



TEXAS A&M UNIVERSITY
Engineering

SemEval 2023 Task 6: Rhetorical Roles Prediction (RR): Legal Document Segmentation

Presented By:
Rohan, Hemal, Shubham, Upasana, Abhishek

- AI applications in the legal domain such as Judgment summarizing, judgment outcome prediction, precedent search, etc require the segmentation of long and unstructured legal documents into semantically coherent text segments.
 - The problem we aim to solve is to segment a given legal document into topically semantic and coherent units and assign a label to those segments which are referred to as Rhetorical Roles (RR), by predicting the RR label for each sentence.
 - The legal document needs to be considered as sequential sentences and the classification of sentences with the prediction of the RR with single label multiple classes for each sentence is needed to be achieved.
-

- In populous countries (e.g., India), pending legal cases have grown exponentially.
- An automated system can help in many intermediate tasks and expedite the process.

IN THE COURT OF THE V ADDL SESSIONS JUDGE, MYSORE. Dated this the 23rd day of May 2013 ... The Petitioner is a businessman and he is permanent resident of Mysore City... On behalf of the Prosecution the learned Public Prosecutor has filed objection to the bail Petition stating that, there ...Now, the points that arise for consideration of the Court are: 1. Whether the Petitioner has made out sufficient grounds to release him on Anticipatory Bail? ... Heard the arguments advanced by the learned advocate for the Petitioner and the learned Public Prosecutor...Considering all these aspects, the Court is of the view that, ...Point No.2: For the foregoing reasons and in view of my above discussions, I proceed to pass the following ...The High Court by its order dated October 26, 1982 set aside the order of the Tribunal and also the assessment on the ground ...The petitioners are falsely implicated and the charge sheet has been filed against the petitioners merely ...My findings on the above points are as follows: Point No.1 : In the Positive Point No.2 : As per final order for the following...In a decision reported in (2013) 1 KCCR 334 case of K.Ramachandra Reddy Vs. State of Karnataka by the Station House Officer...The decision of the Andhra Pradesh High Court ... are not relevant for purposes of deciding the question which has arisen before us...

City... Fact

On behalf of the Prosecution the learned Public Prosecutor has filed objection to the bail Petition stating that, there ... Arg by Respondent

Now, the points that arise for consideration of the Court are: 1. Whether the Petitioner has made out sufficient grounds to release him on Anticipatory Bail? ... Issue

Heard the arguments advanced by the learned advocate for the Petitioner and the learned Public Prosecutor... None

Considering all these aspects, the Court is of the view that, ... Ratio

Point No.2: For the foregoing reasons and in view of my above discussions, I proceed to pass the following ... Ruling by present court

The High Court by its order dated October 26, 1982 set aside the order of the Tribunal and also the assessment on the ground ... Ruling by lower court

The petitioners are falsely implicated and the charge sheet has been filed against the petitioners merely ... Arg by Petitioner

My findings on the above points are as follows: Point No.1 : In the Positive Point No.2 : As per final order for the following... Analysis

In a decision reported in (2013) 1 KCCR 334 case of K.Ramachandra Reddy Vs.

- <https://github.com/Legal-NLP-EkStep/rhetorical-role-baseline>
 - Number of documents in Train dataset: 245
 - Number of documents in Validation dataset size: 30
 - Each document contains sentences which are classified into rhetorical roles
 - Total sentences for Training: 28986
 - Total sentences for Validation: 2890
 - Test dataset will be released in Jan 23
-

1. Preamble: Court name, the details of parties, lawyers and judges' names
 2. Facts: Chronology of events
 3. Ruling by Lower Court: Judgments given by the lower courts
 4. Issues: Such Legal Questions Framed
 5. Argument by Petitioner
 6. Argument by Respondent,
 7. Analysis: Courts discussion on the evidence
 8. Statute: discussion on Established laws
 9. Precedent Relied: discussions relied
 10. Precedent Not Relied: Discussions not relied
 11. Ratio of the decision: result of the analysis by the court.
 12. Ruling by Present Court: Final decision + conclusion + order of the Court
 13. None: Sentence does not fit any of the above
-

