

# Comparative Analysis of RNN Architectures for Sentiment Classification

Mounika Teppola

121012135

DATA641 - Natural Language Processing

Homework 3

## Introduction

Sentiment classification is a fundamental task in Natural Language Processing (NLP) that focuses on identifying the emotional polarity of a given piece of text, such as a movie review, tweet or customer feedback. The goal is to automatically determine whether the expressed sentiment is positive, negative or neutral. This capability is crucial for applications like opinion mining, product recommendation systems and social media analytics, where understanding user attitudes can drive data-driven decision-making.

Traditional machine learning approaches for sentiment classification, such as Support Vector Machines (SVMs) or Naive Bayes classifiers, rely heavily on handcrafted features like term frequencies or n-grams. However, these methods often fail to capture the sequential and contextual dependencies inherent in natural language. To address this limitation, deep learning architectures particularly Recurrent Neural Networks (RNNs), have emerged as powerful tools for sequence modeling.

RNN-based models, including Long Short-Term Memory (LSTM) and Bidirectional LSTM (BiLSTM), are specifically designed to capture temporal dependencies and contextual relationships between words in a sequence. While simple RNNs are capable of learning short-term dependencies, they often suffer from issues such as vanishing or exploding gradients, making them less effective for long text sequences. LSTMs and BiLSTMs mitigate these problems through gating mechanisms and bidirectional processing, thereby achieving more robust performance on sequential data.

This project conducts a comprehensive comparative analysis of multiple RNN architectures—RNN, LSTM and BiLSTM for the task of sentiment classification using the IMDb Movie Review Dataset. The objective is to evaluate how different architectures, activation functions, optimizers and sequence lengths affect model performance and training efficiency. Through this analysis, the project aims to identify the optimal configuration that balances accuracy, F1-score and computational cost under CPU constraints.

## Dataset Summary

The dataset used in this project is the **IMDb Movie Review dataset**, which contains **50,000 movie reviews** evenly split into **25,000 training** and **25,000 testing** samples. Each review is labeled as either **positive** or **negative**, making it suitable for binary sentiment classification.

## Data Preprocessing

The preprocessing pipeline converts raw text into sequences suitable for RNN models. The steps are:

1. **Lowercasing and Cleaning:** All reviews were converted to lowercase, and punctuation and special characters were removed to standardize the text and reduce vocabulary size.
2. **Tokenization and Sequencing:** Using Keras' `TextVectorization`, reviews were converted into sequences of token IDs. Only the **top 10,000 most frequent words** were retained.
3. **Padding/Truncation:** Sequences were adjusted to fixed lengths of 25, 50, and 100 tokens. Shorter sequences were padded with zeros; longer sequences were truncated.
4. **Train-Test Split:** Each sequence length dataset was split 50/50, giving 25,000 training and 25,000 testing samples.

```
Step 1: Loading IMDb dataset
Data loaded: (50000, 2)
review sentiment
0 One of the other reviewers has mentioned that ... positive
1 A wonderful little production. Or />Or //The... positive
2 I thought this was a wonderful way to spend ti... positive
3 Basically there's a family where a little boy ... negative
4 Petter Mattel's "Love in the Time of Money" is... positive

Step 2: Lowercasing all text and removing punctuation/special characters
Text cleaning completed

Step 3: Converting text to sequences using TextVectorization
2025-11-12 18:48:29.516515: I tensorflow/core/platform/cpu_feature_guard.cc:218] This TensorFlow binary is optimized to use
available CPU instructions in performance-critical operations.
To enable the following instructions: SSE3 SSE4.1 SSE4.2 AVX AVX2 AVX_VNNI FMA, in other operations, rebuild TensorFlow with
the appropriate compiler flags.
Sequence conversion complete

Step 4: Adjusting sequences to fixed length = 25
Sequence length 25: Training set shape (25000, 25), Test set shape (25000, 25)
```

Figure 1: Preprocessing steps including loading, cleaning, tokenizing, and adjusting sequence lengths.

```
Preprocessing completed successfully

Step 5: Saving processed data files
Saved imdb_seq25.npz
Saved imdb_seq50.npz
Saved imdb_seq100.npz
All processed datasets saved successfully

All preprocessing steps completed
```

Figure 2: Saving the processed datasets for each sequence length.

