# Precedent-Aware Multi-Agent Retrieval-Augmented Generation in Case Law Analysis

Shatrunjay Kumar Singh[1]

[1]Bloomberg LP

Publication Date: 2026/01/07

**Abstract:** Retrieval-Augmented Generation (RAG) systems promise practical legal assistance by grounding Large Language Models (LLMs) in external authority. However, standard RAG optimizes semantic similarity and often fails to respect common-law constraints such as jurisdictional bindingness, court hierarchy, temporal validity, and negative treatment. We propose Precedent- Aware Multi-Agent RAG (PA-MA-RAG), an agentic architecture that decomposes legal research and writing into specialized agents for issue framing, authority planning, retrieval, precedent ranking, conflict resolution, drafting, and citation verification. Our method introduces an authority- constrained re-ranking objective that prioritizes controlling precedents while penalizing overruled or otherwise negatively treated cases. The verifier agent enforces evidence-grounded generation by requiring each legal proposition to be supported by retrieved holdings and quotations. We describe an evaluation protocol for both precedent retrieval and citation-grounded legal analysis generation, including authority correctness, supported-claim rate, and robustness to conflicting precedent.

*Keywords: Precedent-Aware RAG, Multi-Agent Systems, Legal Information Retrieval, Stare Decisis, Authority Ranking, Citation Networks, CLERC, COLIEE.*

## I. INTRODUCTION

Legal analysis in common-law jurisdictions is precedent-driven: arguments must cite controlling authority, distinguish unfavorable cases, and avoid reliance on invalid or negatively treated precedent. While recent RAG methods improve factuality by retrieving external passages, legal tasks impose additional constraints that are not captured by relevance alone, including jurisdiction, court hierarchy, and subsequent history. As a result, naive RAG can surface persuasive but non- binding cases, miss controlling decisions, or synthesize unsupported claims despite including citations.

This paper introduces Precedent-Aware Multi-Agent RAG (PA-MA-RAG), which treats precedent selection as a structured decision problem and uses a coordinated set of agents to (i) interpret the legal query, (ii) retrieve candidate authorities, (iii) rank and filter them by doctrinal authority, and (iv) generate and verify a citation-grounded legal analysis. Our design aligns retrieval and generation with stare decisis by making authority and validity explicit signals.

➢ *Our Contributions are:*

- Multi-agent legal workflow: We specify an agent decomposition for precedent-centered legal research and writing, including explicit conflict checking and citation verification.
- Authority-constrained ranking: We formalize a precedent authority score that combines bindingness, jurisdictional match, temporal validity, citation-graph centrality, and negative- treatment penalties.
- Evaluation protocol: We propose metrics that test not only relevance but also authority correctness, supported-claim rate, and robustness to conflicting or overruled precedents.

## II. BACKGROUND: THE ANATOMY OF LEGAL PRECEDENT

➢ *The Doctrine of Stare Decisis*

Stare decisis structures common-law reasoning by requiring courts to follow binding precedent from superior courts within the same jurisdiction and, under some conditions, their own prior decisions. Binding force depends on institutional hierarchy (vertical stare decisis) and the court's willingness to adhere to prior decisions (horizontal

stare decisis). For legal assistance systems, the key implication is that two semantically similar cases can carry very different legal weight depending on jurisdiction and court level.

