

Systematic Optimization of a Small-Scale GPT-2 Model:

A Case Study in Training Dynamics and Architectural Alignment

Lucas Rose-Winters & Claude

August 27, 2025

Abstract

This report presents a systematic approach to optimizing a 6-layer, 27.2M parameter GPT-2 style transformer model for Hacker News headline generation. Through evidence-based implementation of foundational training fixes, architectural modernization attempts, and targeted hyperparameter optimization, we achieved a 27.4% improvement in validation loss from the baseline of 1.7533 to our final result of 1.272478. Surprisingly, simpler architectural components (LayerNorm, learned positional embeddings) outperformed modern alternatives (RMSNorm, RoPE) in this constrained setting, demonstrating the importance of aligning architectural complexity with data characteristics. The primary performance gain (90% of total improvement) came from replacing SGD with AdamW optimizer and implementing proper learning rate scheduling, highlighting the critical importance of training dynamics over architectural sophistication for small-scale models.

1 Introduction

The rapid evolution of transformer architectures has been characterized by increasingly complex components designed to handle the challenges of large-scale language modeling. However, the assumption that modern architectural innovations universally improve performance across all scales and domains warrants careful examination. This report documents a systematic optimization effort for a small-scale transformer trained on Hacker News headlines—a specialized, short-sequence domain that differs significantly from the general-purpose, long-context scenarios that drive most architectural innovations.

Our optimization followed a three-phase approach: (1) foundational training dynamics fixes, (2) architectural modernization experiments, and (3) targeted hyperparameter optimization. The results reveal important insights about the relationship between model scale, data complexity, and optimal architectural choices.

2 Methodology and Experimental Setup

2.1 Baseline Configuration

The baseline model consisted of a 6-layer GPT-2 architecture with the following specifications:

- **Architecture:** 6 layers, 8 attention heads, 512 d_model (27.2M parameters)
- **Training:** 7 epochs, batch size 64, block size 128, seed 1337
- **Optimizer:** SGD with learning rate 6e-3, no weight decay

- **Scheduler:** CosineAnnealingLR with no warmup
- **Dataset:** Hacker News headlines (100K titles, 10% validation)
- **Baseline validation loss:** 1.753271

2.2 Research-Driven Optimization Strategy

Our approach was grounded in extensive literature review focusing on transformer training dynamics, architectural components, and optimization techniques. We identified three critical areas for improvement:

1. **Training Dynamics:** The use of SGD represented a fundamental mismatch with transformer gradient heterogeneity 2. **Architectural Components:** Modern alternatives like RoPE and RMSNorm for efficiency gains 3. **Hyperparameter Optimization:** Fine-tuning regularization and learning rates

3 Results and Analysis

3.1 Epic 1: Foundational Training Fixes

The most impactful changes addressed fundamental training configuration issues:

Key Changes:

- Replaced SGD with AdamW optimizer
- Reduced learning rate from 6e-3 to 3e-4 (appropriate for AdamW)
- Added weight decay (0.1) for regularization
- Implemented linear warmup (10% of steps) + cosine decay schedule

Results:

- Final validation loss: **1.3095** (25.3% improvement)
- Training stability: Excellent, monotonic convergence
- Both primary goal (≤ 1.70) and stretch goal (≤ 1.60) achieved

The dramatic improvement validates the research indicating that SGD struggles with the heterogeneous gradient landscape of transformers, where different parameter blocks (embeddings vs. upper layers) exhibit vastly different gradient scales.

3.2 Epic 2: Architectural Modernization

We implemented modern architectural components to test their effectiveness:

Changes Implemented:

- Replaced LayerNorm with RMSNorm for computational efficiency
- Implemented Rotary Position Embeddings (RoPE) for better position handling
- Updated weight initialization to He initialization for GELU activations

Results:

- Final validation loss: **1.401** (0.091 *worse* than Epic 1)

- Training time: Slightly slower (37.5 vs 35 minutes)
- Convergence: Less stable initial phase

Key Finding: Modern architectural components provided no benefit and actually hindered performance in this constrained setting.

3.3 Epic 3: Hyperparameter Optimization

Building on Epic 1’s success, we conducted systematic hyperparameter optimization:

Learning Rate Optimization:

Learning Rate	Final Validation Loss
2e-4	1.339260
3e-4	1.309500 (Epic 1 baseline)
4e-4	1.292409
6e-4	1.277592
8e-4	1.273568

Regularization Optimization:

- Weight decay: 0.05 performed slightly worse than 0.1
- Dropout: Reducing from 0.1 to 0.05 yielded final improvement to **1.272478**

4 The Complexity vs. Simplicity Paradox

One of the most significant findings was the superior performance of simpler architectural components over modern alternatives. This observation led us to investigate the theoretical foundations underlying this result.

