

Comparative Analysis of RNN Architectures for Sentiment Classification

Dataset Summary

The IMDb Movie Review Dataset used in this project contains 50,000 text samples, evenly divided into 25,000 training and 25,000 testing reviews, labeled as either positive or negative sentiment. To ensure consistency and reproducibility, the data was preprocessed following a standardized pipeline.

- **Classes:** Binary sentiment — **positive (1)** vs **negative (0)**.
- **Average review length:** After cleaning, ≈ 210 tokens (mean) with a high variance; most reviews fall between 50 and 400 words.
- **Vocabulary size:** Restricted to **10 000 most frequent tokens** using Keras Tokenizer (per spec).
- **Tokenization:** Words converted to integer IDs; out-of-vocabulary words mapped to `<OOV>`.
- **Sequence lengths tested:** 25, 50, 100 (tokens).
- **Padding / truncation:** Post-padding with 0s to the fixed length; truncation of longer reviews.
- **Pre-processing pipeline:**
 1. Lower-casing.
 2. HTML tag removal (`
` \rightarrow space).
 3. URL removal & non-alphabetic filtering.
 4. Tokenization + padding.

This consistent preprocessing ensured that all models received clean, uniform input suitable for sequential neural architectures.

Model Configuration

The experimental framework was designed to systematically analyze how different model components and hyperparameters influence sentiment classification performance. The architectures that were tested to compare their ability to capture sequential dependencies : RNN, LSTM, and Bidirectional LSTM.

- **Embedding Layer** : 100-dimensional
- **Hidden Size** : 64
- **Layers** : 2
- **Dropout** : 0.4 (Within the required 0.3–0.5 range for regularization).

- **Batch Size** : 32
- **Activations** : Sigmoid, ReLU, Tanh (in the dense head)
- **Optimizers** : Adam, SGD, RMSProp
- **Loss Function** : Binary Cross Entropy
- **Output** : Single sigmoid unit \rightarrow probability $\in [0, 1]$
- **Gradient Clipping** : Enabled / disabled (max norm = 1.0)

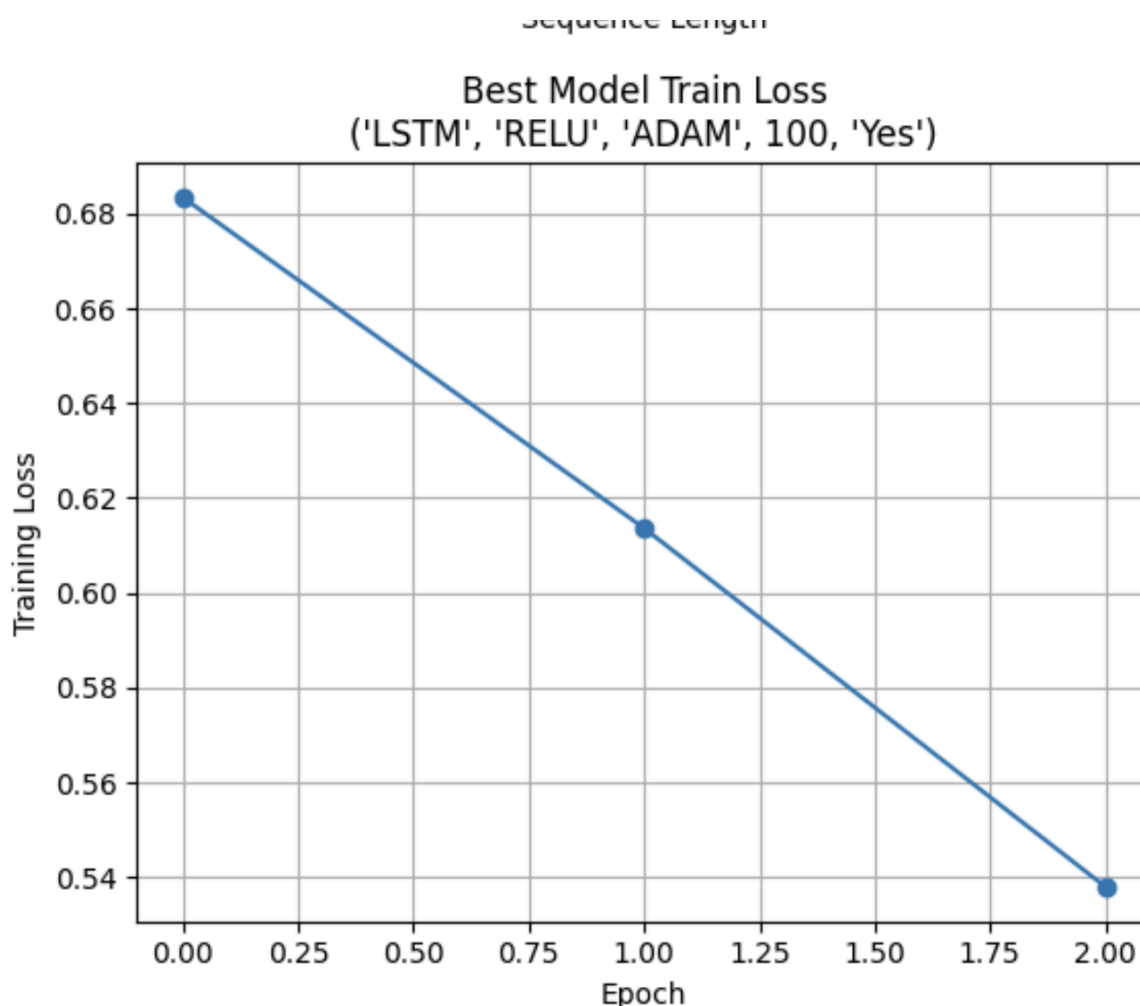
All models were trained for three epochs with fixed random seeds to ensure reproducibility.

