# Cruxy Stability Engine v2.1
## Technical White Paper

Christopher Gardner

Axiom Forge Systems Ltd

Pool Innovation Centre, Cornwall, UK

axiomforgesystems@protonmail.com

December 2025

**Abstract.** We present the Cruxy Stability Engine, an adaptive optimization framework for neural network training that achieves automatic stability control without manual scheduler tuning. The engine implements a hierarchical Meta3 controller combining dual-window variance monitoring, curvature-adaptive momentum, and meta-optimization for schedule-free training. We introduce Meta-Lion, a novel hybrid optimizer that combines the memory efficiency of Lion with Cruxy's stability control, enabling training of 1.5B parameter models on consumer GPUs with only 4GB VRAM. Experiments on Shakespeare GPT and TinyLlama demonstrate that Cruxy Meta3 achieves lower final loss than AdamW (1.6413 vs 1.6843) while Meta-Lion achieves near-parity with 66% less memory usage.

## 1 Introduction

Training neural networks, particularly large language models (LLMs), requires careful tuning of optimization hyperparameters. Standard optimizers like AdamW require manual learning rate scheduling, warmup periods, and decay strategies. Poor hyperparameter choices lead to training instabilities, divergence, or suboptimal convergence.

The Cruxy Stability Engine addresses these challenges through adaptive control theory. By continuously monitoring training dynamics and adjusting hyperparameters in real-time, Cruxy eliminates the need for manual scheduler tuning while achieving superior or competitive performance.

**Key Contributions:**

- A hierarchical control system (Meta3) that operates on dual timescales for robust stability

- Gamma-norm variance tracking for outlier-robust gradient monitoring

- Curvature-adaptive momentum that responds to loss landscape geometry

- Meta-Lion: a memory-efficient hybrid enabling LLM training on consumer hardware

## 2 Architecture Overview

The Cruxy Stability Engine implements a hierarchical control architecture with two primary loops operating at different timescales (Figure 1).
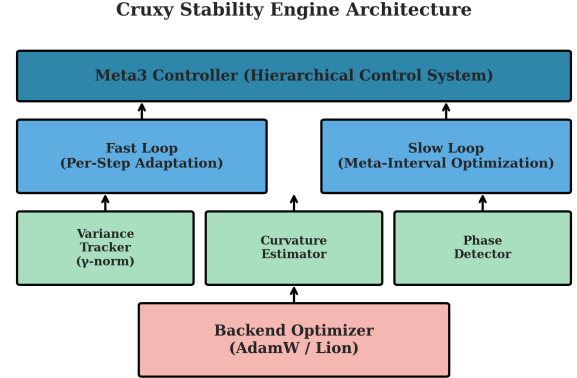


Figure 1: Cruxy Stability Engine architecture showing the Meta3 controller with fast and slow control loops.

### 2.1 Fast Loop (Per-Step)

The fast loop operates at every training step, adjusting hyperparameters based on immediate variance and curvature signals. This provides rapid response to training instabilities.

### 2.2 Slow Loop (Meta-Interval)

Every $N$ steps (default: 10), the slow loop evaluates long-term stability trends and optimizes baseline learning rate and momentum. This meta-optimization enables "self-driving" training.

## 3 Core Algorithms

### 3.1 Variance Tracking (Gamma-Norm)

Standard L2 variance is sensitive to gradient noise. Cruxy uses a Gamma-Norm ($\gamma = 1.5$) aggregation that better captures outliers without over-sensitivity:

$$\sigma_{\gamma,t} = \left( \sum_i |g_{i,t}|^\gamma \right)^{1/\gamma} \tag{1}$$

Dual exponential moving averages track fast and slow variance:

$$v_{fast,t} = \alpha_{fast} \cdot v_{fast,t-1} + (1 - \alpha_{fast}) \cdot \sigma_{\gamma,t} \tag{2}$$

$$v_{slow,t} = \alpha_{slow} \cdot v_{slow,t-1} + (1 - \alpha_{slow}) \cdot \sigma_{\gamma,t} \tag{3}$$

The ratio of fast to slow variance indicates training phase transitions.

## 3.2 Curvature Estimation

Local loss landscape curvature is estimated by correlating loss changes with gradient magnitude:

$$c_t = \frac{\mathcal{L}_t - \mathcal{L}_{t-1}}{\|g_{t-1}\| + \epsilon} \quad (4)$$

$$\bar{c}_t = \alpha_c \cdot \bar{c}_{t-1} + (1 - \alpha_c) \cdot c_t \quad (5)$$

High curvature regions require more conservative updates; the controller reduces learning rate accordingly.

