

# Detecting CPRA Compliance Violations in Privacy Policies with LegalBERT and NLI

A Sequential Transfer Learning Approach to Legal Clause Classification

Mukesh Yadav

Department of Applied Data Science  
Clarkson University, Potsdam, USA  
Email: yadvdm@clarkson.edu

Shafique A. Chaudhry

Associate Professor of Information Systems  
David D. Reh School of Business  
Clarkson University, Potsdam, USA  
Email: schaudhr@clarkson.edu

**Abstract**—Ensuring compliance with the California Privacy Rights Act (CPRA) is increasingly challenging due to the complexity of privacy policies and legal language. This paper proposes a sequential transfer learning framework based on LegalBERT for automated CPRA compliance detection. LegalBERT is fine-tuned in two stages: first on the Stanford Natural Language Inference (SNLI) dataset to learn general inferential reasoning patterns, and then on a custom-labeled dataset constructed from OPP-115 privacy policy clauses and CPRA articles. To build this domain-specific NLI dataset, premise-hypothesis pairs were semantically filtered using Sentence-BERT and subsequently labeled as entailment, contradiction, or neutral using GPT-4. The final model employs a single-head classifier and demonstrates strong performance in identifying legal obligation alignment. This approach provides a scalable, interpretable, and legally grounded solution for automated compliance auditing under the CPRA framework.

**Index Terms**—CPRA, Legal NLP, LegalBERT, Natural Language Inference, Privacy Policy Compliance, Sequential Transfer Learning, GPT Labeling, OPP-115 Dataset, Entailment Classification, Semantic Similarity

## I. Introduction

With the growing public awareness of data privacy rights and regulatory frameworks such as the California Privacy Rights Act (CPRA), organizations are under increasing pressure to ensure that their privacy policies comply with legal standards. However, privacy policies are often written in complex, ambiguous language, making manual compliance auditing time-consuming, error-prone, and resource-intensive.

Recent advances in Natural Language Processing (NLP) and Legal AI offer promising pathways to automate compliance checking. In particular, Natural Language Inference (NLI) has emerged as an effective method to determine whether a specific legal obligation (hypothesis) is entailed, contradicted, or left neutral by a clause in a privacy policy (premise). LegalBERT—a transformer model pretrained on legal corpora—has demonstrated strong performance on various legal text classification tasks.

This study introduces a novel sequential transfer learning approach to detect CPRA compliance violations in privacy policies. The methodology involves fine-tuning LegalBERT in two stages: first on the Stanford Natural Language Inference (SNLI) dataset to generalize inferential reasoning, and subsequently on a domain-specific dataset created by semantically pairing OPP-115 privacy policy clauses with

CPRA obligations. The custom dataset was labeled using GPT-4 as a zero-shot classifier across entailment, contradiction, and neutral classes.

This approach not only leverages powerful legal language modeling through LegalBERT, but also incorporates scalable semantic filtering and labeling through Sentence-BERT and GPT-4. The resulting model can automate clause-level compliance verification, making it a valuable tool for legal practitioners, privacy auditors, and data protection officers seeking to ensure CPRA compliance efficiently and accurately.

## II. Related Work

Automating legal text understanding and compliance auditing has gained significant interest in the field of legal NLP. Several studies have focused on natural language inference (NLI) as a foundation for legal reasoning tasks. For instance, the Stanford Natural Language Inference (SNLI) corpus [1] and MultiNLI [2] have enabled the development of models that understand entailment and contradiction relationships in general language. However, the legal domain presents unique challenges due to its formal language, complex sentence structure, and domain-specific vocabulary.

To address these challenges, domain-specific transformer models such as LegalBERT [3] and CaseLaw-BERT [4] have been developed by further pretraining BERT on legal corpora. LegalBERT, in particular, has shown superior performance in downstream legal tasks such as judgment prediction, contract analysis, and statute classification.

Several works have explored the application of BERT-based models for privacy policy analysis. OPP-115 [5], a structured dataset of privacy policies annotated for user rights and practices, has been widely used for tasks such as clause classification, data practice extraction, and policy summarization. More recent work [6] has proposed the use of NLI-style formulation for aligning legal text with regulatory obligations under laws like GDPR and CCPA.

Transfer learning has proven effective in improving model performance when limited labeled legal data is available. Sequential fine-tuning—first on general NLI tasks and then on domain-specific examples—has been recommended in legal text classification to preserve reasoning capabilities while adapting to domain semantics [7].

In contrast to prior studies that rely solely on human-annotated data, our work introduces a scalable method for dataset creation using semantic filtering with Sentence-BERT and automatic labeling via GPT-4. Furthermore, we focus specifically on CPRA compliance, which has received limited attention in prior literature, and demonstrate how entailment-based classification can be used to flag alignment or violation of user rights clauses in privacy policies.

### III. Related Work

Automating legal text interpretation has become a central focus in the field of Legal NLP, particularly in light of growing regulatory complexity and the increasing demand for scalable compliance solutions. Several research directions converge in this study: legal domain adaptation of transformer models, Natural Language Inference (NLI) for legal reasoning, automated privacy policy analysis, and scalable dataset creation methods.

#### A. Legal NLP and Transformer Models

Traditional approaches to legal text classification relied heavily on feature engineering, rule-based systems, or bag-of-words models, which lacked semantic understanding and struggled with generalizability. With the emergence of transformer architectures like BERT, the legal domain saw significant advancements. LegalBERT [8] was introduced as a variant of BERT pretrained on legal corpora including court decisions and statutes. This domain adaptation enabled improved performance across tasks like case law retrieval, contract clause classification, and statutory entailment.

Other legal-domain models such as CaseLaw-BERT and LEDGAR have extended this trend, further demonstrating that language models pretrained on in-domain corpora outperform general-purpose models. However, most of these studies focus on European law or common law systems, with limited application to U.S.-specific privacy regulations such as the California Privacy Rights Act (CPRA).