- Past works have centered on the creation of annotated corpora and the task of automatic rhetorical role labeling. Venturi (2012) developed a corpus, TEMIS of 504 sentences annotated both syntactically and semantically. The work of *Wyner et al. (2013)* focuses on the process of annotation and conducting inter-annotator studies.
 - One approach was followed using Conditional Random Fields (CRF) with handcrafted features conducted Document segmentation of U.S. court documents to segment the documents into functional and issue-specific parts. Paper *Saravanan et al. (2008)* developed automatic labeling of rhetorical roles, where CRFs were used to label seven rhetorical roles. Bhatia (2014) created Genre Analysis of Legal Texts to create seven rhetorical categories. *Nejadgholi et al. (2017)* developed a method for identification of factual and non factual sentences using fastText.
 - Paper Walker et al. (2019) compares different ML approaches and rule-based scripts for rhetorical role identification. The baseline used is a closely related work by paper *Bhattacharya et al. (2019)*, where they label rhetorical roles in Indian Supreme Court documents by using the BiLSTM-CRF model with sent2vec features. Baseline model is a multi-task learning-based model for RR prediction that outperforms the system of paper *Bhattacharya et al. (2019)*.
-

- The text segmentation and classification problem are modeled as a single-sentence prediction task. Given a legal document D , containing the sentences $[s_1, s_2, \dots, s_n]$, the task of rhetorical role prediction is to predict the label (or role) y_i for each sentence $s_i \in D$.
 - Multiple approaches have been tried and tested on the given dataset, where different tokenizers and encodings were used to generate embeddings of legal text. These embeddings are then classified as single-label multiclass rhetorical roles using various classifiers. Our final model which surpasses the metrics produced by the baseline model used in the paper follows MultiTask Learning.
-

Model - Naive Bayes



TEXAS A&M UNIVERSITY
Engineering

Naive Bayes classifiers is a family of classifiers with each of them sharing a common principle of all feature pairs being classified being independent from each other. These algorithms work on the base principle of the Bayes' Theorem which can be mathematically stated as follows:

$$P(X|Y) = (P(Y|X) * P(X)) / P(Y)$$

$P(X|Y)$: probability of event X given event Y is true.

$P(Y|X)$: probability of event Y given event X is true.

$P(X)$, $P(Y)$: independent probabilities of events X and Y respectively.

- The main motivation of doc2vec is to represent document into numeric value.
 - Doc2Vec model is used to create a vectorized representation of a group of words taken collectively as a single unit, i.e., a document.
 - Documents, unlike words, do not possess logical structures. Hence, we add another vector called paragraph ID/ paragraph vector.
 - Distributed Memory version of Paragraph Vector (PV-DM) works analogous to the CBOW model of word2vec while Distributed Bag of Words version of Paragraph Vector (PV-DBOW) to the skip-gram model of word2vec.
-