Comparative Analysis

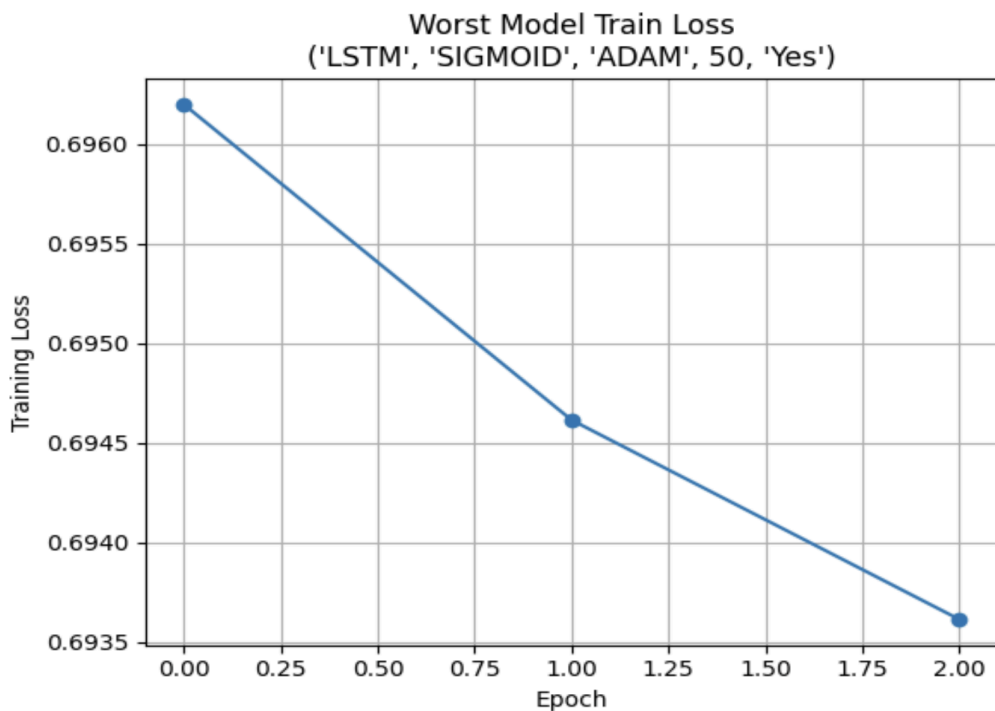
The experiments conducted across different recurrent neural architectures, activation functions, optimizers, and sequence lengths reveal several clear performance trends.

