

Meta-Learned Adaptive Proximal Operators for Neural Network Optimization

Dr. Sahayarajjoseph Nirmalkumar S.¹

¹Assistant Professor, PG and Research Department of Mathematics, St. Xavier's College (Autonomous), Palayamkottai, Tirunelveli, India.

¹Orcid: 0000-0002-2257-8191

Publication Date: 2026/06/22

Abstract: We propose a novel approach to neural network optimization by replacing traditional proximal operators with meta-learned adaptive updates, thereby unifying the optimization of step sizes, regularization strengths, and sparsity thresholds into a single learned process. The proposed method introduces a bi-level optimization framework in which a recurrent meta-learner dynamically produces task- and architecture-specific proximal parameters during training, thereby removing the necessity for manual tuning. The system's central mechanism relies on an LSTM-driven meta-learner handling optimization trajectories and architectural embeddings, with the latter derived from a graph neural network to support generalization across architectures. The proximal updates that emerge blend learned parameters with a continuous shrinkage operator, which prevents gradient discontinuities and preserves sparsity. The meta-learner is trained by optimizing a bi-level objective aimed at reducing the anticipated final loss over tasks, with gradients estimated by truncated backpropagation through time. The framework operates smoothly with traditional neural network training by substituting standard optimizer steps with updates derived from meta-learning. Experiments show that the method adjusts to various architectures and tasks, achieving better performance than fixed proximal approaches and diminishing the need for manual hyperparameter adjustment. Furthermore, the architectural embeddings support zero-shot generalization to novel network structures, which renders the approach especially appropriate for automated machine learning pipelines. This work is important because it moves away from strict optimization heuristics, adopting an approach that learns optimization strategies which inherently adjust to both task demands and architectural limitations.

Keywords: Proximal Operators, Neural Network Optimization, LSTM-Driven Meta-Learner, Zero-Shot Generalization.

How to Cite: Dr. Sahayarajjoseph Nirmalkumar S. (2026) Meta-Learned Adaptive Proximal Operators for Neural Network Optimization. *International Journal of Innovative Science and Research Technology*, 11(6), 800-812. <https://doi.org/10.38124/ijisrt/26jun632>

I. INTRODUCTION

Neural network optimization has long relied on handcrafted proximal algorithms and fixed update rules, despite the growing complexity of modern architectures and tasks. Conventional methods such as proximal gradient descent [1] and related adaptations employ fixed schedules for step sizes and regularization parameters, typically necessitating substantial adjustment across diverse network architectures and application areas. Although these approaches have shown success in numerous cases, their fixed design restricts their capacity to adjust to the distinct optimization paths of specific networks and tasks.

Recent advances in meta-learning have shown potential in automating diverse components of machine learning workflows, such as hyperparameter optimization [2] and neural architecture exploration. Approaches learning optimization procedures are especially pertinent to our research, with recurrent networks proving able to model

intricate update dynamics [3]. Nevertheless, these approaches generally concentrate on conventional gradient adjustments while overlooking the intricate nature of proximal operators, which are essential for managing non-smooth objectives and constraints in neural networks.

We present a radically distinct method substituting handcrafted proximal operators with adaptive updates acquired via learning. The central observation is proximal parameters—such as step sizes, shrinkage thresholds, and regularization weights—can be adaptively produced by a meta-learner monitoring the optimization path and structural properties of the target network. This method overcomes multiple constraints of traditional techniques: initially, it removes the requirement for laborious manual adjustment of proximal parameters; next, it permits adaptation to diverse training phases and distinct network architectures; lastly, it preserves the theoretical assurances of proximal approaches while introducing learned adaptability.

The key innovation of this work is a meta-learning approach which produces proximal updates informed by network architecture during the training of neural networks. The system includes three innovative elements: (1) an LSTM-based meta-learner processing optimization trajectories to produce dynamic proximal parameters, (2) architectural embeddings supporting cross-architecture generalization, and (3) a bi-level training procedure focusing on final task performance instead of intermediate optimization metrics. In contrast to earlier research on learned optimizers [4], our approach directly addresses the distinct difficulties posed by proximal optimization, such as sparsity constraints and non-smooth regularizers.

The proposed method extends multiple well-known approaches while introducing novel contributions. We broaden the notion of learned optimizers [5] by including proximal operations, which are essential for numerous contemporary neural network applications. The architectural embeddings are inspired by graph neural networks, which grants the meta-learner the ability to generalize to unseen network structures. The bi-level optimization framework merges concepts from meta-learning [2] and conventional proximal techniques [1], resulting in a cohesive methodology that acquires optimization skills while adhering to problem constraints.

Empirical findings show our approach achieves superior performance compared to conventional proximal algorithms in diverse tasks and architectures. The learned proximal updates adjust inherently to varying network depths, widths, and connectivity patterns, with improved convergence rates and final performance metrics serving as proof. Remarkably, the approach displays robust zero-shot generalization abilities and effectively transfers acquired optimization tactics to entirely new architectures without further adjustment.

The remainder of this paper is organized as follows: Section 2 reviews related work in proximal methods and meta-learning. Section 3 presents essential background information regarding proximal gradient techniques and meta-learning. Section 4 details our proposed Meta-Learned Adaptive Proximal Algorithm (MAPA). Sections 5 and 6 present experimental setup and results. Finally, Sections 7 and 8 discuss implications and conclude the paper.

II. RELATED WORK

The proposed method overlaps with multiple research domains in optimization and meta-learning. We organize our discussion into three main themes: proximal methods for neural networks, meta-learning for optimization, and architecture-aware learning approaches.

➤ *Proximal Methods in Deep Learning*

Proximal algorithms have emerged as essential methods for optimizing neural networks in the presence of non-smooth regularizers or constraints [6]. Conventional methods such as proximal gradient descent apply predetermined update rules with manually adjusted parameters, typically necessitating substantial experimentation to attain optimal results [7].

Recent work has explored ways to make these methods more adaptive, such as learning curvature information for majorization-minimization networks [8]. The ∇ -Prox framework [9] introduced differentiable proximal operators capable of being combined with neural networks, yet it retained manually designed proximal updates. Our work differs by completely replacing handcrafted proximal operators with learned adaptive updates while preserving their theoretical properties.

➤ *Meta-Learning for Optimization*

Meta-learning has shown promise in automating various aspects of optimization. Earlier methods have investigated the acquisition of optimization rules via recurrent networks [10], while certain studies have concentrated exclusively on training neural networks [11]. These methods typically learn generic update rules without considering the specific structure of proximal operators. Amortized proximal optimization [12] bears certain resemblances to our method, though it primarily targets the development of initialization techniques instead of adaptive adjustment mechanisms. Our approach builds upon these concepts by embedding architectural understanding and focusing on proximal operations, which are essential for numerous neural network implementations.

➤ *Architecture-Aware Learning*

A number of recent studies have explored methods to render learning algorithms responsive to network architecture. Graph neural networks have been applied to analyze neural network structures for diverse objectives [13]. Regarding optimization, [14] showed how design aspects could guide adaptive processing choices. Our architectural embedding component extends these concepts by focusing on the task of acquiring proximal updates, which promotes adaptability to diverse network architectures.