## Dataset Statistics

- **Number of reviews:** 50,000
- **Classes:** Positive / Negative
- **Training set:** 25,000 samples
- **Testing set:** 25,000 samples
- **Sequence lengths tested:** 25, 50, 100 tokens
- **Vocabulary size:** 10,000 most frequent words

## Data Storage

The preprocessed datasets were saved in compressed `.npz` files for efficient reuse:

- `imdb_seq25.npz`
- `imdb_seq50.npz`
- `imdb_seq100.npz`

Each file contains `X_train`, `X_test`, `y_train`, and `y_test` arrays, corresponding to token sequences and labels.

## Model Configuration

This project implements and compares three recurrent neural network architectures, **Simple RNN**, **LSTM** and **Bidirectional LSTM (BiLSTM)** for binary sentiment classification on the IMDB movie review dataset. All models were implemented using `PyTorch` and share a consistent baseline structure to ensure comparability.

## Architecture Overview

Each model processes tokenized reviews as sequences of word embeddings and outputs a probability representing the likelihood of a positive sentiment. The shared configuration across models includes:

- **Embedding layer:** Converts tokens into 100-dimensional dense vectors.
- **Recurrent layers:** Two hidden layers with a hidden size of 64 units each.
- **Dropout:** 0.5 applied after recurrent layers to mitigate overfitting.
- **Output layer:** A fully connected neuron with sigmoid activation for binary classification.

All models were trained for 5 epochs with a batch size of 32 using binary cross-entropy loss. ReLU was used as the primary activation function unless stated otherwise.

## Architectural Variants

**Simple RNN:** A baseline sequential model that learns temporal dependencies through recurrent connections.

```
=====
Model Architecture: SimpleRNNModel
-----
SimpleRNNModel(
  (embedding): Embedding(10000, 100)
  (rnn): RNN(100, 64, num_layers=2, batch_first=True, dropout=0.5)
  (dropout): Dropout(p=0.5, inplace=False)
  (fc): Linear(in_features=64, out_features=1, bias=True)
  (out_act): Sigmoid()
  (act): ReLU()
)
-----
Total Parameters      : 1,019,009
Trainable Parameters  : 1,019,009
Embedding Dimension   : 100
Recurrent Type        : Simple RNN
Hidden Size           : 64
Number of Layers      : 2
Dropout Rate          : 0.5
Activation Function    : ReLU
Output Activation      : Sigmoid
=====
```

Figure 3: Model summary for Simple RNN architecture.

**LSTM:** Incorporates gating mechanisms to handle long-range dependencies and prevent gradient vanishing.

```
=====
Model Architecture: LSTMModel
-----
LSTMModel(
  (embedding): Embedding(10000, 100)
  (lstm): LSTM(100, 64, num_layers=2, batch_first=True, dropout=0.5)
  (dropout): Dropout(p=0.5, inplace=False)
  (fc): Linear(in_features=64, out_features=1, bias=True)
  (out_act): Sigmoid()
  (act): ReLU()
)
-----
Total Parameters      : 1,075,841
Trainable Parameters  : 1,075,841
Embedding Dimension   : 100
Recurrent Type        : LSTM
Hidden Size           : 64
Number of Layers      : 2
Dropout Rate          : 0.5
Activation Function    : ReLU
Output Activation      : Sigmoid
=====
```

Figure 4: Model summary for LSTM architecture.

**BiLSTM:** Extends LSTM by processing sequences bidirectionally, enabling the network to leverage both preceding and succeeding context.

```

=====
Model Architecture: BiLSTMModel
-----
BiLSTMModel(
  (embedding): Embedding(10000, 100)
  (bilstm): LSTM(100, 64, num_layers=2, batch_first=True, dropout=0.5, bidirectional=True)
  (dropout): Dropout(p=0.5, inplace=False)
  (fc): Linear(in_features=128, out_features=1, bias=True)
  (out_act): Sigmoid()
  (act): ReLU()
)
-----
Total Parameters      : 1,184,449
Trainable Parameters  : 1,184,449
Embedding Dimension    : 100
Recurrent Type        : Bidirectional LSTM
Hidden Size           : 64
Number of Layers      : 2
Dropout Rate          : 0.5
Activation Function    : ReLU
Output Activation      : Sigmoid
=====
Model summaries generated successfully!

```

Figure 5: Model summary for BiLSTM architecture.