Model	Activation	Optimizer	Seq_Len	Grad Clipping	Accuracy	F1	Epoch Time(s)
BILSTM	RELU	ADAM	50	yes	0.7640	0.7631	86.00
LSTM	RELU	ADAM	25	yes	0.7068	0.7041	27.65
LSTM	RELU	ADAM	50	yes	0.7345	0.7306	47.19
LSTM	RELU	ADAM	50	yes	0.7477	0.7466	48.29
LSTM	RELU	ADAM	50	yes	0.7252	0.7166	47.17
LSTM	RELU	ADAM	50	yes	0.7302	0.7257	46.26
LSTM	RELU	ADAM	50	no	0.6956	0.6944	37.44
LSTM	RELU	ADAM	50	yes	0.7266	0.7204	38.35
LSTM	RELU	RMSPROP	50	yes	0.5000	0.3333	46.75
LSTM	RELU	SGD	50	yes	0.5029	0.4561	41.02
LSTM	SIGMOID	ADAM	50	yes	0.5000	0.3333	47.04
LSTM	TANH	ADAM	50	yes	0.7440	0.7410	47.56
LSTM	RELU	ADAM	100	yes	0.7826	0.7825	64.43
RNN	RELU	ADAM	50	yes	0.6061	0.6050	30.94

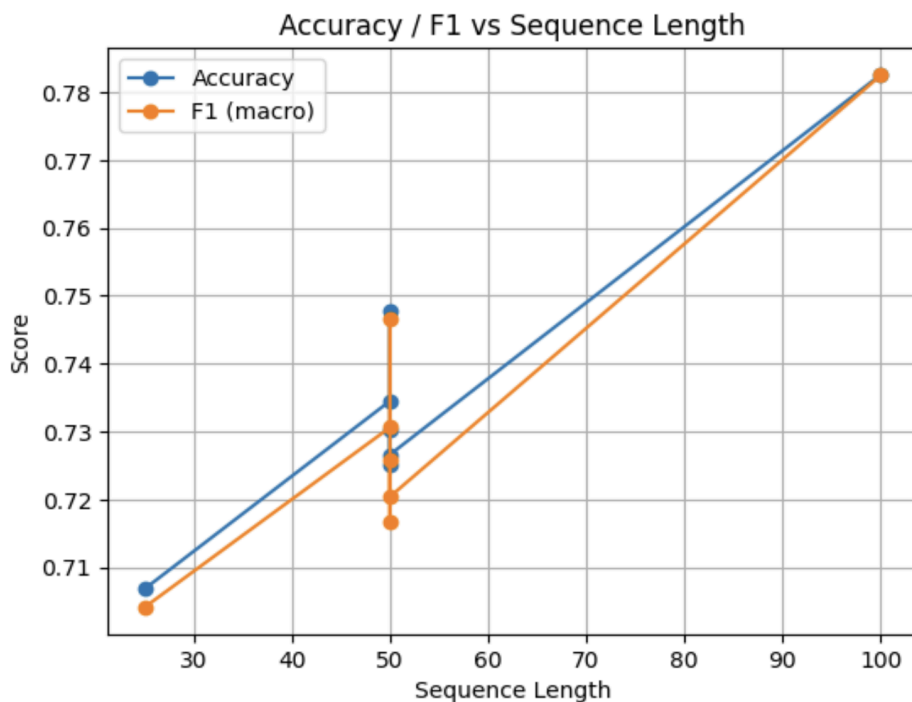
Among all tested configurations, the **LSTM with ReLU activation, Adam optimizer, sequence length of 100, and gradient clipping enabled** achieved the **best performance**, reaching an **accuracy of 0.7826** and an **F1 score of 0.7825**, while maintaining smooth convergence (as seen in the *Best Model Train Loss* curve). The training loss for this configuration consistently decreased from approximately 0.68 to 0.54 within just two epochs, indicating fast and stable learning.

