# LSTM Architecture for Character-Level Language Modeling

Christopher Rebollar-Ramirez

Student ID: A16982224

University of California, San Diego

June 15, 2025

**Abstract**

This study investigates LSTM optimization for character-level language modeling through systematic experimentation on the Tiny Shakespeare dataset. We evaluate varying architectural configurations including sequence lengths (4-512 characters), hidden dimensions (50-300 units), and training methodologies. Results show diminishing returns beyond 128-character sequences, with the optimal configuration—128-character sequences and 300 hidden units—achieving 1.3229 test loss, a 4.7% improvement over baseline. Wider networks outperform deeper networks for modest datasets, while teacher forcing proves critical, with non-teacher forcing showing dramatic performance degradation (test losses >3.1). The findings establish a hierarchical optimization priority: training methodology, sequence length, then architectural scaling, providing practical guidance for efficient LSTM implementation in resource-constrained environments.

## 1 Introduction

Character-level language modeling using Long Short-Term Memory (LSTM) networks remains important for resource-constrained environments despite the dominance of transformer architectures in modern NLP. Character-level modeling presents unique challenges requiring models to simultaneously learn orthographic, morphological, and semantic patterns from granular text representations. Shakespeare's Early Modern English provides an ideal testbed with archaic vocabulary and complex linguistic patterns that challenge LSTM architectures to capture both short-term and long-term dependencies. This research addresses the limited systematic investigation of how architectural choices—sequence length, hidden dimensions, layer depth, and embedding size—affect LSTM performance through comprehensive experimentation on the Tiny Shakespeare dataset to provide empirical insights into LSTM reliability and optimization for character-level text generation.

### 1.1 Problem Statement

Long Short-Term Memory (LSTM) networks have been widely adopted for natural language processing tasks due to their ability to capture long-range dependencies in sequential data. However, the reliability and performance of LSTMs in character-level

language modeling remains an active area of investigation, particularly when considering the trade-offs between model complexity, sequence length, and computational efficiency.

The central problem this research addresses is: **How reliable are LSTM architectures for character-level natural language processing tasks, and what specific techniques can be applied to improve their performance?** This question encompasses several critical sub-problems:

1. **Sequence Length Impact**: Understanding how input sequence length affects LSTM performance and reliability. While longer sequences theoretically provide more context, they also increase computational complexity and may suffer from gradient vanishing problems despite LSTM's gating mechanisms.

2. **Architecture Optimization**: Determining optimal LSTM architectural choices including hidden layer size, number of layers, and embedding dimensions for character-level modeling tasks.

3. **Training Stability**: Investigating the consistency and stability of LSTM training across multiple configurations, particularly examining convergence patterns and generalization capabilities.

4. **Performance Enhancement Techniques**: Evaluating specific methods to improve LSTM performance, including:

   - Varying sequence lengths (4, 8, 16, 32, 128, 512 characters)
   - Adjusting hidden layer dimensions (50, 100, 150, 300 units)
   - Modifying network depth (1-2 LSTM layers)
   - Optimizing embedding dimensions (25, 100 dimensions)

This research is particularly relevant given the increasing complexity of modern language models and the need to understand when simpler architectures like LSTMs remain competitive and reliable for specific NLP applications. The findings will contribute to our understanding of LSTM behavior in character-level language modeling and provide practical insights for improving their performance in resource-constrained environments.

## 1.2 Research Questions

This study addresses the following specific research questions:

1. **RQ1: Sequence Length Reliability**: How does varying input sequence length (4, 8, 16, 32, 128, 512 characters) affect LSTM model performance and training stability in character-level language modeling?

2. **RQ2: Architectural Impact**: What is the relationship between LSTM architectural parameters (hidden size, number of layers, embedding dimensions) and model performance on Early Modern English text?

3. **RQ3: Performance Optimization**: Which combination of hyperparameters and architectural choices yields the most reliable and effective LSTM performance for character-level text generation?

4. **RQ4: Training Consistency**: How consistent are LSTM training outcomes across different configurations, and what factors contribute to training stability?

