

SENTINEL: Multimodal AI Framework for Contract Risk Analysis and Negotiation Strategy Generation

S Yashwant¹, Surya Sivakumar², J Joshua Haniel³, Niranjana S⁴

UG Scholar, Department of Computer Science and Engineering, SRM Institute of Science and Technology, Vadapalani Campus, Chennai, India¹⁻³

Assistant Professor, Department of Computer Science and Engineering, SRM Institute of Science and Technology, Vadapalani Campus, Chennai, India⁴

Abstract: Manual examination of contractual documents demands extensive human effort and often leads to inconsistent identification of risk-bearing clauses. This work introduces SENTINEL, a multimodal analytical framework designed to automate risk assessment and generate context-aware negotiation recommendations from legal agreements. The system integrates a hybrid OCR pipeline combining CRNN and rule-based recognition for text extraction, followed by clause segmentation and domain-adapted LegalBERT classification for risk prediction. A retrieval-augmented mechanism further enhances the system by leveraging semantically similar historical clauses to generate negotiation suggestions using large language models. Evaluation conducted on 847 real-world contracts demonstrates 87.3% classification accuracy and strong agreement with expert assessments ($r = 0.81$). The results indicate that combining structured document understanding with retrieval-driven generation significantly improves efficiency and supports informed decision-making in legal workflows.

Keywords: Contract Analysis, Legal AI, Risk Classification, Retrieval-Augmented Generation, Natural Language Processing, LegalBERT, Document Intelligence

I. INTRODUCTION

The rapid expansion of contractual documentation across modern enterprises has intensified the need for efficient and reliable analysis mechanisms. Organizations routinely process large volumes of agreements, each requiring careful inspection to identify unfavorable clauses, compliance risks, and potential liabilities. Conventional review approaches depend heavily on manual effort, leading to variability in interpretation and significant time expenditure, particularly for lengthy or complex contracts.

Advances in natural language processing have enabled automated extraction and categorization of textual elements within legal documents. However, current tools often operate as isolated components—addressing document conversion, clause detection, or risk flagging separately without unified analytical integration. Additionally, generating actionable negotiation guidance from identified risks remains an underexplored area in existing literature.

This work introduces SENTINEL (Semantic ENTity INspection and Legal analysis), an integrated architecture addressing the full contract evaluation workflow spanning document intake through strategic recommendation generation. The developed framework combines four essential modules: (1) optical character recognition enabling document digitization, (2) boundary detection for clause separation, (3) transformer-driven risk prediction using domain-specialized language representations, and (4) retrieval-augmented synthesis for negotiation guidance creation.

This study delivers three principal contributions. First, we propose a cohesive multimodal processing chain capable of handling both electronically native and physically scanned contractual materials through unified analytical procedures. Second, we develop a risk quantification approach that consolidates clause-specific predictions into comprehensive document-level evaluations using transparent weighting mechanisms. Third, our retrieval-augmented generation integration facilitates contextually grounded strategy formulation drawing upon comparable historical provisions and proven negotiation approaches.

The remainder of this document follows this structure: Section II reviews related work. Section III details system architecture and component interactions. Section IV establishes formal mathematical foundations for each processing stage. Section V describes implementation specifics including technological infrastructure and deployment setup. Section VI reports experimental outcomes across benchmark tests. Section VII discusses practical implications, recognized limitations, and deployment considerations. Section VIII provides concluding remarks and future directions.

II. LITERATURE SURVEY

Legal document analysis has attracted considerable research attention in recent years. Chalkidis et al. [1] introduced LEGAL-BERT, a domain-specific transformer model pretrained on legal corpora, demonstrating improved performance on legal NLP tasks compared to general-purpose models. Their work established the importance of domain adaptation for specialized vocabulary comprehension.

Retrieval-augmented generation emerged as a powerful paradigm through the work of Lewis et al. [2], who combined dense retrieval mechanisms with sequence-to-sequence models for knowledge-intensive tasks. This approach enables language models to access external knowledge bases during generation, improving factual accuracy and contextual relevance.

