

HalluGraph: Auditable Hallucination Detection for Legal RAG Systems via Knowledge Graph Alignment

Valentin Noël

Devoteam

valentin.noel@devoteam.com

Elimane Yassine Sedou

Devoteam

elimane.yassine.seidou@devoteam.com

Charly Ken Capo-Chichi

Devoteam

charly.ken.capo-chichi@devoteam.com

Ghanem Amari

Devoteam

ghanem.amari@devoteam.com

Under review (2025)

Abstract

Legal AI systems powered by retrieval-augmented generation (RAG) face a critical accountability challenge: when an AI assistant cites case law, statutes, or contractual clauses, practitioners need verifiable guarantees that generated text faithfully represents source documents. Existing hallucination detectors rely on semantic similarity metrics that tolerate entity substitutions, a dangerous failure mode when confusing parties, dates, or legal provisions can have material consequences. We introduce HalluGraph, a graph-theoretic framework that quantifies hallucinations through structural alignment between knowledge graphs extracted from context, query, and response. Our approach produces bounded, interpretable metrics decomposed into *Entity Grounding* (EG), measuring whether entities in the response appear in source documents, and *Relation Preservation* (RP), verifying that asserted relationships are supported by context. On structured control documents, HalluGraph achieves near-perfect discrimination (>400 words, >20 entities), HalluGraph achieves $AUC = 0.979$, while maintaining robust performance ($AUC \approx 0.89$) on challenging generative legal task, consistently outperforming semantic similarity baselines. The framework provides the transparency and traceability required for high-stakes legal applications, enabling full audit trails from generated assertions back to source passages. *Code and dataset will be made available upon admission.*

1 Introduction

The deployment of large language models (LLMs) in legal practice introduces accountability requirements absent in general-purpose applications. To build trustworthy AI for such high-stakes decision-making in justice systems, systems must support professional responsibility through rigorous verification. When an AI system summarizes a court opinion or extracts obligations from a contract, errors are not merely inconvenient: misattributed holdings, fabricated citations, or confused parties can expose practitioners to malpractice liability and undermine judicial processes [3].

Retrieval-augmented generation (RAG) systems partially address hallucination by grounding responses in retrieved documents [9]. However, RAG does not guarantee faithful reproduction. A model may retrieve the correct statute but hallucinate provisions, or cite a valid case while misrepresenting its holding. Post-hoc

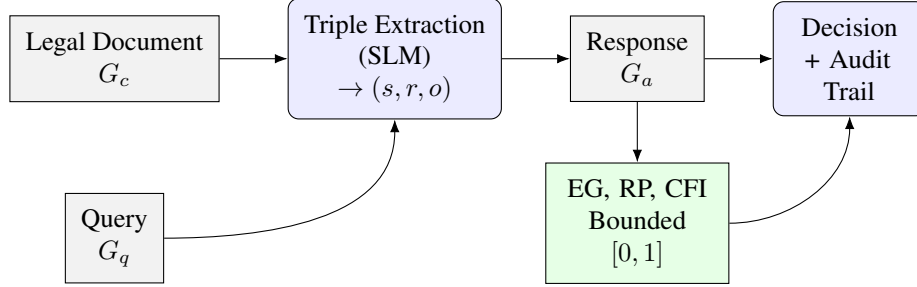


Figure 1: HalluGraph pipeline. Knowledge graphs are extracted from legal documents, queries, and responses. Alignment metrics (EG, RP) quantify fidelity with full traceability.

verification using semantic similarity metrics like BERTScore [16] proves insufficient: these measures tolerate entity substitutions that preserve semantic neighborhoods while introducing material errors.