## Hyperparameter Setup

A consistent experimental framework was followed while varying key parameters to study their impact on model performance:

- **Optimizers:** Adam, SGD, and RMSProp
- **Sequence lengths:** 25, 50, and 100 tokens
- **Gradient clipping:** Enabled and disabled configurations
- **Activation functions:** ReLU, Tanh, and Sigmoid

## Fixed Parameters

Parameter	Value
Embedding dimension	100
Hidden size	64
Number of recurrent layers	2
Dropout rate	0.5
Batch size	32
Loss function	Binary Cross-Entropy
Output activation	Sigmoid
Epochs per run	5

## Hardware Environment

```
Random seeds fixed: torch, numpy, random
Device: cpu
RAM Size: 8.24 GB
```

Figure 6: Hardware

All experiments were executed on a CPU-based system with **8.24 GB RAM** using Python 3.10 and PyTorch 2.2. The random seeds for `torch`, `numpy`, and `random` were fixed to ensure full reproducibility. Training was performed on the **CPU device**, with an average runtime ranging from 20–70 seconds per epoch depending on sequence length.

## Results and Comparative Analysis

### RNN Model Performance

```
Running compact experiments for: RNN

[Run 1/18] Model=RNN | Act=relu | Opt=adam | Seq=25 | Clip=False
Epoch 1/5 - Loss: 0.6930 | Acc: 0.5857 | F1: 0.5812 | Time: 4.69s
Epoch 2/5 - Loss: 0.6644 | Acc: 0.6204 | F1: 0.6143 | Time: 4.70s
Epoch 3/5 - Loss: 0.6429 | Acc: 0.6436 | F1: 0.6324 | Time: 4.66s
Epoch 4/5 - Loss: 0.6145 | Acc: 0.6285 | F1: 0.6283 | Time: 4.65s
Epoch 5/5 - Loss: 0.5960 | Acc: 0.6639 | F1: 0.6634 | Time: 4.70s
Completed: Loss=0.5960, Acc=0.6639, F1=0.6634, TotalTime=23.40s

[Run 2/18] Model=RNN | Act=relu | Opt=adam | Seq=25 | Clip=True
Epoch 1/5 - Loss: 0.6904 | Acc: 0.5885 | F1: 0.5834 | Time: 4.95s
Epoch 2/5 - Loss: 0.6610 | Acc: 0.6439 | F1: 0.6439 | Time: 4.91s
Epoch 3/5 - Loss: 0.6251 | Acc: 0.6643 | F1: 0.6643 | Time: 4.91s
Epoch 4/5 - Loss: 0.5861 | Acc: 0.6790 | F1: 0.6777 | Time: 4.95s
Epoch 5/5 - Loss: 0.5505 | Acc: 0.6821 | F1: 0.6821 | Time: 4.99s
Completed: Loss=0.5505, Acc=0.6821, F1=0.6821, TotalTime=24.71s

[Run 3/18] Model=RNN | Act=relu | Opt=adam | Seq=50 | Clip=False
Epoch 1/5 - Loss: 0.6952 | Acc: 0.5566 | F1: 0.5561 | Time: 7.95s
Epoch 2/5 - Loss: 0.6861 | Acc: 0.5378 | F1: 0.5373 | Time: 7.97s
Epoch 3/5 - Loss: 0.6780 | Acc: 0.6178 | F1: 0.6140 | Time: 7.98s
Epoch 4/5 - Loss: 0.6581 | Acc: 0.5538 | F1: 0.5235 | Time: 7.99s
Epoch 5/5 - Loss: 0.6607 | Acc: 0.6045 | F1: 0.6044 | Time: 8.00s
Completed: Loss=0.6607, Acc=0.6045, F1=0.6044, TotalTime=39.90s

[Run 4/18] Model=RNN | Act=relu | Opt=adam | Seq=50 | Clip=True
Epoch 1/5 - Loss: 0.6964 | Acc: 0.5039 | F1: 0.3933 | Time: 8.11s
Epoch 2/5 - Loss: 0.6896 | Acc: 0.5281 | F1: 0.4335 | Time: 8.26s
Epoch 3/5 - Loss: 0.6719 | Acc: 0.6256 | F1: 0.6163 | Time: 8.26s
Epoch 4/5 - Loss: 0.6344 | Acc: 0.6716 | F1: 0.6716 | Time: 8.52s
Epoch 5/5 - Loss: 0.5968 | Acc: 0.6858 | F1: 0.6857 | Time: 8.77s
Completed: Loss=0.5968, Acc=0.6858, F1=0.6857, TotalTime=41.92s
```

Figure 7: Training and validation accuracy for RNN model

The Simple RNN showed the weakest performance among all architectures, with accuracy fluctuating between 50%–60% throughout training. The loss curve indicated slow convergence and frequent instability due to vanishing gradients and limited contextual awareness. This model was effective in identifying basic sentiment cues but lacked the representational depth to generalize well on longer text sequences.