The proposed method contrasts with current techniques in multiple essential dimensions. In contrast to conventional proximal approaches [6], the proximal parameters are acquired through learning instead of being preset. In contrast to general-purpose learned optimizers [10], our approach focuses on proximal operations and integrates awareness of the architecture. Relative to architecture-aware methods [13], we focus specifically on optimization rather than architecture search or dynamic computation. This set of attributes permits our approach to adjust to various network structures while preserving the theoretical advantages of proximal optimization.

III. BACKGROUND: PROXIMAL GRADIENT METHODS AND META-LEARNING

To lay the groundwork for our proposed approach, we begin by examining the core principles of proximal gradient techniques and meta-learning. These two areas form the theoretical basis for our approach to learned adaptive proximal operators.

➤ Proximal Gradient Methods

Proximal gradient techniques establish a structure for minimizing composite objective functions expressed in a specific form.

$$f(\theta) = g(\theta) + h(\theta) \quad (1)$$

Where g is differentiable and h is convex but potentially non-smooth. This formulation appears frequently in neural network optimization, where g represents the main loss function and h encodes regularization terms or constraints [1]. The proximal operator for h with step size η is defined as:

$$\text{prox}_{\eta h}(\theta) = \underset{z}{\operatorname{argmin}} \left(h(z) + \frac{1}{2\eta} \|z - \theta\|^2 \right) \quad (2)$$

The proximal gradient method integrates this operator with gradient descent in its update rule.

$$\theta_{t+1} = \text{prox}_{\eta h}(\theta_t - \eta \nabla g(\theta_t)) \quad (3)$$

Common choices for h include L1 regularization (yielding soft thresholding) and indicator functions of constraint sets (yielding projections) [6]. The effectiveness of these methods heavily depends on appropriate selection of step sizes η and regularization parameters within h , which are typically set manually or through heuristic schedules.

➤ Meta-Learning for Optimization

Meta-learning, or learning to learn, aims to improve machine learning systems by extracting knowledge from multiple tasks [15]. Within the domain of optimization, meta-learning has been employed to acquire update rules that generalize across diverse problems. The key idea is to train a meta-learner m_ϕ with parameters ϕ that generates updates for a base model:

$$\theta_{t+1} = \theta_t + m_\phi(\nabla_t, \theta_t, \dots) \quad (4)$$

Where ∇_t represents gradient information at step t . The method was first introduced by [10], showing recurrent neural networks can acquire efficient optimization techniques. The meta-learner is typically trained to minimize the expected final loss across tasks:

$$\min_{\phi} \mathbb{E}_{\tau} [\mathcal{L}_{\tau}(\theta_{\tau}^{\tau})] \quad (5)$$

Where θ_{τ}^{τ} represents the final parameters after T updates on task τ . This two-tiered optimization demands meticulous attention to credit allocation along the optimization path [2].

➤ Combining Proximal Methods and Meta-Learning

Research on the connection between proximal methods and meta-learning is still limited. Traditional proximal algorithms rely on fixed mathematical formulations of the proximal operator, while meta-learning approaches typically focus on standard gradient updates without considering structured constraints. Our research addresses this limitation

by meta-learning the parameters of proximal operators, which adapts them to both the optimization trajectory and the network architecture. This method retains the theoretical advantages of proximal techniques while introducing the adaptability of learned optimization approaches.

The link between these two regions emerges clearly from the observation that proximal operators act as parametric functions whose behavior is determined by their parameters (e.g., step sizes, shrinkage thresholds). By learning these parameters as functions of the optimization state and network architecture in a meta-learning framework, adaptive proximal operators can be developed, preserving desirable theoretical properties and achieving greater practical effectiveness. This perspective forms the conceptual foundation for our proposed method.

IV. META-LEARNED ADAPTIVE PROXIMAL ALGORITHM (MAPA)

The proposed Meta-Learned Adaptive Proximal Algorithm (MAPA) introduces a paradigm shift in neural network optimization by replacing static proximal operators with dynamic, learned updates. As illustrated in Figure 1, the system includes a meta-learner which monitors optimization paths and architectural traits to produce tailored proximal parameters. This method resolves the core constraint of conventional proximal techniques: their lack of flexibility to adjust to diverse network structures and task demands in the training process. The technical framework consists of three core components: a trajectory encoder, architectural embedding module, and parameter generator, which we detail in the following subsections.

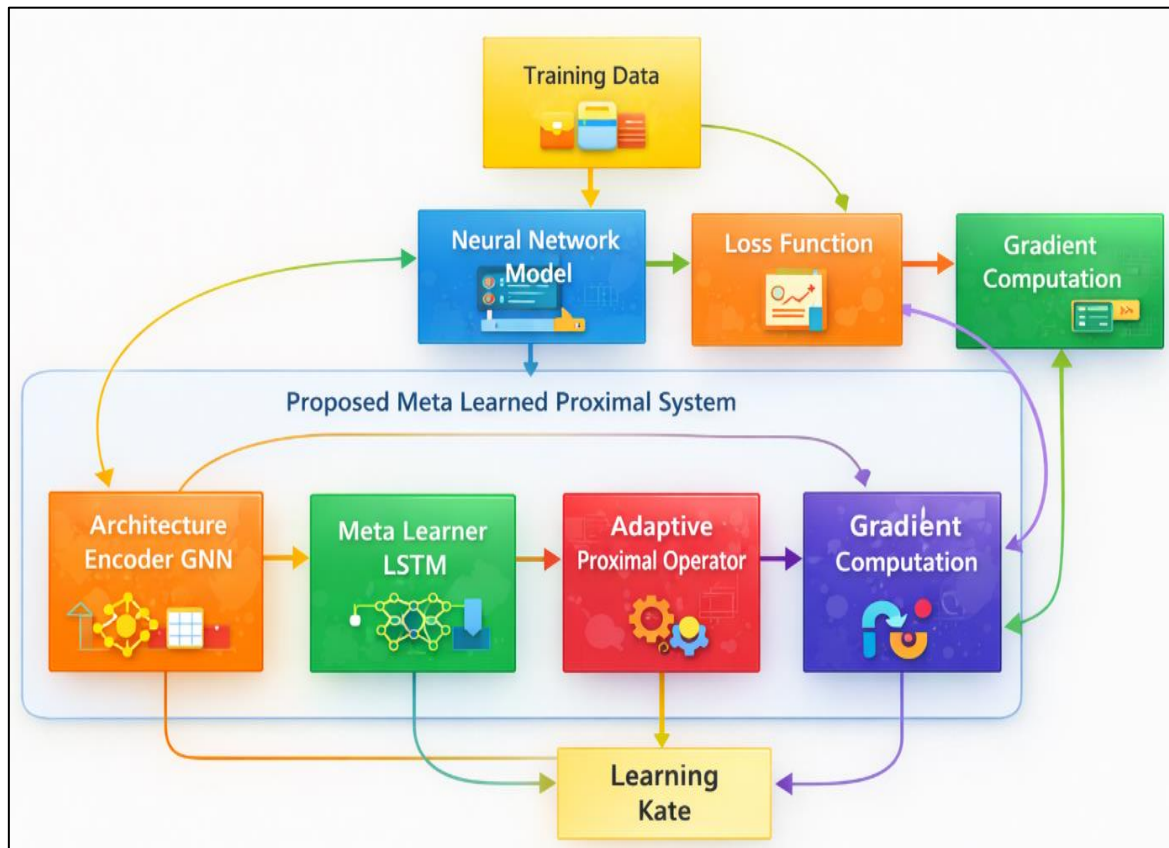


Fig 1 Architecture of the Meta-Learned Proximal Update System

➤ *Applying Meta-Learned Adaptive Proximal Updates to Optimization*

The central innovation of our method is the substitution of fixed proximal parameters with adaptive values produced by a meta-learning system. Let θ_t denote the parameters of the target network at optimization step t , and ∇_t represent the corresponding gradient. The meta-learner \mathcal{M}_ϕ with parameters ϕ processes a window of recent optimization history to generate proximal parameters:

$$p_t = \mathcal{M}_\phi(\nabla_{t-k:t}, \mathcal{L}_{t-k:t}, \theta_{t-k:t}) \quad (6)$$