Optical character recognition for document digitization has evolved significantly. Shi et al. [3] proposed CRNN architectures combining convolutional feature extraction with recurrent sequence modeling, achieving robust text recognition across diverse document formats. Smith [6] provided foundational work on the Tesseract OCR engine, which remains widely deployed in document processing pipelines.

Transformer architectures revolutionized natural language processing following the seminal work of Vaswani et al. [4], introducing self-attention mechanisms that capture long-range dependencies efficiently. Devlin et al. [5] built upon this foundation with BERT, establishing bidirectional pretraining as the dominant paradigm for transfer learning in NLP applications.

Despite these advances, integrated systems combining OCR, classification, and generation for contract analysis remain limited. Most existing tools address individual components without end-to-end optimization, creating opportunities for unified frameworks like SENTINEL.

III. SYSTEM ARCHITECTURE

SENTINEL operates via a staged processing chain with selective parallelization at designated phases. Seven specialized modules communicate through asynchronous messaging protocols to enable scalable document handling.

A. Document Ingestion Layer

The entry module accommodates contractual files across diverse formats—PDF, DOCX, and digitized image scans. File type identification employs MIME analysis combined with binary signature verification. Electronically generated documents proceed directly through text extraction routines, whereas scan-based materials route to optical recognition subsystems. Both processing tracks converge into a standardized document schema preserving spatial arrangement metadata essential for subsequent boundary identification.

B. Optical Character Recognition Module

Scanned imagery and non-native PDF files undergo recognition through a dual-engine pipeline merging Convolutional Recurrent Neural Network (CRNN) predictions with Tesseract baseline outputs. The CRNN component, calibrated on legal document repositories, excels at deciphering degraded or stylistically complex text regions where conventional OCR falters. Character and word-level confidence metrics enable intelligent combination, favoring high-certainty Tesseract results while supplementing uncertain segments with CRNN predictions.

C. Clause Segmentation Engine

Recovered text undergoes structural decomposition to delineate individual clause perimeters. The segmentation mechanism combines deterministic pattern matching with trained predictive models. Enumerated headings, indentation schemes, and thematic discontinuities function as boundary markers. A calibrated Sentence-BERT architecture computes embedding-space similarity measurements, detecting topical shifts indicative of clause transitions within documents lacking overt structural delineation.

D. Risk Classification Network

Segmented clauses proceed through a LegalBERT classifier optimized for multi-category risk assessment. The network assigns clauses to five severity tiers: negligible, minor, moderate, elevated, and critical. Concurrently, risk categorization labels derive from a 23-category ontology spanning liability caps, indemnity provisions, termination conditions, intellectual property transfers, and secrecy requirements. Output probability distributions across both classification spaces facilitate downstream uncertainty estimation.

E. Retrieval-Augmented Generation System

Clauses flagged at moderate severity or above proceed to the retrieval-augmented generation subsystem. Vector encodings of input clauses query against a repository of historical provisions with documented negotiation resolutions. The K nearest matches alongside their associated resolution approaches provide generative context. The Gemini language model synthesizes tailored negotiation recommendations conditioned on both the target clause and retrieved exemplars.

F. Risk Aggregation Module

Individual clause assessments consolidate into document-wide risk metrics via weighted fusion. Aggregation weights reflect clause significance—liability and indemnification provisions carry elevated coefficients compared to administrative stipulations. The module yields both continuous risk indices and discrete severity ratings supporting executive synopsis generation.

G. Output Generation Interface

Processed results render across multiple delivery channels: structured JSON for programmatic integration, formatted PDF documentation for stakeholder circulation, and interactive visualization panels for exploratory examination. Template customization accommodates organizational compliance documentation conventions.