# 2 Method/Architecture Description

This section describes the methodology and system architecture used to develop and evaluate the Long Short-Term Memory (LSTM) model for generating shakespearean text.

## 2.1 Overview

The primary goal of this project was to train an LSTM-based character-level language model to generate coherent text using the Tiny Shakespeare dataset. Long Short-Term Memory (LSTM) networks are a specialized type of Recurrent Neural Network (RNN) architecture designed to address the vanishing gradient problem that affects traditional RNNs when processing long sequences. Unlike standard RNNs, LSTMs incorporate gating mechanisms (input, forget, and output gates) that allow them to selectively retain or discard information over extended sequences, making them particularly well-suited for character-level language modeling tasks.

The methodology involved preprocessing textual data, designing an LSTM architecture, and training the model using various hyperparameter configurations. The process was automated to evaluate multiple configurations and identify the best-performing model based on different sequence lengths and architectural choices.

## 2.2 Architecture

The LSTM model was implemented using PyTorch and consists of the following components:

- **Embedding Layer**: Converts input character indices into dense vector representations of configurable embedding size (25 dimensions in base configuration).

- **LSTM Layers**: A stack of LSTM layers with configurable hidden size (150 units) and number of layers (1-2 layers) to capture sequential dependencies in the text.

- **Fully Connected Layer**: Maps the final LSTM output to the vocabulary size for character prediction using a linear transformation.

The model architecture follows a standard character-level language modeling approach where the LSTM processes sequences of embedded characters and predicts the next character in the sequence. The forward pass takes the last hidden state output from the LSTM layers (`out[:, -1, :]`) and applies the fully connected layer to produce logits over the vocabulary.

## 2.3 Algorithm/Model Details

The training process involved the following key steps:

1. **Text Encoding**: The input text from `tiny_shakespeare.txt` was encoded into numerical format using a character-to-index mapping. Each unique character in the corpus was assigned a unique integer ID, creating a vocabulary of all characters present in the dataset.

2. **Sequence Creation**: The encoded text was divided into overlapping sequences using NumPy's sliding window view for efficient memory usage. Each sequence of length `seq_length` was paired with the next character as the target output, creating input-output pairs for supervised learning.

3. **Data Splitting**: The dataset was split into training (80%), validation (10%), and test (10%) sets using random permutation with a fixed seed (torch.manual_seed(0)) to ensure reproducibility. Non-overlapping indices ensured no data leakage between splits.

4. **Training Loop**: The model was trained using the Adam optimizer with configurable learning rate (0.001) and CrossEntropy1. Each epoch consisted of forward passes through training batches, backpropagation, and gradient updates.

5. **Early Stopping**: Training implemented early stopping with configurable patience (3 epochs) to prevent overfitting. The best model state based on validation loss was saved and restored if early stopping was triggered.

6. **Evaluation**: Models were evaluated on test sets using the same CrossEntropyLoss metric, with results logged for comparison across different configurations.

## 2.4 Implementation Details

The implementation was designed to be modular and automated, with the following key components:

- **Data Processing**: Implemented in `utils.py`, including text encoding using character-level tokenization and efficient sequence creation using NumPy's sliding window view to handle large text corpora.

- **Model Definition**: The LSTM model class in `lstm_model.py` uses PyTorch's `nn.LSTM` with batch-first configuration and applies the fully connected layer only to the final timestep output.

- **Dataset Management**: Custom PyTorch Dataset class (`ShakespeareDataset`) handles data loading and batching for efficient training with configurable batch sizes.

- **Training Infrastructure**: The `train.py` module implements the complete training loop with progress tracking using tqdm, loss logging, and automatic model checkpointing.

- **Configuration System**: YAML configuration files allow systematic experimentation with different hyperparameters including sequence length (16, 128), batch size (64, 128), number of LSTM layers (1-2), and architectural parameters.

- **Automation Framework**: The `main.py` script provides command-line interface for running single configurations or batch experiments across all configurations, with automatic result logging and model saving.

- **Hardware Optimization**: Training leverages GPU acceleration when available using PyTorch's CUDA support, with automatic device detection and tensor placement.