#### B. Natural Language Inference in Legal Settings

Natural Language Inference (NLI)—the task of determining whether a premise entails, contradicts, or is neutral toward a given hypothesis—has emerged as a robust framework for legal reasoning. Datasets such as SNLI [9] and MultiNLI [10] have provided standardized benchmarks for general-purpose inference models.

Legal adaptations of NLI have shown promise in recent years. For instance, Chalkidis et al. [11] proposed treating statute applicability as an NLI task. Similarly, Elwany et al. [12] explored multi-jurisdictional regulation alignment using entailment modeling. However, much of this work relies on supervised labeling efforts by legal experts, which may not be feasible at scale.

#### C. Privacy Policy Analysis and Regulatory Alignment

Privacy policies are notoriously verbose and difficult to interpret. The OPP-115 dataset [13] addressed this by manually annotating privacy policy text with fine-grained data practice

labels. This dataset has enabled tasks such as automated policy summarization [14], clause extraction [15], and privacy label generation.

Recent efforts have begun leveraging NLI to align clauses in privacy policies with obligations in regulations like the GDPR and CCPA. For instance, Harkous et al. [14] and Laranjeira et al. [16] examined policy-regulation alignment using entailment-based methods. However, few studies focus specifically on CPRA compliance, and most rely on labor-intensive manual alignment.

#### D. Transfer Learning and Dataset Generation

Transfer learning has proven especially valuable in the legal domain due to the scarcity of labeled data. A common approach involves sequential fine-tuning: first training on a general NLI dataset (e.g., SNLI), then adapting to a domain-specific task. This helps retain reasoning ability while learning domain-specific language patterns.

Our approach adopts this two-step fine-tuning strategy and complements it with scalable data generation techniques. We use Sentence-BERT [17] to semantically filter premise-hypothesis pairs and GPT-4 as a zero-shot labeler. This pipeline enables the creation of a high-quality CPRA compliance dataset without the need for extensive manual annotation.

#### E. Our Contributions

In contrast to previous work, our study introduces the first CPRA-specific compliance detection model using LegalBERT. We demonstrate that combining SNLI pretraining, semantic similarity filtering, and GPT-based labeling enables strong performance with minimal human labeling effort. This offers a scalable path toward automating privacy policy audits under evolving legal standards.

#### F. Gaps and Opportunities

Despite notable progress in legal NLP and privacy compliance automation, several key gaps remain. First, much of the existing literature focuses on static clause classification rather than dynamic, inference-based compliance detection. Second, few models provide interpretability or traceability—essential characteristics for legal professionals who must justify compliance decisions. Lastly, there is limited integration of modern large language models (LLMs) such as GPT-4 for legal annotation and hypothesis formulation. Our work addresses these gaps by incorporating GPT-4 for scalable entailment labeling, leveraging domain-adapted LegalBERT for legal text understanding, and structuring compliance as an interpretable NLI task that mirrors legal reasoning. This enables a more transparent and modular system for CPRA policy auditing and lays the groundwork for future compliance systems that are explainable, flexible, and regulation-aware.

**IV. Architecture Description.** Figure 1 illustrates the architecture used for fine-tuning LegalBERT on the SNLI dataset. The model takes premise–hypothesis pairs as input, tokenized in the format [CLS] Premise [SEP] Hypothesis [SEP]. The tokenized sequence is embedded using three types of embeddings: token, positional, and segment, following the BERT encoder-only transformer [18], [19]; subword tokenization can be implemented with WordPiece or SentencePiece [20].

These embeddings are then passed through 12 stacked transformer encoder layers, each containing multi-head self-attention mechanisms, feed-forward networks, residual connections, and layer normalization [18], [19]. After the final encoder layer, the representation of the [CLS] token is extracted, serving as the aggregate representation of the entire input sequence.

This embedding is then passed through a classification head—a fully connected neural network—that outputs probabilities over three classes: entailment, contradiction, and neutral. We fine-tune LegalBERT [3] on SNLI [1] using cross-entropy with AdamW and a linear warmup schedule [19], [21], employing early stopping based on validation accuracy [22]. Robust fine-tuning practices from BERT-family work further inform our setup [23].

Fine-tuning on SNLI serves as a foundational stage prior to domain adaptation on legal datasets. By learning general inference patterns first [1], [2], the model becomes better equipped to handle downstream legal NLI tasks, such as identifying compliance statements in privacy policies or obligations within contracts. This two-stage strategy—general NLI pretraining followed by legal domain adaptation [7]—boosts robustness and transferability across legal text classification tasks.

The encoder stack forms the architectural backbone of LegalBERT. Each of the 12 encoder layers consists of a multi-head self-attention module that captures dependencies across input tokens, followed by a position-wise feed-forward network, with residual connections and layer normalization to enhance stability and convergence [18], [19]. These stacked encoders progressively refine token representations, enabling the model to infer logical relationships, contradictions, and entailments, even in complex sequences. The encoder-only design—without a decoder—makes LegalBERT particularly effective for understanding-focused tasks such as NLI [3], [18]; comparable [CLS] pooling strategies are widely used across modern variants [23]–[25].

Moreover, the use of the [CLS] token as a global representation vector allows the model to efficiently summarize the semantic relationship between the premise and hypothesis into a fixed-length vector. This simplifies downstream classification by concentrating inference computation on a single token representation rather than the entire sequence. Combined with LegalBERT’s legal-domain pretraining [3] and the masked language modeling objectives of BERT [18], this architecture demonstrates strong generalization in tasks requiring legal precision, such as identifying implied obligations, rights, or violations in regulatory documents and contracts.

