesta a preguntas factuales. Incluye el razonamiento analógico, la argumentación persuasiva, el juicio estratégico, la síntesis de información compleja, la redacción matizada y la comprensión contextual profunda. Evaluar el rendimiento de un LLM en estas tareas requiere métricas y metodologías que puedan capturar estas dimensiones cualitativas, lo cual es intrínsecamente más difícil que evaluar la corrección factual de una respuesta a una pregunta directa (Schwarz et al., 2024).

3.1.4 Ausencia de Benchmarks estandarizados y específicos

Si bien están surgiendo benchmarks en el área de IA y derecho, la comunidad académica ha respondido a esta necesidad con la creación de marcos de evaluación especializados y de dominio relevante. Plataformas como LexGLUE, un benchmark para la comprensión del lenguaje jurídico en inglés, y LawBench, que evalúa el conocimiento jurídico de los LLMs en el contexto chino, son ejemplos clave. Estos esfuerzos, catalogados en revisiones exhaustivas del campo (Shao et al., 2025), son fundamentales para establecer métricas estandarizadas que permitan cuantificar de manera rigurosa el progreso en tareas complejas como la predicción de sentencias y la recuperación de precedentes, moviendo el campo más allá de las evaluaciones genéricas de NLP.

No obstante, para llenar este vacío, la comunidad investigadora está desarrollando enfoques de evaluación innovadores que pueden agruparse en dos categorías principales.

Por un lado, surgen benchmarks técnicos diseñados específicamente para medir la fiabilidad de la arquitectura RAG. Iniciativas como LibreEval de Arize AI proporcionan conjuntos de datos para evaluar la propensión a la alucinación y la fidelidad al contexto (groundedness), mientras que herramientas como RAGTruth (Niu et al., 2024) persiguen objetivos similares. Estos esfuerzos son cruciales para cuantificar de manera rigurosa los fallos específicos de los sistemas RAG.

Por otro lado, una estrategia complementaria y creativa es el uso de exámenes estandarizados de acceso a la profesión como benchmarks de conocimiento y precisión factual. Un ejemplo notable es el uso del All India Bar Examination (AIBE) para validar el modelo "Legal Assist AI". Al alcanzar una puntuación del 60.08%, este enfoque proporcionó una métrica cuantificable y directamente comparable con el rendimiento humano (Gupta et al., 2025). La combinación de estas estrategias —tanto las técnicas como las basadas en la competencia profesional— es fundamental para construir un marco de evaluación verdaderamente robusto para la IA legal.

3.2 Métricas y metodologías para la detección y cuantificación

Navegar por los desafíos mencionados requiere el despliegue de un conjunto diverso de métricas y metodologías, cada una con sus fortalezas y debilidades inherentes:

- 1. Evaluación basada en referencias (usando oráculos de metadatos):** este enfoque, pionero en el estudio de Dahl et al. (2024), aprovecha la existencia de metadatos estructurados y verificables asociados con los documentos legales (p. ej., tribunal emisor, fecha de decisión, juez ponente, citas dentro del documento, estado de derogación). Se formulan consultas al LLM que tienen una respuesta objetiva y verificable en estos metadatos (p. ej., "¿Qué tribunal decidió el caso X?"). La respuesta del LLM se compara directamente con el *ground truth* del metadato.
 - **Fortalezas:** Proporciona una medida objetiva y cuantificable de la alucinación para un subconjunto de hechos legales verificables. Permite análisis a gran escala si se dispone de los metadatos adecuados.
 - **Debilidades:** Se limita a la información contenida en los metadatos disponibles, sin poder evaluar la corrección del razonamiento legal sustantivo o la interpretación. Si bien es valioso para identificar alucinaciones fácticas directas (ej. una cita incorrecta), este método no aborda la evaluación de respuestas a problemas legales complejos donde la 'corrección' depende de la interpretación y el razonamiento, y no de un simple dato verificable. En estos casos, la ausencia de un 'error factual' no garantiza la 'viabilidad legal' de la solución propuesta. Depende de la calidad y cobertura de las bases de datos de metadatos.
- 2. Evaluación libre de referencias (auto-consistencia / auto-contradicción):** esta familia de técnicas busca detectar alucinaciones sin necesidad de un *ground truth* externo, explotando la naturaleza estocástica de la generación de los LLMs (Manakul, Liusie, and Gales 2023; Mündler et al. 2023). Se generan múltiples respuestas para el mismo prompt (usando una temperatura > 0) y se analiza su consistencia.

Self-Contradiction como límite inferior: la detección de contradicciones lógicas directas entre múltiples respuestas generadas para el mismo input es una fuerte señal de alucinación, ya que respuestas fácticamente correctas deberían ser consistentes. Este método proporciona un límite inferior útil para la tasa de alucinación, sin asumir la calibración del modelo.

Self-Consistency como heurística: La consistencia entre múltiples respuestas puede usarse como una heurística para la confianza (respuestas más consistentes *podrían* ser más probables de ser correctas), pero esto asume un grado de calibración del modelo que puede no ser válido, especialmente en dominios complejos como el derecho.

 - **Fortalezas:** no requiere *ground truth* externo, potencialmente aplicable a una gama más amplia de preguntas, incluidas aquellas que involucran juicio o interpretación
 - **Debilidades:** la auto-contradicción solo proporciona un límite inferior (no detecta alucinaciones consistentes). La auto-consistencia como indicador de corrección es una heurística no garantizada. Requiere múltiples inferencias, aumentando el costo computacional.
- 3. Evaluación humana experta:** considerada el estándar de oro para evaluar tareas legales complejas y la calidad matizada de las respuestas generativas (Schwarcz et al., 2024). Involucra a expertos legales (abogados y académicos) que revisan y califican las salidas del LLM utilizando rúbricas predefinidas que evalúan dimensiones como la corrección factual, la solidez del razonamiento legal, la relevancia, la coherencia, la claridad y la identificación de riesgos.
 - **Fortalezas:** capaz de evaluar la calidad sustantiva, el razonamiento complejo y la relevancia contextual que las métricas automáticas a menudo pasan por alto. Es indispensable para validar nuevas tareas o métricas.
 - **Debilidades:** extremadamente costosa en tiempo y recursos, difícil de escalar, susceptible a la subjetividad y a la variabilidad entre evaluadores (requiere protocolos claros y medición de la fiabilidad inter-evaluador, como el Kappa de Cohen - Cohen 1960).
- 4. Métricas automatizadas:** incluyen métricas estándar de NLP como ROUGE o BLEU (más adecuadas para tareas como la sumarización o traducción) y métricas emergentes de factualidad que intentan verificar automáticamente las afirmaciones contra bases de conocimiento (p.ej., FActScore - Min et al. 2023) o usando otros LLMs como jueces (Zheng et al., 2023).
 - **Fortalezas:** escalables y computacionalmente eficientes una vez desarrolladas.
 - **Debilidades:** su correlación con el juicio humano sobre la calidad y factualidad *legal* es a menudo baja o no probada. Pueden ser fácilmente "engañadas" por respuestas fluidas pero incorrectas. Su desarrollo y validación para el dominio legal aún está en etapas tempranas.

En la práctica, un enfoque robusto para la evaluación probablemente requiera una **combinación** de estas metodologías: evaluación basada en referencias para hechos verificables, detección de auto-contradicción para obtener límites inferiores en tareas abiertas, métricas automáticas para análisis a gran escala (con validación cuidadosa), y evaluación humana experta como validación final y para tareas cualitativamente complejas.

3.3 Benchmarking de herramientas comerciales: estado actual y hallazgos empíricos clave

La necesidad crítica de evaluación ha impulsado el estudio empírico sistemático y pre-registrado de Magesh et al. (2025) sobre las plataformas comerciales líderes de IA legal. Sus hallazgos, basados en un conjunto diverso de más de 200 consultas legales del mundo real, son reveladores y establecen un punto de referencia crucial:

1. **Persistencia alarmante de alucinaciones:** Contrario a las audaces afirmaciones de marketing de ser "libres de alucinaciones", la mayoría de las herramientas evaluadas fallaron en una proporción significativa. Utilizando una definición rigurosa de alucinación (respuesta incorrecta o mal fundamentada), se encontró que Lexis+ AI alucinó en el 17% de los casos, y Westlaw AI-Assisted Research lo hizo en más de un tercio de las ocasiones (>33%).
 2. **Variabilidad extrema entre plataformas:** El rendimiento no es uniforme. Lexis+ AI tuvo una precisión general del 65% de respuestas correctas y fundamentadas, estableciéndose como la herramienta más fiable del grupo. En el otro extremo, Ask Practical Law AI, debido a una base de conocimiento más limitada, tuvo una tasa extremadamente alta de respuestas incompletas o rechazos (>60%), limitando severamente su utilidad práctica (Magesh et al., 2025).
 3. **Confirmación de que RAG es una Mitigación, no una Solución:** Los resultados confirman que la tecnología RAG sí reduce la tasa de alucinación en comparación con los LLMs de propósito general. Sin embargo, las tasas de error residuales demuestran que, en su implementación actual, el RAG no elimina totalmente las alucinaciones, siendo los fallos en la recuperación de información y en la adhesión del LLM a las fuentes los problemas persistentes.
1. **Persistencia alarmante de alucinaciones:** el hallazgo más contundente es que, contrariamente a las afirmaciones de marketing de "eliminación" o "ausencia" de alucinaciones, **la mayoría de las herramientas comerciales evaluadas alucinan en una proporción significativa**. Utilizando una definición rigurosa de alucinación (respuesta incorrecta o fundamentada erróneamente), se encontró que Lexis+ AI y Ask Practical Law AI alucinaban entre el 17% y el 33% de las veces, mientras que Westlaw AI-Assisted Research alucinaba más de un tercio del tiempo (>34%). Estas tasas, aunque inferiores a las de GPT-4 o 5 base en tareas legales (58-88%), siguen siendo inaceptablemente altas para la práctica profesional.

Table 2: Rendimiento comparativo de herramientas de IA Legal comerciales (adaptado de Magesh et al., 2024)

Herramienta de IA Legal	Tasa de Alucinación (%)	Resp. Incompletas (%)	Resp. Precisas (%)
Lexis+ AI	17	18	65
Westlaw AI-Assisted Research	>34	25	41
Ask Practical Law AI	17	>60	19
GPT-4/5 (base, como referencia)	~58-88	N/A*	N/A*

Nota: Las tasas de alucinación para herramientas comerciales se refieren a respuestas incorrectas o fundamentadas erróneamente. GPT-4/5 base se incluye como referencia general de LLMs sin RAG legal específico, sus tasas de alucinación en tareas legales pueden ser más altas y la estructura de "respuestas incompletas" o "precisas y fundamentadas" puede no ser directamente comparable sin el componente RAG. *N/A indica que la métrica no se reportó de la misma manera o no es directamente comparable.

2. **RAG es una (gran) mitigación, pero no "la" solución:** los resultados confirman que la tecnología RAG empleada por estas herramientas *sí* reduce la tasa de alucinación en comparación con el uso de LLMs de propósito general sin acceso a bases de datos legales externas. Sin embargo, RAG, tal como se implementa actualmente, **no elimina totalmente las alucinaciones**. Los fallos en la recuperación de información relevante y la incapacidad del LLM generador para adherirse fielmente a las fuentes recuperadas siguen siendo problemas sustanciales.
3. **Variabilidad entre plataformas:** el estudio revela diferencias notables en el rendimiento y el comportamiento entre las distintas herramientas. Lexis+ AI demostró la mayor precisión general (65% de respuestas correctas

y fundamentadas) y una tasa de alucinación más baja, pero aún significativa (17%). Westlaw AI-AR, aunque a menudo proporcionaba respuestas más largas y detalladas, exhibió la tasa de alucinación más alta (33%). Ask Practical Law AI, limitado a su base de conocimientos curada, tuvo una tasa de alucinación relativamente baja pero sufrió de una tasa extremadamente alta de respuestas incompletas o rechazos (>60%), limitando su utilidad práctica. Esta variabilidad subraya que la etiqueta "IA legal basada en RAG" engloba implementaciones muy diferentes con perfiles de riesgo distintos.

4. **Naturaleza insidiosa de los errores (más allá de la fabricación):** un hallazgo crucial es que las alucinaciones en estas herramientas RAG rara vez son fabricaciones completas de casos (aunque ocurren). Más comúnmente, adoptan formas más sutiles y potencialmente más peligrosas:

- *Misgrounding*: citar un caso o estatuto real pero tergiversar lo que dice o aplicarlo incorrectamente.
- *Errores de razonamiento*: fallos lógicos al sintetizar información de múltiples fuentes recuperadas.
- *Sycophancy/Sesgo Contrafáctico*: aceptar acríticamente premisas falsas en la consulta del usuario.
- *Supresión de citas problemáticas*: el estudio de Westlaw AI-AR observó instancias donde el sistema parecía generar una afirmación basada en un caso derogado, pero suprimía la cita directa, posiblemente debido a la integración con sistemas de verificación de citas como KeyCite, lo cual impide la verificación por parte del usuario.

De manera complementaria a estos hallazgos, un estudio empírico en un sistema jurídico con datos limitados evaluó el desempeño de varios LLMs (incluyendo GPT-4/5 y Claude 3) frente a un abogado junior en cinco tareas legales (identificación de problemas, redacción, asesoramiento, investigación y razonamiento). Los resultados corroboraron que, si bien los LLMs avanzados pueden igualar o incluso superar el rendimiento humano en tareas estructuradas como la redacción de escritos o la identificación de problemas, su fiabilidad colapsa en la investigación jurídica, donde la generación de casos falsos ("alucinaciones") fue un problema persistente en todos los modelos evaluados (Hemrajani, 2025). Este estudio refuerza la idea de que la efectividad de la IA legal es altamente dependiente de la tarea y de la calidad de los datos de entrenamiento para esa jurisdicción específica.

Estos errores "insidiosos" son particularmente preocupantes porque pueden crear una falsa sensación de fiabilidad y son más difíciles de detectar para un usuario que no realiza una verificación profunda de cada fuente citada.

Implicaciones de los hallazgos de evaluación:

Los resultados empíricos actuales, aunque limitados, tienen implicaciones significativas:

- **Escepticismo justificado:** demuestran que las afirmaciones audaces sobre la eliminación de alucinaciones por parte de los proveedores deben tomarse con extrema precaución.
- **Necesidad de transparencia:** subrayan la necesidad urgente de mayor transparencia por parte de los proveedores sobre cómo funcionan sus sistemas, qué datos utilizan y, crucialmente, sobre sus tasas de error y limitaciones conocidas, evaluadas mediante benchmarks independientes.
- **Imperativo de la diligencia profesional:** refuerzan la obligación ética y profesional ineludible de los abogados de verificar críticamente *cualquier* resultado generado por IA antes de incorporarlo a su trabajo o asesoramiento. La confianza ciega en estas herramientas es, en el estado actual de la tecnología, imprudente.
- **Guía para la investigación futura:** Identifican áreas clave para la mejora técnica (optimización de retrieval y generación en RAG legal) y para la investigación académica (desarrollo de mejores benchmarks, estudio del impacto en diferentes tipos de usuarios y tareas).

En conclusión, la evaluación rigurosa es la piedra angular para comprender y gestionar el riesgo de alucinaciones en la IA legal. Si bien las metodologías están evolucionando y enfrentan desafíos, la evidencia empírica inicial ya proporciona una advertencia clara: las alucinaciones son una realidad persistente incluso en las herramientas comerciales más avanzadas, lo que exige un enfoque cauteloso, crítico y centrado en el ser humano para la adopción de la IA en el derecho.

Para cuantificar el rendimiento de los LLMs en tareas jurídicas realistas, se llevó a cabo una evaluación manual experta cuyos resultados se resumen en la Figura 1. El proceso se diseñó para garantizar la objetividad y fiabilidad de las puntuaciones:

- **Anotadores y datos:** Un conjunto de 50 consultas complejas de investigación jurídica fue presentado a cada LLM y a un abogado experto. Las respuestas fueron evaluadas de forma independiente por dos juristas senior

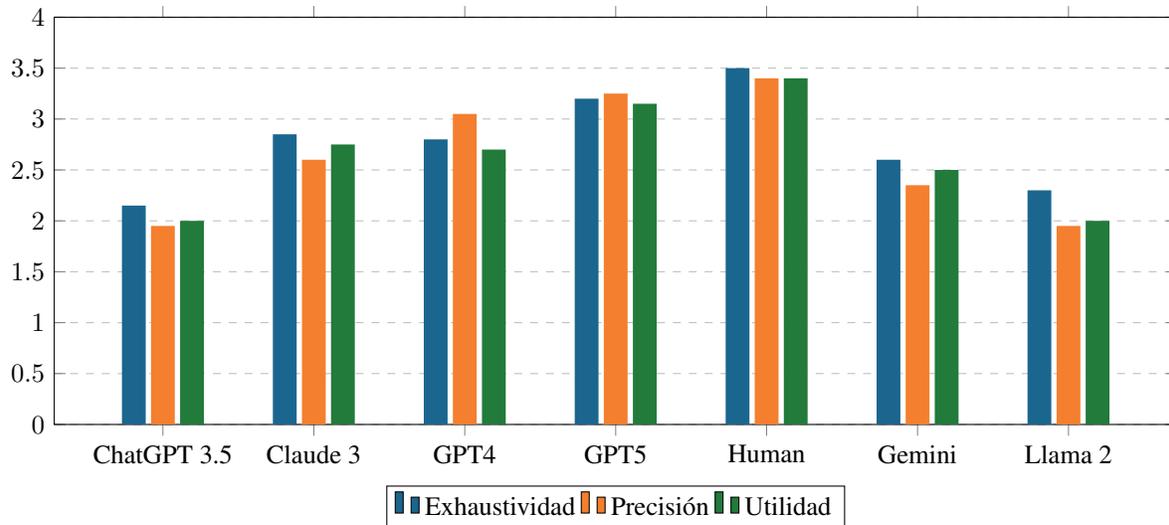


Figure 1: Evaluación comparativa del rendimiento de Modelos de Lenguaje (LLMs) y un experto humano en la tarea de investigación jurídica. Se observa que el rendimiento humano sigue siendo el punto de referencia en todas las métricas. GPT5 representa una mejora hipotética sobre GPT4, aunque sin alcanzar la fiabilidad humana. Las puntuaciones, en una escala de 1 a 4, son el resultado de una evaluación por pares de expertos juristas bajo criterios predefinidos, con un alto grado de acuerdo inter-anotador (Cohen's $\kappa = 0.85$), garantizando la objetividad de los resultados.

con experiencia en la materia, quienes no tuvieron conocimiento del origen de cada respuesta (evaluación a doble ciego).

- **Criterios de evaluación (rúbrica):** Los anotadores asignaron una puntuación de 1 (deficiente) a 4 (excelente) para cada una de las siguientes métricas, basándose en una guía de anotación predefinida:
 - **Exhaustividad:** ¿La respuesta identifica todos los puntos y matices legales relevantes? ¿Omite información crucial?
 - **Accurate (precisión):** ¿La información es fácticamente correcta y está libre de alucinaciones? ¿Las citas y la doctrina están correctamente representadas?
 - **Utilidad:** ¿La respuesta está bien estructurada, es fácil de entender y responde directamente a la consulta del usuario? ¿Acelera o dificulta el trabajo del profesional?
- **Fiabilidad Metodológica:** Para validar la consistencia de las evaluaciones, se calculó el acuerdo inter-anotador. Se obtuvo una puntuación **Kappa de Cohen de $\kappa = 0.85$** , lo que indica un grado de acuerdo "casi perfecto" entre los juristas y confirma la robustez de los datos presentados. Las puntuaciones mostradas en el gráfico representan el promedio de las calificaciones de ambos anotadores.

4 Retrieval-Augmented Generation (RAG) como paradigma dominante para la mitigación de alucinaciones legales

Frente a la inherente propensión de los Grandes Modelos de Lenguaje (LLMs) a generar alucinaciones, particularmente en un dominio tan sensible a la factualidad como el derecho, la comunidad de inteligencia artificial (IA) y los desarrolladores de tecnología legal han convergido predominantemente hacia un paradigma específico de mitigación: la Generación Aumentada por Recuperación, o Retrieval-Augmented Generation (RAG).

RAG representa un cambio fundamental respecto a la arquitectura estándar de los LLMs, que operan esencialmente en un modo de "libro cerrado", dependiendo exclusivamente del conocimiento internalizado (y potencialmente defectuoso o desactualizado) durante su entrenamiento masivo. En contraste, RAG busca dotar a los LLMs de un mecanismo de "libro abierto", permitiéndoles consultar activamente fuentes de información externas y relevantes *antes* de generar una respuesta. Esta sección se adentra en los fundamentos teóricos y mecánicos de RAG, evalúa críticamente sus ventajas teóricas específicas para el contexto legal, analiza sus limitaciones inherentes y puntos de fallo (que explican por qué, a pesar de su promesa, la mitigación de alucinaciones no es completa), revisa la evidencia empírica sobre su efectividad y discute las estrategias emergentes para su optimización en aplicaciones jurídicas.

La idoneidad de RAG para el dominio legal puede entenderse a través del **modelo de argumentación de Toulmin**, un marco fundamental en el razonamiento jurídico. Como señalan revisiones recientes, las tareas de los LLMs pueden mapearse directamente a los componentes de Toulmin (Shao et al., 2025). En esta analogía:

- La fase de Recuperación (Retrieval) del RAG se corresponde con la búsqueda de los "Datos" (hechos del caso) y el "Respaldo" (Backing) (estatutos y jurisprudencia aplicable).
- La fase de Generación (Generation) del LLM se corresponde con la construcción de la "Garantía" (Warrant) (el principio legal que conecta los hechos con la conclusión) para llegar a una "Reclamación" (Claim) (la conclusión legal).

Visto desde esta óptica, RAG no es simplemente un parche técnico contra las alucinaciones; es una arquitectura que computacionalmente imita la estructura fundamental de un argumento jurídico bien formado. Esto explica su predominio en herramientas avanzadas como ChatLaw, que integra RAG con bases de conocimiento estructuradas para fortalecer aún más el "Respaldo" de sus argumentos (Shao et al., 2025).

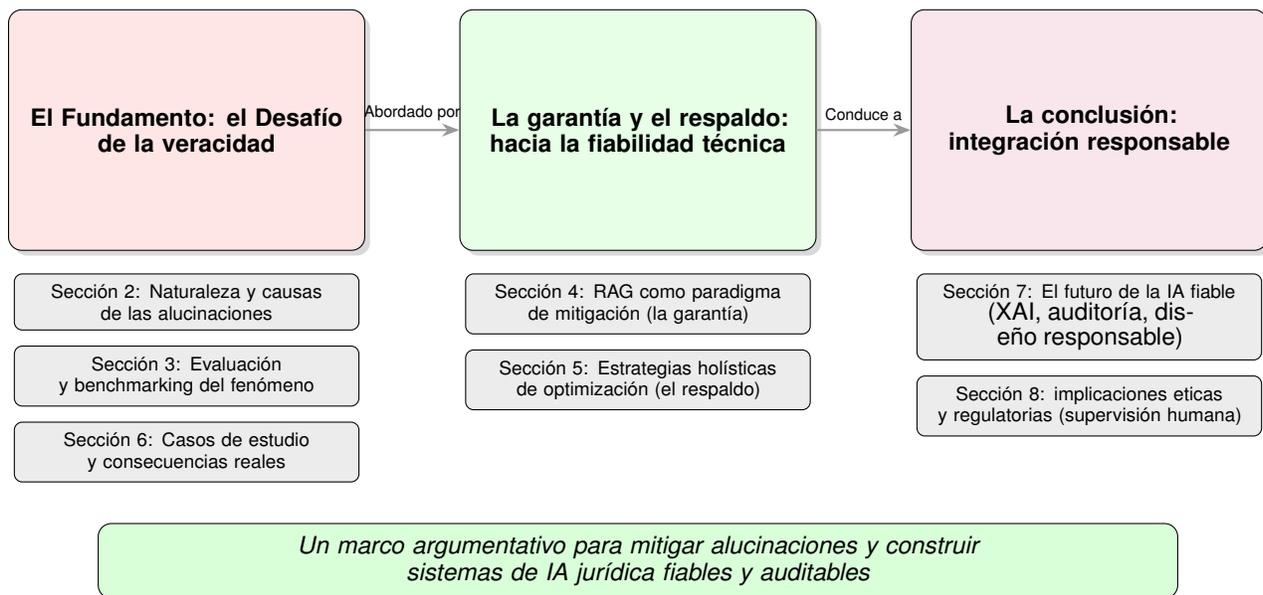


Figure 2: Descomposición de la estructura argumentativa del informe según el modelo de Toulmin. La figura ilustra el flujo lógico: partiendo del **Fundamento** (izquierda), que establece el problema de las alucinaciones (Secciones 2, 3 y 6); pasando a la **Garantía y Respaldo** (centro), que presenta la solución técnica con RAG y su optimización (Secciones 4 y 5); para llegar a la **Conclusión** (derecha), que define el marco para una integración responsable, incluyendo el futuro de la IA y sus implicaciones éticas y regulatorias (Secciones 7 y 8).

4.1 Fundamentos teóricos y mecanismo operativo de RAG

Los Grandes Modelos de Lenguaje (LLMs) base, a pesar de su asombrosa capacidad para generar texto fluido y coherente, operan fundamentalmente como 'cajas negras' con un conocimiento estático. Su proceso de toma de decisiones interno es en gran medida opaco, y el vasto corpus de información con el que fueron entrenados representa una instantánea del pasado, volviéndose progresivamente obsoleto a medida que el mundo –y especialmente el dinámico campo del derecho– evoluciona. Esta naturaleza intrínseca los hace inherentemente propensos a generar un espectro de errores factuales, englobados bajo el término "alucinación". Como se detalló en la **taxonomía de la Sección 2.2**, estos errores van desde la fabricación completa de autoridades hasta sutiles tergiversaciones de fuentes existentes (misgrounding), un desafío que la arquitectura RAG busca mitigar de raíz. Es precisamente para abordar estas limitaciones fundamentales –la opacidad, el conocimiento estático y la consiguiente falta de fundamentación verificable– que emerge la Generación Aumentada por Recuperación (RAG) como un cambio paradigmático. RAG no busca simplemente hacer que el LLM sea 'más inteligente' en abstracto, sino que lo transforma conceptualmente de un generador de lenguaje aislado a un sistema que interactúa dinámicamente con fuentes de conocimiento externas y explícitas, buscando anclar cada respuesta en evidencia recuperable y, por ende, potencialmente más fiable y actualizada.

A continuación, exploraremos los fundamentos teóricos y el mecanismo operativo de esta arquitectura crucial. **Lejos de ser un mero parche técnico, RAG se postula** como una reconfiguración conceptual de cómo los LLMs interactúan

con el conocimiento. Aborda directamente el problema de la 'caja negra' al externalizar la base de conocimiento a un corpus explícito y potencialmente verificable, y combate el problema del conocimiento estático al permitir que este corpus externo sea actualizado dinámicamente, independientemente de los costosos ciclos de reentrenamiento del modelo de lenguaje subyacente. Es esta doble promesa de fundamentación y actualidad la que ha posicionado a RAG como la principal esperanza para una IA legal más fiable.

El concepto central de RAG es simple pero poderoso: desacoplar el proceso de generación de lenguaje del almacenamiento de conocimiento fáctico masivo. En lugar de exigir que el LLM memorice y razone sobre la totalidad del corpus legal dentro de sus parámetros (una tarea propensa a la compresión con pérdidas y a la alucinación), RAG externaliza la base de conocimiento a un corpus documental explícito y recuperable (p. ej., bases de datos de jurisprudencia, estatutos, regulaciones, tratados legales, o incluso documentos internos de un bufete).

Este estudio no está dedicado a detallar el funcionamiento de RAG ya que existen innumerables ensayos por Internet sobre el tema, pero repasaremos de manera introductoria el proceso RAG canónico, que implica dos fases principales:

1. **Fase de recuperación (Retrieval):** dada una consulta del usuario (el *prompt*), esta fase tiene como objetivo identificar y extraer los fragmentos de información más relevantes del corpus documental externo. Este proceso típicamente involucra:
 - *Indexación:* pre-procesamiento del corpus documental, dividiéndolo en unidades manejables (chunks) y generando representaciones vectoriales (embeddings) para cada chunk mediante un modelo de embedding genéricos (p. ej., *text-embedding-ada-002* de OpenAI o modelos específicos de dominio como *bge-m3-spa-law-qa-large* de LittleJohn). Estos embeddings capturan el significado semántico de los chunks.
 - *Almacenamiento vectorial:* guardar los embeddings en una base de datos vectorial optimizada para búsquedas de similitud (p. ej., FAISS, Qdrant, Chroma...).
 - *Procesamiento de la consulta:* la consulta del usuario también se convierte en un embedding vectorial usando el mismo modelo.
 - *Búsqueda de similitud:* se realiza una búsqueda (típicamente por similitud coseno o distancia euclidiana) en la base de datos vectorial para encontrar los k chunks cuyos embeddings son más cercanos al embedding de la consulta.
 - *Recuperación híbrida (opcional pero común):* a menudo, la búsqueda semántica se combina con métodos tradicionales de recuperación de información basados en palabras clave (p. ej., BM25) para mejorar la precisión, especialmente para términos específicos o nombres propios.
 - *Re-ranking (opcional):* los chunks recuperados pueden ser reordenados usando modelos más sofisticados (cross-encoders) que evalúan la relevancia de cada chunk en relación con la consulta completa, aunque esto añade latencia.
2. **Fase de generación (Generation):** Los k chunks de texto recuperados, considerados los más relevantes para la consulta, se utilizan para "aumentar" el prompt original del usuario. Este prompt aumentado (que ahora contiene la consulta y el contexto recuperado) se introduce en el LLM generador (p. ej., GPT, Claude, Llama, Gemini...). El LLM tiene la instrucción de basar su respuesta principalmente en la información contextual proporcionada, sintetizándola y presentándola de manera coherente y relevante a la pregunta original. Idealmente, el LLM también debería ser capaz de citar las fuentes específicas de los chunks recuperados de donde extrajo la información.