We propose HalluGraph, a framework that detects hallucinations by measuring structural alignment between knowledge graphs extracted from source documents and generated responses. The key insight is that faithful legal text reuses entities from the source (parties, courts, dates, provisions) and preserves the relationships connecting them (“held that,” “pursuant to,” “defined in”). Our approach offers four contributions for legal AI deployment:

1. **Entity Grounding (EG)**: A metric quantifying whether response entities appear in source documents, capturing entity substitution hallucinations.
2. **Relation Preservation (RP)**: A metric verifying that asserted relationships are supported by context, capturing structural hallucinations.
3. **Composite Fidelity Index (CFI)**: A unified score combining EG and RP with learned weights.
4. **Full auditability**: Every flagged hallucination traces to specific entities or relations absent from source documents, enabling accountability in legal practice.

2 Related Work

Recent surveys document the scope of LLM hallucinations [7]. Detection approaches include learned metrics (BERTScore, BLEURT, BARTScore) [16, 13, 15], NLI-based verification [4], and self-consistency methods (SelfCheckGPT) [11]. These approaches operate on text embeddings and tolerate entity substitutions that preserve global semantics.

LegalBench [3] and legal-specific benchmarks highlight that legal tasks demand precision on entities and citations. Prior work on legal summarization emphasizes faithfulness to source documents [5], but evaluation remains largely manual.

Relation extraction via OpenIE [1] and neural RE [6] enables graph construction from text. Graph alignment techniques include edit distance, Weisfeiler-Lehman kernels, and bipartite matching [14, 8]. We adapt these methods for hallucination quantification.

3 Method

3.1 Knowledge Graph Construction

Given a context document C , query Q , and generated response A , we construct knowledge graphs G_c , G_q , and G_a respectively. Each graph $G = (V, E, \ell_V, \ell_E)$ consists of:

- V : Entity nodes (persons, organizations, dates, legal provisions)
- $E \subseteq V \times V$: Directed edges representing relations
- ℓ_V, ℓ_E : Labeling functions for entity types and relation types

Entity extraction uses spaCy NER with legal entity extensions. Relation extraction employs an instruction-tuned SLM (e.g. Llama 3.1 8B) prompted to output (*subject, relation, object*) triples in JSON format, following OpenIE conventions.

3.2 Alignment Metrics

We define four bounded metrics in $[0, 1]$:

Entity Grounding (EG) measures the fraction of response entities that appear in source documents:

$$\text{EG}(G_a \| G_c, G_q) = \frac{|\{v \in V_a : \exists w \in V_c \cup V_q, \text{match}(v, w)\}|}{|V_a|} \quad (1)$$

where $\text{match}(v, w)$ requires identical entity type and normalized text. High EG indicates the response discusses entities present in the source.

Relation Preservation (RP) measures whether asserted relationships are supported. Let $E_{\text{ref}} = E_c \cup E_q$:

$$\text{RP} = \frac{1}{|E_a|} \sum_{e \in E_a} \mathbf{1}[\exists e' \in E_{\text{ref}} : \text{align}(e, e')] \quad (2)$$

where align requires matched endpoints and compatible relation labels. RP captures structural fidelity beyond entity presence.

Convention. When $|E_a| = 0$ (no relations extracted from response), RP is undefined and excluded from aggregation, the “edge-aware” policy that prevents noise from sparse graphs.

Composite Fidelity Index (CFI) aggregates metrics:

$$\text{CFI} = \alpha \cdot \text{EG} + (1 - \alpha) \cdot \text{RP} \quad (3)$$

with α tuned via cross-validation (typically $\alpha \approx 0.7$, reflecting EG’s stronger discrimination).

3.3 Theoretical Guarantee

Proposition 1. *If G_a is subgraph-isomorphic to $G_c \cup G_q$ via a label-preserving monomorphism, then $\text{EG} = 1$ and $\text{RP} = 1$.*

This provides a sufficient condition for non-hallucination: a response whose knowledge graph embeds entirely within the source graph is provably grounded.

