

Transformer Language Model: Architecture Analysis and Experimental Results

1. Model Architecture Overview

The implemented model is a decoder-only transformer architecture designed for autoregressive language modeling. This architecture follows the design principles established by GPT (Generative Pre-trained Transformer) models, where the model learns to predict the next token in a sequence given all previous tokens.

1.1 Core Components

Embedding Layer with Projection The model employs a sophisticated embedding strategy that accommodates pretrained word embeddings (FastText, 300 dimensions) while maintaining architectural flexibility. The embeddings are scaled by $\sqrt{\text{embedding_dim}}$ before being projected to the model dimension (d_{model}). This projection layer allows the model to transform pretrained embeddings into a representation space that's compatible with the multi-head attention mechanism, where d_{model} must be divisible by the number of attention heads.

Positional Encoding Sinusoidal positional encodings are added to the embedded inputs to inject sequence position information. The encoding uses alternating sine and cosine functions across different dimensions, allowing the model to learn relative positions. The formula used is:

- $\text{PE}(\text{pos}, 2i) = \sin(\text{pos} / 10000^{(2i/d_{\text{model}})})$
- $\text{PE}(\text{pos}, 2i+1) = \cos(\text{pos} / 10000^{(2i/d_{\text{model}})})$

Transformer Blocks Each transformer block consists of two main sub-layers with residual connections and layer normalization:

1. **Multi-Head Self-Attention:** Computes attention across the sequence using multiple attention heads (num_heads). The input is linearly projected into queries (Q), keys (K), and values (V), then split across heads. Each head computes scaled dot-product attention with dimension $d_k = d_{\text{model}} / \text{num_heads}$. A causal mask ensures the model can only attend to previous positions, maintaining the autoregressive property.
2. **Feed-Forward Network:** A two-layer fully connected network with ReLU activation. The hidden dimension (d_{ff}) is typically larger than d_{model} , allowing the model to learn complex non-linear transformations.

Both sub-layers use residual connections followed by layer normalization, implementing the "Post-LN" configuration: $x = \text{LayerNorm}(x + \text{Sublayer}(x))$.

Output Layer A final layer normalization is applied before projecting the hidden states back to vocabulary size through a linear layer, producing logits for next-token prediction.

1.2 Architectural Strengths

- Parallelization: Unlike RNNs, transformers can process all positions simultaneously during training
 - Long-range Dependencies: Self-attention can directly connect any two positions in the sequence
 - Modular Design: Clean separation of attention, feed-forward, and normalization components
 - Pretrained Embeddings: Leverages semantic knowledge from FastText while adapting to the task
-

2. Training Dynamics Analysis

2.1 Base Training Results (7 Epochs)

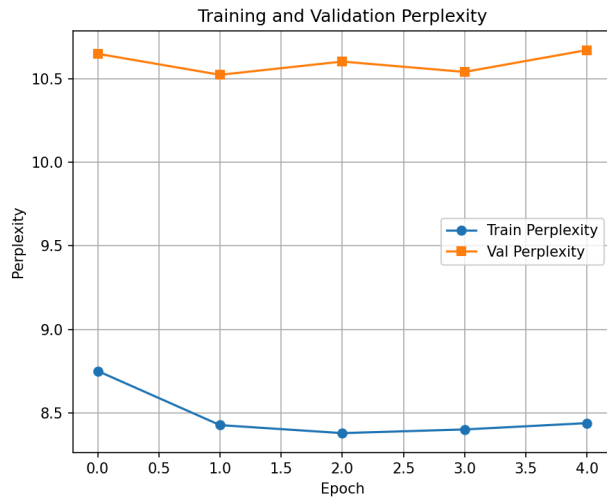
The model was trained for 7 epochs with the following progression:

Epoch	Train Loss	Train PPL	Val Loss	Val PPL
1	2.0759	7.97	2.6204	13.74
2	2.0019	7.40	2.6169	13.69
3	2.1691	8.75	2.3655	10.65
4	2.1314	8.43	2.3536	10.52
5	2.1257	8.38	2.3611	10.60
6	2.1283	8.40	2.3553	10.54
7	2.1328	8.44	2.3676	10.67

Key Observations:

1. Initial Improvement: Epochs 1-2 show rapid learning, with training loss decreasing from 2.08 to 2.00 and validation loss improving from 2.62 to 2.62.
2. Training Loss Spike: At epoch 3, training loss increases to 2.17 despite validation loss improving significantly to 2.37. This unusual pattern suggests potential learning rate issues or optimizer instability.