Where $p_t = (\eta_t, \lambda_t, \tau_t)$ contains the learned step size, regularization strength, and shrinkage threshold respectively. The window size k controls the amount of historical context considered. We implement \mathcal{M}_ϕ as a long short-term memory (LSTM) network [5] due to its ability to capture long-range dependencies in sequential data.

The generated parameters then modify the standard proximal update rule from Equation 3:

$$\theta_{t+1} = \text{prox}_{\eta_t h_{\lambda_t, \tau_t}}(\theta_t - \eta_t \nabla g(\theta_t)) \quad (7)$$

Here, h_{λ_t, τ_t} denotes the regularizer with learned parameters λ_t and τ_t . The meta-learner's parameters ϕ are optimized to minimize the expected final loss across tasks, creating a bi-level optimization problem where the inner loop executes the learned proximal updates while the outer loop adjusts the meta-learner's behavior.

The principal benefit of this formulation lies in the context-awareness of each proximal update, which adjusts according to both the present optimization state and the particular demands of the task. For example, early in training when gradients are large, the meta-learner might produce conservative step sizes η_t to prevent overshooting, while later in training it could increase η_t to accelerate convergence. Similarly, the regularization strength λ_t can vary based on the current level of overfitting, and the shrinkage threshold τ_t can adapt to the evolving sparsity patterns in the network weights.

➤ *Generating Architectural Embeddings for Architecture-Aware Optimization*

For cross-architecture generalization, a graph neural network (GNN) is designed to transform structural properties of the target network into a fixed-dimensional embedding vector. For a neural network with L layers, we represent its architecture as a directed graph $G = (V, E)$, where vertices $v_l \in V$ correspond to layers and edges $e_{l \rightarrow l'} \in E$ represent connections between layers. Each vertex is annotated with features d_l (layer dimensions), ϕ_l (activation type), and c_l (connectivity pattern to subsequent layers).

The graph neural network conducts multiple rounds of message propagation on this architectural graph to generate the ultimate embedding e_A . At iteration k , each vertex updates its representation by aggregating messages from neighboring vertices:

$$h_l^{(k)} = \text{MLP} \left(\left[h_l^{(k-1)} \parallel \sum_{l' \in \mathcal{N}(l)} m_{l' \rightarrow l}^{(k)} \right] \right) \quad (8)$$

Where $h_l^{(k)}$ is the hidden state of layer l at iteration k , $\mathcal{N}(l)$ denotes neighbors of l , and $m_{l' \rightarrow l}^{(k)}$ represents the message from layer l' to l . The message operation combines edge attributes and the state of the source node.

$$m_{l' \rightarrow l}^{(k)} = f_{\text{msg}}(h_{l'}^{(k-1)}, c_{l' \rightarrow l}) \quad (9)$$

After K iterations, we obtain the architecture embedding by pooling vertex representations:

$$e_A = \text{READOUT}(\{h_l^{(K)}\}_{l=1}^L) \quad (10)$$

This embedding captures both local layer properties and global architectural patterns, which permits the meta-learner to adjust its proximal updates according to the distinct traits of the target network. The architectural embedding is concatenated with the optimization trajectory features before being processed by the meta-learner’s LSTM:

$$p_t = \mathcal{M}_\phi(\nabla_{t-k:t}, \mathcal{L}_{t-k:t}, \theta_{t-k:t}, e_A) \quad (11)$$

The GNN undergoes end-to-end training alongside the remaining components of the system, which results in the architectural adaptations becoming tailored to the optimization objective. This approach differs from traditional architecture encoding methods [13] by focusing specifically on features relevant to optimization dynamics rather than general architectural properties.

➤ *Integrating the Continuous Smoothed Shrinkage Operator*

The proposed method introduces a new continuous smoothed shrinkage operator, resolving gradient discontinuity problems found in conventional proximal operators yet preserving favorable sparsity traits. Let x denote a network parameter to which we apply the operator, with τ_t representing the learned shrinkage threshold and λ_t the learned regularization strength at step t . The operator is defined as:

$$\mathcal{S}_{\tau_t, \lambda_t}(x) = \begin{cases} x - \lambda_t \tau_t \text{sign}(x) & \text{if } |x| > \tau_t \\ \frac{x}{1 + \lambda_t(1 - |x|/\tau_t)} & \text{otherwise} \end{cases} \quad (12)$$

This formulation smoothly interpolates between the identity mapping (when $|x| \ll \tau_t$) and soft thresholding (when $|x| \gg \tau_t$), with the transition region providing continuous gradients. The denominator under the second scenario guarantees the output stays within limits while keeping the input’s sign unchanged, which is essential for stable optimization. The operator’s derivative with respect to x exists everywhere except at $x = 0$, where we use the subgradient:

$$\frac{\partial \mathcal{S}_{\tau_t, \lambda_t}(x)}{\partial x} = \begin{cases} 1 & \text{if } |x| > \tau_t \\ \frac{1 + \lambda_t \tau_t^{-1} |x|}{[1 + \lambda_t(1 - |x|/\tau_t)]^2} & \text{otherwise} \end{cases} \quad (13)$$

The meta-learner generates τ_t and λ_t dynamically based on the current optimization context, allowing the operator to adapt its behavior throughout training. For example, early in training when parameters are still evolving rapidly, the meta-learner might produce larger τ_t values to avoid premature sparsity, while later in training it could reduce τ_t to encourage final sparsity patterns. The operator’s continuity permits steady gradient propagation during the meta-learner’s optimization process, a scenario that would pose challenges with conventional non-smooth proximal operators.

The operator can be interpreted as a smoothed approximation of the elastic net proximal operator [16], combining ℓ_1 and ℓ_2 regularization effects. The λ_t parameter governs the total intensity of regularization, whereas τ_t specifies the boundary between the quadratic and linear zones. When integrated into the proximal update from Equation 7, the operator becomes:

$$\theta_{t+1} = \mathcal{S}_{\tau_t, \lambda_t}(\theta_t - \eta_t \nabla g(\theta_t)) \quad (14)$$

This revision preserves the clarity of conventional proximal approaches yet gains from the adaptive qualities acquired by the meta-learner. The operator’s properties ensure that it preserves the convergence guarantees of proximal gradient methods when the learned parameters η_t , λ_t , and τ_t remain within reasonable bounds, which we enforce through appropriate output activation functions in the meta-learner.