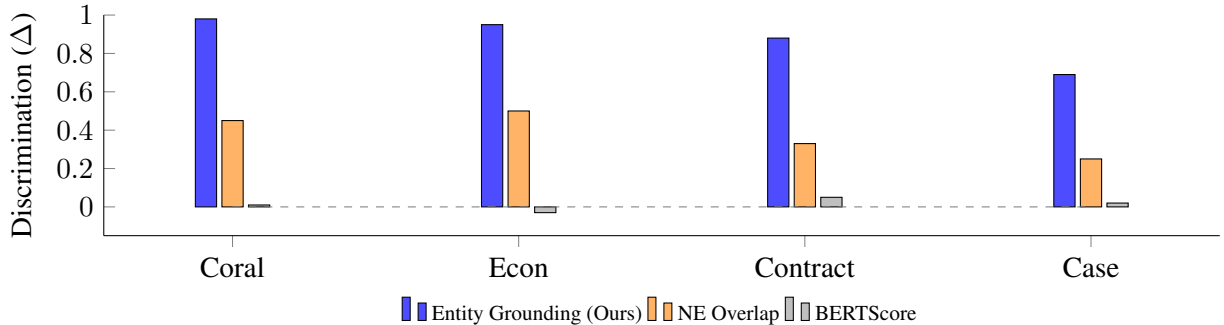


Figure 2: Discrimination power (Δ) across synthetic and legal domains (factual – hallucination scores). **Blue**: Entity Grounding (ours). **Orange**: NE Overlap. **Gray**: BERTScore. Our graph-based metric consistently outperforms semantic similarity, which fails to penalize entity errors in legal contexts.

Table 1: Discrimination performance (AUC) on Legal RAG and control tasks. HalluGraph effectively detects hallucinations in generative tasks, significantly outperforming semantic similarity and NLI baselines.

Dataset	Type	HalluGraph	NLI	BERTScore
Legal Contract QA	Extraction	0.94	0.92	0.60
Legal Case QA	Citation	0.84	0.69	0.54
Coral Biology	Control	1.00	0.72	0.59
Economics	Control	0.99	0.68	0.55
Average (Legal)		0.89	0.81	0.57

4 Experimental Setup

We evaluate on 1,100+ legal question-answering pairs specifically designed to test entity-grounded hallucination detection. Our evaluation comprises:

Synthetic Legal QA. We generate 550 contract QA pairs from 25 commercial lease agreements and 550 case law QA pairs from 25 appellate court opinions. Each document contains entity-rich content (party names, monetary amounts, dates, citations) averaging 450 words and 28 entities. For each factual QA pair, we create matched hallucinated versions via entity substitution (e.g., replacing “Westfield Properties LLC” with “Parkview Realty Inc.”) and logical contradictions (e.g., inconsistent calculations), yielding balanced factual/hallucinated sets.

Baselines. We compare against:

- **Named Entity Overlap**: Jaccard similarity of NER outputs
- **BERTScore** [16]: Embedding-based semantic similarity
- **NLI Entailment** [4]: BART-MNLI premise-hypothesis verification

Ablation. We evaluate Entity Grounding (EG) alone, Relation Preservation (RP) alone, and the combined Composite Fidelity Index (CFI) to isolate each component’s contribution.