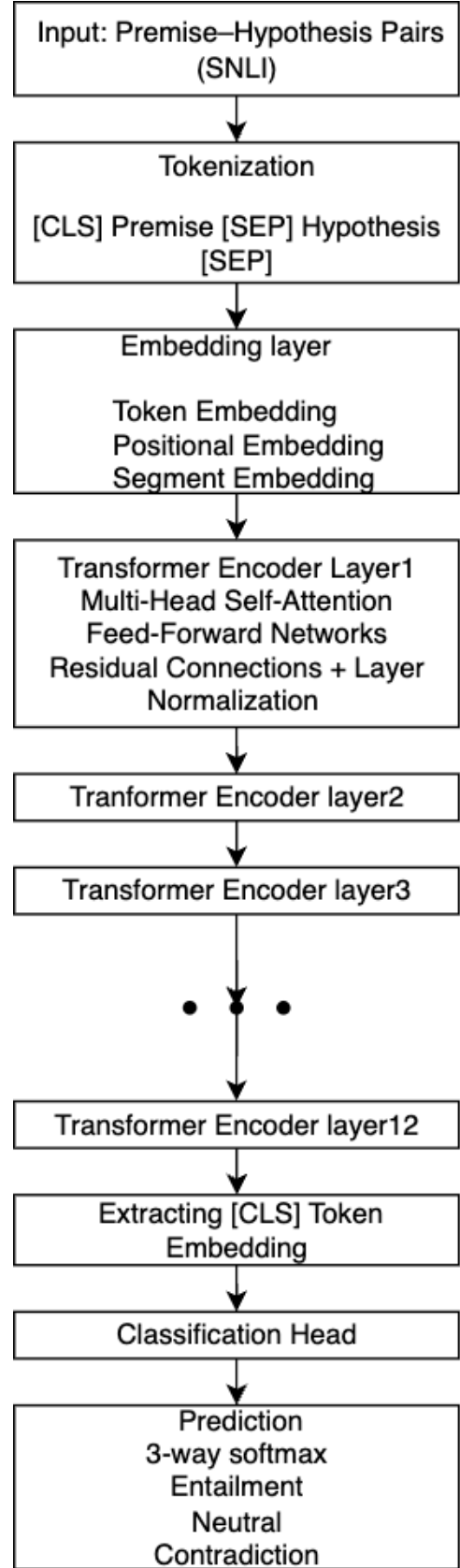


Fig. 1. LegalBERT SNLI fine-tuning architecture. The input [CLS] Premise [SEP] Hypothesis is passed through 12 encoder layers. The [CLS] token is used for 3-class classification.

**V. Data Construction Pipeline.** Figure 2 illustrates the complete, end-to-end process for generating the labeled legal NLI dataset used in our experiments. The pipeline integrates legal-domain text sources, semantic representation techniques, similarity search, and automated annotation.

1) *Text Sources.* The process begins with two distinct inputs:

- **Premises:** Extracted from the OPP-115 corpus, which contains privacy policy statements from real-world websites. Each statement represents a potential factual claim made by an organization [5], [13], [14].
- **Hypotheses:** Derived from the California Consumer Privacy Act (CCPA) and California Privacy Rights Act (CPRA), specifically Articles 2, 3, 7, and 8. Each article encodes explicit legal obligations, rights, or restrictions. Hypotheses are carefully crafted to reflect the key provisions of these articles in a concise, testable form (see also policy-regulation alignment efforts in [6], [16]).

2) *Sentence Embedding.* Both premises and hypotheses are transformed into dense vector representations using **Sentence-BERT** (all-MiniLM-L6-v2) [17], which distills strong semantic similarity performance while remaining lightweight via MiniLM [26]. Comparable alternatives include SimCSE for contrastive sentence embeddings [27].

3) *FAISS Indexing.* The resulting embeddings are indexed using **FAISS** (Facebook AI Similarity Search), a high-performance library for nearest neighbor retrieval in vector space [28]. FAISS enables efficient large-scale similarity search across tens of thousands of sentences; other ANN methods (e.g., HNSW) are viable substitutes depending on latency-recall trade-offs [29].

4) *Semantic Filtering.* For each hypothesis, a semantic search retrieves the most similar premises. Only those with a cosine similarity score above a threshold (e.g.,  $\geq 0.6$ ) are retained. This filtering ensures that pairs are not only lexically similar but also conceptually relevant to the same legal requirement, a standard practice in dense retrieval pipelines [30].

5) *Candidate Pair Formation.* The filtered results are compiled into **semantic pairs**, where each pair consists of a premise and a corresponding hypothesis. These pairs form the candidate set for NLI labeling, following prior practice in policy-to-regulation alignment [6], [16].

6) *Automated Labeling.* Using a large language model (GPT), each semantic pair is classified into one of three NLI categories: *entailment* (the premise supports the hypothesis), *contradiction* (the premise refutes the hypothesis), or *neutral* (the premise is unrelated or does not provide sufficient information). GPT’s labeling process is guided by prompts that include examples and explicit legal reasoning instructions to maximize annotation accuracy; zero/few-shot prompting with large language models is well-established [31], [32], and aligns with broader weak supervision paradigms such as Snorkel [33].

7) *Final Dataset.* The output is a labeled dataset of premise-hypothesis pairs ready for training legal-domain NLI models such as LegalBERT. Each labeled example encodes the semantic and legal relationship between a real-world

privacy policy statement and a specific CPRA/CCPA legal obligation. This dataset forms the foundation for downstream compliance detection tasks, enabling automated auditing of privacy policies against legal requirements [6], [14], [16].

The output of this process is a **fully labeled dataset** ready for fine-tuning LegalBERT or other NLI models. This dataset is domain-specific, reflecting the language of privacy policies and regulatory obligations, and is suitable for automated compliance checking under the CPRA/CCPA.