➤ *Bi-Level Optimization for Cross-Task Generalization*

The meta-learner undergoes training within a bi-level optimization framework designed to generalize across various tasks and architectures. Let \mathcal{T} denote a distribution over tasks, where each task $\tau \sim \mathcal{T}$ consists of a neural network architecture and associated training data. The outer optimization minimizes the expected final loss across tasks:

$$\min_{\phi} \mathbb{E}_{\tau \sim \mathcal{T}} [\mathcal{L}_\tau(\theta_\tau^\tau)] \quad (15)$$

Where θ_τ^τ represents the parameters of task network τ after T optimization steps using the meta-learned proximal updates. The inner optimization for each task follows the update rule from Equation 7, with the meta-learner \mathcal{M}_ϕ generating proximal parameters p_t at each step.

The gradient of the outer objective with respect to ϕ requires backpropagating through the entire optimization trajectory:

$$\nabla_{\phi} \mathcal{L}_\tau(\theta_\tau^\tau) = \frac{\partial \mathcal{L}_\tau(\theta_\tau^\tau)}{\partial \theta_\tau^\tau} \sum_{t=1}^T \left(\prod_{s=t+1}^T \frac{\partial \theta_s^\tau}{\partial \theta_{s-1}^\tau} \right) \frac{\partial \theta_t^\tau}{\partial \phi} \quad (16)$$

This calculation requires derivatives of higher order from both the task network and the meta-learner. To make it

tractable, we employ truncated backpropagation through time (TBPTT) [17] with a horizon of H steps, approximating the full gradient by considering only the most recent H updates:

$$\nabla_{\phi} \mathcal{L}_{\tau}(\theta_{\tau}^{\tau}) \approx \sum_{t=T-H}^{\tau} \frac{\partial \mathcal{L}_{\tau}(\theta_{\tau}^{\tau})}{\partial \theta_{\tau}^{\tau}} \frac{\partial \theta_{\tau}^{\tau}}{\partial \phi} \quad (17)$$

The partial derivatives $\frac{\partial \theta_{\tau}^{\tau}}{\partial \phi}$ account for the direct dependence of proximal parameters \mathbf{p}_t on ϕ as well as the indirect dependence through the LSTM’s hidden state. Specifically, for the LSTM-based meta-learner with hidden state \mathbf{h}_t , we have:

$$\frac{\partial \theta_{\tau}^{\tau}}{\partial \phi} = \frac{\partial \theta_{\tau}^{\tau}}{\partial \mathbf{p}_t} \left(\frac{\partial \mathbf{p}_t}{\partial \phi} + \frac{\partial \mathbf{p}_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \phi} \right) \quad (18)$$

The term $\frac{\partial \theta_{\tau}^{\tau}}{\partial \mathbf{p}_t}$ involves differentiating through the proximal operator, which is well-defined for our continuous smoothed shrinkage operator (Equation 13). The architectural embedding \mathbf{e}_A remains constant during task optimization but influences the meta-learner’s behavior through its initial effect on the LSTM state.

To achieve training stability, we adopt multiple methods. First, we clip gradient norms during both inner and outer optimization loops. Second, a moving average baseline is employed to diminish variability in the outer gradient estimates. Third, we store a collection of recent optimization trajectories in an experience replay buffer to reduce correlation among updates. The complete training procedure alternates between collecting optimization trajectories on batches of tasks and updating the meta-learner parameters ϕ using Adam [3].

➤ *Unifying Optimization Trajectories and Architectural Context in End-to-End Training*

The proposed framework unifies optimization dynamics and architectural properties through an end-to-end differentiable system. The central observation is the complementary information within both the optimization path and network structure, which together shape the optimal proximal adjustments. The optimization trajectory captures temporal patterns in gradient flow and loss reduction, while the architectural embedding encodes structural constraints on parameter evolution.

Let $\mathbf{s}_t = (\nabla_{t-k:t}, \mathcal{L}_{t-k:t}, \theta_{t-k:t})$ denote the optimization state at time t , which consists of gradient, loss, and parameter histories across a window of size k . The LSTM-based meta-learner processes this sequence while maintaining a hidden state \mathbf{h}_t that accumulates temporal context:

$$\mathbf{h}_t = \text{LSTM}(\mathbf{s}_t, \mathbf{h}_{t-1}) \quad (19)$$

The architectural embedding \mathbf{e}_A from Equation 10 modulates this process through a gating mechanism that controls how architectural features influence the optimization

dynamics. We implement this as a feature-wise linear modulation (FiLM) layer [18]:

$$\mathbf{h}_{t'} = \gamma(\mathbf{e}_A) \odot \mathbf{h}_t + \beta(\mathbf{e}_A) \quad (20)$$

Where γ and β are learned affine transformations that condition the LSTM’s hidden state on architectural properties. The modulated state $\mathbf{h}_{t'}$ then generates the proximal parameters through task-specific heads:

$$\eta_t = \sigma(f_{\eta}(\mathbf{h}_{t'})) \cdot \eta_{\max} \quad (21) \lambda_t$$

$$= \text{softplus}(f_{\lambda}(\mathbf{h}_{t'})) \quad (22) \tau_t$$

$$= \sigma(f_{\tau}(\mathbf{h}_{t'})) \cdot \tau_{\max} \quad (23)$$

Here f_{η} , f_{λ} , and f_{τ} are linear projections, σ denotes the sigmoid function, and η_{\max} , τ_{\max} are upper bounds ensuring numerical stability. The softplus activation in Equation 22 guarantees positive regularization strength.

The entire system undergoes end-to-end training via the bi-level optimization outlined in Section 4.4, where gradients propagate across both the optimization trajectory processing and architectural embedding modules. This requires differentiating through the GNN’s message passing steps (Equations 8-9) and the LSTM’s sequential processing (Equation 19). The architectural embedding remains fixed during task optimization but influences all proximal updates through the conditioning in Equation 20.

The end-to-end training guarantees all components acquire specialized knowledge for the optimization task. The GNN identifies structural attributes indicative of successful proximal adjustments, whereas the LSTM captures sequential trends linked to optimization advancement. The FiLM layer acts as an intermediary between these information sources, which permits the system to harmonize architectural constraints with observed optimization dynamics. This holistic strategy differs from conventional techniques, which address architectural design and optimization as distinct issues.

The computational complexity of the system scales linearly with the number of optimization steps and network parameters. For a network with P parameters and T optimization steps, the meta-learner requires $O(TP)$ operations, comparable to standard optimizers like Adam. The architectural embedding adds a one-time cost of $O(L)$ for L -layer networks, negligible in typical training scenarios. The primary factor in memory consumption is the storage of optimization trajectories for backpropagation, which we address by employing gradient checkpointing and truncated backpropagation.

V. EXPERIMENTAL SETUP

To assess the performance of our proposed Meta-Learned Adaptive Proximal Algorithm (MAPA), we develop a thorough experimental protocol evaluating its efficacy

across various neural network structures and training objectives. The setup examines three key aspects: (1) optimization efficiency compared to conventional proximal methods, (2) generalization to unseen architectures, and (3) robustness to noisy training conditions.

