

Natural Language Inference

RoBERTa & BiLSTM Approach

Abstract

- We tackled the Natural Language Inference task using two models: a BiLSTM and a fine-tuned transformer.
- Both were trained to determine if a hypothesis logically follows a given premise.
- The BiLSTM achieved over 70% F1, while the transformer exceeded 90%, highlighting the effectiveness of deep pretrained models.

Introduction

WHAT

- NLI is a task that determines whether a hypothesis logically follows from a given premise.
- It involves classifying the relationship as entailment, neutral, or contradiction (binary in our case).

WHY

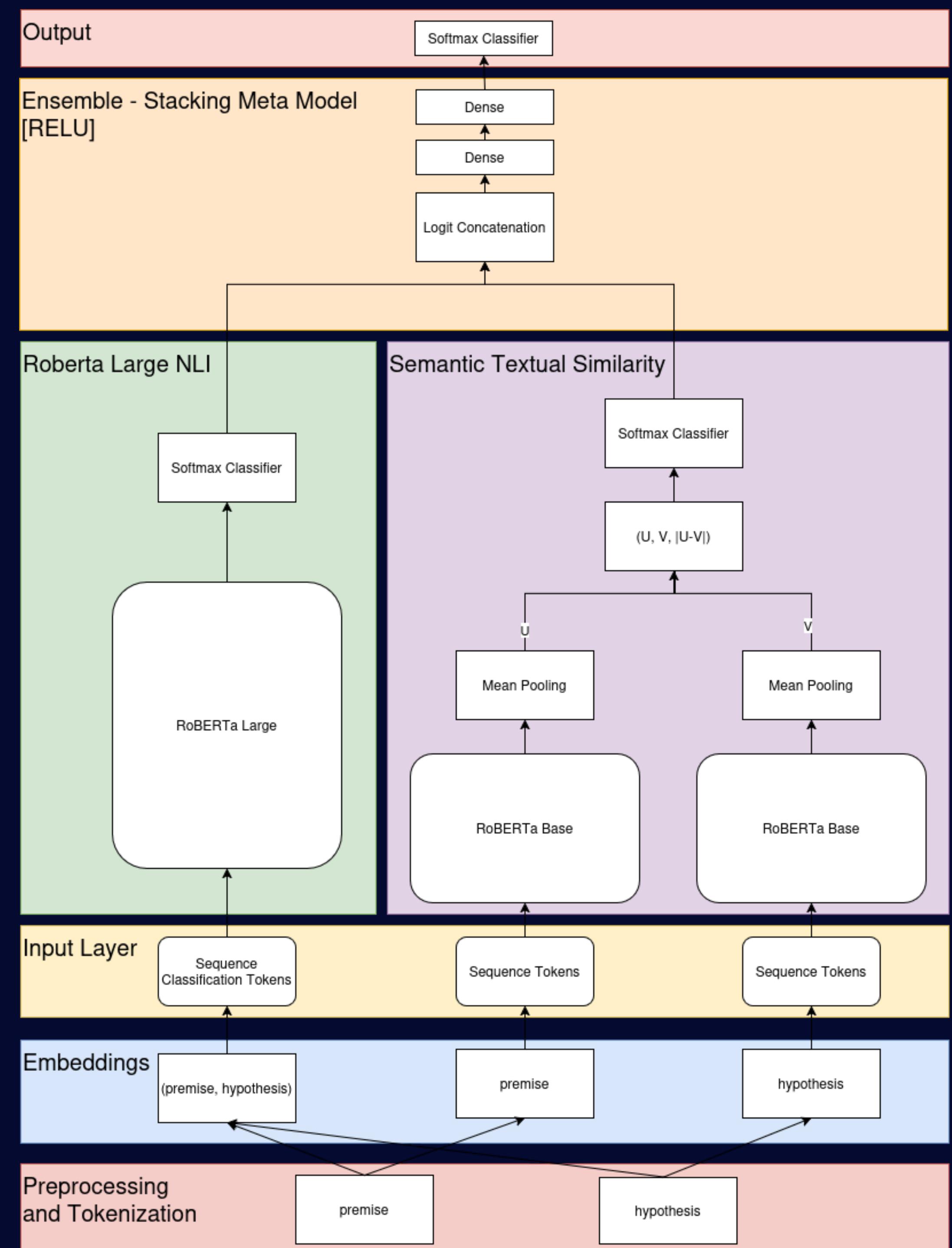
- Comparing a simpler BiLSTM and a deep pretrained transformer helps evaluate trade-offs between efficiency and accuracy.
- This dual approach shows how different architectures understand logical relationships in text.