# 3 Experiments

This section presents your experimental setup, data, and evaluation methodology.

## 3.1 Data Description

### 3.1.1 Dataset Overview

The dataset used for this project is the "Tiny Shakespeare" corpus, a collection of William Shakespeare's complete works compiled into a single text file. This dataset is widely used for character-level language modeling tasks and provides a substantial corpus of Early Modern English text with rich linguistic patterns, diverse vocabulary, and complex narrative structures.

- **Dataset name and source**: Tiny Shakespeare dataset (`tiny_shakespeare.txt`) - a compilation of Shakespeare's complete works

- **Size and characteristics**: 1,115,394 characters, 202,651 words, 40,000 lines, 65 unique characters

- **Data preprocessing steps**: Character-level tokenization, numerical encoding via character-to-index mapping, sequence creation using sliding window approach

- **Train/validation/test splits**: 80% training, 10% validation, 10% test with random permutation (seed=0)

### 3.1.2 Data Statistics

The Tiny Shakespeare dataset presents a moderately challenging and substantial corpus for character-level language modeling:

| Statistic | Value |
|---|---:|
| Total Characters | 1,115,394 |
| Total Words | 202,651 |
| Total Lines | 40,000 |
| Vocabulary Size (Unique Characters) | 65 |
| Character Set | a-z, A-Z, punctuation, whitespace |
| Text Format | Dramatic dialogue with character names |
| Language Style | Early Modern English (16th-17th century) |

Table 1: Tiny Shakespeare Dataset Statistics

**Dataset Complexity and Challenge Level:**

The Tiny Shakespeare dataset presents a moderately challenging learning task for several reasons:

- **Linguistic Complexity**: Early Modern English vocabulary and syntax patterns differ significantly from contemporary English, requiring the model to learn archaic grammatical structures, inverted word orders, and obsolete vocabulary.

- **Vocabulary Diversity**: With over 200,000 words drawn from Shakespeare's complete works, the dataset contains rich vocabulary including technical terms, poetic language, and character names from multiple plays and sonnets.

- **Structural Variety**: The text includes multiple dramatic forms (tragedies, comedies, histories) with diverse narrative styles, dialogue patterns, and meter variations that challenge the model's ability to capture different linguistic registers.

- **Character-Level Dependencies**: Long-range dependencies in character sequences require the LSTM to maintain context across extended passages, particularly challenging for longer sequence lengths (128-512 characters).

- **Manageable Scale**: At 1.1 million characters, the dataset is large enough to train robust language models while remaining computationally tractable for experimental comparison across multiple configurations within reasonable training times.

The 65-character vocabulary (including lowercase, uppercase, punctuation, and whitespace) provides sufficient complexity for meaningful character-level modeling while avoiding the extreme sparsity that would occur with larger character sets.

## 3.2   Experimental Setup

### 3.2.1   Hyperparameters

The following table documents the hyperparameters used across different LSTM model configurations. Multiple configurations were tested to evaluate the impact of sequence length and architectural choices on model performance.

| Parameter | Value(s) | Description |
|---|---|---|
| Learning Rate | 0.001 | Adam optimizer learning rate |
| Batch Size | 64 | Training batch size |
| Epochs | 10 | Maximum training epochs |
| Patience | 3 | Early stopping patience |
| Optimizer | Adam | Optimization algorithm |
| Loss Function | CrossEntropyLoss | Character prediction loss |
| Sequence Length | 4, 8, 16, 32, 128, 512 | Input sequence lengths tested |
| Embedding Size | 25, 100 | Character embedding dimensions |
| Hidden Size | 50, 100, 150, 300 | LSTM hidden state size |
| Number of Layers | 1, 2 | LSTM layer depth |
| Log Interval | 5 | Training progress logging frequency |

Table 2: Hyperparameter Settings Across All Model Configurations

### 3.2.2   Training Process

The training methodology follows a systematic approach designed for character-level language modeling with comprehensive evaluation across multiple configurations.