➤ *Hierarchy*

The architecture of court hierarchy forms a pyramid, with trial courts (District Courts) at the base, handling initial facts and law application; intermediate Courts of Appeals (Circuit Courts) in the middle, reviewing trial decisions; and the Supreme Court at the apex, as the final authority for federal matters, ensuring a system of review for legal correctness and consistent interpretation, both federally and within each state system. This structure allows for specialization (trial vs. appeal) and efficient resource allocation, with appellate courts focusing on legal precedent rather than re-trying facts.
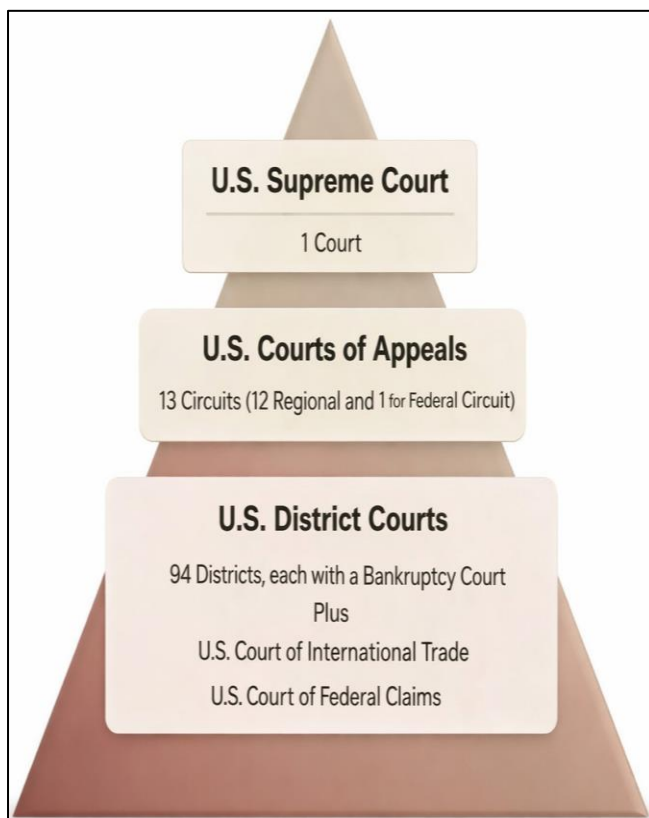


Fig 1 Hierarchy

➤ *Subsequent History and Negative Treatment*

Precedents can lose force through reversal, overruling, abrogation, statutory change, or narrower limitation. Legal researchers therefore track subsequent history and treatment signals (e.g., whether later cases follow, distinguish, or criticize a decision). A precedent-aware retrieval system should explicitly model these signals to avoid citing 'bad law' or relying on outdated doctrine.

➤ *The Language of Citations: Metadata as Legal Signal*

Legal citations provide structured metadata that identifies a case, its court, date, and reporter location. Beyond identification, citations enable a network view of doctrine through inter-case references. These metadata

signals are crucial inputs for precedent-aware ranking.

➤ *Why Standard RAG Fails in Legal Settings*

Standard RAG pipelines retrieve passages by lexical or dense similarity and then prompt an LLM to synthesize an answer. In legal analysis, this can surface persuasive but non-binding cases, miss controlling authorities, or generate plausible yet unsupported statements. The root cause is misalignment between retrieval objectives (semantic relevance) and legal validity constraints (authority and subsequent history).

## III. PRECEDENT-AWARE MULTI-AGENT RAG (PA-MA-RAG)

➤ *Overview*

PA-MA-RAG decomposes precedent-centered legal reasoning into a set of agents that iteratively retrieve, rank, and verify authority. The design goal is to separate authority-sensitive decisions (what to cite) from language generation (how to explain), while enforcing evidence constraints during drafting. Conventional Legal RAG pipelines (e.g., "Dynamic Legal RAG") typically enhance general-purpose retrieval and generation with legal-domain modules such as Legal Entity Recognition (LER), citation parsing, and access to specialized legal knowledge bases. However, these pipelines often lack an explicit, *systematic* mechanism for enforcing precedent doctrine— including court hierarchy, jurisdictional bindingness, temporal validity, and conflicting lines of authority—during both retrieval and generation. PA-MA-RAG fills this gap by operationalizing core common-law reasoning principles directly inside the RAG loop, ensuring that the system prioritizes *controlling authority* and produces analysis that is not only relevant, but also legally well-grounded and internally consistent (Bench-Capon, 2005).

