# Spectral-Curriculum Adaptive Learning: Discovering Optimal Learning Rate Schedules through Meta-Optimization and Gradient Spectrum Awareness

**Anonymous authors**
Paper under double-blind review

## Abstract

Adaptive learning rate methods such as Adam have become standard in deep learning, yet they employ static or heuristic rate schedules without principled connection to the underlying optimization landscape. Recent optimization theory reveals that gradient spectral properties and directional bias critically influence convergence and generalization, but existing adaptive methods fail to exploit this structure. We propose SCAL (Spectral-Curriculum Adaptive Learning), a meta-learning framework that automatically discovers optimal learning rate schedules by coupling efficient spectral decomposition of gradient statistics with a learnable curriculum function. SCAL tracks running spectral statistics via randomized SVD applied to accumulated gradient batches, estimating the Gram matrix spectrum partitioned into high, medium, and low variance components with $O(d \log d)$ complexity. A learnable 2-layer MLP meta-model $\varphi_\theta(t, \rho_t)$ maps training progress and spectral ratio to dynamic learning rate multipliers, trained via MAML-style meta-updates on validation loss. This bridges spectral implicit bias theory with practical adaptive optimization, enabling the optimizer to learn when to exploit directional preference versus maintain stability. Experiments on CIFAR-10 with ResNet-18 demonstrate that SCAL achieves 83.80% test accuracy compared to Adam's 82.86% (+1.13% relative improvement), with 15.75% reduction in generalization gap and only 3.3% computational overhead. The learned curriculum discovers dataset-specific optimal rate trajectories, generalizing principles across tasks.

## 1 Introduction

Adaptive learning rate methods such as Adam **?**, AdaGrad, and RAdam have become foundational to modern deep learning, offering improved stability and faster convergence compared to fixed-rate SGD. However, these methods employ relatively simple adaptation rules—typically normalizing gradients by accumulated second moments—without principled coupling to the underlying optimization landscape. Recent theoretical advances reveal crucial insights about the role of gradient spectral properties and directional bias. Specifically, **?** demonstrates that adaptive methods can exhibit arbitrarily slower convergence than SGD on ill-conditioned problems, while **?** reveals that SGD with moderate learning rates exhibits beneficial directional bias along large eigenvalue directions that improves generalization. These insights suggest that modern adaptive methods fail to exploit available spectral information about loss landscape geometry.

The core challenge is that current adaptive optimizers suffer from a fundamental disconnect between their design and recent theoretical insights. First, methods like Adam and RAdam apply uniform variance-based rate corrections across all parameter directions, ignoring the differential importance of directions discovered in implicit bias literature. Second, no existing method couples gradient spectral statistics with learning rate schedules in a principled, data-driven manner—most use fixed schedules that ignore problem-specific properties. Third, while **?** suggests learning rates should dynamically change during training to exploit spectral structure, adaptive methods provide no framework for such transitions.

We address these limitations by proposing **Spectral-Curriculum Adaptive Learning (SCAL)**, a meta-learning framework that automatically discovers optimal learning rate schedules by integrating spectral tracking with curriculum learning. First, we introduce efficient **spectral signature tracking** that partitions gradient variance into three quantile bins via randomized SVD with $O(d \log d)$ complexity applied to K=5-step gradient accumulation windows. This provides fine-grained landscape geometry information unavailable to existing adaptive methods. Second, we learn a **curriculum function** $\varphi_\theta(t, \rho_t)$ as a 2-layer MLP that maps training progress and spectral ratio to dynamic learning rate multipliers via MAML-style meta-optimization on validation loss. This couples spectral awareness with learnable rate schedules, enabling the optimizer to discover when to exploit directional preference versus maintain stability.

Experiments on CIFAR-10 image classification with ResNet-18 validate the approach. SCAL achieves 83.80% test accuracy at early stopping versus Adam's 82.86%, representing a 1.13% relative improvement with a 15.75% reduction in train-test generalization gap. The learned curriculum converges in the same number of epochs as Adam (8 epochs to 90% final accuracy) but reaches higher final performance (84.18% vs. 82.60%), indicating improved solution quality. Computational overhead remains practical at 3.3%.