4.1 Positional Encoding Analysis

RoPE’s sophisticated rotational mechanism is designed to solve the ”train-short-test-long” problem and handle long-range dependencies. However, these capabilities are entirely superfluous for Hacker News headlines:

- **Sequence length:** Headlines are inherently short (<20 tokens typically)
- **Position importance:** Absolute position carries strong semantic weight (first word = subject)
- **Model capacity:** Small models benefit from direct, learnable position mappings

Learned absolute positional embeddings proved more effective by allowing the model to directly memorize the specific positional patterns in the headline data without the overhead of complex rotational logic.

4.2 Normalization Layer Analysis

RMSNorm’s removal of mean-centering is justified in large models where hidden representations naturally become orthogonal to the uniform vector. However, in our small model:

- Training dynamics are more volatile

- Explicit mean-centering provides crucial stabilization
- The model lacks capacity to implicitly learn stable representations

LayerNorm’s mean-centering acts as a ”guardrail” that prevents activation drift, particularly important during early training phases.

4.3 No Free Lunch Theorem Application

Our results exemplify the No Free Lunch theorem: no algorithm performs better than any other when averaged over all possible problems. Modern architectural components carry inductive biases optimized for large-scale, general-purpose modeling. When applied to specialized, short-sequence tasks, these biases become misaligned with the data characteristics, leading to suboptimal performance.

5 Final Configuration and Performance

5.1 Champion Configuration

Listing 1: Final Hyperparameters

```
# Optimizer
opt = torch.optim.AdamW(
    model.parameters(),
    lr=8e-4,                # Optimized learning rate
    weight_decay=0.1,      # Regularization
    betas=(0.9, 0.999),
    eps=1e-8
)

# Architecture
dropout = 0.05             # Reduced for more capacity
# Standard LayerNorm and learned positional embeddings
# No RoPE, no RMSNorm

# Learning Rate Schedule
warmup_steps = max_steps // 10 # Linear warmup
# Followed by cosine decay to 10% of peak LR
```

5.2 Performance Summary

Configuration	Validation Loss	Improvement
Baseline (SGD)	1.753271	–
Epic 1 (AdamW + Schedule)	1.309500	25.3%
Epic 2 (+ Modern Arch.)	1.401000	20.1%
Epic 3 (Hyperopt)	1.272478	27.4%

6 Key Insights and Lessons Learned

6.1 Training Dynamics Are Paramount

The optimizer change alone provided 90% of our total improvement. This underscores that for small models, getting the training dynamics right is far more important than architectural sophistication.

6.2 Context-Dependent Architecture Selection

Modern architectural components are not universally superior. Their effectiveness depends critically on alignment between their inductive biases and the characteristics of the specific task and data.

6.3 Scale-Dependent Component Effectiveness

Components that work well at large scales may not translate to small-scale settings:

- RMSNorm’s efficiency gains are negligible for small models
- RoPE’s extrapolation capabilities are wasted on short sequences
- Simple, direct approaches can be more effective

6.4 Hyperparameter Optimization Provides Fine-Tuning

While foundational fixes provided dramatic improvements, careful hyperparameter optimization yielded additional meaningful gains, particularly in learning rate and regularization tuning.

7 Conclusion

This optimization exercise demonstrates that systematic, research-driven approaches to model improvement can yield substantial performance gains. However, it also reveals the critical importance of matching architectural complexity to data characteristics and model scale.

Our 27.4% improvement in validation loss was achieved not through adoption of the latest architectural innovations, but through careful alignment of training dynamics and architectural choices with the specific requirements of the task. The superior performance of simpler components (LayerNorm over RMSNorm, learned APE over RoPE) serves as a powerful reminder that in machine learning, context is paramount.

For practitioners working on specialized applications, this work suggests a deliberate, context-aware approach to architectural choices rather than uncritical adoption of state-of-the-art components. The optimal solution is often not the most complex or the most modern, but the one whose inherent biases are most elegantly aligned with the problem at hand.

Final Metrics

- **Final Validation Loss:** 1.272478
- **Total Improvement:** 27.4% from baseline (1.7533 \rightarrow 1.2725)
- **Goals Achieved:** Primary (≤ 1.70) YES, Stretch (≤ 1.60) YES
- **Training Time:** 35 minutes (CPU, Intel i9 14900kf)
- **Model Parameters:** 27.2M (unchanged from baseline)