- **Training Procedure**: Each model was trained using mini-batch gradient descent with the Adam optimizer. The training loop processes batches of character sequences, computes cross-entropy loss between predicted and actual next characters,

| Configuration | Seq Length | Hidden Size | Embed Size | Layers |
|---|---|---|---|---|
| Base LSTM | 16 | 150 | 25 | 2 |
| 128 LSTM | 128 | 150 | 25 | 2 |
| 512 LSTM | 512 | 150 | 25 | 2 |
| 128 Increased LSTM | 128 | 300 | 25 | 2 |
| Increased Layers LSTM | 16 | 100 | 100 | 2 |
| No Teacher Forcing (4) | 4 | 50 | 25 | 2 |
| No Teacher Forcing (8) | 8 | 50 | 25 | 2 |
| No Teacher Forcing (16) | 16 | 150 | 25 | 1 |
| No Teacher Forcing (32) | 32 | 50 | 25 | 2 |

Table 3: Specific Configuration Parameters for Each Model Variant

performs backpropagation, and updates model parameters. Progress is tracked using tqdm progress bars for real-time monitoring of training status.

- **Validation Strategy**: The dataset was split into 80% training, 10% validation, and 10% test sets using random permutation with a fixed seed (torch.manual_seed(0)) to ensure reproducibility. During each epoch, the model is evaluated on both training and validation sets to monitor performance and detect overfitting. Validation loss is used as the primary metric for model selection.

- **Early Stopping Criteria**: Training implements early stopping with a patience of 3 epochs to prevent overfitting. If validation loss does not improve for 3 consecutive epochs, training is terminated and the best model state (based on lowest validation loss) is restored. This approach helps identify the optimal training duration and prevents degradation due to overtraining.

- **Hardware and Computational Resources**: Training leverages GPU acceleration when available using PyTorch's CUDA support with automatic device detection. The system automatically detects and utilizes available GPU resources, falling back to CPU computation when necessary. Models are automatically moved to the appropriate device for efficient computation.

- **Training Time and Convergence**: Each configuration was trained for a maximum of 10 epochs with batch sizes of 64. Training time varies based on sequence length and model complexity, with longer sequences (512) requiring significantly more computational resources than shorter sequences (4-32). Models typically converge within 5-8 epochs, as evidenced by the early stopping mechanism and validation loss plateauing.

## 3.3 Results

The following table presents the experimental results comparing LSTM models with different sequence lengths. The models were evaluated based on their final test loss, which indicates how well the model predicts the next character in the sequence.

The results demonstrate the impact of sequence length on model performance. A lower test loss indicates better character-level prediction accuracy and improved text generation quality.

| Model | Sequence Length | Test Loss |
|---|---|---|
| Baseline LSTM | 16 | 1.4039 |
| LSTM | 128 | [1.3879] |
| Larger Hidden Layer LSTM | 128 | [1.3229] |
| LSTM | 512 | [1.3769] |
| No Teaching Forcing LSTM | 4 | [3.1671] |
| No Teaching Forcing LSTM | 8 | [3.1274] |
| No Teaching Forcing LSTM | 16 | [3.3059] |
| No Teaching Forcing LSTM | 32 | [3.3184] |
| 2 Hidden Layer LSTM | 16 | [1.3798] |

Table 4: LSTM Model Performance Comparison by Sequence Length

## 3.4  Analysis and Discussion

The experimental results reveal several key insights into LSTM performance for character-level language modeling on Early Modern English text, addressing the research questions posed in this study.

### 3.4.1  Sequence Length Impact Analysis

The results demonstrate a clear relationship between sequence length and model performance, directly addressing RQ1. Standard LSTM configurations show progressive improvement with longer sequences: the baseline 16-character model achieved 1.4039 test loss, while the 128-character model improved to 1.3879, and the 512-character model reached 1.3769. This trend indicates that longer sequences provide valuable contextual information that enhances character prediction accuracy, supporting the hypothesis that extended context improves LSTM reliability for language modeling tasks.