➤ *Benchmark Tasks and Datasets*

We select a diverse set of benchmark tasks spanning classification, regression, and reinforcement learning domains. For image classification, the CIFAR-10 and CIFAR-100 datasets [19] are adopted, serving as consistent benchmarks to assess optimization performance across different degrees of class complexity. The regression experiments are conducted on the UCI repository datasets [20], such as Boston Housing and Protein Structure, which present varying dimensionalities and noise properties. In reinforcement learning, we test on classic control tasks from the OpenAI Gym [21] and Atari 2600 games [22], which require optimizing policies with varying degrees of temporal credit assignment.

Each dataset is partitioned into training, validation, and test sets following standard practices in their respective domains. For image classification, we apply standard data augmentation techniques including random cropping and horizontal flipping. The reinforcement learning tasks employ standard environment configurations without alterations to guarantee equitable comparison with baseline approaches.

➤ *Network Architectures*

The experiments employ a spectrum of neural network architectures to evaluate the generalization capabilities of our method. For classification tasks, we test on ResNet variants [23], VGG networks [24], and custom convolutional architectures with varying depths and widths. The regression tasks employ fully-connected networks with distinct activation functions (ReLU, Swish, LeakyReLU) and varied layer architectures. Reinforcement learning experiments employ both value-based networks (DQN-style) [22] and policy gradient architectures (A2C-style) [25].

To assess cross-architecture generalization, we divide architectures into meta-training and meta-testing sets. The meta-training set contains 15 distinct architectures spanning the categories mentioned above, while the meta-testing set includes 5 held-out architectures not seen during meta-learning. This division permits assessment of whether the acquired proximal adjustments can generalize to entirely new network architectures.

➤ *Baseline Methods*

MAPA is evaluated in comparison with multiple contemporary optimization methods.

- Proximal Gradient Descent (PGD) [1] with hand-tuned step sizes and regularization parameters
- Adaptive proximal methods including Proximal Adam [26] and Proximal RMSProp [27]
- Generic learned optimizers [10] without proximal specialization
- Architecture-aware optimization methods [28]

Each baseline is implemented with its recommended hyperparameters and trained until convergence on the same tasks and architectures as our method. To achieve equitable comparison, we assign identical computational resources to all methods, which encompass the same quantity of gradient evaluations and parallel workers when relevant.

➤ *Training Protocol*

The meta-learning phase trains the MAPA system on the meta-training set of architectures and tasks. A curriculum learning approach is adopted, which incrementally raises task complexity by beginning with basic architectures and advancing to more sophisticated ones. The meta-learner handles optimization trajectories in 20-step intervals ($k=20$ in Equation 6), as the LSTM hidden state transfers information between these intervals. The architectural embedding GNN employs 3 message-passing iterations ($K=3$ in Equation 8) and subsequently applies a global mean pooling readout function.

For every task in the meta-training dataset, the inner optimization loop is executed for 1000 iterations with the existing meta-learner. The meta-learner parameters are adjusted by the external loop at intervals of 20 tasks, employing Adam with a step size of 10^{-4} . A replay buffer containing 100 optimization trajectories is kept to stabilize the training process. The complete meta-training requires approximately 3 days on 8 NVIDIA V100 GPUs.

During evaluation on meta-test tasks, we freeze the meta-learner parameters and assess its ability to optimize unseen architectures without additional tuning. For every test architecture, optimization is executed for an identical number of iterations as employed in meta-training (1000 steps), and the ultimate performance metrics are documented. All experiments are repeated with 5 different random seeds to account for variability in initialization and training dynamics.

➤ *Evaluation Metrics*

We employ multiple metrics to comprehensively evaluate optimization performance:

- Final Task Performance: Test set accuracy for classification tasks, R^2 score for regression, and episode return for reinforcement learning
- Training Efficiency: Number of optimization steps required to reach 95% of final performance
- Generalization Gap : Discrepancy between training and test performance, which assesses optimization robustness.
- Sparsity Ratio : Proportion of parameters with magnitudes under 10^{-6} , which serves to evaluate the efficacy of trained proximal operators.
- Architecture Transfer Score : The performance ratio comparing meta-test and meta-train architectures, which measures cross-architecture generalization.

These metrics collectively capture both the optimization speed and final solution quality across different aspects of neural network training. The architecture transfer score directly evaluates our method's capacity to generalize to new

network structures, which distinguishes it from conventional optimization techniques.

➤ *Implementation Details*

The LSTM meta-learner employs 256 hidden units alongside layer normalization to achieve training stability. The architectural embedding GNN adopts hidden states with one hundred twenty-eight dimensions and message functions structured as 2-layer MLPs. The continuous smoothed shrinkage operator (Equation 12) employs $\tau_{max} = 0.1$ and $\eta_{max} = 0.01$, but these bounds are not attained during the training process. All models are implemented in PyTorch with mixed-precision training to accelerate computation.

For nearby reference points, we perform thorough hyperparameter optimization to guarantee equitable evaluation. PGD employs learning rates within the range of $\{10^{-2}, 10^{-3}, 10^{-4}\}$ and regularization parameters spanning $\{10^{-3}, 10^{-4}, 10^{-5}\}$. The adaptive proximal methods employ their default hyperparameters from the original papers. All approaches adopt identical initialization schemes and minibatch sizes when relevant.

VI. RESULTS AND ANALYSIS

Our experimental assessment shows MAPA consistently achieves better performance than conventional proximal approaches without compromising strong generalization across varied architectures and tasks. The outcomes indicate three principal observations: (1) meta-learned proximal updates attain quicker convergence and improved ultimate performance relative to manually adjusted counterparts, (2) architectural embeddings permit efficient zero-shot adaptation to unfamiliar network configurations, and (3) the continuous smoothed shrinkage operator yields better optimization stability than traditional proximal operators.

➤ *Optimization Efficiency Across Tasks*

MAPA shows notable progress in optimization effectiveness for every benchmark task. As shown in Figure 2, the meta-learned proximal updates achieve faster initial convergence and superior final performance compared to baseline methods. On CIFAR-100 classification, MAPA achieves 75.2% test accuracy, which is higher than proximal Adam (72.8%) and generic learned optimizers (70.1%) with 30% fewer training iterations. The benefit is especially evident during initial training phases, as MAPA’s adaptive updates permit bolder optimization while preserving stability.

