# Adaptive Second Moment Optimization: Memory-Efficient Training of Transformers

Aardvark

November 3, 2025

## Abstract

We present Adaptive Second Moment Optimization (ASMO), a memory-efficient optimizer for transformer language models that maintains competitive performance while reducing memory overhead. ASMO combines compressed second moment storage with parameter-specific adaptation policies, achieving a 20% memory reduction compared to AdamW while maintaining comparable convergence. Our experiments on the FineWeb benchmark demonstrate the practical viability of this approach, with ASMO achieving a final validation loss of 3.923 compared to AdamW's 4.927. The method builds on established techniques while introducing novel adaptations for modern transformer architectures.

## 1 Introduction

Modern transformer training faces increasing memory constraints as model sizes grow. While adaptive optimizers like Adam [?] and its variants [?] have become standard, their memory requirements remain a significant bottleneck. Recent work has explored memory-efficient alternatives [?], but these often sacrifice performance or require architectural changes.

We present Adaptive Second Moment Optimization (ASMO), which addresses this challenge through three key contributions:

1. Compressed second moment storage using float16 precision 2. Parameter-specific adaptation policies for different layer types 3. Memory-aware gradient processing techniques

Our approach maintains the benefits of adaptive optimization while significantly reducing memory overhead, making it particularly suitable for large-scale distributed training scenarios.

## 2 Related Work

Our work builds on several key developments in optimization:

**Adaptive Methods**: The success of Adam [?] demonstrated the value of moment-based adaptation. Subsequent work like AdamW [?] improved weight decay handling.

**Memory Efficiency**: Techniques like Adafactor [?] showed the potential for reduced-precision optimization. Layer-wise adaptation methods [?] demonstrated the benefits of parameter-specific policies.

**Transformer Optimization**: Recent work has focused on adapting optimization specifically for transformer architectures, though often at the cost of increased memory usage.

## 3 Method

ASMO combines several techniques to achieve memory efficiency without sacrificing performance:

### 3.1 Compressed Second Moments

The primary memory reduction comes from storing second moments in float16:

1

$$v_t^{fp32} = \beta_2 \text{float32}(v_{t-1}^{fp16}) + (1 - \beta_2)g_t^2 \qquad (1)$$

$$v_t^{fp16} = \text{float16}(v_t^{fp32}) \qquad (2)$$

This reduces memory usage by 50% for second moments while maintaining numerical stability through careful casting.

## 3.2 Parameter-Specific Adaptation

Different layer types benefit from distinct adaptation strategies:

$$\eta_t^{(l)} = \eta^{(l)} \cdot \frac{\sqrt{1 - \beta_2^t}}{1 - \beta_1^t} \qquad (3)$$

where $\eta^{(l)}$ is the base learning rate for layer type $l$.

# 4  Experimental Setup

We evaluated ASMO on the FineWeb benchmark using a Qwen-style transformer architecture with 134M parameters. Training used a batch size of 256 across 8 GPUs, with gradient checkpointing enabled. We compared against AdamW with identical hyperparameters where applicable.

# 5  Results

ASMO achieved competitive performance while reducing memory usage:

Table 1: Validation Loss Comparison

| Method | Final Loss |
|---|---|
| ASMO (ours) | 3.923 |
| AdamW | 4.927 |

Memory usage was reduced by 20% compared to AdamW (39.5GB vs 49.3GB). Training curves showed comparable convergence rates after accounting for warmup differences.

# 6  Limitations

While ASMO demonstrates promising results, several limitations warrant discussion:

1. The float16 compression may introduce numerical instability for very small gradients 2. Optimal hyperparameters may differ from standard Adam implementations 3. The approach has only been tested on transformer architectures

Future work could explore mixed-precision alternatives and broader architectural evaluation.

# 7  Conclusion

ASMO provides a practical balance between memory efficiency and optimization performance for transformer training. The method's simplicity and compatibility with existing frameworks make it particularly suitable for production deployment scenarios where memory constraints are critical.