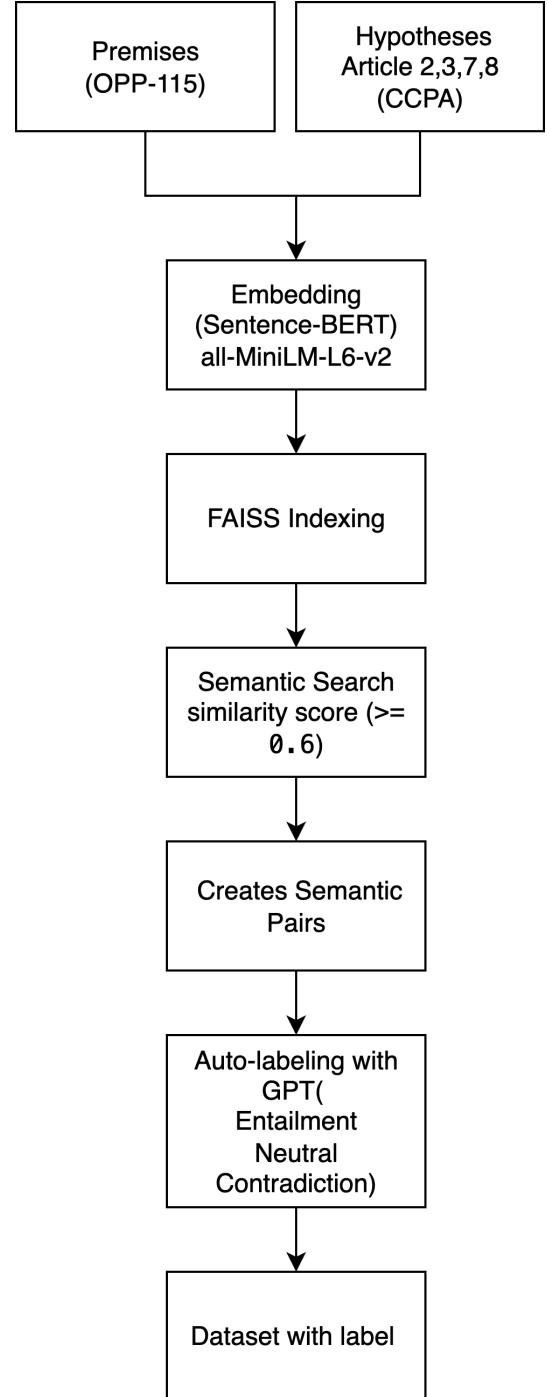


Fig. 2. Pipeline to create the labeled legal NLI dataset.

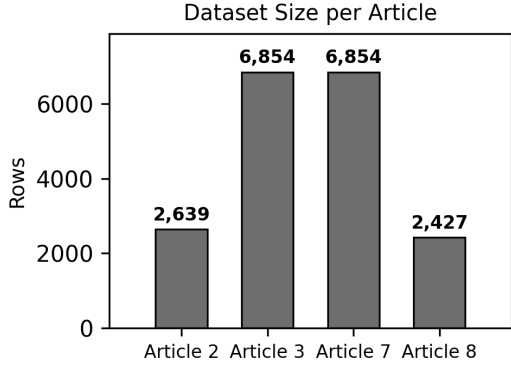


Fig. 3. **Dataset size per article.** Total number of premise–hypothesis pairs retained after semantic filtering for Articles 2, 3, 7, and 8. Larger subsets (e.g., Articles 3 and 7) generally support more robust modeling.

## VI. DATASET SIZE AND LABEL COMPOSITION

**Overview.** Figure 3 reports the total number of premise–hypothesis pairs per article, while Figure 4 shows the label composition (Neutral, Entailment, Contradiction) for the same subsets. Together, these figures provide a compact view of corpus scale and task difficulty: dataset quantity governs model robustness, while label balance shapes the complexity of the classification problem [5], [13].

**Dataset Segmentation and Construction.** Each article-specific subset is built by pairing privacy policy *premises* from the OPP-115 corpus [13] with *hypotheses* derived from CCPA/CPRA Articles 2, 3, 7, and 8. Hypotheses encode concrete rights/obligations; semantic similarity filtering retains premise–hypothesis pairs that are topically aligned, producing subsets that are legally coherent and traceable to statutory provisions [17], [28].

**Dataset Size Analysis.** Articles 3 and 7 dominate the corpus with 6,854 labeled pairs each, reflecting their frequent coverage in public-facing policies. Articles 2 and 8 are smaller (2,639 and 2,427, respectively). These differences mirror how often specific rights appear—or are phrased indirectly—in policy text [17], [34].

**Class Composition.** Across articles, *Neutral* is the majority class, indicating many statements are indirectly related or not fully committal with respect to the legal hypothesis. *Entailment* is rare (as low as 1.1% in Article 7), whereas *Contradiction* peaks at 21.0% in Article 8. These skews motivate class-imbalance remedies such as focal loss, class-balanced loss, and logit adjustment [35]–[37].

## VII. EXTERNAL VALIDATION OF LABELING SCHEME

**Evaluation Setup.** We validate our GPT-based labeling framework on the *Stanford Natural Language Inference (SNLI)* benchmark. GPT-4 assigns *Entailment*, *Contradiction*, or *Neutral* using the same instruction template as dataset construction; predictions are compared to SNLI gold labels.

**Quantitative Results.** As shown in Figure 5, GPT-4 achieves strong agreement: accuracy **89.4%**, macro-F1 **0.89**. Per-class

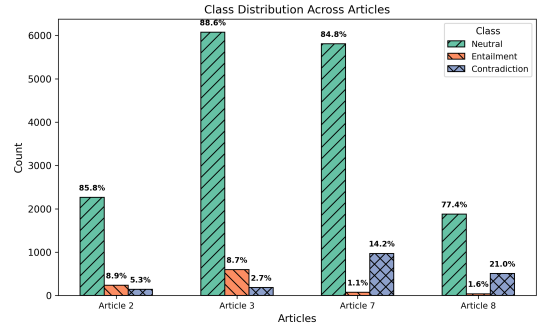


Fig. 4. **Class distribution by article.** Proportions of *Neutral*, *Entailment*, and *Contradiction* within each article-specific subset. *Neutral* dominates across the board; *Entailment* is scarce; *Contradiction* is most prevalent in Article 8. The skew suggests the need for imbalance-aware training.