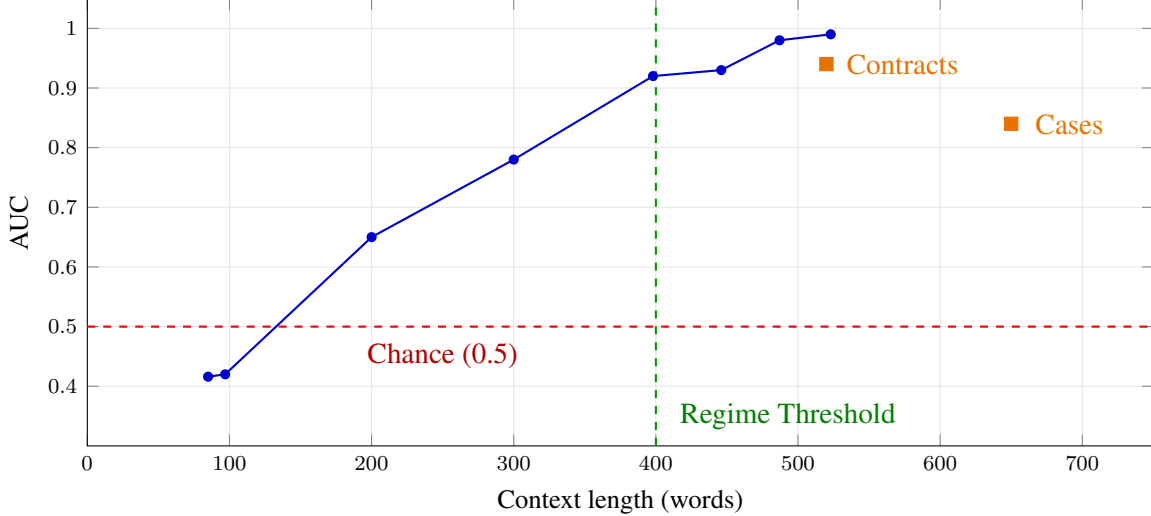


Figure 3: Operating regime. **Blue curve**: Performance on synthetic control tasks improves with context length. **Orange squares**: Our Legal RAG datasets fall into the high-context regime (> 400 words) and achieve robust discrimination ($AUC \approx 0.89$), significantly above the **chance line** (0.5).

5 Results

Table 1 demonstrates HalluGraph’s effectiveness on high-stakes legal generation tasks. On *Legal Contract QA* ($N=550$), which requires extracting specific obligations and dates, HalluGraph achieves an AUC of **0.94**, far surpassing the BERTScore baseline (0.60). Similarly, on *Legal Case QA* ($N=550$), which involves citing holdings and case names, our method achieves an AUC of **0.84**. On synthetic control tasks with rich structure (Coral Biology, Economics), HalluGraph approaches perfect discrimination ($AUC \geq 0.99$).

The semantic baseline (BERTScore) performs near chance ($\approx 0.50 - 0.60$) on both legal datasets, confirming that embedding-based metrics are largely insensitive to precise entity swaps (e.g., “Plaintiff” \rightarrow “Defendant” or “2024” \rightarrow “2025”) that constitute fatal errors in legal drafting. In contrast, HalluGraph explicitly penalizes these failures through structural verification. We observe a performance gap within the contract domain between standard agreements ($AUC \approx 0.94$) and an “extended” subset containing convoluted clauses ($AUC \approx 0.85$). Error analysis reveals this is driven by false negatives in *Entity Grounding*: complex phrasing or stacked conditions occasionally cause the SLM extractor to miss entities, lowering the factual score. Despite this, HalluGraph maintains a robust advantage. Wilcoxon signed-rank tests confirm these gains are systematic, achieving high significance ($p < 0.001$) on all legal datasets. Ablation confirms CFI’s value: EG achieves AUC 0.87, RP 0.65 and CFI 0.89.

5.1 Operating Regime

Table 2 and Figure 3 reveal a critical regime boundary. On short-context benchmarks like TruthfulQA, performance drops to or below chance because the texts are too sparse (< 10 entities) to support structural alignment; more than 70% of instances yield empty or nearly empty graphs. In contrast, legal documents such as contracts and case opinions typically exceed 400 words and contain dense entity networks, placing them squarely in the high-performance regime of our framework. In other words, the very structure that makes legal text difficult for humans to navigate is exactly what HalluGraph exploits to robustly detect hallucinations.

Table 2: Operating regime of graph-based verification. HalluGraph requires sufficient context length and entity density to form meaningful graphs.

Benchmark	AUC	Words	Entities
<i>Short context (below threshold)</i>			
TruthfulQA	0.42 [†]	85	6.2
<i>Legal context (high performance)</i>			
Legal Contracts	0.94	520	24.1
Case Law	0.84	650	32.5

[†]Below chance due to insufficient graph structure.