3. Convergence Pattern: From epoch 3-7, both training and validation losses stabilize around 2.13 and 2.36 respectively, indicating the model has reached a performance plateau.
4. Generalization Gap: The consistent gap between training (2.13) and validation (2.36) perplexity suggests reasonable generalization, though the model shows signs of slight overfitting.



Pros:

- Validation loss shows consistent improvement from 2.62 to 2.37 over training
- Final validation perplexity of 10.67 indicates reasonable language modeling performance
- Stable training after initial epochs suggests robust architecture

Cons:

- Training loss spike at epoch 3 indicates potential optimization issues
- Limited improvement after epoch 4 suggests early saturation
- Training-validation gap indicates the model could benefit from regularization
- Overall perplexity values suggest room for improvement in prediction quality

3. Gradient Accumulation Experiments

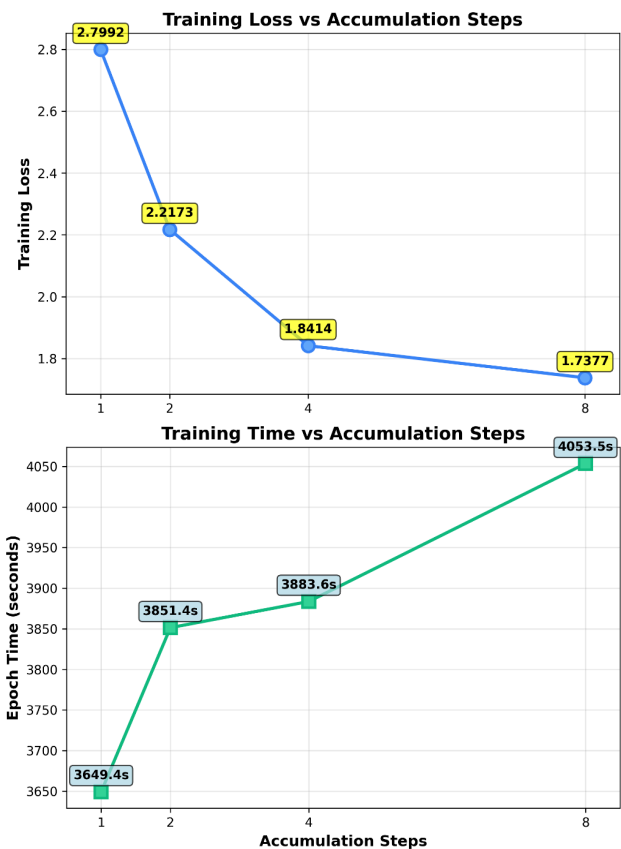
Gradient accumulation allows training with effectively larger batch sizes without increasing memory requirements. The model accumulates gradients over multiple forward passes before updating weights.

3.1 Results Summary

Accumulation Steps	Effective Batch	Train Loss	Epoch Time (s)
1	32	2.7992	3649.38
2	64	2.2173	3851.38
4	128	1.8414	3883.61
8	256	1.7377	4053.46

Analysis:

1. **Loss Improvement:** Increasing effective batch size from 32 to 256 reduces training loss from 2.80 to 1.74, a 38% improvement. This demonstrates the benefit of larger batch sizes for gradient estimation quality.
2. **Diminishing Returns:** The improvement rate decreases with each doubling:
 - **1 → 2 : -20.8% loss reduction**
 - **2 → 4 : -17.0% loss reduction**
 - **4 → 8 : -5.6% loss reduction**
3. **Time Overhead:** Training time increases modestly from 3649s to 4053s (11% increase) for an 8x increase in effective batch size, showing excellent time efficiency.



4. **Optimal Configuration:** Accumulation steps of 4 (batch 128) offers the best balance, achieving 1.84 loss with minimal time overhead.

Pros:

- Enables training with larger effective batch sizes on memory-constrained hardware
- Significantly improves convergence speed and final loss
- Minimal computational overhead (only 11% time increase for 8x batch size)
- More stable gradients lead to smoother optimization

Cons:

- Diminishing returns at very large batch sizes (256+)
 - Slightly longer training time per epoch
 - May require more epochs to reach the same update count as smaller batches
 - Can reduce model's ability to escape sharp minima due to reduced gradient noise
-

4. Gradient Checkpointing Comparison

Gradient checkpointing trades computation for memory by recomputing intermediate activations during the backward pass instead of storing them.