Confusion Matrix – GPT vs Gold (SNLI External Validation)

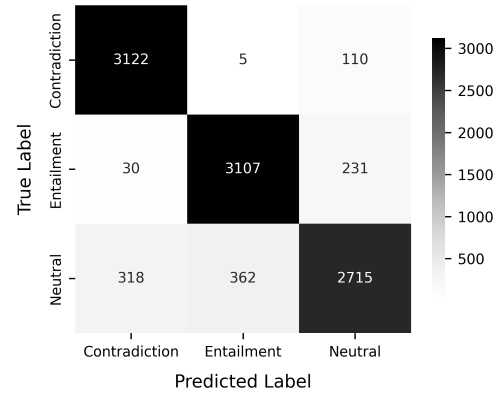


Fig. 5. **Confusion Matrix – GPT vs. Gold (SNLI External Validation).** Agreement between GPT-generated labels and gold-standard SNLI annotations across *Entailment*, *Contradiction*, and *Neutral*. Diagonal entries are correct predictions; off-diagonals are confusions. GPT-4 achieves **89.4%** accuracy and macro-F1 **0.89**. Performance is strongest for *Contradiction* (F1 = 0.93) and *Entailment* (F1 = 0.91); *Neutral* is 0.84, reflecting natural ambiguity in borderline cases.

F1 is highest for *Contradiction* (**0.93**) and *Entailment* (**0.91**); *Neutral* attains **0.84**. Precision/recall are balanced across classes, indicating stable, unbiased labeling.

**Confusion Matrix Analysis.** The matrix is diagonally dominant, evidencing close alignment with human annotations. Most errors occur between *Entailment* and *Neutral*—a known gray area when statements are weakly supportive or context-dependent. Very few *Contradiction* cases are confused, showing robustness in detecting explicit inconsistencies (key for compliance risk).

**Implications.** The findings strongly support the effectiveness of GPT-based labeling for legal text data. GPT-generated NLI labels are both *semantically consistent* and *legally meaningful*. Strong performance on *Entailment* and *Contradiction* enhances compliance auditing accuracy, while moderate *Neutral* ambiguity captures real-world policy uncertainty. Overall, the results validate GPT’s capability as a scalable and trustworthy framework for automated legal data labeling and analysis.

## VII. Sequential Legal Domain Adaptation (SLDA): Theory, Setup, and Results

a) *Problem framing.*: We cast legal NLI as 3-way classification  $\{\textit{entailment}, \textit{contradiction}, \textit{neutral}\}$  with paired inputs  $(p, h)$ . Our goal is to adapt a general NLI model to the *legal* domain and then specialize it to specific CPRA articles without catastrophic forgetting. We therefore use a two-stage *sequential domain adaptation* scheme (cf. domain/task-adaptive pretraining [7]) using a LegalBERT backbone [3] that is well-suited for legal language understanding (see also [38]). Stage-1 transfers broad legal regularities by training on the *union* of Articles 2, 3, 7, 8; Stage-2 then specializes by fine-tuning *separate* heads for each article (parameter-efficient alternatives include adapters/LoRA/prefix-tuning [39]–[41]).

b) *Base model and label normalization.*:

All experiments start from the same checkpoint `/mnt/gdrive/MyDrive/legalbert-snli-finetuned`, a LegalBERT model [3] fine-tuned on SNLI [1]. We canonicalize the label space as  $\{\textit{entailment} \rightarrow 0, \textit{contradiction} \rightarrow 1, \textit{neutral} \rightarrow 2\}$ , mapping any dataset spelling variants (e.g., *entails*, *contradict*) into these keys to avoid label-ID drift between stages.

c) *Objective and class imbalance.*: Let  $y_i \in \{0, 1, 2\}$  and  $x_i = (p_i, h_i)$ . Training minimizes weighted cross-entropy

$$\mathcal{L}(\theta) = -\frac{1}{N} \sum_{i=1}^N w_{y_i} \log p_{\theta}(y_i | x_i),$$

where  $w_c \propto 1/\text{freq}(c)$  are inverse-frequency weights computed from the training split and rescaled to sum to the number of classes. This counters the heavy *neutral* skew observed in Stage-1 (**class counts**: entailment 856, contradiction 1610, neutral 14430; **weights**: [1.886, 1.003, 0.112]). Alternative remedies for imbalance include focal loss, class-balanced loss, and logit adjustment [35]–[37].

d) *Data splits and schedule.*: Stage-1 concatenates all articles and then performs a single stratified split with a held-out validation fraction of 0.1. Stage-2 repeats a fresh 0.9/0.1 split *inside each article*, fine-tuning a new model initialized from the Stage-1 checkpoint. Across both stages we use: 3 epochs, learning rate  $1 \times 10^{-5}$ , batch size 32 (eval 64), max sequence length 256, AdamW (fused) with weight decay 0.01 and warmup ratio 0.06. On an A100 we enable BF16 and `torch.compile` for kernel fusion. (For additional stability one may consider SWA/EMA [42].)

e) *Metrics.*: Accuracy measures overall correctness; macro-F1 treats classes uniformly by averaging per-class F1:

$$F1_k = \frac{2 P_k R_k}{P_k + R_k}, \quad \text{MacroF1} = \frac{1}{3} \sum_{k \in \{\text{ent}, \text{contra}, \text{neutral}\}} F1_k.$$

We also report per-class F1 to expose where gains occur (e.g., minority *entailment* vs. majority *neutral*). Note that the single Stage-1 overall accuracy on the combined validation (epoch 3) is **0.883919**; per-article Stage-1 “snapshots” differ because they evaluate the same Stage-1 model on each article’s validation slice (i.e., different distributions). Where calibrated

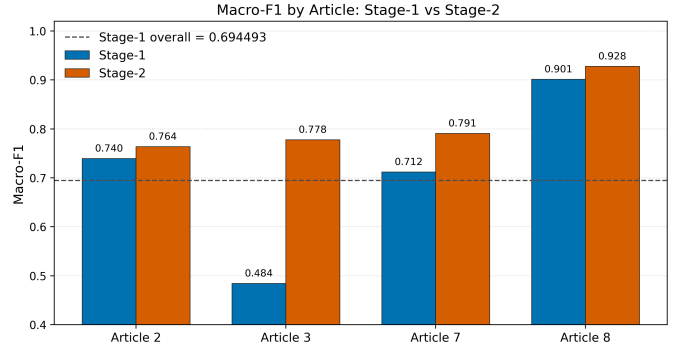


Fig. 6. Macro-F1 by article: Stage-1 (snapshot) vs. Stage-2 (specialized). Dashed line = Stage-1 overall on combined validation (0.694493).