Model - Wiki BM25 Lexical Model



TEXAS A&M UNIVERSITY
Engineering

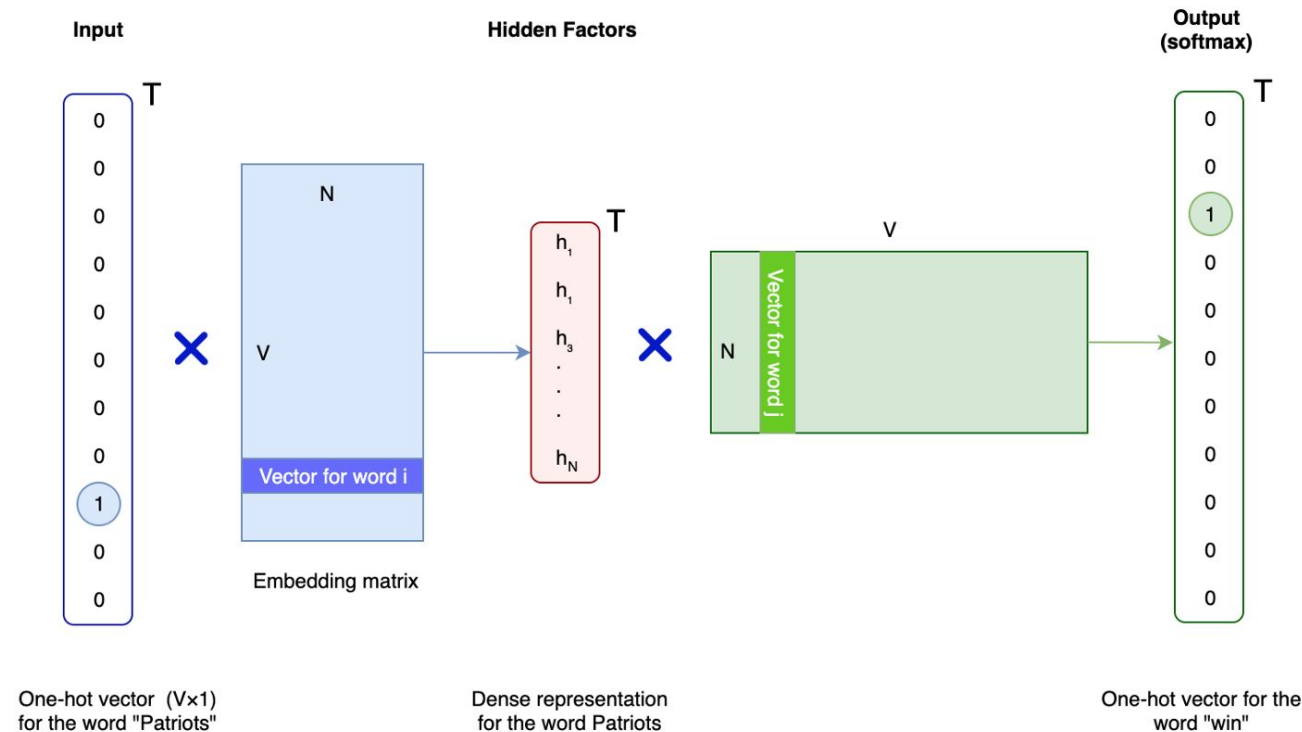
- Tokenizer and query encoder of bm25 improved TF-IDF method, with refinements of TF and IDF components is used to generate the embeddings of the legal text.
 - Classifier Multinomial logistic regression is chosen to perform the classification on word embeddings after training and fitting the model with training data of text embeddings.
 - **Analysis:** As multinomial logistic regression is generalized linear model, it cannot fit non-linear data points. The embedding vector data can be seen as a distribution of non-linear points hence, it does not draw a perfect line as a separation between the data points.
-

Model - GloVe Embedding + FCN



TEXAS A&M UNIVERSITY
Engineering

GloVe Embeddings are a type of word embedding that encode the co-occurrence probability ratio between two words as vector differences.