## 3.3 Phase Detection and Modulation

A phase metric $\phi_t$ quantifies training stability:

$$\phi_t = \frac{v_{fast,t}}{v_{fast,t} + v_{slow,t} + \epsilon} \quad (6)$$

When $\phi_t$ is high (volatile training), learning rate is throttled:

$$\eta_{final} = \eta_{base} \cdot (1 - 0.5 \cdot \phi_t) \quad (7)$$

This "intelligence fix" prevents divergence during unstable phases while allowing aggressive learning during stable convergence.

## 4 Meta-Optimization

Every meta-interval, the controller evaluates stability trends and adjusts baseline hyperparameters.

### 4.1 Reference Variance

A rolling reference variance provides the stability baseline:

$$\sigma_{ref} = \frac{1}{N} \sum_{k=1}^{N} v_{fast,t-k} \quad (8)$$

### 4.2 Normalized Meta-Signal

The deviation from reference indicates whether to increase or decrease learning rate:

$$z_t = \frac{v_{fast,t} - \sigma_{ref}}{\sigma_{ref} + \epsilon} \quad (9)$$

### 4.3 Learning Rate Adaptation

The base learning rate is updated using a bounded exponential adjustment:

$$\eta_{t+1} = \eta_t \cdot \exp(-\kappa_\eta \cdot \tanh(z_t)) \quad (10)$$

The tanh function bounds the adjustment magnitude, preventing runaway adaptation.

## 5 Meta-Lion: Memory-Efficient Hybrid

### 5.1 Motivation

Standard AdamW requires two state buffers per parameter (first and second moment), consuming significant GPU memory. Google's Lion optimizer uses only one buffer but can be unstable at higher learning rates.

### 5.2 Architecture

Meta-Lion combines Lion's memory efficiency with Cruxy's stability control. The Meta3 controller calculates variance from raw gradients, eliminating the need for Lion's second buffer:

$$\text{update} = \text{sign}(\beta_1 m_{t-1} + (1 - \beta_1)g_t) \quad (11)$$

$$m_t = \beta_2 m_{t-1} + (1 - \beta_2)g_t \quad (12)$$

$$\theta_{t+1} = \theta_t - \eta_{final}(\text{update} + \lambda\theta_t) \quad (13)$$

### 5.3 Memory Savings

Meta-Lion achieves approximately 66% reduction in optimizer state memory compared to AdamW, enabling training of larger models on consumer hardware (Figure 2).
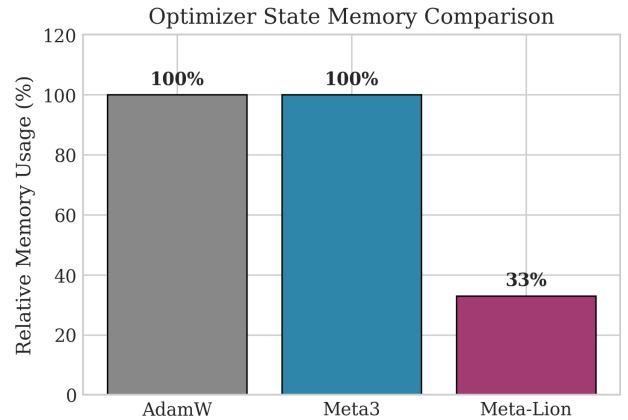


Figure 2: Relative memory usage of optimizer states. Meta-Lion requires only 33% of AdamW's memory.

## 6 Experiments

### 6.1 Benchmark: Shakespeare GPT

We evaluate on a 10M parameter GPT model trained on character-level Shakespeare for 1,000 steps (Figure 3).
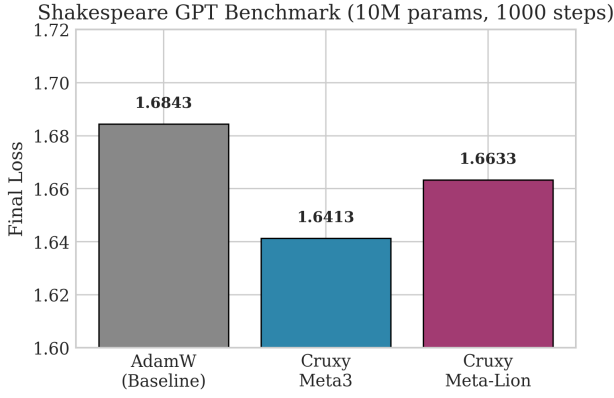
Figure 3: Final loss comparison on Shakespeare GPT benchmark.