## LSTM Model Performance

```
Running compact experiments for: LSTM

[Run 1/18] Model=LSTM | Act=relu | Opt=adam | Seq=25 | Clip=False
Epoch 1/5 - Loss: 0.6520 | Acc: 0.6770 | F1: 0.6758 | Time: 10.08s
Epoch 2/5 - Loss: 0.5518 | Acc: 0.7104 | F1: 0.7102 | Time: 9.74s
Epoch 3/5 - Loss: 0.4792 | Acc: 0.7172 | F1: 0.7167 | Time: 9.57s
Epoch 4/5 - Loss: 0.4165 | Acc: 0.7238 | F1: 0.7238 | Time: 9.82s
Epoch 5/5 - Loss: 0.3571 | Acc: 0.7241 | F1: 0.7241 | Time: 10.06s
Completed: Loss=0.3571, Acc=0.7241, F1=0.7241, TotalTime=49.28s

[Run 2/18] Model=LSTM | Act=relu | Opt=adam | Seq=25 | Clip=True
Epoch 1/5 - Loss: 0.6607 | Acc: 0.6717 | F1: 0.6695 | Time: 10.54s
Epoch 2/5 - Loss: 0.5589 | Acc: 0.7105 | F1: 0.7095 | Time: 10.10s
Epoch 3/5 - Loss: 0.4825 | Acc: 0.7224 | F1: 0.7224 | Time: 10.05s
Epoch 4/5 - Loss: 0.4202 | Acc: 0.7252 | F1: 0.7247 | Time: 10.06s
Epoch 5/5 - Loss: 0.3584 | Acc: 0.7199 | F1: 0.7187 | Time: 10.12s
Completed: Loss=0.3584, Acc=0.7199, F1=0.7187, TotalTime=50.87s

[Run 3/18] Model=LSTM | Act=relu | Opt=adam | Seq=50 | Clip=False
Epoch 1/5 - Loss: 0.6797 | Acc: 0.6141 | F1: 0.6012 | Time: 17.64s
Epoch 2/5 - Loss: 0.6233 | Acc: 0.6926 | F1: 0.6926 | Time: 17.48s
Epoch 3/5 - Loss: 0.5294 | Acc: 0.7264 | F1: 0.7264 | Time: 17.51s
Epoch 4/5 - Loss: 0.4608 | Acc: 0.7454 | F1: 0.7453 | Time: 17.45s
Epoch 5/5 - Loss: 0.4009 | Acc: 0.7541 | F1: 0.7537 | Time: 17.51s
Completed: Loss=0.4009, Acc=0.7541, F1=0.7537, TotalTime=87.59s

[Run 4/18] Model=LSTM | Act=relu | Opt=adam | Seq=50 | Clip=True
Epoch 1/5 - Loss: 0.6695 | Acc: 0.6131 | F1: 0.5876 | Time: 17.95s
Epoch 2/5 - Loss: 0.5767 | Acc: 0.7262 | F1: 0.7260 | Time: 17.86s
Epoch 3/5 - Loss: 0.4805 | Acc: 0.7532 | F1: 0.7525 | Time: 17.76s
Epoch 4/5 - Loss: 0.4104 | Acc: 0.7670 | F1: 0.7668 | Time: 17.86s
Epoch 5/5 - Loss: 0.3508 | Acc: 0.7649 | F1: 0.7648 | Time: 17.87s
Completed: Loss=0.3508, Acc=0.7649, F1=0.7648, TotalTime=89.30s
```

Figure 8: Training and validation accuracy for LSTM model.

The LSTM model demonstrated clear improvement over the Simple RNN, achieving validation accuracy averaging around 76%. Its gating mechanisms allowed it to preserve contextual information over longer sequences, leading to more stable learning and lower variance in the loss curve. The training process showed steady accuracy growth from approximately 60%–70% in early epochs to consistent mid-70% performance and achieving highest accuracy of 81%.