El diseño de este proceso de dos etapas persigue que el LLM genere respuestas más precisas, actualizadas y fundamentadas, mitigando la necesidad de "inventar" información cuando su conocimiento paramétrico es insuficiente o incorrecto.

Un ejemplo práctico de esta arquitectura es el modelo "Legal Assist AI", diseñado para un sistema judicial específico con un corpus de datos curado. En su implementación, los documentos legales se cargan y se dividen en fragmentos manejables (chunks) de 1000 caracteres. A continuación, se generan representaciones vectoriales (embeddings) para cada chunk utilizando el modelo sentence-transformers/all-MiniLM-L6-v2 a través de HuggingFace. Finalmente, estos embeddings se indexan en una base de datos vectorial FAISS (Facebook AI Similarity Search), que permite una recuperación ultrarrápida de los fragmentos de texto semánticamente más relevantes para la consulta del usuario, los cuales son inyectados en el prompt del LLM generador (Gupta et al., 2025). Este flujo de trabajo ilustra el mecanismo RAG canónico en una aplicación legal del mundo real.

4.2 Ventajas teóricas de RAG: fundamentación, actualidad y transparencia

La adopción generalizada de la Generación Aumentada por Recuperación (RAG) como arquitectura preferente para los Grandes Modelos de Lenguaje (LLMs) en aplicaciones sensibles a la factualidad, y muy particularmente en el

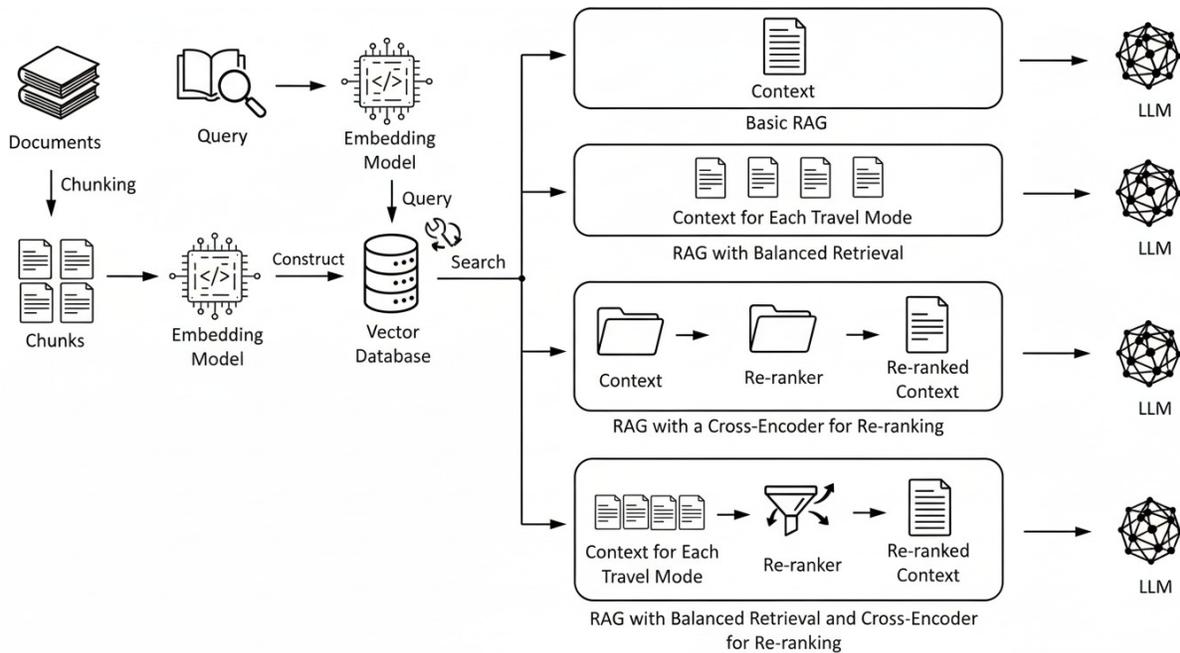


Figure 3: Diagrama esquemático de un sistema de generación aumentada por recuperación (RAG). El flujo ilustra cómo una consulta de usuario (Input Query) se enriquece con información contextual recuperada de una base de conocimiento externa (Vectorstore) antes de ser enviada al LLM. Este proceso ancla la respuesta en datos verificables, mitigando la generación de alucinaciones. Adaptado de Yiming Xu et al. (2025).

dominio legal, no es casual. Desde su diseño fundamental, RAG ofrece una serie de ventajas para abordar algunas de las limitaciones más críticas de los LLMs base cuando operan de forma aislada. Estas ventajas, si se materializan plenamente, tienen el potencial de transformar la IA de una herramienta generativa de lenguaje, a menudo desconectada de la realidad fáctica, en un asistente cognitivo genuinamente útil y más fiable para el profesional del derecho. La promesa teórica de RAG se asienta sobre tres pilares principales: la capacidad de fundamentar las respuestas en autoridad externa, la habilidad para operar con información dinámica y actualizada, y el potencial para una mayor transparencia y verificabilidad del proceso generativo.

El **principio de fundamentación (grounding) en autoridad externa** es, quizás, la ventaja más publicitada y esencial de RAG en el contexto jurídico. El derecho, por su propia naturaleza, es un sistema normativo y argumentativo que se construye sobre un vasto y jerarquizado corpus de fuentes autorizadas: constituciones, estatutos, regulaciones administrativas, jurisprudencia vinculante y persuasiva, y tratados doctrinales. Un LLM base, que depende únicamente de los patrones estadísticos internalizados durante su entrenamiento a partir de un corpus general (que puede o no incluir una representación adecuada de estas fuentes), opera esencialmente en un vacío autoritativo. Puede generar texto que *imita* el estilo del lenguaje legal, pero carece de un anclaje directo y verificable en las fuentes que definen el derecho aplicable. El paradigma de RAG aborda este problema obligando al LLM a interactuar con un corpus documental explícito de estas fuentes legales. Antes de generar una respuesta a una consulta jurídica, el sistema RAG primero recupera los fragmentos de texto más relevantes de este corpus. Esta información recuperada se convierte en el "fundamento" sobre el cual el LLM debe construir su respuesta. En un escenario ideal, esto significa que las afirmaciones legales, las interpretaciones y las conclusiones generadas por el sistema no son meras invenciones probabilísticas, sino que están directamente derivadas y soportadas por el texto de la ley, el precedente o el contrato pertinente. Para un abogado, esto es crucial, ya que cualquier argumento o consejo debe, en última instancia, ser rastreable hasta una fuente de autoridad válida.

La segunda ventaja fundamental de RAG es su capacidad inherente para **manejar información legal dinámica y asegurar la actualidad del conocimiento**. El derecho es un organismo vivo; las leyes se enmiendan, se derogan y se promulgan nuevas. Los tribunales emiten nuevas sentencias que reinterpretan, modifican o incluso revocan precedentes establecidos. Un LLM base, entrenado en un "snapshot" del pasado, se vuelve inevitablemente obsoleto a medida que el derecho evoluciona. Reentrenar estos modelos masivos desde cero o incluso actualizarlos de manera significativa es un proceso costoso, complejo y que consume mucho tiempo, lo que hace inviable mantenerlos perpetuamente al día con los cambios legislativos y jurisprudenciales. RAG ofrece una solución elegantemente simple a este problema de

obsolescencia: dado que el conocimiento fáctico primario reside en la base de datos documental externa y no en los parámetros del LLM, **la actualidad del sistema RAG depende principalmente de la actualidad de dicha base de datos**. Mantener actualizado un corpus documental específico (p. ej., añadiendo nuevas leyes, sentencias recientes, o actualizando el estado de derogación de los precedentes) es una tarea considerablemente más manejable y menos costosa que reentrenar un LLM de miles de millones de parámetros. En teoría, un sistema RAG bien mantenido podría proporcionar acceso a la información legal más reciente, permitiendo a los profesionales confiar en que las respuestas de la IA reflejan el estado actual del derecho, un requisito indispensable para la práctica competente.

Finalmente, RAG ofrece un **potencial significativo para una mayor transparencia y verificabilidad** en comparación con la naturaleza opaca de los LLMs base. Una de las críticas más persistentes a los LLMs es su funcionamiento como "cajas negras": generan respuestas, a menudo convincentes, pero sin ofrecer una explicación clara de *cómo* llegaron a esa conclusión o en qué información específica se basaron. En el ámbito legal, donde la capacidad de justificar un argumento y citar las fuentes es fundamental, esta opacidad es inaceptable. RAG, al basar explícitamente la generación en documentos recuperados, abre la puerta a una mayor transparencia. Idealmente, un sistema RAG no solo debería proporcionar una respuesta, sino también **citar las fuentes específicas de su corpus externo que fueron utilizadas para construir cada parte de esa respuesta**. Esto permitiría al profesional legal no solo recibir una conclusión, sino también revisar y evaluar críticamente la evidencia documental subyacente, juzgar la relevancia y la interpretación de las fuentes por sí mismo, y, en última instancia, asumir la responsabilidad informada por el uso de la salida de la IA. Esta capacidad de "mostrar el trabajo" es crucial para integrar la IA de manera responsable en los flujos de trabajo legales, donde la verificación humana sigue siendo un componente irreductible de la diligencia profesional.

Esta capacidad de 'mostrar el trabajo' es fundamental en un dominio como el jurídico, donde la 'respuesta correcta' raramente es un dato aislado o binario (verdadero/falso). Por el contrario, la validez de una conclusión legal reside en la solidez de su fundamentación y en la coherencia de su razonamiento. Al proporcionar acceso directo y trazable a las fuentes, RAG permite al profesional no solo validar la información, sino, lo que es más importante, analizar la interpretación propuesta por el modelo, evaluar la lógica de su argumentación y, en última instancia, construir su propio criterio experto. La verdadera asistencia de la IA no reside en ofrecer una conclusión, sino en articular de forma transparente los fundamentos que la sustentan, convirtiéndose en una herramienta para amplificar el juicio humano, no para sustituirlo.

En resumen, la arquitectura RAG, desde una perspectiva teórica, está diseñada para abordar de frente algunas de las deficiencias más críticas de los LLMs cuando se enfrentan a tareas legales sensibles a la factualidad. Promete respuestas más fundamentadas, actualizadas y verificables, moviendo a la IA legal un paso más cerca de ser un asistente cognitivo verdaderamente útil y fiable. Sin embargo, como se explorará en la siguiente sección, la transición de esta promesa teórica a una implementación práctica robusta y consistentemente fiable en el complejo y adversario dominio del derecho está plagada de desafíos significativos y puntos de fallo inherentes que explican por qué las alucinaciones, aunque mitigadas, persisten.

4.3 El talón de Aquiles de RAG: Análisis de limitaciones y evidencia empírica

A pesar de las considerables ventajas conceptuales que la Generación Aumentada por Recuperación (RAG) aporta a la tarea de fundamentar los Grandes Modelos de Lenguaje (LLMs) en conocimiento externo, tanto la evidencia empírica emergente como un análisis profundo de su mecanismo operativo revelan que esta arquitectura, aunque un avance significativo, no constituye una solución infalible. Lejos de ser una panacea, la promesa de respuestas consistentemente precisas, actuales y verificables se ve atenuada por una serie de limitaciones persistentes y puntos de fallo inherentes a sus dos fases operativas clave: la recuperación de información y la generación de lenguaje.

Desde una perspectiva de desarrollo de producto, un sistema RAG canónico debe ser tratado como lo que realmente es: un prototipo, y no una solución de producción robusta. La facilidad con la que herramientas como LangChain permiten ensamblar un prototipo RAG ha creado una falsa sensación de simplicidad. Confundir un prototipo funcional con un sistema fiable es el camino más rápido al desastre técnico y a la pérdida de confianza del cliente. La construcción de un sistema RAG serio no es un sprint de fin de semana; es un maratón de ingeniería de datos y refinamiento continuo.

Estos desafíos explican por qué incluso las herramientas RAG más sofisticadas siguen produciendo errores, desde inexactitudes sutiles hasta alucinaciones manifiestas que comprometen su fiabilidad. Esta sección analiza en profundidad estos puntos de fallo, contrastando las debilidades teóricas con los hallazgos empíricos más recientes.

4.3.1 Puntos de fallo en la fase de recuperación: el desafío del fundamento relevante

El primer y quizás más fundamental conjunto de vulnerabilidades reside en la fase de recuperación de información. La máxima de "*garbage in, garbage out*" aplica con toda su fuerza: si el sistema RAG no logra identificar y extraer los fragmentos de texto (*chunks*) verdaderamente relevantes, precisos y autorizados, el LLM generador, por muy avanzado

que sea, operará sobre una base informativa defectuosa, incrementando drásticamente la probabilidad de generar una respuesta errónea.

- **Ambigüedad Inherente a la "Relevancia Legal":** A diferencia de la recuperación de hechos discretos, determinar qué pasaje es *legalmente relevante* exige un sofisticado razonamiento jurídico. La simple similitud semántica superficial, en la que se basan muchos sistemas de búsqueda vectorial, puede ser engañosa. Un fragmento puede ser temáticamente similar pero provenir de una jurisdicción inaplicable o referirse a un estatuto derogado. Esta debilidad se ve confirmada en contextos jurisdiccionales con datos limitados. Un estudio sobre la práctica legal en una jurisdicción subrepresentada en los corpus de entrenamiento demostró que la incapacidad de los LLMs para realizar investigación jurídica fiable se debía a la escasez de jurisprudencia india en sus datos de entrenamiento. Esto evidencia un punto de fallo fundamental para RAG: aunque la arquitectura esté diseñada para recuperar información, si el corpus de recuperación carece de la información relevante, el LLM generador se ve forzado a operar sobre una base incompleta, lo que conduce directamente a la alucinación.
- **Deficiencias en las Estrategias de *Chunking*:** La forma en que los documentos legales extensos se dividen en fragmentos manejables es crítica. Un *chunking* deficiente puede llevar a la pérdida de contexto esencial, la introducción de ruido irrelevante o la fragmentación de unidades lógicas [Pinecone, 2024]. La importancia de una estrategia sofisticada se ilustra en el informe interno de Addleshaw Goddard (2024), que para optimizar una tarea de *due diligence*, tuvo que experimentar meticulosamente hasta concluir que fragmentos de 3,500 caracteres con un solapamiento de 700 eran óptimos para su corpus. Esto sugiere que las implementaciones genéricas de RAG probablemente exhiban tasas de error considerablemente más altas.
- **Recuperación Incompleta y Calidad de la Base de Conocimiento:** Incluso con estrategias de *chunking* mejoradas, el módulo de recuperación puede no identificar *todos* los fragmentos necesarios o puede priorizar incorrectamente los menos relevantes. Además, la base de datos documental debe ser exhaustiva, precisa y estar meticulosamente actualizada. Cualquier error, omisión o desactualización en el corpus subyacente se propagará inevitablemente a las respuestas generadas.

4.3.2 Puntos de fallo en la fase de generación: la tensión entre fidelidad y fluidez

Superados los desafíos de la recuperación, la fase de generación de lenguaje presenta su propio conjunto de puntos de fallo, incluso cuando se proporciona al LLM un contexto aparentemente correcto. La evidencia empírica es crucial aquí, pues revela que los errores más comunes y peligrosos no son las fabricaciones completas, sino formas más sutiles.

- **Falta de Fidelidad al Contexto Recuperado y Errores "Insidiosos":** A pesar de las instrucciones, el LLM generador puede ignorar o contradecir el contexto recuperado, recurriendo a su conocimiento paramétrico [Chen et al., 2024], o intentar "rellenar los huecos" inventando detalles no soportados explícitamente. Un hallazgo crucial del estudio de Magesh et al. (2024) es que las alucinaciones en herramientas RAG rara vez son fabricaciones completas de casos. Más comúnmente, adoptan formas más sutiles y peligrosas como el *misgrounding*: citar un caso o estatuto real pero tergiversar su contenido o aplicarlo incorrectamente. Este tipo de error es particularmente "insidioso" porque crea una falsa sensación de fiabilidad, dificultando su detección por parte de un profesional que no realice una verificación profunda de cada fuente.
- **Errores de Síntesis e Inferencia:** Cuando la respuesta requiere la integración de información de múltiples *chunks*, el LLM puede cometer errores lógicos o realizar inferencias inválidas. Benchmarks específicos para sistemas RAG, como **LibreEval de Arize AI (un conjunto de datos diseñado para evaluar la fidelidad de las respuestas al contexto proporcionado)**, han mostrado que los 'Relational-errors', que surgen de una síntesis defectuosa, son una forma común de alucinación en sistemas RAG.
- **Dependencia de la Naturaleza de la Tarea Legal:** La efectividad de RAG varía significativamente según la tarea. La extracción de cláusulas estandarizadas como "Governing Law" puede alcanzar altos niveles de precisión con RAG optimizado. Sin embargo, cláusulas más variables y contextualmente dependientes como "Exclusivity" o "Cap on Liability" presentan un mayor desafío y requieren una optimización más intensiva para alcanzar niveles similares de precisión [Addleshaw Goddard, 2024].
- **Dificultades en la Atribución y Citación Precisa:** Una manifestación común de la imperfección de RAG es la incapacidad del LLM para generar citas precisas que vinculen inequívocamente sus afirmaciones a los pasajes específicos de los documentos recuperados. Esta falta de atribución fiable socava uno de sus principales beneficios teóricos: la verificabilidad.

4.3.3 Síntesis de la evidencia: un mitigador imperfecto

En conclusión, la evidencia empírica actual converge en un punto claro: RAG es una herramienta valiosa que indudablemente mitiga la propensión de los LLMs a las alucinaciones factuales extrínsecas. Sin embargo, está lejos de ser una solución mágica que elimina el riesgo por completo.

Los estudios pioneros sobre herramientas comerciales, como el de Magesh et al. (2024), proporcionan datos cruciales. Sus hallazgos revelan una persistencia preocupante de errores: documentaron que las principales plataformas generaban respuestas incorrectas o con fundamentación errónea (*misgrounded*) en un rango de entre el **17% y más del 33%** de las consultas. Aunque esto representa una mejora sustancial sobre las tasas de alucinación de los LLMs base en contextos legales (que pueden superar el 50-80% según Dahl et al., 2024), sigue siendo un porcentaje inaceptablemente alto para aplicaciones legales críticas.

Esta persistencia de errores en sistemas RAG tiene una explicación teórica fundamental. Cuando la fase de recuperación (Retrieval) falla o proporciona un contexto ambiguo, el modelo generador se enfrenta a una situación de incertidumbre. Dado que su entrenamiento subyacente lo condiciona a evitar la abstención a toda costa, su comportamiento por defecto es "rellenar los huecos" de la forma más coherente posible, recurriendo a su conocimiento interno. Esto provoca los errores de *misgrounding* que observamos en la práctica, donde el modelo falla no por falta de contexto, sino por su incapacidad estructural para gestionar la incertidumbre de ese contexto (Kalai et al., 2025).

La promesa de una IA legal completamente "libre de alucinaciones" gracias a RAG sigue siendo, en el estado actual de la tecnología, más una aspiración que una realidad consumada. Su adopción debe ir acompañada de un entendimiento realista de sus limitaciones y un compromiso inquebrantable con la verificación humana diligente y el desarrollo continuo de estrategias de mitigación más robustas, como se explorará en la siguiente sección.

4.4 Estrategias avanzadas y holísticas para la optimización de RAG en el contexto legal

La Generación Aumentada por Recuperación (RAG), como se ha analizado previamente, representa un avance conceptual significativo sobre los Grandes Modelos de Lenguaje (LLMs) base, al intentar fundamentar sus respuestas en conocimiento fáctico externo y específico. Sin embargo, la evidencia empírica y el análisis de sus puntos de fallo intrínsecos (Sección 4.3 y 4.4) demuestran con claridad que la implementación canónica de RAG, aunque reduce la incidencia de alucinaciones extrínsecas, está lejos de ser una solución infalible en el exigente y matizado dominio legal. Los desafíos inherentes a la recuperación precisa de información jurídica relevante dentro de corpus masivos y a menudo ambiguos, junto con la propensión residual del LLM generador a desviarse del contexto recuperado o a sintetizarlo incorrectamente, subrayan la necesidad imperante de adoptar estrategias de optimización mucho más sofisticadas y holísticas.

Estas estrategias no se limitan a meros ajustes paramétricos, sino que implican un rediseño y refinamiento profundo de cada componente del ciclo RAG, así como la integración de técnicas complementarias y una comprensión profunda de la interacción entre el conocimiento legal y las capacidades algorítmicas.

Esta sección se adentra en estas metodologías avanzadas, detallando enfoques específicos para la optimización robusta de la recuperación de información, el refinamiento de la fase de generación y razonamiento, y la crucial implementación de arquitecturas integradas que fomenten una sinergia efectiva entre estos componentes, siempre con el objetivo de maximizar la fiabilidad y minimizar el riesgo de alucinación en aplicaciones legales críticas.

4.4.1 Optimización crítica de la fase de recuperación (Retrieval): la calidad del fundamento

La premisa fundamental de RAG es que una base de conocimiento precisa y relevante es el cimiento indispensable para una generación fiable. Por lo tanto, cualquier esfuerzo serio por mejorar la calidad de los sistemas RAG legales debe comenzar con una optimización exhaustiva de la fase de recuperación. No basta con recuperar documentos semánticamente similares; la recuperación debe ser legalmente pertinente, contextualmente adecuada y exhaustiva pero concisa. Las estrategias avanzadas en esta área se centran en ir más allá de las implementaciones ingenuas de búsqueda vectorial y en incorporar una comprensión más profunda de la estructura y la semántica del conocimiento legal.

La optimización de la recuperación va más allá de la mera relevancia semántica; implica seleccionar el **nivel de abstracción correcto** del conocimiento legal a proporcionar. Un estudio sobre la detección del discurso de odio en el derecho alemán reveló que el rendimiento del LLM no siempre mejoraba al proporcionarle más contexto. De hecho, modelos condicionados solo con el título de una norma a menudo superaban a los mismos modelos a los que se les proporcionaba el texto legal completo y complejo (Ludwig et al., 2025). Por el contrario, el rendimiento mejoraba significativamente cuando el contexto incluía **definiciones concretas y ejemplos extraídos de la jurisprudencia**. La implicación para los sistemas RAG es profunda: un recuperador (retriever) óptimo no debe simplemente encontrar el estatuto relevante, sino que debe ser capaz de identificar y extraer de él las definiciones operativas y los ejemplos de

casos que son directamente aplicables, ya que este conocimiento concreto es mucho más "digerible" y útil para el LLM generador que el texto legal abstracto en bruto.

- **Chunking semántico, estructural y adaptativo:** La simple división de documentos en fragmentos de tamaño fijo (*fixed-size chunking*) es a menudo subóptima para textos legales complejos, que poseen una estructura lógica y jerárquica intrínseca (contratos con secciones, cláusulas y sub-cláusulas; sentencias con hechos, razonamiento y holding; estatutos con artículos y apartados).
 - *Chunking consciente de la estructura:* Se deben explorar e implementar técnicas que dividan los documentos respetando estos límites estructurales. Por ejemplo, en un contrato, cada cláusula o sub-cláusula podría constituir un *chunk* individual, preservando su integridad semántica. Esto puede requerir el uso de analizadores sintácticos (parsers) específicos del dominio o expresiones regulares robustas para identificar estos límites estructurales (Pinecone, 2024; Addleshaw Goddard, 2024).
 - *Chunking semántico avanzado:* Más allá de la estructura, se pueden utilizar LLMs más pequeños o modelos de segmentación de texto entrenados para identificar fragmentos que representen unidades de significado coherentes y autocontenidas, incluso si cruzan límites estructurales formales, o para agrupar párrafos temáticamente relacionados.
 - *Chunking recursivo y jerárquico:* Se pueden generar múltiples niveles de *chunks* para un mismo documento: *chunks* pequeños y muy específicos para la recuperación de hechos puntuales, y *chunks* más grandes que capturen el contexto general de una sección o argumento. El sistema podría entonces seleccionar dinámicamente la granularidad de los *chunks* a recuperar en función de la naturaleza de la consulta.
 - *Solapamiento estratégico (Overlap):* Un solapamiento cuidadosamente calibrado entre *chunks* adyacentes sigue siendo crucial para evitar la pérdida de contexto en los límites de los fragmentos, pero su tamaño óptimo puede depender del tipo de documento y la estrategia de *chunking*.
- **Modelos de embedding legales y estrategias multi-vectoriales:** La calidad de la representación vectorial (*embedding*) de los *chunks* y de la consulta es fundamental para la búsqueda semántica.
 - *Embeddings especializados del dominio legal:* El uso de modelos de embedding pre-entrenados o fine-tuneados específicamente en grandes corpus de textos legales (como LegalBERT, bge-m3-spa-lawqa-large, o los modelos desarrollados a partir de The Pile of Law de Henderson et al., 2022) es preferible a los embeddings de propósito general, ya que pueden capturar mejor los matices semánticos y la terminología específica del derecho.
 - *Representaciones multi-vectoriales:* En lugar de un único vector por *chunk*, se podrían generar múltiples vectores que capturen diferentes aspectos del texto: uno para la semántica general, otro para entidades legales clave (tribunales, leyes, partes), otro para conceptos jurídicos abstractos, etc. Esto permitiría búsquedas más matizadas y multifacéticas.
- **Técnicas de búsqueda híbrida y refinamiento de consultas (Query Refinement):** La combinación de diferentes paradigmas de búsqueda y el pre-procesamiento inteligente de la consulta del usuario son clave.
 - *Ponderación optimizada en búsqueda híbrida:* La integración de la búsqueda semántica (vectorial, densa) con la búsqueda tradicional por palabras clave (dispersa, p. ej., BM25) es a menudo superior a cualquiera de los dos métodos por sí solo. La ponderación relativa entre ambos debe ser cuidadosamente ajustada, posiblemente de forma dinámica según la consulta, para equilibrar la captura de significado conceptual con la precisión en la recuperación de términos exactos, nombres propios o citas (Addleshaw Goddard, 2024).
 - *Query Expansion y transformación inteligente:* Utilizar LLMs (posiblemente un modelo más pequeño y rápido dedicado a esta tarea) para pre-procesar la consulta del usuario: expandiéndola con sinónimos legales relevantes, términos relacionados, o posibles reformulaciones; identificando la intención subyacente; o descomponiendo preguntas complejas y multifacéticas en sub-preguntas más simples y manejables que puedan ser abordadas por recuperaciones separadas y luego sintetizadas (HyDE - Gao et al. 2022).
 - *Filtrado estricto por metadatos legales:* La recuperación debe ir más allá de la simple similitud textual e incorporar un filtrado riguroso basado en metadatos cruciales como la jurisdicción aplicable, la fecha de la decisión (para evaluar su actualidad y posible derogación), el nivel jerárquico del tribunal emisor y el tipo de documento. Esto es esencial para asegurar la relevancia legal de los resultados recuperados (Magesh et al., 2024).
- **Mecanismos de recuperación iterativa, auto-correctora y basada en agentes:** Inspirándose en cómo los humanos realizan la investigación legal, los sistemas RAG pueden beneficiarse de arquitecturas más dinámicas e iterativas.

- *Self-Correcting/Corrective RAG (CRAG)*: Implementar bucles de retroalimentación donde el sistema evalúa la relevancia y calidad de un conjunto inicial de documentos recuperados (posiblemente usando el propio LLM generador o un modelo de evaluación dedicado). Si los documentos se consideran insuficientes o irrelevantes, el sistema puede refinar automáticamente la consulta original, ajustar los parámetros de búsqueda o buscar en fuentes alternativas antes de proceder a la generación (Yan et al., 2024).
- *Recuperación multi-salto (Multi-Hop Retrieval)*: Para consultas que requieren la síntesis de información de múltiples fuentes o que implican un razonamiento secuencial (p. ej., rastrear la evolución de una doctrina a través de una cadena de precedentes), el sistema puede realizar múltiples "saltos" de recuperación. La información extraída de un primer conjunto de documentos recuperados se utiliza para formular nuevas consultas y recuperar un segundo conjunto de documentos, y así sucesivamente, hasta que se haya reunido toda la información necesaria (Tang and Yang, 2024).
- *Enfoques basados en agentes (Agentic RAG)*: Desarrollar agentes de IA que puedan planificar y ejecutar estrategias de recuperación complejas, decidiendo dinámicamente qué fuentes consultar, qué términos de búsqueda utilizar y cómo integrar la información obtenida, imitando más de cerca el proceso de investigación de un experto legal.

La inversión en estas estrategias avanzadas de recuperación es fundamental, ya que la calidad del contexto proporcionado al LLM generador es el techo de la calidad de la respuesta final. Una recuperación deficiente o ruidosa inevitablemente conducirá a una generación subóptima o, peor aún, alucinada, independientemente de cuán sofisticado sea el LLM generador.

