


```
"4bbb0629e66146edaf4ac7bde47062fb": {  
  "text": [  
    "The",  
    "factum",  
    "of",  
    "accident",  
    ",",  
    "allegation",  
    "of",  
    "rash",  
    "and",  
    "negligent",  
    "driving",  
    "causing",  
    "death",  
    "of",  
    "Sukendra",  
    "Pal",  
    "Singh",  
    "were",  
    "denied",  
    "."  
  ],  
  "labels": [  
    "O",  
    "O",  
    "O",  
    "O",  
    "O",  
    "O",  
    "O",  
    "O",  
    "O",  
    "O",  
    "O",  
    "O",  
    "O",  
    "O",  
    "O",  
    "B_OTHER_PERSON",  
    "I_OTHER_PERSON",  
    "I_OTHER_PERSON",  
    "O",  
    "O",  
    "O"  
  ]  
},
```

Dataset 2 snippets:

[illegible]

```

    "3": {
      "text": [
        "Great",
        "laptop",
        "that",
        "offers",
        "many",
        "great",
        "features",
        "!"
      ],
      "labels": [
        "0",
        "0",
        "0",
        "0",
        "0",
        "0",
        "0",
        "B",
        "0"
      ]
    },

```

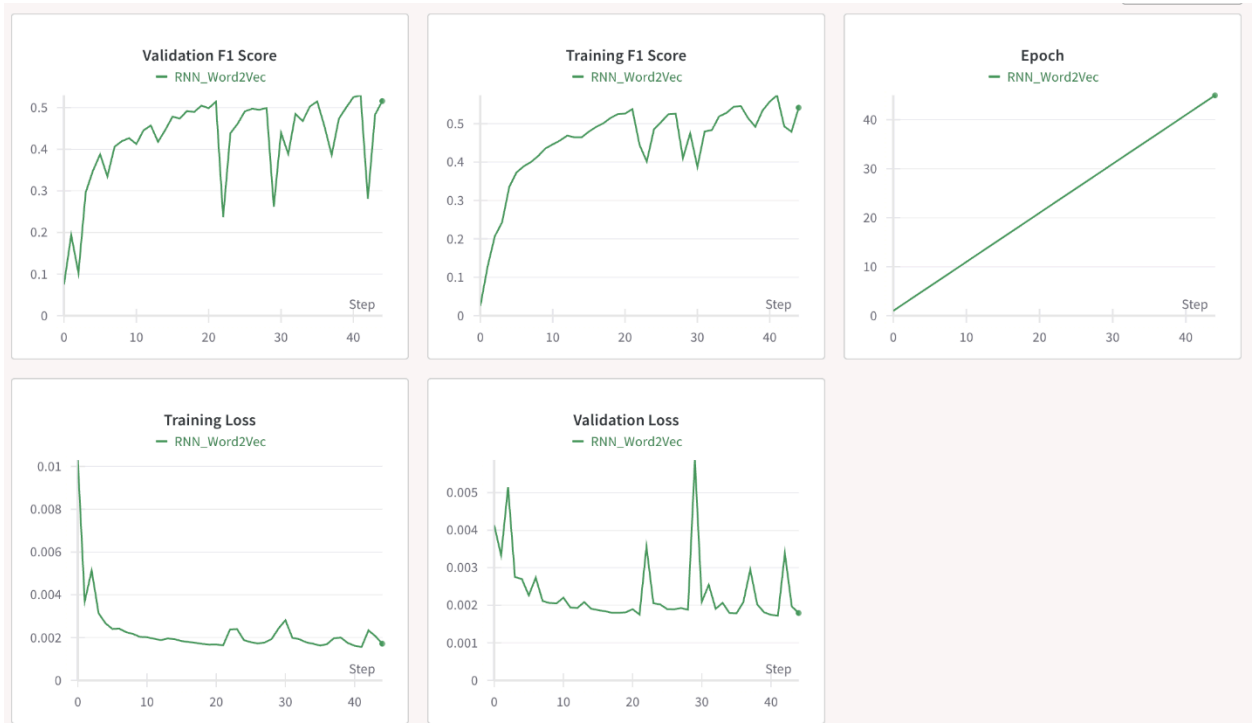
Additional Preprocessing:

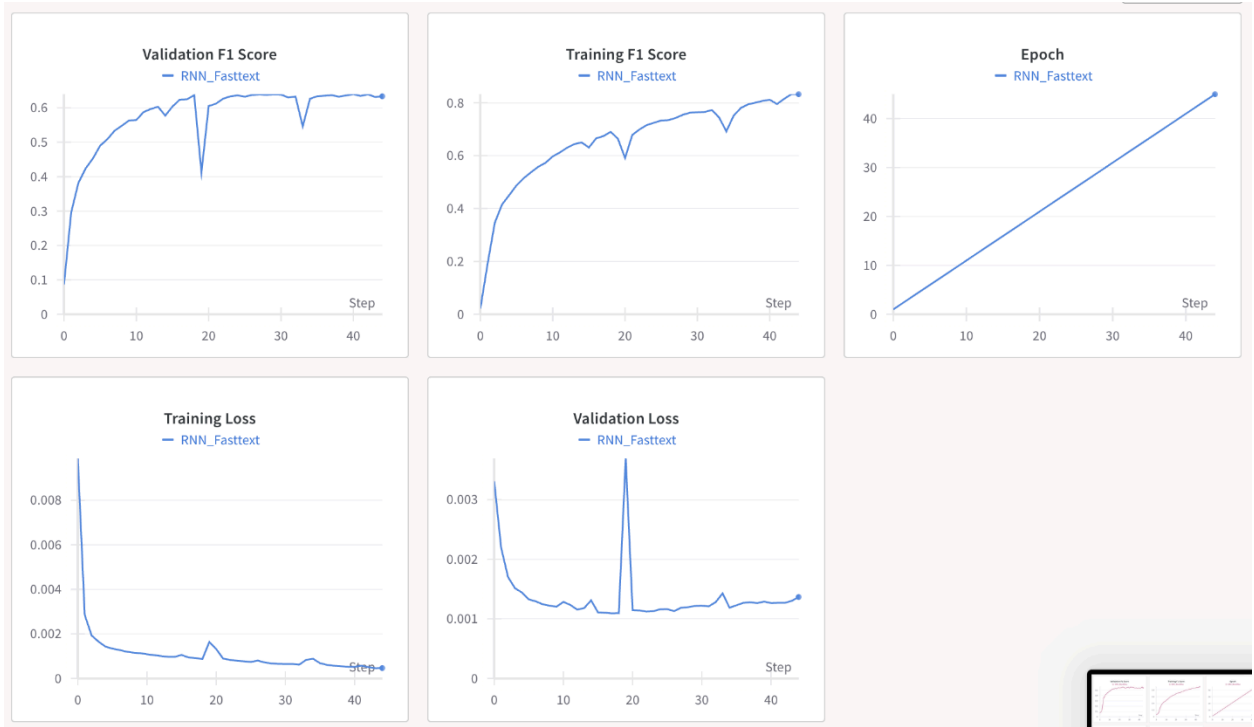
Special emphasis was given on separating tokens that were followed by sentence competition markers like full stop, question mark and other punctuation marks. Rules mentioned in the original json file were applied to generate and label the given texts.

Plots:

Part 2 dataset 1:

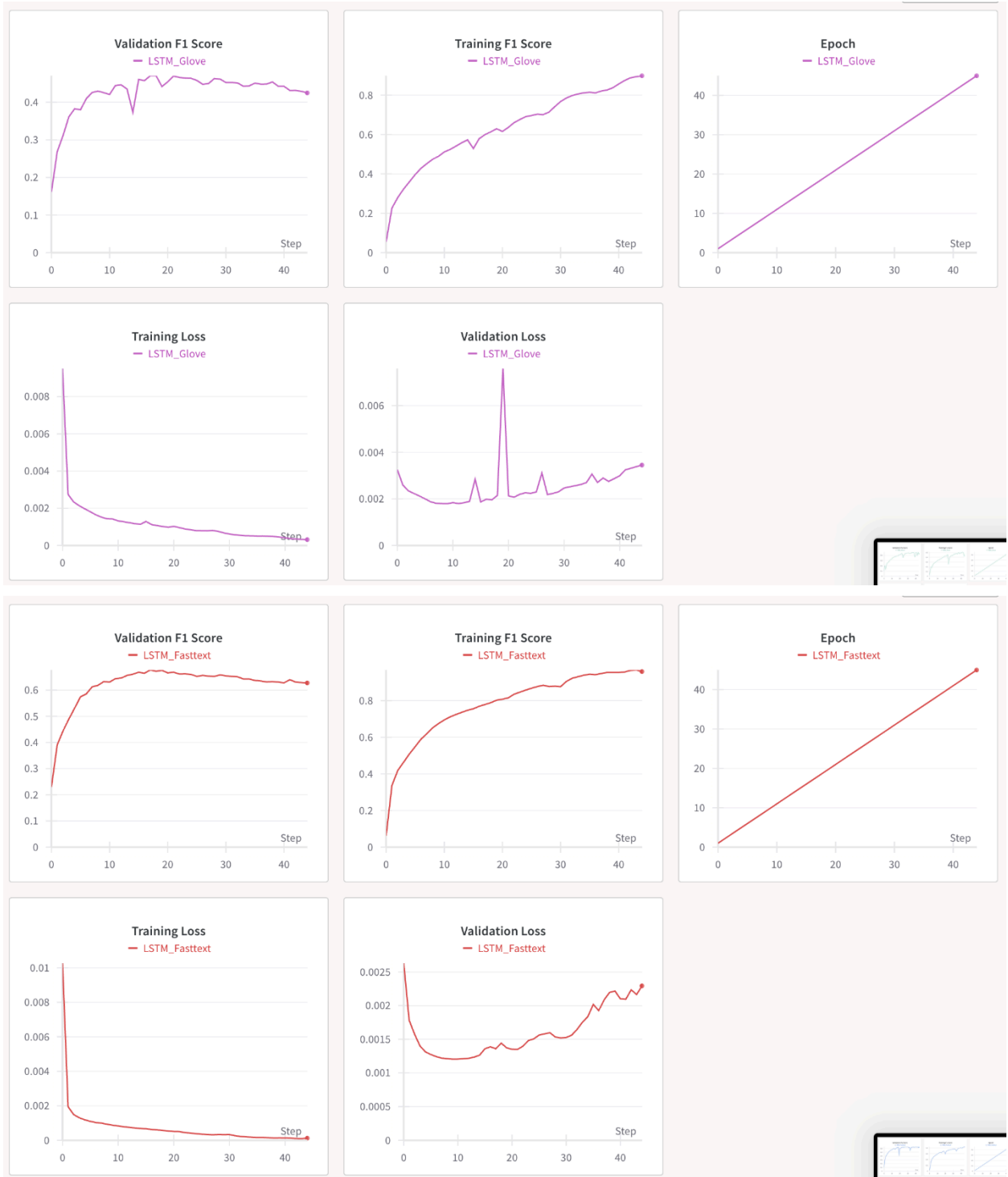
a. RNN





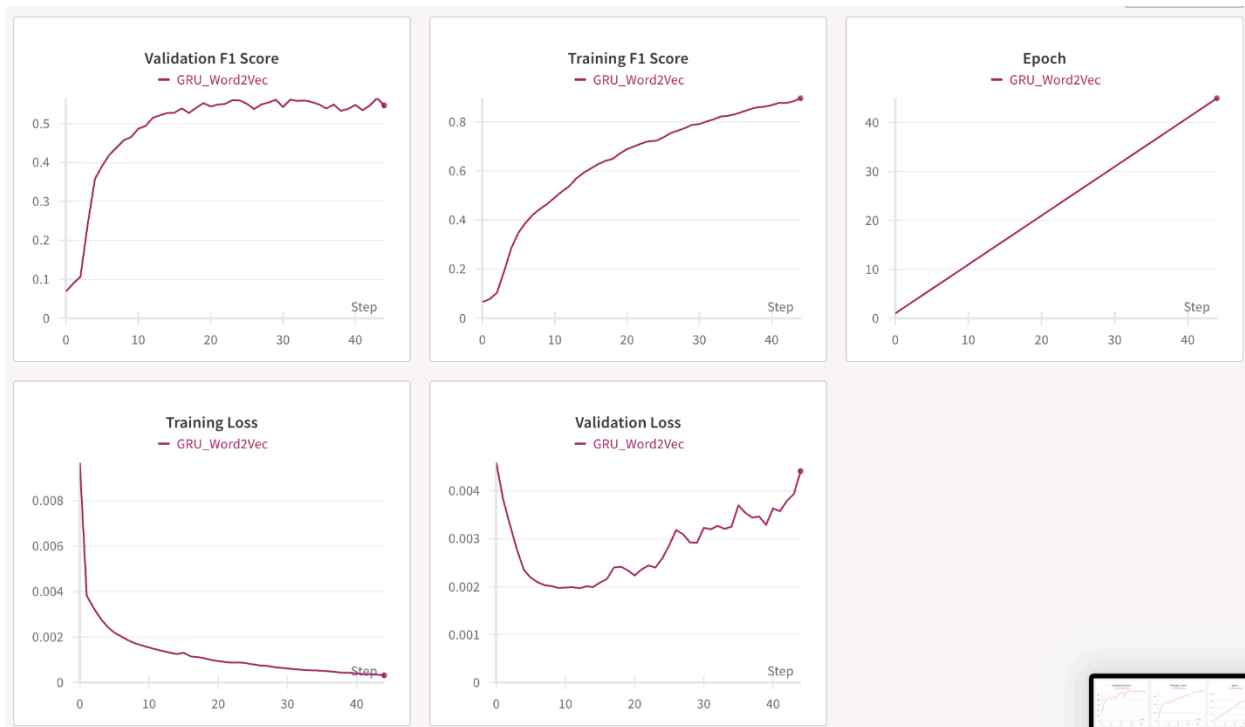
b. LSTM



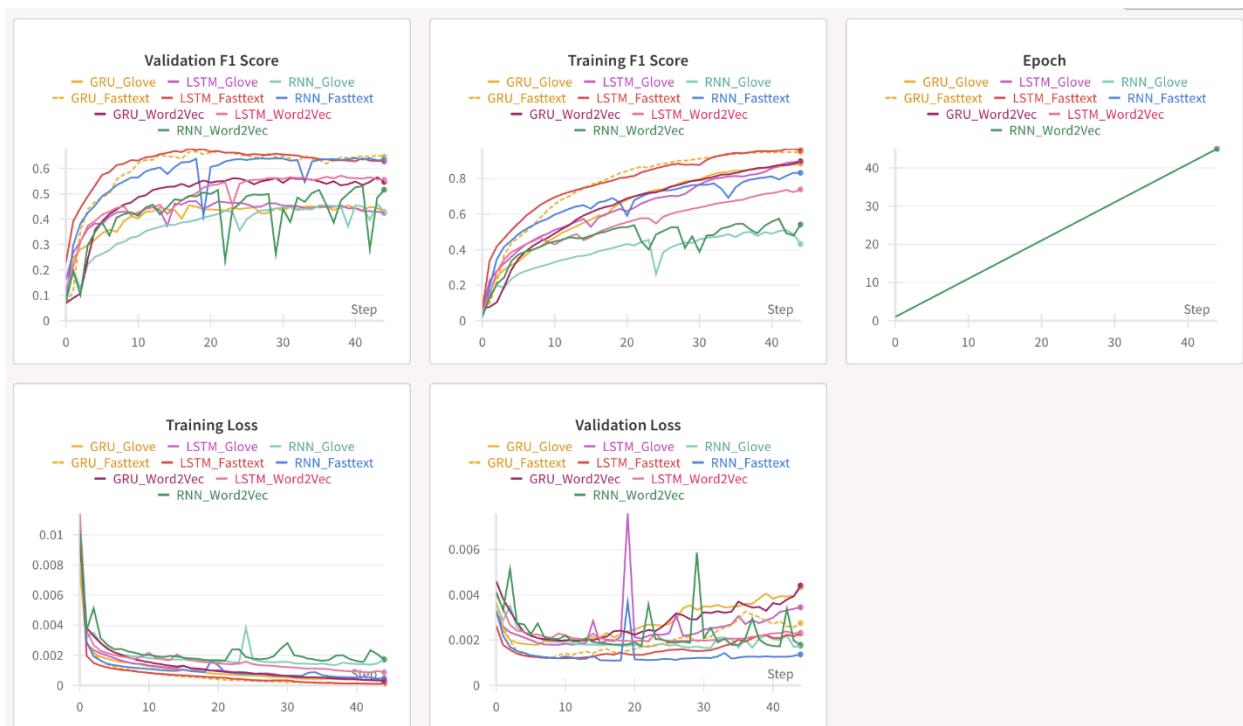


c. GRU



