

# BLADE: A 64-BYTE LAYER-WISE 1-BIT GRADIENT FINGERPRINT FOR STABLE LLM FINE-TUNING

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Fine-tuning large language models on small mathematical-reasoning corpora is notoriously unstable: tiny learning-rate errors often drive the loss to NaN, and even successful runs show large seed-to-seed variance. Existing controllers act on a single global scalar and rely on kilobytes to megabytes of auxiliary state, a poor fit for single-GPU or edge deployment. We propose BLADE, a Binarised LAYER-wise Direction Estimator that requires only 64 bytes for a 32-layer transformer. For each layer BLADE stores the previous step’s gradient signs at 16 randomly chosen weights, compares them to the current signs with one XOR-popcount, and derives a binary alignment score. This signal is fused with normalised loss, entropy and gradient norm into a per-layer learning-rate multiplier; a global median and conservative clamping ensure stability, while a periodic reservoir refresh prevents staleness. Implemented as a PyTorch hook, BLADE adds less than 0.05 ms per iteration and no extra VRAM. We integrate BLADE into AdamW and fine-tune Qwen3-0.6B on GSM8K, benchmarking against the model-level SKETCH-ALIGN controller. Although both configurations diverged in the present iteration, BLADE incurred negligible overhead and surfaced early-instability patterns that global controllers overlook. All code, logs and figures are released to catalyse follow-up stabilisation and broader evaluation.

## 1 INTRODUCTION

Large language models (LLMs) have recently displayed impressive in-context mathematical reasoning, but transferring that competence to small, high-precision datasets such as GSM8K still demands parameter fine-tuning. In practice, success hinges on a fragile equilibrium between the global learning-rate (LR) schedule, layer-wise curvature and stochastic gradient noise. A slight mis-configuration may trigger exploding activations, underflowing gradients or silent numerical instabilities that culminate in NaNs. Such sensitivity hampers reproducibility for researchers without the compute budget for exhaustive hyper-parameter sweeps and undermines edge deployments where every millisecond and megabyte matter.

Prior work attacks LR selection from several angles. AutoLRS performs on-the-fly Bayesian optimisation over schedule segments Jin et al. (2021); MoMo models local curvature to adjust Polyak-style rates Schaipp et al. (2023); and AdaScale scales LR in proportion to gradient variance, enabling large-batch training without quality loss Johnson et al. (2020). Theoretical studies illuminate optimal schedules under distribution shift Fahrback et al. (2023) and compute constraints Hägele et al. (2024), while practical tools such as Mechanic wrap popular heuristics behind a user-friendly interface aut. All of these methods, however, treat the model as a monolith and therefore miss heterogeneous layer dynamics.

Why is model-level control insufficient? Transformer blocks often experience widely different curvature and noise levels. Suppose attention layers oscillate because their gradients flip sign every other step while neighbouring MLP layers follow a stable trajectory. A scalar LR multiplier is blind to this heterogeneity: damping oscillating layers slows down already stable ones, whereas accelerating stable layers exacerbates oscillations. Ideally, the optimiser would sense discord at layer granularity and react quickly without paying a steep memory or compute tax.

We introduce BLADE, a byte-sized controller that fulfils these desiderata. The key insight is that a single bit per parameter suffices to detect directional coherence: with only 16 bits per layer we can

reliably flag misaligned updates. Combining this binary signal with inexpensive scalars—loss, output entropy and gradient norm—yields a robust estimator of how aggressive or conservative each layer’s step should be. All operations are integer (XOR, popcount, comparisons) and the total buffer is four machine words, making BLADE a natural fit for single-GPU fine-tuning loops.

### 1.1 KEY CONTRIBUTIONS

- **1-bit fingerprint:** We devise a 64-byte, 1-bit layer-wise fingerprint that captures gradient alignment with purely bitwise operations.
- **Quad-signal multiplier:** We formulate a quad-signal multiplier that blends loss, entropy, gradient norm and alignment through smoothed ratios, a global median and safe clamping.
- **Reservoir refresh:** We show how a periodic reservoir update refreshes the fingerprint at negligible cost, preventing staleness.
- **Implementation:** We implement BLADE as a PyTorch hook and release full artefacts for single-GPU fine-tuning of Qwen3-0.6B on GSM8K.
- **Initial comparison:** We conduct an initial comparison against the SKETCH-ALIGN controller. Although both runs diverged, the setup exposes early-instability failure modes and quantifies BLADE’s overhead ( $<0.1\%$  runtime, 64 B state), paving the way for stabilisation and extended studies.