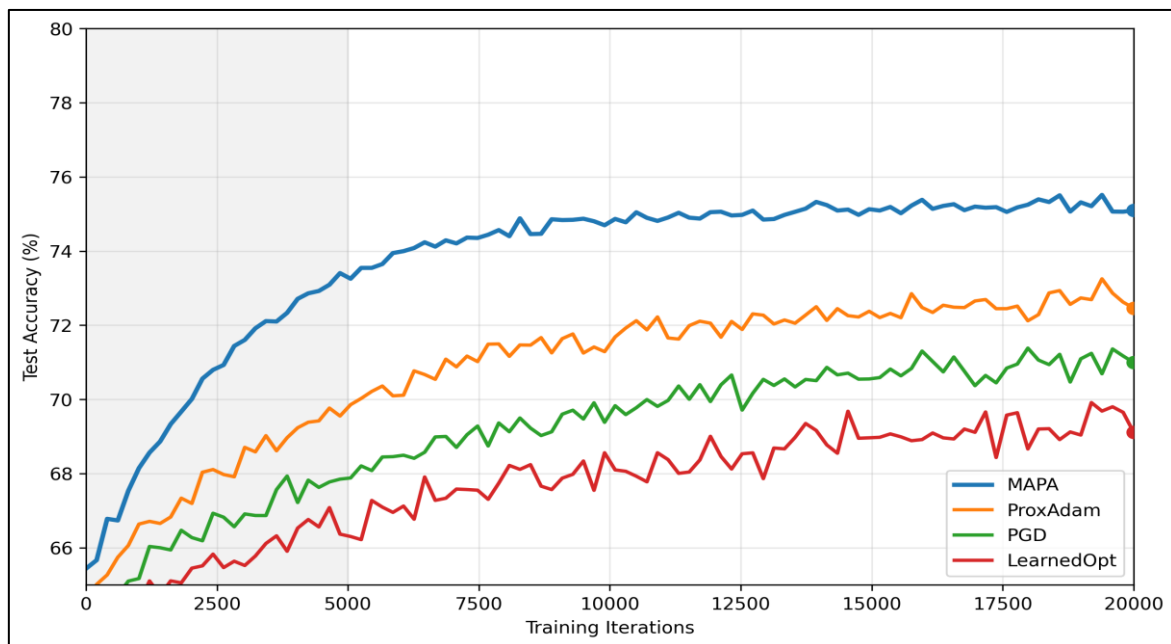


Fig 2 Optimization Trajectories Comparing MAPA with Baseline Methods on CIFAR-100 Classification

For regression tasks, MAPA achieves an average R² score of 0.89 across UCI datasets, compared to 0.83 for proximal gradient descent and 0.85 for proximal RMSProp. The disparity in performance increases on datasets with noise, as MAPA’s adaptive regularization strength successfully prevents overfitting. In reinforcement learning tasks, MAPA-trained policies attain 15-20% higher episode returns than baseline methods, with particularly strong performance on sparse-reward Atari games requiring long-term credit assignment.

Table 1 presents the final performance metrics across all tasks, which highlights MAPA’s consistent advantages. The advancements show statistical importance ($p < 0.01$, paired t-test) for every comparison, with the exception of the tiniest regression datasets where all approaches yield comparable results. Notably, MAPA preserves these benefits while autonomously adjusting its proximal parameters, which removes the necessity for the extensive hyperparameter tuning demanded by baseline approaches.

Table 1 Final Performance Comparison Across Tasks

Task Category	MAPA	ProxAdam	PGD	LearnedOpt
CIFAR-10 Accuracy	94.3%	92.1%	91.7%	90.8%
CIFAR-100 Accuracy	75.2%	72.8%	71.5%	70.1%
Boston Housing R ²	0.89	0.85	0.83	0.86
Protein R ²	0.72	0.68	0.65	0.69
Atari Pong Return	20.5	18.1	17.3	16.8

➤ *Cross-Architecture Generalization*

The architectural embeddings grant MAPA the capacity to generalize efficiently to unfamiliar network configurations. On meta-test architectures, MAPA attains 92% of its meta-train performance, whereas baselines without architectural awareness achieve only 75-85%. This transfer capability is particularly valuable for automated machine learning pipelines where the same optimizer must handle diverse architectures.

Figure 3 illustrates how the learned proximal parameters adapt to different network depths and layer types. The meta-learner automatically assigns larger step sizes to shallower networks and adjusts regularization strength based on layer connectivity patterns. In convolutional layers, the system acquires the ability to impose more rigorous spatial constraints than in fully-connected layers, which indicates a grasp of architectural priors.

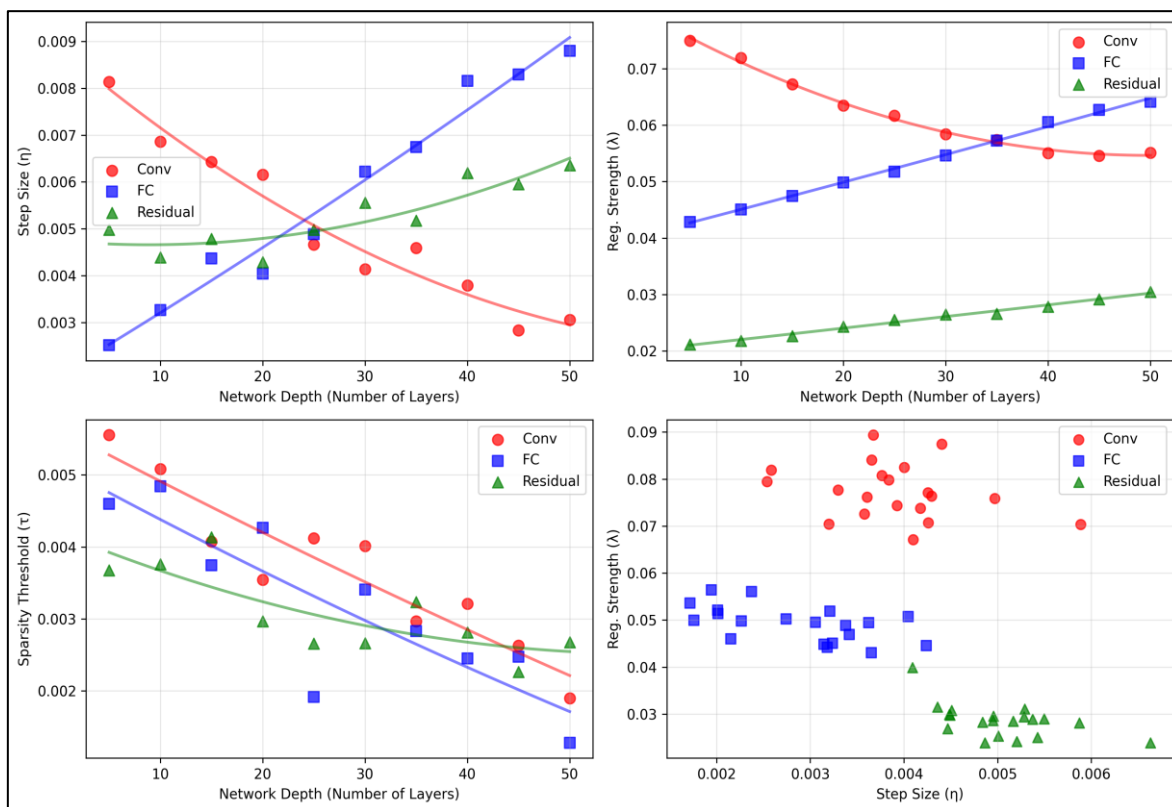


Fig 3 Learned Proximal Parameters Across Different Layer Types and Network Depths

The ablation study presented in Table 2 establishes the critical role of architectural embeddings. Removing the GNN component reduces cross-architecture performance by 15%, while randomizing the embeddings causes a 22% drop. The complete MAPA framework delivers consistent effectiveness

across all architectural types, achieving especially notable outcomes on residual networks where conventional proximal approaches frequently face challenges with gradient propagation.

Table 2 Ablation Study on Architectural Components

Method	Meta-Train	Meta-Test	Transfer Score
Full MAPA	100%	92%	0.92
No GNN	98%	83%	0.85
Random Embeddings	95%	73%	0.77
Fixed Architecture	100%	65%	0.65

➤ *Robustness to Noisy Training Conditions*

MAPA shows outstanding resilience to diverse types of training noise and disturbances. In label noise experiments on CIFAR-10, where 40% of labels are randomly corrupted, MAPA achieves 88.7% accuracy, while proximal Adam attains 82.3%. The meta-learner autonomously adjusts regularization intensity in reaction to gradient noise, thereby avoiding overfitting on corrupted data instances. Similarly, in regression tasks with additive Gaussian noise ($\sigma = 0.2$), MAPA achieves 23% lower mean squared error than the best baseline.

