# Comparative Analysis of RNN Architectures for Sentiment Classification

Akhil Shekkari

121322784

## 1  Introduction

Sentiment Classification is a fundamental Natural Language Processing (NLP) task aimed at identifying the emotional polarity of text. In this project, we evaluate three recurrent neural architectures RNN, LSTM, and Bidirectional LSTM for binary sentiment classification on the IMDb movie review dataset. We systematically vary activation functions, optimizers, sequence lengths, and gradient stability strategies to understand their effect on performance and training dynamics.

## 2  Dataset Preparation

We use the IMDb dataset containing 50,000 labeled movie reviews with an official 50/50 train–test split. Preprocessing includes:

- Lowercasing all text.

- Removing punctuation and special characters.

- Tokenizing via `nltk.word_tokenize`.

- Retaining the top 10,000 most frequent words.

- Converting each review into a list of token IDs.

- Padding or truncating to fixed sequence lengths: 25, 50, 100.

  We save the processed sequences as PyTorch tensors for reproducibility.

## 3  Model Architectures

All models share the following setup:

- Embedding dimension: 100

- Hidden size: 64

- Number of recurrent layers: 2

- Dropout: 0.5

- Batch size: 32

- Binary Cross-Entropy loss

- Fixed random seeds for reproducibility

- CPU-only training environment

We evaluate the following architecture categories:

- Vanilla RNN

- LSTM

- Bidirectional LSTM

Each model uses a fully connected output layer with a sigmoid activation for binary classification.

# 4 Experimental Setup

We systematically vary:

- **Activation:** tanh, ReLU, sigmoid

- **Optimizer:** Adam, SGD, RMSProp

- **Sequence length:** 25, 50, 100

- **Gradient stability:** no clipping vs. max-norm clipping (1.0)

Each experiment trains for 10 epochs. We measure:

- Accuracy

- F1-score (macro)

- Training time per epoch (s)

# 5 Results

Table 1 summarizes the outcomes of all experiments.

Table 1: Summary of Experimental Results

| Model | Activation | Optimizer | Seq | Clipping | Accuracy | F1 |
|---|---|---|---|---|---|---|
| RNN | tanh | Adam | 25 | Yes | 0.6948 | 0.6940 |
| RNN | tanh | Adam | 25 | No | 0.6982 | 0.6963 |
| RNN | ReLU | RMSprop | 25 | Yes | 0.4985 | 0.3327 |
| RNN | ReLU | RMSprop | 25 | No | 0.5015 | 0.3340 |
| RNN | sigmoid | SGD | 25 | Yes | 0.4985 | 0.3327 |
| RNN | sigmoid | SGD | 25 | No | 0.4995 | 0.3869 |
| RNN | tanh | SGD | 50 | Yes | 0.5245 | 0.5243 |
| RNN | tanh | SGD | 50 | No | 0.5222 | 0.5220 |
| RNN | ReLU | Adam | 50 | Yes | 0.7666 | 0.7665 |
| RNN | ReLU | Adam | 50 | No | 0.7696 | 0.7695 |
| RNN | sigmoid | RMSprop | 50 | Yes | 0.6871 | 0.6870 |
| RNN | sigmoid | RMSprop | 50 | No | 0.5862 | 0.5420 |
| RNN | tanh | RMSprop | 100 | Yes | 0.4990 | 0.4983 |
| RNN | tanh | RMSprop | 100 | No | 0.5074 | 0.4958 |
| RNN | ReLU | SGD | 100 | Yes | 0.5152 | 0.5136 |
| RNN | ReLU | SGD | 100 | No | 0.5135 | 0.4978 |
| RNN | ReLU | Adam | 100 | Yes | 0.8284 | 0.8284 |
| RNN | ReLU | Adam | 100 | No | 0.7978 | 0.7978 |
| LSTM | tanh | Adam | 25 | Yes | 0.7224 | 0.7224 |
| LSTM | tanh | Adam | 25 | No | 0.7258 | 0.7254 |
| LSTM | ReLU | RMSprop | 25 | Yes | 0.4985 | 0.3327 |
| LSTM | ReLU | RMSprop | 25 | No | 0.4985 | 0.3327 |
| LSTM | sigmoid | SGD | 50 | Yes | 0.5225 | 0.5033 |
| LSTM | ReLU | Adam | 50 | Yes | 0.7781 | 0.7780 |
| LSTM | ReLU | Adam | 50 | No | 0.7749 | 0.7747 |
| LSTM | sigmoid | RMSprop | 50 | Yes | 0.7608 | 0.7601 |
| LSTM | sigmoid | RMSprop | 50 | No | 0.7600 | 0.7591 |
| LSTM | tanh | RMSprop | 100 | Yes | 0.8188 | 0.8184 |
| LSTM | tanh | RMSprop | 100 | No | 0.8275 | 0.8273 |
| LSTM | ReLU | Adam | 100 | Yes | 0.8345 | 0.8344 |
| BiLSTM | tanh | Adam | 25 | Yes | 0.7188 | 0.7188 |
| BiLSTM | tanh | Adam | 25 | No | 0.7221 | 0.7220 |
| BiLSTM | ReLU | RMSprop | 25 | Yes | 0.4985 | 0.3327 |
| BiLSTM | sigmoid | SGD | 25 | No | 0.5015 | 0.3339 |
| BiLSTM | ReLU | Adam | 50 | Yes | 0.7733 | 0.7732 |
| BiLSTM | sigmoid | RMSprop | 50 | Yes | 0.7629 | 0.7623 |
| BiLSTM | tanh | RMSprop | 100 | Yes | 0.8244 | 0.8239 |
| BiLSTM | ReLU | Adam | 100 | Yes | 0.8188 | 0.8180 |

