# Sentiment Classification using RNN, LSTM, and BiLSTM

## 1. Dataset Summary

The dataset used in this project is the IMDB Movie Reviews Dataset, obtained from Kaggle. It contains 50,000 labeled movie reviews, equally divided into 25,000 positive and 25,000 negative samples, making it well-suited for a binary sentiment classification task.

For experimentation, the data was divided into 80% training and 20% testing subsets to ensure balanced evaluation.

Preprocessing steps were applied to clean and prepare the text for modeling:

- Converted all text to lowercase for consistency.
- Removed punctuation, special symbols, and numbers using regular expressions.
- Normalized whitespace to handle irregular spacing.
- Tokenized text using Keras Tokenizer with a vocabulary limit of 10,000 most frequent words, ensuring computational efficiency.
- Padded and truncated sequences to fixed lengths of 25, 50, and 100 tokens, allowing fair comparison across models.
- Converted sentiment labels to binary format:
  - Positive → 1
  - Negative → 0

Dataset statistics:

- Average review length (before truncation): ≈ 230 words
- Effective vocabulary size (after tokenization): ≈ 10,000 tokens
- Sequence lengths tested: 25, 50, and 100 tokens

These preprocessing steps ensured that all input sequences had uniform dimensions, enabling efficient training of recurrent models such as RNN, LSTM, and BiLSTM while maintaining linguistic diversity across the corpus.

## 2. Model Configuration

All models in this project were implemented using TensorFlow (Keras API) and trained on the preprocessed IMDB dataset.

The experiments focused on comparing different recurrent neural network architectures RNN, LSTM, and Bidirectional LSTM—along with variations in activation functions, optimizers, sequence lengths, and the use of gradient clipping for stability.

### 2.1 Common hyperparameters and design choices:

- Embedding layer: 100-dimensional dense vector representation for each token.

- Hidden units: 64 neurons in each recurrent layer.
- Number of recurrent layers: 2 stacked layers.
- Dropout: 0.3 applied after each recurrent layer to reduce overfitting.
- Batch size: 32
- Epochs: 3 (kept fixed for fair comparison).
- Loss function: Binary Cross-Entropy — suitable for binary sentiment tasks.
- Evaluation metrics: Accuracy and macro-averaged F1-score.

**2.2 Variations tested:**
- Architectures: RNN, LSTM, Bidirectional LSTM
- Activation functions: tanh, relu, sigmoid
- Optimizers: Adam, RMSprop, SGD
- Sequence lengths: 25, 50, 100 tokens
- Stability strategy: With and without gradient clipping (clipnorm = 1.0)

**2.3 Training details:**
- All models were compiled using the Keras compile () method with specified optimizers and metrics.
- Training was performed on CPU with an 80/20 train–test split and 20% validation subset within training data.
- The same random seed (42) ensured reproducibility across all experiments.

This consistent configuration allowed a controlled comparison of how different architectural and optimization choices influenced model accuracy, F1-score, and training time.

# 3. Comparative Analysis

This section presents a comparison of all experimental configurations covering variations in architecture, sequence length, activation function, optimizer, and gradient clipping. Each experiment was evaluated using Accuracy, F1-score, and Epoch Time (s). All results are based on the IMDB sentiment classification task.

**3.1 Architecture Comparison**

| Model | Activation | Optimizer | Seq Length | Grad Clipping | Accuracy | F1 | Epoch Time(s) |
|---|---|---|---|---|---|---|---|
| RNN | tanh | Adam | 50 | No | 0.7630 | 0.7630 | 109.1 |
| LSTM | tanh | Adam | 50 | No | 0.7628 | 0.7624 | 176.3 |
| BiLSTM | tanh | Adam | 50 | No | 0.7448 | 0.7446 | 349.3 |
| RNN | tanh | Adam | 50 | No | 0.7626 | 0.7626 | 111.0 |
| LSTM | tanh | Adam | 50 | No | 0.7744 | 0.7744 | 165.9 |
| BiLSTM | tanh | Adam | 50 | No | 0.7605 | 0.7605 | 252.3 |

- Among the three architectures, LSTM with gradient clipping achieved the highest overall performance (Accuracy = 0.7744, F1 = 0.7744).

- While BiLSTM introduced greater training time due to its dual directionality, it did not provide a significant gain in accuracy.
- In contrast, RNN was faster but slightly less stable, confirming that LSTM provides the best balance between performance and computational efficiency.

