

BONSAI NETWORKS: STRUCTURED PRUNING AND SPARSE TRAINING OF FOUNDATION MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

The recent trend of scaling neural networks to unprecedented sizes demands efficient structured sparsity for practical deployment, yet precisely controlling sparsity levels and patterns for hardware acceleration remains a challenge. This paper introduces **Adaptive Soft-Thresholding Algorithm (ASTRA)**, which achieves a target sparsity by adapting the group regularization strength based on affordable sparsity characterizations. We establish ASTRA’s theoretical foundations, proving the existence of stable regularizations that yield the desired sparsity. We demonstrate sublinear and linear convergence rates for both model parameters and the regularization weight in deterministic settings, and crucially, an almost sure convergence with *mean-square* rate $O(1/t)$ in the practical stochastic gradient setting. Overall, ASTRA offers a theoretically-grounded and versatile method for direct and precise control over structured sparsity, allowing pruning and fine-tuning of foundation models into Bonsai Networks: accelerator-friendly miniatures trained to match the teacher’s output while preserving downstream performance.

1 INTRODUCTION

The groundbreaking performance of modern deep neural networks, such as Large Language Models (LLMs) and Vision Language Models (VLMs), comes at a substantial computational and memory cost, both during the training phase and at inference time (Zhou et al., 2024). While the importance of scale is an ongoing debate (Goel et al., 2025), this cost is projected to continue to increase due to their widespread adoption and the belief that further breakthroughs will require unprecedented scales. This trend has spurred extensive research on model compression, with popular methods including quantization of model weights and activations during (Jacob et al., 2018) or post-training (Frantar et al., 2023).

Another approach to model compression is network pruning, inspired primarily by foundational work such as Optimal Brain Damage (OBD) (LeCun et al., 1989) and Optimal Brain Surgeon (OBS) (Hassibi et al., 1993). A paradigm shift was spurred by the Lottery Ticket Hypothesis (LTH) (Frankle & Carbin, 2019), which demonstrated that dense networks contain sparse “winning tickets” capable of matching the performance of the dense model when trained in isolation. Although LTH validated the potential of highly sparse networks, the computational expense of finding these tickets motivated the development of more efficient sparse training paradigms.

These paradigms include Static Sparse Training (SST), where a fixed sparse topology is identified before training (Lee et al., 2019; Wang et al., 2020), and Dynamic Sparse Training (DST), which adjusts connectivity during a single training run (Evci et al., 2020; Liu et al., 2020). While prominent recent successes in pruning LLMs have come from post-training techniques such as SparseGPT (Frantar & Alistarh, 2023) and Wanda (Sun et al., 2024), the significant cost of pre-training dense models motivates more efficient DST approaches. However, a critical challenge persists for DST methods: they naturally produce unstructured sparsity where individual weights are zeroed out. Such fine-grained patterns do not translate into practical speedups on modern hardware, for which structured sparsity that prunes entire neurons or channels is essential. Yet, existing methods that promote structured sparsity within the DST paradigm often rely on complex, ad-hoc heuristics and lack a unified, theoretically grounded framework.

To address this, we approach the problem from an optimization perspective. Theoretically grounded pruning and sparse training can be viewed through two dual lenses: seeking maximal sparsity for

a given model fidelity (Aghasi et al., 2017), or minimizing training loss for a pre-specified sparsity budget. Our work focuses on the latter, which is more aligned with the practical goals of dynamic training. We propose the Adaptive Soft-Thresholding Algorithm (ASTRA), a theoretically grounded method that induces structured sparsity within a dynamic training paradigm that is applicable to stochastic gradient approaches, thus eliminating the need for dense gradient computations. Our approach dynamically adjusts a group ℓ_1 regularization penalty, guiding iterates towards an equilibrium that satisfies a pre-specified target sparsity for desired structural patterns. Furthermore, the proposed theoretical framework provides a principled understanding of many existing heuristic-based approaches, thereby clarifying their underlying assumptions and limitations. By leveraging the soft-thresholding operator and established parameter grouping strategies (Mairal et al., 2011), ASTRA integrates seamlessly with existing online optimization methods, facilitating broad applicability and enabling hardware-aligned structured pruning.

CONTRIBUTIONS

- We formalize a scalar root-finding view of a target sparsity with first-order oracles.
- We provide a two time-scale scheme that tracks the sparsifying regularization in $O(1/t)$.
- A group-wise extension with closed-form prox that maps to accelerator-friendly patterns.
- Our framework is a local proximal-control of several modern heuristics, opening possibilities for extensions to structured sparsity.

2 PRELIMINARIES: PROXIMAL GRADIENT DESCENT

We consider the task of minimizing a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ under the sparsity constraint $\|\mathbf{w}\|_0 \leq \kappa$ for some integer $\kappa > 0$. A common surrogate is to employ ℓ_1 regularization:

$$\min_{\mathbf{w} \in \mathbb{R}^n} \{F(\mathbf{w}; \lambda) := f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1\}, \quad (1)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex, L -smooth, and μ -strongly convex and $\lambda > 0$ is the regularization weight. Let $\partial\|\cdot\|_1$ denote the subdifferential of the ℓ_1 -norm; an element $\mathbf{g} \in \partial\|\mathbf{w}\|_1$ satisfies $g_i = \text{sgn}(w_i)$ if $w_i \neq 0$, and $g_i \in [-1, 1]$ otherwise. A vector $\mathbf{w}(\lambda)$ is optimal if and only if it satisfies the condition:

$$0 \in \nabla f(\mathbf{w}(\lambda)) + \lambda \partial\|\mathbf{w}(\lambda)\|_1, \quad \text{i.e., } \|\nabla f(\mathbf{w}(\lambda))\|_\infty \leq \lambda. \quad (2)$$

Proximal Gradient Descent (PGD) solves the optimization problem in Equation (1) iteratively via:

$$\mathbf{w}_{t+1} = \text{prox}_{\eta_t \lambda}(\mathbf{w}_t - \eta_t \nabla f(\mathbf{w}_t)), \quad (3)$$

with a step size $\eta_t = \eta < \frac{2}{L}$ to ensure convergence, where the proximal operator associated with the ℓ_1 -norm is the soft-thresholding function, given component-wise by $[\text{prox}_\alpha(\mathbf{z})]_i = \text{sgn}(z_i) \max\{|z_i| - \alpha, 0\}$, $\forall i \in [n]$.

The convergence properties of PGD depend critically on the assumptions regarding f . When f is μ -strongly convex (for $\mu > 0$), PGD achieves a linear (geometric) convergence rate (Beck & Teboulle, 2009; Combettes & Pesquet, 2011) with a contraction factor (if $\eta_t = \frac{1}{L}$) of $\rho = \sqrt{1 - \mu/L}$, i.e.,

$$\|\mathbf{w}_{t+1} - \mathbf{w}(\lambda)\|_2 \leq \rho \|\mathbf{w}_t - \mathbf{w}(\lambda)\|_2. \quad (4)$$

If f is convex but not strongly convex (i.e., $\mu = 0$), PGD converges at a rate of $O(1/t)$ for the objective value. Accelerated methods, such as the Fast Iterative Soft-Thresholding (FISTA) (Beck & Teboulle, 2009), achieve an improved $O(1/t^2)$ convergence rate for the objective value:

$$F(\mathbf{w}_t; \lambda) - F(\mathbf{w}(\lambda); \lambda) \leq \frac{2L \|\mathbf{w}_0 - \mathbf{w}(\lambda)\|_2^2}{(t+1)^2}.$$

In many applications, only some stochastic estimates $\nabla f(\mathbf{w}_t, \boldsymbol{\xi}_t)$ of the gradient $\nabla f(\mathbf{w}_t)$ are accessible, where $\boldsymbol{\xi}_t$ is a random variable. In the strongly convex setting, under standard assumptions on the stochastic gradients (e.g., unbiasedness, and bounded variance), Stochastic Proximal Gradient

Descent (SPGD) achieves a sublinear convergence rate in expectation. This rate is typically (Rosasco et al., 2014):

$$\mathbb{E} [\|\mathbf{w}_t - \mathbf{w}(\lambda)\|_2^2] \leq \frac{C}{\mu t},$$

for some constant C that depends on problem parameters. These standard results are derived and discussed extensively in foundational texts and surveys by Beck & Teboulle (2009); Nesterov (2014); Rosasco et al. (2014); Bottou et al. (2018).

3 METHOD

Building on the surrogate objective in Equation (1), our goal is to find a minimizer \mathbf{w}_κ of f subject to the sparsity constraint $\|\mathbf{w}\|_0 \leq \kappa$, possibly under an additional structured sparsity constraint. The surrogate goal is to choose a regularization weight λ_κ that induces the (unstructured) cardinality constraint $\|\mathbf{w}(\lambda_\kappa)\|_0 \leq \kappa$ (see Section 3.4 for structured patterns).

Throughout, bold symbols denote vectors and vector-valued functions, while regular symbols denote scalars and scalar-valued functions (w_i is component i of \mathbf{w}). We denote the support of $\mathbf{w} \in \mathbb{R}^n$ by $\text{supp}(\mathbf{w}) := \{i \in [n] : w_i \neq 0\}$ and its complement by $\mathcal{Z}(\mathbf{w}) := [n] \setminus \text{supp}(\mathbf{w})$.

Denote the solution map $\mathbf{w}(\lambda) := \arg \min_{\mathbf{w}} F(\mathbf{w}; \lambda)$. We apply PGD under the following assumption, which ensures that $\mathbf{w}(\lambda)$ is single-valued.

Assumption 1. *The function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is μ -strongly convex and L -smooth with $L, \mu > 0$.*

The strong convexity in Assumption 1 yields two properties crucial to our framework (proofs in Appendix B.2):

Lemma 1 (Compactness of the solution path). *Let $\mathcal{W}_\Lambda = \{\mathbf{w}(\lambda) \mid \lambda \geq 0\}$ denote the solution path. Under Assumption 1, \mathcal{W}_Λ is compact: it is a closed set, and $\exists R > 0$ such that $\|\mathbf{w}(\lambda)\|_2 \leq R$ for all $\lambda \geq 0$.*

Lemma 2 (Lipschitz continuity in λ). *The solution map $\mathbf{w} : \mathbb{R}_+ \rightarrow \mathbb{R}^n$ is $L_{\mathbf{w}}$ -Lipschitz with respect to λ in the Euclidean norm, with $L_{\mathbf{w}} = \frac{\sqrt{n}(L+\mu)}{\mu L}$.*

The domain of f can be any closed and convex set \mathcal{W} with $\mathbf{0}_n \in \mathcal{W}$, under the assumption that \mathcal{W}_Λ is included in its interior $\text{int}(\mathcal{W})$, so that we can get the simplified first-order Karush-Kuhn-Tucker (KKT) optimality conditions for $F(\cdot; \lambda)$ as follows:

$$|\nabla_i f(\mathbf{w}(\lambda))| = \lambda, \quad i \in \text{supp}(\mathbf{w}(\lambda)); \quad |\nabla_i f(\mathbf{w}(\lambda))| \leq \lambda, \quad i \in \mathcal{Z}(\mathbf{w}(\lambda)), \quad (5)$$

where $\nabla_i f(\mathbf{w})$ denotes the i -th coordinate of $\nabla f(\mathbf{w})$.

The inequality over $\mathcal{Z}(\mathbf{w}(\lambda))$ need not be strict, as equality may occur. Define $\mathbf{w}_{-i} := (w_1, \dots, w_{i-1}, 0, w_{i+1}, \dots, w_n)$. We prove the following in Appendix B.2:

Lemma (Sparsity characterization). *Under Assumption 1, the following necessary and sufficient sparsity gauge holds:*

$$w_i(\lambda) = 0 \iff |\nabla_i f(\mathbf{w}_{-i}(\lambda))| \leq \lambda. \quad (6)$$

While strong convexity might seem restrictive, it can be enforced without altering the support selected at optimum for a fixed λ . Consider the Tikhonov-perturbed objective $h(\mathbf{w}) = f(\mathbf{w}) + \frac{\gamma}{2} \|\mathbf{w}\|_2^2$ with $\gamma > 0$. Then $\nabla_i h(\mathbf{w}_{-i}(\lambda)) = \nabla_i f(\mathbf{w}_{-i}(\lambda))$, so the gauge equation 6 is unchanged. Consequently, for any $\gamma > 0$, the solution $\mathbf{w}^h(\lambda, \gamma) := \arg \min_{\mathbf{w}} \{F(\mathbf{w}; \lambda) + \frac{\gamma}{2} \|\mathbf{w}\|_2^2\}$ has the same inactive set determined by equation 6. In Appendix C.1 we show that $\mathbf{w}^h(\lambda, \gamma) \rightarrow \arg \min_{\mathbf{w} \in \mathcal{W}(\lambda)} \|\mathbf{w}\|_2$ as $\gamma \rightarrow 0^+$, akin to the Elastic Net (Zou & Hastie, 2005), which approximates the minimum-norm element of $\mathcal{W}(\lambda)$.

