

# Memory Optimization Techniques for Transformer Training: An Experimental Analysis

author: Nicolas Stupak

## 1. Introduction

Training large-scale transformer models is constrained by GPU memory limitations, which directly impacts the maximum batch size and, consequently, training efficiency and model convergence. This report presents various memory optimization techniques applied to transformer-based language modeling. It evaluates the trade-offs between memory consumption, training speed, and model quality across multiple optimization strategies.

## 2. Experimental Setup

### 2.1 Dataset Description

All experiments were conducted on a plwiki dataset filtered to only high-quality articles with the following characteristics:

- **Sequence Length:** 1024 tokens
- **Vocabulary Size:** 30,000 tokens (WordPiece tokenization)
- **Training Objective:** Next-token prediction with cross-entropy loss

### 2.2 Model Architecture

A transformer-based language model with the following specifications was used:

- **Model Parameters:** 24,785,384 parameters
- **Architecture Details:**
  - Number of layers: 12
  - Hidden dimension: 768
  - Number of attention heads: 12
  - Head dimension: 64
  - Feed-forward dimension: 3072
- **Baseline Configuration:** Standard transformer with full-precision (FP32) training

### 2.3 Hyperparameters

Consistent hyperparameters were maintained across all experiments:

- **Optimizer:** Adam
- **Learning Rate:** 5e-4 with linear warmup
- **Training Steps:** 200 iterations (10 epochs with ~20 steps per epoch)
- **Dropout:** 0.1
- **Gradient Clipping:** 1.0

## 2.4 Hardware Configuration

- **GPU:** NVIDIA A100 (40GB HBM)
- **CUDA Version:** 11.x
- **Framework:** PyTorch with custom CUDA kernels for optimized implementations

## 2.5 Optimization Techniques Evaluated

The following memory optimization techniques were evaluated:

1. **BF16 Mixed Precision (BF16-AMP):** Automatic mixed precision training using bfloat16
2. **Flash Attention:** IO-aware exact attention with tiling and recomputation
3. **Local Attention:** Sparse attention with sliding window (window sizes 128 and 256)
4. **Gradient Checkpointing (GradCP):** Trading computation for memory by recomputing activations
5. **Combined Techniques:** Flash Attention + Gradient Checkpointing

## 3. Results

### 3.1 Memory and Performance Comparison

Table 1 presents the results for all optimization techniques at their respective maximum batch sizes:

Technique	Batch Size	Peak Memory (MB)	Time per Step (s)	Total Time (s)	Perplexity
<b>Baseline (FP32)</b>	463	40,798	2.34	467.8	12.67
<b>BF16-AMP (same BS)</b>	-	-	-	-	-
<b>BF16-AMP (max BS)</b>	457	43,867	2.30	460.5	13.00
<b>FlashAttn (same BS)</b>	463	18,401	1.19	238.2	13.02
<b>FlashAttn (max BS)</b>	1,135	44,613	2.37	474.0	12.38
<b>LocalAttn-W128 (same BS)</b>	463	18,401	1.18	236.3	12.78
<b>LocalAttn-W128 (max BS)</b>	1,135	44,613	2.41	482.7	11.26
<b>LocalAttn-W256 (same BS)</b>	463	18,401	1.16	233.0	13.11
<b>LocalAttn-W256 (max BS)</b>	1,135	44,613	2.30	459.7	12.31

Technique	Batch Size	Peak Memory (MB)	Time per Step (s)	Total Time (s)	Perplexity
<b>GradCP (same BS)</b>	463	9,879	2.65	529.3	13.38
<b>GradCP (max BS)</b>	2,038	42,357	13.12	2,624.2	12.23
<b>FlashAttn+GradCP2,038 (max BS)</b>	2,038	41,444	5.62	1,123.0	12.17

### 3.2 Relative Performance Improvements

Table 2 shows the improvements (or degradations) relative to the baseline:

Technique	PPL Change	Speedup	Memory Reduction
BF16-AMP (same BS)	-	-	-
BF16-AMP (max BS)	+2.6%	1.02x	-7.5%
FlashAttn (same BS)	+2.7%	1.96x	+54.9%
FlashAttn (max BS)	-2.3%	0.99x	-9.4%
LocalAttn-W128 (same BS)	+0.8%	1.98x	+54.9%
LocalAttn-W128 (max BS)	-11.2%	0.97x	-9.4%
LocalAttn-W256 (same BS)	+3.5%	2.01x	+54.9%
LocalAttn-W256 (max BS)	-2.9%	1.02x	-9.4%
GradCP (same BS)	+5.6%	0.88x	+75.8%
GradCP (max BS)	-3.5%	0.18x	-3.8%
FlashAttn+GradCP (max BS)	-3.9%	0.42x	-1.6%

### 3.3 Memory Usage Analysis

Figure 1 (conceptual representation in table form) shows peak memory usage at the baseline batch size of 463:

Technique	Peak Allocated (MB)	Peak Reserved (MB)	Memory Efficiency
Baseline	2,944	44,098	6.7%
BF16-AMP	1,659	44,378	3.7%
FlashAttn	1,664	20,500	8.1%
LocalAttn-W128	1,664	20,500	8.1%
GradCP	1,676	12,800	13.1%

## 4. Analysis and Discussion

### 4.1 Memory Reduction Effectiveness

**Gradient Checkpointing** emerges as the most effective memory reduction technique at the same batch size, achieving 75.8% memory savings. This is because gradient checkpointing eliminates the need to store intermediate activations during the forward pass by recomputing them on-demand during the backward pass. The memory footprint becomes  $O(\sqrt{n})$  instead of  $O(n)$  for  $n$  layers.