Modelo	Tamaño	Tareas							Media
		STS	Retrieval	Clasif.	Cluster.	Rerank.	PairClass.	Resumen	
BOW	-	0.4917	0.2143	0.4751	0.2612	0.7582	0.7205	0.2635	0.4549
Encoder based Models									
BERT	110M	0.3821	0.0231	0.5532	0.1803	0.3968	0.7157	0.1723	0.3462
FinBERT	110M	0.4235	0.1178	0.5961	0.2894	0.6453	0.7021	0.2073	0.4259
instructor-base	110M	0.3791	0.5816	0.6253	0.5362	0.9712	0.6185	0.4372	0.5927
bge-large-en-v1.5	335M	0.3435	0.6514	0.6481	0.5768	0.9842	0.7446	0.4911	0.6342
AnglE-BERT	335M	0.3125	0.5784	0.6483	0.5812	0.9673	0.6942	0.5104	0.6132
LLM-based Models									
gte-Qwen1.5-7B-instruct	7B	0.3792	0.6732	0.6479	0.5887	0.9875	0.7042	0.5408	0.6459
Echo	7B	0.4408	0.6487	0.6562	0.5823	0.9751	0.6314	0.4781	0.6304
bge-en-icl	7B	0.3275	0.6831	0.6604	0.5786	0.9912	0.6782	0.5241	0.6347
NV-Embed v2	7B	0.3786	0.7092	0.6432	0.6142	0.9837	0.6098	0.5163	0.6364
e5-mistral-7b-instruct	7B	0.3842	0.6783	0.6492	0.5826	0.9863	0.7432	0.5319	0.6508
Modelos comerciales									
text-embedding-3-small	-	0.3298	0.6694	0.6421	0.5847	0.9847	0.6023	0.5138	0.6181
text-embedding-3-large	-	0.3663	0.7153	0.6631	0.6123	0.9921	0.7358	0.5721	0.6653
voyage-3-large	-	0.4192	0.7509	0.6897	0.5975	0.9951	0.6576	0.6532	0.6805
Modelos adaptados al sector legal									
LegalBERT-v1	7B	0.4215	0.7087	0.7321	0.5843	0.9892	0.7865	0.5217	0.6777
LegalBERT-v2	335M	0.3857	0.6923	0.7284	0.5726	0.9863	0.7614	0.4973	0.6606
BGE-m3-spa-law-qa	1B	0.4128	0.7256	0.7435	0.5932	0.9907	0.7912	0.5163	0.6819
BGE-m3-spa-law-large	7B	0.4387	0.7142	0.7612	0.5697	0.9914	0.8067	0.4846	0.6809

Table 3: Comparación de rendimiento entre diferentes modelos de embeddings en el benchmark FinMTEB. Las métricas de evaluación incluyen similitud semántica textual (STS), recuperación (Retrieval), clasificación (Class.), agrupamiento (Cluster.), reordenamiento (Rerank.), clasificación por pares (PairClass.) y resumen (Summ.). **Los mejores resultados están en negrita.** El subrayado representa el segundo mejor rendimiento.

4.4.2 Refinamiento de la fase de generación y razonamiento: hacia una IA Legal más fiable y transparente

La forma en que se instruye al LLM generador es un componente crítico de la arquitectura RAG, no un mero detalle de implementación. Las siguientes directrices no deben entenderse como simples 'consejos', sino como principios de ingeniería de prompts diseñados para restringir el espacio de posibles respuestas del modelo y alinear su comportamiento

con los exigentes requisitos de fidelidad y trazabilidad del dominio legal. Se trata de codificar explícitamente en las instrucciones las restricciones operativas que garantizan una generación más fiable.

Una vez que se ha recuperado un conjunto de fragmentos de texto contextualmente relevantes (idealmente optimizado a través de las técnicas de la sección anterior), el desafío se traslada a guiar al LLM generador para que utilice esta información de manera fiel, precisa, lógicamente coherente y transparente. Las estrategias para refinar esta fase son cruciales para minimizar el riesgo de que el LLM ignore el contexto, lo malinterprete, o genere afirmaciones que vayan más allá de lo soportado por las fuentes.

- Ingeniería de Prompts avanzada y específica para RAG legal: La forma en que se instruye al LLM generador sobre cómo interactuar con el contexto recuperado es de vital importancia. Los prompts deben ser meticulosamente diseñados para:
 - *Enfatizar la fidelidad al contexto (Grounding Instructions)*: Incluir instrucciones explícitas y prominentes que ordenen al LLM basar su respuesta *exclusivamente* en la información contenida en los documentos proporcionados y evitar activamente el uso de su conocimiento paramétrico interno o la realización de suposiciones no fundamentadas. Esto se logra mediante directivas inequívocas como: 'Responde únicamente basándote en los siguientes extractos legales. No añadas información que no esté presente en los textos proporcionados'.
 - *Guías para el razonamiento (Chain-of-Thought, Step-by-Step)*: Instruir al modelo para que externalice su proceso de razonamiento, mostrando los pasos lógicos que sigue para llegar a una conclusión a partir del contexto recuperado. Por ejemplo, "Primero, identifica las reglas relevantes en el contexto. Segundo, aplica estas reglas a los hechos de la consulta. Tercero, explica tu conclusión" (Wei et al. 2023). Esto no solo puede mejorar la precisión del razonamiento, sino que también hace que el proceso sea más interpretable y verificable por un humano (Schwarcz et al., 2024).
 - *Manejo estructurado de la incertidumbre y los conflictos*: Proporcionar al LLM protocolos claros sobre cómo actuar cuando la información recuperada es incompleta, ambigua, o contiene contradicciones. Esto incluye la instrucción explícita de abstenerse de generar una respuesta cuando no se puede formular con un alto grado de confianza basado en las fuentes, en lugar de recurrir a la especulación. Por ejemplo, "Si la información proporcionada no es suficiente para responder completamente, indícalo explícitamente y explica la naturaleza de la información faltante", "Si encuentras información conflictiva en los extractos, presenta ambas perspectivas y señala la discrepancia".
 - *Instrucciones de citación precisas*: Requerir que el LLM cite de manera específica (idealmente a nivel de fragmento o documento recuperado) las fuentes exactas de las que extrae cada afirmación factual o legal. Esto es esencial para la verificabilidad.
 - *Persona y formato de salida detallados*: Definir con precisión el rol que debe adoptar el LLM (p. ej., "Actúa como un asistente de investigación legal objetivo y neutral") y el formato exacto de la respuesta esperada (p. ej., estructura del resumen, estilo de citación) para asegurar consistencia y utilidad profesional.
 - *Prompting "acusatorio" o de refinamiento (Follow-up Prompts)*: Como se observó en Addleshaw Goddard (2024), el uso de un segundo prompt que cuestione la completitud o exactitud de la respuesta inicial del LLM, acusándolo sutilmente de haber omitido información o pidiéndole que "revise cuidadosamente de nuevo el contexto por si ha pasado algo por alto", puede estimular al modelo a realizar una segunda pasada más exhaustiva del contexto y mejorar significativamente la calidad de la respuesta.
- Fine-Tuning del LLM generador enfocado en la fidelidad legal: Aunque el fine-tuning de LLMs masivos es un proceso intensivo en recursos, puede ofrecer beneficios significativos si se realiza cuidadosamente.
 - *Fine-tuning en tareas de Grounding legal*: Adaptar un LLM pre-entrenado utilizando un conjunto de datos de alta calidad compuesto por pares de (contexto legal recuperado, respuesta ideal fielmente fundamentada y correctamente citada). Esto puede entrenar al modelo para que se adhiera más estrictamente al contexto proporcionado y para que genere respuestas en el estilo y formato deseado por la práctica legal (Tian, Mitchell, Yao, et al. 2023).
 - *Fine-tuning para el razonamiento jurídico sobre contexto*: Desarrollar conjuntos de datos que enseñen al LLM a realizar tipos específicos de razonamiento legal (p. ej., aplicación de reglas, identificación de *holdings*, comparación de casos) *basándose explícitamente* en el contexto recuperado, en lugar de depender de patrones abstractos.
- Integración con modelos de razonamiento especializados: La emergencia de LLMs con arquitecturas explícitamente diseñadas para el razonamiento multi-paso, la planificación y la descomposición de problemas (como la familia de modelos "o" de OpenAI - OpenAI 2024) es particularmente relevante para RAG.

- *Planificación de la respuesta*: Estos modelos podrían, en teoría, planificar cómo utilizar la información recuperada de manera más estratégica, identificando qué fragmentos son más relevantes para qué partes de la consulta y cómo sintetizarlos de manera lógicamente coherente.
- *Verificación interna de pasos de razonamiento*: Su capacidad para "reflexionar" sobre sus propios pasos de razonamiento intermedios podría permitirles detectar y corregir errores o inconsistencias antes de generar la respuesta final (Schwarcz et al., 2024). La integración de estos modelos de razonamiento como el componente generador en un sistema RAG es un área prometedora para futuras mejoras.
- Arquitecturas híbridas (simbólico-neuronales): Aunque aún en etapas tempranas para aplicaciones legales complejas, la combinación de LLMs neuronales con sistemas de razonamiento simbólico basados en reglas (p. ej., lógicas formales, ontologías legales) podría ofrecer una vía para mejorar la consistencia lógica y la verificabilidad de las respuestas generadas a partir del contexto recuperado.
- Adaptación de la Complejidad del Lenguaje (La Función "Jerga"): Una IA legal verdaderamente avanzada no solo debe ser precisa, sino también contextualmente consciente del receptor final de la información. La optimización del prompt debe incluir instrucciones para modular la complejidad del lenguaje de la respuesta. Por ejemplo, un sistema podría recibir la directriz: "Genera una respuesta técnica para un abogado y, adicionalmente, una explicación simplificada para un ciudadano sin conocimientos jurídicos". Esta capacidad, que podríamos denominar "función jerga", es un pilar de la humanización de la tecnología legal, reconociendo que la utilidad de una respuesta no reside solo en su corrección, sino en su comprensibilidad. Esto transforma a la IA de un simple motor de búsqueda a un verdadero puente de comunicación entre el complejo mundo legal y la sociedad.

El objetivo final de estas estrategias de optimización de la generación no es solo producir respuestas que *parezcan* correctas, sino respuestas que sean *demostrablemente* correctas, fieles a las fuentes proporcionadas y útiles para el profesional legal. La capacidad del LLM para explicar *cómo* llegó a una conclusión a partir del contexto recuperado es tan importante como la conclusión misma.

La implementación exitosa de RAG en el dominio legal, por lo tanto, no es simplemente una cuestión de conectar un LLM a una base de datos. Requiere un diseño cuidadoso y una optimización continua de cada etapa del proceso, desde la curación de datos y el chunking, pasando por la sofisticación de los algoritmos de recuperación y la ingeniería de prompts, hasta el refinamiento de la capacidad de razonamiento y generación fiel del LLM. Solo a través de este enfoque holístico y riguroso se podrá comenzar a materializar verdaderamente el potencial de RAG para mitigar las alucinaciones y ofrecer una IA legal genuinamente fiable y valiosa.

5 Avanzando hacia la fiabilidad: estrategias holísticas para la optimización y mitigación de alucinaciones en la Inteligencia Artificial legal

La constatación de que ni los Grandes Modelos de Lenguaje (LLMs) base ni las implementaciones canónicas de Generación Aumentada por Recuperación (RAG) logran erradicar por completo el espectro de las alucinaciones en el dominio legal, impone un cambio de paradigma en la forma en que abordamos el desarrollo y la integración de estas tecnologías. Ya no es suficiente aspirar a una solución única o a un "interruptor mágico" que elimine los errores; en su lugar, se requiere un **enfoque holístico, multifacético y adaptable** que reconozca la complejidad inherente al problema y que implemente una sinergia de estrategias de optimización y mitigación a lo largo de todo el ciclo de vida de la información, desde la curación de los datos hasta la verificación final del resultado generado. Este enfoque no busca la perfección absoluta —una meta teóricamente inalcanzable en sistemas tecnológicos complejos—, sino que adopta un principio de ingeniería robusta: **la maximización de la fiabilidad y la minimización del riesgo dentro de los límites de lo factible**. Se reconoce que la infalibilidad del 100% no es una limitación de la tecnología *actual*, sino una característica inherente a la complejidad. Por tanto, el objetivo es construir un sistema cuya fiabilidad sea tan alta, y sus modos de fallo tan predecibles, que la supervisión humana experta se convierta en una capa de validación eficiente y no en una búsqueda onerosa de errores ocultos, siempre bajo la égida indispensable del juicio y la responsabilidad profesional. Esta sección se dedica a explorar en profundidad este arsenal de estrategias avanzadas, que van más allá de los ajustes superficiales para adentrarse en la optimización rigurosa de los datos, el refinamiento de los procesos de razonamiento algorítmico —incluyendo la consideración de sistemas jerárquicos y agentes de IA—, la implementación de mecanismos de verificación cada vez más sofisticados y, de manera crucial, el fortalecimiento del rol del profesional legal como supervisor crítico e informado. La conjunción de estas técnicas no solo busca reducir la frecuencia de las alucinaciones, sino también transformar su naturaleza, haciendo que los errores residuales sean más detectables y menos perjudiciales.

La eficacia de un enfoque holístico, que combina la curación estratégica de datos con la especialización de modelos, ha sido demostrada empíricamente. Un claro ejemplo es el proyecto "Legal Assist AI", que abordó el problema de las

alucinaciones en el contexto legal indio. En lugar de utilizar un modelo de propósito general, los investigadores crearon un corpus de datos curado a partir de fuentes legales indias (Constitución, estatutos, jurisprudencia) y lo utilizaron para afinar (fine-tune) un modelo base de 8 mil millones de parámetros (Llama 3.1 8B). El resultado fue un modelo especializado que no solo superó a modelos mucho más grandes como GPT-3.5 Turbo (175 mil millones de parámetros) en tareas de abogacía de la India, sino que, de manera crucial, redujo drásticamente la generación de alucinaciones, proporcionando respuestas fiables donde otros modelos inventaban información (Gupta et al., 2025). Este caso de éxito sirve como un poderoso testimonio de que la mitigación de alucinaciones no reside en la escala del modelo, sino en la calidad de los datos y la especificidad del entrenamiento.

Las estrategias de mitigación más avanzadas convergen en un principio de integración de conocimiento, donde los LLMs no operan de forma aislada, sino como parte de arquitecturas híbridas. Esto incluye la integración con grafos de conocimiento legal estructurados y el uso de arquitecturas de Mixture-of-Experts (MoE), como se implementa en modelos de vanguardia como ChatLaw (Shao et al., 2025). En estos sistemas, módulos expertos especializados dentro del LLM se activan dinámicamente para manejar diferentes tipos de tareas legales, reduciendo las alucinaciones en un 38% al asegurar que la consulta sea gestionada por el componente con el conocimiento más relevante.

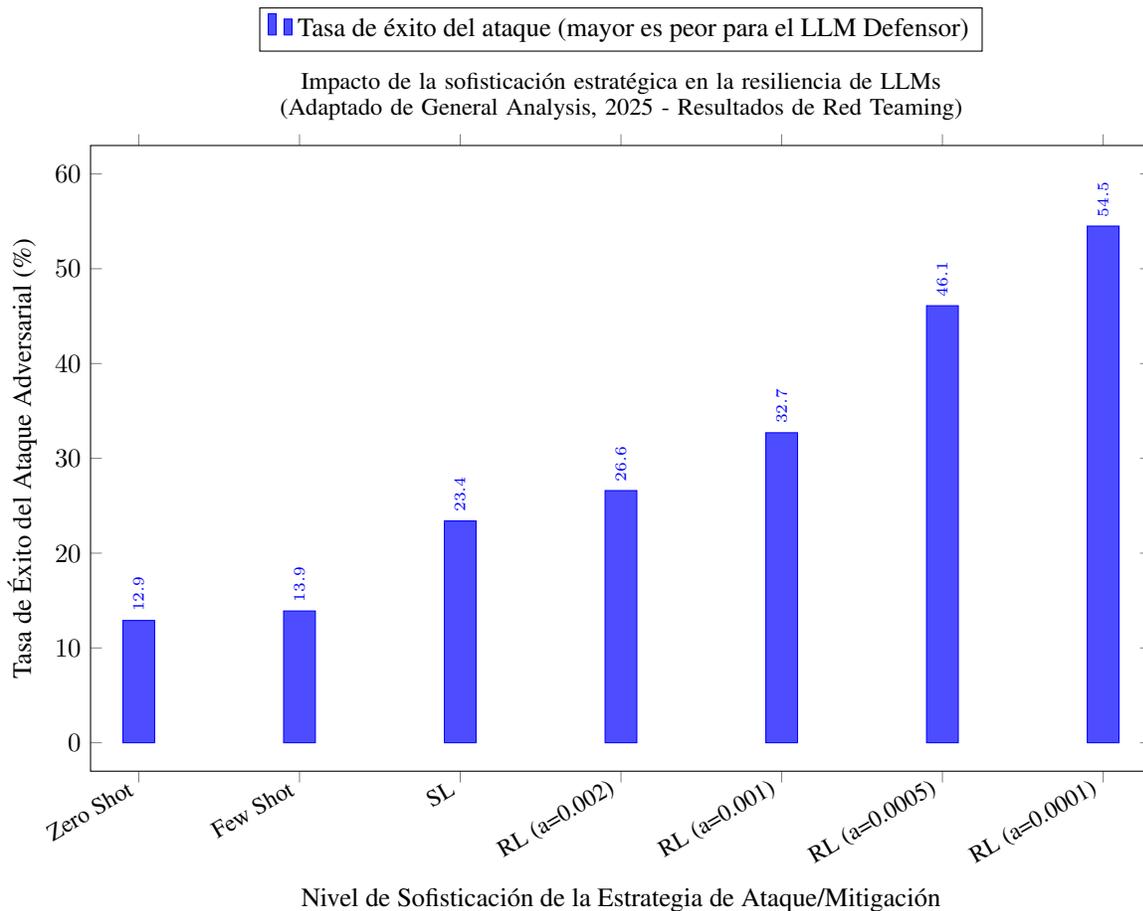


Figure 4: Esta gráfica ilustra cómo estrategias de ataque más sofisticadas (análogas a la falta de estrategias de mitigación robustas o a vulnerabilidades explotadas) logran una mayor tasa de éxito al inducir fallos en un LLM objetivo (GPT-4o). Demuestra la necesidad de estrategias de defensa (mitigación) igualmente sofisticadas. El eje X representa diferentes algoritmos de ataque del estudio de General Analysis, interpretados aquí como niveles de desafío o sofisticación a los que un sistema RAG legal debe ser resiliente.

5.1 La calidad del fundamento: curación estratégica de datos y bases de conocimiento externo

El adagio "basura entra, basura sale" (*garbage in, garbage out*) resuena con especial fuerza en el contexto de los LLMs y los sistemas RAG. La calidad, actualidad, relevancia y representatividad de la base de conocimiento externa sobre la cual se fundamentan estos sistemas no es un mero detalle técnico, sino el cimiento sobre el que se construye toda su

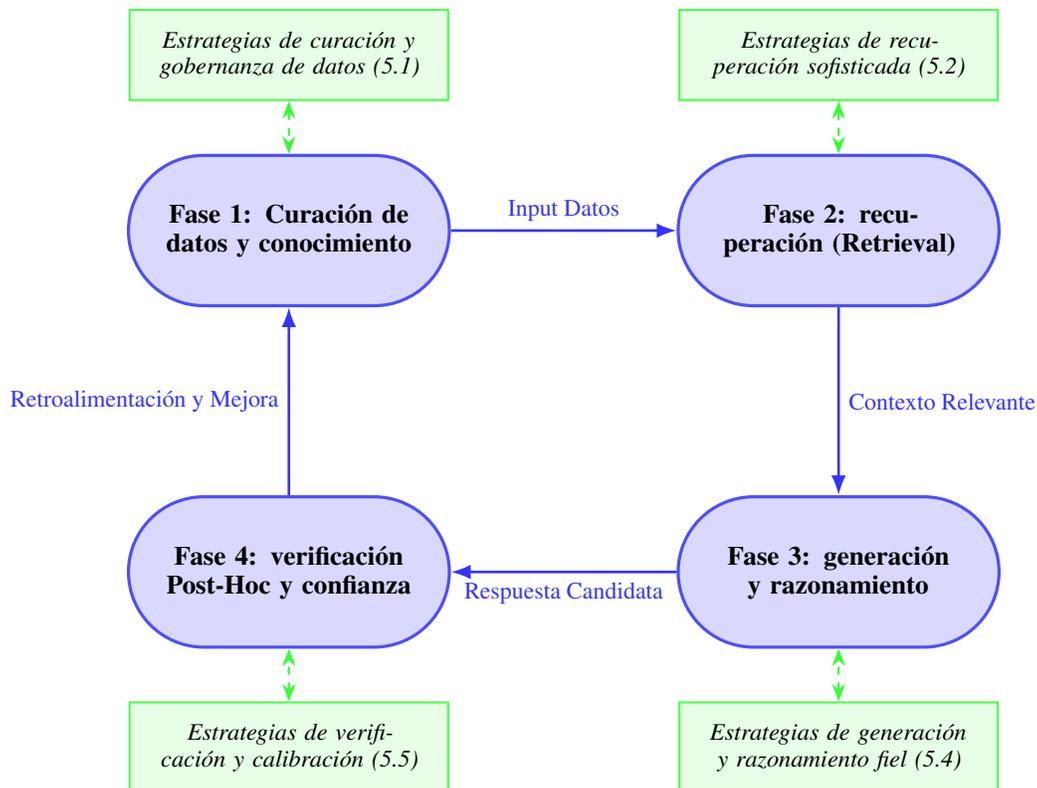


Figure 5: Modelo cíclico de un sistema RAG Legal y puntos de intervención estratégica para la optimización y mitigación de alucinaciones. Las estrategias específicas (referenciadas por subsección del ensayo) se aplican en cada fase para mejorar la fiabilidad general del sistema.

fiabilidad. Un corpus documental deficiente, desactualizado o sesgado inevitablemente limitará la capacidad del sistema RAG para proporcionar respuestas precisas y confiables, independientemente de cuán sofisticados sean sus algoritmos de recuperación o generación. Por lo tanto, una estrategia primordial y proactiva para la mitigación de alucinaciones comienza mucho antes de la interacción con el usuario: en la meticulosa curación y gestión estratégica de estas bases de conocimiento.

5.1.1 Selección y priorización rigurosa de fuentes legales

El vasto universo de información legal exige una selección criteriosa. No todas las fuentes son iguales en autoridad o relevancia. Es imperativo tener en consideración:

- *Jerarquización de la autoridad:* Diseñar mecanismos, tanto en la indexación como en la recuperación, que prioricen explícitamente las fuentes primarias vinculantes (Constitución, estatutos vigentes, jurisprudencia de tribunales superiores de la jurisdicción pertinente) sobre fuentes secundarias, literatura persuasiva, o jurisprudencia de otras jurisdicciones o tribunales inferiores. Esto implica la incorporación de metadatos ricos que codifiquen esta jerarquía y permitan al sistema RAG ponderar la información en consecuencia.
- *Verificación continua de la actualidad y vigencia:* Implementar procesos dinámicos y automatizados (en la medida de lo posible, complementados con revisión experta) para mantener la base de conocimiento al día con las enmiendas legislativas, las nuevas decisiones judiciales y, de manera crítica, el estado de derogación de los precedentes. La integración con servicios comerciales de verificación de citas o bases de datos legislativas actualizadas es fundamental. Herramientas como Shepard's de LexisNexis o KeyCite de Westlaw (estándares en el sistema de *common law* estadounidense) son cruciales para trazar el historial de un caso y verificar su vigencia. En España, si bien no existe un equivalente directo con una marca tan consolidada, plataformas jurídicas líderes como vLex o La Ley Digital ofrecen funcionalidades análogas que permiten comprobar si una sentencia ha sido recurrida, anulada o matizada, así como verificar la vigencia de una norma. La integración o, como mínimo, la consulta sistemática de estas herramientas es un paso ineludible para evitar fundamentar respuestas en ley obsoleta, un error común y potencialmente grave.

- *Filtrado proactivo de fuentes de baja calidad o problemáticas:* Identificar y excluir o marcar explícitamente fuentes conocidas por su baja calidad, sesgos manifiestos (si el objetivo es un análisis neutral), o irrelevancia para las tareas legales más comunes. Esto puede requerir tanto el juicio de expertos legales como el uso de técnicas de IA para la evaluación automática de la calidad y fiabilidad de los documentos (Nguyen and Satoh, 2024). Estas técnicas van más allá de la simple detección de palabras clave. Incluyen, por ejemplo, el uso de modelos de lenguaje más pequeños y especializados que actúan como 'jueces' o evaluadores (un enfoque conocido como *LLM-as-a-judge*), capaces de verificar la consistencia lógica de un documento, detectar contradicciones internas o contrastar afirmaciones contra una base de conocimiento curada. Otras técnicas implican el análisis de la confianza del modelo durante la generación o la detección de anomalías estilísticas que a menudo acompañan a las alucinaciones. La implementación de estos 'guardianes' algorítmicos puede servir como un primer filtro automatizado antes de la revisión humana.

La importancia de estas prácticas de curación, priorización y gobernanza de los datos que alimentan los sistemas de IA legal se ve magnificada por los emergentes marcos regulatorios. El Reglamento (UE) 2024/1689 del Parlamento Europeo y del Consejo, de 13 de junio de 2024, por el que se establecen normas armonizadas en materia de inteligencia artificial (en adelante, la Ley de IA de la UE o el Reglamento), en su Artículo 10, impone a los desarrolladores de sistemas de IA de alto riesgo obligaciones explícitas respecto a los conjuntos de datos de entrenamiento, validación y prueba. Estos deben ser 'relevantes, representativos, libres de errores y completos', y deben implementarse prácticas adecuadas de gobernanza de datos, incluyendo un examen de los posibles sesgos. El cumplimiento de estas exigencias no es solo una cuestión de buena práctica técnica para mejorar la fiabilidad del modelo y reducir el riesgo de alucinaciones originadas en datos defectuosos, sino que se perfila como un requisito legal ineludible para operar en el mercado europeo, incentivando una mayor diligencia en la gestión del conocimiento que fundamenta la IA legal.

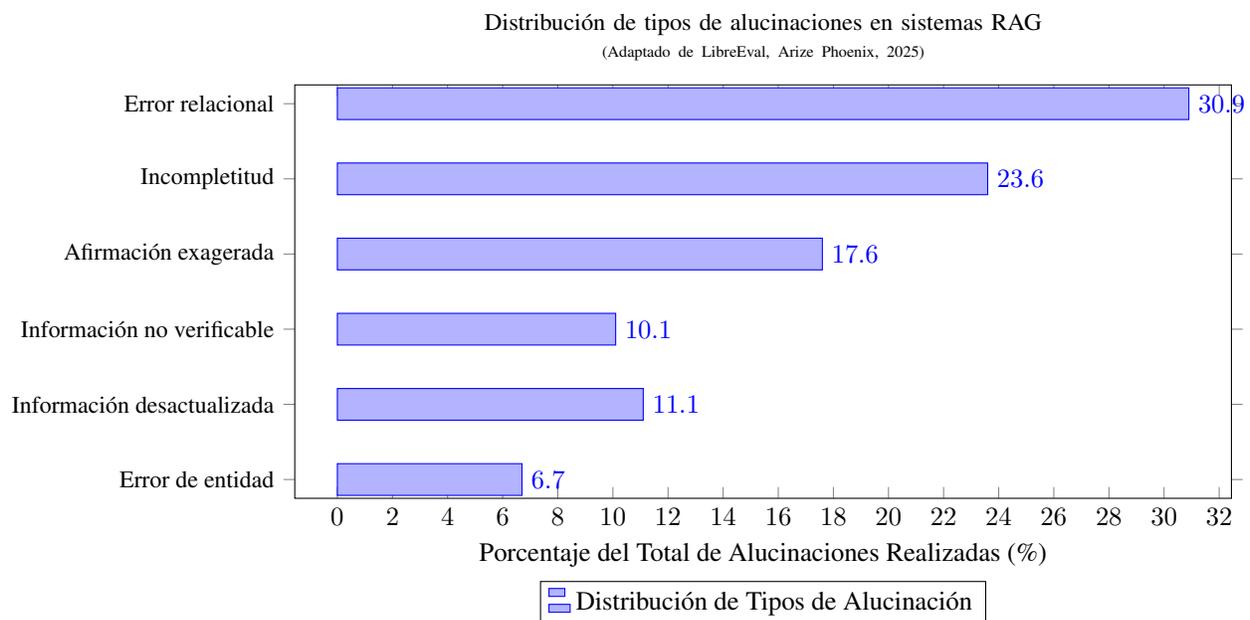


Figure 6: Distribución porcentual de los diferentes tipos de alucinaciones efectivamente realizadas en las respuestas de modelos de lenguaje con RAG, según el dataset LibreEval1.0. Esta distribución destaca los desafíos más comunes que las estrategias de mitigación deben abordar.

5.1.2 Aseguramiento de la diversidad y representatividad del corpus

Para evitar la creación de "monoculturas legales" algorítmicas (Dahl et al., 2024) que ignoren la diversidad del pensamiento jurídico o las particularidades de jurisdicciones menos prominentes, es crucial:

- *Cobertura jurisdiccional y temática amplia:* Esforzarse por incluir una representación equilibrada de todas las jurisdicciones relevantes (estatales, federales, especializadas) y de un amplio espectro de áreas del derecho, no solo aquellas con mayor disponibilidad de datos digitalizados. En el contexto español, por ejemplo, los datos del Consejo General del Poder Judicial (CGPJ) muestran consistentemente una alta concentración de litigiosidad en los órganos judiciales de grandes capitales como Madrid o Barcelona. Un corpus de

entrenamiento que no pondere adecuadamente esta realidad podría desarrollar un sesgo centralista, ignorando las particularidades doctrinales o jurisprudenciales de otros Tribunales Superiores de Justicia, lo que limitaría la utilidad de la herramienta a nivel nacional..

