

Legal Clause Similarity — Baseline Report

Dataset & Splits

- Train: 101,905 clauses across 394 labels (203,810 pairs)
- Validation: 21,846 clauses across 394 labels (43,692 pairs)
- Test: 22,061 clauses across 394 labels (44,122 pairs)
- Negative/positive sampling ratio per split: 1.0:1

Architectures

1) BiLSTM Siamese

- Embedding: vocab 30k, dim 128, padding idx 0
- Encoder: Bidirectional LSTM (hidden 128 each direction) → mean pooling
- Projection: LayerNorm → Linear → ReLU → Dropout(0.3)
- Comparator: $\text{concat}([h_1, h_2, |h_1 - h_2|, h_1 \odot h_2]) \rightarrow \text{MLP}(256 \rightarrow 1)$

2) Attentive BiGRU Siamese

- Embedding: same as above
- Encoder: Bidirectional GRU (hidden 128) + additive attention (masked softmax) to get a single vector
- Comparator: $\text{concat}([h_1, h_2, |h_1 - h_2|, h_1 \odot h_2]) \rightarrow \text{MLP}(256 \rightarrow 1)$

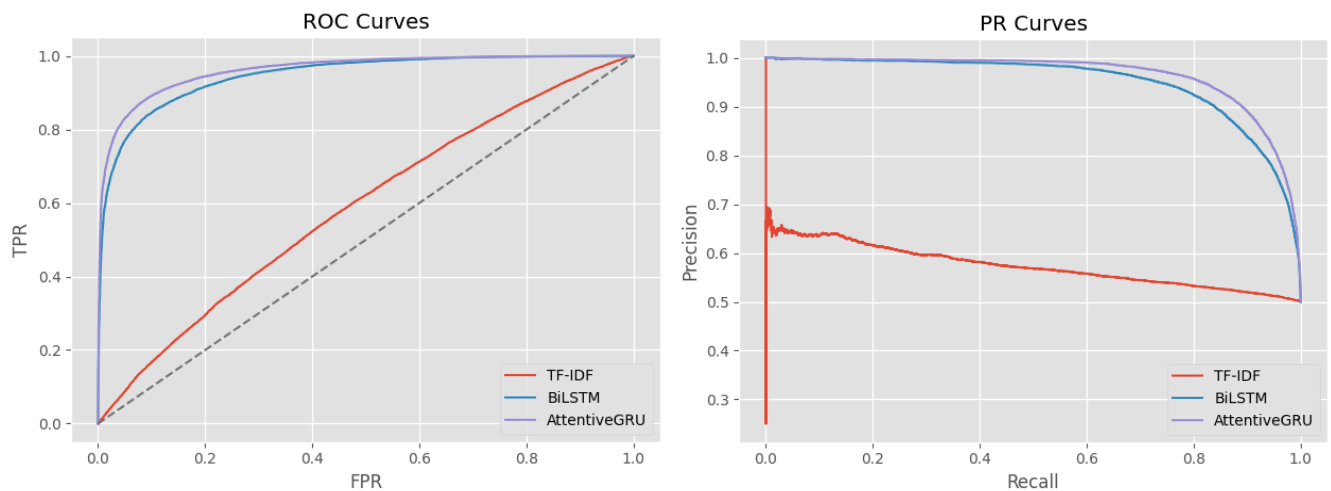
3) TF-IDF + Logistic Regression (baseline)

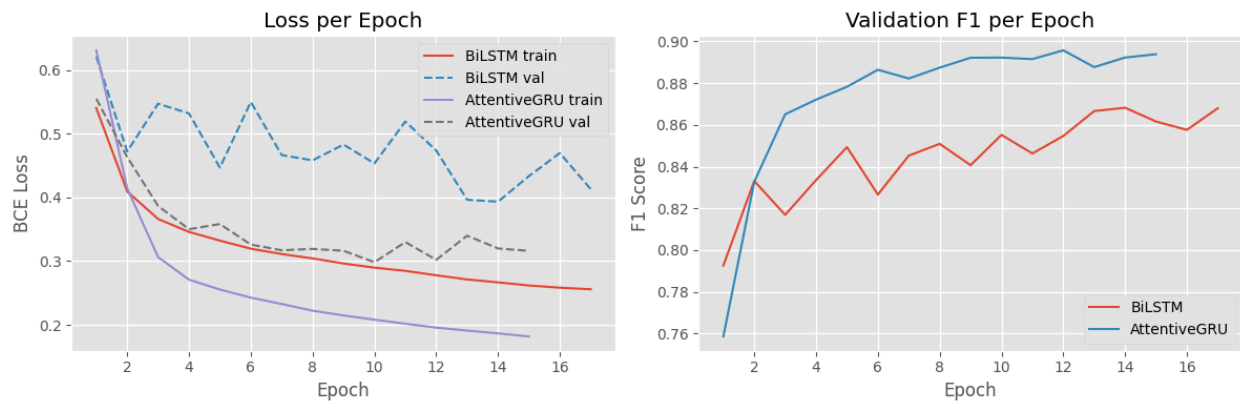
- Text featurization: word 1–2 grams, max 50k features, lower-cased, token pattern `[A-Za-z ']+`
- Classifier: Logistic Regression (liblinear), C tuned on validation

Training Setup

- Max sequence length: 150
- Batch sizes: train 64, eval 128
- Optimizer: Adam (lr=0.001, weight_decay=0.0001), gradient clip=2.0
- Training schedule: up to 50 epochs with early stopping (patience=3) on validation F1
- Loss: BCEWithLogitsLoss; thresholds chosen by best validation F1

Training Graphs





Test Performance

Model	Acc	Prec	Rec	F1	ROC-AUC	PR-AUC	Thr	Train Time (s)	Loss
BiLSTM	0.868	0.857	0.883	0.870	0.946	0.949	0.681	9490.024	0.392
	614	923	550	548	642	848	313	658	995
AttentiveGRU	0.895	0.902	0.887	0.895	0.961	0.964	0.645	17230.37	0.299
RU	993	432	992	154	654	579	903	6958	011
TFIDF	0.507	0.503	0.988	0.667	0.586	0.573	0.217	38.75307	—
	411	776	668	452	827	302	144	8	

Comparative Analysis

- By F1:

1. AttentiveGRU (0.895154),
2. BiLSTM (0.870548),
3. TFIDF (0.667452).

- Relative F1 lift over TF-IDF: BiLSTM = **30.4%**, AttentiveGRU = **34.2%**.

- By accuracy: AttentiveGRU (0.895993) > BiLSTM (0.868614) » TFIDF (0.507411).

- **By area metrics:** AttentiveGRU leads on ROC-AUC (0.961654) and PR-AUC (0.964579); BiLSTM is close; TF-IDF lags far behind.
- **By training time:** TF-IDF is the fastest baseline ($\approx 1\times$). BiLSTM is $\sim 244.8\times$ TF-IDF; AttentiveGRU is $\sim 444.8\times$ TF-IDF. AttentiveGRU vs BiLSTM train-time ratio $\approx 1.82\times$.
- **Error tendencies:** false positives from shared boilerplate with different scope/carve-outs; false negatives from modality shifts (“may” \rightarrow “shall”) or reordered conditions. Attention helps recover long-range cues, improving recall and PR-AUC.

Takeaways

- Attention pooling on BiGRU yields the strongest overall ranking and classification performance on legal clauses, at the cost of $\sim 1.82\times$ the BiLSTM training time.
- The BiLSTM Siamese remains competitive and more compute-efficient; good default when resources are tight.
- Pure lexical TF-IDF fails on paraphrases and structural semantics; keep it only as a sanity baseline.