3.1 STABLE REGULARIZATIONS

Once $\mathbf{w}(\lambda)$ is computed, λ is adjusted so that $|\nabla_i f(\mathbf{w}_{-i}(\lambda))| \leq \lambda$ is satisfied by at most κ components. However, perturbing λ causes the coordinates of $\mathbf{w}(\lambda)$ to change, creating a circular dependency. Theorem 1 (proof in Appendix B.3) shows that λ can be continuously adjusted to cross the

curve of any non-negative continuous function, establishing the existence of what we term *stable regularizations*.

Theorem 1 (Stable regularizations). *Let $\phi : \mathbb{R}^n \rightarrow \mathbb{R}_+$ be a continuous nonnegative function. Then there exists $\lambda \geq 0$ such that $\lambda = \phi \circ \mathbf{w}(\lambda)$. Define the compact set of ϕ -stable regularizations w.r.t to f as:*

$$\Lambda(\phi) := \{\lambda \geq 0 \mid \phi \circ \mathbf{w}(\lambda) = \lambda\}. \quad (7)$$

It should be noted that the set of ϕ -stable regularizations is not necessarily a singleton, and even for simple quadratic f , the Lasso path can be non-monotone.

Let $a_{[j]}$ denote the j -th largest element of the vector $\mathbf{a} \in \mathbb{R}^n$, and define ϕ_κ as:

$$\phi_\kappa(\mathbf{w}) := [|\nabla_1 f(\mathbf{w}_{-1})|, |\nabla_2 f(\mathbf{w}_{-2})|, \dots, |\nabla_n f(\mathbf{w}_{-n})|]_{[\kappa+1]}. \quad (8)$$

The map ϕ_κ is nonnegative and is continuous since both $\mathbf{w} \rightarrow \nabla_i f(\mathbf{w}_{-i})$ and $\mathbf{a} \mapsto |\mathbf{a}|_{[\kappa]}$ are continuous. Theorem 1 guarantees the existence of a stable regularization $\lambda \in \Lambda(\phi_\kappa)$ for which $\phi_\kappa(\mathbf{w}(\lambda)) = \lambda$. Applying Equation (6), we get:

$$\begin{aligned} w_i(\lambda) = 0 &\iff |\nabla_i f(\mathbf{w}_{-i}(\lambda))| \leq \lambda \\ &\iff |\nabla_i f(\mathbf{w}_{-i}(\lambda))| \leq \phi_\kappa(\mathbf{w}(\lambda)) \end{aligned} \quad (9)$$

By definition, $\phi_\kappa(\mathbf{w}(\lambda))$ is the κ -th largest entry among $(|\nabla_i f(\mathbf{w}_{-i}(\lambda))|)_{1 \leq i \leq n}$ and the inequality (9) is satisfied by at least $n - \kappa$ components, i.e., at most κ components are active. Thus, the task of identifying λ_κ can be reduced to finding a root of the scalar map $\lambda \mapsto \phi_\kappa(\mathbf{w}(\lambda)) - \lambda$.

In a black-box setting, computing ϕ_κ requires $|\text{supp}(\mathbf{w}(\lambda))|$ gradient evaluations, which is costly when the support is large. We therefore introduce a cheaper surrogate ψ_κ (proof in Appendix B.4):

Theorem 2 (Practical characterization). *Let $\mathbf{w} = \mathbf{w}(\lambda)$. For any $i \in [n]$, if $\exists \alpha > 0$ such that $|\nabla_i f(\mathbf{w}) - \alpha w_i| \leq \lambda$ then $w_i = 0$. Consequently, for any $\alpha \in \mathbb{R}_{>0}^n$, the function $\psi_{\kappa, \alpha} : \mathbf{w} \mapsto |\nabla f(\mathbf{w}) - \alpha \odot \mathbf{w}|_{[k+1]}$, where \odot is the element-wise product, satisfies:*

$$\lambda \in \Lambda(\psi_{\kappa, \alpha}) \implies \|\mathbf{w}(\lambda)\|_0 \leq \kappa.$$

In the remainder of the paper, the hyperparameter $\alpha \in \mathbb{R}_{>0}^n$ is fixed and dropped from notations. Note that the sparsity characterization from Theorem 2 does not require the strict convexity of f , and ψ_κ has the nice property of being Lipschitz continuous:

Lemma 3 (Lipschitzness of ψ_κ). *ψ_κ is L_ψ -Lipschitz in the Euclidean norm with $L_\psi := L + \max_i \alpha_i$.*

As shown in the proof in Appendix B.2, for all $\lambda \geq \lambda_{\max} := \|\nabla f(\mathbf{0}_n)\|_\infty$, we have $\mathbf{w}(\lambda) = \mathbf{0}_n$. Therefore, to find a stable regularization for ϕ_κ or ψ_κ , we could apply a bisection method initialized with the interval $[0, \lambda_{\max}]$: at each iteration, one solves the subproblem at the midpoint and subsequently shrinks the interval. This method yields an ϵ -accurate root with a total complexity of $O(\log^2(1/\epsilon))$, requiring $O(\log(1/\epsilon))$ linear-rate PGD subproblems. [In the next section, we introduce a flexible scheme that avoids fully computing \$\mathbf{w}\(\lambda\)\$ for each trial of \$\lambda\$ and operates effectively in the stochastic setting.](#)

3.2 ADAPTIVE SOFT-THRESHOLDING ALGORITHM (ASTRA)

Given the integer $\kappa > 0$ for the targeted sparsity level (which is then omitted from notations), set:

$$\Psi(\lambda) := \psi(\mathbf{w}(\lambda)), \quad \text{and} \quad \Phi(\lambda) := \Psi(\lambda) - \lambda.$$

We propose an Adaptive Soft-Thresholding Algorithm (ASTRA) for exact gradients:

$$\lambda_{t+1} = \Pi_{[0, \lambda_{\max}]} [(1 - \beta_t)\lambda_t + \beta_t \psi(\mathbf{w}_t)], \quad (10)$$

$$\mathbf{w}_{t+1} = \text{prox}_{\eta \lambda_{t+1}}(\mathbf{w}_t - \eta \nabla f(\mathbf{w}_t)), \quad (11)$$

where $\Pi_{\mathcal{C}}$ is the projection operator onto a set \mathcal{C} , and β_t is a suitably chosen parameter.

Using the tracking errors $\delta_t := \|\mathbf{w}_t - \mathbf{w}(\lambda_t)\|$ and $\epsilon_t := \psi(\mathbf{w}_t) - \Psi(\lambda_t)$, the λ -update becomes:

$$\lambda_{t+1} = \Pi_{[0, \lambda_{\max}]} [\lambda_t + \beta_t (\Phi(\lambda_t) + \epsilon_t)]. \quad (12)$$

This is similar to a Robbins-Monro scheme with deterministic perturbation ϵ_t . For any convergence guarantees of the proposed algorithm, the iterates are required to be bounded, which is proven in Appendix B.5 in two steps:

Lemma 4 (Decay Rate of δ_t). *Let $\rho = \sqrt{1 - \mu/L}$ and $\eta \in (0, 1/L)$. The error δ_t satisfies:*

$$\delta_{t+1} \leq \rho(1 + \beta_t L_{\mathbf{w}} L_{\psi}) \delta_t + \rho L_{\mathbf{w}} \beta_t |\Phi(\lambda_t)| \quad (13)$$

Corollary 1 (Boundedness of iterates). *Let $\bar{\beta} = \sup_t \beta_t$ such that $\rho(1 + \bar{\beta} L_{\mathbf{w}} L_{\psi}) < 1$, then the sequence $(\mathbf{w}_t)_t$ generated by the iteration from Equation (11) is bounded, and $\delta_t = O(\beta_t)$.*

Our strategy of using ℓ_1 penalty to induce sparsity comes at the price of introducing a bias in the solution. By targeting the minimal stable regularization, λ_* , this bias is minimized:

$$\lambda_* := \min_{\lambda \in \Lambda(\psi)} \lambda = \min \{ \lambda \in \mathbb{R}_+ \mid \Phi(\lambda) = 0 \}. \quad (14)$$

We will adopt the following stability assumption around λ_* :

Assumption 2 (Asymptotic stability λ_*). *We assume that λ_* is asymptotically stable for the ordinary differential equation (ODE) $\dot{\lambda} = \Phi(\lambda)$: there exist a constant $c > 0$ and a neighborhood \mathcal{N} of λ_* such that for all $\lambda \in \mathcal{N}$:*

$$(\lambda - \lambda_*)(\Phi(\lambda) - \Phi(\lambda_*)) \leq -c(\lambda - \lambda_*)^2 \quad (15)$$

Assumption 2 ensures that the equilibrium λ_* of the ODE $\dot{\lambda} = \Phi(\lambda)$ is locally attractive around λ_* . This condition implies that the curve of Φ must properly cross the λ -axis at λ_* with a negative slope, locally, and that $|\Phi(\lambda)|$ is bounded below by $c|\lambda - \lambda_*|$ in a neighborhood of λ_* . If, for instance, $\Phi(\lambda)$ only “touched” zero at λ_* from one side ($\Phi(\lambda_*) = 0$ but $\Phi(\lambda)$ does not change sign), or if $\Phi(\lambda)$ crossed zero but too “flatly” (e.g., if $\Phi'(\lambda_*) = 0$ in a differentiable case), then the stability and convergence behavior around λ_* would be dictated by higher-order, non-linear terms, which would result in considerably slower rates. Our assumption, a standard for achieving $O(1/t)$ rates (Polyak & Juditsky, 1992; Borkar, 2008), guarantees robust and sufficiently fast local convergence.

The sequence $(\beta_t)_t$ is chosen to satisfy the Robbins-Monro conditions: $\sum_{t=0}^{\infty} \beta_t = \infty$ and $\sum_{t=0}^{\infty} \beta_t^2 < \infty$, to constitute a non-uniform infinite timeline and ensure that the perturbation of the ODE is controllable, i.e., $\sum_t \beta_t \delta_t^2 < \infty$, while being small enough to not skip the stable neighborhood of λ_* . **Since we will use the initialization $\lambda_0 = 0$, the attractive neighborhood is reached in finite time, and with a suitable choice of (β_t) , we can guarantee that the iterates remain in such neighborhood.**

Theorem 3 (Sublinear Convergence Rate). *Let $\beta_t = \beta_0/(t + t_0)$ for $t \geq 0$, with $\beta_0 > 0$ and $t_0 > 0$. Assume λ_* satisfies Assumption 2. Then, for a sufficiently large β_0 , and t_0 chosen large enough so that $\beta_t < \min \left\{ 1, \frac{\rho-1}{\rho L_{\mathbf{w}} L_{\psi}} \right\}$, then the iterates converge to the solution with rates*

$$\|\lambda_t - \lambda_*\|^2 = O(1/t), \quad \|\mathbf{w}_t - \mathbf{w}(\lambda_*)\|^2 = O(1/t). \quad (16)$$

See Appendix B.6 for the proof.

ASTRA dynamics are mainly controlled by the convergence of the λ -update in $O(1/t)$, even if the \mathbf{w} -update can be exponentially faster. We show in Appendix C.2 that if Ψ is contractive around λ_* , the following intermittent schedule converges linearly to a neighborhood of λ_* of radius $O(\rho^T)$:

- $\lambda_{k+1} = (1 - \beta)\lambda_k + \beta\psi(\mathbf{w}_{kT})$, using constant step size β .
- For $0 \leq j < T - 1$: $\mathbf{w}_{kT+j+1} = \text{prox}_{\eta\lambda_{k+1}}(\mathbf{w}_{kT+j} - \eta\nabla f(\mathbf{w}_{kT+j}))$.