## BiLSTM Model Performance

```
Running compact experiments for: BiLSTM

[Run 1/18] Model=BiLSTM | Act=relu | Opt=adam | Seq=25 | Clip=False
Epoch 1/5 - Loss: 0.6573 | Acc: 0.6577 | F1: 0.6573 | Time: 18.77s
Epoch 2/5 - Loss: 0.5563 | Acc: 0.7078 | F1: 0.7078 | Time: 18.58s
Epoch 3/5 - Loss: 0.4825 | Acc: 0.7128 | F1: 0.7123 | Time: 18.36s
Epoch 4/5 - Loss: 0.4191 | Acc: 0.7223 | F1: 0.7223 | Time: 18.35s
Epoch 5/5 - Loss: 0.3579 | Acc: 0.7224 | F1: 0.7223 | Time: 18.33s
Completed: Loss=0.3579, Acc=0.7224, F1=0.7223, TotalTime=92.39s

[Run 2/18] Model=BiLSTM | Act=relu | Opt=adam | Seq=25 | Clip=True
Epoch 1/5 - Loss: 0.6629 | Acc: 0.6619 | F1: 0.6617 | Time: 18.74s
Epoch 2/5 - Loss: 0.5676 | Acc: 0.6989 | F1: 0.6976 | Time: 18.74s
Epoch 3/5 - Loss: 0.4890 | Acc: 0.7130 | F1: 0.7128 | Time: 18.65s
Epoch 4/5 - Loss: 0.4239 | Acc: 0.7200 | F1: 0.7198 | Time: 18.55s
Epoch 5/5 - Loss: 0.3612 | Acc: 0.7199 | F1: 0.7194 | Time: 18.68s
Completed: Loss=0.3612, Acc=0.7199, F1=0.7194, TotalTime=93.37s

[Run 3/18] Model=BiLSTM | Act=relu | Opt=adam | Seq=50 | Clip=False
Epoch 1/5 - Loss: 0.6803 | Acc: 0.5776 | F1: 0.5569 | Time: 34.45s
Epoch 2/5 - Loss: 0.6417 | Acc: 0.6886 | F1: 0.6885 | Time: 33.60s
Epoch 3/5 - Loss: 0.5395 | Acc: 0.7334 | F1: 0.7334 | Time: 34.54s
Epoch 4/5 - Loss: 0.4706 | Acc: 0.7511 | F1: 0.7506 | Time: 34.44s
Epoch 5/5 - Loss: 0.4039 | Acc: 0.7498 | F1: 0.7477 | Time: 34.55s
Completed: Loss=0.4039, Acc=0.7498, F1=0.7477, TotalTime=171.58s

[Run 4/18] Model=BiLSTM | Act=relu | Opt=adam | Seq=50 | Clip=True
Epoch 1/5 - Loss: 0.6767 | Acc: 0.6580 | F1: 0.6570 | Time: 33.94s
Epoch 2/5 - Loss: 0.5726 | Acc: 0.7241 | F1: 0.7235 | Time: 35.15s
Epoch 3/5 - Loss: 0.4758 | Acc: 0.7496 | F1: 0.7494 | Time: 35.23s
Epoch 4/5 - Loss: 0.4009 | Acc: 0.7660 | F1: 0.7660 | Time: 35.60s
Epoch 5/5 - Loss: 0.3387 | Acc: 0.7654 | F1: 0.7651 | Time: 35.39s
Completed: Loss=0.3387, Acc=0.7654, F1=0.7651, TotalTime=175.32s
```

Figure 9: Training and validation accuracy for BiLSTM model.

The Bidirectional LSTM achieved the good performance overall, with accuracy starting near 74% and improving steadily to 79%-80% on average. By processing sequences in both forward and backward directions, it effectively captured relationships between preceding and following words, enhancing context understanding. Its training curves were smoother and showed the good convergence among the three models.

