# Text Classification using CNN, RNN and HAN
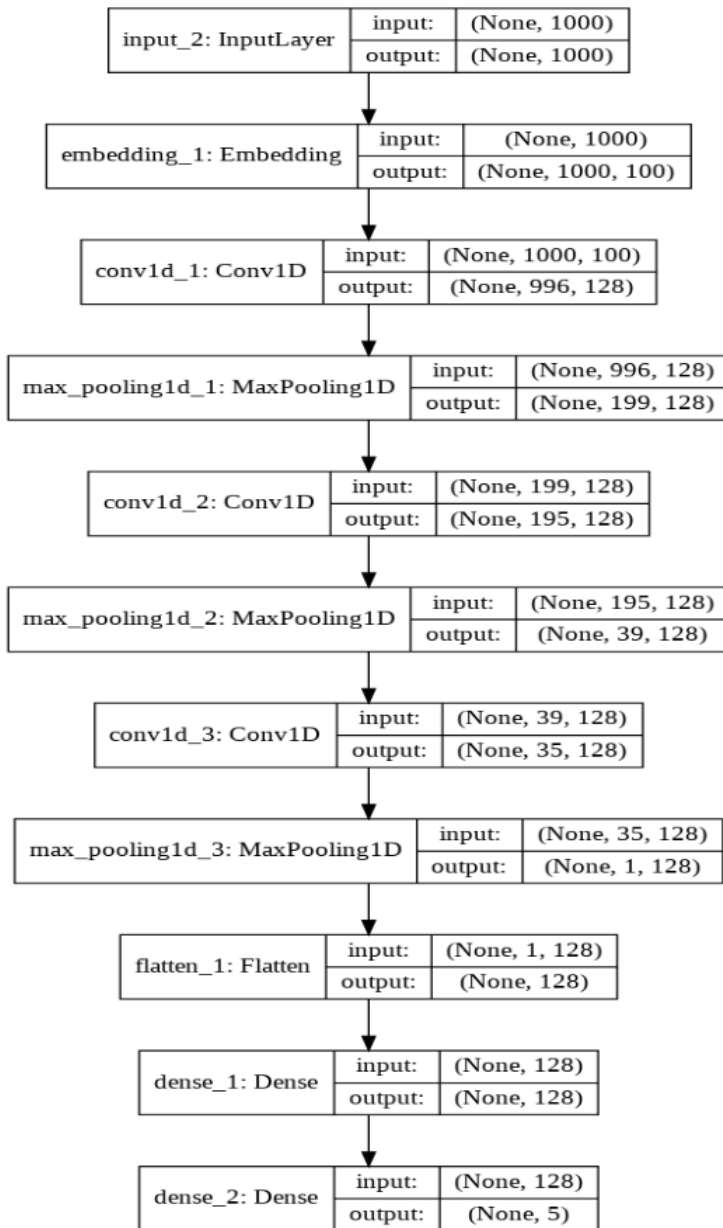## Authors: Tejas Gupta, Tushar Dahibhate, Shraddha Dhyade