It should be noted that to properly implement the intermittent schedule, we also need to know some properties of f (i.e., μ and L), which is not always accessible, especially in a black-box setting. Furthermore, the necessary conditions for Ψ to be contractive are not clear; however a sufficient condition is the separability of f , i.e., $f(\mathbf{w}) = \sum_{i=1}^n f_i(w_i)$.

3.3 STOCHASTIC ADAPTIVE SOFT-THRESHOLDING ALGORITHM (SASTRA)

We now consider the case where only stochastic estimates of $\nabla f(\mathbf{w}_t)$ are available. Following the formulation from Equation (12), the stochastic extension is seen as two-timescale stochastic approximation of the ODE $(\dot{\lambda}, \dot{\mathbf{w}}) = (\Phi(\lambda), G(\mathbf{w}))$ where $G_t(\mathbf{w}) = \mathbf{w}_t - \text{prox}(\mathbf{w}_t - \eta_t f(\mathbf{w}))$.

The main challenge facing a stochastic ASTRA (SASTRA) is the nonlinearity of the order statistics, which introduces as bias in the estimation of our sparsity gauge $\psi(\mathbf{w}_t)$. We will assume that we have access to an unbiased mean map $\bar{\psi}_t$, which can be approximated practically via exponential moving average of the gradients. Let (\mathcal{F}_t) be the natural filtration and assume the following:

Assumption 3 (Unbiasedness). *The stochastic gradient g_t is unbiased with bounded variance: $\mathbb{E}[g_t | \mathcal{F}_t] = \nabla f(x_t)$ and $\mathbb{E}[\|g_t - \nabla f(x_t)\|^2 | \mathcal{F}_t] \leq \sigma^2$.*

Assumption 4 (Bounded variance). *There is a mean map $\bar{\psi}_t$ and noise ξ_t such that $\bar{\psi}_t = \psi(\mathbf{w}_t) + \xi_t$, $\mathbb{E}\xi_t | \mathcal{F}_t = 0$, $\mathbb{E}[\xi_t^2 | \mathcal{F}_t] \leq \sigma_\psi^2$.*

Consider the two-time-scale iteration

$$\lambda_{t+1} = \Pi_{[0, \lambda_{\max}]} [(1 - \beta_t)\lambda_t + \beta_t \psi_t], \quad \mathbf{w}_{t+1} = \text{prox}_{\eta_t \lambda_{t+1}}(\mathbf{w}_t - \eta_t g_t), \quad (17)$$

The Proximal Stochastic Gradient Descent (ProxSGD) step (\mathbf{w} -update) targets the slowly moving $\mathbf{w}(\lambda_t)$, for which the tracking error $E[\|\mathbf{w}_t - \mathbf{w}(\lambda_t)\|^2]$ is typically $O(1/t)$.

Theorem 4 (Convergence with log-slowed β_t). *Let $\eta_t = \frac{\eta_0}{(t+t_0)}$ and:*

$$\beta_t = \frac{\beta_0}{(t+t_0)(\log(t+t_0))^q}, \quad q > \frac{1}{2},$$

then $\sum_t \beta_t^2 / \eta_t < \infty$ and $\sum_t \eta_t^2 < \infty$. Then $(\mathbf{w}_t, \lambda_t)$ converges a.s with mean rates:

$$\mathbb{E}[\|\mathbf{w}_t - \mathbf{w}(\lambda_t)\|^2] = O(1/t), \quad \mathbb{E}[\|\lambda_t - \lambda_\star\|^2] = O(\beta_t).$$

We break down the result in Appendix B.7, which is largely inspired by Theorem 3 and relies on the Robbins-SigmundRobbins & Siegmund (1971) lemma, which requires $\sum_t \beta_t^2 / \eta_t < \infty$. To fulfill such a condition, we opt for logarithmic slowed β_t . The other option is to simply choose $\eta_t = \Theta(t^a)$ for $a \in (\frac{1}{2}, 1)$, but this implies slower convergence of $(\mathbf{w}_t)_t$.

3.4 STRUCTURED SPARSE TRAINING

Although unstructured sparsity can achieve high *theoretical* compression rates, its irregular nature offers limited practical speedups on modern hardware (Jaiswal et al., 2023). Structured sparsity overcomes this by removing entire groups of parameters (Xie et al., 2023), which makes the resulting models amenable to hardware acceleration and efficient memory access. ASTRA naturally extends to structured sparsity by applying its adaptive mechanism to groups of parameters rather than individual weights. This is achieved by employing a group ℓ_1 regularization penalty (Mairal et al., 2011) to encourage sparsity at the group level (Bach et al., 2012), using ℓ_2 -norm within a group.

In order to represent the different sparsity patterns, we introduce a Structured Sparsity Algebra (SSA) in Appendix C.3. Here, we formulate the regularization structure that is adapted to SSA.

Let $\mathcal{B} = \{\mathbf{b}_1, \dots, \mathbf{b}_{pq}\}$ be a partition of $[n]$ into disjoint pq blocks. We then define the group grid $\mathcal{G} = \{\mathbf{g}_1, \dots, \mathbf{g}_q\}$, consisting of q disjoint groups $\mathbf{g}_i = \{\mathbf{b}_{(i-1)q+1}, \dots, \mathbf{b}_{iq+p}\}$. To extend our framework to induce κ -sparse groups, we replace our unstructured ℓ_1 regularization with the gauge function

$$\Omega_{\mathcal{B}, \mathcal{G}}(\mathbf{w}) := \sum_{i=1}^q \lambda_i \sum_{\mathbf{b} \in \mathbf{g}_i} \|\mathbf{w}_{\mathbf{b}}\|,$$

then λ_i is dynamically adjusted using ASTRA to keep at most κ blocks per group. The regularization is now symmetric across groups. We consider the case with a single group ($q = 1$).

The proximal operator for such block-wise regularization is

$$\text{prox}_\alpha(\mathbf{w}_{\mathbf{b}_j}) = \max \left\{ 0, 1 - \frac{\alpha}{\|\mathbf{w}_{\mathbf{b}_j}\|_2} \right\} \mathbf{w}_{\mathbf{b}_j}.$$

The optimality conditions are then extended as $\nabla_{\mathbf{b}_j} f(\mathbf{w}^*) \in \lambda \partial \|\mathbf{w}_{\mathbf{b}_j}^*\|_2$, with $\partial \|\mathbf{w}_{\mathbf{b}}\|_2 = \frac{\mathbf{w}_{\mathbf{b}}}{\|\mathbf{w}_{\mathbf{b}}\|_2}$ for an active group, and $\partial \|\mathbf{w}_{\mathbf{b}}\|_2 \in B_{|\mathbf{b}|}(0, 1)$ for inactive ones.

The stable regularizations extend smoothly; the characterization from Equation (6) becomes:

$$\|\mathbf{w}_b(\lambda)\| = 0 \iff \|\nabla_b f(\mathbf{w}_{-b}(\lambda))\| \leq \lambda. \quad (18)$$

We write \mathcal{S}_{++}^B for the set of block-diagonal SPD matrices conformable with the block partitioning:

$$\mathcal{S}_{++}^B := \left\{ \mathbf{A} = \text{blkdiag}(\mathbf{A}_1, \dots, \mathbf{A}_m) \mid \mathbf{A}_j \in \mathbb{R}^{|b_j| \times |b_j|}, \mathbf{A}_j \succ 0 \right\}$$

Theorem 5 (Practical characterization 2). *If $\exists \mathbf{A} \in \mathcal{S}_{++}^B$ such that $\|\nabla_b f(\mathbf{w}(\lambda)) - \mathbf{A}\mathbf{w}_b(\lambda)\| \leq \lambda$ then $\|\mathbf{w}_b\| = 0$. Hence, it follows that for $\psi_\kappa^{(\mathbf{A})} : \mathbf{w} \mapsto [\mathbf{u}(\mathbf{w})]_{[\kappa+1]}$ with $\mathbf{u}_j(\mathbf{w}) = \|\nabla_{b_j} f(\mathbf{w}) - \mathbf{A}_j \mathbf{w}_{b_j}\|$ and $\mathbf{A} \in \mathcal{S}_{++}^B$, if $\lambda \in \Lambda(\psi_\kappa^{(\mathbf{A})})$ then the solution $\mathbf{w}(\lambda)$ has at most κ active blocks, i.e. $\text{nnz}(\{\|\mathbf{w}_{b_1}(\lambda)\|, \dots, \|\mathbf{w}_{b_p}(\lambda)\|\}) \leq \kappa$.*

While Theorem 5 is more general than the unstructured characterization of Theorem 2, we can still use $\mathbf{A} = \text{diag}(\boldsymbol{\alpha}) \in \mathcal{S}_{++}^B$ for $\boldsymbol{\alpha} \in \mathbb{R}_{>0}^n$ while keeping the same guarantees (proof in Appendix B.8).

4 RELATED WORK

A first-order Taylor expansion of the optimality condition in Equation (6) around \mathbf{w} gives rise to a set of saliency scores. Specifically, using H , the Hessian of f , the saliency of a weight w_i is evaluated as:

$$s_i := |\nabla_i f(\mathbf{w}_{-i})| \approx |\nabla_i f(\mathbf{w}) - H_{i,i}(\mathbf{w})w_i|, \quad (19)$$

where

Connections to Foundational Pruning Methods. Our framework recovers and generalizes several classic pruning criteria. For example, the OBD (LeCun et al., 1989) criterion, which estimates the change in the objective function δF_i from pruning a single weight w_i , is given by:

$$\delta F_i \approx - \left(\nabla_i f(\mathbf{w})w_i - \frac{1}{2} H_{i,i} w_i^2 \right). \quad (20)$$

The derivative of this saliency with respect to w_i is precisely our score, that is, $s_i = \left| \frac{d}{dw_i} \delta F_i \right|$. Crucially, our derivation of s_i does not require the strong assumptions of OBD, namely that the model is fully converged ($\nabla f(\mathbf{w}) = \mathbf{0}$) or that the Hessian is diagonal. Similarly, OBS (Hassibi et al., 1993), which analytically computes the change in objective after optimally re-adjusting all other weights, can be seen as a more complex prune-finetune heuristic. While OBS relies on inverting the Hessian to find the optimal compensation, ASTRA achieves a similar result iteratively during training based on gradients.

Unifying Modern Pruning Heuristics. ASTRA provides a theoretical grounding for modern empirically driven methods. For example, Rigging The Lottery (RigL) (Evcı et al., 2020) implements a dynamic train-prune-grow cycle where weights with the lowest magnitude are pruned and inactive weights with the largest gradient are grown. This heuristic can be interpreted as an instance of ASTRA: when the weight multiplier α_i in our framework is large such that $|\alpha_i w_i| \gg |\nabla_i f(\mathbf{w})|$, the saliency scores for active weights are dominated by their magnitude. For inactive weights ($w_i = 0$), the scores reduce to the gradient magnitudes $(|\nabla_i f(\mathbf{w})|)_i$, mirroring RigL’s grow phase.

Furthermore, our framework provides a theoretical grounding for recent one-shot pruning methods for LLMs. For the standard layer-wise reconstruction objective, $f(\mathbf{W}) = \|\mathbf{X}\mathbf{W} - \mathbf{Y}\|^2$, used by methods like Wanda (Sun et al., 2024) and SparseGPT (Frantar & Alistarh, 2023), the Hessian diagonals are $H_{j,j} = \|\mathbf{X}_j\|_2^2$. Under typical initialization conditions where $\nabla f(\mathbf{W}_0) = \mathbf{0}$, our general saliency score from Equation (19) simplifies to $s_{ij} = |W_{ij}| \|\mathbf{X}_j\|_2^2$. This is a direct variant of the empirical Wanda heuristic, which uses $|W_{ij}| \|\mathbf{X}_j\|_2$. SparseGPT, in contrast, uses a more complex iterative process based on OBS. It greedily removes weights and then analytically updates the remaining ones to compensate for the change in the objective based on the inverse of the Hessian, while ASTRA uses gradient descent to continuously adapt weights throughout the training process.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

Algorithm 1 Stochastic ASTRA Algorithm

Require: Loss function \mathcal{L} , optimizer opt , parameter $\mathbf{w} \in \mathbb{R}_n$, sparsity level κ , $\alpha \in \mathbb{R}_{>0}^n$, grids \mathcal{B}, \mathcal{G} , timescales $(\eta_t)_t, (\beta_t)_t$, EMA decay ρ .