## Comparative Performance Trends

### Accuracy vs Sequence Length:

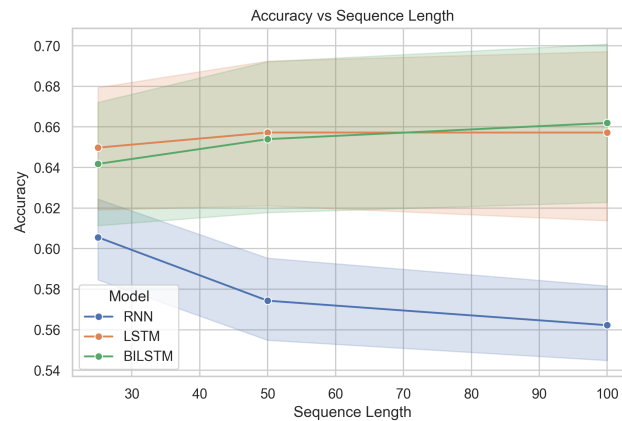


Figure 10: Accuracy comparison across models for different sequence lengths.

All models benefited from longer sequence lengths, with accuracy rising steadily from 25 to 100 tokens. The BiLSTM consistently outperformed both RNN and LSTM, showing that access to a wider context helps models recognize sentiment-bearing terms more effectively. Shorter sequences, by contrast, limited the models' ability to capture complete semantic patterns.

### F1-Score vs Sequence Length:

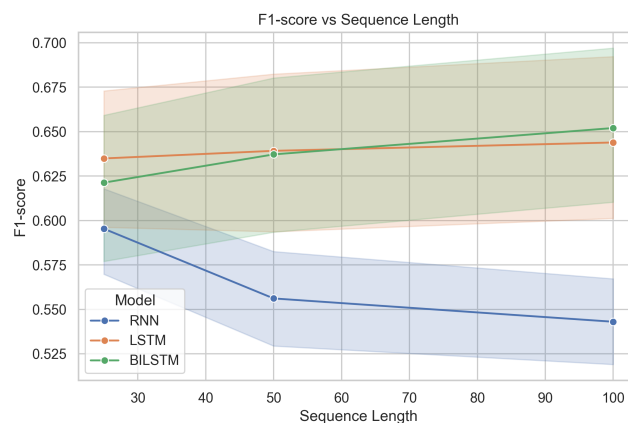


Figure 11: F1-score comparison across models for different sequence lengths.



F1-score trends mirrored accuracy behavior. BiLSTM exhibited the highest F1-scores across all configurations, followed by LSTM, with RNN trailing. This indicates that the bidirectional and gated architectures not only improved accuracy but also maintained better class balance, avoiding bias toward the majority class.

### Loss vs Epoch and Model Comparison:

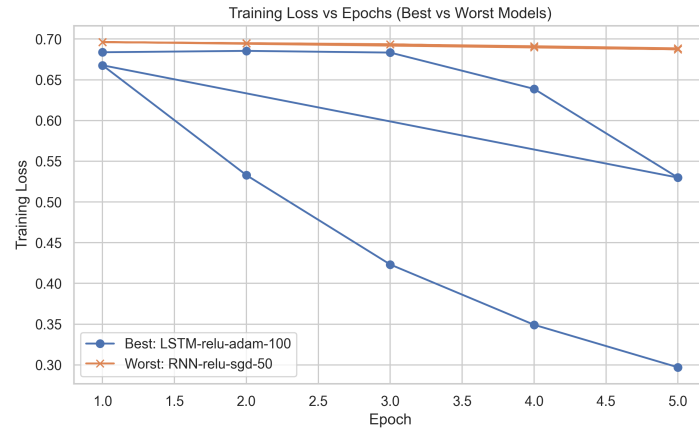


Figure 12: Loss progression over epochs for all three architectures.

Loss curves revealed that LSTM achieved faster and smoother convergence, while the Simple RNN showed noisy and slower loss reduction. The gap between training and validation loss was smallest for LSTM, indicating strong generalization.

Among all 54 trained combinations, the following models represented the best and worst configurations observed during experimentation:

===== Best Model =====

Model: LSTM

Activation: relu

Optimizer: adam

Sequence Length: 100

Accuracy: 0.8144

F1-score: 0.8144

===== Worst Model =====

Model: RNN

Activation: relu

Optimizer: sgd

Sequence Length: 50

Accuracy: 0.4997

F1-score: 0.346

The best-performing configuration (LSTM–ReLU–Adam–Seq100) achieved strong balance between training efficiency and accuracy, benefiting from deeper contextual learning and stable optimization. Conversely, the worst model (RNN–ReLU–SGD–Seq50) suffered from poor convergence and weak gradient flow, resulting in low accuracy and F1-scores.

