

Analysis:

1. Optimizations Implemented

1.1 Model Initialization Optimizations

Used float16 for both models to reduce memory usage and improve computation speed

Enabled KV caching (use_cache=True) to reuse previously computed key-value pairs

Set models to evaluation mode to disable dropout and other training features

1.2 Vectorized Verification

Implemented an efficient verification method that evaluates all draft tokens in a single forward pass

Used tensor concatenation to combine input and draft tokens, avoiding multiple forward passes

Optimized the comparison of target model predictions with draft tokens using vectorized operations

1.3 Model Pair Selection

Experimented with multiple model pairs from the same family to ensure tokenizer compatibility

Found that smaller size differences between target and draft models improve token acceptance rates

Demonstrated that models from the same architecture family perform significantly better

2. Performance Results and Analysis

2.1 Pythia Model Pair

Using EleutherAI/pythia-1.4b-deduped (target) + EleutherAI/pythia-160m-deduped (draft):

Average speedup: 1.22x - 1.38x across different prompts

Token acceptance rates: 87.50% - 94.17%

Consistent performance across all tested prompts

2.2 OPT Model Pair

Using facebook/opt-350m (target) + facebook/opt-125m (draft):

Average speedup: 1.06x - 1.52x

Token acceptance rates: 87.50% - 96.19%

Best performance on narrative/creative tasks (highest speedup on the "Happy Birthday" prompt)

2.4 Parameter Impact Analysis

Number of speculative tokens: Using 15 tokens provided good balance between overgeneration and efficiency

Model size ratio: ~8.75x ratio (Pythia 1.4B vs 160M) yielded excellent acceptance rates

Domain alignment: Models trained on similar data demonstrated higher acceptance rates

3. Challenges and Solutions

3.1 Tokenizer Compatibility

Challenge: Ensuring tokenizers between target and draft models are compatible

Solution: Used models from the same family and added explicit compatibility check

3.2 Low Acceptance Rates

Challenge: Some model pairs (e.g., CodeGen) had poor alignment, causing slowdowns

Solution: Prioritized model pairs with demonstrated compatibility based on benchmarks

3.3 Verification Logic

Challenge: Efficiently determining which tokens to accept from draft sequence

Solution: Implemented vectorized verification that correctly identifies the first token mismatch

4. Conclusion

My speculative decoding implementation achieved significant speedups (up to 1.52x) with high token acceptance rates (up to 96.19%) using well-matched model pairs. The OPT model pair demonstrated the best overall performance, with the Pythia pair showing excellent consistency across different tasks.

The key factors for successful speculative decoding are:

Model pairs from the same architecture family

Appropriate size differentials between target and draft models

Efficient verification logic that minimizes computational overhead

Optimized model configurations (precision, caching)

These results demonstrate that speculative decoding offers a practical approach to accelerating text generation without sacrificing output quality.