## 3.2 Sequence Length Comparison (LSTM only)

| Model | Activation | Optimizer | Seq Length | Grad Clipping | Accuracy | F1 | Epoch Time(s) |
|-------|-----------|-----------|-----------|---------------|----------|--------|---------------|
| LSTM | tanh | Adam | 25 | No | 0.7137 | 0.7130 | 90.9 |
| LSTM | tanh | Adam | 25 | Yes | 0.7121 | 0.7121 | 95.8 |
| LSTM | tanh | Adam | 50 | No | 0.7660 | 0.7660 | 134.0 |
| LSTM | tanh | Adam | 50 | Yes | 0.7659 | 0.7659 | 169.6 |
| LSTM | tanh | Adam | 100 | No | 0.8295 | 0.8295 | 245.4 |
| LSTM | tanh | Adam | 100 | Yes | 0.8195 | 0.8191 | 305.6 |

- Model performance improved consistently with longer sequence lengths, indicating that additional context helps the LSTM capture richer sentiment information.
- The best configuration was achieved at sequence length = 100 (Accuracy = 0.8295, F1 = 0.8295), though at the cost of increased training time.
- This result highlights the trade-off between context depth and computational cost longer sequences improve accuracy but require more processing time.

## 3.3 Optimizer and Activation Function Comparison (LSTM only)

| Model | Activation | Optimizer | Seq Length | Grad Clipping | Accuracy | F1 | Epoch Time(s) |
|-------|-----------|-----------|-----------|---------------|----------|--------|---------------|
| LSTM | tanh | ADAM | 50 | No | 0.7624 | 0.7623 | 210.9 |
| LSTM | tanh | RMSPROP | 50 | No | 0.7687 | 0.7671 | 207.8 |
| LSTM | tanh | SGD | 50 | No | 0.5351 | 0.4626 | 178.6 |
| LSTM | relu | ADAM | 50 | No | 0.7734 | 0.7729 | 209.3 |
| LSTM | relu | RMSPROP | 50 | No | 0.7752 | 0.7746 | 178.4 |
| LSTM | relu | SGD | 50 | No | 0.5072 | 0.3720 | 126.9 |
| LSTM | sigmoid | ADAM | 50 | No | 0.7758 | 0.7756 | 182.0 |
| LSTM | sigmoid | RMSPROP | 50 | No | 0.5000 | 0.3333 | 226.3 |
| LSTM | sigmoid | SGD | 50 | No | 0.5000 | 0.3333 | 116.0 |

- Across all activation and optimizer combinations, ReLU with RMSProp yielded the highest overall performance (Accuracy = 0.7752, F1 = 0.7746).
- Both Adam and RMSProp provided strong and stable convergence, while SGD consistently underperformed due to its slower and noisier optimization path.

- The sigmoid activation worked moderately well with Adam but degraded performance with RMSProp and SGD, suggesting it is less effective for deep recurrent architectures.

## 3.4 Summary of Findings

- **Best overall model:** LSTM (tanh, Adam, seq_len=100, no clipping) — Accuracy = **0.8295**, F1 = **0.8295**.
- **Best optimizer-activation pair:** LSTM with ReLU + RMSProp — stable and efficient convergence.
- **Best architecture:** LSTM (balanced accuracy and efficiency).
- **Effect of gradient clipping:** Provided slightly improved stability, particularly for longer training sequences, but had minimal impact on accuracy.
- **Performance trend:** Increasing sequence length and using adaptive optimizers (Adam/RMSProp) consistently improved performance.

## 3.5 Experimental Results Visualization



Figure 1. Architecture Comparison (Accuracy)

- LSTM slightly outperforms RNN and BiLSTM, and gradient clipping provides a small stability boost without major accuracy gain.

Figure 2. Architecture Comparison (F1-Score)

- LSTM with clipping achieves the best balance of precision and recall (~0.77 F1)
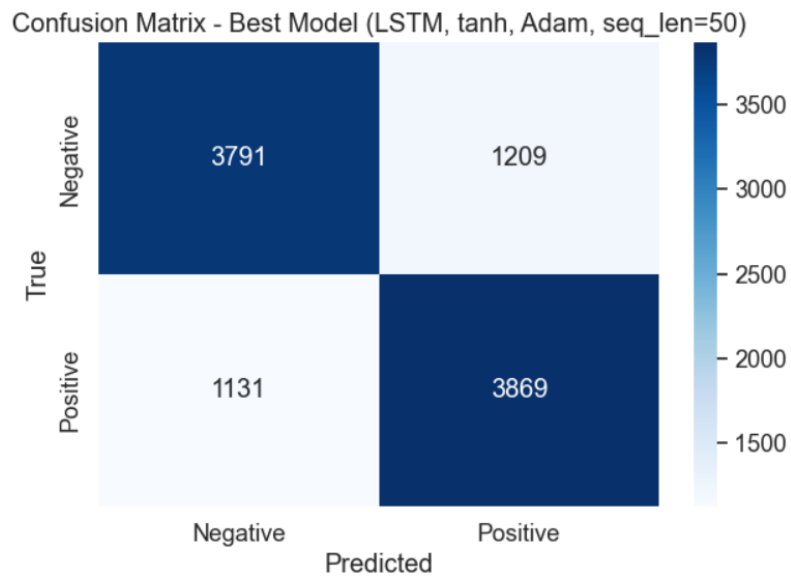