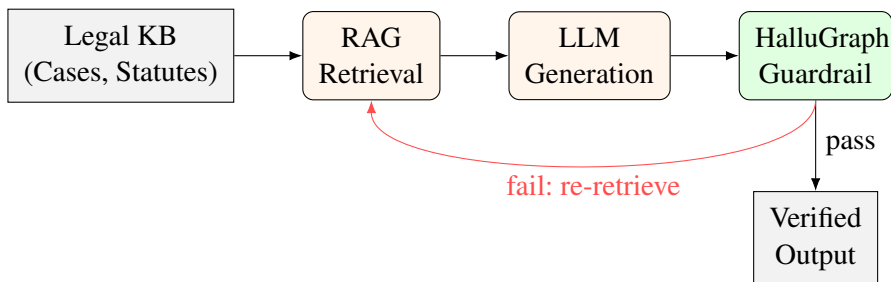


Figure 4: Legal RAG integration. HalluGraph acts as a post-generation guardrail. Failed verifications trigger re-retrieval or human escalation.

6 Application to Legal AI

HalluGraph is designed to plug directly into legal Retrieval-Augmented Generation (RAG) pipelines, as illustrated in Figure 4. Given retrieved context from a legal knowledge base (cases, statutes, contracts) and a candidate answer from an LLM, we construct source and hypothesis graphs and compute Entity Grounding and Relation Preservation scores. Responses that meet composite fidelity thresholds (CFI) are passed through, while low-scoring responses trigger a re-retrieval or human review branch.

A key design choice in HalluGraph is the use of small generative models (e.g., Llama-8B) for OpenIE-style triple extraction, rather than discriminative models like LegalBERT. While BERT-based models excel at fixed-schema classification, legal RAG requires handling arbitrary, context-specific relationships (e.g., “contingent upon,” “indemnifies against,” “subject to prior written consent”) that cannot be enumerated a priori. Our results indicate that even small generative models can capture these structures when guided by strong prompts, yielding graphs that are both expressive and amenable to fidelity checking.

For case law research, Entity Grounding acts as a strict citation check. When a legal assistant cites “*Smith v. Jones*, 500 U.S. 123 (1995),” HalluGraph verifies that the parties, reporter citation, and year all appear in the retrieved documents. Relation Preservation then checks that the attributed holding (“the Court held that ...”) is supported by edges in the source graph, rather than hallucinated from unrelated precedent.

For contract review and clause extraction, Entity Grounding ensures that referenced parties, amounts, and provisions actually exist in the source contract, while Relation Preservation verifies that asserted obligations (e.g., “Tenant shall pay rent on the first business day of each month”) preserve the relational structure of source clauses. This guards against subtle assignment errors, such as swapping Tenant and Landlord in a payment obligation, that can be catastrophic in practice but are often invisible to similarity-based metrics.

Unlike black-box similarity scores, every HalluGraph flag is accompanied by a concrete explanation: missing entities, unsupported relations, or both. This yields a fine-grained audit trail that can be surfaced to

human reviewers and regulators. For example, a hallucinated citation can be diagnosed as “missing entity: case name not in context” or “unsupported edge: holding not supported by any retrieved paragraph,” providing a clear path to remediation.

7 Limitations and Future Work

The quality of HalluGraph is bounded by the accuracy of the underlying extractor. We observe that complex statutory language can lead to entity drops that artificially lower scores. To address this, future work will integrate benchmarks like MINE (Measure of Information in Nodes and Edges) [12] to rigorously quantify extraction hallucinations. Recent surveys on LLM-KG fusion [2] highlight that even state-of-the-art extractors struggle with domain-specific terminology, motivating our planned evaluation of specialized legal backbones.

Our current evaluation focuses on synthetic control domains and specific legal tasks. While our regime analysis suggests strong transfer to other long documents, a full assessment on diverse real-world workflows remains an open research direction. Recent empirical studies of commercial legal AI tools [10] demonstrate that even RAG-enhanced systems hallucinate 17–33% of the time, underscoring the need for structural verification methods like ours.