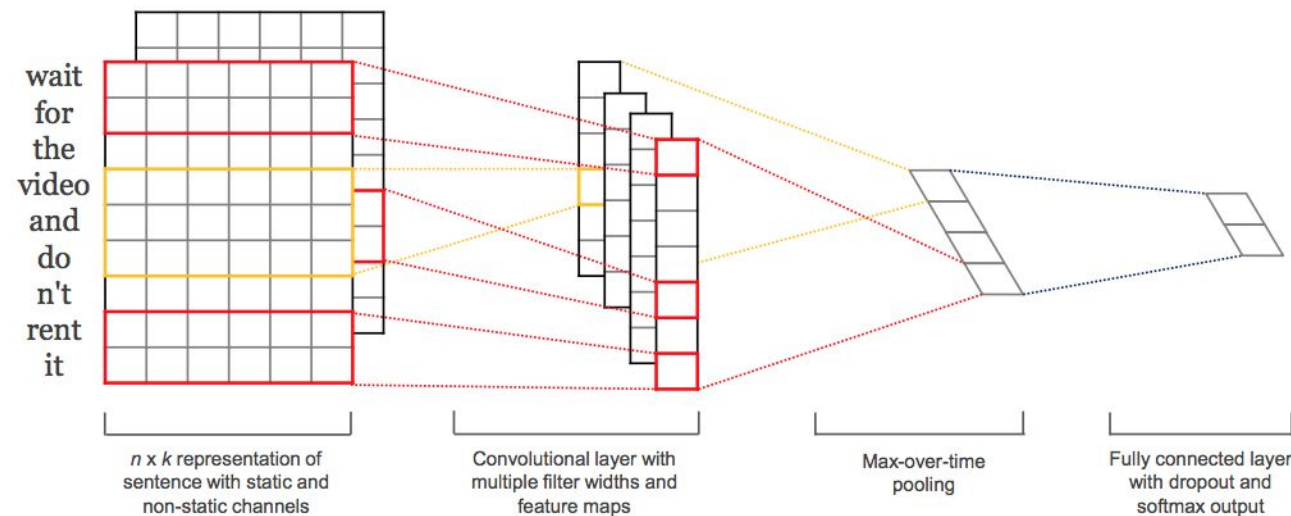


Model - Glove Embedding + 1DCNN



TEXAS A&M UNIVERSITY
Engineering

- convolutional kernel is sliding window, only its job is to look at embeddings for multiple words, rather than small areas of pixels in an image
- similar to representing an n-gram



T5, or Text-to-Text Transfer Transformer, is a [Transformer](#) based architecture that uses a text-to-text approach. Every task – including translation, question answering, and classification – is cast as feeding the model text as input and training it to generate some target text. This allows for the use of the same model, loss function, hyperparameters, etc. across our diverse set of tasks. The changes compared to [BERT](#) include:

- adding a *causal* decoder to the bidirectional architecture.
- replacing the fill-in-the-blank cloze task with a mix of alternative pre-training tasks.

For our application we have used T5 encoder of several sizes including T5-base, T5-large, T5-xl, and T5-xxl.

FAISS (Facebook AI Similarity Search) is a library that allows developers to quickly search for embeddings of multimedia documents that are similar to each other. It solves limitations of traditional query search engines that are optimized for hash-based searches, and provides more scalable similarity search functions.

T5 Encoder + (PLDA)



1. PLDA is a generative model where we assume that the data samples X of a class are generated from a Gaussian distribution. The mean of Gaussian represents the class variable y is generated from another Gaussian distribution called as prior.
 1. For the task of recognition, we can compare two examples of unseen classes using PLDA scores by comparing likelihood of examples from same class vs likelihood of examples from different.
 2. We can perform clustering of examples into classes using PLDA scores between all pairs of examples from entire set of examples.
-

T-5 encoder with BiLSTM



TEXAS A&M UNIVERSITY
Engineering

- Used T-5-base model to get embeddings for each sentence and used a bidirectional LSTM model followed by a linear layer to get labels.

```
LSTM(768, 64, num_layers=3, batch_first=True, bidirectional=True)  
Linear(in_features=128, out_features=13, bias=True)
```