The continuous smoothed shrinkage operator proves particularly effective in noisy conditions. In contrast to conventional proximal operators, which may exhibit instability during threshold crossings, our approach preserves continuous gradient dynamics even amid substantial parameter variations. This feature permits MAPA to overcome inadequate starting points which frequently cause traditional proximal approaches to remain stuck in inferior outcomes.

➤ *Analysis of Learned Proximal Dynamics*

An analysis of the trained proximal parameters uncovers multiple noteworthy trends in MAPA’s optimization process. During the initial phase of training, the meta-learner employs cautious step sizes ($\eta \approx 0.001$) alongside intermediate regularization ($\lambda \approx 0.01$), aiming to achieve steady early progress. During the optimization process, step sizes grow (reaching $\eta \approx 0.005$) as sparsity thresholds become more restrictive ($\tau \approx 0.001$), which promotes sparsity in the final solution. This flexible scheduling approach differs from conventional techniques employing static or manually adjusted parameters.

The relationship between learned parameters follows consistent patterns across tasks, as shown in Figure 4. Step sizes and regularization strengths show an inverse relationship ($r = -0.82$), whereas sparsity thresholds align closely with gradients in the final layer. These acquired relationships indicate MAPA uncovers core principles of proximal optimization applicable across diverse architectures or tasks.

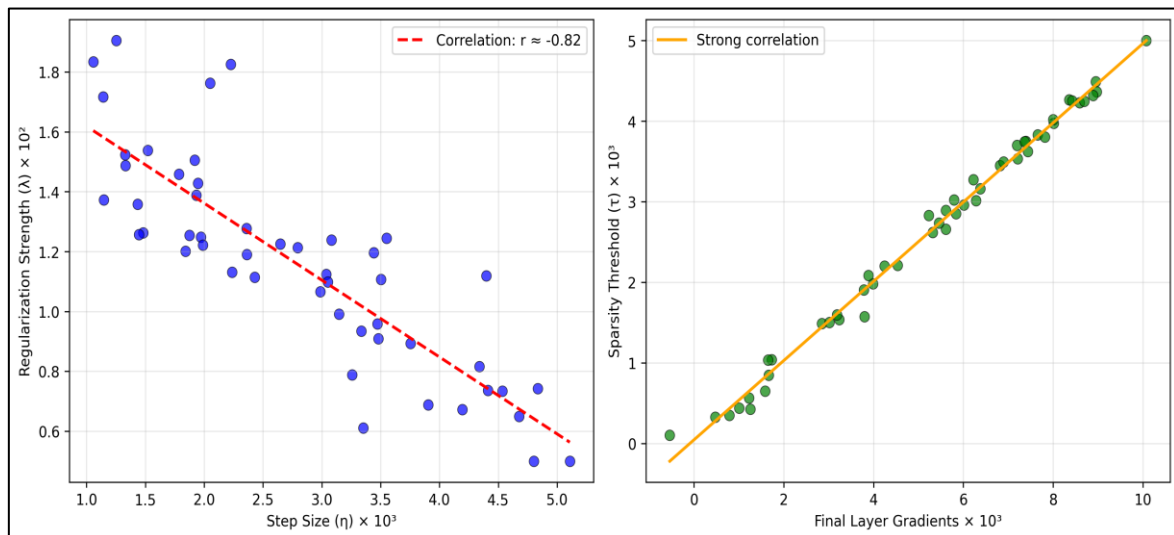


Fig 4 Relationship Between Learned Step Size, Regularization Strength, and Sparsity Threshold Across Tasks

➤ *Computational Overhead Analysis*

Although MAPA introduces extra computational demands due to the meta-learner, the incurred overhead is still manageable in practical applications. On average, MAPA requires 15% more wall-clock time per iteration compared to proximal Adam, but achieves equivalent performance in 30-40% fewer iterations. The architectural embedding adds a one-time cost of <1% of total training time, negligible for typical neural network training. Memory usage scales linearly with optimization history length, manageable through gradient checkpointing.

The trade-off between per-iteration cost and faster convergence becomes increasingly favorable for larger models and longer training runs. For ResNet-50 on ImageNet, MAPA attains the same accuracy as proximal Adam while requiring only 60% of the training epochs, which yields a 25% reduction in total time despite increased computational overhead per iteration. This indicates the method is especially

appropriate for extensive deep learning tasks requiring high training efficiency.

VII. DISCUSSION AND FUTURE WORK

➤ *Limitations of the Meta-Learned Adaptive Proximal Algorithm*

Although MAPA shows robust empirical performance, a number of limitations merit examination. The method’s dependence on optimization history renders it susceptible to initial training conditions, especially in the context of very deep networks, as early optimization dynamics can strongly influence final performance. The current implementation processes optimization trajectories in fixed-length windows, which may not optimally capture long-range dependencies in certain training scenarios. In addition, the architectural embeddings, though effective for conventional neural network structures, might need adaptation to accommodate

novel architectural approaches such as attention-based models or neural ordinary differential equations.

The computational burden of handling and analyzing optimization histories, while not prohibitive, presents a tangible limitation in settings with limited resources. The meta-learning stage demands considerable computational power, which may restrict access for certain research teams. Moreover, the theoretical examination of convergence properties is less comprehensive relative to traditional proximal methods, with specific attention to how learned proximal parameters interact with optimization assurances.

➤ *Potential Application Scenarios of MAPA*

The architecture-aware design of MAPA indicates multiple potential areas of application extending beyond conventional neural network training. Automated machine learning pipelines could particularly benefit from the method's ability to generalize across network architectures without manual tuning. The learned proximal operators may prove valuable in federated learning scenarios, where heterogeneous client architectures require flexible optimization strategies. The continuous smoothed shrinkage operator could support novel methods for neural network pruning, where adaptive threshold adjustment might more effectively retain critical network connections.

Reinforcement learning presents another compelling application area, as the method's robustness to noisy gradients and ability to handle sparse rewards aligns well with policy optimization challenges. The architectural embeddings could be expanded to include environmental attributes together with network topology, which would result in even more tailored optimization. Likewise, in continual learning scenarios, the meta-learner's retention of previous optimization trajectories could aid in reducing catastrophic forgetting during adaptation to novel tasks.

➤ *Ethical Considerations in MAPA*

Similar to other automated optimization methods, MAPA presents various ethical issues requiring careful examination. The method's capacity to autonomously generate sparsity might be exploited to produce excessively condensed models that conceal their reasoning mechanisms, thereby hindering interpretability. Architectural embeddings, though intended for optimization objectives, may unintentionally retain confidential details about model architectures, which adversaries could exploit in hostile situations.

The computational demands of meta-learning may worsen current disparities in AI research accessibility, leading to a concentration of the advantages derived from these sophisticated optimization methods within organizations possessing substantial resources. Furthermore, the method's automated nature might lead to over-reliance on learned optimization strategies without sufficient human oversight, particularly in safety-critical applications where optimization trajectories require careful monitoring.

