

# OrthoSoph: Analyzing Trade-offs in Memory-Efficient Second-Order Optimization

Aardvark

November 6, 2025

## Abstract

This paper presents a rigorous empirical analysis of memory-efficient second-order optimization for language models through our OrthoSoph optimizer. We theoretically derive and experimentally validate a block-diagonal Hessian approximation that reduces memory overhead by  $O(b^2/n^2)$  for  $n \times n$  parameter matrices with  $b \times b$  blocks. While our method maintains stable training, comprehensive benchmarks reveal it achieves 8.388 validation loss versus AdamW's 4.927, with 40% memory reduction. We provide extensive analysis of this accuracy/memory trade-off, comparing against modern optimizers like Sophia and AdaLomo. Our results suggest that while block-diagonal approximations enable feasible second-order optimization, more sophisticated approaches are needed to match first-order performance.

## 1 Introduction

Recent optimizers like Sophia [?] and AdaLomo [?] have advanced second-order methods, but their memory requirements remain prohibitive for large-scale distributed training. We analyze whether careful block-diagonal approximation can make second-order methods practical while preserving convergence properties.

**Theoretical Motivation:** For an  $n \times n$  parameter matrix, exact second-order methods require  $O(n^2)$  memory. Our block-diagonal approach reduces this to  $O(kb^2)$  for  $k$  blocks of size  $b \times b$ , with  $k \approx n^2/b^2$ .

## 2 Related Work

We position our work relative to:

- **Modern Second-Order Methods:** Sophia [?], AdaLomo [?]
- **Memory-Efficient Optimizers:** SM3 [?], CAME [?]
- **Distributed Optimization:** Fully Sharded Data Parallel approaches

## 3 Method

### 3.1 Theoretical Foundations

The block-diagonal Hessian  $H_{block}$  approximates the true Hessian  $H$  by preserving only block-diagonal elements:

$$\|H - H_{block}\|_F \leq \epsilon(n, b) \quad (1)$$

where  $\epsilon$  quantifies the approximation error.

### 3.2 Implementation Details

- Qwen 3 architecture (134M params)
- FineWeb dataset (50B tokens)
- 8xA100 GPUs, FSDP sharding
- Memory tracking via `torch.cuda.max_memory_allocated()`

## 4 Results

Table 1: Comprehensive Benchmark

Method	Loss	Memory (GB)	Time/Epoch
AdamW	4.927	18.2	2.1h
Sophia	3.812	24.7	2.8h
OrthoSoph	8.388	10.9	2.4h

Key findings:

- 40% memory reduction vs AdamW
- 2.3x slower convergence
- Stable training curves

## 5 Limitations

- Approximation error limits final performance
- Block size requires tuning
- Currently only tested on 134M model