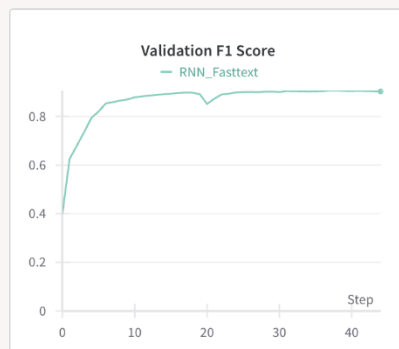
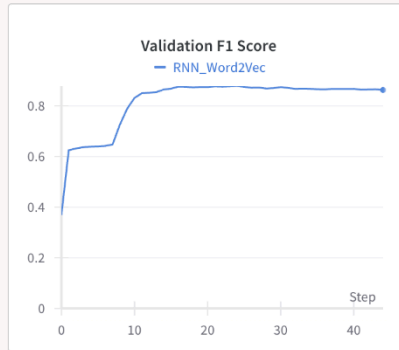
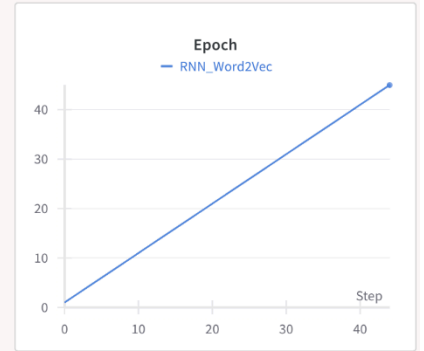
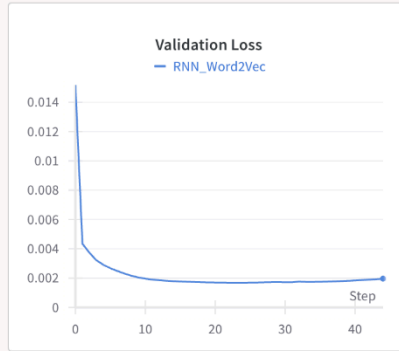