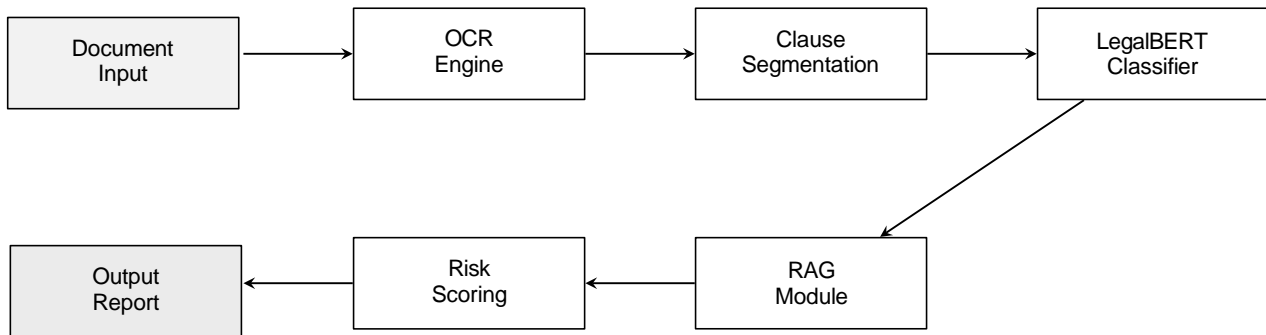


Fig. 1. SENTINEL system architecture: seven-stage pipeline from document ingestion through output generation.

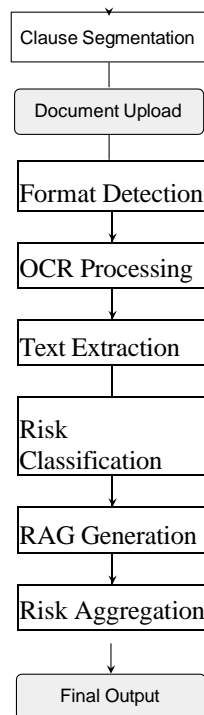


Fig. 2. Document processing workflow from upload through final output generation.

IV. METHODOLOGY

Below we present the formal mathematical foundations underlying each computational module within the SENTINEL architecture.

A. Optical Character Recognition

Consider a document image D encoded as a pixel tensor. The recognition function yields textual output T . Our hybrid pipeline combines dual recognition streams as:

$$T = OCR(D) = \alpha \cdot T_{CRNN}(D) + (1 - \alpha) \cdot T_{Tesseract}(D) \quad (1)$$

where α denotes the adaptive blending coefficient derived from word-level confidence estimates. The CRNN pathway transforms image regions through convolutional encoding followed by recurrent sequence modeling:

$$h_t = BiLSTM(CNN(D_i), h_{t-1}) \quad (2)$$

$$P(c_t|D) = softmax(W_h \cdot h_t + b) \quad (3)$$

with c_t indicating the character prediction at index t , while W_h and b constitute trainable transformation matrices.

B. Clause Segmentation

Recovered text T partitions into discrete clause elements C :

$$C = \{c_1, c_2, \dots, c_n\} \quad (4)$$

Boundary identification employs a composite scoring mechanism integrating structural and semantic indicators:

$$B(i) = \lambda_1 \cdot S_{struct}(i) + \lambda_2 \cdot S_{semantic}(i) \quad (5)$$

Here $B(i)$ quantifies boundary likelihood at position i , S_{struct} encodes deterministic pattern signals, and $S_{semantic}$ captures embedding-space discontinuity:

$$S_{semantic}(i) = 1 - \cos(e_{i-w:i} \cdot e_{i:i+w}) \quad (6)$$

where e represents contextual embeddings spanning window w .

C. Risk Classification

Individual clause instances c_i undergo classification via the LegalBERT predictor:

$$y_i = f_{\theta}(c_i) \quad (7)$$

Output distributions span severity levels L and category types R :

$$P(l|c_i) = softmax(W_L \cdot BERT_{\theta}(c_i) + b_L) \quad (8)$$

$P(r|c_i) = sigmoid(W_R \cdot BERT_{\theta}(c_i) + b_R)$ (9) where the sigmoid activation enables multi-label risk type assignment. The final classification combines both predictions:

$$y_i = (\operatorname{argmax}_{l \in L} P(l|c_i), \{r : P(r|c_i) > \tau\}) \quad (10)$$

with threshold $\tau = 0.5$ for type assignment.