$t \leftarrow 0$

Initialize: $\mathbf{m}_t \leftarrow \mathbf{0}_n, \lambda_t^{\mathcal{G}} \leftarrow \mathbf{0}_q$.

while not terminated **do**

 Sample batch b_t .

 Get (pseudo) gradient $g_t = opt(\mathcal{L}(b_t))$

$\mathbf{m}_t \leftarrow (1 - \rho)\mathbf{m}_{t-1} + \rho g_t$

 compute $\psi_{\mathcal{G}} = \psi_{\kappa}^{(\alpha)}(\mathbf{w}, \mathbf{m}_t)$ per group

$\lambda_t \leftarrow (1 - \beta_t)\lambda_{t-1} + \beta_t \psi_{\mathcal{G}}$

$\mathbf{w} \leftarrow \text{prox}_{\eta_t \lambda_t, \mathcal{G}}(\mathbf{w} - \eta_t g_t)$

$t \leftarrow t + 1$

end while

Algorithm 2 Bonsai Networks

Require: Neural network \mathcal{N}_{θ} parameterized by θ , SGD optimizer opt , warm-up time T_w , freeze time T_f .

for $t = 1 \dots T_w$ **do**

 Update θ using SGD optimizer.

 Update SASTRA states \mathbf{m}_t and $\lambda_t^{\mathcal{G}}$.

end for

for $t = T_w + 1 \dots T_f$ **do**

 Update θ using SGD optimizer.

 Update SASTRA states \mathbf{m}_t and λ_t .

 Soft-threshold using λ_t

end for

Freeze $\text{supp}(\theta)$ using magnitude pruning.

for $t = T_f + 1 \dots T$ **do**

 Apply SGD steps to $\text{supp}(\theta)$

end for

Figure 1: The core algorithms used to train sparse neural networks with provided sparsity Structure using Stochastic ASTRA iterations.

5 EXPERIMENTS

Our adaptive regularization approach is, at its core, a sparse-coding–inspired method. We therefore first benchmark it against its closest counterpart, Iterative Hard Thresholding (IHT). In particular, our ASTRA procedure of identifying a stable support and then updating the active weights without further thresholding; can be viewed as an IHT-style heuristic. We report results in Appendix A.1, comparing ASTRA to IHT and Optimal Brain Damage (OBD), which is equivalent to WANDA in the one-dimensional output case. The results show that ASTRA delivers consistently strong performance across settings, especially at high sparsity.

5.1 BENCHMARKING SPARSE TRAINING

We extend our approach to sparsely train a *Bonsai* ResNet-32 on CIFAR-10 and CIFAR-100 at high sparsity, following the iterations detailed in Figure 1. We use CIFAR-10/100 (50k train / 10k test, 32×32). Per-channel normalization (dataset-specific mean/std) is applied. No external data are used; augmentations are limited to random horizontal flip and random crop with 4-pixel padding. Test-time evaluation uses a single center crop (the identity at 32×32).

Our primary backbone is ResNet-32 as introduced by Wang et al. (2020). All convolutional and linear layers are candidates for pruning in the unstructured case. In the channel structured sparse training, where we sparsify at the filter level, we additionally exclude the first convolution, the final classifier.

We consider two variants:

- **Unstructured:** block size (1, 1, 1, 1) with global coupling across network weights as defined by our Structured Sparsity Algebra, excluding batch-normalization parameters. This *unstructured Bonsai ResNet-32* outperforms several established sparse-training approaches.
- **Structured:** channel-wise pruning that removes entire filters at once (excluding the input convolution and downsampling layers), trading flexibility for hardware alignment; this variant trails slightly due to the reduced degrees of freedom when selecting the sparsity mask.

We compare against various sparsity inducing techniques, including Magnitude Pruning (MP), OBD, and Lottery Ticket (LT), Gradient Signal Preservation Criterion (GraSP) (Wang et al., 2020),

Sparse Evolutionary Training(SET)(Mocanu et al., 2017), Dynamic Sparse Reparameterization (DSR)(Mostafa & Wang, 2019), and the classical Iterative Hard Thresholding (IHT),

Table 1: Test accuracy of ResNet32 on CIFAR-10 and CIFAR-100. Starred methods are reproduced results. Non-starred methods have scores as were self-reported in the literature (Wang et al., 2020).

Paradigm	Method	CIFAR-10		CIFAR-100	
		90%	95%	90%	95%
Initial	SNIP	92.59	91.01	68.89	68.89
	GraSP*	92.38	91.39	69.24	66.59
Pruning	OBD	94.17	93.29	71.46	68.73
	MP Prune	94.21	93.02	72.34	67.38
	LT	93.31	91.06	68.99	66.12
DST	DSR	92.97	91.61	69.63	68.20
	SET	92.30	90.75	69.66	67.10
	IHT*	92.29	91.60	68.64	67.03
	Bonsai	93.05	92.72	71.89	71.60
	Bonsai-Channel	91.11	90.23	69.53	67.49
Dense	Baseline*	94.32		74.44	

5.2 BENCHMARKING STRUCTURED PRUNING

For structured pruning, we test on pruning the LM head of Qwen 3 - 8B LLM (Yang et al., 2025) on wikitext data. To do so, we take 16000 tokens from wikitext data that we project into the vocabulary space. The goal then is to prune the head in a structured way to gain on memory and space. We target 75% sparsity, under the criterion of minimizing the KL divergence between the original output logits and those of our pruned layer. Our goal is to achieve a 4:16 block-sparsity with block-size $16 \times x 16$ along the K dimension. In this case, the layer weight shape is 4096×151936 , for a total of 622M parameters and 2.4M blocks. Our goal is to leverage the Dense \times 4:16 Block-Sparse kernel in Appendix C.4, which achieve a $\times 1.64$ speed up on Nvidia A100 GPU when the input is column major. Using Wanda to prune the layer (without structure) yields a mean sample KL divergence of 0.016 versus 0.08 in unstructured block- sparse and 0.10 for our 4:16 Block-Sparse structure. In contrast, the unstructured block-sparse pattern at 75% sparsity achieves $\times 0.7$ performance compared to $\times 1.64$ using our structure w.r.t to the dense matrix multiplication.

6 CONCLUSION AND DISCUSSION

We introduced ASTRA, an adaptive soft-thresholding framework that drives a model toward a *target* sparsity pattern by updating the regularization strength online. By casting target sparsity selection as a scalar root-finding problem over the regularization weight and tracking it with a two time-scale recursion, we obtain principled control of sparsity with provable rates. In the deterministic setting we established sublinear (and, under an intermittent schedule, local linear) convergence of both parameters and regularization, while in the stochastic setting we proved almost-sure convergence with mean-square rates consistent with the step-size schedule. These results provide a unifying, proximal view that clarifies the assumptions behind several heuristic sparse-training methods.

Beyond theory, we demonstrated that *Bonsai* networks trained with SASTRA achieve competitive accuracy at high sparsity on CIFAR-10/100, matching or surpassing strong pruning and DST baselines, and we provided a structured LLM case study showing that hardware-aligned 4:16 block sparsity can deliver practical speedups with controlled output drift.

Limitations and future work. Our analysis relies on convexity, while modern deep networks are nonconvex; extending the guarantees to these relaxed conditions is a natural next step. The stochastic theory assumes an unbiased surrogate for the order-statistic gauge; quantifying robustness under bias would further strengthen the framework. Finally, while our structured experiments validate a 4:16 kernel on a single LLM head, broader system-level evaluations (end-to-end latency/throughput across kernels and hardware) and scaling experiments on larger backbones are promising directions.

REFERENCES

- 486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
- Ilan Adler, Zhiyue T. Hu, and Tianyi Lin. New proximal newton-type methods for convex optimization. In *2020 59th IEEE Conference on Decision and Control (CDC)*, pp. 4828–4835, 2020. doi: 10.1109/CDC42340.2020.9304423.
- Alireza Aghasi, Afshin Abdi, Nam Nguyen, and Justin Romberg. Net-trim: Convex pruning of deep neural networks with performance guarantee. In *31st Conference on Neural Information Processing Systems (NIPS 2017)*, Long Beach, CA, 2017. URL <https://arxiv.org/abs/1611.05162>.
- Francis R. Bach, Rodolphe Jenatton, Julien Mairal, and Guillaume Obozinski. *Optimization with Sparsity-Inducing Penalties*. NOW: Foundations and Trends in Machine Learning, 2012. ISBN 9781601985101. doi: 10.1561/22000000015.
- Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009. doi: 10.1137/080716542.
- Stephen Becker, Jalal Fadili, and Peter Ochs. On quasi-newton forward–backward splitting: Proximal calculus and convergence, 2018. URL <https://arxiv.org/abs/1801.08691>.
- V.S. Borkar. *Stochastic Approximation: A Dynamical Systems Viewpoint*. Cambridge University Press, 2008. ISBN 9780521515924. URL <https://books.google.com/books?id=QLxIvgAACAAJ>.
- Léon Bottou, Frank E. Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning, 2018.
- Sébastien Bubeck. Convex optimization: Algorithms and complexity. *Found. Trends Mach. Learn.*, 8(3–4):231–357, November 2015. ISSN 1935-8237. doi: 10.1561/22000000050. URL <https://doi.org/10.1561/22000000050>.
- Patrick L. Combettes and Jean-Christophe Pesquet. Proximal splitting methods in signal processing. In Heinz H. Bauschke, Regina S. Burachik, Patrick L. Combettes, Veit Elser, D. Russell Luke, and Henry Wolkowicz (eds.), *Fixed-Point Algorithms for Inverse Problems in Science and Engineering*, pp. 185–212. Springer, New York, 2011. doi: 10.1007/978-1-4419-9569-8_10.
- Utku Evci, Trevor Gale, Jacob Menick, Pablo Samuel Castro, and Erich Elsen. Rigging the lottery: Making all tickets winners. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*, 2020.
- Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations (ICLR)*, 2019. URL <https://openreview.net/forum?id=rJl-b3RcF7>.
- Elias Frantar and Dan Alistarh. SparseGPT: Massive Language Models Can Be Accurately Pruned in One-Shot. In *International Conference on Machine Learning (ICML)*, 2023.
- Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. Gptq: Accurate post-training quantization for generative pre-trained transformers, 2023. URL <https://arxiv.org/abs/2210.17323>.
- Yash Goel, Ayan Sengupta, and Tanmoy Chakraborty. Position: Enough of scaling llms! lets focus on downscaling. In *Forty-second International Conference on Machine Learning Position Paper Track, ICML, 2025*. URL <https://openreview.net/forum?id=CyJlJgEzZs>.
- B. Hassibi, D.G. Stork, and G.J. Wolff. Optimal brain surgeon and general network pruning. In *IEEE International Conference on Neural Networks*, pp. 293–299 vol.1, 1993. doi: 10.1109/ICNN.1993.298572.
- Yasutoshi Ida, Sekitoshi Kanai, Atsutoshi Kumagai, Tomoharu Iwata, and Yasuhiro Fujiwara. Fast iterative hard thresholding methods with pruning gradient computations. In *Advances in Neural Information Processing Systems*, volume 37, pp. 52836–52857, 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/5eaa54503005d9125ad6aa3044e912d8-Paper-Conference.pdf.