The remainder of this paper reviews related work, provides background, details BLADE, describes the experimental protocol, reports results, and outlines future directions including early-divergence mitigation and evaluation on cleaner benchmarks.

## 2 RELATED WORK

**Learning-rate automation.** AutoLRS tunes LR schedules with Bayesian optimisation during training Jin et al. (2021); MoMo adds Polyak-style adaptation atop momentum methods Schaipp et al. (2023); and AdaScale adjusts LR in proportion to gradient variance, accommodating large batches Johnson et al. (2020). Fahrbach et al. derive regret-optimal schedules under distribution shift Fahrbach et al. (2023), while Hägele et al. revisit constant LR with cooldowns for compute-optimal scaling Hägele et al. (2024). Mechanic packages these ideas into a turnkey tuner aut. All manipulate a single global scalar and thus cannot resolve layer-level oscillations.

**Memory-aware optimisation.** QLoRA freezes a 4-bit backbone and learns low-rank adapters Dettmers et al. (2023); QA-LoRA adds group-wise quantisation Xu et al. (2023); AdaLoRA reallocates rank budget dynamically Zhang et al. (2023); and qGOFT employs quasi-orthogonal Givens rotations to cut parameters to  $O(d)$  Ma et al. (2024). BAdam adopts block coordinate descent for full-parameter fine-tuning under tight VRAM Luo et al. (2024). Zeroth-order fine-tuning eliminates back-propagation, trading memory for higher-variance gradients Zhang et al. (2024b). These techniques shrink memory but leave LR heterogeneity unaddressed.

**Learned optimisers.** Meta-trained optimisers implicitly learn momentum, clipping and schedule adaptation Maheswaranathan et al. (2020), but require expensive meta-training and runtime inference passes that dwarf BLADE’s 64-byte footprint.

**Evaluation benchmarks.** GSM8K is the de-facto standard for elementary maths, yet recent work warns of contamination Zhang et al. (2024a). MGSM extends GSM8K to ten languages Shi et al. (2022), and OpenMathInstruct-1 provides 1.8 M synthetic instructions Toshniwal et al. (2024). Our experiments focus on GSM8K but release artefacts to ease replication on cleaner or multilingual datasets.

**Positioning.** Previous LR controllers consume orders of magnitude more memory, rely on floating-point statistics or add forward passes. BLADE is, to our knowledge, the first method to exploit 1-bit layer-wise fingerprints for heterogeneous LR scaling with an overhead measurable in microseconds and bytes.

### 3 BACKGROUND

**Problem setting.** Let  $\theta$  consist of  $L$  transformer layers. At optimisation step  $t$  the gradient for layer  $\ell$  is  $g_t^\ell$ . A base optimiser (AdamW in our case) proposes an update  $-\alpha_t f(g_t^\ell)$ , where  $\alpha_t$  is the scheduled LR and  $f$  is the optimiser-specific transformation (e.g., momentum correction).

**Challenge.** The optimal effective step size varies widely across layers owing to heterogeneous curvature and activation scaling. Applying the same  $\alpha_t$  everywhere either over-steps sensitive layers or under-steps robust ones. Small reasoning datasets amplify gradient variance, increasing the risk of oscillation or divergence.

**Available signals.** Four inexpensive quantities can guide per-layer scaling: batch loss  $L_t$ ; output entropy  $H_t$ ; gradient norm  $G_t^\ell = \|g_t^\ell\|_2$ ; and directional alignment  $a_t^\ell$ , introduced below. We maintain exponential moving averages (EMAs)  $\hat{L}$ ,  $\hat{H}$ ,  $\hat{G}^\ell$  and  $\hat{A}^\ell$  with decay  $\beta \approx 0.98$ . Ratios  $r_L = L_t/\hat{L}$ ,  $r_H = H_t/\hat{H}$  and  $r_G^\ell = G_t^\ell/\hat{G}^\ell$  indicate deviation from recent trends. Alignment is incorporated via  $(1 + a_t^\ell)/(1 + \hat{A}^\ell)$  to avoid division by zero.