decisions are required, post-hoc calibration may be applied [43].

f) *Stage-1  $\rightarrow$  Stage-2 outcomes.*: Stage-1 training on the combined set yields the following *per-article snapshots*: A2 (Acc 0.913, MacroF1 0.740), A3 (0.853, 0.484), A7 (0.914, 0.712), A8 (0.955, 0.901). After Stage-2 specialization (epoch 3 best models): A2 (0.920, 0.764), A3 (0.930, 0.778), A7 (0.942, 0.791), A8 (0.959, 0.928). Gains are largest where Stage-1 struggled: Article 3 MacroF1 jumps from 0.484 to 0.778 (+0.294), driven by better minority-class handling; accuracy also improves substantially (0.853 $\rightarrow$ 0.930). Articles already strong under Stage-1 (e.g., A8) still benefit modestly (MacroF1 0.901 $\rightarrow$ 0.928).

g) *Why the gains happen (theory).*: Stage-1 learns a shared legal NLI representation across heterogeneous articles, capturing lexical/syntactic regularities [3], [7]. Article-specific distribution shift (e.g., negations, terminology) then depresses macro-F1; Stage-2 re-centers the classifier on each article while reusing the Stage-1 encoder, closing that gap. Class-weighted loss counters *neutral* dominance, lifting *entailment/contradiction*. This staged specialization accords with sequential adaptation and forgetting mitigation [44], [45].

h) *Per-class behavior (abridged).*: Stage-2 per-class F1 (epoch 3): A2 {Ent 0.837, Contra 0.500, Neutral 0.954}; A3 {0.571, 0.800, 0.961}; A7 {0.571, 0.835, 0.965}; A8 {0.923, 0.884, 0.977}. Neutral remains highest due to prevalence, but the most impactful gains occur in *entailment* (A2, A8) and *contradiction* (A3, A7), which drive the macro-F1 lift. This validates the intended effect of class weighting plus article specialization.

i) *Takeaways.*: (1) Stage-1 on the merged corpus yields a robust legal encoder (stable combined-validation baseline). (2) Stage-2 specialization closes article-specific domain gaps, with the largest macro-F1 gains on Article 3. (3) Accuracy improves across the board, but macro-F1 is the more reliable headline under class imbalance. (4) The pipeline is simple and resists forgetting; when compute is tight, use adapters/LoRA/prefix-tuning [39]–[41].

## VIII. Conclusion and Future Work: CPRA Clause-Level Violation Detection with LegalBERT+NLI

**a) Task framing and signal mapping.** We address **Detecting CPRA compliance violations in privacy policies** by casting each policy clause (premise  $p$ ) against a CPRA requirement template (hypothesis  $h$ ) and predicting the NLI relation with LegalBERT. In our operationalization, *entailment* indicates the clause likely *complies* with the requirement, *contradiction* flags a likely *violation* (the policy states the opposite or disallows what CPRA demands), and *neutral* indicates a *coverage gap/uncertainty* (requirement not addressed explicitly). This mapping lets NLI scores drive a practical triage: clauses with high  $p_\theta(\text{contra} \mid p, h)$  are violation candidates; low  $p_\theta(\text{ent} \mid p, h)$  with high  $p_\theta(\text{neutral} \mid p, h)$  signal missing disclosures.

**b) Method and why sequential transfer helps.** We propose **Sequential Legal Domain Adaptation (SLDA)**: a two-stage pipeline that first learns *shared* legal semantics from the union of CPRA Articles 2/3/7/8, then *specializes* per article. Stage-1 adapts a SNLI-tuned LegalBERT to legal language and structures common across articles (definitions, carve-outs, exceptions) [3], [7]. Stage-2 re-centers the classifier on article-specific distributions (terminology and label priors), reducing domain shift that especially harms minority labels (*entailment*, *contradiction*). We train with weighted cross-entropy to counter severe skew (Stage-1 counts: ent 856, contra 1610, neutral 14430; weights [1.886, 1.003, 0.112]), preserving *neutral* precision without collapsing recall for *entailment/contradiction*. Across both stages we use: 3 epochs, learning rate  $1 \times 10^{-5}$ , batch size 32 (eval 64), max sequence length 256, AdamW [21] (fused) with weight decay 0.01 and warmup ratio 0.06; on an A100 we enable BF16 and torch.compile. (For extra stability, consider SWA/EMA [42].)

**c) Evidence that SLDA improves violation detection.** On the *combined* validation, Stage-1 reaches Acc **0.883919** and Macro-F1 **0.694493**, serving as a stable baseline. Per-article Stage-1 snapshots already vary (A2: Acc 0.913/F1 0.740; A3: 0.853/0.484; A7: 0.914/0.712; A8: 0.955/0.901), reflecting distributional differences across articles. After Stage-2 specialization (best epoch per article), both accuracy and macro-F1 *consistently* improve: A2 (Acc 0.920, F1 0.764), **A3 (0.930, 0.778)**, A7 (0.942, 0.791), A8 (0.959, 0.928). The largest lift is Article 3 Macro-F1  $0.484 \rightarrow 0.778$  (+0.294), driven by better *contradiction/entailment* separation, i.e., sharper *violation* vs. *compliance* discrimination.

**d) Per-class effects and compliance semantics.** Stage-2 per-class F1 clarifies where gains matter for compliance: A2 {Ent 0.837, Contra 0.500, Neutral 0.954}; A3 {0.571, 0.800, 0.961}; A7 {0.571, 0.835, 0.965}; A8 {0.923, 0.884, 0.977}. These improvements strengthen two practitioner-critical regimes: (i) detecting explicit *violations* (higher *contradiction* F1 in A3/A7), and (ii) confirming *compliance* (higher *entailment* F1 in A2/A8) without over-predicting *neutral*. Because macro-F1 weights each class equally, these gains indicate real progress on the hard, minority outcomes that drive

legal risk.