Figure 1 & 2 presents the PA-MA-RAG workflow, where multiple specialized agents collaborate to convert an initial legal question into a precedent-grounded answer. The Issue Framer Agent extracts jurisdictional and doctrinal constraints (e.g., forum, court level, time window), after which the Retriever Agent executes hybrid retrieval over case law and supporting materials. Retrieved items are then processed by an Authority & Validity Ranker Agent, which re-ranks candidates using a structured precedent authority score (hierarchy, jurisdiction, recency, and treatment signals), explicitly reducing reliance on persuasive or outdated authorities when binding precedent is available. Next, a Conflict Checker Agent identifies contradictions among candidate precedents and resolves them via hierarchy and later-in-time priority rules or escalates uncertainty when conflicts cannot be safely reconciled. Finally, the Drafting Agent produces a citation- anchored legal analysis, while a Citation Verifier Agent enforces claim-level support by requiring that each key proposition be traceable to retrieved holdings or quoted passages—triggering targeted re-retrieval when gaps are detected. This multi-agent design extends the "legal-domain RAG" foundation by embedding precedent governance as a first-class constraint across retrieval, ranking, conflict resolution, and generation (Hinkle, 2015).
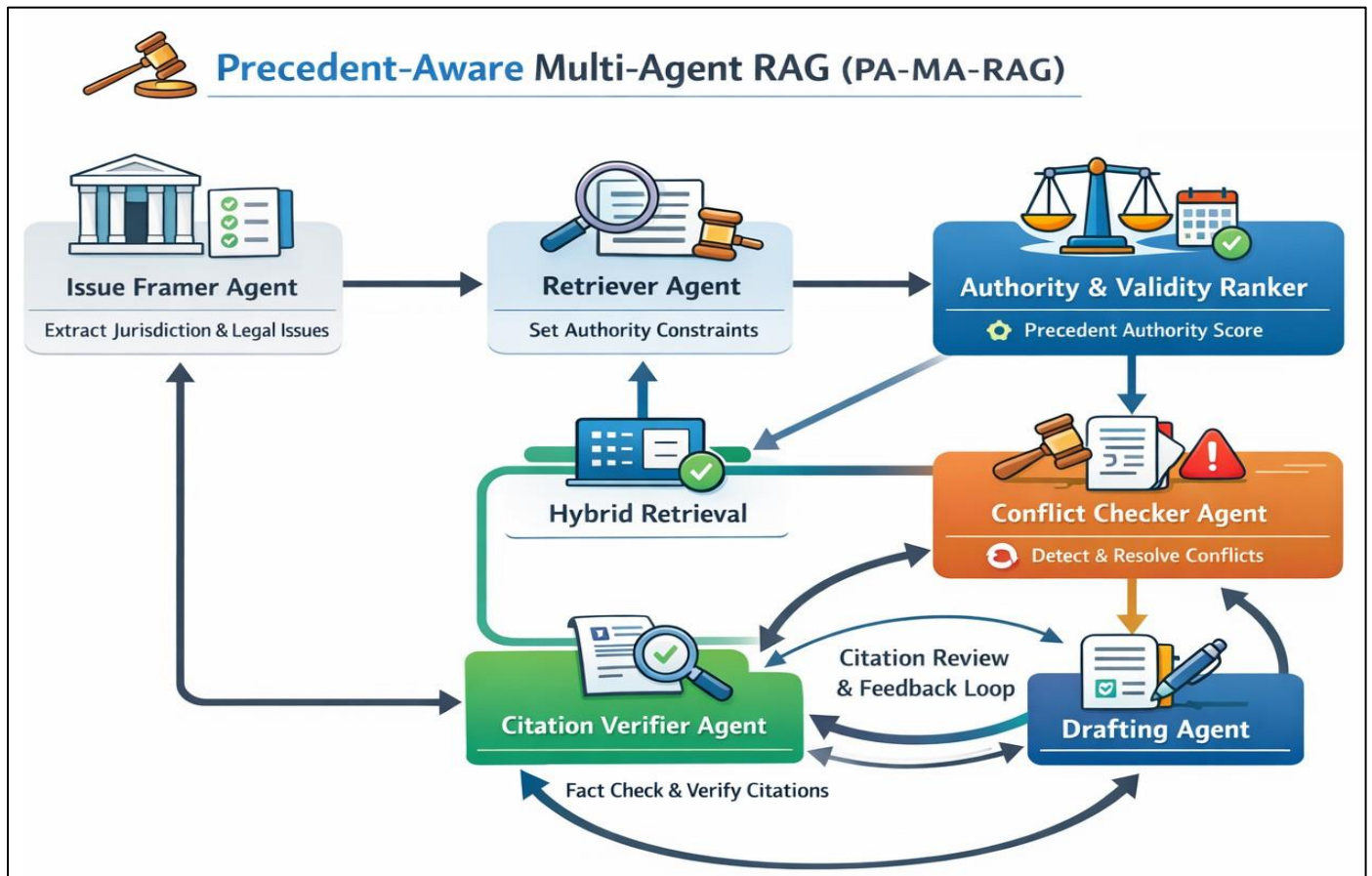
Fig 2 Precedent-Aware Multi-Agent RAG (PA-MA-RAG) Workflow Illustrating Agent Coordination, Authority-Aware Ranking, Conflict Handling, and Citation Verification (Your Name, 2025).
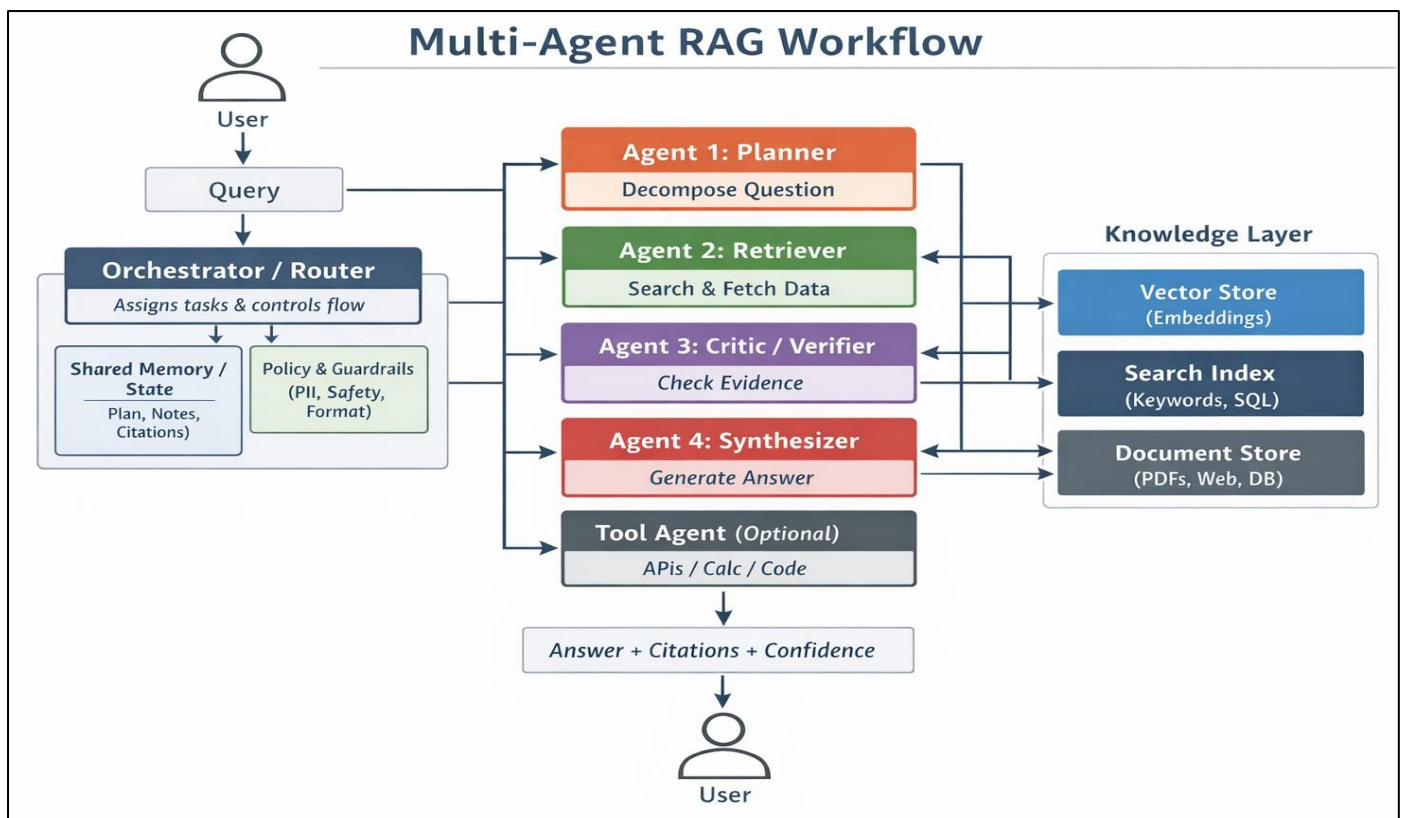