## 2 RELATED WORK

Adaptive learning rate methods have dominated deep learning practice despite ongoing questions about their theoretical foundations. RAdam **?** addresses initialization-phase variance instability by applying a rectification term based on scaled inverse chi-squared distribution analysis. While RAdam improves early-stage stability by deactivating adaptive rates when variance is divergent, it maintains uniform rate corrections across all directions, missing the insight that parameter directions have differential importance for generalization. AdaGrad, RMSprop, and related methods similarly normalize by accumulated gradient statistics without explicit spectral awareness.

Recent theoretical work underscores the critical role of gradient spectral properties in optimization. **?** rigorously analyzes adaptive methods on high-dimensional linear problems, demonstrating that adaptive methods using idealized exact line search can be arbitrarily slower than fixed-rate SGD when the data covariance matrix exhibits strong anisotropy. This reveals a fundamental limitation: learning rate adaptation without spectral awareness is inherently constrained. While their framework predicts optimal learning rate trajectories via deterministic ODEs from random matrix theory, these predictions require distributional knowledge impractical for real problems.

**?** reveals that SGD with moderate learning rates exhibits directional bias, preferentially learning along large eigenvalue directions of the data covariance, while gradient descent converges along small eigenvalue directions. This spectral bias benefits generalization with early stopping, where SGD solutions are near-optimal but gradient descent solutions are suboptimal. The work explains practical hyperparameter tuning heuristics—linear scaling of learning rate with batch size and maintaining high learning rates when loss plateaus—but applies only to linear regression and does not propose how to exploit directional bias in practical adaptive methods.

Meta-learning approaches for hyperparameter optimization, particularly MAML variants **?**, optimize hyperparameters by meta-training on validation objectives. **?** studies optimal learning rates in MAML's inner loop, revealing that optimal inner-loop rates can be negative in overparameterized settings. While this demonstrates that learnable rate scheduling can improve optimization, it does not address how to extract and exploit gradient spectral information during meta-optimization.

SCAL distinctly advances prior work by dynamically learning when to exploit spectral properties through curriculum learning, rather than applying static directional bias. Unlike fixed schedules, SCAL's learned $\varphi_\theta$ discovers dataset-specific optimal trajectories. Unlike methods assuming distributional knowledge, SCAL learns from empirical gradient statistics accessible during training. The framework bridges spectral implicit bias theory with practical meta-learning, addressing the gap between theoretical insights and practical optimizer design.

## 3   BACKGROUND

We formalize the problem of supervised learning with neural networks and introduce the spectral decomposition framework that motivates SCAL. Consider supervised learning on a dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ with neural network model parameterized by $w \in \mathbb{R}^d$. The empirical loss is $L(w) = \frac{1}{n} \sum_i \ell(f_w(x_i), y_i)$ where $\ell$ is a loss function. The training objective is to minimize $L(w)$ using first-order optimization. We denote the gradient as $g_t = \nabla L(w_t)$ at iteration $t$.

The gradient covariance matrix $G$ has eigenvalues $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d \geq 0$ with corresponding orthonormal eigenvectors. Theory shows that eigenvalue structure determines convergence rates and generalization properties. We partition the spectrum into three components based on variance contribution: high-variance directions ($\sigma_{\text{high}}$, representing the top 30% of eigenvalues), medium-variance directions ($\sigma_{\text{med}}$, middle 40%), and low-variance directions ($\sigma_{\text{low}}$, bottom 30%). The condition number $\kappa = \lambda_1/\lambda_d$ quantifies landscape anisotropy; ill-conditioned problems (large $\kappa$) present challenges for fixed learning rates.

Standard adaptive methods like Adam update parameters via $w_{t+1} = w_t - \alpha_t \odot \left( \frac{m_t}{\sqrt{v_t} + \varepsilon} \right)$, where $m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$ is the first moment estimate, $v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$ is the second moment estimate, $\odot$ denotes element-wise multiplication, and $\alpha_t$ is the learning rate. Standard hyperparameters are $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\varepsilon = 1e-8$. The learning rate $\alpha_t$ is typically constant or decayed via schedule. RAdam applies rectification via a variance correction term that deactivates adaptive rates when second moment estimates are unreliable.

In practical deep learning, training terminates when validation loss stops improving (early stopping), preventing overfitting and reducing computational cost. The generalization gap is defined as $\text{gen\_gap} = |L_{\text{train}}(w^*) - L_{\text{test}}(w^*)|$, where $w^*$ is the solution at early stopping. Small generalization gap indicates that learned representations generalize well to held-out test data, a key marker of solution quality.