Observations made in Baseline



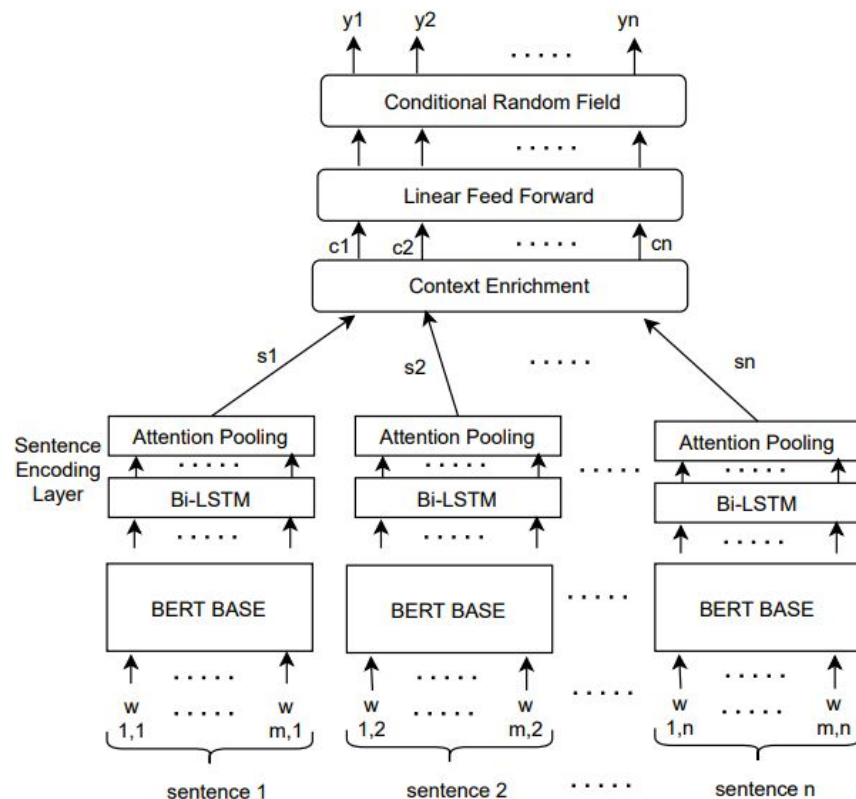
TEXAS A&M UNIVERSITY
Engineering

- Rhetorical role labels do not change abruptly across sentences in a document, and the text tends to maintain topical coherence. Given the label y for a sentence s_i in the document, we hypothesize that the chances of shift (change) in the label for the next sentence s_{i+1} are low.
 - LegalBERT fine-tuned on legal documents performs better than Bert-base-uncased for this task as it is trained on domain specific documents.
-

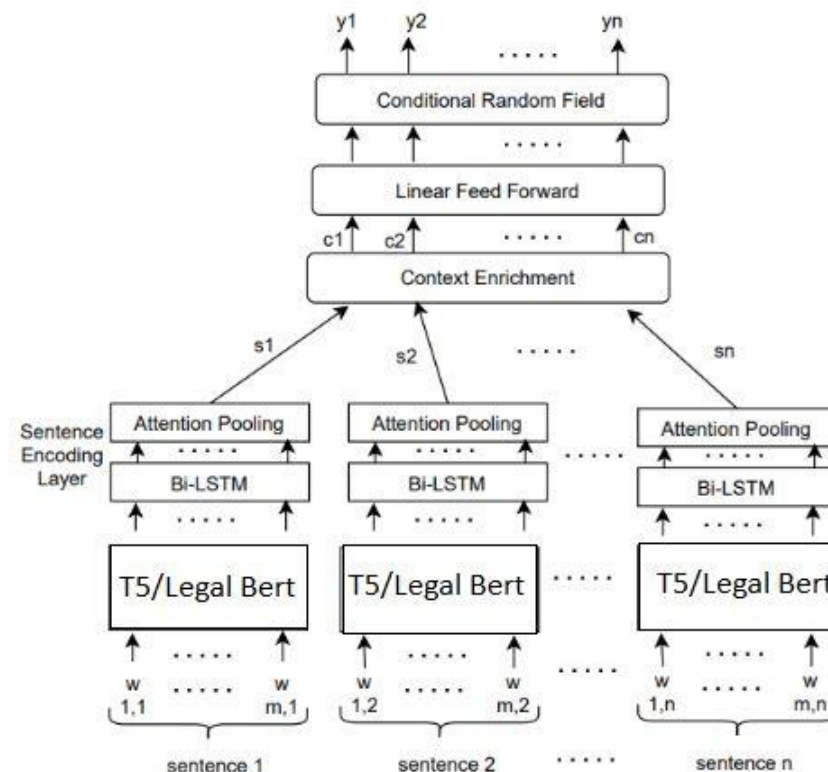
Baseline and Our Approach



TEXAS A&M UNIVERSITY
Engineering



Baseline



Our Approach

The key idea behind the baseline preprocessing approach was to ensure that the sentences of a single document stays in a single batch. This was based on the understanding that the next role is 88% of the times same as that of the previous role.