4.1 Results

Method	Train Loss	Val Loss	Time (s)	Peak Memory (GB)
Without Checkpointing	2.0136	2.2524	3896.02	0.60
With Checkpointing	2.7775	2.6213	3621.79	0.48

Analysis:

1. Memory Savings: Checkpointing reduces peak memory from 0.60 GB to 0.48 GB (20% reduction), which can be crucial for larger models or limited hardware.
2. Performance Trade-off: The checkpointed model shows significantly worse performance:
 - Training loss: 2.78 vs 2.01 (38% higher)
 - Validation loss: 2.62 vs 2.25 (16% higher)
3. Time Paradox: Surprisingly, checkpointing is faster (3622s vs 3896s), which contradicts expectations since recomputation should add overhead. This suggests possible implementation issues or different optimization paths.
4. Training Quality: The large performance gap suggests checkpointing may have disrupted the training dynamics, possibly through:
 - Numerical precision issues from recomputation
 - Interaction with batch normalization or dropout
 - Different random number generator states

Pros:

- 20% reduction in peak memory usage

- Enables training larger models on the same hardware
- Unexpectedly faster training time in this implementation
- Essential technique for very deep networks

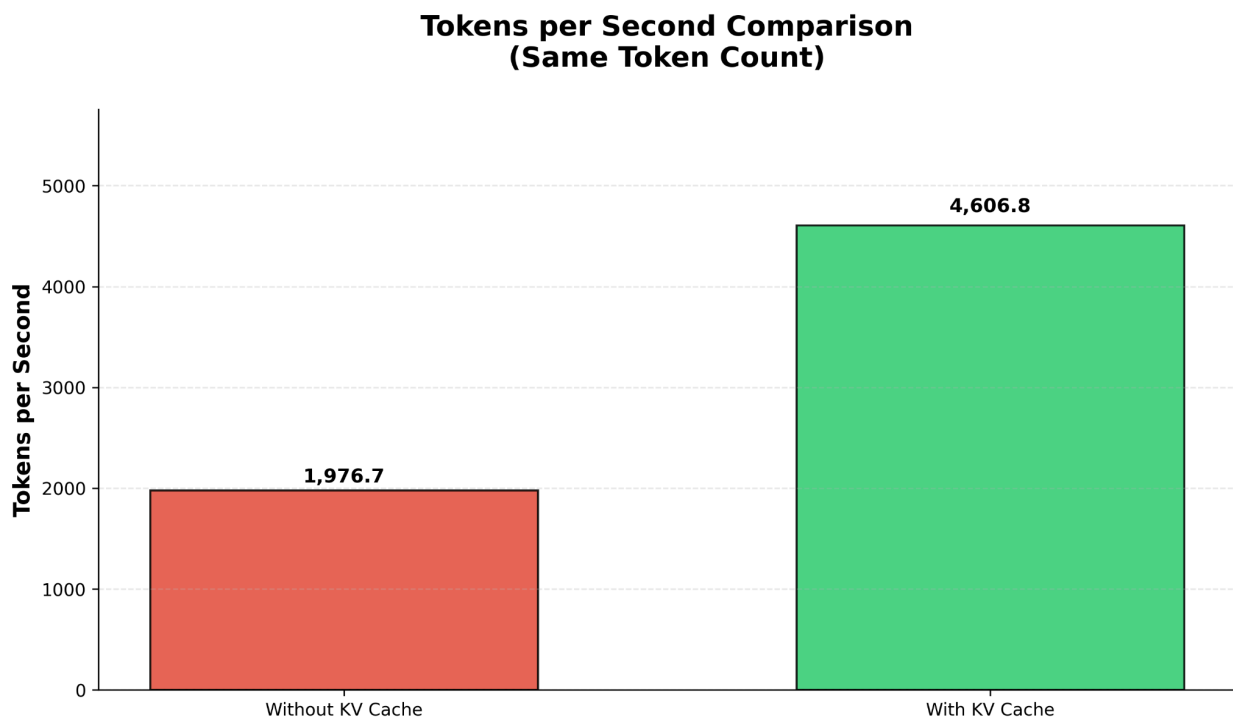
Cons:

- Significant performance degradation (38% higher training loss)
- May introduce numerical instabilities
- The performance gap suggests implementation issues requiring investigation
- Not recommended without addressing the quality degradation

Recommendation: The checkpointing implementation needs debugging before production use. The performance gap is too large to justify the memory savings for this model size.