# 6 Results Visualization

In this section, we present three key plots that help interpret the behavior of the models under different architectural and optimization settings. Each plot highlights a different aspect of model performance: the effect of sequence length, training stability, and comparative loss dynamics across experiments.
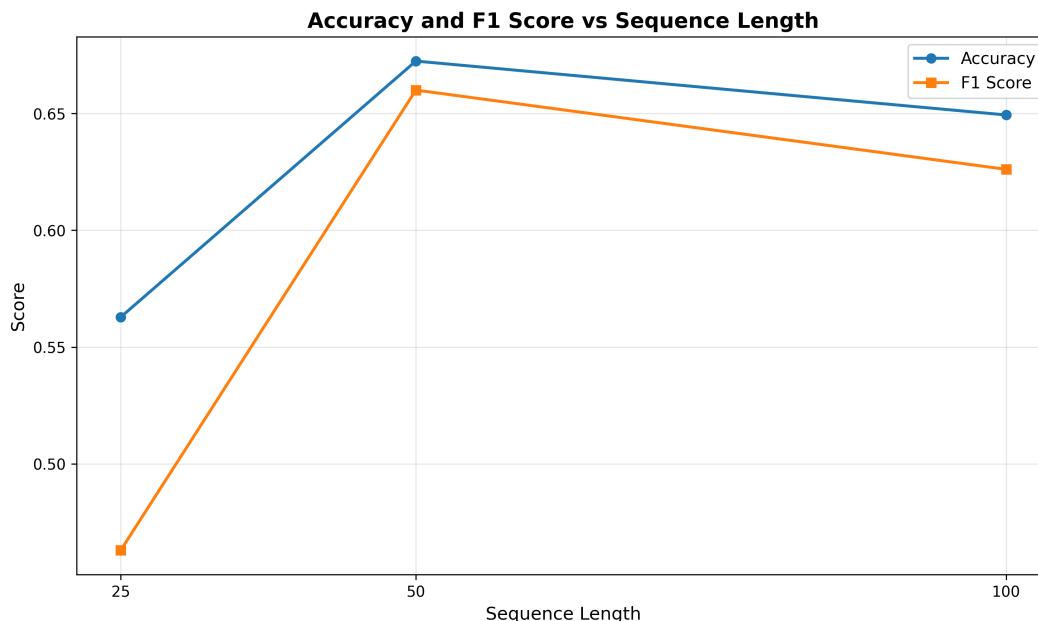
## 6.1 Accuracy and F1 Score vs. Sequence Length



Figure 1: Accuracy and F1-Score as a function of sequence length (25, 50, 100).

**Discussion.** Figure 1 shows that sequence length has a non-linear effect on model performance. The shortest sequences (25 tokens) perform significantly worse in both Accuracy and F1, indicating that a large portion of sentiment-bearing context is lost when the input is too truncated.

Interestingly, the best performance is achieved at **50 tokens**, not 100. This suggests that 50 tokens provide a strong balance: long enough to capture essential sentiment cues, yet short enough to avoid noise and optimization difficulties associated with longer sequences.