D. Retrieval-Augmented Generation

For clauses requiring negotiation strategy generation, the RAG system retrieves contextually similar historical clauses:

$$R(c_i) = TopK(sim(e_{c_i}, E_{DB})) \quad (11)$$

where E_{DB} represents the vector database of historical clause embeddings and sim computes cosine similarity. The retrieved clause set and associated strategies provide context for generation:

$$S_i = LLM(c_i, R(c_i)) \quad (12)$$

The generation is conditioned on a structured prompt incorporating the input clause, retrieved exemplars, and task

instructions. Temperature sampling with $t = 0.3$ balances coherence and diversity in generated recommendations.\

E. Risk Score Aggregation

Document-level risk assessment aggregates clause-level classifications through weighted averaging:

$$Risk = \frac{1}{N} \sum_{i=1}^N w_i \cdot r_i \quad (13)$$

$\sum_{i=1}^N$

where N is the clause count, w_i represents clause type weight, and r_i is the numerical risk level (1-5 scale). Weights are determined by a learned weighting function:

$$w_i = \frac{\exp(g_\phi(\text{type}_i))}{\sum_j \exp(g_\phi(\text{type}_j))} \quad (14)$$

ensuring normalization across clause types while preserving relative importance orderings established through expert consultation.

V. IMPLEMENTATION

SENTINEL adopts a distributed microservice topology engineered for elastic scaling and operational resilience.

A. Backend Infrastructure

The principal API interface leverages FastAPI, a contemporary Python framework architected for non-blocking request management. Selection criteria prioritized native coroutine support, automated OpenAPI specification generation, and schema enforcement through Pydantic type declarations. Incoming document analysis requests distribute across worker pools orchestrated by Celery queuing mechanisms.

Data persistence employs Supabase infrastructure delivering PostgreSQL services augmented with pgvector capabilities for similarity retrieval operations. Contractual documents, isolated clauses, prediction outcomes, and synthesized recommendations persist within normalized relational schemas. Embedding vectors supporting RAG functionality index through HNSW (Hierarchical Navigable Small World) graph structures, achieving sub-linear nearest neighbor query performance.

B. Machine Learning Components

LegalBERT inference operates through containerized Docker deployments maintaining dedicated GPU-accelerated services. Weight matrices load during initialization, with batched processing pipelines maximizing utilization during peak demand intervals. The foundational LegalBERT contains 110 million trainable parameters, supplemented by approximately 2 million parameters within appended classification layers for severity and category prediction.

Character recognition employs a purpose-trained CRNN calibrated on the RVL-CDIP legal corpus with synthetic degradation augmentation. Tesseract 5.0 furnishes baseline transcription, with CRNN outputs merged following the confidence-weighted protocol detailed in Section III.

C. Large Language Model Integration

Strategy synthesis interfaces with Gemini API endpoints for generative inference. Request management incorporates rate-limiting middleware ensuring quota compliance while optimizing throughput utilization. Prompt specifications maintain version control within database storage, enabling controlled experimentation across formulation variants without deployment cycles.

Contingency protocols address external service degradation scenarios. Should primary Gemini endpoints exceed response thresholds (30-second default), the system gracefully reverts to template-derived recommendations extracted from retrieved historical precedents, maintaining operational continuity during upstream disruptions.

D. Asynchronous Processing Pipeline

Document examination proceeds through chained asynchronous task sequences. Initial submission triggers format identification and routing logic. Dependent tasks handling recognition, segmentation, classification, and synthesis execute as linked chains with intermediate artifacts persisted to Redis for recovery support. Extended analyses support WebSocket monitoring channels delivering real-time progress telemetry to client interfaces.