**Text classification problem:**
Using deep learning, we aim to perform a sentiment analysis of the yelp online reviews. The reviews are in the form of star ratings(1-5) and textual format. Thus, we use the pretrained GloVe vectors dataset and perform a sentiment analysis to later perform comparison of the aforementioned models.
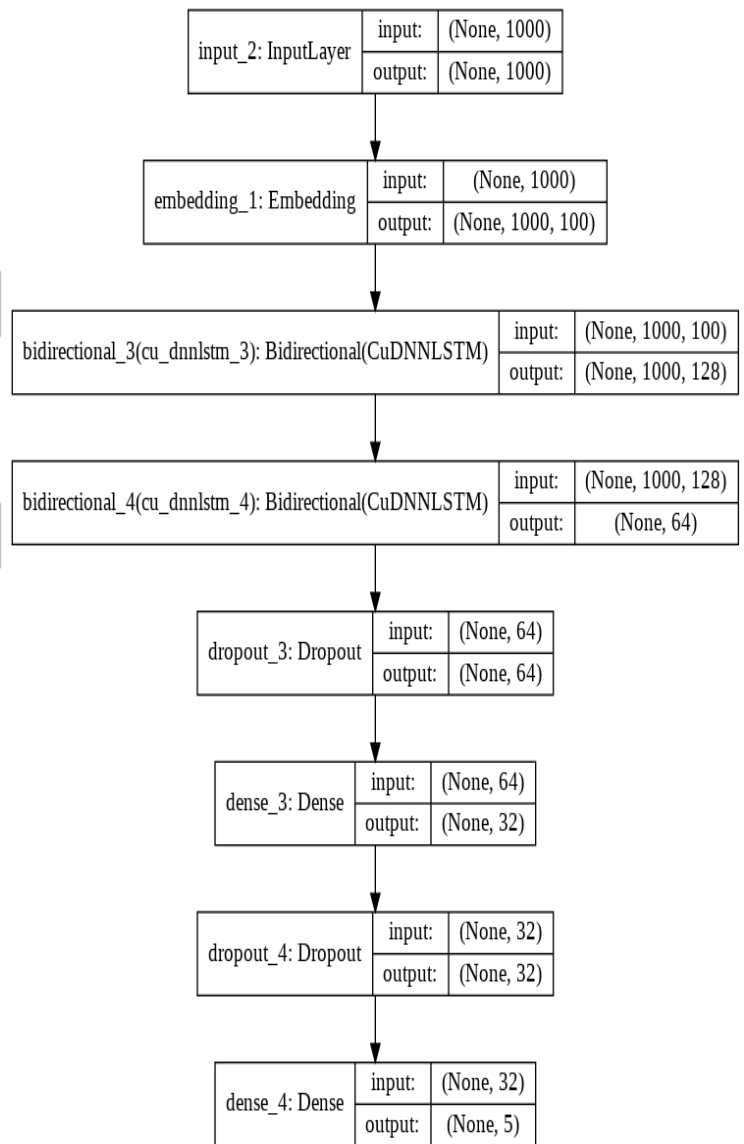
**Description of the data sets:**
The yelp dataset consists of reviews with star ratings 1-5 and textual reviews and other information such as user id, review id usefulness of the review, business id, date, review id, user id and comments on the review, given by other users. However, the ratings and reviews are of prime importance to us in this analysis. The dataset has 10000 samples which were divided into train test and validation datasets as follows: The training set contains 6000 samples, validation data has 2000 samples and test data has 2000 samples.
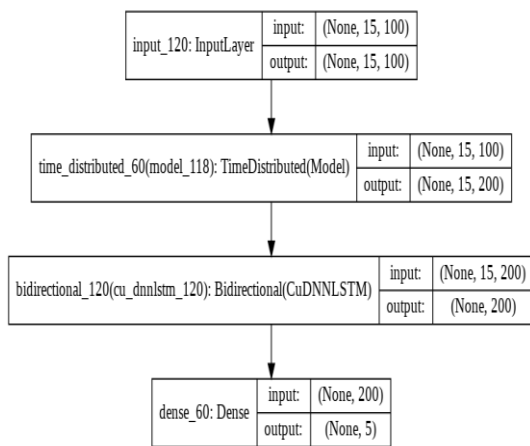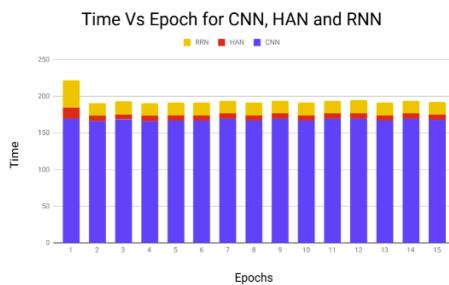
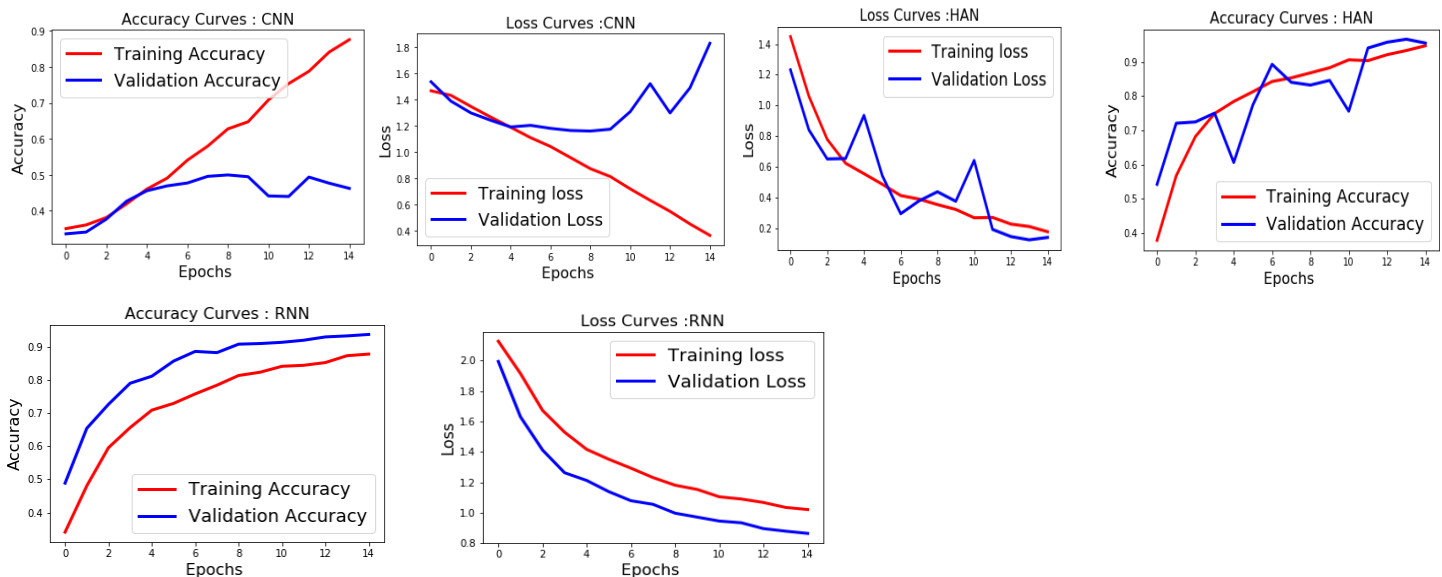**Architectures: Visual Graphs**