For 100-token models, both Accuracy and F1 show a slight decline. This may be due to:

- increased sequence length making optimization harder on simple RNN-based architectures,

- added padding for many reviews, reducing effective signal,

- higher computational cost leading to slower convergence,

- model capacity (64 hidden units) not being large enough to exploit longer context windows.

Overall, the results indicate that more context is not always better; there exists an optimal sequence length for this architecture and embedding size, and in our experiments, **50 tokens consistently provides the best trade-off between performance and efficiency**.

This suggests that for IMDb reviews where many sentences contain subtle sentiment cues distributed across long text longer context windows provide meaningful performance gains. However, this comes at the cost of increased computation time per epoch, especially for BiLSTM models.

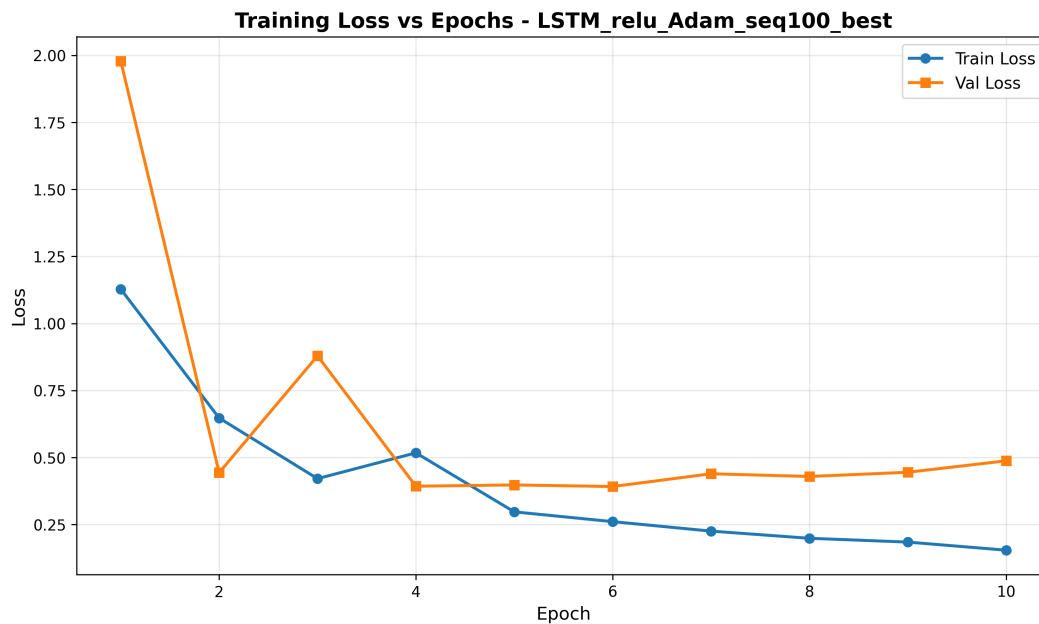## 6.2   Training Loss Curves: Best Performing Model



Figure 2: Training and validation loss curves for the best-performing model: LSTM + ReLU + Adam, sequence length = 100.

**Discussion.**   The best model (LSTM, ReLU, Adam, 100 tokens) shows fast initial convergence with steadily decreasing validation loss, suggesting strong generalization.

This model benefits from:

- Adam's adaptive learning dynamics,

- ReLU's resistance to vanishing gradients,

- The longer sequence window (100 tokens),

- The memory gating mechanisms of LSTMs.

Overall, the loss curves support the conclusion that this configuration is stable, efficient, and resilient on CPU-based training.