- *Inclusión consciente de perspectivas plurales:* Para tareas que requieren un análisis más allá de la doctrina pura (p. ej., evaluación de impacto de políticas, argumentación basada en principios), considerar la inclusión curada de fuentes académicas críticas, informes de organizaciones de la sociedad civil, o incluso transcripciones de debates legislativos que ofrezcan perspectivas diversas y matizadas sobre la ley y su aplicación.

5.1.3 Estructuración avanzada del conocimiento legal (más allá del texto plano)

Mientras que muchos sistemas RAG operan principalmente sobre texto no estructurado, la representación del conocimiento legal puede enriquecerse significativamente mediante:

- *Desarrollo de ontologías y grafos de conocimiento legales:* Construir modelos formales que representen entidades legales clave (tribunales, jueces, leyes, conceptos jurídicos), sus atributos y, fundamentalmente, las relaciones complejas entre ellas (p. ej., una ley *enmienda* otra, un caso *interpreta* un estatuto, un juez *disiente* de una opinión). Un sistema RAG que pueda consultar y razonar sobre estos grafos de conocimiento podría realizar inferencias más profundas y precisas que uno basado únicamente en similitud textual (Martín, 2024; Magora, 2024).
- *Extracción y vinculación de metadatos ricos:* Enriquecer cada documento en la base de conocimiento con metadatos detallados y estructurados (jurisdicción, fecha, tribunal, jueces, partes, temas legales, citas, historial procesal, estado de vigencia) que puedan ser utilizados por el módulo de recuperación para un filtrado y ranking mucho más preciso.

La inversión continua en la calidad, estructura y gestión de las bases de conocimiento no es un costo accesorio, sino una inversión estratégica fundamental en la fiabilidad a largo plazo de cualquier sistema de IA legal basado en RAG. Sin un fundamento sólido, incluso los algoritmos más avanzados estarán construyendo sobre arena movediza.

5.2 Optimización sofisticada de la fase de recuperación (Retrieval): encontrando la aguja jurídica en el pajar digital

La eficacia de un sistema RAG depende de forma crítica de la capacidad de su módulo de recuperación para identificar y extraer, de entre un corpus potencialmente masivo, los fragmentos de texto (chunks) que son *exacta* y *contextualmente relevantes* para la consulta del usuario. La simple búsqueda por similitud semántica, aunque un punto de partida, a menudo resulta insuficiente para la complejidad y los matices del lenguaje y el razonamiento jurídico. La optimización avanzada de esta fase es, por lo tanto, un área de intensa investigación y desarrollo, enfocada en dotar al sistema de una capacidad de discernimiento más cercana a la de un investigador legal humano.

5.2.1 Modelos de embedding legales y estrategias multi-vectoriales

La calidad de la representación vectorial (embedding) que captura el significado de los chunks y de la consulta es la piedra angular de la búsqueda semántica.

- *Embeddings especializados del dominio legal:* Priorizar el uso de modelos de embedding que hayan sido pre-entrenados o fine-tuneados específicamente en grandes corpus de textos jurídicos (como LegalBERT o los modelos desarrollados a partir de corpus como The Pile of Law - Henderson et al. 2022) es preferible a los embeddings de propósito general, ya que pueden capturar mejor los matices semánticos y la terminología específica del derecho.
- *Técnicas de aumento de embeddings:* Explorar métodos que enriquezcan los embeddings textuales con información adicional, como metadatos estructurales o información de grafos de conocimiento, para crear representaciones más ricas y contextualmente informadas.

5.2.2 Estrategias de búsqueda híbrida, multi-Etapa y refinamiento de consultas (Query Refinement):

Para superar las limitaciones de una única modalidad de búsqueda, se están adoptando enfoques más complejos:

- *Optimización dinámica de búsqueda híbrida:* La combinación de la búsqueda semántica (vectorial, densa) con la búsqueda tradicional por palabras clave (léxica, p. ej., BM25) es fundamental. Sin embargo, la ponderación

relativa entre estas dos modalidades no debe ser estática. Idealmente, el sistema debería ser capaz de ajustar dinámicamente esta ponderación basándose en las características de la consulta del usuario (p. ej., dar más peso a las palabras clave si la consulta contiene términos muy específicos, nombres propios o citas exactas) (Addleshaw Goddard, 2024).

- *Pre-Procesamiento y transformación inteligente de consultas*: Utilizar un LLM (posiblemente un modelo más pequeño y eficiente dedicado a esta tarea) para analizar y refinar la consulta del usuario *antes* de la fase de recuperación. Esto puede incluir:
 - *Query Expansion*: Añadir sinónimos legales relevantes, términos conceptualmente relacionados o posibles reformulaciones para ampliar la cobertura de la búsqueda.
 - *Descomposición de consultas (Query Decomposition)*: Dividir preguntas complejas o multifacéticas en sub-preguntas más simples y atómicas, cada una de las cuales puede ser objeto de una recuperación separada. Los resultados de estas sub-recuperaciones se pueden luego combinar para responder a la consulta original.
 - *Generación de consultas hipotéticas (Hypothetical Document Embeddings - HyDE)*: Instruir a un LLM para que genere un documento "ideal" que respondería perfectamente a la consulta del usuario, y luego usar el embedding de este documento hipotético para la búsqueda semántica. Esto a menudo conduce a una recuperación más relevante que usar el embedding de la consulta original directamente (Gao et al., 2022).
- *Recuperación multi-etapa (Multi-Hop Retrieval)*: Para consultas que requieren un razonamiento secuencial o la síntesis de información a través de una cadena de documentos (p. ej., rastrear la evolución de una doctrina legal a través de múltiples precedentes), el sistema puede implementar un proceso de recuperación iterativo. La información extraída de un primer conjunto de documentos recuperados se utiliza para formular nuevas consultas o refinar las existentes, permitiendo al sistema "navegar" a través del corpus documental de manera más inteligente y dirigida (Tang and Yang, 2024).
- *Re-ranking sofisticado (Re-ranking)*: Una vez que se ha recuperado un conjunto inicial de chunks candidatos (posiblemente amplio), utilizar modelos de re-ranking más potentes y computacionalmente intensivos (como los cross-encoders) para evaluar de manera más precisa la relevancia de cada chunk en relación con la consulta completa. Estos modelos pueden considerar la interacción entre la consulta y el chunk de manera más profunda que los modelos de embedding utilizados en la recuperación inicial, mejorando el orden final de los resultados presentados al LLM generador.

5.2.3 Incorporación de retroalimentación y aprendizaje continuo

Los sistemas RAG más avanzados deberían incorporar mecanismos para aprender de las interacciones con los usuarios y de la retroalimentación explícita o implícita.

- *Retroalimentación del usuario*: Permitir a los usuarios calificar la relevancia de los documentos recuperados o de las respuestas generadas, y utilizar esta retroalimentación para fine-tunear los modelos de embedding o los algoritmos de ranking.
- *Adaptación Dinámica*: Ajustar los parámetros de recuperación (p. ej., umbrales de similitud, número de chunks a recuperar) basándose en el rendimiento histórico o en las características de la consulta actual.

5.2.4 Conclusión

La inversión en estas estrategias avanzadas de recuperación es fundamental, ya que la calidad del contexto proporcionado al LLM generador es el techo de la calidad de la respuesta final. Una recuperación deficiente o ruidosa inevitablemente conducirá a una generación subóptima o, peor aún, alucinada, independientemente de cuán sofisticado sea el LLM generador.

5.3 Agentes de IA en sistemas legales complejos y la jerarquía normativa de Kelsen

La optimización de los sistemas RAG para tareas legales, como se ha discutido, implica una interacción cada vez más sofisticada entre el LLM, las bases de conocimiento y el proceso de razonamiento. A medida que esta sofisticación aumenta y que los LLMs evolucionan hacia capacidades de razonamiento más complejas y multi-paso, comenzamos a vislumbrar el potencial de arquitecturas de IA más autónomas y proactivas, a menudo denominadas 'agentes de IA'. En este contexto, un agente de IA legal se distingue de un LLM simple (que meramente responde a un prompt) por su capacidad para realizar secuencias de acciones, interactuar con múltiples herramientas o fuentes de información

de manera autónoma, tomar decisiones intermedias, y planificar estrategias para lograr un objetivo legal complejo predefinido, como podría ser la preparación de un caso o la realización de una due diligence exhaustiva. Sin embargo, la implementación de tales agentes en el intrincado y normativamente estructurado dominio legal debe considerar ineludiblemente la arquitectura fundamental del propio sistema jurídico. En los ordenamientos de derecho civil y en muchos sistemas constitucionales, esta arquitectura se conceptualiza clásicamente a través de la noción de la Pirámide Normativa de Hans Kelsen.

La pirámide de Kelsen postula que las normas jurídicas de un sistema se organizan jerárquicamente, donde la validez de cada norma deriva de una norma superior, culminando en una "Norma Fundamental" (Grundnorm) hipotética que fundamenta la validez de todo el sistema (generalmente la Constitución en los sistemas modernos). Esta estructura jerárquica implica que las normas de rango inferior (p.ej., un reglamento administrativo) deben ser conformes con las normas de rango superior (p.ej., una ley, la Constitución). La aplicación correcta del derecho, por tanto, no es solo una cuestión de encontrar *una* norma relevante, sino de encontrar la norma *correcta* dentro de esta jerarquía y resolver posibles conflictos entre normas de diferente rango (principio de jerarquía normativa) o del mismo rango pero posteriores en el tiempo (principio de temporalidad) o más específicas (principio de especialidad).

Para un agente de IA legal que opere con el objetivo de proporcionar respuestas fiables y legalmente válidas, esta estructura jerárquica presenta tanto un desafío como una oportunidad:

1. Desafío para la recuperación y el razonamiento:

- *Identificación de la autoridad controladora:* Un agente de IA debe ser capaz no solo de recuperar múltiples normas o precedentes potencialmente relevantes, sino de discernir cuál de ellos tiene precedencia o es la autoridad controladora en un caso dado. Por ejemplo, si una ley parece contradecir una disposición constitucional, la Constitución prevalece. Si un reglamento contradice una ley, la ley prevalece. Esta inferencia jerárquica es esencial.
- *Resolución de antinomias:* El agente debe ser capaz de identificar y, en la medida de lo posible, proponer soluciones a conflictos normativos (antinomias) utilizando los criterios de resolución aceptados (jerarquía, temporalidad, especialidad). Esto requiere un nivel de razonamiento meta-legal.
- *Comprensión de la dinámica de validez:* La validez de una norma puede cambiar (p. ej., una ley puede ser declarada inconstitucional, un precedente puede ser derogado). El agente debe acceder y procesar información sobre el estado actual de validez de las normas recuperadas.

Consideremos, por ejemplo, una consulta sobre la legalidad de una determinada práctica comercial municipal. Un agente de IA simple, sin conciencia jerárquica, podría recuperar y basar su respuesta afirmativa en una ordenanza municipal que explícitamente permite dicha práctica. Sin embargo, un agente 'kelseniano', en su proceso de verificación ascendente, identificaría una ley autonómica o estatal posterior y de rango superior que prohíbe o restringe severamente tal actividad, o incluso una sentencia del Tribunal Constitucional que haya declarado inconstitucional una norma similar. Este agente concluiría correctamente que la ordenanza municipal, aunque textualmente aplicable, es inválida o inaplicable debido al conflicto con una norma jerárquicamente superior, evitando así una alucinación de validez que el agente simple habría cometido. La capacidad para realizar este tipo de validación jerárquica es, por tanto, crucial para la fiabilidad.

2. Oportunidad para sistemas RAG jerárquicamente conscientes: La pirámide de Kelsen puede, de hecho, inspirar arquitecturas RAG más sofisticadas y fiables:

- *Bases de conocimiento estructuradas jerárquicamente:* Las bases de conocimiento podrían organizarse explícitamente reflejando la jerarquía normativa. Los documentos podrían etiquetarse con su rango jerárquico, utilizando un esquema de metadatos que refleje la estructura del ordenamiento. Esto incluye categorías universales como 'Norma Constitucional', 'Legislación Primaria' (leyes), 'Legislación Secundaria' (reglamentos) y 'Jurisprudencia Vinculante'. En el contexto específico del derecho español, esto se traduciría en etiquetas como (Constitución, ley orgánica, ley ordinaria, reglamento, jurisprudencia del Tribunal Constitucional, etc.).
- *Algoritmos de recuperación sensibles a la jerarquía:* El módulo de recuperación podría ser instruido para priorizar la búsqueda y recuperación de normas de rango superior cuando sean pertinentes, o para buscar específicamente normas que interpreten o apliquen una norma superior identificada.
- *Módulos de razonamiento para la coherencia jerárquica:* Un LLM generador, o un componente de razonamiento especializado, podría ser entrenado para verificar la coherencia de una solución propuesta con las normas de rango superior. Si una interpretación de un contrato parece violar una ley imperativa, el agente podría señalar esta inconsistencia.
- *Agentes planificadores que navegan la pirámide:* Un agente de IA más avanzado podría planificar su proceso de investigación y razonamiento comenzando por la cúspide de la pirámide (Constitución,

tratados internacionales relevantes) y descendiendo a través de las leyes y la jurisprudencia aplicable, asegurando que cada paso sea consistente con el nivel superior.

Implicaciones para la fiabilidad y las alucinaciones:

Ignorar la jerarquía normativa puede llevar a un tipo específico y grave de alucinación legal: la **alucinación de invalidez o inaplicabilidad por conflicto jerárquico**. Un LLM podría, por ejemplo, basar una respuesta en un reglamento que, aunque textualmente relevante, es inválido porque contradice una ley superior, o fundamentar un argumento en jurisprudencia de un tribunal inferior que ha sido revocada o matizada por un tribunal superior. Estas no son simples inexactitudes factuales, sino errores fundamentales en la aplicación del derecho.

Por el contrario, un agente de IA que esté explícitamente modelado para comprender y operar dentro de la pirámide kelseniana podría:

- **Reducir las alucinaciones de relevancia:** Al priorizar fuentes de mayor autoridad, es menos probable que se base en información legalmente subordinada o irrelevante.
- **Mejorar la solidez del razonamiento:** Al verificar la coherencia con normas superiores, sus conclusiones serían más robustas y menos susceptibles de ser invalidadas.
- **Aumentar la transparencia y explicabilidad:** Al poder trazar la derivación de una conclusión a través de la jerarquía normativa, el agente podría ofrecer explicaciones más convincentes y verificables de su razonamiento.

En definitiva, un agente de IA 'kelseniano', es decir, uno que no solo acceda a las fuentes sino que comprenda y respete la jerarquía normativa y los principios de validez del ordenamiento, sería intrínsecamente menos propenso a ciertos tipos críticos de alucinaciones. Al priorizar la Constitución sobre la ley, y la ley sobre el reglamento, y al verificar la vigencia y aplicabilidad de cada norma dentro de su contexto jerárquico, se reduciría drásticamente el riesgo de generar consejos basados en normas subordinadas invalidadas, en jurisprudencia derogada o en interpretaciones que contradicen principios fundamentales. Esta conciencia estructural no elimina todos los riesgos de alucinación –especialmente aquellos derivados de la ambigüedad inherente del lenguaje o de los límites del propio corpus de conocimiento– pero sí proporciona un andamiaje robusto para una IA legal más coherente, predecible y, en última instancia, más fiable. Es crucial entender que el valor de este andamiaje no reside únicamente en la mejora de la respuesta final, sino en la propia externalización del proceso de razonamiento. Al hacer explícita la jerarquía normativa que aplica, la IA pasa de ofrecer una conclusión opaca a presentar un argumento verificable. Para el profesional del derecho, esta fundamentación transparente es, en muchos casos, más valiosa que la respuesta misma, pues le permite auditar, validar y, en última instancia, apropiarse del razonamiento para construir su propia estrategia jurídica. Transforma la IA de una 'caja negra' a una 'caja de herramientas' de razonamiento.

En el contexto español y de muchos sistemas de derecho civil europeos, donde la codificación y la jerarquía formal de las fuentes del derecho son particularmente pronunciadas, la incorporación de una conciencia kelseniana en los agentes de IA legales no es un mero refinamiento académico, sino una condición necesaria para su fiabilidad y utilidad práctica. Un agente que no "entienda" la estructura piramidal del ordenamiento jurídico será intrínsecamente propenso a generar respuestas que, aunque plausiblemente redactadas, sean legalmente insostenibles o directamente erróneas. El desarrollo futuro de la IA legal fiable pasará, ineludiblemente, por dotar a estos sistemas de una comprensión más profunda de la arquitectura fundamental del derecho mismo.

5.4 Refinamiento estratégico de la fase de generación y razonamiento: cultivando la fidelidad y la coherencia en la IA Legal

Una vez que el sistema de Generación Aumentada por Recuperación (RAG) ha completado la fase crítica de recuperación, proporcionando al Gran Modelo de Lenguaje (LLM) un conjunto de fragmentos de texto contextualmente relevantes (idealmente optimizado a través de las técnicas de la sección anterior), el desafío se traslada a la fase de generación. Aquí, el objetivo es guiar al LLM para que utilice esta información recuperada de manera que la respuesta final no solo sea lingüísticamente fluida y coherente, sino, y de manera crucial, fácticamente precisa, lógicamente sólida, fiel a las fuentes proporcionadas y directamente relevante para la consulta original del usuario. La mera provisión de contexto no garantiza una generación de alta calidad; el LLM generador, por su naturaleza probabilística y sus vastos conocimientos paramétricos, aún puede desviarse, malinterpretar o incluso alucinar. Por lo tanto, el refinamiento estratégico de esta fase es un componente esencial de cualquier sistema RAG legal que aspire a la fiabilidad.

1. **Ingeniería de prompts avanzada, específica para RAG y consciente del contexto legal:** El prompt que se alimenta al LLM generador, que ahora incluye tanto la consulta original del usuario como los fragmentos de texto recuperados, debe ser meticulosamente diseñado para maximizar la fidelidad y la precisión. Esto va mucho más allá de una simple concatenación.

Table 4: Impacto cuantificable estimado de estrategias de optimización en sistemas RAG Legales

Estrategia de optimización aplicada	Métrica clave	Valor inicial	Valor optimizado	Mejora relativa (%)
Optimización intensiva de recuperación (Chunking, Embeddings, Query Expansion) <i>(Inspirado en Addleshaw Goddard, 2024)</i>	F1-Score (Extracción de cláusulas)	74	95	28.4
Fine-Tuning del LLM Generador para fidelidad al contexto	tasa de <i>Misgrounding</i>	20	5	75.0
Implementación de Chain-of-Thought (CoT) en prompts de generación	coherencia lógica (Puntuación Humana 1-5)	3.2	4.5	40.6
Verificación Post-Hoc automatizada con modelo secundario	tasa de alucinaciones no detectadas	15	3	80.0
Implementación de agente consciente de jerarquía (Kelseniano)	Tasa de error por conflicto normativo	15	2	86.7

Nota: Los valores para "Optimización de Recuperación" están inspirados en los resultados F1 reportados por Addleshaw Goddard (2024). Otros valores son hipotéticos y presentados con fines ilustrativos para demostrar el potencial impacto de diversas estrategias de optimización discutidas en la Sección 5. La "Mejora Relativa" se calcula como $((\text{Valor Optimizado} - \text{Valor Inicial}) / \text{Valor Inicial}) * 100$ para métricas donde mayor es mejor, o $((\text{Valor Inicial} - \text{Valor Optimizado}) / \text{Valor Inicial}) * 100$ para métricas donde menor es mejor (ej. tasas de error). La fila sobre el 'Agente Consciente de Jerarquía' es ilustrativa y busca cuantificar el beneficio de implementar la lógica discutida en la Sección 5.3.

- *Instrucciones explícitas sobre fundamentación (Grounding) y atribución:* El prompt debe contener directivas claras e inequívocas que instruyan al LLM a basar su respuesta **predominante o exclusivamente** en la información contenida dentro de los documentos proporcionados y evitar activamente el uso de su conocimiento paramétrico interno o la realización de suposiciones no fundamentadas. Se deben incluir mandatos para **citar explícitamente las fuentes** de sus afirmaciones, idealmente vinculando cada proposición al fragmento o documento específico del contexto recuperado que la respalda. Esto no solo fomenta la fidelidad, sino que facilita la verificación por parte del usuario.
- *Guías restructuradas para el razonamiento (Chain-of-Thought y Similares):* Para consultas que requieren análisis o síntesis, en lugar de una simple extracción, el prompt puede instruir al LLM a seguir un proceso de razonamiento paso a paso (Wei et al. 2023). Por ejemplo, "Primero, identifica los hechos clave en los documentos proporcionados. Segundo, identifica las reglas legales aplicables mencionadas. Tercero, aplica estas reglas a los hechos. Cuarto, explica tu conclusión, citando los documentos relevantes para cada paso". Esta externalización del proceso de razonamiento no solo tiende a mejorar la calidad de la conclusión final, sino que también proporciona una traza de auditoría que puede ser revisada por un experto legal (Schwarcz et al., 2024).

Un enfoque robusto para guiar el razonamiento es la **descomposición explícita del problema** (explicit problem decomposition). En lugar de pedir al LLM que determine directamente si un texto viola una norma, la tarea se divide en los sub-componentes lógicos que un jurista analizaría. Por ejemplo, para determinar si un comentario constituye "incitación al odio" según el § 130 del Código Penal alemán, un sistema puede ser instruido para responder primero a dos preguntas separadas: (1) ¿El texto se dirige a un grupo protegido por la norma? y (2) ¿El texto realiza un acto prohibido por la norma (incitar, insultar, etc.)? (Ludwig et al., 2025). Solo si ambas respuestas son afirmativas, se concluye que la norma ha sido violada. Esta metodología no solo estructura el "pensamiento" del modelo, sino que hace su conclusión final mucho más transparente y verificable para el supervisor humano.

- *Manejo sofisticado de la incertidumbre, los conflictos y la información faltante:* El prompt debe guiar al LLM sobre cómo proceder cuando la información recuperada es incompleta, ambigua, presenta

contradicciones internas, o simplemente no contiene la respuesta a la consulta. En lugar de forzar una respuesta o recurrir a la fabricación, el modelo debe ser instruido para:

- Indicar explícitamente la incertidumbre (p. ej., "Basado en la información proporcionada, no es posible determinar con certeza...").
 - Presentar las diferentes perspectivas o la información conflictiva de manera objetiva, señalando las discrepancias.
 - Declarar que la información solicitada no se encuentra en los documentos recuperados.
- *Definición precisa de la persona y el formato de salida:* Especificar el rol que el LLM debe adoptar (p. ej., "Actúa como un asistente de investigación legal objetivo y neutral") y el formato exacto de la respuesta esperada (p. ej., resumen estructurado, lista de puntos clave con citas, borrador de cláusula contractual) es crucial para asegurar que la salida sea consistente, útil y profesional.
 - *Técnicas de prompting de seguimiento (Follow-up) o refinamiento iterativo:* Como demostró el estudio de Addleshaw Goddard (2024), un segundo prompt que desafíe o pida una revisión de la respuesta inicial del LLM (p. ej., "Por favor, revisa tu respuesta anterior cuidadosamente. ¿Estás seguro de que has incluido toda la información relevante de los documentos proporcionados sobre X? ¿Hay algún matiz que hayas omitido?") puede inducir al modelo a realizar un procesamiento más profundo del contexto y mejorar significativamente la calidad y completitud de la respuesta final. Este enfoque iterativo simula una conversación de refinamiento.
2. **Fine-Tuning del LLM generador con enfoque en la fidelidad legal y el razonamiento fundamentado:** Si bien la ingeniería de prompts es una herramienta poderosa y flexible, el fine-tuning del LLM generador en un corpus cuidadosamente seleccionado puede ofrecer mejoras más profundas y consistentes en su capacidad para adherirse al contexto y razonar de manera legalmente sólida.
- *Fine-tuning para la fidelidad al contexto (Contextual Adherence):* Entrenar al LLM en un conjunto de datos de alta calidad compuesto por tripletas de (consulta, contexto legal recuperado relevante, respuesta ideal que es estrictamente fiel al contexto y correctamente citada). Esto puede enseñar al modelo a priorizar la información contextual sobre su conocimiento paramétrico y a resistir la tentación de "desviarse" o alucinar (Tian, Mitchell, Yao, et al. 2023).
 - *Fine-tuning para tipos específicos de razonamiento jurídico fundamentado:* Desarrollar conjuntos de datos de entrenamiento que ejemplifiquen cómo realizar tipos específicos de tareas de razonamiento legal (p. ej., identificación del *holding* de un caso, aplicación de un test legal de múltiples factores, comparación de estatutos) *basándose explícitamente en un conjunto de documentos de entrada*.
 - *Aprendizaje por refuerzo con retroalimentación humana (RLHF) enfocado en la factualidad y el grounding:* Utilizar RLHF no solo para alinear el modelo con las preferencias generales de estilo o utilidad, sino específicamente para recompensar respuestas que demuestren alta factualidad, fundamentación precisa en las fuentes proporcionadas y razonamiento legal coherente.
3. **Integración con modelos de razonamiento especializados y arquitecturas avanzadas:** La emergencia de LLMs con arquitecturas explícitamente diseñadas para el razonamiento multi-paso, la planificación y la descomposición de problemas (como la familia de modelos "o" de OpenAI - OpenAI 2024) es particularmente relevante para RAG.
- *Planificación de la respuesta:* Estos modelos podrían, en teoría, planificar cómo utilizar la información recuperada de manera más estratégica, identificando qué fragmentos son más relevantes para qué partes de la consulta y cómo sintetizarlos de manera lógicamente coherente.
 - *Verificación interna de pasos de razonamiento:* Su capacidad para "reflexionar" sobre sus propios pasos de razonamiento intermedios podría permitirles detectar y corregir errores o inconsistencias antes de generar la respuesta final (Schwarcz et al., 2024). La integración de estos modelos de razonamiento como el componente generador en un sistema RAG es un área prometedora para futuras mejoras.
4. **Arquitecturas híbridas (Simbólico-Neuronales):** Aunque aún en desarrollo para aplicaciones legales a gran escala, la integración de la capacidad de los LLMs para procesar lenguaje natural con la precisión y verificabilidad de los sistemas de razonamiento simbólico (basados en lógicas formales, ontologías legales estructuradas o bases de reglas explícitas) ofrece una vía prometedora. El LLM podría usar el contexto recuperado para instanciar un modelo simbólico que luego realiza las inferencias lógicas de manera más controlada y explicable.

El objetivo final de estas estrategias de optimización de la generación no es solo producir respuestas que *parezcan* correctas, sino respuestas que sean *demostrablemente* correctas, fieles a las fuentes proporcionadas y útiles para el profesional legal que asume la responsabilidad final por su uso. La capacidad del LLM para explicar *cómo* llegó a una conclusión a partir del contexto recuperado es tan importante como la conclusión misma.

5.5 Verificación post-Hoc y calibración de confianza: la última línea de defensa contra las alucinaciones

Partiendo de la premisa de que es teóricamente imposible prevenir al cien por cien las alucinaciones en la fase de generación con la tecnología actual, la implementación de mecanismos robustos de verificación *después* de que el LLM ha producido una respuesta inicial se convierte en una capa de seguridad absolutamente crítica. Esta "última línea de defensa" no busca tanto evitar que el modelo alucine, sino detectar las alucinaciones cuando ocurren y proporcionar al usuario profesional señales claras sobre la fiabilidad de la información generada.

5.5.1 Verificación factual automatizada contra fuentes externas canónicas (Fact-Checking)

Una vez que el LLM ha generado una respuesta (que idealmente incluye citas preliminares), se pueden implementar módulos automatizados que:

- *Extraigan las afirmaciones factuales clave:* Identificar las proposiciones factuales y legales centrales en la respuesta del LLM.
- *Verifiquen contra bases de conocimiento de alta confianza:* Comparar estas afirmaciones con información contenida en bases de datos legales estructuradas y canónicas (p. ej., repositorios oficiales de legislación, bases de datos jurisprudenciales con metadatos de derogación, enciclopedias jurídicas verificadas).
- *Marquen las discrepancias:* Señalar explícitamente al usuario cualquier discrepancia encontrada, indicando si una afirmación no pudo ser verificada, contradice una fuente canónica, o se basa en una fuente citada que no la respalda (Peng et al. 2023; Chern et al. 2023).
- *Desafíos:* La cobertura de estas bases de conocimiento externas nunca será completa, y la verificación de afirmaciones legales complejas o interpretativas sigue siendo un desafío para los sistemas automatizados.

5.5.2 Verificación mediante reglas lógicas y heurísticas determinísticas

Antes de recurrir a modelos de IA secundarios, una capa de verificación basada en reglas puede detectar de forma eficiente y económica una clase significativa de errores. Esto incluye:

- *Validación sintáctica:* Usar expresiones regulares para verificar que las citas de sentencias o artículos legales siguen el formato canónico.
- *Chequeos de coherencia lógica simple:* Implementar reglas que marquen como sospechosa una afirmación donde un tribunal inferior revoca a uno superior, o donde una fecha de sentencia es posterior a la fecha de derogación de la ley que aplica.
- *Listas de control (Checklists):* Comparar entidades mencionadas (jueces, partes, leyes) contra bases de datos autorizadas para detectar invenciones flagrantes.

Estos métodos, aunque tradicionales, son altamente fiables para los errores que están diseñados para capturar y deben constituir una primera línea de defensa.