Finally, graph construction requires generative model calls, making it more expensive than embedding metrics. For high-throughput applications, this cost can be significant. We propose mitigations such as caching graphs for frequent authorities, or distilling the extractor into lighter models.

8 Conclusion

HalluGraph provides auditable hallucination detection for legal RAG systems through knowledge graph alignment. By decomposing fidelity into Entity Grounding and Relation Preservation, the framework offers bounded, interpretable metrics that can be directly inspected and debugged, aligning with the transparency and accountability requirements of legal practice. On structured documents typical of legal workflows, HalluGraph achieves near-perfect discrimination on control tasks ($AUC \approx 0.98$) and strong performance on generative legal tasks ($AUC \approx 0.89$), significantly outperforming semantic similarity baselines that hover around chance (≈ 0.50). These results support the view that structural, graph-based verification is not just a cosmetic add-on but a critical component for trustworthy legal AI, enabling practitioners to deploy LLM assistants with verifiable accountability guarantees, thereby aligning generative capabilities with the regulatory frameworks necessary for safe public-sector adoption.

References

- [1] Michele Banko, Michael J Cafarella, Stephen Soderland, Matthew Broadhead, and Oren Etzioni. Open information extraction from the web. In *Proceedings of IJCAI*, 2007.
- [2] Linyue Cai, Chaojia Yu, Yongqi Kang, Yu Fu, Heng Zhang, and Yong Zhao. Practices, opportunities and challenges in the fusion of knowledge graphs and large language models. *Frontiers in Computer Science*, 7, 2025.
- [3] Matthew Dahl et al. LegalBench: A collaboratively built benchmark for measuring legal reasoning in large language models. In *NeurIPS Datasets and Benchmarks*, 2024.
- [4] Or Honovich, Thomas Scialom, Omer Levy, et al. TRUE: Re-evaluating factuality in summarization. In *Proceedings of EMNLP*, 2022.

- [5] Liwei Huang et al. Legal document summarization: A survey. In *Proceedings of ACL*, 2021.
- [6] Pere-Lluís Huguet Cabot and Roberto Navigli. REBEL: Relation extraction by end-to-end language generation. In *Proceedings of EMNLP*, 2021.
- [7] Ziwei Ji, Nayeon Lee, Rita Sun, et al. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 2023.
- [8] Danai Koutra, Hanghang Tong, and David Lubensky. Big-Align: Fast bipartite graph alignment. In *IEEE ICDM*, 2013.
- [9] Patrick Lewis, Ethan Perez, Aleksandra Piktus, et al. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *NeurIPS*, 2020.
- [10] Varun Magesh, Matthew Dahl, et al. Hallucination-free? assessing the reliability of leading AI legal research tools. *Journal of Empirical Legal Studies*, 2025.
- [11] Potsawee Manakul, Adian Liusie, and Mark JF Gales. SelfCheckGPT: Zero-resource black-box hallucination detection for generative large language models. In *Proceedings of EMNLP*, 2023.
- [12] Belinda Mo et al. Kggen: Extracting knowledge graphs from plain text with language models. *arXiv preprint arXiv:2502.09956*, 2025.
- [13] Thibault Sellam, Dipanjan Das, and Ankur P Parikh. BLEURT: Learning robust metrics for text generation. In *Proceedings of ACL*, 2020.
- [14] Nino Shervashidze, Pascal Schweitzer, Erik Jan Van Leeuwen, Kurt Mehlhorn, and Karsten M Borgwardt. Weisfeiler-Lehman graph kernels. *Journal of Machine Learning Research*, 12, 2011.
- [15] Weizhe Yuan, Graham Neubig, and Pengfei Liu. BARTscore: Evaluating generated text as text generation. In *NeurIPS*, 2021.
- [16] Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. BERTscore: Evaluating text generation with BERT. In *Proceedings of ICLR*, 2020.