- 540 Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard,
541 Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for
542 efficient integer-arithmetical-only inference. In *Proceedings of the IEEE Conference on Computer
543 Vision and Pattern Recognition (CVPR)*, June 2018.
- 544 Ayush Jaiswal, Shiwei Liu, Tianlong Liu, Mingbao Lin, and Zhangyang Wang. Dynamic
545 Sparse Training with Structured Sparsity. In *ICLR 2024 Workshop on LLM Agents*, 2023.
546 arXiv:2305.02299.
- 547 Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. In D. Touretzky
548 (ed.), *Advances in Neural Information Processing Systems*, volume 2. Morgan-Kaufmann,
549 1989. URL [https://proceedings.neurips.cc/paper_files/paper/1989/
550 file/6c9882bba1c7093bd25041881277658-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/1989/file/6c9882bba1c7093bd25041881277658-Paper.pdf).
- 551 Namhoon Lee, Thalaiyasingam Ajanthan, and Philip H. S. Torr. Snip: Single-shot network pruning
552 based on connection sensitivity. In *International Conference on Learning Representations*, ICLR,
553 2019. URL <https://openreview.net/forum?id=B1VZqjAcYX>.
- 554 Jun Liu and Ye Yuan. Almost sure convergence rates analysis and saddle avoidance of stochastic
555 gradient methods. *Journal of Machine Learning Research*, 25(271):1–40, 2024. URL [http://
556 jmlr.org/papers/v25/23-1436.html](http://jmlr.org/papers/v25/23-1436.html).
- 557 Junjie Liu, Zhe Xu, Runbin Shi, Ray C. C. Cheung, and Hayden K. H. So. Dynamic sparse training:
558 Find efficient sparse network from scratch with trainable masked layers, 2020. URL [https://
559 openreview.net/forum?id=SJ1bGJrtDB](https://openreview.net/forum?id=SJ1bGJrtDB).
- 560 Julien Mairal, Rodolphe Jenatton, Guillaume Obozinski, and Francis Bach. Convex and network
561 flow optimization for structured sparsity. *J. Mach. Learn. Res.*, 12:2681–2720, November 2011.
562 ISSN 1532-4435.
- 563 Peter Melchior, Rémy Joseph, and Fred Moolekamp. Proximal adam: Robust adaptive up-
564 date scheme for constrained optimization, 2020. URL [https://arxiv.org/abs/1910.
565 10094](https://arxiv.org/abs/1910.10094).
- 566 Decebal Constantin Mocanu, Elena Mocanu, Peter Stone, Phuong H. Nguyen, Madeleine Gibescu,
567 and Antonio Liotta. Evolutionary training of sparse artificial neural networks: A network science
568 perspective. *CoRR*, abs/1707.04780, 2017. URL <http://arxiv.org/abs/1707.04780>.
- 569 Hesham Mostafa and Xin Wang. Parameter efficient training of deep convolutional neural networks
570 by dynamic sparse reparameterization, 2019. URL [https://arxiv.org/abs/1902.
571 05967](https://arxiv.org/abs/1902.05967).
- 572 Yurii Nesterov. Introductory lectures on convex optimization - a basic course. In *Applied Optimiza-
573 tion*, 2014. URL <https://api.semanticscholar.org/CorpusID:62288331>.
- 574 Boris T. Polyak and Anatoli B. Juditsky. Acceleration of stochastic approximation by averaging.
575 *SIAM Journal on Control and Optimization*, 30(4):838–855, 1992. doi: 10.1137/0330046.
- 576 H. Robbins and D. Siegmund. A convergence theorem for non-negative almost supermartin-
577 gales and some applications. In Jagdish S. Rustagi (ed.), *Optimizing Methods in Statistics*,
578 pp. 233–257. Academic Press, 1971. ISBN 978-0-12-604550-5. doi: [https://doi.org/10.1016/
579 B978-0-12-604550-5.50015-8](https://doi.org/10.1016/B978-0-12-604550-5.50015-8). URL [https://www.sciencedirect.com/science/
580 article/pii/B9780126045505500158](https://www.sciencedirect.com/science/article/pii/B9780126045505500158).
- 581 Lorenzo Rosasco, Silvia Villa, and Bang Công Vũ. Convergence of stochastic proximal gradient
582 algorithm. *Appl. Math. Optim.*, 82:891–917, 2014. doi: 10.1007/s00245-019-09617-7.
- 583 Mingjie Sun, Zhuang Liu, Anna Bair, and J. Zico Kolter. A simple and effective pruning approach
584 for large language models. In *The Twelfth International Conference on Learning Representations*,
585 ICLR, 2024. URL <https://openreview.net/forum?id=PxoFut3dWW>.
- 586 Chaoqi Wang, Guodong Zhang, and Roger Grosse. Picking winning tickets before training by
587 preserving gradient flow. In *International Conference on Learning Representations*, ICLR, 2020.
588 URL <https://openreview.net/forum?id=SkgsACVKPH>.

594 Xiaoru Xie, Mingyu Zhu, Siyuan Lu, and Zhongfeng Wang. Efficient layer-wise n:m sparse cnn
595 accelerator with flexible spec: Sparse processing element clusters. *Micromachines*, 14, 2023.
596 URL <https://api.semanticscholar.org/CorpusID:257238483>.
597

598 An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang
599 Gao, Chengen Huang, Chenxu Lv, Chuji Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu,
600 Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin
601 Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang,
602 Le Yu, Lianghao Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui
603 Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang
604 Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yinger
605 Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan
606 Qiu. Qwen3 technical report, 2025. URL <https://arxiv.org/abs/2505.09388>.

607 Zixuan Zhou, Xuefei Ning, Ke Hong, Tianyu Fu, Jiaming Xu, Shiyao Li, Yuming Lou, Luning
608 Wang, Zhihang Yuan, Xiuhong Li, Shengen Yan, Guohao Dai, Xiao-Ping Zhang, Yuhan Dong,
609 and Yu Wang. A survey on efficient inference for large language models, 2024. URL <https://arxiv.org/abs/2404.14294>.
610

611 Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of*
612 *the Royal Statistical Society Series B: Statistical Methodology*, 67(2):301–320, 03 2005. ISSN
613 1369-7412. doi: 10.1111/j.1467-9868.2005.00503.x. URL [https://doi.org/10.1111/](https://doi.org/10.1111/j.1467-9868.2005.00503.x)
614 [j.1467-9868.2005.00503.x](https://doi.org/10.1111/j.1467-9868.2005.00503.x).
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

648 A NOTES ON REPRODUCIBILITY AND RESULTS

649 A.1 SPARSE-CODING BENCHMARK

650 Since our algorithm behaves like a pseudo-IHT, we benchmark it against IHT. Our theory requires
 651 the strong convexity assumption, but in practice that is not guaranteed. We therefore adopt the
 652 accelerated FISTA for the w -update instead of the plain ISTA iteration, and use the accelerated
 653 point to perform the λ -update. We adopt the same benchmark as Ida et al. (2024) using gisette,
 654 ledgar, real-sim and epsilon from LIBSVM. For each density level κ , the table reports the score of
 655 method \mathcal{M} as:

$$656 \text{Score}(\mathcal{M}) = \frac{MSE(\mathcal{M})}{MSE(\text{best method})} - 1$$

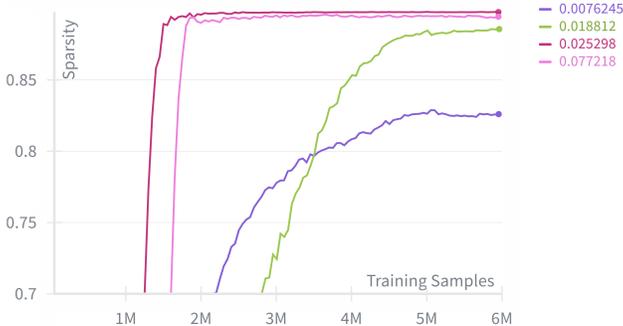
657 As such, the best method always has score 0. When two methods show the same score, the one in
 658 bold is the best one, meaning that its error is within 1% of the best approach.

Dataset	Method	κ				
		1%	10%	25%	50%	75%
Gisette	IHT	.005	0.091	.111	.298	0.84
	OBD	77.50	25.00	8.16	2.89	2.73
	ASTRA	.000	.000	.000	.000	.000
Epsilon	IHT	.060	.012	.002	.002	.001
	OBD	.327	.027	.013	.000	.000
	ASTRA	.000	.000	.000	.000	.000
Real-Sim	IHT	.084	.162	.161	.150	.050
	OBD	.325	.062	.000	.000	.000
	ASTRA	.000	.000	.018	.051	.023
Edgar	IHT	.085	.052	.068	.066	.055
	OBD	1.441	.226	.089	.000	.000
	ASTRA	.000	.000	.000	.025	.003

659 Table 2: Sparse Coding Benchmark: relative error w.r.t to the best performing method.

660 The table shows that ASTRA performs consistently better than OBD (Wanda) and IHT in high
 661 sparsity modes.

662 A.2 EFFECT OF LEARNING RATE



663 Figure 2: Reaching the target sparsity faster requires a higher learning rate, but that also results in
 664 worse performance in Neural Network training due to convergence to sub-optimal region that favors
 665 minimizing the ℓ_1 norm over the loss function

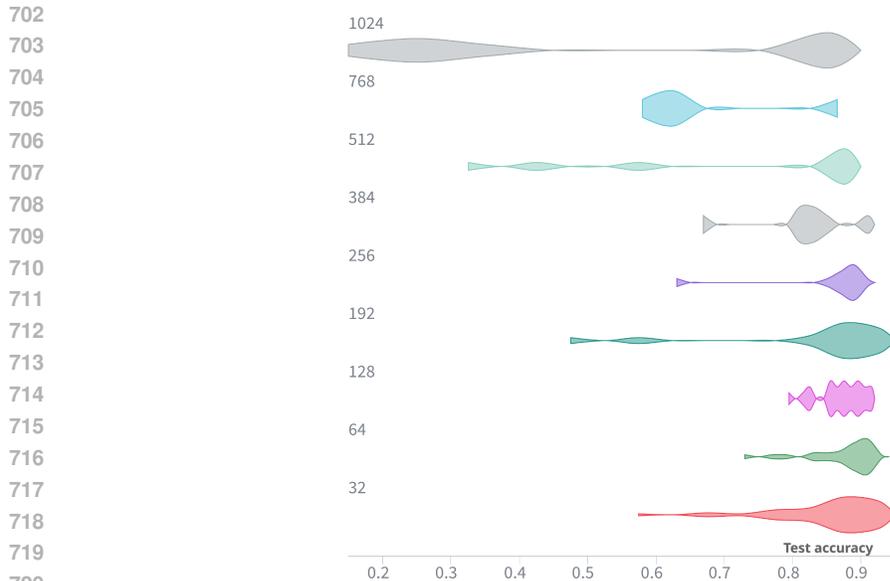


Figure 3: Distribution of the test accuracy across different batch sizes, estimated on 212 runs for arbitrary choices of hyperparameters. The batch size is the most significant factor that affects SAS-TRA performance, lower batches yield higher performance, with a -0.503 correlation with the final test accuracy.

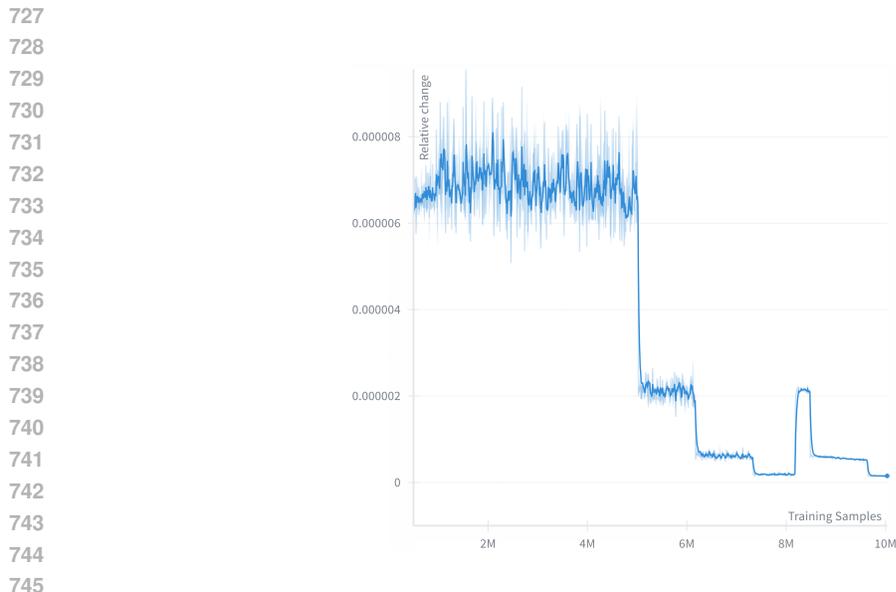


Figure 4: At training sample 8M, the support of the layer weight is frozen. In this case, The layer weight was 35% sparse before T_f and 87% sparse afterwards, yet, the relative delta norm is less than 0.0003%, showing that the components of the parameter were already hovering around 0.0.

A.3 HYPER-PARAMETERS FOR REPORTED RESULTS

These values are also provided in the configuration files of the reproducibility code. The reader can refer to the implementation for more details.