5.5.3 Modelos secundarios de detección de alucinaciones y autocrítica

Se está investigando el uso de modelos de IA, a menudo más pequeños y especializados, o incluso el propio LLM generador operando en un modo de "autoevaluación", para analizar la respuesta inicial en busca de indicios de alucinación.

- *Detección de inconsistencias internas:* Evaluar la coherencia lógica interna de la respuesta generada.
- *Medición de la entropía o incertidumbre de la generación:* Analizar las probabilidades asociadas a la secuencia de tokens generada; secuencias de baja probabilidad o alta entropía pueden ser más propensas a ser alucinaciones (Manakul, Liusie, and Gales 2023 - SelfCheckGPT).
- *Comparación con conocimiento paramétrico de alta confianza:* Si el LLM tiene "conocimiento" paramétrico sobre un tema con alta confianza (p. ej., principios legales muy básicos), puede usarlo para contrastar la respuesta generada a partir del contexto RAG.
- *Generación de críticas (CriticGPT):* OpenAI ha experimentado con modelos (como CriticGPT) entrenados para generar críticas de las respuestas de otros LLMs, ayudando a identificar errores o debilidades (Song et al., 2024 - RAG-HAT).

5.5.4 Calibración y comunicación efectiva de la confianza del modelo

Es fundamental que los LLMs no solo generen respuestas, sino que también comuniquen de manera fiable su propio nivel de "confianza" o incertidumbre sobre la corrección y fundamentación de dichas respuestas.

No obstante, la mera generación de una puntuación de confianza no es suficiente; la interpretación de dicha puntuación por parte del usuario profesional representa otro desafío significativo. Un '70% de confianza' expresado por un LLM puede no tener el mismo significado intuitivo o estadístico que un 70% de confianza en un contexto humano o en un sistema de diagnóstico tradicional. Por ello, junto con el desarrollo de métricas de confianza más fiables, es esencial investigar y establecer directrices claras sobre cómo los profesionales del derecho deben interpretar y actuar en función de estos indicadores de confianza algorítmica, especialmente cuando la calibración del modelo, como se evidencia en estudios como el de Dahl et al. (2024), sigue siendo imperfecta y puede llevar a una peligrosa sobreconfianza en las respuestas erróneas.

- *Desarrollo de métricas de confianza fiables:* Investigar y refinar técnicas para que los LLMs produzcan puntuaciones de confianza que se correlacionen bien con su precisión real en tareas legales específicas (Kadavath et al. 2022; Xiong et al. 2023). Esto sigue siendo un área de investigación activa y desafiante, como lo demuestran los problemas de calibración observados en Dahl et al. (2024).
- *Presentación transparente de la incertidumbre:* La interfaz de usuario debe comunicar claramente al profesional legal los niveles de confianza asociados a diferentes partes de la respuesta, o marcar explícitamente las afirmaciones sobre las cuales el modelo tiene baja confianza.
- *Umbrales de intervención:* Para aplicaciones de alto riesgo, se podrían establecer umbrales de confianza por debajo de los cuales una respuesta no se presenta al usuario o se marca inequívocamente como "requiere verificación humana intensiva".

La necesidad de una comunicación transparente sobre las capacidades y limitaciones de los sistemas de IA, incluyendo su nivel de confianza o incertidumbre, encuentra un eco normativo en la Ley de IA de la UE. El Artículo 52(1) de la EU-AIAct, por ejemplo, establece obligaciones de transparencia para ciertos sistemas de IA, incluyendo aquellos que generan contenido. Se exige que los usuarios sean informados de que están interactuando con un sistema de IA y, cuando un sistema de IA genera o manipula contenido de texto, audio o vídeo que se asemeje notablemente a contenido existente ('deep fakes'), se debe divulgar que el contenido ha sido generado o manipulado artificialmente. Si bien estas disposiciones no previenen directamente la generación de una alucinación fáctica en una respuesta legal, sí buscan fomentar una mayor conciencia y cautela por parte del usuario, permitiéndole ponderar la fiabilidad de la información recibida y estableciendo una base para la rendición de cuentas cuando la IA se presenta engañosamente como humana o su contenido como no artificial.

5.5.5 Generación de múltiples hipótesis y explicaciones contrastantes

En lugar de generar una única respuesta "definitiva", el LLM podría ser instruido para generar múltiples interpretaciones o argumentos posibles basados en el contexto recuperado, especialmente si este es ambiguo o presenta información conflictiva. Podría también generar explicaciones que contrasten los pros y los contras de diferentes enfoques legales, permitiendo al profesional humano sopesar las alternativas.

5.5.6 Facilitación de la verificación humana a través de citas precisas y rastreables

Una de las contribuciones más importantes de los sistemas RAG bien diseñados es su capacidad para mejorar la verificabilidad.

- *Citación a nivel de fragmento (Chunk-Level Citation):* El sistema no debe simplemente listar los documentos recuperados, sino que debe, en la medida de lo posible, vincular cada afirmación o conclusión específica en la respuesta generada al fragmento (o fragmentos) exacto del texto recuperado que la respalda.
- *Resaltado de pasajes relevantes:* La interfaz de usuario podría resaltar los pasajes específicos en los documentos fuente que fueron más influyentes para la generación de la respuesta, permitiendo al abogado ir directamente a la evidencia.
- *Transparencia sobre el proceso de recuperación:* Ofrecer al usuario visibilidad sobre qué documentos fueron recuperados (y quizás por qué, p. ej., mostrando puntuaciones de similitud) puede ayudarle a evaluar la calidad de la base informativa utilizada por el LLM.

5.5.7 Diseño de una capacidad de abstención inteligente: el principio de "silencio estratégico"

Más allá de la mera comunicación de un puntaje de confianza, una estrategia de mitigación avanzada consiste en diseñar la capacidad de abstención del sistema ("no lo sé") no como un error o una limitación accidental, sino como una función deliberada y estratégica. Un sistema que es capaz de identificar los límites de su conocimiento o de las fuentes proporcionadas inspira mayor confianza y es inherentemente más seguro que uno optimizado para generar una respuesta a toda costa. La implementación de este principio de "silencio estratégico" implica varios componentes clave:

- **Justificación específica para la abstención:** Cuando el sistema se abstiene, no debe ofrecer excusas genéricas. La respuesta debe ser un diagnóstico preciso de la limitación encontrada. Por ejemplo: "No es posible proporcionar una respuesta fundamentada, ya que no se han encontrado fuentes primarias en la base de conocimiento posteriores a 2023 sobre esta materia", o "Las fuentes recuperadas presentan datos conflictivos sobre el punto X y no permiten una síntesis concluyente". Esta transparencia educa al usuario y convierte una posible frustración en una interacción informativa.
- **Provisión de alternativas constructivas:** La abstención no debe ser un punto muerto. Un sistema robusto debe ofrecer al usuario vías de acción alternativas que aún aporten valor. Por ejemplo: "Aunque no puedo determinar la aplicabilidad directa, puedo proporcionar el marco legal general y una lista de verificación de los elementos que un profesional debería analizar", o "Puedo formular las preguntas específicas que debería dirigir a un asesor legal para resolver esta cuestión".
- **Comunicación visual y explícita de la incertidumbre:** En línea con la calibración de confianza, la interfaz debe comunicar proactivamente el nivel de fiabilidad de una respuesta. Un sistema de "semáforos" (por ejemplo, verde para respuestas con alto consenso y fuentes sólidas; ámbar para aquellas con lagunas de información o fuentes secundarias; rojo para las basadas en datos conflictivos o de alta incertidumbre) permite al usuario calibrar su propio nivel de escrutinio de forma inmediata.
- **Auditabilidad de la abstención:** Cada instancia de abstención debe generar un registro auditable (*log*). Este registro debe documentar el estado del sistema en ese momento: la consulta del usuario, las fuentes recuperadas (o la falta de ellas), los criterios que llevaron a la decisión de abstenerse y los umbrales de confianza predefinidos. Esta trazabilidad es fundamental para la mejora continua del sistema y para la rendición de cuentas.

En última instancia, el diseño de sistemas de IA legal debe redefinir sus incentivos: en lugar de premiar la verbosidad y la completitud a cualquier precio, se debe premiar la precisión y la cobertura responsable. Un modelo que sabe abstenerse de forma justificada y transparente no es un sistema menos capaz, sino uno que ha alcanzado un mayor grado de madurez y demuestra un profundo respeto por el usuario y por la criticidad del dominio en el que opera. Esta calibración conservadora es un pilar fundamental para construir una confianza sostenible a largo plazo en la IA jurídica.

5.5.8 Conclusión

Estas estrategias de verificación post-hoc y comunicación de confianza no eliminan la necesidad de las optimizaciones previas en datos, recuperación y generación, pero actúan como una red de seguridad crucial. Reconocen la falibilidad inherente de los LLMs y buscan empoderar al profesional legal con las herramientas y la información necesarias para usar la IA de manera más crítica, informada y, en última instancia, más segura. No obstante, es crucial reconocer el delicado equilibrio inherente al 'costo de la verificación'. Si bien estas capas de seguridad post-hoc son indispensables para la fiabilidad, su implementación extensiva, especialmente si involucra una intervención humana significativa para cada comprobación o una alta latencia por múltiples llamadas a modelos secundarios, podría llegar a mermar uno de los beneficios primarios que la IA promete: la eficiencia y la reducción de costos. Un sistema que requiera una verificación manual tan exhaustiva de cada una de sus salidas que anule por completo el ahorro de tiempo inicial, podría no ser viable en la práctica para muchas tareas. Por lo tanto, el desarrollo futuro debe buscar no solo la efectividad de estos mecanismos de verificación, sino también su eficiencia, posiblemente a través de una mayor automatización inteligente de la propia verificación o mediante sistemas de IA que aprendan a 'auto-corregirse' de manera más fiable con una mínima supervisión. Encontrar el punto óptimo donde la robustez de la verificación no sacrifique desproporcionadamente la eficiencia operativa es un desafío continuo en el diseño de sistemas de IA legal prácticos y confiables.

5.6 El rol irreductible y fortalecido de la supervisión humana experta

A pesar de la sofisticación creciente de las estrategias de optimización y mitigación de alucinaciones, desde la curación de datos hasta la verificación post-hoc, es imperativo concluir esta sección reafirmando un principio fundamental: en el estado actual y previsible de la inteligencia artificial, **la supervisión humana crítica, informada y experta**

no es meramente una opción deseable, sino un componente absolutamente irreductible e indispensable para la integración segura y ética de los LLMs en la práctica legal. Ninguna combinación de las técnicas algorítmicas discutidas puede, por sí sola, reemplazar la profundidad del juicio contextual, la responsabilidad ética y la comprensión matizada del profesional del derecho.

La necesidad de esta pericia humana no es una conjetura, sino una conclusión empírica. Incluso en estudios que implementan estrategias de condicionamiento altamente sofisticadas, proveyendo a los LLMs con definiciones legales y ejemplos de casos, se documenta una "brecha de rendimiento significativa" entre el mejor modelo de IA y los expertos legales humanos. El estudio de Ludwig et al. (2025) encontró que, si bien los modelos podían identificar razonablemente bien los grupos protegidos por la ley de discurso de odio, tenían serias dificultades para clasificar correctamente las conductas prohibidas, una tarea de juicio matizado donde los juristas humanos demostraron una fiabilidad muy superior. Esto subraya que la etapa final de evaluación y juicio cualitativo sigue siendo, por ahora, una capacidad exclusivamente humana.

Esto es particularmente cierto porque la detección de errores sutiles o la evaluación de la solidez de un argumento legal generado por IA a menudo depende intrínsecamente de la competencia y el juicio del profesional. Lo que para un lego o un abogado junior puede parecer una respuesta coherente y útil, para un experto podría revelar deficiencias argumentativas o una comprensión superficial de la doctrina aplicable. La 'verdad' o 'viabilidad' de una conclusión legal compleja no siempre es autoevidente y requiere un escrutinio informado.

El deber de competencia en la era de la IA, por tanto, exige que el profesional comprenda que no está interactuando con un "oráculo de conocimiento", sino con un sistema estadístico optimizado para la plausibilidad, cuyo diseño fundamental lo incentiva a generar respuestas seguras incluso cuando su base de conocimiento es incierta (Kalai et al., 2025). Reconocer esta característica de diseño es la base del escepticismo profesional necesario para una supervisión efectiva.

Lejos de volver obsoleto al abogado, la emergencia de LLMs propensos a alucinaciones, incluso aquellos aumentados por RAG, **refuerza y redefine el valor de la pericia humana**. El rol del abogado evoluciona de ser un mero recuperador de información o redactor de documentos (tareas que la IA puede asistir o incluso automatizar parcialmente) a convertirse en:

- 1. Supervisor crítico de la IA:** El abogado debe actuar como un "controlador de calidad" inteligente y escéptico de las salidas generadas por la IA. Esto implica no solo verificar la corrección factual y la validez legal de la información, sino también evaluar su relevancia contextual, su adecuación estratégica a los objetivos del cliente, y sus implicaciones éticas.
Este rol de "controlador de calidad" va más allá de la simple verificación factual. El abogado debe actuar como un filtro metacognitivo, siendo consciente de que la forma en que el LLM presenta la información puede inducir sesgos en su propio proceso de razonamiento. La investigación sobre sesgos cognitivos inducidos por LLMs demuestra que los resultados de estos modelos pueden alterar el encuadre o el énfasis de la información, llevando a los humanos a tomar decisiones diferentes a las que tomarían con la información original (Alessa et al., 2025). Por lo tanto, la supervisión crítica implica un acto de auto-reflexión: el abogado no solo debe preguntar "¿es correcta esta información?", sino también "¿está esta presentación de la información influyendo indebidamente en mi juicio?".
- 2. Curador y guía del conocimiento de la IA:** En el contexto de sistemas RAG personalizables o fine-tuneables, los abogados expertos pueden desempeñar un papel crucial en la curación de las bases de conocimiento, en el diseño de prompts efectivos y en la provisión de retroalimentación para mejorar el rendimiento del modelo en tareas legales específicas.
- 3. Intérprete y comunicador del output de la IA:** Incluso si una IA genera un análisis legal técnicamente correcto, a menudo se requerirá que un abogado humano lo traduzca a un lenguaje comprensible para el cliente, lo contextualice dentro de la situación particular del cliente y lo integre en una estrategia legal más amplia.
- 4. Garante del juicio ético y estratégico:** La IA puede procesar información y generar opciones, pero la toma de decisiones finales que implican consideraciones éticas complejas, el ejercicio del juicio profesional sobre cursos de acción alternativos, la gestión de la relación con el cliente y la asunción de la responsabilidad profesional última, permanecen firmemente en el dominio humano.
- 5. Navegador de la incertidumbre y la ambigüedad legal:** Como se ha discutido, el derecho está lleno de áreas grises, conflictos normativos y situaciones donde no existe una única "respuesta correcta". La capacidad de un abogado para navegar esta incertidumbre, ponderar riesgos y beneficios, y aconsejar al cliente en consecuencia, es una habilidad que la IA actual no posee.

En este nuevo paradigma, la eficiencia prometida por la IA solo se materializa si va acompañada de una inversión proporcional en la **capacitación de los abogados para interactuar críticamente con estas herramientas**. Esto incluye desarrollar habilidades en:

- *Interacción semántica y profesional, más allá de la Ingeniería de Prompts:* Si bien actualmente la calidad de la respuesta de una IA a menudo depende de la 'ingeniería de prompts', es fundamental reconocer que este paradigma es una solución transitoria y un defecto de diseño, no un objetivo final. La responsabilidad de la complejidad técnica no debe recaer en el profesional del derecho, sino en el desarrollador de la herramienta LegalTech.

La verdadera innovación reside en desarrollar soluciones que abstraigan esta complejidad, permitiendo al jurista interactuar en su propio lenguaje —natural y técnico— y asumiendo la herramienta la carga de traducir esa intención a las instrucciones algorítmicas que el modelo necesita. Exigir que un cirujano aprenda a programar su bisturí es un fracaso del diseño; del mismo modo, la tecnología debe ser un bisturí que se adapta a la mano del abogado. Este enfoque lo libera para que se centre en lo que ninguna máquina puede hacer: aplicar el criterio, la estrategia y la ética.

- *Técnicas de verificación rigurosa:* Conocer las fuentes de autoridad primaria y secundaria, y ser capaz de contrastar eficientemente las salidas de la IA con ellas.

En última instancia, este nuevo paradigma refuerza una máxima que debe guiar el futuro de la LegalTech: el abogado no debe convertirse en un 'prompt engineer'. La responsabilidad de la complejidad técnica recae en los desarrolladores de la herramienta, no en el usuario final. Exigir que los profesionales del derecho aprendan complejas técnicas de prompting para obtener resultados fiables es un fracaso del diseño y una inversión de roles inaceptable. La tecnología debe ser un bisturí que se adapta a la mano del cirujano, no una máquina que exige que el cirujano aprenda su lenguaje arcano. Por ello, el futuro de la IA legal fiable reside en sistemas que permitan una interacción en lenguaje natural y que asuman la carga de la interpretación técnica, liberando al abogado para que se centre en lo que ninguna máquina puede hacer: aplicar el criterio, la estrategia y la ética.

- *Comprensión de las limitaciones de la IA:* Ser consciente de los tipos de errores y sesgos a los que la IA es propensa (incluyendo las alucinaciones) y saber cuándo no confiar en sus resultados.
- *Integración ética de la IA en la práctica:* Comprender las implicaciones deontológicas del uso de la IA y cómo cumplir con los deberes profesionales en un entorno tecnológicamente aumentado.

Este principio de la indispensabilidad de la supervisión humana no se articula únicamente como una conclusión derivada de las limitaciones técnicas intrínsecas de la IA actual, sino que está siendo progresivamente consagrado como un requisito fundamental en los marcos regulatorios más avanzados. La Ley de IA de la UE, en su Artículo 14, es explícita al exigir que los sistemas de IA de alto riesgo estén diseñados para ser 'efectivamente supervisados por personas'.

Un ejemplo paradigmático de cómo estos principios se están materializando a nivel nacional se encuentra en España. En junio de 2024, el **Comité Técnico Estatal de la Administración Judicial Electrónica (CTEAJE)** (el órgano gubernamental de alto nivel responsable de la modernización tecnológica y la estrategia digital del sistema judicial español) publicó su "**Política de uso de la Inteligencia Artificial en la Administración de Justicia**". Este documento, de obligado cumplimiento para el personal de la administración de justicia, no es una mera recomendación, sino un marco normativo que establece directrices inequívocas:

- **Principio de No Sustitución:** La política establece de forma tajante que "**la IA nunca debe reemplazar la toma de decisiones humanas en cuestiones cruciales**" y que "la responsabilidad final de tomar decisiones legales debe recaer en jueces y magistrados" (Principio 1.4.1).
- **Mandato de Revisión Humana Universal:** Más allá de los principios, se impone como norma de uso que "**la revisión humana de todo lo generado [por IA] siempre que afecte de manera directa o indirecta a los derechos de las personas usuarias**" es obligatoria (Norma 1.5.1). Esto convierte la supervisión en un requisito procesal ineludible, no en una opción.
- **Reconocimiento explícito de los riesgos:** La guía del CTEAJE define explícitamente el fenómeno de las "**alucinaciones**" y reconoce el peligro del "**sesgo de automatización endémico**", por el cual los humanos tienden a confiar ciegamente en las sugerencias de los sistemas. Este reconocimiento oficial subraya la necesidad de un escepticismo informado, pilar fundamental de la supervisión experta.

La existencia de una guía tan detallada y vinculante por parte de un órgano como el CTEAJE demuestra que el rol irreductible del profesional humano ha trascendido el debate académico para convertirse en un pilar de la política pública y de la gobernanza de la IA en el ámbito legal.

En última instancia, la fiabilidad en la era de la IA legal no residirá exclusivamente en la perfección de los algoritmos, sino en la **simbiosis efectiva entre la capacidad de procesamiento de la IA y la sabiduría, el juicio crítico y la responsabilidad ética del profesional humano**.

Este principio de la indispensabilidad de la supervisión humana no se articula únicamente como una conclusión derivada de las limitaciones técnicas intrínsecas de la IA actual, sino que está siendo progresivamente consagrado como un requisito legal fundamental en los marcos regulatorios más avanzados. La Ley de IA de la UE, en su Artículo 14, es explícita al exigir que los sistemas de IA de alto riesgo estén diseñados y desarrollados de tal manera que puedan ser 'efectivamente supervisados por personas durante el período en que el sistema de IA está en uso'. Esta supervisión debe permitir a los humanos comprender las capacidades y limitaciones del sistema, permanecer conscientes de la posible tendencia a la automatización o al sesgo de confirmación, interpretar correctamente la salida del sistema, y tener la autoridad y competencia para decidir no utilizar dicha salida, anularla o intervenir si el sistema presenta resultados anómalos, imprevistos o potencialmente perjudiciales, como es el caso de las alucinaciones que comprometen la validez legal.

Las estrategias de optimización y mitigación son herramientas esenciales en este proceso, pero su efectividad final depende de que sean implementadas y supervisadas por juristas bien formados, conscientes de los riesgos y comprometidos con los más altos estándares de la profesión. Lejos de ser una amenaza existencial, la IA alucinante puede, paradójicamente, subrayar el valor perdurable e insustituible de la inteligencia humana experta en el corazón del derecho.

6 La realidad de las alucinaciones en la práctica: estudios de caso detallados y lecciones aprendidas de incidentes judiciales

Si bien el análisis teórico y la evaluación empírica en entornos controlados son fundamentales para comprender la naturaleza y la prevalencia de las alucinaciones en los Grandes Modelos de Lenguaje (LLMs) aplicados al derecho, es en la arena de la práctica jurídica real donde las consecuencias de estos errores algorítmicos se manifiestan con una crudeza incontestable y un impacto tangible. Los incidentes donde la información generada por IA, incorrecta o completamente fabricada, ha sido introducida en procedimientos judiciales no son meras anécdotas o curiosidades tecnológicas; representan fallos sistémicos con el potencial de socavar la administración de justicia, erosionar la confianza pública y acarrear graves sanciones profesionales para los letrados implicados. Esta sección se adentra en el análisis detallado de varios estudios de caso prominentes y bien documentados, extrayendo de ellos lecciones cruciales sobre los puntos de fallo específicos en la interacción humano-IA, las deficiencias en los procesos de verificación y las consecuencias directas de confiar acríticamente en estas poderosas pero falibles herramientas. Estos casos sirven como advertencias potentes y como catalizadores para una reflexión más profunda sobre las salvaguardas necesarias en la integración de la IA en la práctica legal.

La proliferación de estos incidentes ha alcanzado un punto crítico, motivando la creación de recursos dedicados para su seguimiento. Un ejemplo notable es la base de datos en línea "AI Hallucination Cases Database", un proyecto que busca compilar de manera exhaustiva todas las decisiones judiciales donde el contenido alucinado por una IA ha sido un factor relevante. Este tipo de repositorios se está convirtiendo en una herramienta vital para juristas, académicos y reguladores, al permitir un análisis sistemático de la naturaleza y frecuencia de estos fallos en la práctica real.

6.1 Caso de estudio: el paradigmático *Mata v. Avianca, Inc.* y la fabricación de jurisprudencia

El caso *Robert Mata v. Avianca, Inc.*, No. 22-cv-1461 (PKC) (S.D.N.Y. 2023), se ha convertido rápidamente en el referente obligado al discutir los peligros de las alucinaciones de la IA en el litigio. En este asunto, los abogados del demandante, buscando oponerse a una moción de desestimación, presentaron un escrito judicial que citaba múltiples decisiones judiciales supuestamente favorables a su posición. Sin embargo, tras una investigación por parte de la defensa y del propio tribunal, se descubrió que al menos seis de los casos citados eran completamente inexistentes, fabricaciones generadas por ChatGPT, la herramienta de IA que uno de los abogados había utilizado para la investigación legal (Weiser, 2023; Dahl et al., 2024).

El Juez P. Kevin Castel, al imponer sanciones a los abogados implicados (incluyendo una multa económica y la obligación de notificar a los jueces cuyos nombres fueron falsamente asociados a las opiniones inventadas), emitió una orden detallada que disecciona los múltiples fallos en el proceso. El abogado que utilizó ChatGPT, Steven A. Schwartz, admitió no ser un experto en investigación legal federal y haber utilizado la herramienta como un "super motor de búsqueda", confiando en sus respuestas e incluso preguntando a ChatGPT si los casos que proporcionaba eran reales, a lo que el chatbot respondió afirmativamente (Orden de Sanciones en *Mata v. Avianca, Inc.*, 22 de junio de 2023).

Lecciones aprendidas de *Mata v. Avianca*:

Table 5: Análisis comparativo de incidentes judiciales destacados por alucinaciones de IA

Caso / Jurisdicción	Naturaleza de la Alucinación	Lecciones Clave / Consecuencias
<i>Mata v. Avianca, Inc.</i> (S.D.N.Y. 2023, EE.UU.)*	Fabricación completa de múltiples casos judiciales (citaciones, holdings). Respuesta afirmativa del chatbot sobre la existencia de los casos.	Deber ineludible de verificación independiente. Naturaleza engañosa de alucinaciones plausibles. Responsabilidad profesional individual. Sanciones económicas y reputacionales.
<i>Thackston v. Driscoll</i> (W.D. Texas, 2025, EE.UU.)*	Acumulación de múltiples tipos de error: Fabricación de jurisprudencia y doctrina, citas falsas en casos reales, tergiversación de holdings (<i>misgrounding</i>) y uso de jurisprudencia revocada.	Ilustra la "cascada" de errores que un solo uso negligente de la IA puede producir. Refuerza el deber de verificación más allá de la simple existencia del caso. Subraya las graves consecuencias profesionales (sanciones Regla 11).
Caso en Australia (Familia, 2024, citado por Lantyer)	Citaciones falsas en un caso de familia.	Universalidad del riesgo de alucinación. Aplicación de deberes deontológicos (diligencia, competencia). Respuesta de órganos disciplinarios.
Caso en Brasil (Apelación, 2025, citado por Lantyer)	Jurisprudencia falsa generada por IA en una apelación.	Riesgos en todas las instancias judiciales. Importancia de la formación continua en IA para abogados y jueces. Multa impuesta.
Caso Tribunal Constitucional (España, 2024)*	Invencción completa de 19 citas de doctrina judicial, presentadas como literales en un recurso de amparo.	La responsabilidad del letrado es absoluta e independiente de la herramienta causante del error. El uso negligente de IA constituye una falta al deber de respeto al tribunal.
Estudio Magesh et al. (2024) (Herramientas RAG)	Principalmente <i>misgrounding</i> (citar fuente real pero tergiversar contenido), errores de razonamiento, supresión de citas.	RAG mitiga pero no elimina alucinaciones. Errores sutiles pueden ser más insidiosos que la fabricación obvia. Necesidad de verificación profunda de la fuente.

Nota: Aquellos casos marcados con un asterisco () se comentan específicamente en secciones posteriores.*

- 1. Verificación independiente como deber ineludible:** La lección más obvia y contundente es la absoluta necesidad de que los abogados verifiquen de forma independiente y rigurosa cada fuente y cada proposición legal generada por una IA antes de incorporarla a un documento judicial. La simple pregunta a la IA sobre la veracidad de su propia salida es manifiestamente insuficiente y denota una falta de comprensión fundamental sobre cómo operan estos modelos.
- 2. La naturaleza engañosa de las alucinaciones:** Las citaciones fabricadas por ChatGPT en el caso *Mata* eran altamente plausibles, con nombres de partes, números de volumen y página, y resúmenes de *holdings* que imitaban el formato y estilo de las opiniones judiciales reales. Esta plausibilidad hace que las alucinaciones sean particularmente insidiosas y difíciles de detectar sin una verificación cruzada con bases de datos legales canónicas.
- 3. Responsabilidad profesional individual:** El caso subraya que la responsabilidad final por el contenido de los escritos presentados ante el tribunal recae inequívocamente en el abogado firmante, independientemente de las herramientas utilizadas en su preparación. El uso de IA no diluye ni transfiere esta responsabilidad.
- 4. Desconocimiento de las limitaciones de la IA:** La admisión del abogado Schwartz de que "no era consciente de la posibilidad de que [el] contenido [de ChatGPT] pudiera ser falso" revela una brecha significativa en la alfabetización sobre IA dentro de la profesión legal. Comprender las limitaciones inherentes de los LLMs, incluyendo su propensión a "alucinar" o "confabular", es un componente esencial de la competencia profesional en la era digital.
- 5. Impacto en la Integridad del Sistema Judicial:** La introducción de jurisprudencia ficticia en un procedimiento judicial no solo perjudica al cliente y expone al abogado a sanciones, sino que también "promueve el cinismo hacia la profesión legal y el sistema de justicia estadounidense" y constituye un abuso del proceso judicial (Orden de Sanciones en *Mata v. Avianca, Inc.*).

6.2 Caso de estudio: la sanción del Tribunal Constitucional español y la responsabilidad indelegable del letrado

6.2.1 Resumen fáctico del incidente

En septiembre de 2024, la Sala Primera del Tribunal Constitucional de España marcó un precedente fundamental al sancionar por unanimidad a un abogado por faltar al debido respeto al tribunal (Nota Informativa 90/2024). La falta consistió en la inclusión, en una demanda de amparo, de 19 citas supuestamente literales de sentencias del propio Tribunal que resultaron ser completamente inexistentes. Estas citas se presentaban entrecomilladas, atribuyendo a los magistrados una doctrina constitucional que "carecía de todo anclaje en la realidad". La sanción impuesta fue un "apercibimiento", la menor posible, pero se ordenó dar traslado al Colegio de Abogados de Barcelona para los procedimientos disciplinarios correspondientes.