**Flash Attention and Local Attention** techniques achieve substantial 54.9% memory reduction at the same batch size. These methods reduce memory by avoiding the materialization of the full  $N \times N$  attention matrix. Instead of computing and storing the entire attention matrix in high-bandwidth memory (HBM), they use tiling to process smaller blocks in fast on-chip SRAM. According to the Flash Attention paper, the IO complexity is reduced from  $\theta(Nd + N^2)$  to  $\theta(N^2d^2M^{-1})$ , where  $M$  is the SRAM size.

**BF16 Mixed Precision** The baseline batch size of 463 cannot fit in memory with BF16-AMP.

### 4.2 Training Speed Trade-offs

**Flash Attention and Local Attention** variants achieve approximately  $2x$  speedup at the same batch size. This speedup comes from reduced HBM accesses through kernel fusion and tiling. Despite performing additional FLOPs due to recomputation, the reduction in memory-bound operations dominates performance. The slightly better performance of LocalAttn-W256 ( $2.01x$  speedup) compared to LocalAttn-W128 ( $1.98x$  speedup) suggests that the window size affects the granularity of memory access patterns.

**Gradient Checkpointing** shows a slowdown ( $0.88x$  at same batch size,  $0.18x$  at max batch size) because it trades computation for memory. Every activation that was stored during the forward pass must now be recomputed during the backward pass, effectively doubling the forward pass computation. At maximum batch size (2,038), the computational overhead becomes severe (13.12s per step vs. 2.34s baseline).

**Combined FlashAttn+GradCP** achieves a middle ground, showing  $0.42x$  speed ( $2.4x$  slowdown) at maximum batch size. The Flash Attention component mitigates some of the gradient checkpointing overhead through efficient recomputation.

### 4.3 Model Quality and Perplexity

The perplexity results reveal important insights about the approximation quality of different techniques:

**Exact methods** (BF16-AMP, Flash Attention, Gradient Checkpointing) maintain relatively stable perplexity:

- BF16-AMP at max BS: 13.00 (+2.6%)
- FlashAttn at max BS: 12.38 (-2.3%)
- GradCP at max BS: 12.23 (-3.5%)
- FlashAttn+GradCP: 12.17 (-3.9%)

These techniques are mathematically equivalent to the baseline (or in the case of BF16, use a different numerical precision that doesn't significantly impact convergence). The variations in perplexity are primarily due to batch size differences affecting optimization dynamics rather than algorithmic approximations.

**Local Attention** shows varied results depending on window size:

- LocalAttn-W128 at max BS: 11.26 (-11.2% improvement)
- LocalAttn-W256 at max BS: 12.31 (-2.9%)

The surprising improvement with LocalAttn-W128 suggests that the local attention pattern may act as a beneficial inductive bias for this particular task, potentially reducing overfitting. The smaller window forces the model to focus on more local dependencies, which might align well with the linguistic structure of the data.

#### 4.4 Memory-Speed-Quality Pareto Frontier

Analyzing the three-way trade-off reveals several Pareto-optimal configurations:

1. **Best Speed** (same batch size): LocalAttn-W256 ( $2.01\times$  speedup) but with quality degradation
2. **Best Memory Efficiency**: GradCP at same BS (75.8% reduction) with acceptable quality (-5.6% PPL)
3. **Best Quality**: LocalAttn-W128 at max BS (11.26 PPL, -11.2%) with slight slowdown
4. **Best Balanced**: FlashAttn at same BS ( $1.96\times$  speedup, 54.9% memory reduction, +2.7% PPL)

#### 4.5 Combining Techniques: Synergies and Conflicts

The **FlashAttn+GradCP** combination demonstrates that:

1. **Memory benefits stack**: Achieves batch size 2,038 ( $4.4\times$  baseline), similar to GradCP alone
2. **Speed penalties partially mitigate**:  $0.42\times$  speed vs.  $0.18\times$  for GradCP alone ( $2.3\times$  faster than pure GradCP)
3. **Quality is preserved**: 12.17 PPL, best among all techniques

This suggests a multiplicative rather than additive relationship for speed, where Flash Attention's efficient recomputation significantly reduces gradient checkpointing's overhead.

#### Potential unexplored combinations:

- **BF16-AMP + FlashAttn**: Could provide both precision and attention optimization benefits
- **BF16-AMP + GradCP**: Maximum memory reduction for very large models
- **LocalAttn + GradCP**: Might achieve even larger batch sizes with potential quality improvements

#### 4.6 Practical Recommendations

Based on the findings:

1. **For memory-constrained scenarios**: Use Gradient Checkpointing (75.8% reduction) or combine FlashAttn+GradCP for better speed
2. **For speed-critical applications**: Use Flash Attention or Local Attention ( $2x$  speedup with minimal quality loss)
3. **For maximum model quality**: Use LocalAttn-W128 at maximum batch size (surprising 11.2% improvement)
4. **For production training**: FlashAttn at maximum batch size provides the best balance ( $2.45x$  batch size increase, similar speed, -2.3% better perplexity)

### 5. Conclusions

This experimental study demonstrates that memory optimization techniques for transformer training involve complex trade-offs between memory consumption, computational speed, and model quality. No single technique dominates across all metrics, and the optimal choice depends on the specific constraints and priorities of the training scenario.

Key findings include:

1. **Gradient checkpointing** provides the highest memory reduction but with significant speed penalties
2. **Flash Attention** offers an excellent balance of memory savings and speed improvements without quality degradation
3. **Local Attention** surprisingly improved model quality, suggesting architectural inductive biases can benefit from memory optimization constraints
4. **Combined techniques** (FlashAttn+GradCP) can achieve multiplicative benefits, enabling batch sizes  $4.4x$  larger than baseline