Component	Hyper-parameter	CIFAR-10	CIFAR-100
Optimizer	type	SGD	SGD
	Initial learning rate	0.0125	0.02
	Momentum	0.9437	0.977
	Weight decay	6.1×10^{-4}	2.52×10^{-5}
LR Schedule	type	Multi-step	Multi-step
	decay LR every	95 epochs	47
	decay rate	0.33	0.33
SASTRA	α_i	7.66	9.14
	β schedule	constant	constant
	β_t	0.065	0.005
	λ_{\max}	0.0015	0.0013
	EMA ρ	0.018	0.004
	T_w	8	11
	T_f	150	132
Epochs		200	150
Batch size		192	32

Table 3: Hyper-parameters for SASTRA on CIFAR Datasets with ResNet-32

B PROOFS

Notes For expositional clarity when analyzing the solution for a *fixed* λ , we assume, without loss of generality, that all entries in its support are positive, i.e., $w_i(\lambda) > 0$ for all $i \in \text{supp}(\mathbf{w}(\lambda))$. This is justified by a simple change of variables $z_i = -w_i$ for negative components $w_i(\lambda)$, which preserves the ℓ_1 -norm term ($\lambda|w_i| = \lambda|z_i|$) and the convexity of $F(\mathbf{w}; \lambda)$.

B.1 PRELIMINARIES

We will make several uses of the following lemmas (see Nesterov (2014, Thm. 2.1.12) for a proof).

Lemma B.1 (Strong monotonicity). *For a μ -strongly convex function f with L -Lipschitz continuous gradient, its gradient ∇f satisfies, for all $\mathbf{w}_1, \mathbf{w}_2$ in its domain,*

$$\begin{aligned} (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^\top (\mathbf{w}_1 - \mathbf{w}_2) &\geq \frac{\mu L}{\mu + L} \|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 + \frac{1}{\mu + L} \|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2^2 \\ &\geq \frac{\mu L}{\mu + L} \|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 \end{aligned} \quad (21)$$

In particular,

$$\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \leq \frac{1}{\mu} \|\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2)\|_2 \quad (22)$$

which holds trivially if $\mathbf{w}_1 = \mathbf{w}_2$ otherwise it follows from the Cauchy-Schwarz inequality.

Lemma B.2 (Zero solutions beyond). *If $\lambda \geq \|\nabla f(\mathbf{0}_n)\|_\infty$, then $\mathbf{w}(\lambda) = \mathbf{0}_n$.*

Proof. Let $\lambda \geq \|\nabla f(\mathbf{0}_n)\|_\infty$, then:

$$\begin{aligned} F(\mathbf{w}; \lambda) - F(\mathbf{0}_n; \lambda) &= (f(\mathbf{w}) - f(\mathbf{0}_n)) + \lambda \|\mathbf{w}\|_1 \\ &\geq \langle \nabla f(\mathbf{0}_n), \mathbf{w} \rangle + \lambda \|\mathbf{w}\|_1 && \text{using convexity of } f \\ &\geq -\|\nabla f(\mathbf{0}_n)\|_\infty \|\mathbf{w}\|_1 + \lambda \|\mathbf{w}\|_1 && \text{H\"older's inequality with } (p, q) = (1, \infty) \\ &= (\lambda - \|\nabla f(\mathbf{0}_n)\|_\infty) \|\mathbf{w}\|_1 \geq 0 \end{aligned}$$

Therefore $F(\mathbf{w}, \lambda) \geq F(\mathbf{0}_n, \lambda)$ for all $\mathbf{w} \in \mathbb{R}^n$. If, in addition, f is μ -strongly convex for $\mu > 0$, then the minimizer is unique; thus $\mathbf{w}(\lambda) = \mathbf{0}_n$. \square

Theorem B.1 (Robbins & Siegmund (1971, Thm. 1)). *Let (\mathcal{F}_t) be a filtration and let (Y_t) , (a_t) , (b_t) , (c_t) be sequences of nonnegative random variables adapted to (\mathcal{F}_t) such that*

$$\mathbb{E}[Y_{t+1} \mid \mathcal{F}_t] \leq (1 + a_t)Y_t - b_t + c_t, \quad t \geq 0,$$

with $\sum_t a_t < \infty$ and $\sum_t c_t < \infty$ almost surely. Then Y_t converges almost surely to a finite random variable Y_∞ , and $\sum_t b_t < \infty$ almost surely.

B.2 PROOF OF LEMMA 1

Lemma 1 (Compactness of the solution path). *Let $\mathcal{W}_\Lambda = \{\mathbf{w}(\lambda) \mid \lambda \geq 0\}$ denote the solution path. Under Assumption 1, \mathcal{W}_Λ is compact: it is a closed set, and $\exists R > 0$ such that $\|\mathbf{w}(\lambda)\|_2 \leq R$ for all $\lambda \geq 0$.*

Proof. We start by showing boundedness. Let $\lambda_{\max} = \|\nabla f(\mathbf{0}_n)\|_\infty$, then by inspection of the optimality conditions (and uniqueness of the solution, guaranteed by strong convexity), we see that $\mathbf{w}(\lambda) = \mathbf{0} \forall \lambda \geq \lambda_{\max}$.

Hence $\mathcal{W}_\Lambda = \{\mathbf{w}(\lambda) \mid 0 \leq \lambda \leq \lambda_{\max}\}$. For f a strongly convex function, then ∇f is (continuously) invertible which follows from eq. (22), so $\mathbf{w}(\lambda) = (\nabla f)^{-1}(-\lambda s)$ for a subgradient s . As subgradients are bounded $\|s\|_\infty \leq 1$ and $\lambda \in [0, \lambda_{\max}]$ is bounded, and since $(\nabla f)^{-1}$ is continuous, it follows \mathcal{W}_Λ is bounded.

We now show closedness, i.e., that \mathcal{W}_Λ contains all its limit points. Let (\mathbf{w}_k) be a sequence in \mathcal{W}_Λ that converges to $\bar{\mathbf{w}}$, and we wish to show $\bar{\mathbf{w}} \in \mathcal{W}_\Lambda$. If $\bar{\mathbf{w}} = \mathbf{0}_n$, then $\bar{\mathbf{w}} \in \mathcal{W}_\Lambda$ since $\mathbf{0} = \mathbf{w}(\lambda_{\max})$, so from now on assume $\bar{\mathbf{w}} \neq \mathbf{0}_n$. Thus for large enough k , $\mathbf{w}_k \neq \mathbf{0}$. Because $\mathbf{w}_k \in \mathcal{W}_\Lambda$, we can write $\mathbf{w}_k = \mathbf{w}(\lambda_k)$ for some $\lambda_k \geq 0$. For nonzero \mathbf{w}_k , the optimality conditions eq. (2) are equivalent to $\|\nabla f(\mathbf{w}_k)\|_\infty = \lambda_k$ (cf. Lemma B.2). Since (\mathbf{w}_k) converges and since ∇f is continuous (since it is Lipschitz) and norms are continuous, it follows that the sequence (λ_k) also converges; denote its limit by $\bar{\lambda}$. Because $\lambda_k \geq 0$ and $[0, \infty)$ is a closed set, then $\bar{\lambda} \geq 0$ as well. Furthermore, by the same continuity argument, we have $\|\nabla f(\bar{\mathbf{w}})\|_\infty = \bar{\lambda}$, which is the necessary and sufficient condition for optimality, so it follows $\bar{\mathbf{w}} = \mathbf{w}(\bar{\lambda})$ hence $\bar{\mathbf{w}} \in \mathcal{W}_\Lambda$.

□

Lemma 2 (Lipschitz continuity in λ). *The solution map $\mathbf{w} : \mathbb{R}_+ \rightarrow \mathbb{R}^n$ is $L_{\mathbf{w}}$ -Lipschitz with respect to λ in the Euclidean norm, with $L_{\mathbf{w}} = \frac{\sqrt{n}(L+\mu)}{\mu L}$.*

Proof. For any subgradients $s_1 \in \partial\|\mathbf{w}_1\|_1$ and $s_2 \in \partial\|\mathbf{w}_2\|_1$, we have monotonicity:

$$(\mathbf{s}_1 - \mathbf{s}_2)^\top (\mathbf{w}_1 - \mathbf{w}_2) \geq 0 \tag{23}$$

For proofs of this inequality and a broader background on convex optimization, see Bubeck (2015).

Let $\lambda_1, \lambda_2 \geq 0$, and $\mathbf{w}_1 = \mathbf{w}(\lambda_1)$ and $\mathbf{w}_2 = \mathbf{w}(\lambda_2)$. The first-order optimality conditions imply the existence of subgradients $s_1 \in \partial\|\mathbf{w}_1\|_1$ and $s_2 \in \partial\|\mathbf{w}_2\|_1$ such that:

$$\begin{aligned} \nabla f(\mathbf{w}_1) &= -\lambda_1 \mathbf{s}_1 \\ \nabla f(\mathbf{w}_2) &= -\lambda_2 \mathbf{s}_2. \end{aligned}$$

Using the strong convexity of f and Equation (21),

$$\begin{aligned} \frac{\mu L}{\mu + L} \|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 &\leq (\nabla f(\mathbf{w}_1) - \nabla f(\mathbf{w}_2))^\top (\mathbf{w}_1 - \mathbf{w}_2) \\ &= (-\lambda_1 \mathbf{s}_1 - (-\lambda_2 \mathbf{s}_2))^\top (\mathbf{w}_1 - \mathbf{w}_2) \\ &= (\lambda_2 \mathbf{s}_2 - \lambda_1 \mathbf{s}_1)^\top (\mathbf{w}_1 - \mathbf{w}_2) \\ &= -(\lambda_1 \mathbf{s}_1 - \lambda_2 \mathbf{s}_2)^\top (\mathbf{w}_1 - \mathbf{w}_2). \end{aligned}$$

Let's analyze the term $\lambda_1 \mathbf{s}_1 - \lambda_2 \mathbf{s}_2$:

$$\lambda_1 \mathbf{s}_1 - \lambda_2 \mathbf{s}_2 = \lambda_1 \mathbf{s}_1 - \lambda_1 \mathbf{s}_2 + \lambda_1 \mathbf{s}_2 - \lambda_2 \mathbf{s}_2 = \lambda_1 (\mathbf{s}_1 - \mathbf{s}_2) + (\lambda_1 - \lambda_2) \mathbf{s}_2$$

Substituting this back:

$$\begin{aligned} \frac{\mu L}{\mu + L} \|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 &\leq -[\lambda_1(\mathbf{s}_1 - \mathbf{s}_2) + (\lambda_1 - \lambda_2)\mathbf{s}_2]^\top (\mathbf{w}_1 - \mathbf{w}_2) \\ &= -\lambda_1(\mathbf{s}_1 - \mathbf{s}_2)^\top (\mathbf{w}_1 - \mathbf{w}_2) - (\lambda_1 - \lambda_2)\mathbf{s}_2^\top (\mathbf{w}_1 - \mathbf{w}_2). \end{aligned}$$

Using Equation 23, we have $(\mathbf{s}_1 - \mathbf{s}_2)^\top (\mathbf{w}_1 - \mathbf{w}_2) \geq 0$. As $\lambda_1 \geq 0$, the first term $-\lambda_1(\mathbf{s}_1 - \mathbf{s}_2)^\top (\mathbf{w}_1 - \mathbf{w}_2)$ is non-positive. Therefore:

$$\frac{\mu L}{\mu + L} \|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 \leq -(\lambda_1 - \lambda_2)\mathbf{s}_2^\top (\mathbf{w}_1 - \mathbf{w}_2) = (\lambda_1 - \lambda_2)\mathbf{s}_2^\top (\mathbf{w}_2 - \mathbf{w}_1)$$

Applying the Cauchy-Schwarz inequality:

$$\frac{\mu L}{\mu + L} \|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 \leq |\lambda_1 - \lambda_2| \cdot |\mathbf{s}_2^\top (\mathbf{w}_1 - \mathbf{w}_2)| \leq |\lambda_1 - \lambda_2| \cdot \|\mathbf{s}_2\|_2 \|\mathbf{w}_1 - \mathbf{w}_2\|_2$$

Since $\|\mathbf{s}_2\|_\infty \leq 1$ by the definition of $\partial\|\cdot\|_1$, it follows $\|\mathbf{s}_2\|_2 \leq \sqrt{n}$. Substituting, if $\mathbf{w}_1 \neq \mathbf{w}_2$, we can divide by $\|\mathbf{w}_1 - \mathbf{w}_2\|_2$ (if $\mathbf{w}_1 = \mathbf{w}_2$, the inequality holds trivially):

$$\frac{\mu L}{\mu + L} \|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 \leq \sqrt{n}|\lambda_1 - \lambda_2| \|\mathbf{w}_1 - \mathbf{w}_2\|_2.$$

Hence

$$\|\mathbf{w}_1 - \mathbf{w}_2\|_2 \leq \frac{\sqrt{n}(L + \mu)}{\mu L} |\lambda_1 - \lambda_2|.$$

This demonstrates that $\mathbf{w}(\lambda)$ is Lipschitz continuous with constant $\frac{\sqrt{n}(L + \mu)}{\mu L}$. \square

PROOF OF SPARSITY CHARACTERIZATION

Lemma (Sparsity characterization). *Under Assumption 1, the following necessary and sufficient sparsity gauge holds:*

$$w_i(\lambda) = 0 \iff |\nabla_i f(\mathbf{w}_{-i}(\lambda))| \leq \lambda. \quad (6)$$

Proof. Let $\mathbf{w} = \mathbf{w}(\lambda)$ for some $\lambda \geq 0$. We fix an index $i \in \{1, \dots, n\}$ and consider the function $h(t) = \zeta(t) + \lambda|t|$, where $\zeta(t) = f(w_1, \dots, w_{i-1}, t, w_{i+1}, \dots, w_n)$.