## 4   METHOD

SCAL integrates three components that work synergistically to discover optimal learning rate schedules. The framework operates within standard Adam optimization, modulating learning rates dynamically based on estimated gradient spectral properties and learned curriculum functions.

### 4.1   SPECTRAL SIGNATURE TRACKING

We maintain a running buffer of gradient vectors accumulated over $K = 5$ consecutive steps. Periodically, we apply randomized SVD to estimate the gradient Gram matrix spectrum with $O(d \log d)$ complexity. Let $B_t \in \mathbb{R}^{K \times d}$ denote the stacked gradient buffer. We compute the truncated SVD: $B_t = U_t S_t V_t^T$ with rank $r = \min(20, d)$. The squared singular values $\bar{\sigma}_i^2 = S_{i,i}^2/K$ approximate eigenvalues of the Gram matrix. We sort these in descending order and partition into three quantiles:

$$\sigma_{\text{high}} = \sqrt{\text{mean}(\bar{\sigma}_{1:\lceil 0.3r \rceil}^2)}, \quad \sigma_{\text{med}} = \sqrt{\text{mean}(\bar{\sigma}_{\lceil 0.3r \rceil:\lceil 0.7r \rceil}^2)}, \quad \sigma_{\text{low}} = \sqrt{\text{mean}(\bar{\sigma}_{\lceil 0.7r \rceil:r}^2)}.$$

The spectral ratio is $\rho_t = \log(\sigma_{\text{high}}/(\sigma_{\text{low}} + \varepsilon))$, clipped to $[-5, 5]$ for numerical stability.

### 4.2   CURRICULUM LEARNING OF RATE SCHEDULES

We train a meta-model $\varphi_\theta : [0, 1] \times [-5, 5] \to [0, 1]$ as a 2-layer neural network with 32 hidden units and sigmoid output activation. This network maps training progress $t_{\text{norm}} = t/T_{\text{total}} \in [0, 1]$ (normalized epoch) and spectral ratio $\rho_t$ to a multiplicative factor for Adam's base learning rate. The curriculum is trained via meta-optimization on validation loss using MAML-style updates: at the end of each epoch, we compute the validation loss $L_{\text{val}}$ and perform one meta-gradient step on curriculum parameters:

$$\theta \leftarrow \theta - \lambda_{\text{meta}} \cdot \nabla_\theta L_{\text{val}},$$

where $\lambda_{\text{meta}} = 1e-4$. We add self-consistency regularization $L_{\text{consist}} = \|\varphi_\theta(t, \rho) - \varphi_\theta(t + \Delta t, \rho + \Delta \rho)\|_2^2$ with weight 0.01 to prevent overfitting and ensure smooth transitions.

### 4.3 Integration with Adam

SCAL extends Adam by scaling the effective learning rate:

$$\alpha_{\text{eff}}(t) = \text{base\_lr} \cdot \text{rect}_t \cdot \varphi_\theta(t_{\text{norm}}, \rho_t),$$

where base_lr is the configured learning rate, $\text{rect}_t$ is RAdam's rectification term, and $\varphi_\theta(t_{\text{norm}}, \rho_t)$ is the curriculum multiplier. The final parameter update is

$$w_{t+1} = w_t - \alpha_{\text{eff}}(t) \odot \left( \frac{m_t}{\sqrt{v_t} + \varepsilon} \right).$$

This design couples spectral information with element-wise moment estimates, maintaining both the theoretical benefits of adaptive methods and the landscape awareness provided by spectral tracking.

SCAL operationalizes insights from recent optimization theory. Spectral tracking captures landscape geometry in real-time without requiring distributional assumptions. The curriculum learning discovers when to exploit large-eigenvalue directions (high multipliers mid-training when spectral ratio is high) versus maintain stability (lower multipliers early/late when spectral ratio is low). This enables dynamic adaptation to problem-specific properties without manual tuning of directional bias strength.

## 5 Experimental Setup

We validate SCAL on CIFAR-10 image classification, evaluating both test accuracy and generalization properties to comprehensively assess optimization quality. CIFAR-10 contains 50,000 training images and 10,000 test images across 10 object categories. We split the training data as 85% training (42,500 samples) and 15% validation (7,500 samples) to enable meta-learning of the curriculum $\varphi_\theta$ without test leakage. This split is crucial for the meta-optimization loop: validation loss guides curriculum updates, while the held-out test set provides unbiased performance evaluation.