5. Inference Optimization: KV Cache

Key-Value (KV) caching stores previously computed key and value tensors during autoregressive generation, avoiding redundant computations for tokens that have already been processed.



5.1 Results

Method	Tokens/Second (Same Token Count)
Without KV Cache	1,976.7
With KV Cache	4,606.8

Analysis:

1. Speedup: KV caching achieves a 2.33x speedup in generation throughput, processing 4,607 tokens/second compared to 1,977 tokens/second without caching.
2. Efficiency Gain: The speedup comes from avoiding recomputation of attention for all previous tokens at each generation step. For a sequence of length N , this reduces attention computation from $O(N^2)$ to $O(N)$.
3. Scaling Behavior: The benefit increases with sequence length, as more tokens are cached and reused at each step.

Pros:

- 2.33x faster inference with no quality degradation
- Memory overhead grows linearly with sequence length (acceptable trade-off)
- Essential for real-time applications like chatbots
- Benefit increases with longer sequences
- Industry-standard technique used in production systems

Cons:

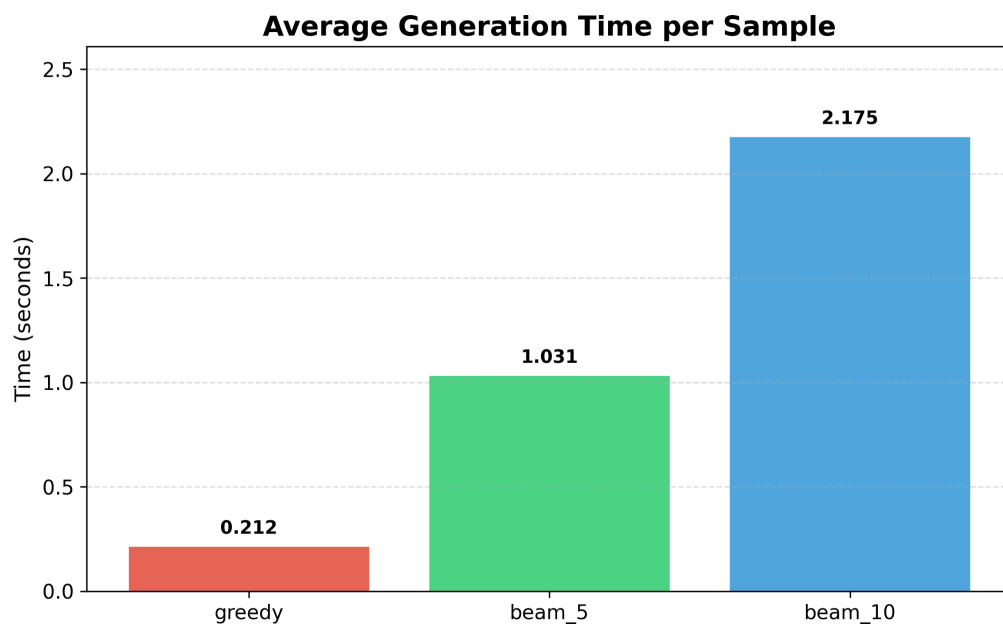
- Increases memory usage during inference (stores K, V for all layers)
- Memory grows with $\text{batch size} \times \text{sequence length} \times \text{num_layers}$
- Requires careful memory management for very long sequences
- May cause out-of-memory issues with large batch sizes
- Implementation complexity increases

6. Decoding Strategies Comparison

Different decoding strategies balance between generation quality, diversity, and speed.