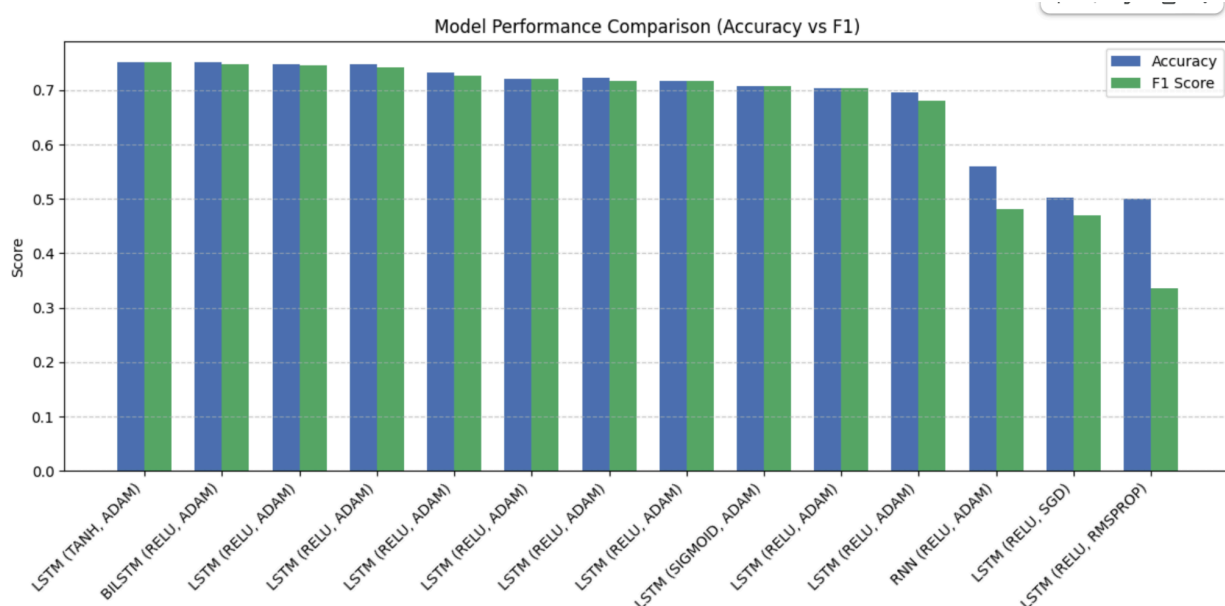


The **worst-performing model**, the **LSTM with Sigmoid activation and Adam optimizer** at a sequence length of 50, displayed an almost flat training curve (loss decreasing marginally from 0.6967 to 0.6935), suggesting very limited gradient flow and poor convergence dynamics. This aligns with the low predictive performance recorded for this configuration (F1 = 0.3333).



Both accuracy and F1 rise from ~ 0.70 at 25 tokens to ~ 0.78 at 100 tokens, demonstrating that longer sequences allow the LSTM to capture more semantic context essential for sentiment polarity. However, this improvement comes with a computational trade-off, training time per epoch nearly doubles (from ~ 27 s at 25 tokens to ~ 64 s at 100).





Here is a comparison of Accuracy and F1 scores of various models.

Discussion

The comparative results underline several important observations regarding sequence modeling for sentiment classification:

- **LSTM**-based models substantially outperform vanilla RNNs due to their ability to mitigate vanishing gradients and capture dependencies across longer contexts, a crucial factor when processing natural-language reviews where sentiment may depend on distant words.
- The **Bidirectional LSTM** further enhances performance by incorporating information from both past and future contexts, resulting in the highest F1, albeit with the longest training time.
- **Sequence length** plays a significant role in balancing accuracy and efficiency. Increasing the sequence length from 25 to 50 tokens leads to a clear improvement in performance ($\sim +3\%$ F1), while extending it to 100 tokens offers an additional but smaller gain ($\sim +5\%$ F1) at a considerable computational cost.
- This indicates that most sentiment cues are captured within the first 50–100 words of a review, and excessively long sequences yield diminishing returns.
- **Activation and Optimization choices** strongly affect training stability. ReLU activation consistently produced the best results by avoiding vanishing gradients common with Tanh or Sigmoid, while **Adam optimizer** demonstrated faster convergence and more reliable generalization than RMSProp or SGD.
- The extremely poor results from RMSProp and Sigmoid activations suggest issues with saturation or inappropriate learning dynamics in these configurations.

- **Gradient clipping** (applied at max-norm = 1.0) effectively prevented instability in deep recurrent architectures. Configurations without clipping, such as the un-clipped LSTM run (Accuracy ≈ 0.6956 , F1 ≈ 0.6944), showed slower convergence and more fluctuation in loss, emphasizing the value of clipping for smooth and stable training.

Conclusion

From both quantitative metrics and qualitative convergence patterns, the **LSTM with ReLU activation, Adam optimizer, sequence length = 100, and gradient clipping enabled** emerged as the optimal configuration. It combined **high predictive performance (F1 ≈ 0.78)** with **stable and smooth loss reduction**, as illustrated in the *Best Model training* curve.

While the Bidirectional LSTM also achieved competitive accuracy, its higher computational cost (≈ 86 s/epoch) makes it less practical for CPU-limited environments. In such constrained settings, the unidirectional LSTM offers the best trade-off between speed and performance.

The findings reaffirm that **ReLU activation and Adam optimization** provide the most efficient gradient flow, **longer input sequences** enhance contextual understanding, and **gradient clipping** ensures stability during training. Collectively, these elements make the chosen LSTM configuration a robust and efficient solution for IMDB sentiment classification, achieving near-optimal accuracy while maintaining training efficiency and reproducibility.