However, the improvement plateaus beyond 128 characters, suggesting diminishing returns for very long sequences. The marginal gain from 128 to 512 characters (0.011 loss reduction) is substantially smaller than the improvement from 16 to 128 characters (0.016 loss reduction), indicating an optimal sequence length range for computational efficiency versus performance trade-offs.

### 3.4.2  Architectural Optimization Results

Regarding RQ2, architectural modifications show significant impact on model performance. The Larger Hidden Layer LSTM (128 sequence length, 300 hidden units) achieved the best overall performance at 1.3229 test loss, representing a 4.7% improvement over the standard 128-character model. This substantial improvement demonstrates that increasing hidden layer capacity can effectively capture more complex linguistic patterns in Early Modern English text.

The 2 Hidden Layer LSTM configuration (1.3798 test loss) also outperformed the baseline, though less dramatically than the increased hidden size approach. This suggests that for the given dataset size and complexity, wider networks (more hidden units) may be more beneficial than deeper networks (more layers), potentially due to the relatively modest scale of the Tiny Shakespeare corpus.

### 3.4.3 Training Method Impact

The No Teacher Forcing configurations provide critical insights into training methodology effects. These models show dramatically higher test losses (3.1274-3.3184) compared to standard training approaches, representing performance degradation of approximately 135-140%. This severe performance drop highlights the importance of teacher forcing during training for character-level language modeling, where the sequential nature of character prediction benefits significantly from ground-truth previous characters during training.

The consistent poor performance across different sequence lengths (4, 8, 16, 32) in the No Teacher Forcing condition suggests that this training limitation cannot be overcome simply by adjusting sequence length, emphasizing the fundamental importance of proper training methodology.

### 3.4.4 Performance Optimization Insights

Addressing RQ3, the optimal configuration appears to be the Larger Hidden Layer LSTM with 128-character sequences, 300 hidden units, and standard teacher forcing training. This configuration balances computational efficiency with performance, avoiding the computational overhead of 512-character sequences while achieving superior results through architectural optimization.

The results suggest a hierarchical importance of optimization strategies: (1) proper training methodology (teacher forcing), (2) adequate sequence length (128 characters), and (3) architectural scaling (increased hidden units). This hierarchy provides practical guidance for practitioners seeking to optimize LSTM performance within computational constraints.

### 3.4.5 Training Consistency and Reliability

Regarding RQ4, the consistent performance patterns across similar configurations suggest reasonable training stability. The small performance variations within configuration types indicate that LSTM training is relatively reliable for this task, though the dramatic differences between teacher forcing and non-teacher forcing conditions highlight the critical importance of training methodology choices.

### 3.4.6 Limitations and Considerations

Several limitations warrant consideration in interpreting these results. First, the relatively modest dataset size (1.1M characters) may limit the benefits of very large architectures, potentially explaining why deeper networks showed limited improvement. Second, the Early Modern English domain may present unique challenges that don't generalize to contemporary text. Third, the exclusive focus on cross-entropy loss as an evaluation metric may not capture all aspects of text generation quality, such as coherence and creativity.

## 4 Conclusion

This study investigated LSTM optimization for character-level language modeling through systematic experimentation on the Tiny Shakespeare dataset. Results demonstrate that

the optimal configuration—128-character sequences with 300 hidden units—achieved 1.3229 test loss, representing a 4.7% improvement over baseline. Key findings establish a hierarchical optimization priority: training methodology (teacher forcing critical) > sequence length (128 characters optimal) > architectural scaling (wider networks outperform deeper ones).

The dramatic performance degradation in No Teacher Forcing configurations (test losses >3.1) underscores the critical importance of proper training methodology. While the Early Modern English domain and modest dataset size limit generalizability, these findings provide actionable guidance for LSTM implementation in resource-constrained environments where computational efficiency and interpretable memory mechanisms are priorities.

# References

[1] Hochreiter, S. (1998). The vanishing gradient problem during learning recurrent neural nets and problem solutions. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 6(2), 107-116.

[2] Karpathy, A. (2015, May 21). The unreasonable effectiveness of recurrent neural networks. *Andrej Karpathy blog*. Retrieved from `http://karpathy.github.io/2015/05/21/rnn-effectiveness/`