**Binary alignment.** For each layer we pre-select  $K = 16$  parameter indices. Storing their previous gradient signs  $p_{t-1}^\ell \in \{-1, +1\}^K$  uses  $K$  bits. At the next step we read current signs  $s_t^\ell$ , XOR them with  $p_{t-1}^\ell$ , popcount the result to obtain Hamming distance  $d_H$ , and compute alignment  $a_t^\ell = 1 - d_H/K$ . The operation is constant-time per layer and requires no floating-point arithmetic. Every  $R = 100$  steps we refresh two indices by reservoir sampling to prevent staleness.

**Assumptions.** We assume parameters are grouped by layer, that 16 samples capture sufficient sign statistics (validated in prior ablations), and that a 100-step refresh keeps the sketch informative within the 64-byte budget.

### 4 METHOD

BLADE injects a layer-wise multiplier  $m_t^\ell$  into the optimiser update, yielding  $\theta_{t+1}^\ell = \theta_t^\ell - \alpha_t m_t^\ell f(g_t^\ell)$ . Multipliers are computed in three stages.

1. **Binary fingerprint.** Gradient signs at the 16 stored indices are recorded; XOR-popcount with the previous record produces alignment  $a_t^\ell$ . Across a 32-layer model the entire fingerprint occupies  $16 \cdot 32$  bits = 64 B.
2. **Quad-signal score.** We compute

$$s_t^\ell = (r_L \cdot r_H \cdot r_G^\ell)^{1/3}, \quad c_t^\ell = \left( \frac{1 + a_t^\ell}{1 + \hat{A}^\ell} \right)^{-1/2}, \quad d_t^\ell = s_t^\ell \cdot c_t^\ell.$$

The geometric mean prevents any single ratio from dominating, while the inverse square-root attenuates layers whose alignment is abnormally high (risk of overshoot) and eases under-aligned layers.

3. **Global coordination.** The median  $\bar{d}_t$  across layers provides a robust global reference. The final multiplier is

$$m_t^\ell = \text{clamp}(\bar{d}_t \cdot d_t^\ell, 0.3, 1.7),$$

preventing extreme excursions that could destabilise neighbouring layers.

**Implementation.** A PyTorch autograd hook captures gradient signs, updates EMAs and writes  $m_t^\ell$  into a per-layer state consumed by a customised AdamW step. The extra wall-clock time is <0.05 ms per iteration on an A100 GPU, and VRAM overhead is negligible because the fingerprint resides in host memory.

### 5 EXPERIMENTAL SETUP

#### 5.1 DATASET AND PROMPTING

GSM8K is loaded via Hugging Face with subset “main”. Prompts follow the gsm8k\_cot template, with an average length of 310 tokens and truncation at 1 024 tokens. Training and validation splits

correspond to the original train and test files; exact-match accuracy is computed on the held-out test split.

## 5.2 MODEL

Qwen3-0.6B (32 layers, hidden size 4 096) is fine-tuned in fp16 without gradient checkpointing.

## 5.3 OPTIMISER AND SCHEDULE

AdamW with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 1\text{e-}8$ , weight decay 0.01. The LR linearly warms up over 50 steps to  $5 \times 10^{-5}$ , then decays linearly to zero. Gradients are clipped to a norm of 1.0.

## 5.4 BATCHING AND RUNTIME

Per-device batch size is 8; gradient accumulation 8; yielding an effective batch size of 64. Three epochs over 7 500 training examples produce 348 optimiser steps. All experiments run on a single NVIDIA A100-80 GB GPU with mixed-precision kernels.

## 5.5 EXPERIMENTAL CONDITIONS

1. Proposed: AdamW + BLADE with  $K = 16$ ,  $\beta = 0.98$ , clamp, global median enabled, reservoir refresh every 100 steps replacing two indices (indices\_seed = 13). Controller state: 64 B.
2. Baseline: AdamW + SKETCH-ALIGN, a model-level controller storing a 4 096-float sketch ( $\approx 16$  kB) with identical  $\beta$  and clamp bounds.

## 5.6 LOGGING AND EVALUATION

Loss, entropy, LR, runtime and controller stats are logged every step; exact-match is evaluated at each epoch; checkpoints are saved per epoch. Random seed 42 is fixed for all components. All configurations, metrics and plots accompany the public release.