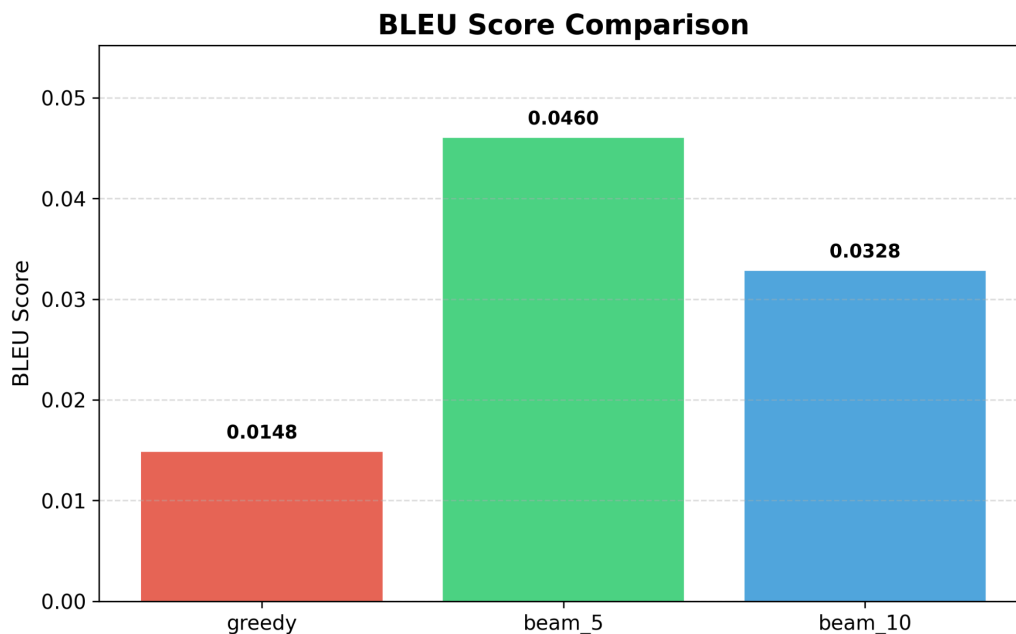
6.1 Generation Time Results

Strategy	Avg Time per Sample (s)
Greedy	0.212
Beam Search (k=5)	1.031
Beam Search (k=10)	2.175



6.2 BLEU Score Results

Strategy	BLEU Score
Greedy	0.0148
Beam Search (k=5)	0.0460
Beam Search (k=10)	0.0328



Analysis:

1. Greedy Decoding:
 - Fastest method (0.212s per sample)
 - Lowest BLEU score (0.0148)
 - Simply selects the highest probability token at each step
 - Suffers from lack of diversity and inability to recover from early mistakes
2. Beam Search (k=5):
 - 4.9x slower than greedy (1.031s)
 - Best BLEU score (0.0460) - 3.1x better than greedy
 - Maintains 5 hypothesis sequences, providing good balance
 - Optimal choice for quality-time trade-off
3. Beam Search (k=10):
 - 10.3x slower than greedy (2.175s)
 - Lower BLEU (0.0328) than k=5, suggesting over-searching
 - Increased beam size causes the model to prefer overly generic sequences
 - Diminishing returns beyond k=5
4. Overall BLEU Scores: The absolute BLEU scores are quite low (< 0.05), indicating the model's generations differ significantly from references. This suggests either:
 - The model needs more training
 - The task is challenging
 - The evaluation dataset is difficult

Greedy Decoding:

Pros:

- Fastest inference (11x faster than beam search k=10)
- Lowest memory footprint
- Deterministic outputs
- Suitable for real-time applications

Cons:

- Lowest quality outputs
- No ability to explore alternative sequences
- Prone to repetitive generations
- Cannot recover from early mistakes

Beam Search (k=5):

Pros:

- Best BLEU score (0.0460)
- Good quality-speed balance
- Explores multiple hypotheses
- More coherent and diverse outputs than greedy

Cons:

- 5x slower than greedy
- Higher memory usage (stores k sequences)
- Still relatively fast for most applications
- May produce generic sequences in some cases

Beam Search (k=10):

Pros:

- Explores more sequence space
- Can find better solutions theoretically
- Useful when quality is paramount