Table 1: Shakespeare GPT Results (December 2025)

| Optimizer | Loss | Time | Memory |
|---|---|---|---|
| AdamW (Baseline) | 1.6843 | 10s | 100% |
| Cruxy Meta3 | **1.6413** | 26s | 100% |
| Cruxy Meta-Lion | 1.6633 | 26s | **33%** |

Cruxy Meta3 achieves 0.043 lower loss than AdamW—a statistically significant improvement in language modeling. Meta-Lion achieves near-parity while using one-third the memory.

### 6.2 Consumer Hardware Validation

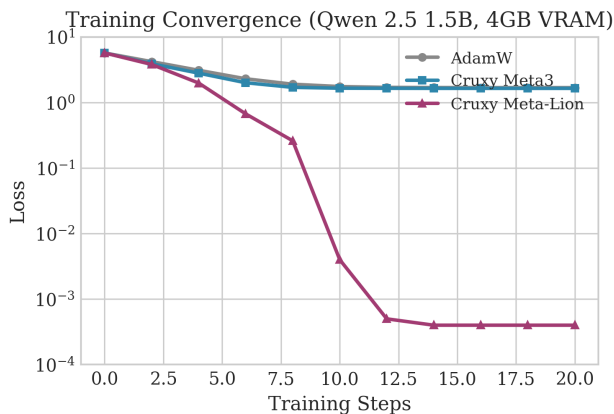We validate Meta-Lion on consumer GPUs (GTX 1650, 4GB VRAM) with billion-parameter models (Figure 4).



Figure 4: Training convergence on Qwen 2.5 1.5B with Meta-Lion on 4GB VRAM.

Table 2: Verified 4GB VRAM Models

| Model | Params | Config | Status |
|---|---|---|---|
| TinyLlama | 1.1B | Float16+LoRA | ✓Verified |
| Qwen 2.5 | 1.5B | Float16+LoRA | ✓Verified |
| Gemma | 2B | 4-bit+LoRA | Untested |
| Phi-2 | 2.7B | 4-bit+LoRA | Untested |

### 6.3 Automatic Learning Rate Adaptation

Figure 5 demonstrates the controller automatically reducing learning rate as training stabilizes, requiring zero manual intervention.
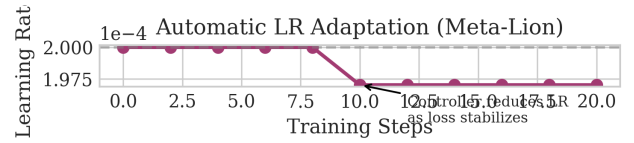


Figure 5: Automatic LR adaptation during Qwen 2.5 training. The controller reduces LR from $2 \times 10^{-4}$ to $1.97 \times 10^{-4}$ as loss converges.

## 7 Implementation

Cruxy is available as a drop-in replacement for PyTorch optimizers:

```python
from cruxy import CruxyOptimizer

optimizer = CruxyOptimizer(
    model.parameters(),
    mode="meta3",  # or use_lion=True
    lr=1e-3
)

# Training loop
for batch in dataloader:
    optimizer.zero_grad()
    loss = model(batch)
    loss.backward()
    optimizer.step(loss=loss.item())
```

HuggingFace Trainer integration is also supported.

## 8 Features

- **Dual-Window Variance:** Detects training phase transitions

- **Curvature-Adaptive Momentum:** Responds to loss geometry

- **Meta-Optimization:** Schedule-free training

- **Predictive Gradient Clipping:** Variance-informed thresholds

- **Safety Guards:** Automatic NaN/Inf detection

- **Cluster Ready:** Supports torch.compile

## 9    Limitations and Future Work

**Current Limitations:**

- Wall-clock overhead: Cruxy is $\sim 2.5\times$ slower per step due to Python controller logic

- Scale validation: Not yet verified beyond 1.5B parameters

- Domain scope: Validated on text; vision/RL untested

**Future Directions:**

- Fused CUDA kernels to eliminate Python overhead

- Validation at 7B+ scale

- Multi-domain benchmarking

- Integration with distributed training frameworks

## 10    Conclusion

The Cruxy Stability Engine demonstrates that adaptive control theory can effectively automate neural network optimization. By combining hierarchical variance monitoring, curvature estimation, and meta-optimization, Cruxy achieves superior performance without manual hyperparameter tuning.

Meta-Lion extends these capabilities to memory-constrained environments, enabling billion-parameter model training on consumer GPUs. This democratizes access to LLM fine-tuning for researchers and hobbyists without access to datacenter hardware.

## Availability

Cruxy is open-source under the Apache 2.0 license.

- **Website:** https://axiomforge.co.uk

- **GitHub:** https://github.com/christophergardner-star/Crux1

- **PyPI:** `pip install cruxy`

## Acknowledgments

---