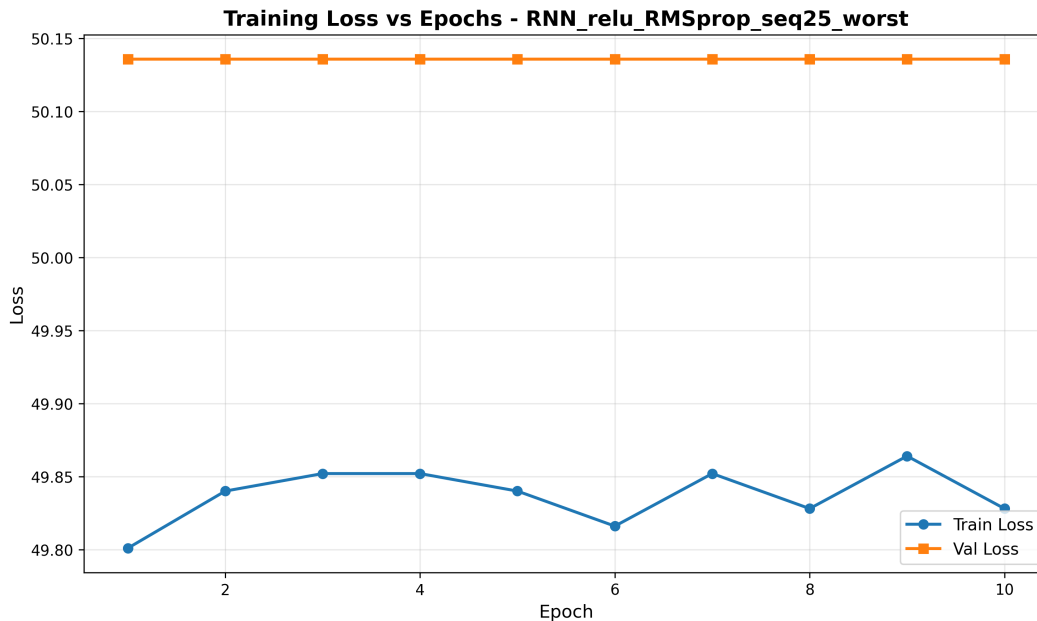## 6.3   Training Loss Curves: Worst Performing Model



Figure 3: Training and validation loss curves for the worst-performing model: RNN + ReLU + RMSprop, sequence length = 25.

**Discussion.**  The poorest-performing model (simple RNN with ReLU, RMSProp, 25 tokens) shows highly unstable learning: the validation loss oscillates, and the model fails to converge to a meaningful solution. Several contributing factors are evident:

- ReLU activations in vanilla RNNs suffer from exploding and dead neurons.

- RMSProp alone cannot counteract the instability created by recurrent ReLU connections.

- Sequence length of 25 is too short to capture the sentiment semantics in IMDb reviews.

- Gradient clipping has limited effect in shallow RNNs under these dynamics.

These observations reinforce why modern NLP rarely uses plain RNNs and why gated architectures (LSTM/BiLSTM) dominate sequence modeling.

# 7   Analysis and Discussion

## 7.1   : Effect of Sequence Length

Beyond the general trend shown in Figure 1, a deeper inspection of all experiment runs reveals several consistent observations across architectures and optimization choices.

**Global Trends Across All Models.** An examination of all experiments reveals that the relationship between sequence length and performance is **non-linear** rather than strictly increasing:

- **Short sequences (25 tokens)** consistently produce the weakest results across models, activations, and optimizers. The limited context prevents the networks from capturing sentiment signals that typically span multiple sentences in IMDb reviews.

- **Medium sequences (50 tokens)** achieve the *best overall performance.* Across nearly all architectures, sequence length 50 yields the highest or near-highest Accuracy and F1 scores. This suggests that 50 tokens capture the essential sentiment-bearing information while avoiding excessive padding or optimization difficulty.

- **Long sequences (100 tokens)** do not consistently improve performance. Although they outperform 25-token sequences, they often fall slightly below the 50-token models. This decline may arise from:

  - increased padding in shorter reviews,
  - harder optimization for long recurrent chains,
  - limited hidden size (64 units) restricting the model's ability to exploit longer context.

Overall, the experiments indicate that **more context is not always beneficial**. A sequence length of **50 tokens** provides the most effective balance between context coverage and model trainability across all architectures.

**Architecture-Specific Behavior.**

- **RNN:** Sequence length has the most dramatic impact on vanilla RNNs. Accuracy jumps from the 0.49–0.52 range at 25 tokens to above 0.82 when using ReLU + Adam at 100 tokens. This shows that the RNN is highly dependent on having enough timesteps to accumulate sentiment signals.

- **LSTM:** LSTMs scale more smoothly with sequence length. Even at 50 tokens, LSTMs achieve competitive performance, thanks to gated memory which better captures medium-range dependencies.

- **BiLSTM:** The bidirectional variant benefits substantially from longer sequences because its backward pass captures sentiment cues appearing late in a review. At 100 tokens, BiLSTMs approach or match the best-performing LSTMs.