HOW

- Non-Transformer Captures word order and context using sequential processing
- Will struggle with long dependencies
- Deep Learning Uses self-attention to understand full sentence context in parallel
- Handling long-range dependencies better.

References

DEEP Learning Transformer



INPUT DATA

- Minimal preprocessing: maintain contextual integrity, with stop words and punctuation deliberately retained.

SEMANTIC TEXTUAL SIMILARITY

- Compare the similarity of semantically meaningful sentence embeddings [2].

RoBERTa LARGE MODEL

- RoBERTa-large model fine-tuned in a pairwise classification setup [3].

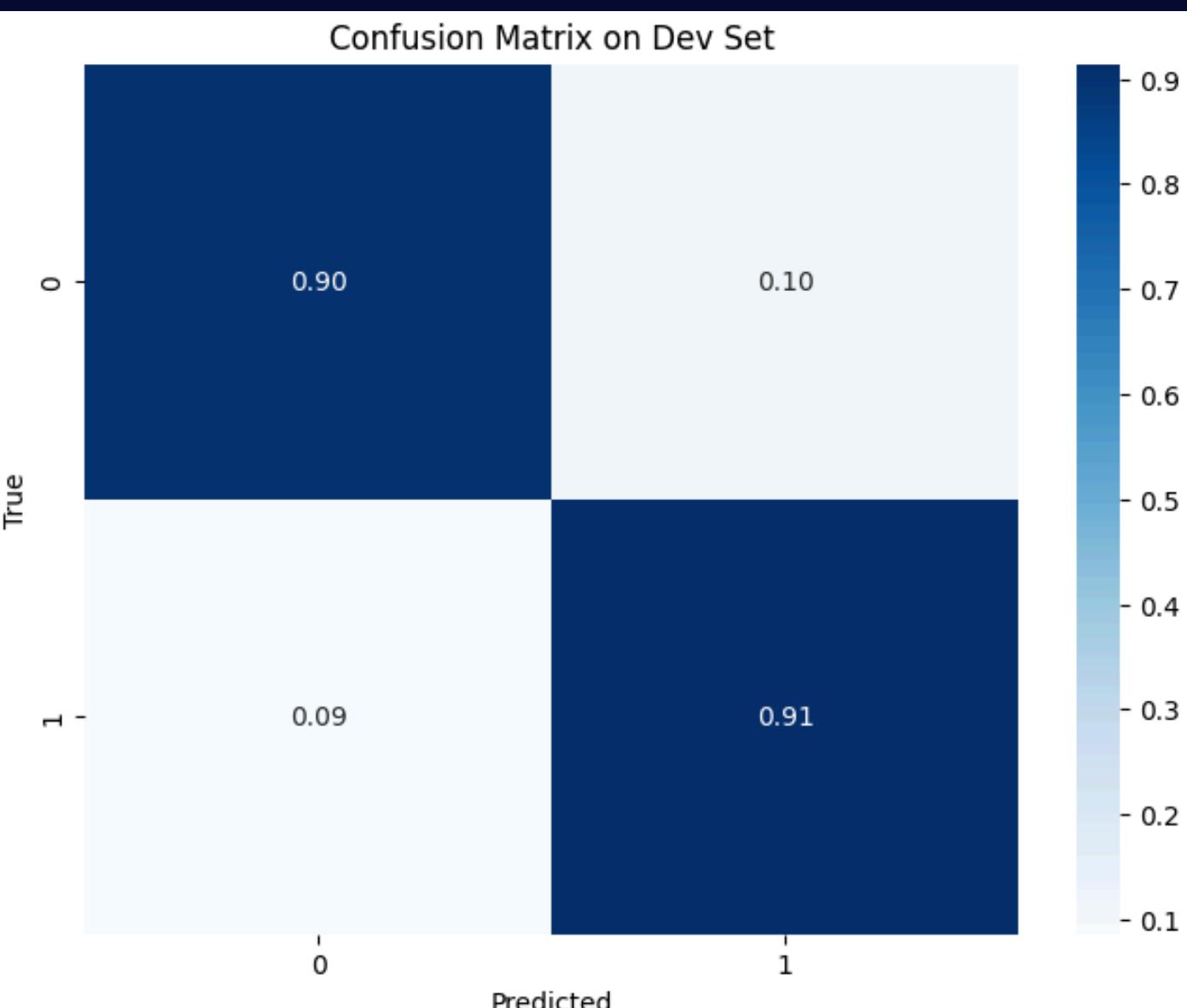
ENSEMBLE META MODEL

- Leverages strengths of different architectures.
- Meta-classifier learns optimal weighting of individual model predictions.
- Improves robustness and generalization.