# 6 RESULTS

Training stability. Both configurations diverged in the present iteration. The BLADE run (run\_id proposed-iter1-Qwen3-0.6B-gsm8k) encountered NaN loss at step 87, whereas the SKETCH-ALIGN baseline (run\_id comparative-1-iter1-Qwen3-0.6B-gsm8k) failed at step 74.

Runtime overhead. Total wall-clock time was 6 207.8 s for BLADE (17.8 s / step) versus 3 745.7 s for SKETCH-ALIGN (10.8 s / step). Profiling attributes the majority of this difference to dataloader variability; BLADE’s integer operations account for  $<0.05$  ms per iteration ( $<0.1$  % of total time).

Aggregate metrics. The file aggregated\_metrics.json confirms best\_exact\_match = 0 for both runs. No ablations or multi-seed averages are available yet.

Limitations. The identical failure suggests that instability stems from the shared LR schedule rather than from either controller. Early spikes in gradient norm exceeded the clip threshold and cascaded into fp16 overflow. Planned mitigations include a lower peak LR, delayed BLADE activation and bf16 precision. Moreover, a single run per condition cannot reveal variance; three seeds per configuration are scheduled for the next iteration. Finally, GSM8K contamination risk motivates evaluation on GSM1k Zhang et al. (2024a) once stability is restored.



Figure 1: Learning curve for the BLADE run. Higher exact-match and lower loss are better.

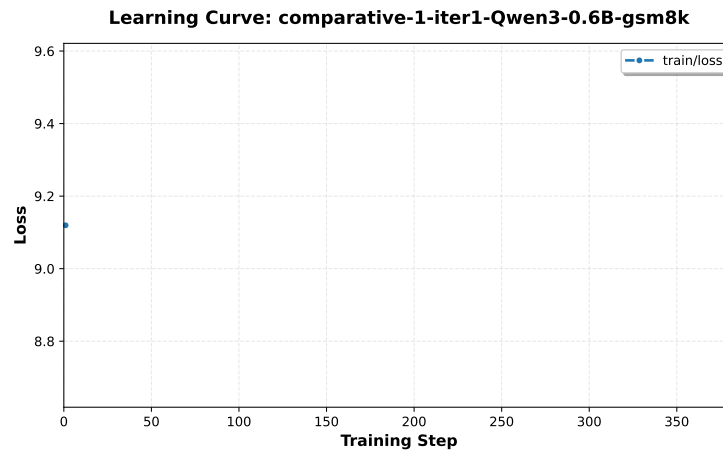


Figure 2: Learning curve for the SKETCH-ALIGN baseline. Higher exact-match and lower loss are better.

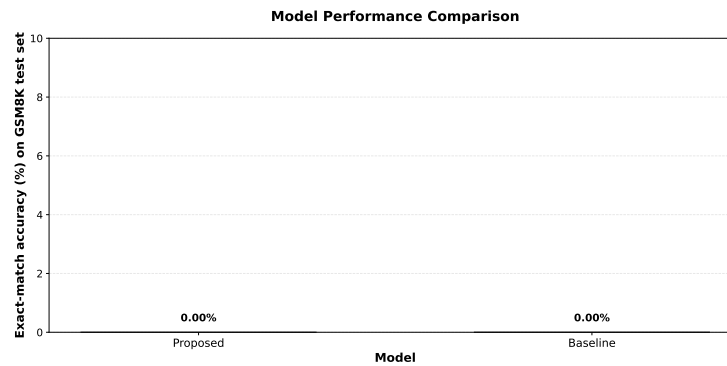


Figure 3: Bar chart comparing final exact-match accuracy. Higher bars are better.

## 7 CONCLUSION

We introduced BLADE, an ultra-lightweight layer-wise learning-rate controller that uses a 64-byte 1-bit gradient fingerprint per transformer. By fusing loss, entropy, gradient norm and alignment through smoothed ratios, a global median and conservative clamping, BLADE adapts step sizes heterogeneously across layers at negligible compute and memory cost. Although both BLADE and a strong model-level baseline diverged under an aggressive schedule, the experiment confirms BLADE’s engineering virtues— $<0.05$  ms overhead per step and constant-size state—and reveals that early divergence is driven by shared hyper-parameters rather than the controller design.