Results



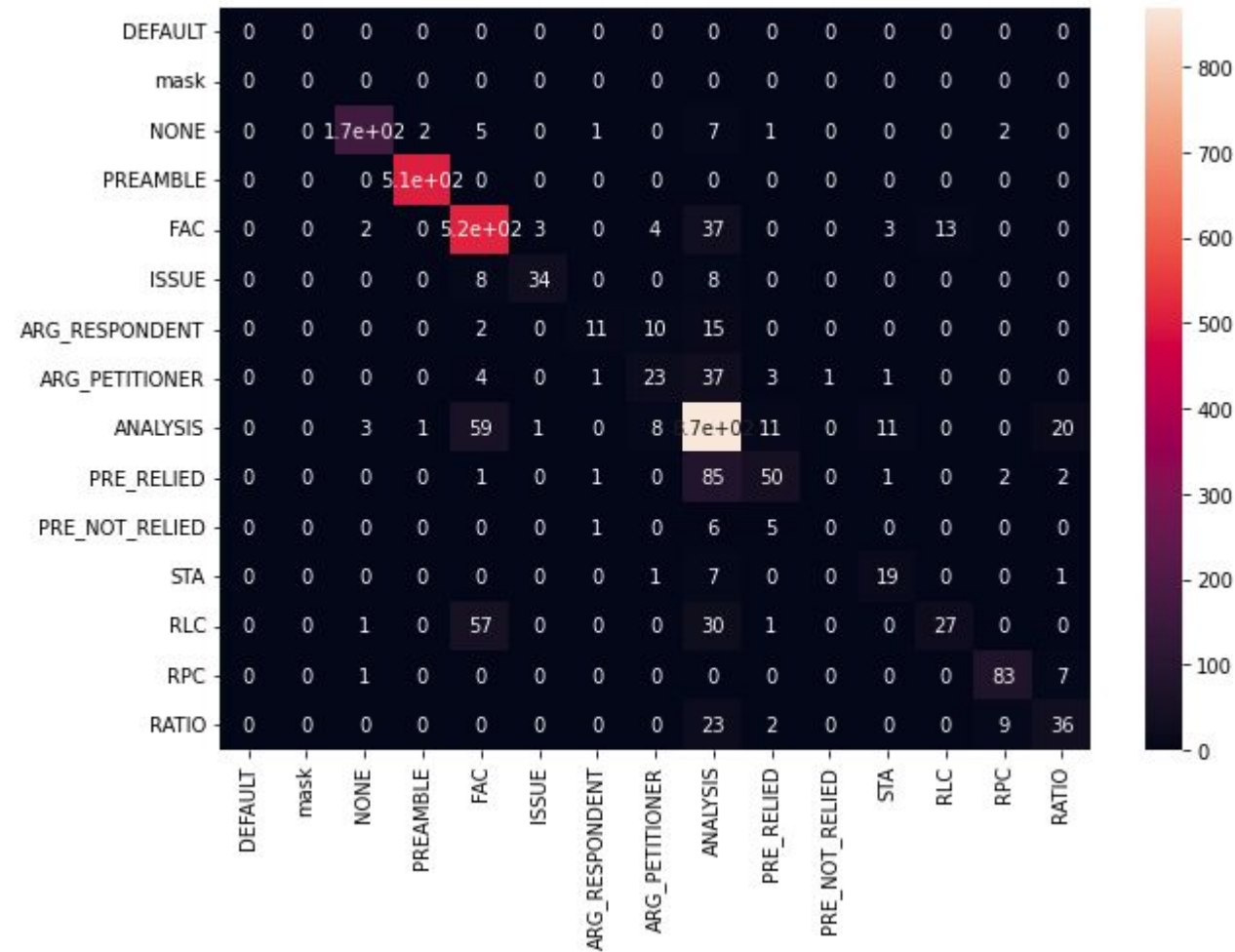
TEXAS A&M UNIVERSITY
Engineering

Model	Macro F1 on Validation	Precision on Validation
Naive Bayes	0.3202008672595919	0.5217426043956311
Doc2Vec	0.35678367302630026	0.46079523508796877
Wiki BM25 Lexical Model + LR	0.382966	0.51591
Glove Embedding + FCN	0.2694453017795601	0.3142378360218918
Glove Embedding + 1DCNN	0.3136997664692862	0.35628853795417964
T5 Encoder + (FAISS)	0.3524619251870554	0.4806228373702422
T5 Encoder + (PLDA)	0.4319042539617168	0.5
T5-base + LSTM	0.3923875427246094	0.4703375424246297
BaseLine	0.5962514492279818	0.8082667592914207
Our Approach (Best) - Baseline modified with T5-large tokenizer and encoder	0.6161397982433457	0.8166029871483154