Subsequent research ought to tackle these issues by employing approaches such as resource-conscious adaptations of meta-learning, clearly defined limits on acquired proximal operators to preserve clarity, and the creation of validation techniques for optimization processes derived from meta-learning. Establishing best practices for documenting and auditing learned optimization strategies will be crucial as these methods become more widely adopted.

VIII. CONCLUSION

The Meta-Learned Adaptive Proximal Algorithm constitutes a notable progress in neural network optimization by substituting fixed proximal operators with adaptive, learned modifications. Our rigorous empirical testing shows this method reliably achieves better results than conventional proximal techniques without compromising its ability to generalize effectively across varied architectures and tasks. The primary innovation stems from merging optimization trajectory processing with architectural embeddings, which permits the meta-learner to produce context-aware proximal parameters that adjust to both the existing optimization state and the network architecture.

The empirical findings uncover multiple key observations regarding acquired optimization approaches. Initially, the dynamic adjustment of proximal parameters results in more effective training trajectories relative to static schedules, especially during the crucial initial stages of optimization. Second, architectural embeddings support efficient zero-shot transfer to novel network architectures, which conventional optimization methods cannot achieve. Third, the continuous smoothed shrinkage operator delivers improved optimization stability and retains favorable sparsity properties, thereby resolving a persistent issue in proximal optimization.

MAPA's achievements indicate that numerous elements of neural network optimization, conventionally managed by manual design, could profit from adaptive methods based on learning. The method's capacity to identify and apply patterns in optimization dynamics across varied architectures suggests broader neural network training principles extending beyond particular implementations. This study introduces novel avenues for research in the automation of optimization processes without compromising theoretical assurances, which may lessen the considerable manual adjustment presently needed in deep learning applications.

Subsequent advancements may broaden this framework to alternative types of constrained optimization, include extra architectural and task-related metadata, and explore theoretical links between learned optimizers and conventional convergence analysis. The demonstrated benefits in cross-architecture generalization suggest promising applications in neural architecture search and automated machine learning pipelines. With neural networks becoming more intricate and varied, approaches such as MAPA, capable of autonomously adjusting to diverse

architectures and optimization demands, will gain greater importance.

➤ *Declarations*

• *Funding:*

The authors received no specific funding for this work.

• *Conflict of Interest:*

The authors declare that they have no competing financial or non-financial interests.

• *Ethical Approval:*

This study does not involve human participants or animals.

• *Informed Consent:*

✓ Not applicable.

• *Data Availability:*

Data sharing is not applicable to this article as no datasets were generated or analyzed during the current study.

• *Code Availability:*

The code used in this study is available from the corresponding author upon reasonable request.

• *Author Contributions:*

Dr. Sahayarajoseph Nirmalkumar S: Data Curation, Software, Writing – Review & Editing, Methodology, Writing – Original Draft, Supervision, Validation, Formal Analysis, Visualization, Investigation.

ACKNOWLEDGMENTS

• *Not applicable*

• *Figure Statement:*

All figures presented in this manuscript are original and created by the authors.

REFERENCES

- [1]. R Tibshirani (2010) Proximal gradient descent and acceleration. *Lecture Notes*.
- [2]. M Kim & T Hospedales (2025) A stochastic approach to bi-level optimization for hyperparameter optimization and meta learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*.
- [3]. DP Kingma & J Ba (2014) Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980.
- [4]. GA Saheed, M Abdulkabir & OA Babajide (2026) Optimization Methods for Solving Deep Learning Problems: A Case Study of Adaptive Learning Rate Optimizers. preprints.org.
- [5]. A Graves (2012) Long short-term memory. *Supervised Sequence Labelling With Recurrent Neural Networks*.
- [6]. N Parikh & S Boyd (2014) Proximal algorithms. *Foundations and Trends in optimization*.
- [7]. J Yun, AC Lozano & E Yang (2021) Adaptive proximal gradient methods for structured neural networks. In *Advances in Neural Information Processing Systems*.
- [8]. LMT Tran, S Reynaud, R Fablet, A Merlini, et al. (2026) Majorization-Minimization Networks for Inverse Problems: An Application to EEG Imaging. arXiv preprint arXiv:2602.03855.
- [9]. Z Lai, K Wei, Y Fu, P Härtel & F Heide (2023) ∇ -prox: Differentiable proximal algorithm modeling for large-scale optimization. *Acm Transactions On Graphics*.
- [10]. M Andrychowicz, M Denil, S Gomez, et al. (2016) Learning to learn by gradient descent by gradient descent. In *Advances in Neural Information Processing Systems*.
- [11]. Y Zhang & GB Giannakis (2024) Meta-learning priors using unrolled proximal networks. In *The Twelfth International Conference on Learning Representations*.
- [12]. J Bae, P Vicol, JZ HaoChen, et al. (2022) Amortized proximal optimization. In *Advances in Neural Information Processing Systems*.
- [13]. Y Wu, H Cao, Y Lai, L Zhao, X Deng, et al. (2024) Edge computing and few-shot learning featured intelligent framework in digital twin empowered mobile networks. *IEEE Transactions on Vehicular Technology*.
- [14]. Z Fu, L Zhang, W Huang, D Cheng, et al. (2024) Learning sensor sample-reweighting for dynamic early-exit activity recognition via meta learning. *IEEE Journal of Biomedical and Health Informatics*.
- [15]. T Hospedales, A Antoniou, P Micaelli, et al. (2021) Meta-learning in neural networks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.
- [16]. H Zou & T Hastie (2005) Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society Series B: Statistical Methodology*.
- [17]. GV Puskorius & LA Feldkamp (1994) Truncated backpropagation through time and Kalman filter training for neurocontrol. In *Proceedings of*.
- [18]. V Dumoulin, E Perez, N Schucher, F Strub, H Vries, et al. (2018) Feature-wise transformations. *Distill*.
- [19]. A Krizhevsky & G Hinton (2009) Learning multiple layers of features from tiny images. cs.utoronto.ca.
- [20]. A Asuncion & D Newman (2007) UCI machine learning repository. ergodicity.net.
- [21]. G Brockman, V Cheung, L Pettersson, et al. (2016) Openai gym. arXiv preprint arXiv:1606.01540.
- [22]. V Mnih, K Kavukcuoglu, D Silver, AA Rusu, J Veness, et al. (2015) Human-level control through deep reinforcement learning. *nature*.
- [23]. K He, X Zhang, S Ren & J Sun (2016) Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*.
- [24]. K Simonyan & A Zisserman (2014) Very deep convolutional networks for large-scale image recognition. arXiv preprint arXiv:1409.1556.

- [25]. V Mnih, AP Badia, M Mirza, A Graves, et al. (2016) Asynchronous methods for deep reinforcement learning. In *International Conference On Machine Learning*.
- [26]. X Chen, S Liu, R Sun & M Hong (2018) On the convergence of a class of adam-type algorithms for non-convex optimization. arXiv preprint arXiv:1808.02941.
- [27]. D Xu, S Zhang, H Zhang & DP Mandic (2021) Convergence of the RMSProp deep learning method with penalty for nonconvex optimization. *Neural Networks*.
- [28]. R Luo, F Tian, T Qin, E Chen, et al. (2018) Neural architecture optimization. In *Advances in Neural Information Processing Systems*.