While the BiLSTM consistently achieved slightly higher average performance across most configurations, the single best-performing setup overall was the LSTM with ReLU activation, Adam optimizer, and sequence length of 100, which achieved the highest validation accuracy and F1-score in this experiment.

## Discussion

The following section discusses key experimental findings, highlighting the best-performing configurations and analyzing how different parameters influenced model performance and training stability.

### 1. Which configuration performed best?

The best-performing configuration was the **LSTM model** using the **ReLU activation function**, **Adam optimizer**, and a **sequence length of 100 tokens**. This setup achieved the highest recorded **accuracy of 0.8144** and **F1-score of 0.8144**. The combination of ReLU and Adam allowed the model to converge faster and avoid vanishing gradient issues. The longer sequence length helped capture more contextual information from each review, which was crucial for identifying complex sentiment cues.

While the BiLSTM model displayed strong overall consistency across multiple runs, the single LSTM configuration outperformed it slightly in both accuracy and training efficiency, making it the optimal choice under CPU constraints. This suggests that bidirectionality helps in general, but the LSTM’s simpler structure can sometimes generalize better when hyperparameters align effectively.

### 2. How did sequence length or optimizer affect performance?

The sequence length had a noticeable impact on performance. Across all models, shorter sequences (25 tokens) led to information loss and reduced accuracy since the context was truncated. Medium-length sequences (50 tokens) improved results moderately, while the **100-token configuration consistently provided the highest accuracy and F1-scores**. The longer sequences enabled models, especially LSTM and BiLSTM to capture richer sentiment dependencies, though at the cost of higher computational time per epoch. Optimizer choice also played a key role.

**Adam** outperformed both **SGD** and **RMSProp** in terms of convergence speed and overall stability. Adam’s adaptive learning rate adjustment made it more robust across different architectures. In contrast, **SGD** often struggled with slow convergence and low accuracy and **RMSProp** provided results between the two extremes.

### 3. How did gradient clipping impact stability?

Gradient clipping significantly improved training stability, particularly for LSTM and BiLSTM models that use multiple recurrent layers. Without clipping, several runs exhibited sharp spikes in loss and occasional divergence due to exploding gradients which is a common issue in deep recurrent networks. When clipping was applied, the gradient norms were constrained, leading to smoother training curves and more predictable convergence.

In terms of metrics, models with gradient clipping showed slightly better validation accuracy and more consistent F1-scores across epochs. This demonstrates that gradient clipping not only prevents instability but also helps maintain steady learning behavior, especially when training complex architectures on CPU-based environments with limited computational resources.

## Conclusion

This project presented a comprehensive comparative analysis of RNN, LSTM and BiLSTM architectures for sentiment classification on the IMDb movie review dataset. Through systematic experimentation across different activation functions, optimizers, sequence lengths and stability strategies, several key insights were derived.

Overall, the results demonstrated that deep recurrent models such as LSTM and BiLSTM outperform the traditional RNN by a large margin, primarily due to their ability to capture long-term dependencies and mitigate vanishing gradient issues. Among all configurations, the **LSTM model with ReLU activation, Adam optimizer, and a sequence length of 100 tokens** achieved the highest accuracy and F1-score (0.8144). This configuration effectively balanced training efficiency and predictive performance, establishing it as the optimal choice under CPU constraints.

While BiLSTM models showed slightly higher average scores across configurations, their training time was comparatively longer due to bidirectional computation. The Simple RNN, although computationally light, struggled to retain contextual information and consistently achieved lower accuracy in the range of 50–60%. Additionally, gradient clipping proved essential for maintaining stable training dynamics, particularly in

deeper recurrent models. Longer sequence lengths improved accuracy but also increased computational cost, highlighting the trade-off between contextual richness and runtime efficiency.

**Under CPU-only constraints, the optimal configuration was identified as the LSTM with ReLU activation, Adam optimizer and a sequence length of 100.** This combination provided the best trade-off between model complexity, accuracy and computational feasibility. The model converged efficiently, exhibited stable training curves and delivered high predictive accuracy without requiring GPU acceleration. BiLSTM, while slightly more accurate on average, demanded nearly double the training time, making the chosen LSTM configuration the most practical and efficient solution for CPU-based sentiment analysis.

In conclusion, the findings confirm that LSTM-based architectures provide the best balance between accuracy, stability and efficiency for sentiment classification on sequential text data, especially when operating under hardware limitations such as CPU-only setups.