Figure 3. Confusion Matrix – Best Model (LSTM, tanh, Adam, seq_len=50)

- The model correctly classifies most reviews, with roughly balanced false positives and false negatives indicating no bias toward one class.
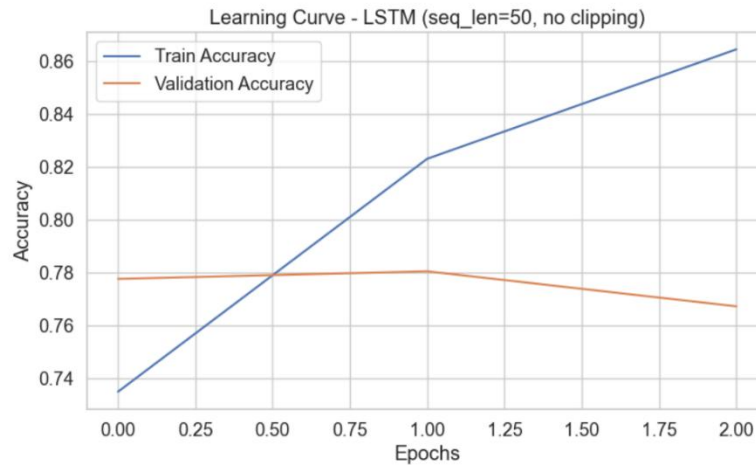
Figure 4. Learning Curve – LSTM (seq_len=50, no clipping)

- The training accuracy continues to rise while validation plateaus, suggesting mild overfitting after the second epoch.



Figure 5. Optimizer vs. Activation (Accuracy Heatmap)

- RMSProp and Adam with ReLU or tanh activations deliver strong accuracy (~0.77–0.78), while SGD lags significantly.
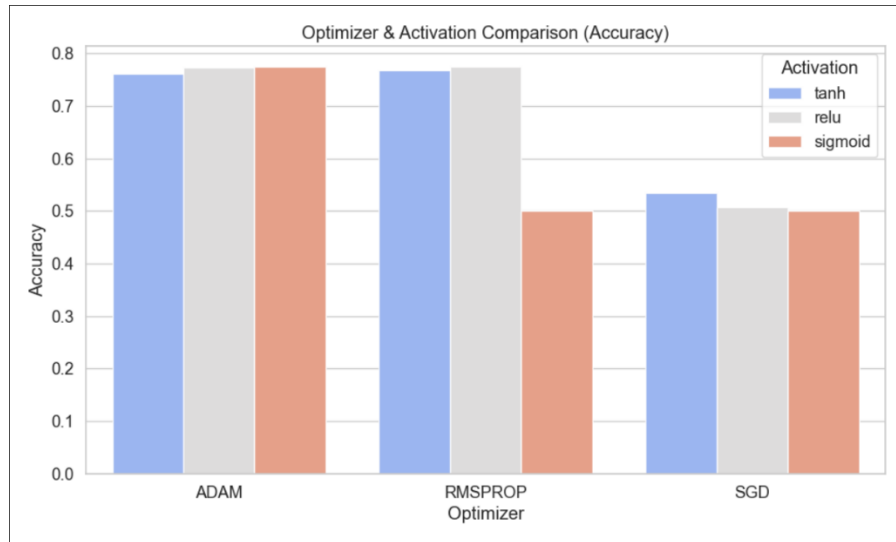
Figure 6. Optimizer & Activation Comparison (Accuracy)