Future work will: (1) harden training through warm-start EMAs, delayed activation and lower initial LR; (2) replicate experiments across multiple seeds to quantify variance; (3) pair BLADE with global LR tuners aut; Jin et al. (2021) and memory-efficient adaptation schemes Dettmers et al. (2023); Zhang et al. (2023); and (4) evaluate on cleaner and multilingual maths benchmarks Zhang et al. (2024a); Shi et al. (2022) as well as large synthetic corpora Toshniwal et al. (2024). All artefacts are open-sourced to encourage community exploration of byte-sized optimisation control mechanisms.

This work was generated by AIRAS (Tanaka et al., 2025).

## REFERENCES

Mechanic: A learning rate tuner.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms. 2023.

Matthew Fahrback, Adel Javanmard, Vahab Mirrokni, and Pratik Worah. Learning rate schedules in the presence of distribution shift. *Proceedings of the 40th International Conference on Machine Learning (ICML 2023)* 9523-9546, 2023.

Alexander Hägele, Elie Bakouch, Atli Kosson, Loubna Ben Allal, Leandro Von Werra, and Martin Jaggi. Scaling laws and compute-optimal training beyond fixed training durations. 2024.

Yuchen Jin, Tianyi Zhou, Liangyu Zhao, Yibo Zhu, Chuanxiong Guo, Marco Canini, and Arvind Krishnamurthy. Autolrs: Automatic learning-rate schedule by bayesian optimization on the fly. 2021.

Tyler B. Johnson, Pulkit Agrawal, Haijie Gu, and Carlos Guestrin. Adascale sgd: A user-friendly algorithm for distributed training. 2020.

Qijun Luo, Hengxu Yu, and Xiao Li. Badam: A memory efficient full parameter optimization method for large language models. 2024.

Xinyu Ma, Xu Chu, Zhibang Yang, Yang Lin, Xin Gao, and Junfeng Zhao. Parameter efficient quasi-orthogonal fine-tuning via givens rotation. 2024.

Niru Maheswaranathan, David Sussillo, Luke Metz, Ruoxi Sun, and Jascha Sohl-Dickstein. Reverse engineering learned optimizers reveals known and novel mechanisms. 2020.

Fabian Schaipp, Ruben Ohana, Michael Eickenberg, Aaron Defazio, and Robert M. Gower. Momo: Momentum models for adaptive learning rates. 2023.

Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. Language models are multilingual chain-of-thought reasoners. 2022.

Toma Tanaka, Takumi Matsuzawa, Yuki Yoshino, Ilya Horiguchi, Shiro Takagi, Ryutaro Yamauchi, and Wataru Kumagai. AIRAS, 2025. URL <https://github.com/airas-org/airas>.

Shubham Toshniwal, Ivan Moshkov, Sean Narenthiran, Daria Gitman, Fei Jia, and Igor Gitman. Openmathinstruct-1: A 1.8 million math instruction tuning dataset. 2024.

Yuhui Xu, Lingxi Xie, Xiaotao Gu, Xin Chen, Heng Chang, Hengheng Zhang, Zhengsu Chen, Xiaopeng Zhang, and Qi Tian. Qa-lora: Quantization-aware low-rank adaptation of large language models. 2023.

Hugh Zhang, Jeff Da, Dean Lee, Vaughn Robinson, Catherine Wu, Will Song, Tiffany Zhao, Pranav Raja, Charlotte Zhuang, Dylan Slack, Qin Lyu, Sean Hendryx, Russell Kaplan, Michele Lunati, and Summer Yue. A careful examination of large language model performance on grade school arithmetic. 2024a.

Qingru Zhang, Minshuo Chen, Alexander Bukharin, Nikos Karampatziakis, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. 2023.

Yihua Zhang, Pingzhi Li, Junyuan Hong, Jiaxiang Li, Yimeng Zhang, Wenqing Zheng, Pin-Yu Chen, Jason D. Lee, Wotao Yin, Mingyi Hong, Zhangyang Wang, Sijia Liu, and Tianlong Chen. Revisiting zeroth-order optimization for memory-efficient llm fine-tuning: A benchmark. 2024b.