Results

TRANSFORMER MODEL PERFORMANCE

Ensemble Model Component Test: Accuracy Comparison			
Model	Train	Dev	Test
RoBERTa Large	99.41%	90.88%	91.09%
Similarity Model	97.07%	78.29%	78.03%
Ensemble Model	99.61%	90.87%	91.02%

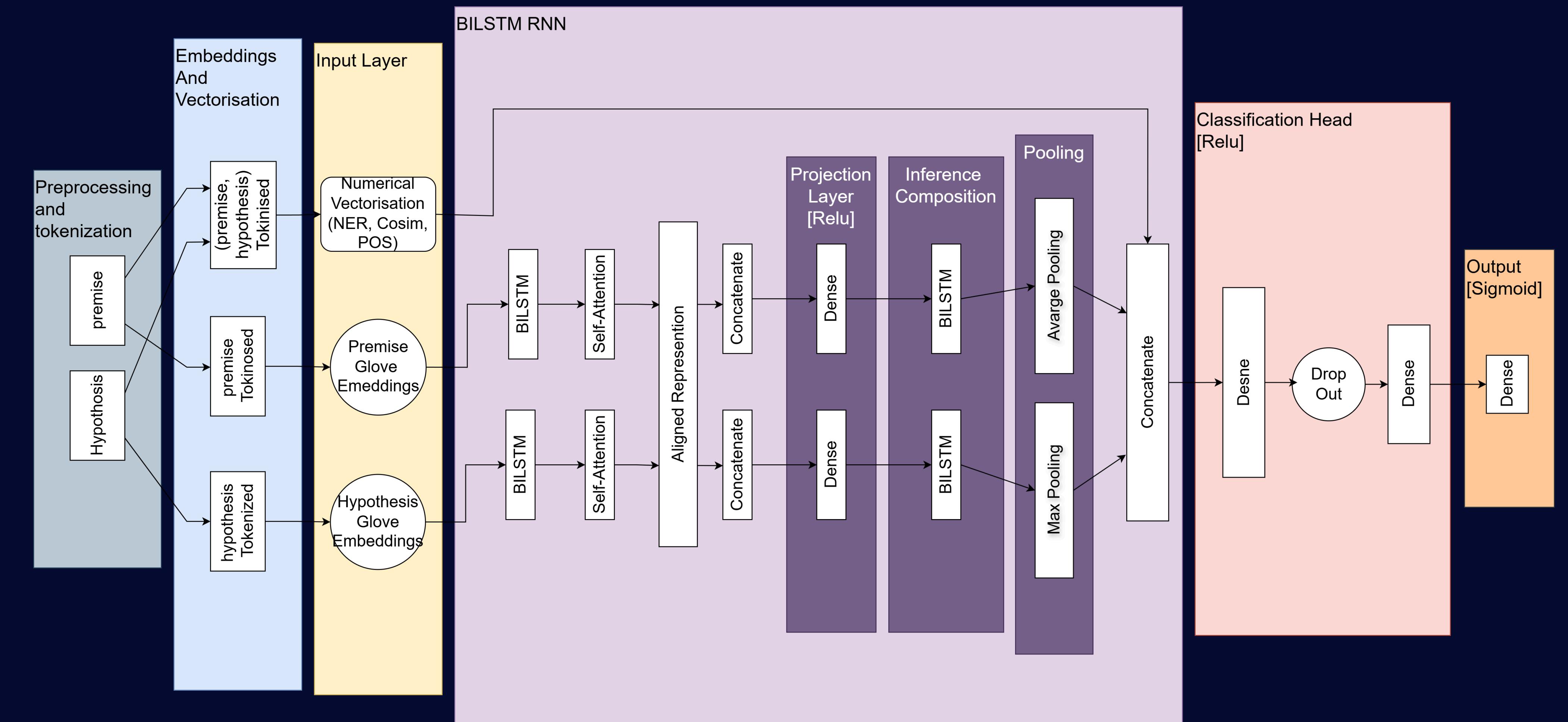


BiLSTM RNN MODEL PERFORMANCE

BiLSTM RNN Performance on Dev Set	
MCC	39.9%
Weighted Precision	69.9%
Weighted Recall	69.9%
Weighted F1-score	69.9%
Accuracy	69.9%