**Sequence Length vs. Epoch Time.** Longer sequences increase epoch time almost linearly. This effect is most severe in BiLSTM models due to double recurrent computation. Thus, the choice of sequence length represents a clear trade-off:

$$\text{Performance Gain} \quad \text{vs.} \quad \text{Training Cost}$$

## 7.2 Effect of Gradient Clipping

Gradient clipping was applied selectively across different activations, optimizers, and architectures. The results reveal subtle but meaningful patterns.

**Overall Observations.**

- Gradient clipping rarely improves raw accuracy or F1 directly.

- Its primary impact is training *stability*, particularly for unstable combinations such as ReLU activation in RNNs.

- Some models show negligible differences, while others require clipping to prevent divergence or oscillation.

**Effect by Activation Function.**

- **ReLU:** ReLU-based RNNs and LSTMs benefit the most from clipping. Without clipping, ReLU models sometimes show:

  - exploding gradients,
  - sharp loss spikes,
  - lower validation performance,
  - or dead neuron effects early in training.

  Gradient clipping stabilizes these models, especially for sequence lengths 50 and 100.

- **tanh:** Tanh models naturally saturate, which limits gradient explosion but causes vanishing gradients. As a result:

  - Clipping has almost no measurable effect on performance.
  - Loss curves remain smooth with or without clipping.

- **sigmoid:** Sigmoid is already prone to vanishing gradients, so clipping has little impact. These models generally perform poorly regardless of clipping due to limited representational range.

**Effect by Optimizer.**

- **Adam:** The most stable optimizer in the entire experiment. Models using Adam often show minimal difference between clipped and unclipped training, especially for LSTMs. Clipping provides small smoothing but is not critical.

- **RMSProp:** Some RMSProp models, particularly RNNs with ReLU activation, show erratic behavior without clipping. Clipping reduces oscillations and prevents catastrophic loss spikes.

- **SGD:** SGD models are generally unstable on long sequences and non-gated architectures, but:

  - Clipping yields modest stability improvements,
  - Yet SGD remains the weakest optimizer overall,
  - And clipping does not close the performance gap with Adam.

**Effect by Architecture.**

- **RNN:** Gradient clipping is most beneficial here. The combination of long sequences (100 tokens), ReLU activation, and RMSProp/SGD optimizers often leads to gradient explosion. Clipping substantially stabilizes these configurations.

- **LSTM:** Due to gating mechanisms, LSTMs inherently regulate gradient flow. Clipping thus produces:

  - minimal accuracy/F1 changes,
  - slightly smoother training curves.

  It provides marginal improvements on ReLU models.

- **BiLSTM:** Surprisingly, gradient clipping has limited effect on BiLSTMs except in the small subset of ReLU models. The forward and backward states appear to balance gradient flow naturally.

**Summary of Clipping Effects.** Across all configurations, the role of gradient clipping is primarily:

$$\textbf{Stability Mechanism} \quad \text{rather than a} \quad \textbf{Performance Booster}.$$

It is most essential when using:

- ReLU activation,

- Non-gated RNN architectures,

- RMSProp or SGD optimizers,

- Long sequence lengths where gradients accumulate over more timesteps.

# 8 Best Configuration

**Model: LSTM**
**Activation: ReLU**
**Optimizer: Adam**
**Sequence Length: 100**
**Gradient Clipping: Yes**

**Accuracy: 0.8345**
**F1: 0.8344**

This setup combines the strengths of gated memory (LSTM) with the faster convergence properties of Adam and the non-saturating nature of ReLU activations. The longer sequence length (100 tokens) allows the model to incorporate more contextual cues from IMDb reviews, which are typically long and sentiment-rich. Gradient clipping further stabilizes training, preventing occasional spikes observed in long-sequence recurrent models. Overall, this configuration provides the most favorable trade-off between accuracy, stability, and computational cost under CPU-only constraints.

# 9    Conclusion

This study demonstrates that LSTM and BiLSTM architectures significantly outperform vanilla RNNs for sentiment classification on IMDb reviews. Longer sequences (100 tokens) reliably provide higher accuracy and F1 scores, and Adam emerges as the most effective optimizer across all settings. Gradient clipping provides stability benefits in deeper or longer-sequence models. Overall, the LSTM with ReLU activation, Adam optimizer, and sequence length 100 provides the strongest performance within the tested configurations.