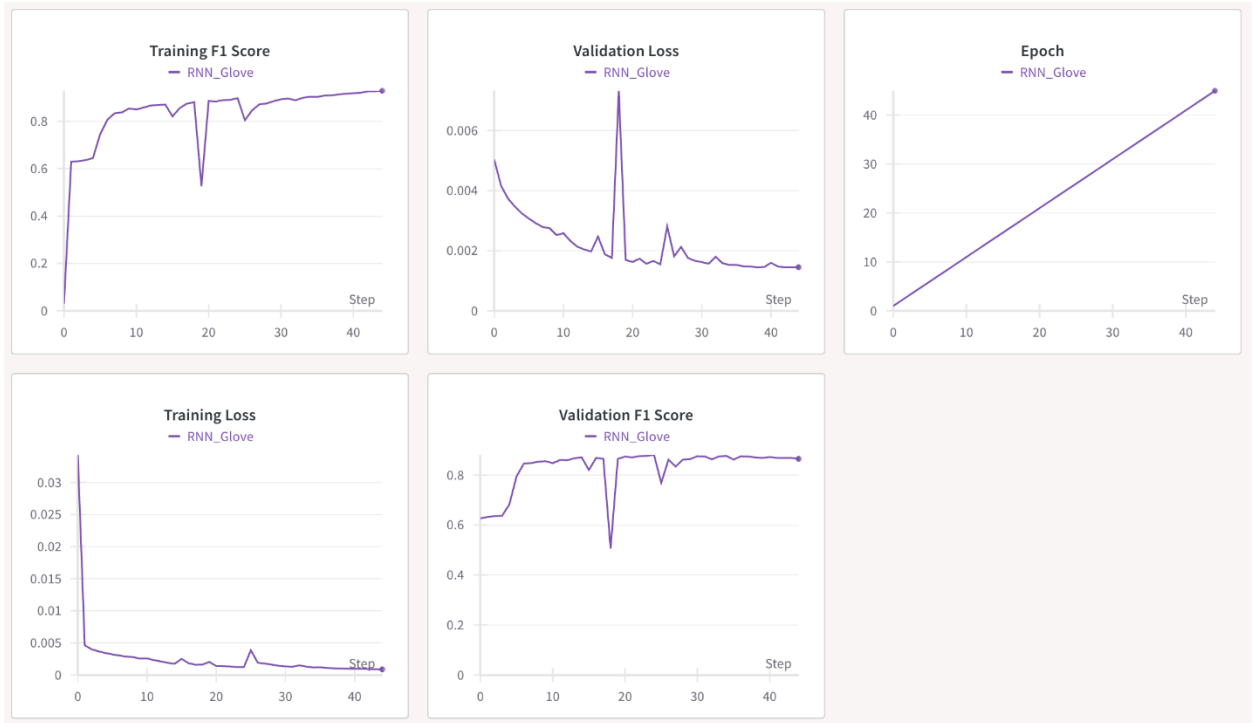
d. Combined



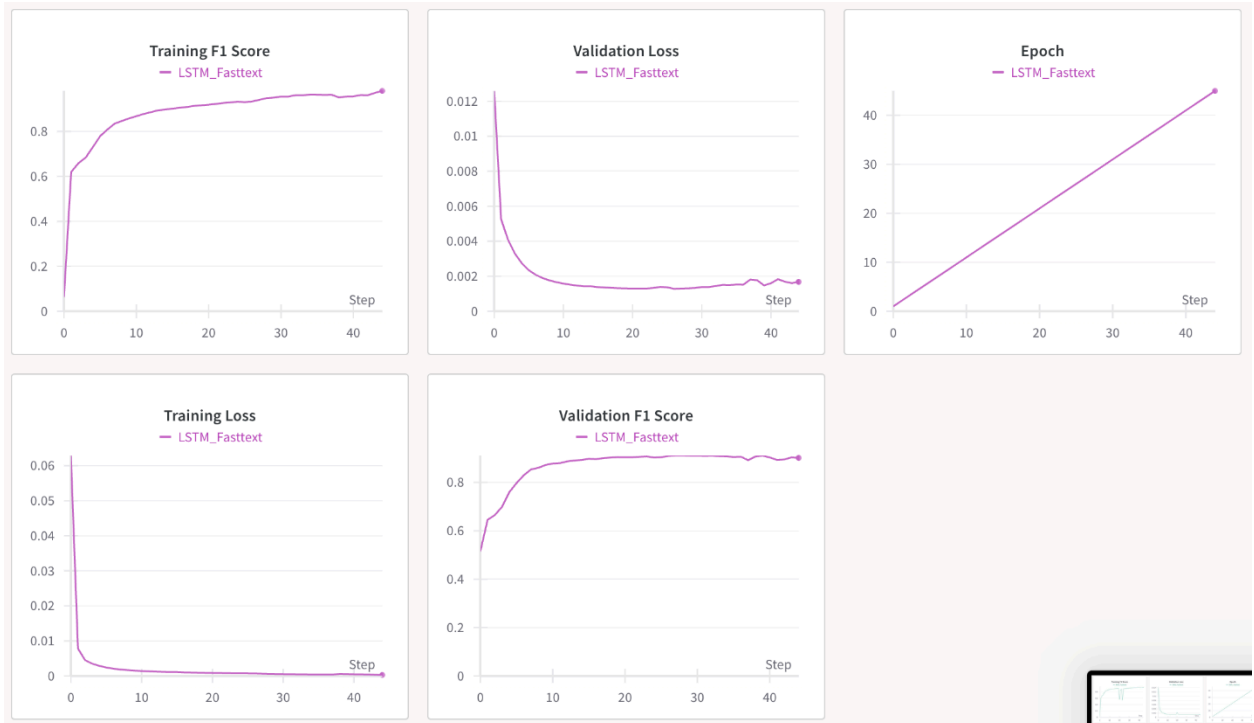
Part 2 dataset 2:

a. RNN



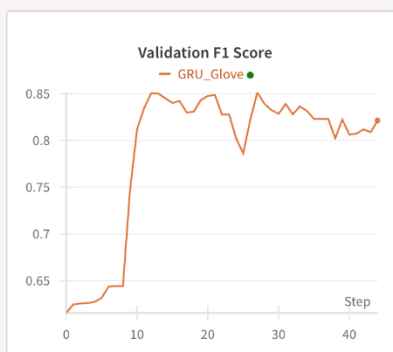
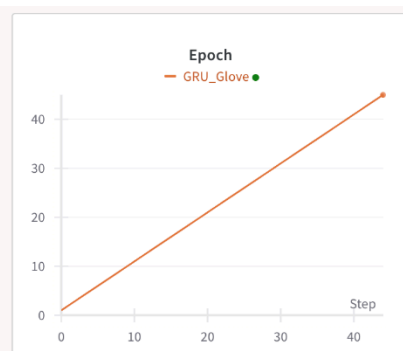
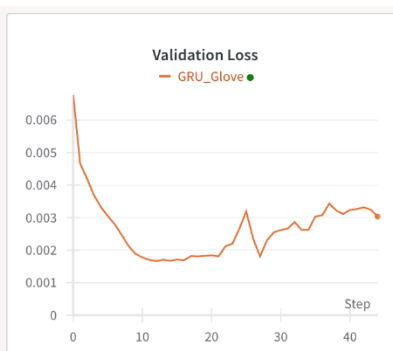
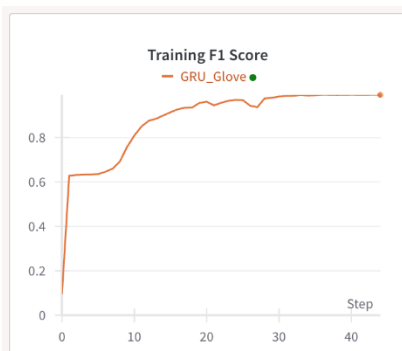
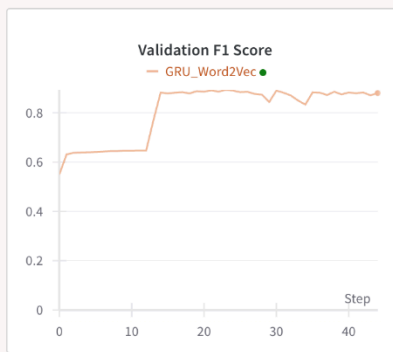
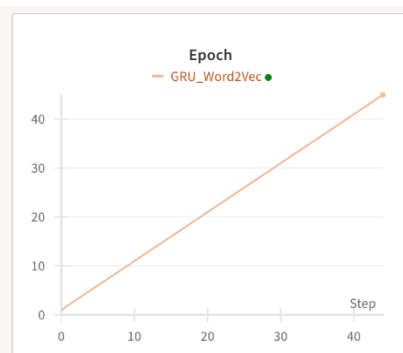
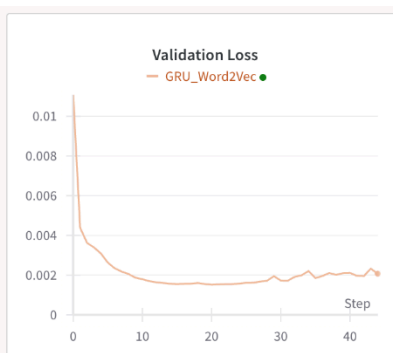