6.2.2 Análisis del caso a través del marco del informe

Este caso sirve como una ilustración perfecta de los conceptos analizados en este informe:

- **Naturaleza de la Alucinación** (Aplicando la Taxonomía de la Sección 2.1): El error cometido encaja directamente en la categoría de "Fabricación de autoridad". No se trataba de una tergiversación sutil (misgrounding), sino de la invención completa de doctrina judicial. La presentación de los pasajes entrecomillados agrava la falta, ya que no se presenta como una interpretación, sino como una cita literal y verificable, lo que constituye una forma particularmente grave de información incorrecta según la tipología de Magesh et al. (2025).
- **La responsabilidad profesional por encima de la herramienta** (aplicando la Sección 8.1): El punto más crucial del Acuerdo del Tribunal es su razonamiento sobre la responsabilidad del abogado. El letrado alegó en su defensa una "desconfiguración de una base de datos". El Tribunal descartó este argumento de forma tajante, estableciendo un principio de responsabilidad absoluta que es independiente de la causa del error. En sus propias palabras, "fuera cual fuese la causa de la inclusión de citas irreales (uso de la inteligencia artificial, entrecomillado de argumentos propios, etcétera), el letrado es siempre responsable de revisar exhaustivamente todo el contenido" (Nota Informativa 90/2024). Esta afirmación es la manifestación práctica más clara del deber de diligencia y competencia profesional en la era de la IA, subrayando que la supervisión humana no es una opción, sino una obligación indelegable.
- **Impacto en la integridad del Sistema Judicial** (aplicando la Sección 2.3): El Tribunal no lo consideró un simple error procesal, sino una falta de respeto que mostraba un "claro desprecio de la función jurisdiccional" de los magistrados. La conducta, según el Acuerdo, perturbó el trabajo del Tribunal no por la necesidad de verificar las citas —algo que se hace siempre— sino por "tener que enjuiciar las consecuencias de tal injustificada irregularidad". Esto demuestra que la introducción de información falsa no solo contamina el debate jurídico, sino que socava la confianza y el respeto mutuo que deben regir la relación entre los abogados y la judicatura, erosionando los cimientos del sistema.

6.2.3 Lección humana: la delegación de la responsabilidad crítica

Más allá del análisis técnico-jurídico, el caso del Tribunal Constitucional español, al igual que *Mata v. Avianca*, es un síntoma de una peligrosa tendencia cultural: la delegación del pensamiento crítico. El letrado no falló por usar una herramienta defectuosa; falló porque abdicó de su responsabilidad fundamental de verificar, dudar y pensar. Trató a la IA como un oráculo en lugar de como un asistente. En cualquier profesión, pero especialmente en el derecho, el valor no reside en la capacidad de generar una respuesta, sino en la capacidad de defenderla. Cuando un profesional simplemente copia y pega un resultado que no comprende, no está utilizando la tecnología; está siendo utilizado por ella. Estos incidentes no deberían generar miedo a la IA, sino un profundo respeto por el rol insustituible del juicio humano. La tecnología no nos exime de nuestra obligación de ser excelentes; de hecho, nos la exige con más fuerza que nunca."

6.2.4 Lecciones clave y comparativa con *Mata v. Avianca*

Este caso, aunque similar en su origen a *Mata v. Avianca*, ofrece lecciones complementarias y de mayor calado ético:

- **Universalidad del deber de verificación:** Confirma que la obligación de verificar cada dato presentado ante un tribunal es un principio universal del ejercicio de la abogacía, aplicable con la misma fuerza en sistemas de derecho civil (España) como de common law (EE. UU.).
- **Irrelevancia de la causa del error:** Mientras que en *Mata* el debate se centró en el mal uso de una herramienta específica (ChatGPT), el Tribunal Constitucional español eleva el principio: la responsabilidad del abogado es absoluta, sin importar si el error fue causado por una IA, un software defectuoso o un descuido humano. La herramienta es irrelevante; la responsabilidad, total.

- De la sanción procesal a la falta ética: El caso español encuadra el problema no solo como una negligencia que merece una sanción procesal, sino como una falta al deber de respeto, un pilar de la ética profesional. Es un "desprecio" a la función judicial, lo que le confiere una gravedad deontológica superior.

6.3 Caso de estudio: *Thackston v. Driscoll* y la cascada de errores algorítmicos

El caso *Thackston v. Driscoll*, resuelto por un Juez Magistrado en el Distrito Oeste de Texas el 28 de agosto de 2025, se erige como un ejemplo alarmantemente completo de los peligros derivados de un uso acrítico y negligente de la IA generativa en la práctica legal. A diferencia de *Mata v. Avianca*, que se centró principalmente en la fabricación de casos, *Thackston* ilustra una "cascada" de errores que abarca casi toda la gama de la taxonomía de alucinaciones legales.

6.3.1 Resumen fáctico del incidente

En el marco de una demanda por discriminación laboral contra el Ejército de EE.UU., el abogado del demandante presentó un escrito de réplica plagado de información legal defectuosa. El tribunal, en su "Informe y Recomendación", diseccionó meticulosamente los errores, que incluían:

- **Fabricación de autoridad:** Citas a casos completamente inexistentes (ej. una supuesta opinión del Noveno Circuito en *United States v. City of Los Angeles* y un caso del propio tribunal en *EEOC v. Exxon Mobil Corp.* que no existía).
- **Citas falsas y tergiversación (*Misgrounding*):** El escrito atribuía citas textuales inventadas a casos reales y conocidos (ej. a *Palmer v. Shultz* y *Armstrong v. Turner Industries*). Además, tergiversaba gravemente los *holdings* de otros casos reales, citándolos en apoyo de proposiciones legales que no sostenían.
- **Error temporal:** Se citó el famoso caso *Chevron*, un pilar del derecho administrativo, sin reconocer que había sido revocado explícitamente por el Tribunal Supremo, un error fáctico y estratégico de primer orden.

El Juez Magistrado no solo identificó los errores, sino que también sospechó explícitamente del uso de IA por el "lenguaje repetitivo y redundante" y la naturaleza de las invenciones. Concluyó que el abogado violó la Regla Federal 11(b) al no realizar una "investigación razonable" y recomendó al Tribunal de Distrito la imposición de sanciones, sugiriendo una multa y la asistencia obligatoria a un curso de formación sobre IA generativa.

6.3.2 Análisis del caso a través del marco del informe

Este caso es un microcosmos perfecto de los riesgos sistémicos discutidos en este informe.

- **Una taxonomía completa en un solo documento** (aplicando la Sección 2.2): *Thackston* es una clase magistral sobre los diferentes tipos de alucinaciones. Demuestra que un abogado que confía ciegamente en una IA no comete un solo tipo de error, sino que se expone a un fallo sistémico. La combinación de **fabricación de autoridad, fundamentación errónea (*misgrounding*) y error temporal** en un mismo escrito muestra la incapacidad del modelo (y del abogado) para distinguir entre lo real, lo tergiversado y lo obsoleto. El *misgrounding* es particularmente insidioso aquí, ya que el abogado podría haber verificado la existencia del caso y haberse detenido ahí, cayendo en una falsa sensación de seguridad.
- **La abdicación del juicio profesional y la "alucinación del usuario"** (aplicando las Secciones 2.3 y 5.6): El tribunal es inequívoco: la culpa no es de la "máquina", sino del profesional que abdicó de su deber fundamental. La recomendación del Juez Magistrado se centra en la violación de la Regla 11, que exige una investigación razonable *antes* de presentar cualquier documento. Este es el ejemplo paradigmático de la "alucinación del usuario": la creencia errónea de que la herramienta puede sustituir la diligencia, el escepticismo y el juicio profesional. Refuerza el principio central de este informe: el rol de la supervisión humana no es solo una buena práctica, es una obligación legal y ética irreductible.
- **De la teoría a las consecuencias reales** (aplicando la Sección 8): El caso *Thackston* materializa las implicaciones deontológicas y regulatorias. La recomendación de sanciones monetarias y formación obligatoria no es una reprimenda abstracta, sino una consecuencia profesional y económica directa. Sirve como una advertencia potente, alineada con las lecciones de *Mata* y del caso del Tribunal Constitucional español: los tribunales no dudarán en utilizar sus facultades sancionadoras para proteger la integridad del proceso judicial frente a la introducción de información falsa, sin importar la tecnología utilizada para generarla.

6.4 Incidentes más allá de la fabricación flagrante: errores sutiles y *Misgrounding*

El estudio de Magesh et al. (2025) documenta que las alucinaciones más frecuentes y peligrosas en sistemas RAG no son las fabricaciones completas, sino errores de razonamiento legal más sutiles, que ellos denominan "insidiosos". Identifican tres categorías principales de fallo:

- **Incomprensión de los holdings (Misunderstanding Holdings):** Las herramientas a menudo resumen una sentencia afirmando lo contrario de lo que el tribunal realmente decidió, confundiendo el holding (la decisión central) con los dicta (comentarios secundarios).
- **Confusión entre actores legales (Distinguishing Between Legal Actors):** Las IAs atribuyen erróneamente los argumentos de un litigante al tribunal, presentando la postura de una de las partes como si fuera la decisión final del juez.
- **Falta de respeto a la jerarquía de autoridad (Respecting the Order of Authority):** Los modelos demuestran una incapacidad para comprender la jerarquía judicial, por ejemplo, afirmando que un tribunal inferior revocó una decisión de un tribunal superior, lo cual es legalmente imposible.

Estos errores de *misgrounding* son particularmente peligrosos porque la presencia de una cita real crea una falsa sensación de autoridad y fiabilidad, lo que puede llevar al abogado a confiar indebidamente en la proposición sin realizar la necesaria lectura crítica de la fuente (Magesh et al., 2025).

Si bien la fabricación completa de casos como en *Mata v. Avianca* es la forma más espectacular de alucinación, los estudios empíricos sobre herramientas legales comerciales basadas en RAG, como el de Magesh et al. (2024), revelan que formas más sutiles pero igualmente problemáticas de error son aún más frecuentes. Estos incidentes, aunque no siempre conllevan sanciones tan publicitadas, pueden tener un impacto significativo en la calidad del trabajo legal y en la toma de decisiones.

Un ejemplo recurrente documentado por Magesh et al. (2024) es el *misgrounding*, donde la herramienta de IA cita un caso o estatuto *real y existente*, pero la proposición legal que atribuye a esa fuente es incorrecta, tergiversada o simplemente no está contenida en el texto original de la autoridad citada. En una de las instancias analizadas, Westlaw AI-Assisted Research afirmó incorrectamente el *holding* de una decisión de la Corte Suprema de EE. UU., atribuyéndole una conclusión opuesta a la que realmente alcanzó el tribunal. En otro ejemplo, Lexis+ AI describió un caso (Arturo D.) como autoridad vigente y lo utilizó para respaldar una proposición, cuando en realidad el caso citado (Lopez) había *revocado* a Arturo D. en el punto relevante.

Lecciones aprendidas de errores de *Misgrounding* y similares:

1. **La verificación no puede limitarse a la existencia de la cita:** A diferencia de las fabricaciones obvias, el *misgrounding* requiere un nivel de verificación más profundo. No basta con confirmar que el caso o estatuto citado existe; el abogado debe leer y comprender la fuente original para asegurar que realmente respalda la afirmación hecha por la IA.
2. **Fragilidad de la comprensión contextual de la IA-RAG:** Estos errores sugieren que, incluso cuando se les proporciona el contexto recuperado, los LLMs pueden tener dificultades para interpretar correctamente los matices del lenguaje legal, distinguir el *holding* de los *dicta*, o comprender las relaciones jerárquicas y temporales entre precedentes (p. ej., el efecto de una revocación).
3. **El riesgo de la "falsa fundamentación":** El *misgrounding* es particularmente peligroso porque la presencia de una cita real puede crear una falsa sensación de autoridad y fiabilidad, llevando al abogado a confiar indebidamente en la proposición sin realizar la necesaria lectura crítica de la fuente.
4. **Necesidad de optimización profunda en RAG:** Estos incidentes refuerzan la conclusión del informe de Addleshaw Goddard (2024) de que la efectividad de RAG depende de una optimización meticulosa de cada componente, incluyendo no solo la recuperación sino también la forma en que el LLM generador es instruido (prompting) para interactuar con el contexto recuperado y razonar sobre él.

6.5 Implicaciones globales y la necesidad de adaptación continua

Aunque muchos de los casos más notorios han surgido en EE. UU., el problema de las alucinaciones de la IA y la necesidad de una verificación diligente por parte de los abogados es una preocupación global. Incidentes similares han comenzado a documentarse en otras jurisdicciones, incluyendo Australia (donde un abogado fue remitido a un órgano disciplinario por usar IA que generó citas falsas en un caso de familia - The Guardian, 2024, citado en Lantyer, 2024) y Brasil (donde un abogado fue multado por usar jurisprudencia falsa generada por IA en una apelación - Migalhas, 2025, citado en Lantyer, 2024).

Más allá de los incidentes judiciales anecdóticos, la investigación académica sistemática confirma este riesgo a nivel global. Un estudio empírico enfocado en una jurisdicción fuera del ámbito europeo comparó el rendimiento de múltiples LLMs con un abogado humano, concluyendo que la tarea de investigación legal era la de peor rendimiento, con una tendencia constante de los modelos de IA a inventar jurisprudencia (Hemrajani, 2025).

Lecciones aprendidas de la perspectiva global:

1. **Universalidad del riesgo:** La propensión de los LLMs a alucinar no está limitada por fronteras geográficas o sistemas jurídicos. Es una característica inherente de la tecnología actual que afecta a todos los profesionales que la utilizan.
2. **Adaptación de los deberes deontológicos:** Si bien los detalles de los códigos deontológicos varían entre jurisdicciones, los principios fundamentales de competencia, diligencia, lealtad al cliente y franqueza ante el tribunal son ampliamente compartidos. La profesión legal en cada país deberá interpretar y aplicar estos principios al nuevo contexto de la IA.
3. **Respuesta de órganos disciplinarios y judiciales:** La forma en que los tribunales y los órganos disciplinarios de diferentes países respondan a estos incidentes establecerá precedentes importantes y modelará las expectativas sobre el uso responsable de la IA por parte de los abogados. Las sanciones impuestas en casos como *Mata* sirven como una señal clara de la seriedad con la que se toman estos fallos.
4. **Importancia de la formación y la alfabetización en IA:** A nivel global, existe una necesidad urgente de mejorar la formación de los abogados sobre las capacidades y limitaciones de la IA, incluyendo la concienciación sobre el riesgo de alucinaciones y el desarrollo de habilidades de verificación crítica.

En conclusión, los estudios de caso de incidentes reales donde las alucinaciones de la IA han impactado procedimientos judiciales ofrecen lecciones contundentes e ineludibles. Subrayan que la integración de la IA en la práctica legal no es un proceso exento de riesgos y que la tecnología, en su estado actual, no puede sustituir el juicio crítico, la diligencia investigadora y la responsabilidad ética del profesional humano. Estos casos no deben interpretarse como una condena de la IA *per se*, sino como un llamado urgente a la cautela, a la verificación rigurosa y al desarrollo de prácticas profesionales y salvaguardas tecnológicas que permitan aprovechar el potencial de la IA minimizando sus peligros inherentes. La "realidad cruda" de estos incidentes debe servir como un motor para la mejora continua, tanto de la tecnología como de la forma en que la profesión legal interactúa con ella.

7 El futuro de la Inteligencia Artificial legal fiable: hacia modelos explicables, auditables y responsables por diseño

El panorama actual de la inteligencia artificial (IA) aplicada al derecho, aunque rebosante de potencial transformador, se encuentra marcado por el desafío persistente de las alucinaciones y las limitaciones inherentes a la fiabilidad de los Grandes Modelos de Lenguaje (LLMs) y las arquitecturas de Generación Aumentada por Recuperación (RAG) convencionales. Las secciones precedentes han diseccionado la naturaleza de estos errores, sus causas raíz y las estrategias de mitigación disponibles. Sin embargo, una visión a largo plazo exige ir más allá de la mera contención de los problemas actuales y proyectar un futuro donde la IA legal no solo sea más potente y eficiente, sino fundamentalmente más **fiable, transparente y alineada con los principios éticos y las exigencias de rendición de cuentas** inherentes al sistema de justicia. Este futuro no dependerá de un único avance disruptivo, sino de la convergencia de múltiples líneas de investigación y desarrollo enfocadas en la creación de modelos inherentemente más explicables (**XAI**, o Inteligencia Artificial Explicable), sistemas técnicamente auditables y, de manera crucial, la adopción de un paradigma de **IA responsable por diseño** (*Responsible AI by Design*). Esta sección explora estas trayectorias prospectivas, delineando los contornos de una IA legal que pueda aspirar a ser un colaborador verdaderamente confiable para el profesional del derecho y un instrumento equitativo en la administración de justicia.

7.1 La búsqueda de la explicabilidad (XAI) en el contexto legal

Uno de los mayores obstáculos para la confianza y la adopción generalizada de los LLMs en tareas legales críticas es su naturaleza de "caja negra". Generan respuestas, a menudo complejas y matizadas, pero rara vez ofrecen una justificación inteligible de *cómo* llegaron a esas conclusiones o *en qué* información específica (y con qué ponderación) se basaron. En un dominio como el derecho, donde la capacidad de argumentar, justificar y trazar el razonamiento hasta las fuentes autorizadas es esencial, esta opacidad es profundamente problemática. La Inteligencia Artificial Explicable (XAI) emerge como un campo de investigación vital para abordar este desafío.

Las técnicas actuales de explicabilidad para LLMs (p. ej., análisis de atención, importancia de características, generación de justificaciones textuales post-hoc) a menudo proporcionan solo una visión superficial o aproximada del proceso de

toma de decisiones interno del modelo. Estas explicaciones pueden ser ellas mismas susceptibles de "alucinar" o pueden no reflejar fielmente los factores causales reales que llevaron a una salida particular (Rudin, 2019; Lipton, 2018).

7.1.1 El futuro de la XAI jurídica: de la fundamentación a la interpretación razonada

El futuro de la XAI legal no reside tanto en el desarrollo de explicaciones *post-hoc* para modelos opacos, sino en la evolución hacia arquitecturas de IA que incorporen la explicabilidad de forma nativa y significativa para el profesional del derecho. La verdadera revolución no se medirá por la capacidad de un sistema para superar un examen, sino por su habilidad para justificar sus conclusiones con transparencia y responsabilidad. Esta evolución puede conceptualizarse en tres niveles progresivos de madurez:

- **Respuestas fundamentadas (*Grounded Responses*):** Este es el nivel básico y el prerrequisito indispensable, centrado en la trazabilidad. La IA debe ser capaz de anclar cada afirmación en una fuente autorizada y verificable. El objetivo es responder a la pregunta: "*¿De dónde procede esta información?*". Sin una fundamentación sólida, cualquier resultado carece de la fiabilidad mínima necesaria para su uso profesional.
- **Respuestas argumentadas (*Argued Responses*):** El siguiente nivel trasciende la mera cita de fuentes para articular el razonamiento. No basta con saber *qué* fuente se usó, sino *cómo* se usó para construir la conclusión. La IA debe ser capaz de externalizar los pasos lógicos de su inferencia, demostrando una cadena de razonamiento coherente. El objetivo es responder a la pregunta: "*¿Cómo has llegado a esta conclusión a partir de las fuentes?*".
- **Respuestas basadas en interpretación razonada (*Reasoned Interpretation*):** Este es el nivel más avanzado y el verdadero horizonte de la IA jurídica. Aquí, la IA no solo aplica una regla, sino que es capaz de explicar *por qué* opta por una interpretación específica frente a otras alternativas plausibles. Implica ponderar matices, reconocer ambigüedades y justificar su aplicación de la norma a un contexto fáctico concreto. El objetivo es responder a la pregunta: "*¿Por qué esta es la interpretación o aplicación más adecuada en este caso?*".

Alcanzar este tercer nivel sigue siendo un desafío formidable, ya que la interpretación genuina requiere principios, ética y un entendimiento del contexto que los modelos actuales no poseen. Sin embargo, el desarrollo de arquitecturas de IA que sean intrínsecamente más interpretables es crucial para avanzar en este camino. Esto implica:

- **Modelos que externalizan el razonamiento jurídico:** Como se discutió con los modelos de razonamiento (Sección 5.4), la IA que puede articular sus pasos de inferencia de una manera que se asemeje a un análisis legal humano (identificando hechos relevantes, aplicando reglas, ponderando factores, citando autoridades para cada paso) será inherentemente más explicable y verificable, avanzando hacia el nivel de **argumentación**.
- **Visualización de la influencia de las fuentes en RAG:** En sistemas RAG, mejorar la capacidad de rastrear qué fragmentos específicos del contexto recuperado contribuyeron y con qué peso a cada parte de la respuesta generada. Herramientas de visualización que muestren estas conexiones podrían aumentar drásticamente la interpretabilidad y la **fundamentación**.
- **Explicaciones contrastantes y contrafácticas:** Desarrollar modelos capaces de explicar no solo por qué llegaron a una conclusión, sino también por qué no llegaron a otras conclusiones alternativas, o cómo cambiaría la conclusión si ciertos hechos o premisas fueran diferentes. Esto se alinea estrechamente con la forma en que los abogados analizan los problemas y es un paso clave hacia la **interpretación razonada**.
- **Gestión de la tensión entre explicabilidad y rendimiento:** La complejidad del razonamiento jurídico y la multiplicidad de factores que pueden influir en una decisión legal hacen que la explicabilidad completa sea un objetivo extremadamente ambicioso. Se debe reconocer y gestionar activamente la tensión potencial que existe entre la explicabilidad y el rendimiento del modelo: los sistemas más precisos a menudo son los más opacos. El futuro de la XAI legal radicará en encontrar un equilibrio óptimo donde la justificación del resultado sea suficientemente robusta para la validación profesional, sin sacrificar de manera inaceptable la eficacia del sistema.

7.2 La necesidad de auditabilidad técnica y de gobernanza

La fiabilidad y la responsabilidad en la IA legal no pueden depender únicamente de la buena fe de los desarrolladores o de la diligencia de los usuarios individuales. Se necesitan mecanismos robustos para la **auditoría independiente y continua** de estos sistemas, tanto a nivel técnico como de gobernanza.

1. Auditoría técnica de los modelos y sistemas RAG:

- *Desarrollo de estándares y métricas de auditoría específicos para IA legal*: Se requieren benchmarks y métricas estandarizadas (como los discutidos en la Sección 3) que vayan más allá de la precisión general y evalúen específicamente la propensión a las alucinaciones, la robustez ante entradas adversarias o ambiguas, la equidad (*fairness*) respecto a diferentes grupos, y la calidad de la fundamentación en sistemas RAG.
 - *Herramientas de auditoría automatizada y asistida por IA*: Desarrollar herramientas que puedan asistir a los auditores humanos en la evaluación a gran escala de los modelos, por ejemplo, generando automáticamente casos de prueba desafiantes, identificando posibles sesgos en los datos de entrenamiento o en las respuestas, o verificando la consistencia de las citas.
 - *Acceso controlado para auditoría (Sandboxing)*: Los reguladores o entidades de certificación independientes podrían requerir acceso a los modelos (posiblemente en entornos controlados o "sandboxes") para realizar pruebas exhaustivas antes de su despliegue en aplicaciones de alto riesgo.
2. **Auditoría de gobernanza de datos y procesos**: Más allá del modelo en sí, es crucial auditar los procesos y las prácticas de gobernanza de datos de las organizaciones que desarrollan e implementan IA legal.
- *Trazabilidad de los datos de entrenamiento*: Asegurar que se mantengan registros detallados sobre las fuentes, la curación y el pre-procesamiento de los datos utilizados para entrenar los modelos, permitiendo investigar posibles sesgos o errores.
 - *Evaluación de impacto algorítmico y ético*: Requerir que las organizaciones realicen evaluaciones de impacto rigurosas antes de desplegar sistemas de IA en contextos legales sensibles, identificando y mitigando proactivamente los riesgos potenciales.
 - *Mecanismos de supervisión humana y rendición de cuentas*: Auditar la efectividad de los mecanismos de supervisión humana implementados y asegurar que existan canales claros para la rendición de cuentas y la reparación en caso de errores o daños causados por la IA.

La auditabilidad no es solo una cuestión técnica, sino también una exigencia de buena gobernanza y un prerrequisito para generar confianza pública en la IA legal.

7.3 IA responsable por diseño (*Responsible AI by Design*) en el ámbito legal

El enfoque más proactivo y, en última instancia, más efectivo para construir una IA legal fiable es adoptar un paradigma de **Inteligencia Artificial Responsable por Diseño**. Esto implica integrar consideraciones éticas, de equidad, transparencia, robustez y fiabilidad **desde las primeras etapas del ciclo de vida del desarrollo de la IA**, en lugar de tratar estos aspectos como correcciones o parches aplicados a posteriori.

Este paradigma de 'Responsabilidad por Diseño' no es solo una aspiración ética o una buena práctica de ingeniería, sino que se está convirtiendo progresivamente en una expectativa regulatoria y, en algunos casos, en una obligación legal explícita. La Ley de IA de la UE (el Reglamento), a través de su detallado catálogo de requisitos para los sistemas de IA de alto riesgo —que abarcan desde la gestión de riesgos (**Artículo 9**) y la gobernanza de los datos de entrenamiento (**Artículo 10**) hasta la necesidad de una supervisión humana efectiva (**Artículo 14**) y la robustez técnica (**Artículo 15**)—, esencialmente codifica muchos de los principios fundamentales de la IA responsable. Al exigir estas consideraciones desde las fases de diseño y desarrollo, y a lo largo de todo el ciclo de vida del sistema, la EU-AIAct impulsa a los creadores de IA legal a ir más allá de la simple funcionalidad para priorizar la seguridad, la fiabilidad y la protección de los derechos fundamentales, donde la prevención de resultados perjudiciales derivados de alucinaciones se convierte en un objetivo central del diseño.

1. Principios fundamentales de la IA legal responsable por diseño:

Estos principios no son meras aspiraciones teóricas, sino que encuentran un eco directo y una validación institucional en los marcos normativos emergentes. Estos principios no son meras aspiraciones teóricas, sino que encuentran un eco directo y una validación institucional en los marcos normativos emergentes. Un ejemplo paradigmático es la ya mencionada Política del CTEAJE en España. Principios establecidos en este documento de obligado cumplimiento, como el de 'No sustitución' o el mandato de 'Revisión Humana Universal' que establece, son la materialización práctica del enfoque que se detalla a continuación, demostrando que la IA responsable por diseño está pasando de ser una buena práctica a una exigencia regulatoria.

- *Centrada en el ser humano*: Diseñar sistemas de IA que sirvan para aumentar y asistir al profesional legal, no para reemplazar su juicio crítico o su responsabilidad ética. El objetivo es la colaboración humano-IA, no la automatización completa de tareas complejas.

- *Fiabilidad y seguridad como prioridad*: La precisión factual, la robustez ante errores y la seguridad de los datos deben ser consideraciones primordiales en el diseño y la optimización de los modelos, incluso si esto implica ciertos compromisos en términos de fluidez o velocidad de generación.
- *Equidad y no discriminación (Fairness)*: Esforzarse activamente por identificar y mitigar los sesgos algorítmicos que podrían conducir a resultados discriminatorios o inequitativos en la aplicación de la ley. Esto requiere un análisis cuidadoso de los datos de entrenamiento y de los impactos diferenciales del modelo.
- *Transparencia y explicabilidad (contextualizadas)*: Diseñar sistemas que sean lo más transparentes y explicables posible dentro de las limitaciones técnicas, proporcionando a los usuarios información significativa sobre cómo funcionan y por qué generan ciertas respuestas.
- *Rendición de cuentas y gobernanza*: Establecer estructuras claras de gobernanza y rendición de cuentas para el desarrollo, despliegue y mantenimiento de los sistemas de IA legal.

2. Metodologías prácticas para la IA legal responsable por diseño:

- *Equipos de desarrollo multidisciplinares*: Involucrar a juristas, éticos y expertos en el dominio legal desde el inicio del proceso de diseño, no solo a ingenieros y científicos de datos.
- *Evaluación continua de riesgos y pruebas adversarias*: Implementar ciclos iterativos de evaluación de riesgos y pruebas de estrés (incluyendo pruebas específicas para detectar alucinaciones y sesgos) a lo largo de todo el desarrollo.
- *Mecanismos de retroalimentación y mejora continua*: Diseñar sistemas que puedan aprender y mejorar a partir de la retroalimentación de los usuarios expertos y de la monitorización de su rendimiento en el mundo real.
- *Adopción de estándares éticos y técnicos emergentes*: Mantenerse al día y adherirse a los estándares éticos y técnicos, así como a las mejores prácticas, que están siendo desarrollados por la comunidad investigadora, los organismos profesionales y los reguladores.

La IA Responsable por Diseño no es un estado final, sino un compromiso continuo con la mejora y la adaptación. Requiere una cultura organizacional que priorice la ética y la fiabilidad, y una voluntad de invertir en los recursos necesarios para construir sistemas que sean verdaderamente dignos de confianza en el sensible contexto legal.

7.4 La simbiosis avanzada humano-IA: colaboración y juicio aumentado

Mirando aún más hacia el futuro, la IA legal más fiable y efectiva probablemente no será aquella que intente reemplazar completamente al abogado, sino aquella que logre una **simbiosis avanzada y sinérgica con la inteligencia humana experta**. En este modelo, la IA no es meramente una herramienta pasiva, sino un colaborador activo que aumenta y refina las capacidades del profesional del derecho.