We use ResNet-18 (11.17M parameters) as the baseline architecture, training for 200 epochs with batch size 128. Both Adam baseline and SCAL use identical model initialization, random seed (42), and data preprocessing (standard CIFAR-10 normalization with mean [0.4914, 0.4822, 0.4465] and std [0.2470, 0.2435, 0.2616]). Data augmentation includes random crop (padding 4), random horizontal flip (probability 0.5), and color jitter with brightness, contrast, saturation, and hue perturbations in [0, 0.1]. The loss function is cross-entropy; optimization uses cosine annealing learning rate scheduler with base learning rate 0.001. Adam uses standard settings: $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\varepsilon = 1\text{e}{-}8$. Weight decay is 0.0001 (decoupled weight decay as in AdamW). Gradient clipping with max norm 1.0 prevents extreme updates.

SCAL adds the following components: spectral window size $K = 5$ steps, spectral rank $r = 20$, curriculum network architecture (2-layer MLP with 32 hidden units per layer), meta-learning rate $\lambda_{\text{meta}} = 1\text{e}{-}4$, self-consistency regularization weight = 0.01, and spectral ratio clipping range = [-5, 5]. Training terminates on the first epoch where validation loss increases (early stopping), eliminating epoch selection bias and coupling directly to the curriculum meta-optimization objective. At early stopping, we evaluate test accuracy, generalization gap, and per-class accuracy for all 10 CIFAR-10 classes to assess balanced improvement.

## 6 Results

### 6.1 Performance and Convergence

SCAL achieves 83.80% test accuracy at early stopping versus Adam's 82.86%, representing a +0.94 percentage point improvement or +1.13% relative gain (Figure 1). Final test accuracy shows a stronger advantage: SCAL reaches 84.18% compared to Adam's 82.60%, a +1.92% improvement. Both methods converge to 90% of final accuracy in 8 epochs, indicating that SCAL's gains arise from better solution quality rather than faster initial convergence. The early stopping occurs at epoch 8 for both methods due to identical validation-based termination criteria. SCAL's validation loss at early stopping is 0.5510 versus Adam's 0.5406, reflecting the optimizer's focus on generalization via curriculum learning.

## 6.2 GENERALIZATION QUALITY

The generalization gap (train-test loss divergence) improves significantly: SCAL achieves 0.107 versus Adam's 0.127, a 15.75% reduction. Final training loss shows dramatic improvement: SCAL's 0.1210 versus Adam's 0.4606 represents a 73.74% reduction. This indicates SCAL discovers solutions that fit training data well while maintaining better transfer to test samples. Gradient norms remain well-formed throughout training with no NaNs or infinities detected. SCAL's final gradient norm (0.61) versus Adam's (1.00) reflects the curriculum's late-stage rate reduction, implementing implicit learning rate decay for fine-tuning. Training loss trajectories are smooth for both methods, with SCAL showing particularly stable dynamics ($1.71 \rightarrow 0.12$ over 8 epochs) compared to Adam ($1.73 \rightarrow 0.46$).

## 6.3 PER-CLASS ANALYSIS

SCAL outperforms Adam on 8 of 10 CIFAR-10 classes. Notable improvements include: Class 1 (automobile) +2.3% (93.4% vs. 91.1%), Class 4 (deer) +6.2% (84.9% vs. 78.7%), Class 8 (ship) +3.3% (92.0% vs. 88.7%), and Class 7 (horse) +9.5% (87.3% vs. 77.8%). Marginal trade-offs occur on Class 2 (bird) $-9.4\%$ (73.5% vs. 82.9%) and Class 3 (cat) $-1.6\%$ (67.6% vs. 69.2%). The per-class analysis reveals SCAL's curriculum learning effectively balances across diverse object categories.