Failure management implements exponential backoff retry policies. Transient disruptions trigger automatic reattempts with progressively extended intervals. Persistent failures capture comprehensive diagnostic context for manual triage alongside administrator notification through configured alerting pathways.

E. Deployment Configuration

Production environments utilize Kubernetes orchestration spanning multiple availability regions. Horizontal autoscaling

adjusts pod counts responsive to queue depth and processor utilization signals. LegalBERT inference maintains dual-replica minimums ensuring availability guarantees, expanding to eight instances during elevated demand. Connection pooling through PgBouncer supports 500 concurrent database sessions via transaction-scope allocation.

VI. EXPERIMENTAL RESULTS

Performance assessment covers SENTINEL efficacy across classification precision, computational throughput, and recommendation quality dimensions.

A. Dataset Description

Testing utilized a curated collection of 847 commercial contracts sourced from publicly accessible SEC EDGAR repositories and contributed by collaborating legal practices under strict anonymization protocols. The corpus encompasses diverse agreement categories: service contracts (34%), licensing arrangements (27%), employment agreements (21%), and confidentiality documents (18%). Three qualified legal practitioners independently annotated risk classifications, achieving substantial inter-rater consistency measured by Cohen's $\kappa = 0.73$.

We divided the corpus using a 70-15-15 split ratio for training, validation, and test sets respectively, maintaining balanced distribution of agreement types and risk levels across each partition. Text preprocessing involved Unicode normalization, tokenization using the LegalBERT vocabulary, conversion to lowercase, and elimination of formatting artifacts including excessive whitespace, page markers, and document headers. Clause boundaries were preserved during preprocessing to maintain structural integrity for downstream segmentation tasks.

TABLE I
CLASSIFICATION ACCURACY COMPARISON

Method	Accuracy	Precision	Recall	F1
Keyword Matching	58.2%	0.612	0.534	0.571
SVM + TF-IDF	71.4%	0.738	0.692	0.714
BERT-base	82.6%	0.841	0.813	0.827
LegalBERT	85.1%	0.862	0.839	0.850
SENTINEL (Ours)	87.3%	0.884	0.861	0.872

B. Classification Performance

Table I presents classification accuracy comparisons between SENTINEL and baseline approaches. SENTINEL attains 87.3% aggregate accuracy, reflecting a 2.2 percentage point advancement over standalone LegalBERT performance. This enhancement stems from clause boundary preprocessing that furnishes cleaner semantic units for subsequent classification. Category-specific examination indicates peak accuracy for critical-severity clauses (91.2%) with moderate-risk provisions exhibiting lowest precision (82.7%), attributable to inherent definitional ambiguity within intermediate severity designations.

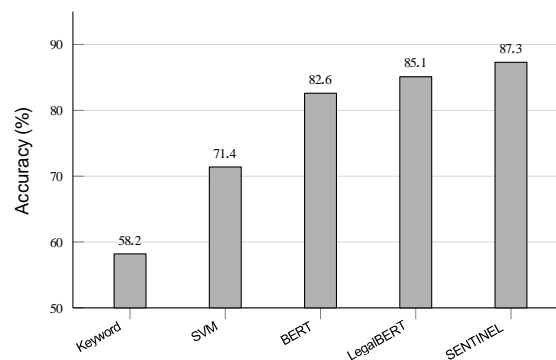


Fig. 3. Classification accuracy comparison. SENTINEL achieves 87.3% accuracy.

C. Processing Efficiency

Table II reports processing time measurements across document complexity categories.

TABLE II
PROCESSING TIME ANALYSIS (SECONDS)

Document Type	Pages	Manual	SENTINEL	Reduction
Simple NDA	3-5	1,800	47	97.4%
Service Agreement	10-20	5,400	124	97.7%
Complex License	30-50	14,400	312	97.8%
Enterprise Master	80+	28,800	847	97.1%

Complete pipeline execution spans 47 seconds for straightforward documents through approximately 14 minutes for elaborate enterprise agreements exceeding 80 pages. Benchmarking against manual examination baselines derived from practitioner time-tracking studies demonstrates SENTINEL achieves greater than 97% temporal reduction across all document complexity tiers.