Confusion Matrix for Best Model



TEXAS A&M UNIVERSITY
Engineering



Different Configurations



TEXAS A&M UNIVERSITY
Engineering

Model	Macro F1 on Validation	Precision on Validation
Sentence Transformer (paraphrase-xlm-r-multilingual-v1)+ PLDA	0.39319976447977234	0.46401384083044983
Sentence Transformer (paraphrase-xlm-r-multilingual-v1)+ FAISS	0.3329157145830798	0.4598615916955017
Baseline modified with T5-base encoder + BiLSTM with 1024 hidden layer	0.6103711640920119	0.8127822160472387
T5-base encoder + BiLSTM with three hidden layer	0.5796309101679527	0.7964571031608197
Baseline modified with LegalBert	0.587383359920286	0.8089614449461618
Baseline modified with LegalBert Large 1.7m	0.5777515872781539	0.7919416464050018
Our Approach (Best) - Baseline modified with T5-large tokenizer and encoder + BiLSTM	0.6161397982433457	0.8166029871483154

- Traditional methods like Naive Bayes and embedding based approaches like Glove, Doc2Vec, Wiki-bm25 performed decently.
 - Next trials with 1D-Conv and LSTM captured long term dependencies leading to better performance
 - By focusing on improving the baseline, T5-large encoder led to best performance.
 - Legal-BERT/Legal-BERT-1.7m did not perform better than T5-encoder. This might be attributed to the fact that LEGAL-BERT (trained on EU legal documents, which also has European competition law) is not trained on Indian IT law documents. Using T5-encoder for embeddings with BiLSTM-CRF provides better results.
-

We considered the task of classifying sentences in Legal Documents based on their respective rhetorical roles. We tried all kinds of approaches ranging from Naive Bayes, Information Retrieval, LSTMs Large Language Model Embedding. Baseline modified with large language model embedding (T5) outperformed the baseline results.

From each approach:

- TF-IDF is based on the bag-of-words model, therefore it does not capture position in text, semantics, co-occurrences in different documents. It does not capture word semantics in the text hence only useful as a lexical level feature.
 - GLOVE embedding does not contain contextual information which could hurt performance.
 - 1D Conv and FCN does not capture long range dependencies.
 - In our best approach the labels with per-label-precision and per-label-F1 score are Argument by Petitioner and Argument by Respondents are low because they are very similar to each other.
 - Also, the model was not able to identify Precedent Not Relied class and classified them as Precedent Relied and Analysis. This can be attributed to incorrect identification of boundary because they are almost similar to one another.
-



TEXAS A&M UNIVERSITY
Engineering

Q & A



TEXAS A&M UNIVERSITY
Engineering

Thank You



TEXAS A&M UNIVERSITY
Engineering

Appendix

Model Parameters



TEXAS A&M UNIVERSITY
Engineering

```
"model": "t5-large",  
"bert_trainable": false,  
"dropout": 0.5,  
"word_lstm_hs": 758,  
"att_pooling_dim_ctx": 200,  
"att_pooling_num_ctx": 15,  
"lr": 3e-05, "batch_size": 32,  
"max_seq_length": 128,  
"max_epochs": 20,  
"early_stopping": 5
```

```
"weighted-f1":  
0.8007992216999676,  
"weighted-precision":  
0.806932444734725,
```

Results Analysis



TEXAS A&M UNIVERSITY
Engineering

"labels": ["DEFAULT", "mask", "NONE", "PREAMBLE", "FAC", "ISSUE",
"ARG_RESPONDENT", "ARG_PETITIONER", "ANALYSIS", "PRE_RELIED",
"PRE_NOT_RELIED", "STA", "RLC", "RPC", "RATIO"],

"per-label-f1": [0.0, 0.0, 0.9322493224932249, 0.9970559371933267,
0.839546191247974, 0.7727272727272727, 0.4150943396226416,
0.39655172413793105, 0.8250355618776671, 0.4651162790697673, 0.0,
0.603174603174603, 0.34615384615384615, 0.8877005347593583,
0.5294117647058822],

"per-label-precision": [0.0, 0.0, 0.9608938547486033, 0.9941291585127201,
0.7920489296636085, 0.8947368421052632, 0.7333333333333333, 0.5,
0.7733333333333333, 0.684931506849315, 0.0, 0.5428571428571428, 0.675,
0.8645833333333334, 0.5454545454545454],

Confusion Matrix for Best Model



TEXAS A&M UNIVERSITY
Engineering