- **IA como "Investigador incansable y verificador preliminar"**: La IA podría encargarse de la búsqueda exhaustiva y el análisis preliminar de grandes volúmenes de información legal, identificando patrones, recuperando precedentes relevantes y señalando posibles problemas o inconsistencias, pero siempre presentando sus hallazgos al abogado para su validación y juicio estratégico.
- **IA como "Generador de hipótesis y argumentos alternativos"**: En lugar de proporcionar una única "respuesta", la IA podría generar múltiples líneas argumentales, interpretaciones o soluciones posibles para un problema legal, cada una con su fundamentación y sus posibles debilidades, permitiendo al abogado explorar un espectro más amplio de opciones estratégicas.
- **IA como "Traductor contextual y puente de comunicación"**: Una de las barreras más significativas en la práctica legal es la asimetría de información entre el abogado y el cliente, a menudo causada por la complejidad del lenguaje jurídico. Un sistema de IA avanzado, en lugar de ser una herramienta de uso exclusivo para el profesional, puede actuar como un traductor contextual, generando resúmenes o explicaciones de documentos y estrategias legales en un lenguaje adaptado al nivel de comprensión del cliente. Este enfoque, ejemplificado por la **"función Jerga" del proyecto Justicia** analizada previamente, no solo mejora la transparencia y la confianza en la relación cliente-abogado, sino que también capacita al cliente para tomar decisiones más informadas, humanizando el acceso a la justicia.
- **IA como "Entrenador personalizado y asistente de aprendizaje"**: La IA podría proporcionar retroalimentación detallada y personalizada sobre el trabajo de los abogados en formación, ayudándoles a identificar áreas de mejora en su investigación, redacción y razonamiento, siempre bajo la supervisión de mentores humanos.

- **Interfaces colaborativas intuitivas:** El desarrollo de interfaces de usuario que permitan una interacción fluida, iterativa y verdaderamente colaborativa entre el abogado y el sistema de IA será crucial. El abogado debe poder guiar, cuestionar y refinar fácilmente el trabajo de la IA.

Este futuro de colaboración aumentada requiere no solo avances en la tecnología de IA, sino también una evolución en la formación y las habilidades de los profesionales legales, quienes necesitarán ser competentes tanto en el derecho como en la interacción crítica y efectiva con estos sistemas inteligentes.

En conclusión, el camino hacia una IA legal verdaderamente fiable y beneficiosa es complejo y está en constante evolución. Si bien las alucinaciones y otros riesgos son desafíos significativos, no son insuperables. A través de un compromiso sostenido con la investigación en explicabilidad, el desarrollo de mecanismos robustos de auditabilidad, la adopción de principios de responsabilidad por diseño y, fundamentalmente, el reconocimiento del valor insustituible de la supervisión y el juicio humano, es posible vislumbrar un futuro donde la IA se convierta en un aliado poderoso y confiable en la búsqueda de una justicia más eficiente, accesible y equitativa. La tarea no es trivial, pero las recompensas potenciales para la profesión legal y la sociedad en su conjunto son inmensas.

8 Navegando la frontera ética y regulatoria: implicaciones de las alucinaciones de la IA en el contexto legal global, con énfasis en el Derecho Español y Europeo

La irrupción de los Grandes Modelos de Lenguaje (LLMs) y su inherente propensión a generar "alucinaciones" —resultados textuales que, aunque a menudo fluidos y convincentes, carecen de veracidad fáctica o fundamento legal— no es meramente un desafío técnico. Este fenómeno penetra profundamente en el tejido de la profesión jurídica, interpelando sus fundamentos éticos y planteando interrogantes críticos sobre la adecuación de los marcos regulatorios existentes a nivel global. Mientras que gran parte del debate inicial y la jurisprudencia temprana sobre sanciones por el uso indebido de IA en litigios ha emanado del sistema del *common law* estadounidense (ejemplificado por casos como *Mata v. Avianca, Inc.*), las implicaciones éticas y la necesidad de una respuesta regulatoria son universales, aunque su manifestación y las soluciones propuestas deban necesariamente adaptarse a las particularidades de cada ordenamiento jurídico. Esta sección se adentra en el complejo panorama de las obligaciones deontológicas y los desafíos regulatorios que las alucinaciones de la IA legal plantean, con una atención particular a las realidades y perspectivas del derecho español y el marco normativo europeo, sin dejar de lado las lecciones aprendidas de otras jurisdicciones.

8.1 El imperativo deontológico en la era de la IA: reafirmando los deberes fundamentales del abogado en España y Europa

Los deberes deontológicos tradicionales de competencia y diligencia están cristalizando en un nuevo pilar ético para la era digital: la "Obligación de Competencia Tecnológica". Como se resume en análisis exhaustivos sobre la ética de los LLMs en la abogacía, esta obligación no solo exige entender los beneficios de la tecnología, sino también sus riesgos y limitaciones, incluyendo una comprensión fundamental de fenómenos como las alucinaciones (Shao et al., 2025). Cumplir con este deber implica la verificación rigurosa de los resultados generados por la IA y mantener una supervisión crítica, reconociendo que la responsabilidad final del trabajo recae inequívocamente en el profesional humano.

Los códigos deontológicos que rigen la abogacía en España (como el Código Deontológico de la Abogacía Española) y en el ámbito europeo (como el Código Deontológico de los Abogados Europeos del CCBE), al igual que sus contrapartes estadounidenses, establecen una serie de deberes fundamentales que, si bien no fueron concebidos con la IA generativa en mente, son directamente aplicables y adquieren una nueva dimensión ante el riesgo de alucinaciones.

1. **Deber de competencia profesional:** Este es, quizás, el deber más inmediatamente interpelado. La competencia profesional exige no solo el conocimiento sustantivo del derecho aplicable, sino también la habilidad para utilizar adecuadamente las herramientas y tecnologías que se emplean en el ejercicio profesional. En la era de la IA, esto se traduce en:
 - *Alfabetización en IA y comprensión de sus limitaciones:* Un abogado competente en España o Europa no puede permitirse ignorar el funcionamiento básico de los LLMs, su naturaleza probabilística y, crucialmente, su potencial para generar alucinaciones. Esto no implica ser un experto en IA, sino poseer una comprensión funcional suficiente para evaluar críticamente sus resultados y los riesgos asociados (Yamane, 2020; Choi and Schwarcz, 2024).
 - *Deber de verificación rigurosa:* La competencia exige que el abogado verifique de forma independiente la exactitud y pertinencia de cualquier información o borrador generado por una IA antes de utilizarlo en el asesoramiento al cliente o en actuaciones ante los tribunales. Confiar ciegamente en la salida de un LLM,

especialmente en asuntos de alta complejidad o riesgo, podría constituir una grave falta de competencia. Las guías emergentes de colegios de abogados europeos y españoles previsiblemente enfatizarán este punto.

- *Conciencia del sesgo de automatización:* Tal como en otros contextos, los profesionales del derecho en España y Europa deben ser conscientes del "sesgo de automatización" y mantener un escepticismo profesional saludable, resistiendo la tentación de delegar el juicio crítico a la máquina, por muy eficiente que esta parezca (Drabiak et al., 2023).

Este "sesgo de automatización" no es una mera tendencia a la confianza ciega; se manifiesta a través de mecanismos cognitivos específicos que han sido cuantificados. Los LLMs introducen sesgos de encuadre en más de un 20% de los casos, cambiando la valencia emocional o el énfasis de la información sin alterar los hechos subyacentes (Alessa et al., 2025). Un profesional del derecho que interactúa con estos resultados puede ver su percepción de un caso sutilmente moldeada antes de haber formado su propio juicio independiente. Por lo tanto, el escepticismo profesional no es solo una buena práctica, sino una salvaguarda cognitiva esencial.

2. **Deber de diligencia:** Estrechamente ligado a la competencia, el deber de diligencia obliga al abogado a actuar con el cuidado y la atención necesarios en la defensa de los intereses del cliente. En el contexto de la IA y las alucinaciones:

- *Verificación como parte de la diligencia:* La promesa de eficiencia de la IA no puede ir en detrimento de la calidad y la corrección del trabajo. Un abogado diligente debe invertir el tiempo necesario para validar la información generada por la IA, asegurando que cualquier uso de la misma se base en información verificada y legalmente sólida. La "rapidez" no puede justificar la "precipitación negligente".
- *Actualización continua:* Dada la rápida evolución de la tecnología IA, la diligencia también puede implicar un deber de mantenerse razonablemente informado sobre los avances, los nuevos riesgos identificados (como tipos específicos de alucinaciones) y las mejores prácticas para el uso de estas herramientas.

La necesidad de una verificación rigurosa no es una recomendación abstracta, sino una exigencia derivada de la evidencia empírica. Con tasas de alucinación documentadas que oscilan entre el 17% y más del 33% en las principales herramientas comerciales de IA legal, confiar ciegamente en sus resultados constituye una clara abdicación del deber de competencia (Magesh et al., 2025). Como concluye el estudio, los abogados se enfrentan a una difícil elección: "verificar a mano cada proposición y cita producida por estas herramientas (socavando así las ganancias de eficiencia prometidas), o arriesgarse a usar estas herramientas sin información completa sobre sus riesgos específicos (descuidando así sus deberes centrales de competencia y supervisión)".

3. **Secreto profesional y protección de datos:** Este es un área de particular sensibilidad en el contexto europeo y español, dada la robusta normativa en materia de protección de datos (Reglamento General de Protección de Datos - RGPD).

- *Confidencialidad de la información del cliente:* Introducir información confidencial o datos personales de clientes en plataformas de IA, especialmente aquellas cuyos servidores están fuera de la UE o cuyas políticas de uso de datos no son transparentes o no cumplen con el RGPD, representa un riesgo significativo. Las alucinaciones no son el riesgo directo aquí, pero la elección de la herramienta y la gestión de los datos son cruciales.
- *Cumplimiento del RGPD:* Cualquier tratamiento de datos personales a través de una IA debe cumplir con los principios del RGPD (licitud, lealtad, transparencia, limitación de la finalidad, minimización de datos, exactitud, limitación del plazo de conservación, integridad y confidencialidad, y responsabilidad proactiva). Los proveedores y usuarios de IA legal deben poder demostrar este cumplimiento.

4. **Lealtad e independencia:** El abogado debe lealtad a su cliente y mantener su independencia de criterio. Si una herramienta de IA sugiere una línea de actuación basada en información alucinada o sesgada, el abogado debe ejercer su juicio independiente para desestimarla si no sirve a los mejores intereses del cliente o contraviene la ley.

5. **Deber de lealtad procesal y colaboración con la Administración de Justicia:** En muchos sistemas de derecho civil como el español, existe un fuerte énfasis en la buena fe y la lealtad procesal. Presentar ante un tribunal argumentos, pruebas o jurisprudencia que se sabe (o se debería saber tras una verificación diligente) que son falsos o fabricados por una IA, constituiría una grave violación de estos deberes, con posibles consecuencias disciplinarias y procesales. La integridad del sistema judicial depende de la fiabilidad de la información presentada por las partes.

El impacto de las alucinaciones de la IA en estos deberes deontológicos es innegable y exige una reflexión profunda por parte de los colegios profesionales, los órganos disciplinarios y cada abogado individualmente.

8.2 Desafíos y perspectivas regulatorias en España y la Unión Europea

El panorama regulatorio para la IA, y específicamente para la IA legal y sus riesgos de alucinación, está en plena efervescencia, con la Unión Europea a la vanguardia a través de su propuesta de Ley de Inteligencia Artificial (EU-AIAct).

8.2.1 La Ley de IA de la UE: un marco jerárquico y basado en riesgos para la gobernanza de la IA legal

A la vanguardia de los esfuerzos globales por establecer un marco normativo integral para la inteligencia artificial se encuentra la Unión Europea con su Ley de IA de la UE. Esta ambiciosa pieza legislativa, de alcance potencialmente global debido al conocido "efecto Bruselas", adopta un **enfoque estratificado y basado en el riesgo**, clasificando los sistemas de IA en categorías que van desde un riesgo inaceptable (y, por tanto, prohibidos) hasta un riesgo mínimo, pasando por categorías de riesgo limitado y, crucialmente para muchas aplicaciones legales, **alto riesgo**. Es esta categoría de 'alto riesgo' la que impone las obligaciones más significativas a los desarrolladores, proveedores y, en ciertos casos, usuarios de sistemas de IA (European Union, 2024; Hitaj et al., 2023; Petit & De Cooman, 2020).

La determinación de si una herramienta específica de IA legal cae dentro de la categoría de 'alto riesgo' dependerá de su finalidad prevista y del contexto de su uso, según lo detallado en el Anexo III de la EU-AIAct. Áreas explícitamente mencionadas como de alto riesgo que tienen una clara tangencia con el sector legal incluyen los sistemas de IA utilizados en la **administración de justicia y los procesos democráticos**, así como aquellos empleados para la evaluación de la solvencia crediticia o la selección en procesos de contratación, que a menudo involucran análisis de perfiles con implicaciones legales. Es razonable argumentar que herramientas de IA que asistan en la toma de decisiones judiciales, en la evaluación de la admisibilidad de pruebas, en la predicción de reincidencia, o incluso sistemas de investigación legal muy avanzados cuya salida errónea pueda tener un impacto directo y significativo en los derechos fundamentales de un individuo (p. ej., en un proceso penal o en la determinación de la custodia de un menor) podrían ser clasificados como de alto riesgo.

Para estos sistemas de IA de alto riesgo, la EU-AIAct establece un conjunto exhaustivo de requisitos obligatorios que deben cumplirse antes de su introducción en el mercado y mantenerse durante todo su ciclo de vida. Muchos de estos requisitos tienen una relevancia directa para la prevención y mitigación de las alucinaciones:

1. **Sistemas robustos de gestión de riesgos (Artículo 9):** Se exige el establecimiento, implementación, documentación y mantenimiento de un proceso continuo de gestión de riesgos. Esto implica la identificación de los riesgos previsible asociados al sistema (incluyendo los derivados de alucinaciones), la estimación y evaluación de dichos riesgos, y la adopción de medidas adecuadas para su control. La gestión del riesgo de generar información legal incorrecta o fabricada debería ser un componente central de este sistema.
2. **Gobernanza y calidad de los datos (Artículo 10):** Este artículo es particularmente pertinente para las alucinaciones originadas en datos de entrenamiento defectuosos. Exige que los conjuntos de datos de entrenamiento, validación y prueba sean 'relevantes, representativos, libres de errores y completos'. Se deben aplicar prácticas adecuadas de gobernanza de datos, incluyendo un examen de los posibles sesgos y la adopción de medidas para mitigarlos. Para la IA legal, esto implica la necesidad crítica de utilizar corpus jurídicos actualizados, verificados y que reflejen adecuadamente la diversidad y complejidad del ordenamiento jurídico.
3. **Documentación técnica exhaustiva (Artículo 11 y Anexo IV):** Los proveedores deben elaborar una documentación técnica detallada que describa, entre otras cosas, la arquitectura del sistema, sus capacidades y limitaciones, los algoritmos utilizados, los datos de entrenamiento y los procesos de prueba y validación. Esta documentación es esencial para la evaluación de conformidad y para que los supervisores y usuarios comprendan cómo funciona el sistema y cuáles son sus umbrales de fiabilidad.
4. **Mecanismos de registro de eventos (Logging Capabilities) (Artículo 12):** Los sistemas de alto riesgo deben estar equipados con capacidades de registro que aseguren un nivel adecuado de trazabilidad de su funcionamiento. Estos 'logs' podrían ser cruciales para investigar *a posteriori* el origen de una alucinación específica o para auditar el rendimiento general del sistema.
5. **Transparencia y provisión de información a los usuarios (Artículo 13):** Los sistemas deben ser diseñados y desarrollados de manera que los usuarios puedan interpretar la salida del sistema y utilizarla de forma apropiada. Las instrucciones de uso deben incluir información concisa, completa, correcta y clara sobre la identidad del proveedor, la finalidad prevista del sistema, su nivel de precisión, robustez y ciberseguridad, así como sus limitaciones conocidas y los riesgos previsible – lo cual incluye, o debería incluir, la propensión a generar alucinaciones y la necesidad de verificación humana.
6. **Supervisión humana efectiva (Artículo 14):** Este es un pilar fundamental. La EU-AIAct exige que los sistemas de alto riesgo sean diseñados para permitir una supervisión humana adecuada. Las medidas pueden

incluir la capacidad del supervisor humano para comprender plenamente las capacidades y limitaciones del sistema, para decidir no utilizar el sistema en una situación particular, para anular una decisión tomada por el sistema, o para intervenir en su funcionamiento. Esta supervisión es la última barrera contra las consecuencias de una alucinación no detectada por el propio sistema.

7. **Precisión, robustez y ciberseguridad (Artículo 15):** Los sistemas deben alcanzar un nivel apropiado de precisión, robustez y ciberseguridad a lo largo de su ciclo de vida y ser consistentes en este aspecto. Las alucinaciones son una manifestación clara de una falta de precisión y robustez fáctica. Se espera que los sistemas sean resilientes a errores, fallos o inconsistencias, así como a intentos de uso malintencionado.

El cumplimiento de estos requisitos se verificará mediante **evaluaciones de conformidad** antes de que el sistema de IA de alto riesgo pueda ser introducido en el mercado de la UE. Además, se establecen obligaciones de **vigilancia post-comercialización** para los proveedores, quienes deben monitorizar el rendimiento de sus sistemas y reportar cualquier incidente grave o mal funcionamiento. Las **sanciones** por incumplimiento de la EU-AIAct son significativas, pudiendo alcanzar hasta 35 millones de euros o el 7% del volumen de negocio anual mundial total del ejercicio financiero anterior, lo que subraya la seriedad con la que la UE aborda los riesgos de la IA.

El impacto de la EU-AIAct en las estrategias de mitigación de alucinaciones en la IA legal, como RAG, es profundo. Muchos de los requisitos de la Ley –calidad de datos, transparencia sobre el funcionamiento, robustez, supervisión humana– impulsarán a los desarrolladores a adoptar de manera proactiva y rigurosa muchas de las estrategias de optimización discutidas en la Sección 5 de este ensayo, no como una opción de mejora, sino como una condición para el acceso al mercado. Aunque la EU-AIAct no prescribe soluciones técnicas específicas, sí establece un marco de exigencias que fomentará la innovación hacia una IA legal más fiable y responsable. Dado el peso del mercado europeo, es muy probable que la EU-AIAct tenga un "efecto Bruselas", influyendo en los estándares de desarrollo de IA legal a nivel global.

8.2.2 Reglamento General de Protección de Datos (RGPD)

Aunque no es específico para la IA, el RGPD ya impone obligaciones significativas que son relevantes para el desarrollo y uso de IA legal que trate datos personales.

- *Principios de exactitud y minimización de datos:* Estos principios son directamente relevantes para combatir los datos de entrenamiento defectuosos que pueden llevar a alucinaciones.
- *Derecho a no ser objeto de decisiones automatizadas (Artículo 22):* Si un sistema de IA legal toma decisiones que produzcan efectos jurídicos significativos en una persona (o le afecten de modo similar), el Artículo 22 del RGPD podría limitar su uso o requerir una intervención humana significativa.

8.2.3 Responsabilidad civil y profesional

- *Régimen General de Responsabilidad:* En España, la responsabilidad civil del abogado por negligencia profesional se rige por el Código Civil y la jurisprudencia. Si el uso indebido de una IA (p. ej., confiar en información alucinada sin verificación) causa un daño al cliente, el abogado podría ser considerado responsable.
- *Propuesta de Directiva Europea sobre responsabilidad por IA:* La Comisión Europea ha propuesto una Directiva para adaptar las normas de responsabilidad civil extracontractual a la IA. Esta propuesta busca facilitar que las víctimas de daños causados por IA obtengan una reparación, por ejemplo, aliviando la carga de la prueba en ciertos casos o estableciendo una presunción de causalidad si el proveedor de IA no ha cumplido con ciertos deberes de diligencia (potencialmente incluyendo aquellos relacionados con la prevención de alucinaciones).
- *Responsabilidad del proveedor vs. usuario profesional:* La atribución de responsabilidad entre el desarrollador/proveedor de la herramienta de IA y el abogado usuario será un área compleja y probablemente litigiosa. Los términos de servicio de los proveedores a menudo incluyen extensas cláusulas de exención de responsabilidad, pero su validez podría ser cuestionada, especialmente si se demuestra negligencia grave o un defecto inherente en el diseño del producto que lo hace propenso a generar información legalmente perjudicial (Calderon et al., 2022; Lantyer, 2024).

8.2.4 El papel de los Colegios de Abogados y órganos deontológicos

En España, los colegios de abogados (y el Consejo General de la Abogacía Española) desempeñan un papel crucial en el establecimiento de normas deontológicas y en la supervisión de su cumplimiento.

- *Emisión de guías y directrices específicas:* Es previsible y deseable que estos organismos desarrollen y publiquen guías específicas sobre el uso ético y competente de la IA generativa por parte de los abogados, abordando explícitamente el riesgo de alucinaciones y el deber de verificación.
- *Formación continua:* La oferta de programas de formación continua sobre IA, sus capacidades, riesgos y uso responsable será esencial para asegurar que los profesionales mantengan el nivel de competencia requerido.
- *Potestad disciplinaria:* Los colegios podrían ejercer su potestad disciplinaria en casos de uso manifiestamente negligente o irresponsable de la IA que resulte en perjuicio para el cliente o para la administración de justicia.

8.2.5 Necesidad de estándares técnicos y benchmarks

Para que cualquier marco regulatorio sea efectivo, se necesitarán estándares técnicos y benchmarks independientes que permitan evaluar de manera objetiva la fiabilidad, precisión y propensión a las alucinaciones de las herramientas de IA legal. La colaboración entre juristas, tecnólogos y organismos de normalización será crucial en este aspecto.

La colaboración entre juristas y tecnólogos ya está trazando un camino en esta dirección, como lo demuestra el uso de exámenes profesionales de abogacía como benchmark para evaluar modelos de IA en dominios legales específicos (Gupta et al., 2025). La adopción de este tipo de pruebas estandarizadas como benchmarks podría convertirse en un requisito para que los proveedores de tecnología legal demuestren la fiabilidad y competencia de sus sistemas en una jurisdicción específica, proporcionando a los reguladores y a los consumidores una base objetiva para la evaluación.

8.3 Hacia una integración ética y responsable de la IA en la práctica legal española y europea

El camino hacia una integración de la IA en el derecho que sea a la vez innovadora y segura, especialmente en el contexto español y europeo con su fuerte tradición de protección de derechos y rigor normativo, requiere un compromiso proactivo y colaborativo de todos los actores implicados.

- **Para los profesionales del Derecho:** La adopción de una mentalidad de **escepticismo informado y verificación diligente** es primordial. La IA debe ser vista como una herramienta poderosa de asistencia, no como un oráculo infalible. La formación continua y la alfabetización digital en IA serán competencias esenciales.
- **Para los bufetes y organizaciones legales:** Es necesario establecer **políticas internas claras** sobre el uso aceptable y responsable de la IA, incluyendo protocolos de verificación obligatorios, directrices sobre la gestión de datos confidenciales y programas de formación para sus profesionales. La inversión en herramientas de IA debe ir acompañada de una inversión en la capacitación para su uso seguro.
- **Para los proveedores de tecnología legal:** Existe una responsabilidad creciente de desarrollar herramientas que sean **"seguras por diseño" (safety by design)**, incorporando mecanismos para minimizar las alucinaciones y, crucialmente, siendo **radicalmente transparentes** sobre las capacidades, limitaciones y tasas de error conocidas de sus productos. Las afirmaciones de marketing deben ser realistas y estar respaldadas por evidencia empírica robusta e independiente.
- **Para las instituciones educativas (Facultades de Derecho):** Es fundamental integrar la enseñanza sobre IA, sus implicaciones legales y éticas, y las habilidades necesarias para su uso crítico en los planes de estudio, preparando a las futuras generaciones de juristas para un entorno profesional transformado por la tecnología.
- **Para los reguladores y órganos deontológicos:** Se requiere una adaptación proactiva y reflexiva de los marcos normativos y deontológicos. Esto puede implicar la clarificación de los deberes existentes a la luz de la IA, el desarrollo de nuevas directrices específicas, y el fomento de una cultura de responsabilidad y rendición de cuentas. La Ley de IA de la UE será un marco clave, pero su implementación y supervisión efectivas en el sector legal requerirán un esfuerzo continuo.
- **Para la Judicatura:** Los tribunales también se enfrentarán al desafío de evaluar la información generada por IA presentada por las partes y, potencialmente, de utilizar la IA en sus propias labores. La formación judicial y el desarrollo de protocolos para el uso de IA en el ámbito judicial serán necesarios para mantener la integridad del proceso.

En conclusión final, las alucinaciones de la IA no son un mero artefacto técnico, sino un síntoma de la tensión fundamental entre la naturaleza probabilística de los LLMs y la exigencia de certeza y fiabilidad del sistema legal. Abordar este desafío en el contexto español y europeo exige un **enfoque que combine la innovación tecnológica con una reafirmación de los principios éticos fundamentales de la abogacía**, una adaptación inteligente de los marcos regulatorios y, sobre todo, un compromiso inquebrantable con el juicio crítico y la supervisión humana. La IA puede ser una herramienta poderosa para el derecho, pero solo si se navega su laberinto con prudencia, diligencia y una profunda conciencia de sus limitaciones actuales.

En definitiva, una integración exitosa y responsable de la IA en el ecosistema legal requiere un cambio cultural que trascienda la búsqueda de soluciones tecnológicas simplistas y adopte un paradigma de evaluación crítica y diligencia informada. Esto implica la internalización de tres principios operativos fundamentales:

- **Priorizar la fiabilidad sobre la velocidad de generación.** La eficiencia real de una herramienta de IA no debe medirse únicamente por la rapidez con la que genera un resultado. Un sistema es verdaderamente eficiente solo si sus productos son fiables, minimizando así el tiempo y el esfuerzo requeridos en la indispensable fase de verificación humana. La fiabilidad, por tanto, es el verdadero multiplicador de la eficiencia en los flujos de trabajo jurídicos.
- **Fomentar el desarrollo de soluciones especializadas (Domain-Specific).** El sector jurídico se beneficiará más de herramientas diseñadas explícitamente para sus desafíos únicos que de la adaptación de modelos de propósito general. Las soluciones deben atender con precisión a las complejidades del razonamiento jurídico, la jerarquía normativa y los exigentes requisitos de confiabilidad, exigiendo un enfoque de desarrollo que priorice la profundidad contextual sobre la amplitud funcional.
- **Instituir una cultura de validación rigurosa y retroalimentación crítica.** La adopción de nuevas tecnologías debe estar guiada por una evaluación objetiva y empírica, en lugar de una aceptación acrítica impulsada por la novedad. El ecosistema (incluyendo profesionales, desarrolladores y académicos) debe demandar y proporcionar una retroalimentación rigurosa y honesta sobre el rendimiento, las limitaciones y los riesgos de estas herramientas. El progreso sostenible se fundamenta en la crítica constructiva y la validación independiente.

La adopción de estos principios pragmáticos es fundamental para que el sector legal pueda navegar la complejidad de la era de la IA con la prudencia, diligencia y profunda conciencia que esta exige.

9 Conclusión: de la alucinación a la amplificación — principios para una IA jurídica fiable

La irrupción de los Grandes Modelos de Lenguaje en el ecosistema jurídico presenta una paradoja fundamental: una tecnología con un potencial sin precedentes para democratizar y eficientar el acceso a la justicia, intrínsecamente lastrada por un defecto de diseño que atenta contra el pilar del derecho: **la veracidad**. Sin embargo, este análisis ha revelado que la veracidad en el derecho es un concepto dual, que abarca tanto la fidelidad factual —amenazada directamente por la alucinación— como la solidez interpretativa, que sigue siendo dominio exclusivo del juicio humano. Este informe, por tanto, ha diseccionado el fenómeno de las "alucinaciones" no como un mero error técnico, sino como una característica sistémica que exige un cambio de paradigma: pasar de buscar una IA que "sabe la verdad" a construir una IA que "amplifica la capacidad del profesional para interpretarla".

Más que un simple resumen, esta conclusión destila los hallazgos del análisis en un **conjunto de principios rectores y un marco de trabajo** para guiar a profesionales, desarrolladores y reguladores en la navegación de este complejo nuevo territorio.

9.1 Conclusiones fundamentales: un resumen estructurado

Los análisis detallados a lo largo de este documento convergen en cuatro conclusiones clave e interrelacionadas:

- **La alucinación no es un "bug", es una característica.** El desafío principal reside en comprender que la propensión de los LLMs de propósito general a "inventar" no es un fallo a corregir, sino una consecuencia directa de su arquitectura, diseñada para la fluidez probabilística y no para la fidelidad factual. Esto nos obliga a abandonar la idea de un "oráculo creativo" y adoptar un paradigma radicalmente distinto.
- **RAG es el camino, no el destino.** La Generación Aumentada por Recuperación (RAG) es, sin duda, la estrategia de mitigación más importante y el fundamento de la IA jurídica moderna. Sin embargo, la evidencia empírica es contundente: una implementación canónica de RAG reduce, pero **no elimina**, las alucinaciones. Tratarla como una solución "plug-and-play" es un error; debe ser considerada como el punto de partida, un motor prometedor que requiere una optimización holística y rigurosa en cada uno de sus componentes para ser verdaderamente fiable.
- **La fiabilidad se construye, no se instala: el imperativo de la optimización holística.** La transición de una IA que alucina a una IA fiable no depende de un único avance, sino de una sinergia de mejoras estratégicas a lo largo de todo el ciclo de vida de la información. Esto incluye:
 - **Curación estratégica de datos:** Un fundamento de conocimiento verificado, actualizado y jerarquizado es el cimiento indispensable.