Ensemble Performance on Dev Set	
MCC	81.8%
Weighted Precision	90.9%
Weighted Recall	90.9%
Weighted F1-score	90.9%
Accuracy	90.9%

Deep Learning Non Transformer



INPUT

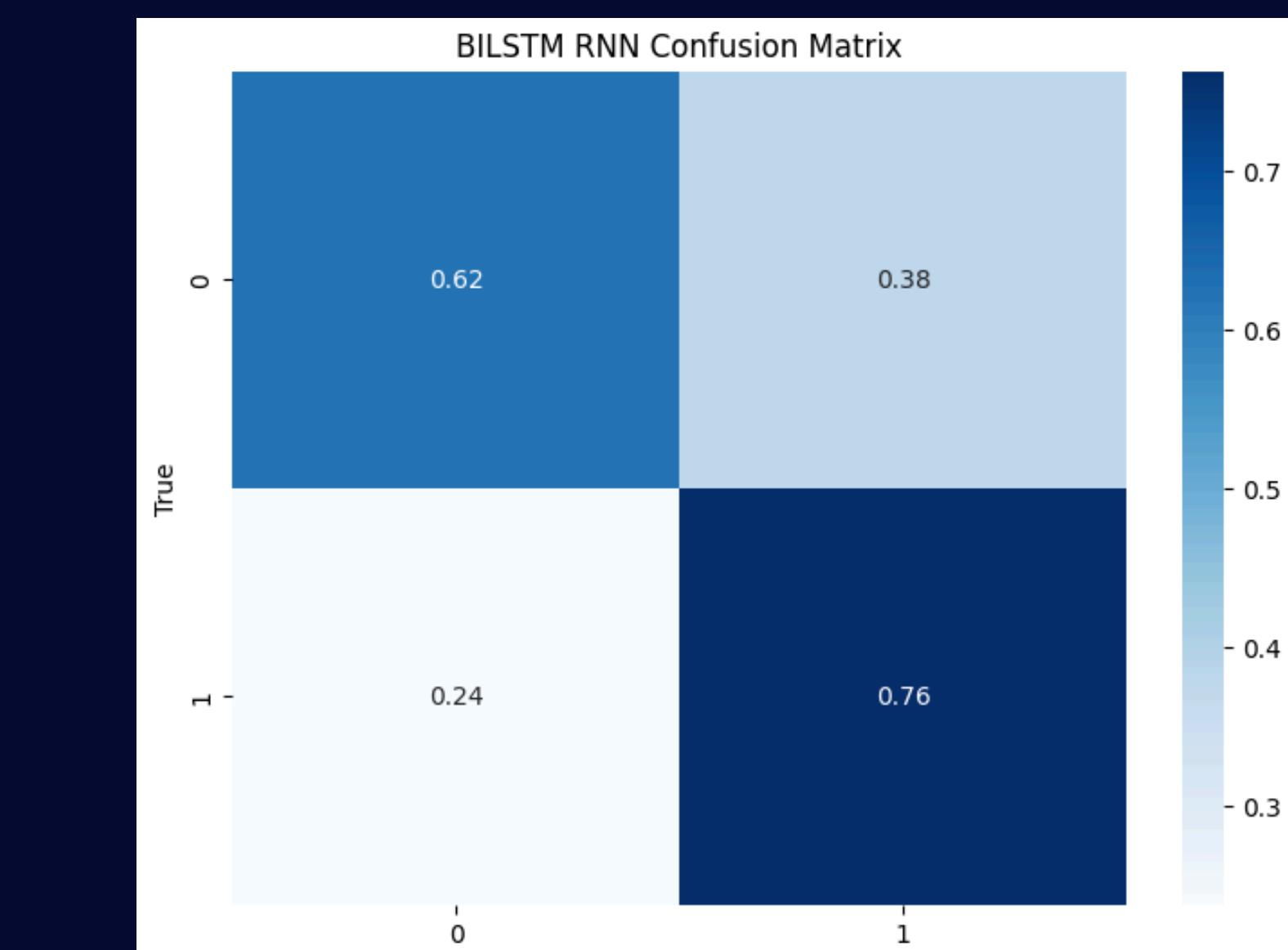
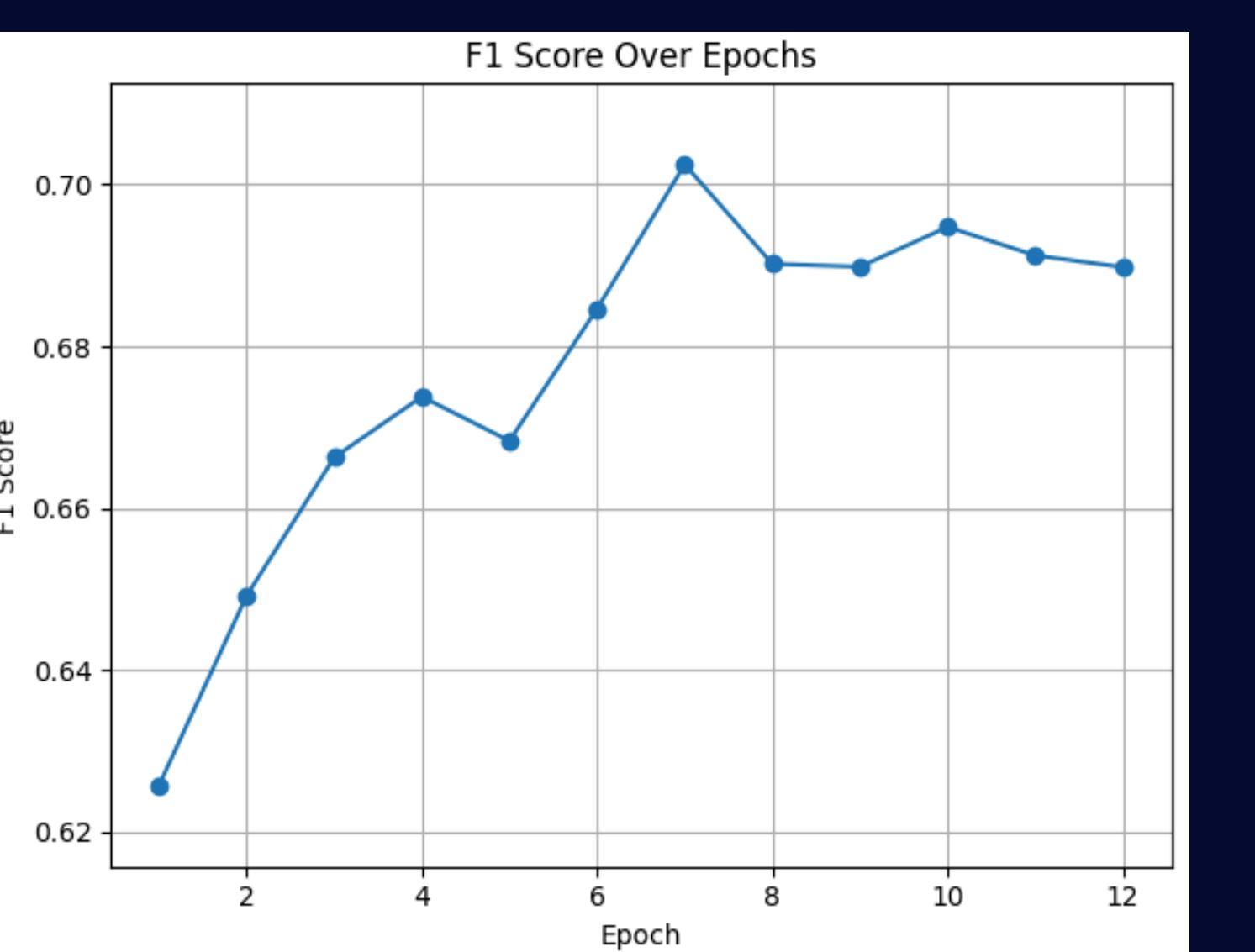
- Text pairs (premise + hypothesis) cleaned and tokenized.
- Tokens are converted into Glove Embeddings (NER) and POS tagging are applied.

BiLSTM RNN ESIM

- Separate BiLSTMs encode each sentence.
- Intermediate layers refine embeddings through projection and composition [1].
- Inference is modeled by comparing encoded sentence representations [1].
- Max pooling aggregates important features across time steps [1].

CLASSIFIER

- Final vector passed through dense layers with dropout.
- Output layer uses sigmoid to predict class label (entailment or not).



Conclusions

- BiLSTM for NLI is lightweight and easy to train
- Transformer model offers higher robustness and better overall performance.

FUTURE WORK

BOTH

- Data set augmentation.

BiLSTM

- More hyper Parameter tuning .
- Fine tuning encoder as well as model.

Transformer

- Explore additional models to ensemble, to improve robustness.