Fig 3 Precedent-Aware Multi-Agent RAG (PA-MA-RAG) Illustrating a Team of Specialized AI Agents Working Together so the Final Answer is Grounded in Retrieved Documents, Not Just "what the Model Remembers."

➢ *Agent-Based Decomposition*

The system consists of eight collaborating agents described in Table 1.

Table 1 Agents in PA-MA-RAG and their Responsibilities.

| Agent | Inputs | Outputs / Responsibility |
|---|---|---|
| **Issue Framer** | User query; optional facts | Extract legal issues, jurisdiction cues, and time constraints; produce a structured query frame. |
| **Authority Planner** | Query frame; court maps | Set binding-first constraints, court ladder, and stop conditions for retrieval. |
| **Retriever** | Planner constraints; corpus index | Hybrid retrieval (BM25 + dense), query decomposition, and citation expansion; returns candidate cases. |
| **Precedent Ranker** | Candidates; metadata; citation graph | Re-rank by authority score; filter negative treatment; select top-k authorities. |
| **Holding Extractor** | Selected cases | Extract holdings / ratio snippets to serve as evidence units for generation. |
| **Conflict Checker** | Authorities; holdings | Detect contradictions; resolve by hierarchy and time; produce a consistent authority set or flagged uncertainty. |
| **Drafting Agent** | Authority set; evidence units | Generate structured analysis with inline citations and quotation-backed claims. |
| **Citation Verifier** | Draft; evidence units | Validate each proposition; request revision or additional retrieval when unsupported. |

Coordination proceeds in loops. The Authority Planner sets constraints; the Retriever proposes candidates; the Ranker selects controlling precedent; the Conflict Checker tests consistency; and the Verifier enforces evidence. If the Verifier rejects a claim, it triggers either targeted retrieval or a rewrite with narrower scope and explicit uncertainty.

➢ *Authority-Constrained Precedent Ranking*

Let q denote the query frame and c a candidate case. We score candidates with an authority- constrained objective:

$$S(c \mid q) = wA \cdot A(c,q) + wJ \cdot J(c,q) + wT \cdot T(c,q) + wC \cdot C(c) - wN \cdot N(c)$$

A(c,q) measures binding strength from court hierarchy; J(c,q) measures jurisdictional match; T(c,q) captures temporal validity; C(c) estimates citation centrality; and N(c) penalizes negative treatment. Weights can be tuned on held-out judgments or set to enforce a binding-first policy.

- Authority score: Encodes court level and bindingness relative to q, favoring controlling courts.
- Jurisdictional score: Rewards exact jurisdiction matches; uses persuasive authority only when necessary.
- Temporal score: Downweights outdated precedents and accounts for subsequent history.
- Citation centrality: Uses the citation graph to prefer influential, widely relied-on cases.
- Negative-treatment penalty: Penalizes overruled, reversed, or abrogated decisions to avoid bad law.

➢ *Conflict Detection and Resolution*

The Conflict Checker compares extracted holdings to identify contradictory tests or standards. Resolution follows a precedence rule set: binding over persuasive, higher over lower courts, and later-in-time decisions over earlier ones

when authority level is equal. Unresolvable conflicts are surfaced explicitly, enabling the Drafting Agent to present competing views with clear jurisdictional assumptions.

➢ *Evidence-Grounded Drafting and Citation Verification*

The Drafting Agent is constrained to evidence units produced by the Holding Extractor. Each paragraph is generated with an explicit mapping from propositions to supporting quotations and citations. The Citation Verifier enforces a supported-claim policy: any claim without a retrieved support span triggers revision or additional retrieval. This reduces hallucinated legal propositions and misattributed holdings.

➢ *Relation to Prior Work*

PA-MA-RAG extends standard RAG by incorporating legal authority constraints and extends multi-agent RAG by specializing agents for precedent selection, conflict resolution, and citation verification.

## IV. IMPLEMENTATION CONSIDERATIONS & CHALLENGES

➢ *Data Acquisition and Curation Requirements*

A deployable precedent-aware system requires more than case text. It needs structured metadata and relational signals to compute authority and validity.

- Court hierarchy mappings: Machine-readable jurisdictional and hierarchical relations between courts.
- Temporal metadata: Decision dates and subsequent history, including reversal or overruling events when available.
- Citation networks: Inter-case citation edges, enriched with citation context and treatment labels where possible.