CNN

RNN

| input_120: InputLayer | input: | (None, 15, 100) |
|---|---|---|
| | output: | (None, 15, 100) |

| time_distributed_60(model_118): TimeDistributed(Model) | input: | (None, 15, 100) |
|---|---|---|
| | output: | (None, 15, 200) |

| bidirectional_120(cu_dnnlstm_120): Bidirectional(CuDNNLSTM) | input: | (None, 15, 200) |
|---|---|---|
| | output: | (None, 200) |

| dense_60: Dense | input: | (None, 200) |
|---|---|---|
| | output: | (None, 5) |

**HAN**

| Architectures: Hyperparameter Table | | | |
|---|---|---|---|
| | CNN | RNN | HAN |
| Epochs | 15 | 15 | 15 |
| Batch Size | 128 | 64 | 256 |
| Learning rate | 0.001 | 0.0001 | 0.001 |
| L1 Regularization | default | 0.01 | 0.00001 |

**Time/Epochs graph:**



**Training and Validation Accuracy and Loss over Epochs:**



## Performance Analysis:

Test accuracy of 0.9355 was achieved in HAN model, while RNN gave 0.8529 accuracy and CNN gave 0.7319 accuracy. Thus, we can say that HAN performed better among the three.

## Hyperparameter Tuning: Choices, rationale, observed impact on the model performance:

We tuned the model by tweaking following hyperparameters-

- **Learning rate-** When the model was overfitting, decreasing the learning rate would increase the performance while on underfitting, the increasing the learning rate would increase the model performance.
- **Regularization parameter–** decreasing the L1 regularization parameters improved the model performance
- **Batch size-** Increasing the batch size improved the model performance.