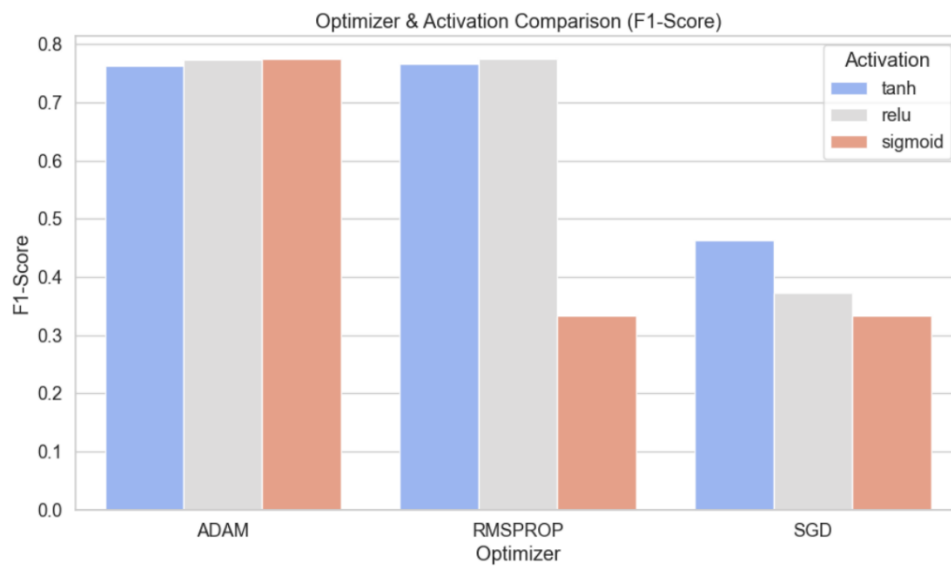- RMSProp + ReLU yields the top accuracy; SGD consistently underperforms.



Figure 7. Optimizer & Activation Comparison (F1-Score)

- Similar trend to accuracy, ReLU and tanh activations perform best with adaptive optimizers (Adam, RMSProp).
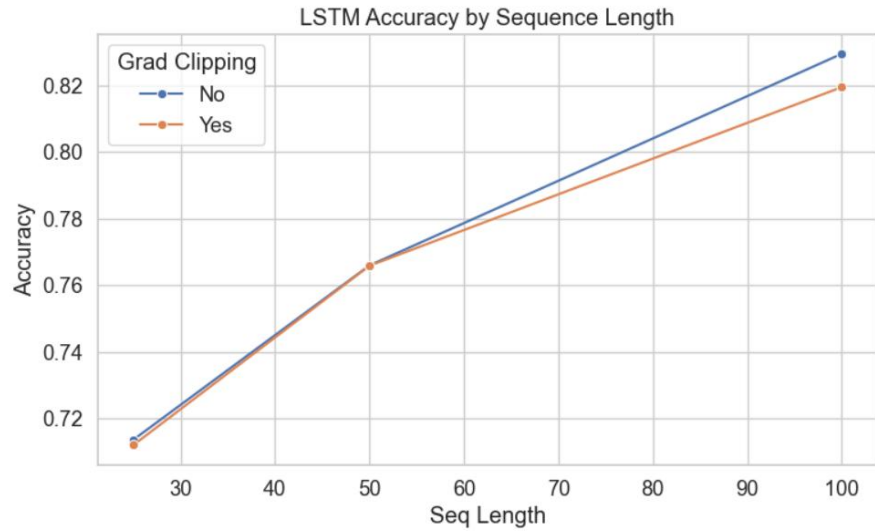
Figure 8. LSTM Accuracy by Sequence Length

- Performance of LSTM across sequence lengths (25, 50, 100).
  Observation: Accuracy rises steadily with longer input context; sequence length = 100 gives the best result.
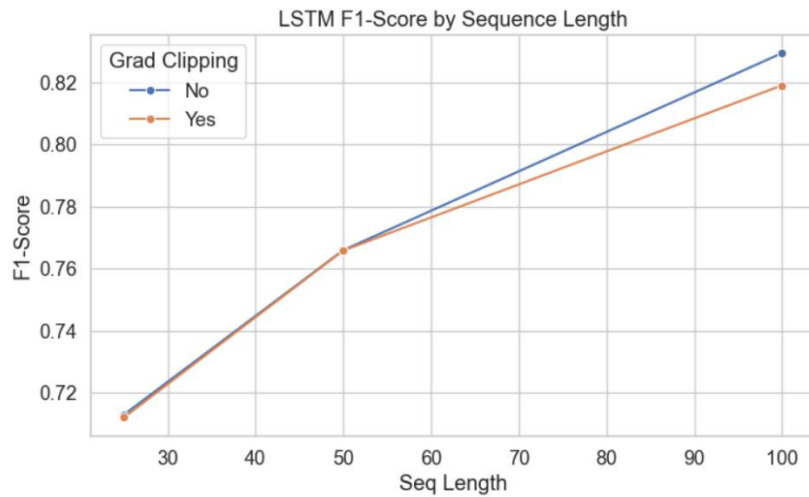