Since f is μ -strongly convex, $\zeta(t)$ is μ -strongly convex in the scalar variable t and the value w_i must minimize $h(t)$ with respect to t .

We remind the reader on the assumption made at the beginning of Appendix B about $w_i \geq 0$. The first-order optimality condition for minimizing $\zeta(t)$ is $0 \in \partial h(w_i)$. We have $\zeta'(t) = \frac{\partial f}{\partial w_i}(w_1, \dots, t, \dots, w_n) = \nabla_i f(w_1, \dots, t, \dots, w_n)$.

The optimality condition becomes:

$$-\zeta'(w_i) \in \lambda \partial|w_i|$$

If $w_i = 0$, then $\partial|w_i| = [-1, 1]$. The condition becomes $\zeta'_i(0) \in [-\lambda, \lambda]$, which is equivalent to $|\zeta'(0)| = |\nabla_i f(x_{-i})| \leq \lambda$. This establishes: $w_i = 0 \implies |\nabla_i f(x_{-i})| \leq \lambda$.

Now we prove the converse: assume $|\nabla_i f(x_{-i})| = |\zeta'(0)| \leq \lambda$.

Since ζ is μ -strongly convex, then its derivative is strictly increasing. Let's suppose, for contradiction, that $w_i > 0$, then $\zeta'(w_i) > \zeta'_i(0) \geq -\lambda$. The optimality condition requires that $\zeta'(w_i) = -\lambda$, which contradicts the supposition. Therefore, $w_i = 0$.

Hence, $w_i = 0$ if and only if $|\zeta'(0)| = |\nabla_i f(x_{-i})| \leq \lambda$. \square

B.3 PROOF OF THEOREM 1

Theorem 1 (Stable regularizations). *Let $\phi : \mathbb{R}^n \rightarrow \mathbb{R}_+$ be a continuous nonnegative function. Then there exists $\lambda \geq 0$ such that $\lambda = \phi \circ \mathbf{w}(\lambda)$. Define the compact set of ϕ -stable regularizations w.r.t to f as:*

$$\Lambda(\phi) := \{\lambda \geq 0 \mid \phi \circ \mathbf{w}(\lambda) = \lambda\}. \quad (7)$$

Proof. Define the function $g : \mathbb{R}_+ \rightarrow \mathbb{R}$ as $g(\lambda) = \phi(\mathbf{w}(\lambda)) - \lambda$. By Lemma 2, the mapping $\lambda \mapsto \mathbf{w}(\lambda)$ is continuous on its domain, and ϕ is continuous, hence their composition is continuous and so is g .

At $\lambda = 0$, $g(0) = \phi(\mathbf{w}(0)) - 0 \geq 0$, since ϕ is nonnegative.

Recall that $\mathbf{w}(\lambda) \in \mathcal{W}_\Lambda$ for all λ and \mathcal{W}_Λ is compact by Lemma 1. As ϕ is continuous on the compact set \mathcal{W}_Λ , it is bounded on \mathcal{W}_Λ , hence $M = \max_{\mathbf{w} \in \mathcal{W}_\Lambda} \phi(\mathbf{w}) < \infty$ is well-defined and finite. This gives the uniform bound $\phi(\mathbf{w}(\lambda)) \leq M$ for all $\lambda \geq 0$.

Consider the limit of $g(\lambda)$ as $\lambda \rightarrow \infty$: $g(\lambda) = \phi(\mathbf{w}(\lambda)) - \lambda \geq M - \lambda$. Thus, $\lim_{\lambda \rightarrow \infty} g(\lambda) = -\infty$.

To summarize, we have:

- g is continuous on \mathbb{R}_+ ,
- $g(0) \geq 0$,
- $\lim_{\lambda \rightarrow \infty} g(\lambda) = -\infty$,

thus the Intermediate Value Theorem (IVT) guarantees that there exists at least one $\lambda_* \in \mathbb{R}_+$ such that $g(\lambda_*) = 0$, i.e., a value λ_* where $\phi(\mathbf{w}(\lambda_*)) = \lambda_*$.

Since g is continuous and $\Lambda(\phi) = g_{-1}(0)$, the compactness of $\Lambda(\phi)$ follows. \square

B.4 PROOF OF THEOREM 2

Theorem 2 (Practical characterization). *Let $\mathbf{w} = \mathbf{w}(\lambda)$. For any $i \in [n]$, if $\exists \alpha > 0$ such that $|\nabla_i f(\mathbf{w}) - \alpha w_i| \leq \lambda$ then $w_i = 0$. Consequently, for any $\alpha \in \mathbb{R}_{>0}^n$, the function $\psi_{\kappa, \alpha} : \mathbf{w} \mapsto |\nabla f(\mathbf{w}) - \alpha \odot \mathbf{w}|_{[k+1]}$, where \odot is the element-wise product, satisfies:*

$$\lambda \in \Lambda(\psi_{\kappa, \alpha}) \implies \|\mathbf{w}(\lambda)\|_0 \leq \kappa.$$

Proof. Assume we have such an α , and we'll show w_i must equal 0. Consider the case $w_i > 0$, then from the optimality conditions $\nabla f_i(x) = -\lambda$ is negative, and so is $-\alpha w_i$ so

$$-\lambda - \alpha w_i = -|\lambda + \alpha w_i| = -|-\lambda - \alpha w_i| = -|\nabla_i f(\mathbf{w}) - \alpha w_i| \geq -\lambda$$

implying $-\alpha w_i \geq 0$, a contradiction. Similarly if $w_i < 0$ then $\lambda - \alpha w_i = |\lambda - \alpha w_i| = |\nabla_i f(\mathbf{w}) - \alpha w_i| \leq \lambda$ implying $-\alpha w_i \leq 0$ i.e., $w_i \geq 0$, also a contradiction.

Hence the only remaining case is $w_i = 0$.

Consequently, the $(k+1)$ -st largest magnitude in ϕ_κ or ψ_κ is of an inactive component i s.t $|\nabla_i f(\mathbf{w}_{-i})|$ and $|\nabla_i f(\mathbf{w}) - \alpha w_i|$ are equal (to $|\nabla_i f(\mathbf{w})|$). \square

B.5 PROOF OF LEMMA 4 AND COROLLARY 1

Lemma 4 (Decay Rate of δ_t). *Let $\rho = \sqrt{1 - \mu/L}$ and $\eta \in (0, 1/L)$. The error δ_t satisfies:*

$$\delta_{t+1} \leq \rho(1 + \beta_t L_{\mathbf{w}} L_\psi) \delta_t + \rho L_{\mathbf{w}} \beta_t |\Phi(\lambda_t)| \quad (13)$$

Proof. \mathbf{w}_{t+1} comes from one iteration of PGD for λ_{t+1} :

$$\begin{aligned} \delta_{t+1} &\leq \rho \|\mathbf{w}_t - \mathbf{w}(\lambda_{t+1})\| && \text{using Equation (4)} \\ &\leq \rho(\|\mathbf{w}_t - \mathbf{w}(\lambda_t)\|_2 + \|\mathbf{w}(\lambda_t) - \mathbf{w}(\lambda_{t+1})\|_2) \\ &\leq \rho(\delta_t + L_{\mathbf{w}} |\lambda_{t+1} - \lambda_t|) && \text{using Lemma 3.} \end{aligned}$$

From the iteration in Equation (10) we have:

$$\begin{aligned} |\lambda_{t+1} - \lambda_t| &\leq \beta_t |\psi(\mathbf{w}_t) - \lambda_t| && \text{inequality due to non-expansiveness of } \Pi_{[0, \lambda_{\max}]} \\ &\leq \beta_t (\epsilon_t + |\Psi(\lambda_t) - \lambda_t|) && \text{recall } \Psi(\lambda) = \psi(\mathbf{w}(\lambda)) \\ &\leq \beta_t L_\psi \delta_t + \beta_t |\Phi(\lambda_t)| \end{aligned}$$

Combining the two results:

$$\delta_{t+1} \leq \rho(1 + \beta_t L_{\mathbf{w}} L_\psi) \delta_t + \rho L_{\mathbf{w}} \beta_t |\Phi(\lambda_t)|.$$

□

Corollary 1 (Boundedness of iterates). *Let $\bar{\beta} = \sup_t \beta_t$ such that $\rho(1 + \bar{\beta} L_{\mathbf{w}} L_\psi) < 1$, then the sequence $(\mathbf{w}_t)_t$ generated by the iteration from Equation (11) is bounded, and $\delta_t = O(\beta_t)$.*

Proof. Define $\bar{\beta} := \sup_t \beta_t$ then:

$$\delta_{t+1} \leq \rho(1 + \bar{\beta} L_{\mathbf{w}} L_\psi) \delta_t + \rho L_{\mathbf{w}} \bar{\beta} |\Phi(\lambda_t)|$$

Writing $q := \rho(1 + \bar{\beta} L_{\mathbf{w}} L_\psi)$ and $C_0 = \max_{\lambda \in [0, \lambda_{\max}]} L_{\mathbf{w}} \bar{\beta} |\Phi(\lambda)|$, we get:

$$\delta_{t+1} \leq q \delta_t + C_0.$$

Hence, if $q < 1$, which is equivalent to $\bar{\beta} < \frac{\rho-1}{\rho L_{\mathbf{w}} L_\psi}$, then unrolling this recurrence from $i = 0$ to $t - 1$:

$$\delta_t \leq q^t \delta_0 + C_0 \sum_{i=1}^{t-1} q^i \leq q \delta_0 + \frac{C_0}{1-q}.$$

Thus \mathbf{w}_t is always a bounded distance away from $\mathbf{w}(\lambda_t)$, and $\mathbf{w}(\lambda_t) \in \mathcal{W}_\Lambda$ and \mathcal{W}_Λ was shown to be bounded in Lemma 1, hence all the \mathbf{w}_t stay within a bounded set. □

B.6 PROOF OF THEOREM 3

Proof. We can write the λ -update as an update towards a root of $\Phi(\lambda) = \Psi(\lambda) - \lambda$:

$$\begin{aligned} \lambda_{t+1} &= \lambda_t + \beta_t (\psi(\mathbf{w}_t) - \lambda_t) \\ &= \lambda_t + \beta_t ((\Psi(\lambda_t) - \lambda_t) + (\psi(\mathbf{w}_t) - \Psi(\lambda_t))) \\ &= \lambda_t + \beta_t \Phi(\lambda_t) + \beta_t \epsilon_t, \end{aligned}$$

where $\epsilon_t := \psi(\mathbf{w}_t) - \Psi(\lambda_t) = \psi(\mathbf{w}_t) - \psi(\mathbf{w}(\lambda_t))$.

The function ψ is L_ψ -Lipschitz (Lemma 3). Thus,

$$|\epsilon_t| \leq L_\psi \|\mathbf{w}_t - \mathbf{w}(\lambda_t)\|_2 = L_\psi \delta_t. \quad (24)$$

From Lemma 4, with $\beta_t = \beta_0/(t + t_0) = O(1/t)$, we have $\delta_t = O(\beta_t) = O(1/t)$. Therefore, $\epsilon_t = O(1/t)$.