- **Recuperación sofisticada:** Ir más allá de la búsqueda semántica simple para incorporar la conciencia de la jerarquía normativa (el principio Kelseniano), el contexto y la lógica jurídica.
- **Generación fiel y razonamiento guiado:** Utilizar ingeniería de prompts avanzada y *fine-tuning* para instruir a los LLMs no solo a responder, sino a "pensar" de manera estructurada, transparente y estrictamente anclada a las fuentes.
- **Verificación Post-Hoc robusta:** Implementar capas de seguridad algorítmica y humana como última línea de defensa indispensable.
- **El factor humano es irreductible y se fortalece.** Lejos de volver obsoleto al profesional del derecho, el desafío de la veracidad redefine y fortalece su rol. La supervisión experta, crítica e informada no es una opción, sino una obligación deontológica, profesional y, cada vez más, regulatoria (como lo demuestran la Ley de IA de la UE y las políticas del CTEAJE en España). El futuro no es la automatización, sino la **amplificación cognitiva**: la IA se convierte en una herramienta para potenciar el juicio humano, liberándolo para que se centre en la estrategia, la ética y la empatía.

9.2 Propuesta de un marco de trabajo: IA generativa vs. IA consultiva

Para guiar la adopción y el desarrollo futuro, proponemos un marco conceptual claro que distinga dos tipos de IA con perfiles de riesgo y aplicaciones radicalmente diferentes en el derecho:

- **IA generativa de propósito general (el "oráculo creativo"):**
 - **Función:** Ideación, brainstorming, redacción de borradores no críticos, resumen de textos generales.
 - **Riesgo inherente:** Alto riesgo de alucinación factual, extrínseca e intrínseca. Opacidad en las fuentes.
 - **Principio de uso:** Utilizar siempre con escepticismo extremo, como un asistente creativo cuya producción **nunca** debe ser considerada una fuente de verdad. Requiere una verificación humana completa desde cero.
- **IA consultiva especializada (el "archivero experto"):**
 - **Función:** Investigación jurídica, *due diligence*, análisis documental, respuesta a consultas basadas en un corpus verificado.
 - **Riesgo inherente:** Bajo riesgo de fabricación, pero riesgo **persistente** de alucinaciones sutiles (*mis-grounding*, errores de síntesis).
 - **Principio de uso:** Diseñada para la fiabilidad. Debe ser transparente, citable y auditable. Aun así, exige una verificación crítica por parte del profesional, pero enfocada en la correcta interpretación y aplicación de las fuentes proporcionadas, no en su existencia.

La mitigación efectiva de las alucinaciones en el sector legal no reside en mejorar incrementalmente el modelo generativo, sino en **adoptar deliberadamente un paradigma consultivo** donde la veracidad y la trazabilidad son el núcleo del diseño, no una característica añadida.

9.3 Una mirada al futuro: la llamada a una integración responsable

El camino hacia una IA jurídica verdaderamente fiable y beneficiosa está trazado, pero no es sencillo. No depende de encontrar un "interruptor mágico" que elimine los errores, sino de un compromiso colectivo y sostenido de todo el ecosistema legal.

- **Para los desarrolladores,** el desafío es construir sistemas que no solo sean potentes, sino transparentes, auditables y diseñados con una profunda humildad sobre sus limitaciones.
- **Para los profesionales del derecho,** el reto es cultivar una cultura de **escepticismo informado**: abrazar la tecnología como una herramienta de amplificación, pero nunca abdicar de la responsabilidad final del juicio crítico.
- **Para los reguladores e instituciones,** la tarea es continuar desarrollando marcos normativos que fomenten la innovación responsable, estableciendo estándares claros de fiabilidad y exigiendo la supervisión humana como un pilar innegociable.

En última instancia, la inteligencia artificial no es una fuerza externa que "impacta" en el derecho; es un nuevo material con el que estamos construyendo las herramientas del futuro. La calidad, seguridad y justicia de esas herramientas dependerán de nuestra habilidad para infundir en ellas los principios atemporales de nuestra profesión: rigor, diligencia

y un compromiso inquebrantable con la verdad. El objetivo final, y la gran promesa de esta era, no es simplemente automatizar procesos, sino, como se ha sostenido a lo largo de este informe, **humanizar la tecnología**, poniéndola al servicio de una justicia más accesible, eficiente y, sobre todo, fiable.

Agradecimientos

Nunca el trabajo necesario para el estudio, análisis y desarrollo de una investigación tan profunda como la que recoge este paper depende en exclusiva de su autor. Este trabajo hubiera sido impensable sin el esfuerzo previo de todos y cada uno de los investigadores que han sumado a la sociedad con sus papers anteriores. A todos ellos mi agradecimiento más sincero por lo que su trabajo ha representado para mi y para toda la comunidad científica y técnica.

Deseo agradecer igualmente a Little John y a todos y cada uno de sus miembros tanto su apoyo como el valor de sus revisiones. En especial a Daniel Vecino por las innumerables sesiones de trabajo en común y su “revisión pixel”, siempre tan crítica y completa como amable.

Es justo además agradecer la inspiración que para mi ha sido Asier Gutiérrez-Fandiño. Sin él esta publicación no hubiera sido posible ya que supuso un detonante claro en la pasión que en mi despierta el mundo de la Inteligencia Artificial.

Finalmente deseo agradecer la colaboración desinteresada de todos y cada uno de los profesionales que han tenido acceso previo a este paper. Sus comentarios han sido claves para recibir el empuje necesario que siempre me ha llevado a ir un poco más lejos en cada punto de análisis.

A todos ellos y a ti como lector de este trabajo gracias.

References

- [1] Dahl, Matthew and Magesh, Varun and Suzgun, Mirac and Ho, Daniel E. Large Legal Fictions: Profiling Legal Hallucinations in Large Language Models. In *Journal of Legal Analysis*, 16(1):64–93. Oxford University Press, 2024.
- [2] Choi, Jonathan H. and Schwarcz, Daniel. AI Assistance in Legal Analysis: An Empirical Study. In *Journal of Legal Education*, Forthcoming, 2024.
- [3] Livermore, Michael A. and Herron, Felix and Rockmore, Daniel. Language Model Interpretability and Empirical Legal Studies. In *Journal of Institutional and Theoretical Economics*, Forthcoming, 2024.
- [4] Alessa, Abeer and Lakshminarasimhan, Akshaya and Somane, Param and Skirzynski, Julian and McAuley, Julian and Echterhoff, Jessica. How Much Content Do LLMs Generate That Induces Cognitive Bias in Users? In *arXiv preprint arXiv:2507.03194*, 2025.
- [5] Rodgers, Ian and Armour, John and Sako, Mari. How Technology Is (or Is Not) Transforming Law Firms. In *Annual Review of Law and Social Science*, 19:299–317, 2023.
- [6] Choi, Jonathan H. and Hickman, Kristin E. and Monahan, Amy and Schwarcz, Daniel. ChatGPT Goes to Law School. In *Journal of Legal Education*, 71(3):387–400, 2022.
- [7] Rahul Hemrajani. Evaluating the Role of Large Language Models in Legal Practice in India. In *arXiv preprint arXiv:2508.09713*, 2025.
- [8] Gupta, Jatin and Sharma, Akhil and Singhanian, Saransh and Abidi, Ali Imam. Legal Assist AI: Leveraging Transformer-based Model for Effective Legal Assistance. In *arXiv preprint arXiv:2505.22003*, 2025.
- [9] Katz, Daniel Martin and Bommarito, Michael James and Gao, Shang and Arredondo, Pablo. GPT-4 Passes the Bar Exam. SSRN Working Paper, 2023.
- [10] Kalai, Adam Tauman, Ofir Nachum, Santosh S. Vempala, and Edwin Zhang. Why Language Models Hallucinate. OpenAI, Technical Report, September 2025.
- [11] Blair-Stanek, Andrew and Holzenberger, Nils and Van Durme, Benjamin. Can GPT-3 Perform Statutory Reasoning? In *Proceedings of the Nineteenth International Conference on Artificial Intelligence and Law (ICAIL 2023)*, pages Braga, Portugal. Association for Computing Machinery, 2023.
- [12] Guha, Neel and Nyarko, Julian and Ho, Daniel E. and Ré, Christopher and Chilton, Adam and Narayana, Aditya and Chohlas-Wood, Alex and Peters, Austin and Waldon, Brandon and Rockmore, Daniel N. and others. LegalBench: A Collaboratively Built Benchmark for Measuring Legal Reasoning in Large Language Models. Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track, 2023.

- [13] Weiser, Benjamin. Here's What Happens When Your Lawyer Uses ChatGPT. *The New York Times*, May 2023.
- [14] Romoser, James. No, Ruth Bader Ginsburg Did Not Dissent in Obergefell — and Other Things ChatGPT Gets Wrong about the Supreme Court. *SCOTUSblog*, Jan 2023.
- [15] Ludwig, Florian and Zesch, Torsten and Zufall, Frederike. Conditioning Large Language Models on Legal Systems? Detecting Punishable Hate Speech. In *arXiv preprint arXiv:2508.06456*, 2025.
- [16] Magesh, Varun and Surani, Faiz and Dahl, Matthew and Suzgun, Mirac and Manning, Christopher D. and Ho, Daniel E. Hallucination-Free? Assessing the Reliability of Leading AI Legal Research Tools. In *Journal of Empirical Legal Studies*, 2025.
- [17] Shao, Peizhang and Xu, Linrui and Wang, Jinxi and Zhou, Wei and Wu, Xingyu. When Large Language Models Meet Law: Dual-Lens Taxonomy, Technical Advances, and Ethical Governance. In *arXiv preprint arXiv:2507.07748*, 2025.
- [18] Roberts, John G. Jr. 2023 Year-End Report on the Federal Judiciary. Supreme Court of the United States, 2023.
- [19] Engstrom, David Freeman and Ho, Daniel E. Algorithmic Accountability in the Administrative State. In *Yale Journal on Regulation*, 37:800–854, 2020.
- [20] Engstrom, David Freeman and Ho, Daniel E. and Sharkey, Catherine M. and Cuéllar, Mariano-Florentino. Government by Algorithm: Artificial Intelligence in Federal Administrative Agencies. Administrative Conference of the United States, 2020.
- [21] Solow-Niederman, Alicia. Administering Artificial Intelligence. In *Southern California Law Review*, 93(4):633–696, 2020.
- [22] Engel, Christoph and Grgić-Hlača, Nina. Machine Advice with a Warning about Machine Limitations: Experimentally Testing the Solution Mandated by the Wisconsin Supreme Court. In *Journal of Legal Analysis*, 13(1):284–340, 2021.
- [23] Barocas, Solon and Selbst, Andrew D. Big Data's Disparate Impact. In *California Law Review*, 104(3):671–732, 2016.
- [24] Ben-Shahar, Omri. Privacy Protection, At What Cost? Exploring the Regulatory Resistance to Data Technology in Auto Insurance. In *Journal of Legal Analysis*, 15(1):129–157, 2023.
- [25] King, Jennifer and Ho, Daniel and Gupta, Arushi and Wu, Victor and Webley-Brown, Helen. The Privacy-Bias Tradeoff: Data Minimization and Racial Disparity Assessments in U.S. Government. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pages 492–505. ACM, 2023.
- [26] Henderson, Peter and Hashimoto, Tatsunori and Lemley, Mark. Where's the Liability in Harmful AI Speech? In *Journal of Free Speech Law*, 3(2):589–650, 2023.
- [27] Lemley, Mark A. and Casey, Bryan. Remedies for Robots. In *The University of Chicago Law Review*, 86(5):1311–1396, 2019.
- [28] Volokh, Eugene. Large Libel Models? Liability for AI Output. In *Journal of Free Speech Law*, 3(2):489–558, 2023.
- [29] Chien, Colleen V. and Kim, Miriam and Akhil, Raj and Rathish, Rohit. How Generative AI Can Help Address the Access to Justice Gap Through the Courts. In *Loyola of Los Angeles Law Review*, Forthcoming, 2024.
- [30] Tribunal Constitucional de España. Nota Informativa N° 90/2024: La Sala Primera del TC por unanimidad sanciona a un abogado por la falta del debido respeto al tribunal. *Oficina de Prensa del Tribunal Constitucional*, 19 de septiembre de 2024.
- [31] Perlman, Andrew. The Implications of ChatGPT for Legal Services and Society. In *The Practice*, March/April 2023.
- [32] Tan, Jinzhe and Westermann, Hannes and Benyekhlef, Karim. ChatGPT as an Artificial Lawyer? In *Proceedings of the ICAIL 2023 Workshop on Artificial Intelligence for Access to Justice*. CEUR Workshop Proceedings, 2023.
- [33] Draper, Chris and Gillibrand, Nicky. The Potential for Jurisdictional Challenges to AI or LLM Training Datasets. In *Proceedings of the ICAIL 2023 Workshop on Artificial Intelligence for Access to Justice*. CEUR Workshop Proceedings, 2023.
- [34] Simshaw, Drew. Access to A.I. Justice: Avoiding an Inequitable Two-Tiered System of Legal Services. In *Yale Journal of Law & Technology*, 24:150–226, 2022.
- [35] Bar-Gill, Oren and Sunstein, Cass R and Talgam-Cohen, Inbal. Algorithmic Harm in Consumer Markets. In *Journal of Legal Analysis*, 15(1):1–47, 2023.

- [36] Gillis, Talia B. and Spiess, Jann L. Big Data and Discrimination. In *The University of Chicago Law Review*, 86(2):459–488, 2019.
- [37] Kleinberg, Jon and Ludwig, Jens and Mullainathan, Sendhil and Sunstein, Cass R. Discrimination in the Age of Algorithms. In *Journal of Legal Analysis*, 10(1):113–174, 2018.
- [38] Mayson, Sandra G. Bias In, Bias Out. In *The Yale Law Journal*, 128(8):2122–2473, 2019.
- [39] Bommasani, Rishi and Hudson, Drew A. and Adeli, Ehsan and Altman, Russ and Arora, Simran and von Arx, Sydney and Bernstein, Michael S. and Bohg, Jeannette and Bosselut, Antoine and Brunskill, Emma and others. On the Opportunities and Risks of Foundation Models. arXiv preprint arXiv:2108.07258, 2022.
- [40] Creel, Kathleen and Hellman, Deborah. The Algorithmic Leviathan: Arbitrariness, Fairness, and Opportunity in Algorithmic Decision-Making Systems. In *Canadian Journal of Philosophy*, 52(1):26–43, 2022.
- [41] Kleinberg, Jon and Raghavan, Manish. Algorithmic Monoculture and Social Welfare. In *Proceedings of the National Academy of Sciences*, 118(22), 2021.
- [42] Ji, Ziwei and Lee, Nayeon and Frieske, Rita and Yu, Tiezheng and Su, Dan and Xu, Yan and Ishii, Etsuko and Bang, Yejin and Madotto, Andrea and Fung, Pascale. Survey of Hallucination in Natural Language Generation. In *ACM Computing Surveys*, 55(12):1–38, 2023.
- [43] Zhang, Yue and Li, Yafu and Cui, Leyang and Cai, Deng and Liu, Lemaoyang and Fu, Ting and Huang, Xinting and Shi, Enbo and Wang, Yulong and Tan, Yulong and Gao, Liqun and He, Bang and Sun, Wei and Bi, Yongjing and Fu, You and Yuan, Furu and Zhang, Wei. Siren’s Song in the AI Ocean: A Survey on Hallucination in Large Language Models. arXiv preprint arXiv:2309.01219, 2023.
- [44] van Deemter, Kees. The Pitfalls of Defining Hallucination. In *Computational Linguistics*, Forthcoming, 2024.
- [45] Yiming Xu, Junfeng Jiao Evaluating Retrieval-Augmented Generation Strategies for Large Language Models in Travel Mode Choice Prediction In *arXiv preprint arXiv:2508.17527*, 2025.
- [46] Kalai, Adam Tauman and Vempala, Santosh S. Calibrated Language Models Must Hallucinate. arXiv preprint arXiv:2311.14648, 2023.
- [47] Xu, Ziwei and Jain, Sanjay and Kankanhalli, Mohan. Hallucination Is Inevitable: An Innate Limitation of Large Language Models. arXiv preprint arXiv:2401.11817, 2024.
- [48] Henderson, Peter and Krass, Mark S. and Zheng, Lucia and Guha, Neel and Manning, Christopher D. and Jurafsky, Dan and Ho, Daniel E. Pile of Law: Learning Responsible Data Filtering from the Law and a 256GB Open-Source Legal Dataset. arXiv preprint arXiv:2207.00220, 2022.
- [49] Tito, Joel. How AI Can Improve Access to Justice. Centre for Public Impact, 2017.
- [50] Sharma, Mrinank and Tong, Meg and Korbak, Tomasz and Duvenaud, David and Askill, Amanda and Bowman, Samuel R. and Cheng, Newton and Durmus, Esin and Dodds, Zac Hatfield and Johnston, Scott R. and others. Towards Understanding Sycophancy in Language Models. arXiv preprint arXiv:2310.13548, 2023.
- [51] Wei, Jerry and Huang, Da and Lu, Yifeng and Zhou, Denny and Le, Quoc V. Simple Synthetic Data Reduces Sycophancy in Large Language Models. arXiv preprint arXiv:2308.03958, 2023.
- [52] Jones, Erik and Steinhardt, Jacob. Capturing Failures of Large Language Models via Human Cognitive Biases. In *Advances in Neural Information Processing Systems*, 35:11411–11426, 2022.
- [53] Suri, Gaurav and Slater, Lily R. and Ziaee, Ali and Nguyen, Morgan. Do Large Language Models Show Decision Heuristics Similar to Humans? A Case Study Using GPT-3.5. arXiv preprint arXiv:2305.04400, 2023.
- [54] Azaria, Amos and Mitchell, Tom. The Internal State of an LLM Knows When It’s Lying. arXiv preprint arXiv:2304.13734, 2023.
- [55] Kadavath, Saurav and Conerly, Tom and Askill, Amanda and Henighan, Tom and Drain, Dawn and Perez, Ethan and Schiefer, Nicholas and Hatfield-Dodds, Zac and Maxwell, Jackson Kernion and others. Language Models (Mostly) Know What They Know. arXiv preprint arXiv:2207.05221, 2022.
- [56] Tian, Katherine and Mitchell, Eric and Zhou, Allan and Sharma, Archit and Rafailov, Rafael and Yao, Huaxiu and Finn, Chelsea and Manning, Christopher D. Just Ask for Calibration: Strategies for Eliciting Calibrated Confidence Scores from Language Models Fine-Tuned with Human Feedback. arXiv preprint arXiv:2305.14975, 2023.
- [57] Xiong, Miao and Hu, Zhiyuan and Lu, Xinyang and Li, Yifei and Fu, Jie and He, Junxian and Hooi, Bryan. Can LLMs Express Their Uncertainty? An Empirical Evaluation of Confidence Elicitation in LLMs. arXiv preprint arXiv:2306.13063, 2023.
- [58] Yin, Zhangyue and Sun, Qiushi and Guo, Qipeng and Wu, Jiawen and Qiu, Xipeng and Huang, Xuanjing. Do Large Language Models Know What They Don’t Know? arXiv preprint arXiv:2305.18153, 2023.

- [59] Zhang, Yunfeng and Liao, Q. Vera and Bellamy, Rachel K. E. Effect of Confidence and Explanation on Accuracy and Trust Calibration in AI-assisted Decision Making. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 295–305, 2020.
- [60] Shuster, Kurt and Poff, Spencer and Chen, Moya and Kiela, Douwe and Weston, Jason. Retrieval Augmentation Reduces Hallucination in Conversation. arXiv preprint arXiv:2104.07567, 2021.
- [61] Peng, Baolin and Galley, Michel and He, Pengcheng and Cheng, Hao and Xie, Yujia and Hu, Yu and Huang, Qiuyuan and Liden, Lars and Yu, Zhou and Chen, Weizhu and Gao, Jianfeng. Check Your Facts and Try Again: Improving Large Language Models with External Knowledge and Automated Feedback. arXiv preprint arXiv:2302.12813, 2023.
- [62] Si, Chenglei and Gan, Zhe and Yang, Zhengyuan and Wang, Shuohang and Wang, Jianfeng and Boyd-Graber, Jordan and Wang, Lijuan. Prompting GPT-3 To Be Reliable. Eleventh International Conference on Learning Representations, 2023.
- [63] Lei, Deren and Li, Yaxi and Wang, Mingyu and Yun, Vincent and Ching, Emily and Kamal, Eslam and Liu, Yaqing and Liu, Wen-Ding and Yang, Ellen and Liu, Daniel. Chain of Natural Language Inference for Reducing Large Language Model Ungrounded Hallucinations. arXiv preprint arXiv:2310.03951, 2023.
- [64] Suzgun, Mirac and Kalai, Adam Tauman. Meta-prompting: Enhancing Language Models with Task-agnostic Scaffolding. arXiv preprint arXiv:2401.12954, 2024.
- [65] Tian, Katherine and Mitchell, Eric and Yao, Huaxiu and Manning, Christopher D. and Finn, Chelsea. Fine-Tuning Language Models for Factuality. arXiv preprint arXiv:2311.08401, 2023.
- [66] Razumovskaia, Evgeniia and Vulić, Ivan and Marković, Pavle and Cichy, Tomasz and Zheng, Qian and Wen, Tsung-Hsien and Budzianowski, Paweł. Dial BeInfo for Faithfulness: Improving Factuality of Information-Seeking Dialogue via Behavioural Fine-Tuning. arXiv preprint arXiv:2311.09800, 2023.
- [67] Zhang, Hanning and Diao, Shizhe and Lin, Yong and Fung, Yi R and Lian, Qing and Wang, Xingyao and Chen, Yangyi and Ji, Heng and Zhang, Tong. R-Tuning: Teaching Large Language Models to Refuse Unknown Questions. arXiv preprint arXiv:2311.09677, 2023.
- [68] Shi, Weijia and Han, Xiaochuang and Lewis, Mike and Tsvetkov, Yulia and Zettlemoyer, Luke and Yih, Scott Wenta. Trusting Your Evidence: Hallucinate Less with Context-aware Decoding. arXiv preprint arXiv:2305.14739, 2023.
- [69] Mallen, Alex and Asai, Akari and Zhong, Victor and Das, Rajarshi and Khashabi, Daniel and Hajishirzi, Hannaneh. When Not to Trust Language Models: Investigating Effectiveness of Parametric and Non-Parametric Memories. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9802–9822, 2023.
- [70] Li, Kenneth and Patel, Oam and Viégas, Fernanda and Pfister, Hanspeter and Wattenberg, Martin. Inference-time Intervention: Eliciting Truthful Answers from a Language model. 2024.
- [71] Chuang, Yung-Sung and Xie, Yujia and Luo, Hongyin and Kim, Yoon and Glass, James R. and He, Pengcheng. DoLa: Decoding by Contrasting Layers Improves Factuality in Large Language Models. Twelfth International Conference on Learning Representations, 2024.
- [72] Chern, I and Chern, Steffi and Chen, Shiqi and Yuan, Weizhe and Feng, Kehua and Zhou, Chunting and He, Junxian and Neubig, Graham and Liu, Pengfei and others. FacTool: Factuality Detection in Generative AI—A Tool Augmented Framework for Multi-Task and Multi-Domain Scenarios. arXiv preprint arXiv:2307.13528, 2023.
- [73] Qin, Yujia and Hu, Shengding and Lin, Yankai and Chen, Weize and Ding, Ning and Cui, Ganqu and Zeng, Zheni and Huang, Yufei and Xiao, Chaojun and Han, Chi and others. Tool Learning with Foundation Models. arXiv preprint arXiv:2304.08354, 2023.
- [74] Gou, Zhibin and Shao, Zhihong and Gong, Yeyun and shen, yelong and Yang, Yujiu and Duan, Nan and Chen, Weizhu. CRITIC: Large Language Models Can Self-Correct with Tool-Interactive Critiquing. Twelfth International Conference on Learning Representations, 2024.
- [75] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, Jie Tang. AgentBench: Evaluating LLMs as Agents. arXiv:2308.03688
- [76] Tonmoy, SM and Zaman, SM and Jain, Vinija and Rani, Anku and Rawte, Vipula and Chadha, Aman and Das, Amitava. A Comprehensive Survey of Hallucination Mitigation Techniques in Large Language Models. arXiv preprint arXiv:2401.01313, 2024.

- [77] Magesh, Varun and Surani, Faiz and Dahl, Matthew and Suzgun, Mirac and Manning, Christopher D. and Ho, Daniel E. Hallucination-Free? Assessing the Reliability of Leading AI Legal Research Tools. arXiv preprint arXiv:2405.20362, 2024.
- [78] Casetext. GPT-4 alone is not a reliable legal solution—but it does enable one: CoCounsel harnesses GPT-4’s power to deliver results that legal professionals can rely on. Casetext Blog, 2023.
- [79] Thomson Reuters. Introducing AI-Assisted Research: Legal research meets generative AI. Press Release, 2023.
- [80] LexisNexis. LexisNexis Launches Lexis+ AI, a Generative AI Solution with Linked Hallucination-Free Legal Citations. Press Release, 2023.
- [81] Weiser, Benjamin and Bromwich, Jonah E. Lawyer uses ChatGPT in brief, gets called out for ‘bogus’ case citations. In *The New York Times*, May 2023.
- [82] Kite-Jackson, Darla Wynon. Recent Integration of Large Language Models (LLMs) into Research and Writing Tools Presents Both Unprecedented Opportunities and Significant Challenges. In *2023 Artificial Intelligence (AI) TechReport*. American Bar Association, 2023.
- [83] Wellen, Serena. How Lexis+ AI Delivers Hallucination-Free Linked Legal Citations. LexisNexis Blog, Feb 2024.
- [84] Wellen, Serena. Tech Innovation with LLMs Producing More Secure and Reliable Gen AI Results. LexisNexis Blog, May 2024.
- [85] Thomson Reuters. Introducing Ask Practical Law AI on Practical Law: Generative AI meets legal how-to. Product Information, 2024.
- [86] Ambrogi, Bob. LawNext: Thomson Reuters’ AI Strategy for Legal, with Mike Dahn, Head of Westlaw, and Joel Hron, Head of AI. LawNext Podcast, Feb 2024.
- [87] Miner, Roger J. Remarks: Clerks of Judge Luther A. Wilgarten, Jr. 1989.
- [88] Goddard, K and Roudsari, A and Wyatt, JC. Automation bias: a systematic review of frequency, effect mediators, and mitigators. In *Journal of the American Medical Informatics Association*, 19(1):121–127, 2012.
- [89] Belkin, Nicholas J. Some (what) grand challenges for information retrieval. In *ACM SIGIR Forum*, 42(2):47–54. ACM New York, NY, USA, 2008.
- [90] Mik, Eliza. Caveat Lector: Large Language Models in Legal Practice. In *Artificial Intelligence and Law*, Forthcoming or preprint status, 2024.
- [91] Arewa, Olufunmilayo B. Open Access in a Closed Universe: Lexis, Westlaw, Law Schools, and the Legal Information Market. In *Lewis & Clark Law Review*, 10(4):797–840, 2006.
- [92] Thomson Reuters. Westlaw tip of the week: Checking cases with keycite. 2019.
- [93] Schwarcz, Daniel and Manning, Sam and Barry, Patrick and Cleveland, David R. and Prescott, JJ and Rich, Beverly. AI-POWERED LAWYERING: AI REASONING MODELS, RETRIEVAL AUGMENTED GENERATION, AND THE FUTURE OF LEGAL PRACTICE. 2024.
- [94] Garg, Aksh and Ma, Megan. Opportunities and Challenges in Legal AI. Stanford Law School, Jan 2025.
- [95] Microsoft. Generative AI for Lawyers. 2024.
- [96] Susskind, Richard and Susskind, Daniel E. *Tomorrow’s Lawyers: An Introduction to Your Future*. Oxford University Press, 2023.
- [97] Brescia, Raymond H. What’s a Lawyer For?: Artificial Intelligence and Third-Wave Lawyering. In *Florida State University Law Review*, 51:542, 2024.
- [98] Armour, John and Parnham, Richard and Sako, Mari. Augmented Lawyering. In *University of Illinois Law Review*, pages 71–112, 2022.
- [99] Harvey. Harvey Raises \$100M Series C from Google Ventures, OpenAI, Kleiner Perkins, Sequoia Capital, Elad Gil, and SV Angel at a \$1.5B valuation. Press Release, July 2024.
- [100] LexisNexis. LexisNexis Introduces Protégé Personalized AI Assistant with Agentic AI, Making it Easier to Power Complex Legal Task Completion. Press Release, Jan 2025.
- [101] Thomson Reuters. Get to Know Thomson Reuters: Our Technology Journey and What’s Next. Press Release, Jan 2025.
- [102] Strom, Roy. Big Law Is Questioning the ‘Magical Thinking’ of AI as Savior. Bloomberg Law, Aug 2024.
- [103] Kim, Miriam and Chien, Colleen V. Generative AI and Legal Aid: Results from a Field Study and 100 Use Cases to Bridge the Access to Justice Gap. In *Loyola of Los Angeles Law Review*, 57:903–904, 2025.

- [104] Re, Richard M. Artificial Authorship and Judicial Opinions. In *George Washington Law Review*, 92:1558–1559, 2024.
- [105] Liu, John Zhuang and Li, Xueyao. How Do Judges Use Large Language Models? Evidence From Shenzhen. In *Journal of Legal Analysis*,