- Holding extraction: Segmentation methods to reliably identify holdings and supporting rationale.

➢ *Technical and Computational Complexities*

Authority-aware ranking adds overhead: metadata joins, graph traversals, and iterative retrieval loops. Efficient indexing and caching are needed for interactive use. Precomputing authority features and using approximate graph algorithms can reduce latency.

➢ *Legal Knowledge Representation*

Holdings are often fact-sensitive and may contain multiple sub-rules. Treatment signals can be noisy or incomplete. Systems should expose traceable justifications, allow user control over jurisdictional assumptions, and support abstention when evidence is insufficient.

## V. EVALUATION METHODOLOGY

➢ *Tasks*

We evaluate two backbone tasks: (i) precedent retrieval, where the system must retrieve controlling or highly relevant cases for a query, and (ii) retrieval-augmented legal analysis generation, where the system drafts a short analysis grounded in cited precedents.

➢ *Benchmark Datasets*

CLERC supports both citation retrieval and retrieval-augmented legal analysis generation. CaseHOLD evaluates holding identification in a multiple-choice format. COLIEE provides shared tasks for legal case retrieval and entailment.

➢ *Baselines and Ablations*

Baselines include single-agent RAG and generic multi-agent RAG without authority constraints. Ablations remove authority scoring, negative-treatment filtering, conflict checking, and citation verification.

➢ *Metrics*

We recommend the following metrics:

- Recall@k / nDCG@k: Standard IR metrics for retrieval quality.

- Authority correctness: Whether binding authority is prioritized when present.

- Citation precision / recall: Whether cited cases match gold citations or expert judgments.

- Supported-claim rate: Fraction of propositions supported by retrieved evidence spans.

- Hallucination rate: Rate of fabricated cases, citations, or holdings.

➢ *Human Evaluation*

Expert evaluation is recommended for final validity judgments, scoring (i) correctness of cited authority, (ii) faithfulness to holdings, (iii) completeness, and (iv) transparency about uncertainty and jurisdictional assumptions.

## VI. ETHICAL CONSIDERATIONS & LIMITATIONS

Fluent legal language and citations can create over-trust. PA-MA-RAG should be framed as research assistance, not legal advice. Systems should display citation provenance, quote supporting passages, and log intermediate steps for auditability.

Limitations include incomplete treatment metadata, ambiguous jurisdictions, corpus gaps, and cost/latency from iterative agent loops. Authority-aware ranking improves alignment with stare decisis but does not guarantee correctness without high-quality data and oversight.

## VII. CONCLUSION AND FUTURE WORK

We introduced PA-MA-RAG, a precedent-aware multi-agent retrieval-augmented generation framework for case law analysis. By combining authority-constrained ranking, conflict-aware precedent resolution, and citation verification, the system better matches common-law reasoning than standard RAG. Future work includes learning authority weights from expert feedback, integrating richer treatment taxonomies, and evaluating on practitioner datasets across jurisdictions.

## REFERENCES

[1]. Hou, A. B., Weller, O., Qin, G., Yang, E., Lawrie, D., Holzenberger, N., Blair-Stanek, A., & Van Durme, B. (2024). CLERC: A Dataset for Legal Case Retrieval and Retrieval-Augmented Analysis Generation. arXiv:2406.17186.

[2]. Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.-t., Rocktäschel, T., Riedel, S., & Kiela, D. (2020). Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. NeurIPS 2020.

[3]. Nguyen, T., Chin, P., & Tai, Y.-W. (2025). MA-RAG: Multi-Agent Retrieval-Augmented Generation via Collaborative Chain-of-Thought Reasoning. arXiv:2505.20096.

[4]. Rabelo, J., Goebel, R., Kim, M.-Y., Kano, Y., Yoshioka, M., & Satoh, K. (2024). Overview and Discussion of the Competition on Legal Information Extraction/Entailment (COLIEE) 2023. Review of Socionetwork Strategies, 18(1), 27-47.

[5]. Zheng, L., Guha, N., Anderson, B. R., Henderson, P., & Ho, D. E. (2021). When Does Pretraining Help? Assessing Self-Supervised Learning for Law and the CaseHOLD Dataset. Proceedings of the 18th International Conference on Artificial Intelligence and Law (ICAIL).