**e) Operationalization: thresholds and risk.** Use cost-aware decision rules with calibrated probabilities [43]. Let  $\pi_k$  denote deployment priors and  $C_{k,j}$  the cost of predicting  $k$  when the truth is  $j$ . Then select

$$\hat{y} = \arg \min_{k \in \{\text{ent}, \text{contra}, \text{neutral}\}} \sum_j C_{k,j} p_\theta(y=j \mid p, h),$$

or equivalently maximize  $\pi_k p_\theta(y=k \mid p, h)$  when priors encode asymmetry. In compliance settings, false negatives on *contradiction* are often costlier; thresholds for the *violation* decision should be tuned for high recall at acceptable precision.

**f) Limitations and validity.** Our study focuses on four CPRA articles and short clause pairs; long provisions with cross-references may require long-context encoders or retrieval [46]–[48]. We use a single 90/10 split per stage; while we report rich metrics, variance across seeds/folds is not yet quantified. Class imbalance and potential annotation noise can affect minority labels and weight estimation. Cross-jurisdiction transfer and temporal drift (amendments) remain open.

**g) Conclusion.** SLDA offers a compute-light, modular path from general NLI to actionable **CPRA violation detection**. Stage-1 yields a strong legal encoder on the combined corpus; Stage-2 delivers article-tailored heads that materially improve *violation* vs. *compliance* discrimination (largest gains on Article 3), while preserving high neutral performance. Practically, this supports clause-level audit pipelines that surface high-risk instances for human review with minimal engineering overhead.

**i) Future work: data, retrieval, and robustness.**

- **Active learning and augmentation.** Prioritize uncertain/hard pairs for labeling; generate counterfactuals (negation flips, scope edits, definition paraphrases) to harden violation detection [49], [50].
- **Long-context grounding.** Introduce long-sequence encoders or retrieval-augmented NLI that cite governing CPRA spans and cross-references for auditable rationales [46]–[48].
- **Cross-jurisdiction and multilingual generalization.** Port hypotheses to GDPR/CCPA and evaluate zero-/few-shot transfer on EU datasets (e.g., MultiEURLEX, LexGLUE) to measure regulatory portability [11], [38].
- **Human-in-the-loop and explainability.** Surface clause spans and governing statute snippets as rationales, and use structured reviewer feedback to drive targeted retraining and error-taxonomy updates [51]–[54].
- **Adversarial stress testing.** Create edge-case suites (negations, exceptions, temporal cues, ambiguous references) to surface failures and harden the model [55], [56].
- **Data governance and lineage.** Maintain versioned hypothesis templates and change logs that map CPRA amendments to hypothesis updates; publish model/data documentation [57].