b. LSTM



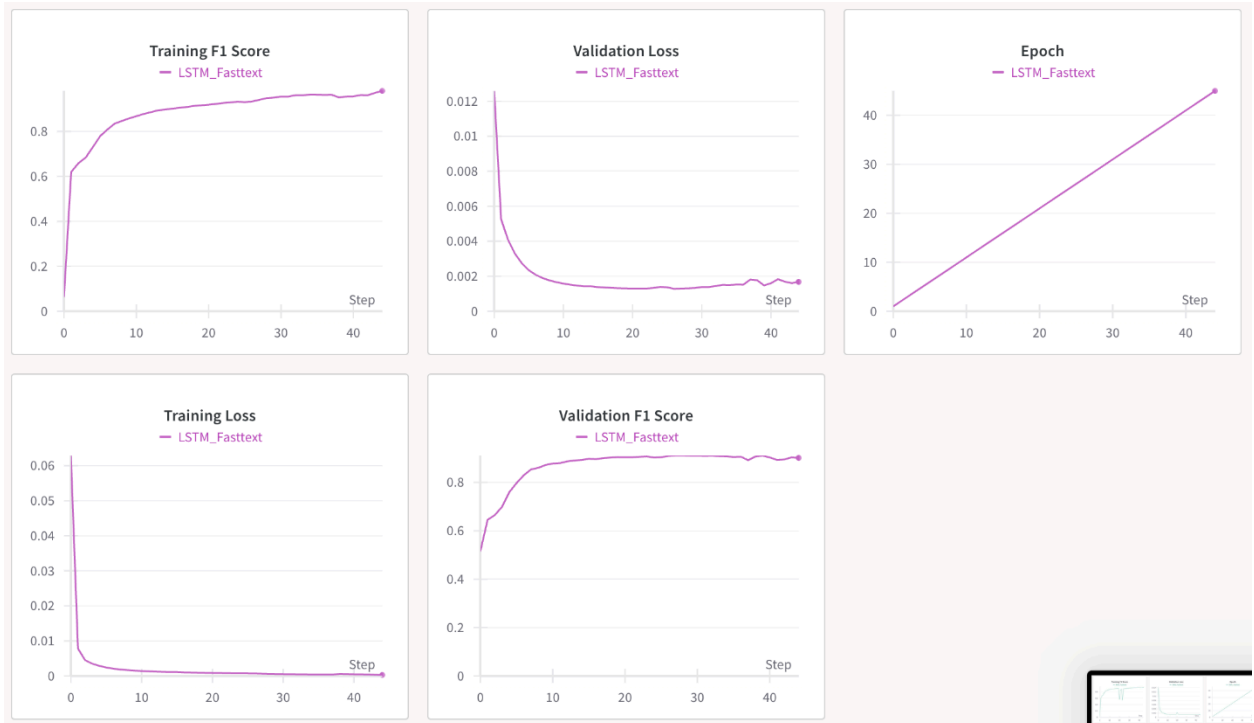


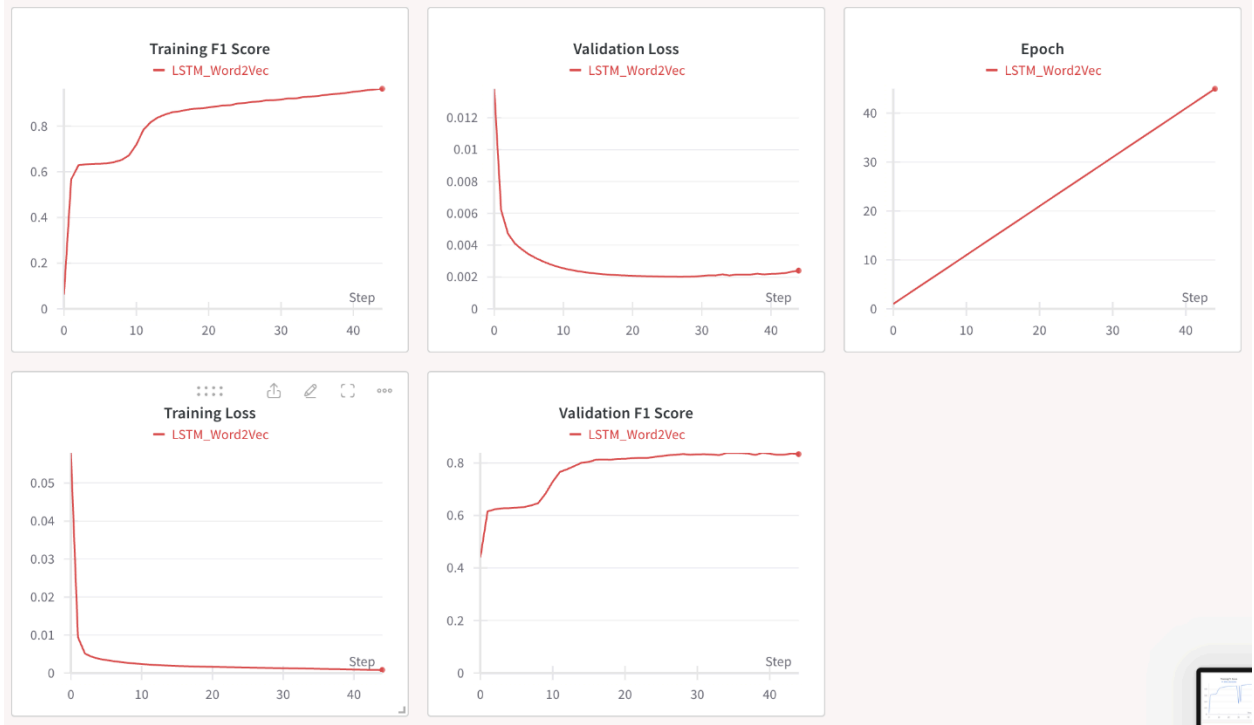
c. GRU





d. BiLSTM





d. Combined



General Analysis and Comparison of word embeddings

1. Word2Vec:

a. Pros:

- i. Effective at capturing semantic relationships between words.
- ii. Can handle out-of-vocabulary words to some extent.

b. Cons:

- i. May not capture subword information.
- ii. Fixed-size embeddings.

2. GloVe (Global Vectors for Word Representation):

a. Pros:

- i. Utilizes global statistics, capturing word co-occurrence information effectively.
- ii. Generally performs well on semantic tasks.

b. Cons:

- i. May struggle with rare or out-of-vocabulary words.
- ii. Fixed-size embeddings.

3. FastText:

a. Pros:

- i. Considers subword information, effective for handling rare words.
- ii. Can handle out-of-vocabulary words by breaking them down into subword components.

b. Cons:

- i. May generate larger embeddings.
- ii. May not capture long-range semantic relationships as well.

Analysis of the plots:

1. We observe that the f1 scores of all the embeddings follow the order: fasttext > word2vec > glove in the dataset 1.
2. We observe that the f1 scores of all the embeddings follow the order: fasttext > glove > word2vec in the dataset 2.
3. The better performance of fasttext in both datasets can be attributed to the fact that it works very well in smaller context and does not consider longer dependencies and

semantics. Since our problem was not focused on deriving or considering long term semantics, addition of the same would have hindered the performance.