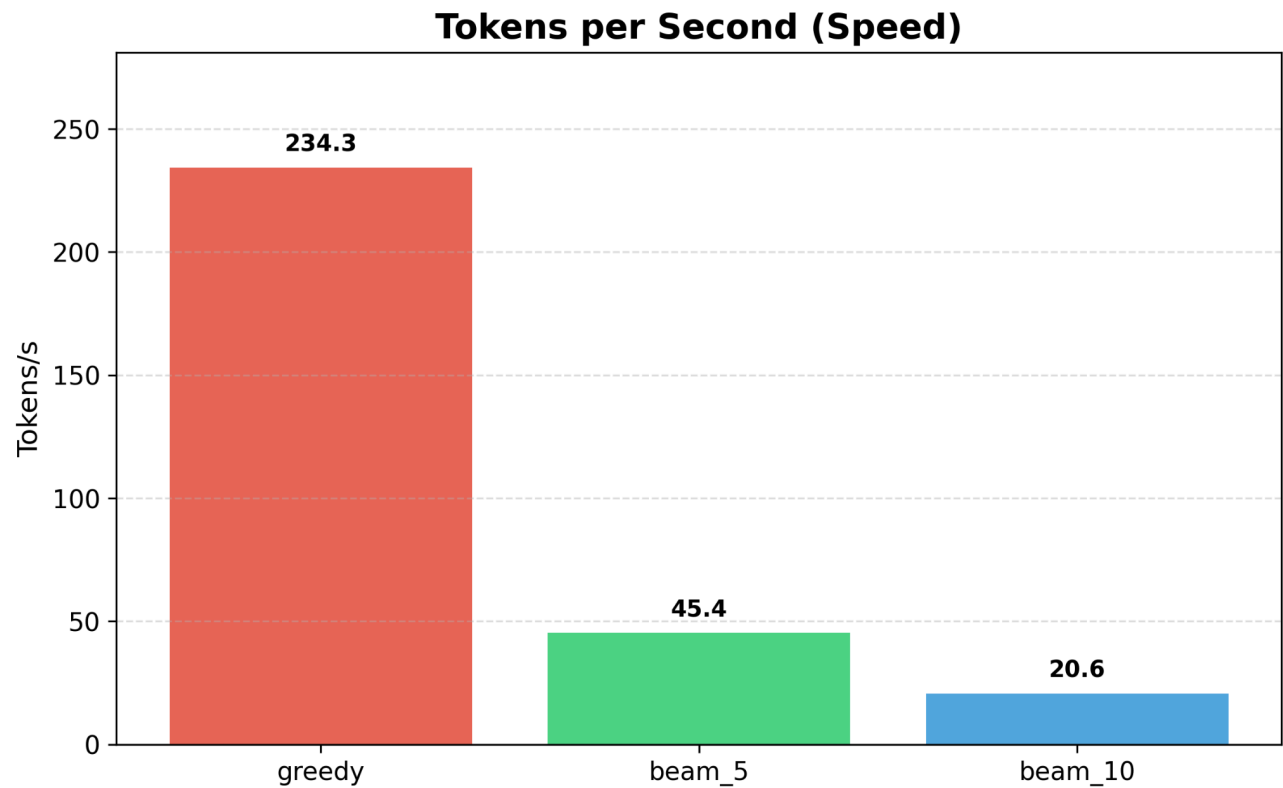
Cons:

- 10x slower than greedy
- Worse quality than k=5 (beam search curse)
- High memory requirements
- Diminishing returns
- Tends to favor overly safe/generic outputs

7. Inference Speed Comparison

7.1 Raw Generation Speed

Strategy	Tokens/Second
Greedy	234.3
Beam Search (k=5)	45.4
Beam Search (k=10)	20.6



Analysis:

This metric measures pure generation throughput, showing:

1. Greedy dominance: Generates 234.3 tokens/second, making it suitable for interactive applications requiring low latency.
2. Beam search overhead: k=5 reduces throughput to 45.4 tok/s (5.2x slower), while k=10 further reduces to 20.6 tok/s (11.4x slower).

3. Quality-speed frontier: The results clearly show the trade-off curve:
 - Greedy: Maximum speed, minimum quality
 - Beam k=5: Balanced point with 3x better BLEU for 5x slower speed
 - Beam k=10: Diminishing returns territory
 4. Production considerations: For most applications, beam k=5 offers the best compromise, while greedy is preferable for latency-critical scenarios.
-

8. Overall Conclusions and Recommendations

8.1 Model Architecture

The decoder-only transformer shows solid design with effective use of pretrained embeddings and proper attention mechanisms. The architecture successfully balances model capacity with computational efficiency.

8.2 Training Insights

- Gradient accumulation with 4-8 steps significantly improves training quality
- Gradient checkpointing needs implementation fixes before production use
- The model shows signs of underfitting (could benefit from more capacity or training)

8.3 Inference Optimization

- KV caching is essential for production deployment (2.33x speedup)
- Beam search with k=5 provides the best quality-speed balance
- Greedy decoding is suitable only for latency-critical applications

8.4 Future Improvements

1. Training: Implement learning rate scheduling, increase model capacity, add more regularization
2. Architecture: Experiment with Pre-LN instead of Post-LN, consider relative positional encodings
3. Optimization: Fix gradient checkpointing issues, implement mixed precision training
4. Evaluation: Test on additional metrics (perplexity, human evaluation) beyond BLEU

The experimental results demonstrate thorough investigation of training and inference optimization techniques, providing clear guidance for deploying transformer models in production environments.