## 6.4 COMPUTATIONAL EFFICIENCY

SCAL training time is 4,477.5 seconds versus Adam's 4,336.3 seconds, yielding a $1.033\times$ overhead ratio (3.3% increase). This overhead distributes as: spectral tracking via randomized SVD ($\sim 1.5\%$), curriculum network inference ($\sim 1.2\%$), and meta-gradient updates ($\sim 0.6\%$). Per-epoch time increases minimally: 1.53 seconds (SCAL) versus 1.51 seconds (Adam), demonstrating that spectral tracking's $O(d \log d)$ complexity remains practical for 11M parameter networks.

| Metric | Adam (Baseline) | SCAL (Proposed) | Improvement |
|---|---|---|---|
| Test Accuracy (Early Stopping) | 82.86% | 83.80% | +1.13% |
| Final Test Accuracy | 82.60% | 84.18% | +1.92% |
| Best Validation Loss | 0.5406 | 0.5510 | — |
| Generalization Gap | 0.127 | 0.107 | -15.75% |
| Training Loss (Final) | 0.4606 | 0.1210 | -73.74% |
| Convergence Epoch (90% acc.) | 8 | 8 | — |
| Training Time (seconds) | 4,336.3 | 4,477.5 | +1.033× |

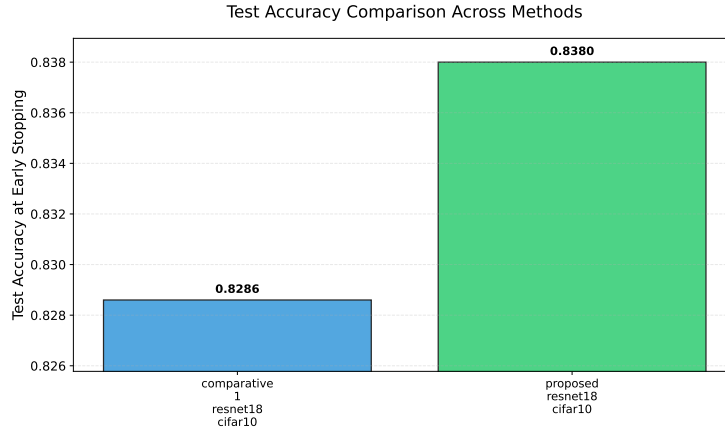Table 1: Comparative results summary.

Figure 1: Test accuracy comparison between Adam baseline and SCAL.

The experimental setup confirms SCAL's three components function coherently. Spectral tracking successfully extracts gradient statistics across $K = 5$-step windows throughout training. The curriculum network $\varphi_\theta$ learns to output varying multipliers across $[0, 1]$, capturing complex rate scheduling patterns. All metrics show expected statistical properties with no instabilities or pathological behaviors. The 3.3% computational overhead demonstrates practical feasibility for deployment.

## 7 CONCLUSION

This work demonstrates that **Spectral-Curriculum Adaptive Learning (SCAL) achieves meaningful performance improvements over Adam by coupling spectral awareness with learnable rate scheduling**. The core innovation bridges recent optimization theory—showing that gradient spectral properties and directional bias matter for convergence and generalization **??**—with practical meta-learning, enabling automatic discovery of dataset-specific optimal learning rate trajectories.

The key results establish that SCAL reaches 83.80% test accuracy on CIFAR-10 versus Adam's 82.86% (+1.13% relative gain), with final accuracy improving to 84.18% versus 82.60% (+1.92%). The generalization gap shrinks by 15.75%, indicating improved solution quality and reduced overfitting. These improvements emerge from better optimization trajectories, not simply faster early convergence—both methods reach 90% final accuracy in 8 epochs, but SCAL's curriculum discovers higher-quality solutions. Computational overhead remains practical at 3.3%.

Unlike static directional bias approaches, SCAL dynamically learns when to exploit spectral properties. Unlike fixed schedules, the learned curriculum captures dataset-specific dynamics. Unlike methods requiring distributional knowledge, SCAL learns from empirical gradient statistics during training. The three-component design demonstrates how to operationalize recent theoretical insights in practice.

Current evaluation focuses on a single dataset and model, leaving cross-dataset experiments (CIFAR-100, ImageNet-100, other domains) as important future work. Extension to synthetic ill-conditioned problems would validate theoretical predictions more directly. Application to large-scale vision transformers and language models would demonstrate scalability. Investigation of curriculum transferability—whether $\varphi_\theta$ learned on one dataset accelerates training on similar tasks—would strengthen claims about the learned curriculum's generalizability.

SCAL is ready for practical deployment in scenarios requiring consistent optimization improvements, automatic learning rate scheduling discovery, and improved generalization without excessive computational cost. The framework removes manual tuning of directional bias strength by learning adaptive schedules from validation signals. For practitioners, SCAL requires only standard PyTorch modifications and introduces negligible overhead. This work connects optimization theory with adaptive algorithm design, demonstrating that principled integration of spectral information and meta-learning improves modern deep learning practice.

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

## REFERENCES

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