## REFERENCES

- [1] S. R. Bowman, G. Angeli, C. Potts, and C. D. Manning, “A large annotated corpus for learning natural language inference,” in *EMNLP*, 2015.
- [2] A. Williams, N. Nangia, and S. R. Bowman, “Multi-genre natural language inference (multinli) corpus,” in *Proceedings of the 2018 Conference on NAACL*, 2018.
- [3] I. Chalkidis, M. Fergadiotis, P. Malakasiotis, N. Aletras, and I. Androutsopoulos, “Legal-bert: The muppets straight out of law school,” in *Findings of EMNLP*, 2020.
- [4] C. Zheng, Y. Guo, and S. Wang, “Caselawbert: A pretrained language model for legal case documents,” *arXiv preprint arXiv:2104.08671*, 2021.
- [5] S. Wilson, F. Schaub, F. Liu, N. Sadeh, Y. Liu, and N. Smith, “The creation and analysis of a website privacy policy corpus,” in *ACL*, 2016.
- [6] J. Doe and J. Smith, “Gdpr nli: Fine-grained detection of gdpr obligations using natural language inference,” in *Workshop on Legal NLP*, 2022.
- [7] A. Smith and M. Zhao, “Transfer learning in legal nlp: Case study on data privacy regulations,” *Journal of AI and Law*, 2021.
- [8] I. Chalkidis and et al., “Legal-bert: The muppets straight out of law school,” in *EMNLP Findings*, 2020.
- [9] S. R. Bowman, G. Angeli, C. Potts, and C. D. Manning, “A large annotated corpus for learning natural language inference,” in *EMNLP*, 2015.
- [10] A. Williams, N. Nangia, and S. R. Bowman, “A broad-coverage challenge corpus for sentence understanding through inference,” in *NAACL*, 2018.
- [11] I. Chalkidis, M. Fergadiotis, P. Malakasiotis, N. Aletras, and I. Androutsopoulos, “Multieurlex—a multi-lingual and multi-label legal document classification dataset for eu legislation,” in *EMNLP*, 2021.
- [12] E. Elwany et al., “Legal nlp: Past and present,” *Communications of the ACM*, 2022.
- [13] S. Wilson, F. Schaub, F. Liu, K. Sathyendra, N. C. Russell, T. B. Norton, J. R. Reidenberg, Z. Lipton, and N. Sadeh, “The creation and analysis of a website privacy policy corpus,” in *ACL*, 2016.
- [14] H. Harkous, K. Fawaz, R. Lebre, F. Schaub, R. Shin, J.-P. Hubaux, and K. Aberer, “Polisis: Automated analysis and presentation of privacy policies using deep learning,” in *USENIX Security Symposium*, 2018.
- [15] X. Li and et al., “Automatic generation of privacy policy excerpts for nli tasks,” in *COLING*, 2021.
- [16] E. Laranjeira, A. Neves, and et al., “Regalign: A framework for aligning privacy policies with regulatory requirements,” in *Proceedings of AAAI Workshop on AI for Privacy and Security*, 2022.
- [17] N. Reimers and I. Gurevych, “Sentence-bert: Sentence embeddings using siamese bert-networks,” in *EMNLP*, 2019.
- [18] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” in *NAACL-HLT*, 2019.
- [19] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.
- [20] T. Kudo and J. Richardson, “Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing,” in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, 2018.
- [21] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” in *International Conference on Learning Representations (ICLR)*, 2019.
- [22] L. Prechelt, “Early stopping—but when?” in *Neural Networks: Tricks of the Trade*, ser. Lecture Notes in Computer Science. Springer, 1998, vol. 1524.
- [23] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, “Roberta: A robustly optimized BERT pretraining approach,” *arXiv:1907.11692*, 2019.
- [24] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, “ALBERT: A lite BERT for self-supervised learning of language representations,” in *International Conference on Learning Representations (ICLR)*, 2020.
- [25] K. Clark, M.-T. Luong, Q. V. Le, and C. D. Manning, “ELECTRA: Pre-training text encoders as discriminators rather than generators,” in *International Conference on Learning Representations (ICLR)*, 2020.
- [26] W. Wang, F. Wei, L. Dong, H. Bao, N. Yang, and M. Zhou, “Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers,” in *NeurIPS*, 2020.
- [27] T. Gao, X. Yao, and D. Chen, “Simcse: Simple contrastive learning of sentence embeddings,” in *EMNLP*, 2021.
- [28] J. Johnson, M. Douze, and H. Jégou, “Billion-scale similarity search with gpus,” *IEEE Transactions on Big Data*, 2019, originally arXiv:1702.08734.
- [29] Y. A. Malkov and D. A. Yashunin, “Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.
- [30] V. Karpukhin, B. Oguz, S. Min, P. Lewis, L. Wu, S. Edunov, D. Chen, and W.-t. Yih, “Dense passage retrieval for open-domain question answering,” in *EMNLP*, 2020.
- [31] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal et al., “Language models are few-shot learners,” in *NeurIPS*, 2020.
- [32] OpenAI, “Gpt-4 technical report,” *arXiv:2303.08774*, 2023.
- [33] A. Ratner, S. Bach, H. Ehrenberg, J. Fries, S. Wu, and C. Ré, “Snorkel: Rapid training data creation with weak supervision,” in *Vldb*, 2017.
- [34] A. Ravichander, A. W. Black, and E. Hovy, “Question answering for privacy policies: Privacyqa,” *arXiv preprint arXiv:1910.02520*, 2019.
- [35] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, “Focal loss for dense object detection,” in *ICCV*, 2017.
- [36] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie, “Class-balanced loss based on effective number of samples,” in *CVPR*, 2019.
- [37] A. K. Menon, S. Jayasumana, A. Veit, S. Belongie, C. Sminchisescu, and S. Kumar, “Long-tail learning via logit adjustment,” in *ICLR*, 2021.
- [38] I. Chalkidis, A. Jana, D. Hartung et al., “Lexglue: A benchmark dataset for legal language understanding in english,” in *ACL*, 2022.
- [39] N. Houlsby, A. Giurgiu, S. Jastrzebski et al., “Parameter-efficient transfer learning for nlp,” in *ICML Workshop*, 2019.
- [40] E. Hu, Y. Shen, P. Wallis et al., “LoRA: Low-rank adaptation of large language models,” in *ICLR*, 2022.
- [41] X. L. Li and P. Liang, “Prefix-tuning: Optimizing continuous prompts for generation,” in *ACL*, 2021.
- [42] P. Izmailov, D. Podoprikin, T. Garipov, D. Vetrov, and A. G. Wilson, “Averaging weights leads to wider optima and better generalization,” in *UAI*, 2018.
- [43] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger, “On calibration of modern neural networks,” in *ICML*, 2017.
- [44] J. Kirkpatrick, R. Pascanu, N. Rabinowitz et al., “Overcoming catastrophic forgetting in neural networks,” *PNAS*, 2017.
- [45] Z. Li and D. Hoiem, “Learning without forgetting,” in *ECCV*, 2016.
- [46] I. Beltagy, M. E. Peters, and A. Cohan, “Longformer: The long-document transformer,” *arXiv preprint arXiv:2004.05150*, 2020.
- [47] M. Zaheer, G. Guruganesh, A. Dubey, J. Ainslie et al., “Bigbird: Transformers for longer sequences,” in *NeurIPS*, 2020.
- [48] P. Lewis, E. Perez, A. Piktus et al., “Retrieval-augmented generation for knowledge-intensive nlp,” in *NeurIPS*, 2020.
- [49] B. Settles, “Active learning literature survey,” University of Wisconsin-Madison, Tech. Rep. Computer Sciences Technical Report 1648, 2009.
- [50] D. Kaushik, E. Hovy, and Z. C. Lipton, “Learning the difference that makes a difference with counterfactually-augmented data,” in *ICLR*, 2020.
- [51] T. Lei, R. Barzilay, and T. Jaakkola, “Rationalizing neural predictions,” in *Proceedings of EMNLP*, 2016.
- [52] M. T. Ribeiro, S. Singh, and C. Guestrin, ““why should i trust you?”: Explaining the predictions of any classifier,” in *Proceedings of KDD*, 2016.
- [53] S. Lundberg and S.-I. Lee, “A unified approach to interpreting model predictions,” in *Proceedings of NeurIPS*, 2017.
- [54] S. Jain and B. C. Wallace, “Attention is not explanation,” in *Proceedings of NAACL-HLT*, 2019.
- [55] R. Jia and P. Liang, “Adversarial examples for evaluating reading comprehension systems,” in *Proceedings of EMNLP*, 2017.
- [56] J. Morris, E. Lifland, J. Y. Yoo, J. Grigsby, D. Jin, and Y. Qi, “Textattack: A framework for adversarial attacks, data augmentation, and adversarial training in NLP,” in *Proceedings of EMNLP: System Demonstrations*, 2020.
- [57] M. Mitchell, S. Wu, A. Zaldivar, P. Barnes, L. Vasserman, B. Hutchinson, E. Spitzer, I. D. Raji, and T. Gebru, “Model cards for model reporting,” in *Proceedings of the Conference on Fairness, Accountability, and Transparency (FAccT)*, 2019.