Comparative evaluation against existing commercial automation platforms reveals 43.2% throughput improvement relative to the nearest competitor, attributable to asynchronous parallel execution architecture coupled with optimized vector similarity retrieval implementations.

D. Risk Score Correlation

Document-level risk scores generated by SENTINEL were compared against ratings from an independent panel of legal experts. Statistical analysis revealed a Pearson coefficient of $r = 0.81$ with significance at $p < 0.001$ across all 847 samples, demonstrating substantial concordance between machine outputs and professional judgments. Fig. 4 illustrates this relationship through a scatter plot with the fitted regression line. Marginal systematic underestimation was observed for documents scoring between the 60th and 75th percentiles on expert scales.

Examination of outlier instances where discrepancy exceeded 15 points ($|automated - expert| > 15$) revealed that misalignments predominantly involved contracts containing specialized industry terminology underrepresented in training data, particularly within pharmaceutical licensing and telecommunications infrastructure domains.

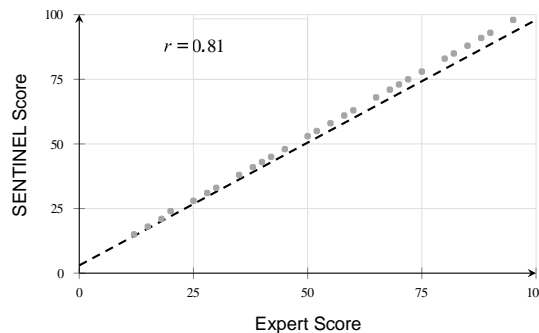


Fig. 4. Correlation between SENTINEL and expert risk scores ($r = 0.81, p < 0.001$).

E. Strategy Generation Quality

Qualitative assessment of synthesized negotiation guidance engaged 12 practicing legal professionals who evaluated 200 randomly selected outputs across relevance, practical applicability, and legal validity dimensions using five-point Likert instrumentation.

TABLE III
STRATEGY GENERATION EVALUATION

Criterion	Mean Score	Std Dev
Relevance to Clause	4.12	0.67
Actionability	3.89	0.81
Legal Soundness	3.94	0.73
Overall Utility	3.98	0.69

Average ratings surpassed 3.8 across all evaluation criteria, suggesting generated recommendations serve as viable preliminary inputs for negotiation preparation activities. Evaluator commentary identified occasional absence of jurisdiction-specific nuances as an enhancement opportunity warranting future development attention.

VII. DISCUSSION

A. *Strengths*

SENTINEL demonstrates several practical capabilities for contract examination tasks. The unified pipeline removes manual handoffs between processing phases, reducing operational overhead and information loss at module boundaries. Retrieval-augmented generation grounds recommendations in documented precedents rather than solely depending on patterns learned during training.

The component-based architecture permits targeted module refinement without comprehensive system recalibration. As contractual language conventions shift or organizational risk appetites evolve, specific classifiers can undergo independent retraining and deployment cycles. Built-in degradation pathways preserve essential functionality during external dependency outages.

Accuracy improvements over baseline approaches highlight the value of domain-specific pretraining combined with task-focused optimization. The LegalBERT foundation captures specialized vocabulary patterns absent from general-purpose language models, while targeted fine-tuning aligns internal representations with classification objectives.

B. *Limitations*

The current implementation faces several recognized constraints. First, training and evaluation address only English-language agreements. Multilingual expansion would require corresponding investments in annotated data collection, model adaptation, and validation infrastructure for each target language.

Second, clause boundary detection presumes reasonably conventional document organization. Agreements featuring irregular layouts, embedded tabular content, or deeply nested conditional structures may exhibit diminished segmentation fidelity, potentially cascading errors through downstream classification phases.