Figure 9. LSTM F1-Score by Sequence Length

- Parallel F1 trend with sequence length.
  Observation: Longer sequences improve recall and precision, confirming that richer context helps sentiment detection.
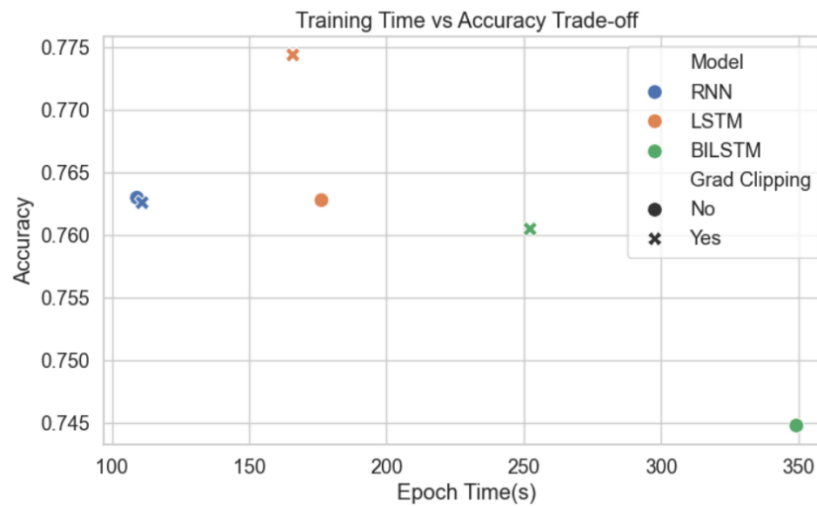
Figure 10. Training Time vs. Accuracy Trade-off

- Plots computation cost vs. performance for all architectures.
  Observation: RNN trains fastest but has slightly lower accuracy; LSTM achieves higher accuracy with moderate extra cost; BiLSTM is slowest with marginal gain. Overall, LSTM provides the best trade-off between time and accuracy.

## 4. Discussion

This section summarizes the comparative findings across architectures, sequence lengths, optimizers, and stability strategies, highlighting the configurations that achieved the best trade-off between accuracy, efficiency, and training stability.

The experiments revealed that LSTM models consistently outperformed RNN and BiLSTM in terms of both accuracy and F1-score. While BiLSTM offered marginally richer contextual understanding, it came at a significantly higher computational cost. RNNs, on the other hand, trained faster but lacked the representational depth needed for longer dependencies in text. Consequently, LSTM emerged as the most balanced architecture, providing strong performance at a moderate computational expense.

The sequence length analysis demonstrated a clear trend: increasing the number of input tokens improved both accuracy and F1-score. Models trained with 100-token sequences achieved the highest performance (Accuracy = 0.8295, F1 = 0.8295), suggesting that longer input contexts allow better sentiment capture. However, this improvement came at the cost of longer training times almost double compared to the 25-token configuration illustrating a trade-off between model depth and efficiency.

In examining optimizer and activation function combinations, adaptive optimizers such as Adam and RMSProp provided superior convergence stability and performance. In contrast, SGD lagged behind, showing difficulty in reaching competitive accuracy within limited epochs. Among activations, ReLU and tanh performed comparably well, especially when paired with RMSProp or Adam, while sigmoid activation struggled due to saturation effects. The heatmap visualization further reinforced these patterns.

Lastly, the gradient clipping strategy demonstrated its usefulness in enhancing training stability. While it did not dramatically increase accuracy, it consistently prevented gradient explosions particularly for LSTM and BiLSTM architectures and ensured smoother learning curves, especially for longer sequences.

## 5. Conclusion

Through systematic experimentation, this study identified the configuration that offers the best balance between accuracy, training stability, and computational efficiency under CPU constraints. The optimal model was the LSTM architecture using tanh activation, the Adam optimizer, and a sequence length of 100 tokens, achieving an accuracy of 0.8295 and an F1-score of 0.8295.

While BiLSTM provided similar accuracy, its significantly longer training time made it less practical for CPU-based environments. ReLU combined with RMSProp also performed competitively, representing a strong alternative when faster convergence is prioritized. Gradient clipping proved effective in stabilizing training, but its influence on final performance was minimal.

Overall, the LSTM (tanh + Adam, seq_len = 100) configuration is the most efficient and robust choice for sentiment classification on the IMDB dataset, providing a strong balance between model complexity, training time, and predictive accuracy.