4. We also observed that word2vec outperformed glove embedding in dataset 1 whereas opposite was true for dataset 2.
5. This can be attributed to the fact that in first problem of NER, many word labels were out of vocabulary which is a weak point of glove embedding, hence word2vec outperformed it.
6. In the second dataset of ATE, glove performed better than word2vec because it captures co-occurrence information quite well that can be quite useful in associating terms with a single entity like emotion or object.

Table comparing the performance of all the models:

Dataset 1

Model number	Embedding Used	Accuracy	Macro f1
RNN word2vec	word2vec	94.634	0.538
RNN glove	glove	93.678	0.441
RNN fasttext	fasttext	95.961	0.640
LSTM word2vec	word2vec	94.674	0.571
LSTM glove	glove	93.978	0.470
LSTM fasttext	fasttext	95.261	0.676
GRU word2vec	word2vec	94.634	0.539
GRU glove	glove	93.678	0.451
GRU fasttext	fasttext	94.961	0.647
BiLSTM word2vec	word2vec	95.122	0.653
BiLSTM glove	glove	94.672	0.537
BiLSTM fasttext	fasttext	96.012	0.701

Dataset 2

Model number	Embedding Used	Accuracy/100	Macro f1
RNN word2vec	word2vec	0.977	0.879
RNN glove	glove	0.983	0.744
RNN fasttext	fasttext	0.986	0.900
LSTM word2vec	word2vec	0.981	0.846
LSTM glove	glove	0.984	0.751
LSTM fasttext	fasttext	0.986	0.904
GRU word2vec	word2vec	0.982	0.866
GRU glove	glove	0.982	0.695
GRU fasttext	fasttext	0.987	0.908
BiLSTM word2vec	word2vec	0.985	0.854
BiLSTM glove	glove	0.987	0.821
BiLSTM fasttext	fasttext	0.988	0.924

Contribution of members:

Each member contributed equally to the project and most of the things were done in sync overlapping with each other.