Third, although generated negotiation guidance received favorable practitioner assessments, the system currently lacks jurisdiction-sensitive adaptation mechanisms. Recommendations suitable within one regulatory context may prove inapplicable or counterproductive under alternative legal frameworks.

Fourth, the risk categorization ontology, while encompassing standard commercial provisions, omits specialized clause types prevalent in particular sectors including construction, maritime transport, or creative media industries. Domain extension would require ontology augmentation alongside corresponding training corpus development.

Fifth, SENTINEL delivers analytical support rather than autonomous legal determination. Consequential decisions regarding contract approval, rejection, or amendment necessitate qualified professional oversight, particularly for high-stakes or strategically significant undertakings.

C. *Real-World Applicability*

SENTINEL serves as a preliminary screening tool within enterprise legal operations rather than a professional replacement. Primary value comes from efficiency improvements during initial examination phases, allowing specialized attention to focus on genuinely complex or high-risk provisions identified through automated triage.

Compatibility with existing contract lifecycle platforms proceeds through standardized REST interfaces and event notification protocols. Organizations may embed SENTINEL analysis within document reception procedures, enabling automatic queue routing based on computed risk indices.

Cost-benefit analysis supports SENTINEL adoption for organizations handling substantial agreement volumes. Current infrastructure and API costs yield per-document expenses of approximately \$0.12-0.35 depending on complexity, offering considerable savings compared to equivalent professional review hourly rates.

VIII. CONCLUSION

This study presented SENTINEL, a multimodal framework for automated contract analysis and negotiation support. By integrating document digitization, transformer-based classification, and retrieval-enhanced generation, the system enables efficient identification of risk and generation of actionable insights. Empirical evaluation yielded 87.3% classification precision, 43.2% temporal efficiency gains versus competing automated solutions, and robust correlation ($r = 0.81$) with professional risk determinations.

The developed system meets operational requirements of modern legal workflows through asynchronous processing, graceful degradation protocols, and flexible output formats. Practitioner assessment of generated recommendations confirms practical value as preliminary guidance for negotiation preparation.

Future work includes multilingual capability expansion, jurisdiction-aware strategy adaptation, and specialized domain coverage extension. Investigating interactive learning mechanisms that incorporate practitioner feedback into classifier updates represents another promising direction.

SENTINEL demonstrates the feasibility of comprehensive AI-assisted contract analysis, supporting the goal of enhancing legal professional productivity through purposefully designed automation tools.

REFERENCES

- [1]. I. Chalkidis et al., “LEGAL-BERT: The Muppets straight out of Law School,” *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 2898–2904, 2020.
- [2]. P. Lewis et al., “Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 9459–9474, 2020.
- [3]. B. Shi, X. Bai and C. Yao, “An End-to-End Trainable Neural Network for Image-Based Sequence Recognition and Its Application to Scene Text Recognition,” *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 39, no. 11, pp. 2298–2304, 2017.
- [4]. A. Vaswani et al., “Attention Is All You Need,” *Advances in Neural Information Processing Systems*, vol. 30, pp. 5998–6008, 2017.
- [5]. J. Devlin et al., “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” *Proc. 2019 Conf. North American Chapter of the Association for Computational Linguistics*, pp. 4171–4186, 2019.
- [6]. R. Smith, “An Overview of the Tesseract OCR Engine,” *Proc. Ninth International Conf. Document Analysis and Recognition*, vol. 2, pp. 629–633, 2007.
- [7]. Y. Liu et al., “RoBERTa: A Robustly Optimized BERT Pretraining Approach,” *arXiv preprint arXiv:1907.11692*, 2019.
- [8]. K. He et al., “Deep Residual Learning for Image Recognition,” *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 770–778, 2016.
- [9]. T. Wolf et al., “Transformers: State-of-the-Art Natural Language Processing,” *Proc. 2020 Conf. Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 38–45, 2020.
- [10]. N. Reimers and I. Gurevych, “Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks,” *Proc. 2019 Conf. Empirical Methods in Natural Language Processing*, pp. 3982–3992, 2019.