Let $\Delta_t := |\lambda_t - \lambda_\star|$. The update for Δ_t is:

$$\begin{aligned} \Delta_{t+1} &= |\lambda_{t+1} - \lambda_\star| \\ &= |\lambda_t + \beta_t \Phi(\lambda_t) + \beta_t \epsilon_t - \lambda_\star| \\ &= |(\lambda_t - \lambda_\star) + \beta_t (\Phi(\lambda_t) - \Phi(\lambda_\star)) + \beta_t \epsilon_t|, \end{aligned}$$

since $\Phi(\lambda_\star) = \Psi(\lambda_\star) - \lambda_\star = 0$.

Consider the Lyapunov function $V_t = \frac{\Delta_t^2}{2} = \frac{1}{2}(\lambda_t - \lambda_\star)^2$.

$$\begin{aligned} V_{t+1} &= \frac{1}{2}(\Delta_t + \beta_t (\Phi(\lambda_t) - \Phi(\lambda_\star)) + \beta_t \epsilon_t)^2 \\ &= V_t + \beta_t \Delta_t [\Phi(\lambda_t) - \Phi(\lambda_\star)] + \beta_t \Delta_t \epsilon_t + \frac{\beta_t^2}{2} (\Phi(\lambda_t) - \Phi(\lambda_\star) + \epsilon_t)^2. \end{aligned}$$

Ψ is L_Ψ -Lipschitz with $L_\Psi = L_\psi L_w$. Thus, $\Phi(\lambda) = \Psi(\lambda) - \lambda$ is $L_\Phi := (L_\Psi + 1)$ -Lipschitz. Since λ_\star is an asymptotically stable equilibrium for $\dot{\lambda} = \Phi(\lambda)$, for λ_t in a neighborhood of λ_\star , we have:

$$\Delta_t(\Phi(\lambda_t) - \Phi(\lambda_\star)) = (\lambda_t - \lambda_\star)(\Phi(\lambda_t) - \Phi(\lambda_\star)) \leq -c(\lambda_t - \lambda_\star)^2 = -cV_t.$$

The term $\beta_t \Delta_t \epsilon_t$ satisfies:

$$\beta_t \Delta_t \epsilon_t \leq 2\beta_t \sqrt{V_t} L_\psi \delta_t.$$

The final term satisfies:

$$\frac{\beta_t^2}{2} (\Phi(\lambda_t) - \Phi(\lambda_\star) + \epsilon_t)^2 \leq \beta_t^2 [(\Phi(\lambda_t) - \Phi(\lambda_\star))^2 + |\epsilon_t|^2] \leq \beta_t^2 (L_\Phi |\Delta_t|^2 + |\epsilon_t|^2),$$

Since the algorithm ensures λ_t remains bounded (as established in Lemma 4), this term is $O(\beta_t^2)$.

So we get:

$$V_{t+1} \leq V_t - 2c\beta_t V_t + 2\beta_t \sqrt{V_t} (1 + L) \delta_t + O(\beta_t^2).$$

Using $\delta_t = O(\beta_t)$:

$$V_{t+1} \leq (1 - 2c\beta_t + O(\beta_t^2))V_t + O(\beta_t^2 \sqrt{V_t}) + O(\beta_t^2). \quad (25)$$

This form of recurrence is standard in stochastic approximation theory. Using Robbins-Siegmund theorem and its corollary (Liu & Yuan, 2024, Lemma 3), although we're operating in a deterministic case, for the rate $V_t = O(1/t)$ (or $\mathbb{E}[V_t] = O(1/t)$ in stochastic settings), one typically needs $2c\beta_0 > 1$. Then the terms $O(\beta_t^2 \sqrt{V_t})$ and $O(\beta_t^2)$ are $O(t^{-2} \sqrt{V_t})$ and $O(t^{-2})$, so we get:

$$\begin{aligned} V_{t+1} &\leq \left(1 - \frac{2c\beta_0}{t} + O\left(\frac{1}{t^2}\right)\right)V_t + O\left(\frac{1}{t^2} \sqrt{V_t}\right) + O\left(\frac{1}{t^2}\right) \\ &\leq \left(1 - \frac{2c\beta_0}{t}\right)V_t + O(t^{-2.5}) + O(t^{-2}) \end{aligned}$$

For a self-contained proof, we can show by induction that if β_0 is large enough and $\beta_t = \frac{\beta_0}{t+t_0}$, then $V_t = O(1/t)$: assume that $V_t \leq \frac{C}{t}$ for some $t > T$, then $V_{t+1} \leq \frac{C}{t} (1 - \frac{2c\beta_0}{t}) + O(t^{5/2}) + O(\beta_t^2)$. Using Taylor expansion $\frac{1}{t+1} = \frac{1}{t} - \frac{1}{t^2} + o(t^{-3})$, we get $V_{t+1} \leq \frac{C}{t+1}$ is equivalent to $2c\beta_0 > 1$.

If $2c\beta_0 > 1$, then for large t , $V_t = O(1/t)$. This gives $\Delta_t^2 = O(1/t)$.

Now we want to bound $\|\mathbf{w}_t - \mathbf{w}\|^2 = \|\mathbf{w}_t - \mathbf{w}(\lambda_\star)\|^2$. Using the triangle inequality:

$$\begin{aligned} \|\mathbf{w}_t - \mathbf{w}\|_2 &\leq \|\mathbf{w}_t - \mathbf{w}(\lambda_t)\|_2 + \|\mathbf{w}(\lambda_t) - \mathbf{w}(\lambda_\star)\|_2 \\ &\leq \delta_t + L_w \Delta_t, \end{aligned}$$

where L_w is the Lipschitz constant of $\mathbf{w}(\cdot)$.

We have $\delta_t = O(\beta_t) = O(1/t)$, and $\Delta_t^2 = O(1/t)$. Thus, $\|\mathbf{w}_t - \mathbf{w}\| = O(1/t) + O(1/\sqrt{t}) = O(1/\sqrt{t})$. This establishes Equation (16) and completes the proof. \square

B.7 CONVERGENCE OF SASTRA (THEOREM 4)

We provide a self-contained proof. The structure and the logic are similar to the deterministic case, although several of the arguments we use here are known results under various names, especially the notion of a two-timescale stochastic ODE from Borkar (2008).

Our goal is to prove that $\mathbb{E}[V_{t+1} | \mathcal{F}_t] \leq (1 - C\eta_t)V_t + O(\eta_t)$ for some constant $C > 0$, which would guarantee that δ_t converges a.s. to 0 with $\mathbb{E}[\delta_t^2] = O(\eta_t)$.

Lemma 5 (Tracking a moving target). *With $\eta_t \leq 1/L$, define the tracker $\Delta_t = |\lambda_{t+1} - \lambda_t|$, the Lyapunov $V_t = \frac{1}{2}\delta_t^2 = \frac{1}{2}\|\mathbf{w}_t - \mathbf{w}(\lambda_t)\|^2$ obeys:*

$$\mathbb{E}[V_{t+1} | \mathcal{F}_t] \leq \left(1 - \frac{\mu\eta_t}{2}\right)V_t + \frac{L_w^2}{\mu\eta_t} \mathbb{E}[\Delta_t^2 | \mathcal{F}_t] + \frac{\sigma^2}{2}\eta_t^2 \quad (26)$$

Proof. Let \mathcal{F}_t be the filtration up to time t . Using the non-expansiveness of the proximal operator and the contractive SGD step with $\eta_t \leq 1/L$ we get:

$$\mathbb{E}[V_{t+1} | \mathcal{F}_t] \leq \frac{1}{2}(1 - \mu\eta_t)\mathbb{E}[\|\mathbf{w}_t - \mathbf{w}(\lambda_{t+1})\|^2 | \mathcal{F}_t] + \frac{1}{2}\eta_t^2\mathbb{E}[\|g_t - \nabla f(\mathbf{w}_t)\|^2 | \mathcal{F}_t] \quad (27)$$

In order to get that $(1 - C\eta_t)$ factor, we use Young's inequality:

$$\|\mathbf{u} + \mathbf{v}\|^2 \leq (1 + \tau)\|\mathbf{u}\|^2 + (1 + \frac{1}{\tau})\|\mathbf{v}\|^2, \forall \tau > 0,$$

with $\mathbf{u} = \mathbf{w}_t - \mathbf{w}(\lambda_t)$, $\mathbf{v} = \mathbf{w}(\lambda_t) - \mathbf{w}(\lambda_{t+1})$, and $\tau = \frac{\mu\eta_t}{2(1-\mu\eta_t)}$ we get:

$$(1 - \mu\eta_t)(1 + \tau) = 1 - \frac{\mu\eta_t}{2}, \quad (1 - \mu\eta_t)\left(1 + \frac{1}{\tau}\right) \leq \frac{2}{\mu\eta_t}. \quad (28)$$

Plugging Equation (28) into Equation (27), with $\mathbf{v} \leq L_{\mathbf{w}}\Delta_t$ (from Lemma 2), we get:

$$\mathbb{E}[V_{t+1} | \mathcal{F}_t] \leq (1 - \frac{\mu\eta_t}{2})V_t + \frac{L_{\mathbf{w}}^2}{\mu\eta_t}\mathbb{E}[(\lambda_{t+1} - \lambda_t)^2 | \mathcal{F}_t] + \frac{\sigma^2}{2}\eta_t^2 \quad (29)$$

□

Lemma 6 (Drift bound for Δ_t). *Let $e_t := \lambda_t - \lambda_*$ and $\delta_t = \|\mathbf{w}_t - \mathbf{w}(\lambda_t)\|$. Then, for a universal constant C :*

$$\mathbb{E}[\Delta_t^2 | \mathcal{F}_t] \leq C\beta_t^2\left(|\Phi(\lambda_t) - \lambda_t|^2 + L_{\psi}^2\delta_t^2 + \sigma_{\psi}^2\right). \quad (30)$$

Proof. Projection is 1-Lipschitz, so $\Delta_t \leq \beta_t|\psi_t - \lambda_t|$. Write $\psi_t = \psi(\mathbf{w}_t) + \xi_t = \Phi(\lambda_t) + (\psi(\mathbf{w}_t) - \Phi(\lambda_t)) + \xi_t$, use $|\psi(\mathbf{w}_t) - \Phi(\lambda_t)| \leq L_{\psi}\delta_t$, $(a+b+c)^2 \leq 3(a^2+b^2+c^2)$, and $\mathbb{E}[\xi_t^2 | \mathcal{F}_t] \leq \sigma_{\psi}^2$. □

Lemma 7 (State mean-square drift). *Under equation 15, there exist constants $C_1, C_2 > 0$ such that*

$$\mathbb{E}[e_{t+1}^2 | \mathcal{F}_t] \leq (1 - 2c\beta_t)e_t^2 + C_1\beta_t^2\delta_t^2 + C_2\beta_t^2 \quad (31)$$

Consequently, for $\beta_t = \frac{\beta_0}{(t+t_0)(\log(t+t_0))^q}$ with any $q > 0$,

$$\mathbb{E}[e_t^2] = O(\beta_t). \quad (32)$$

Proof. From the update:

$$e_{t+1} = e_t + \beta_t(\psi_t - \lambda_t) = e_t + \beta_t(\Phi(\lambda_t) - \lambda_t) + \beta_t b_t + \beta_t \xi_t, \quad b_t := \psi(\mathbf{w}_t) - \Phi(\lambda_t), \quad |b_t| \leq L_{\psi}\delta_t.$$

Expand and take $\mathbb{E}[\cdot | \mathcal{F}_t]$:

$$\mathbb{E}[e_{t+1}^2 | \mathcal{F}_t] = e_t^2 + 2\beta_t e_t(\Phi(\lambda_t) - \lambda_t) + 2\beta_t e_t b_t + \beta_t^2 \mathbb{E}[(\Phi(\lambda_t) - \lambda_t + b_t + \xi_t)^2 | \mathcal{F}_t].$$

Use $\mathbb{E}[\xi_t] = 0$, $\mathbb{E}[\xi_t^2] \leq \sigma_{\psi}^2$, Young's inequality $2e_t b_t \leq e_t^2 + b_t^2$, and the stability from Equation (15):

$$e_t(\Phi(\lambda_t) - \lambda_t) = e_t(\Phi(\lambda_t) - \Phi(\lambda_*)) - e_t^2 \leq -ce_t^2 - e_t^2 = -(1+c)e_t^2.$$

Therefore, the $2\beta_t$ term contributes at most $-2(1+c)\beta_t e_t^2 + \beta_t e_t^2 + \beta_t b_t^2 = -(1+2c)\beta_t e_t^2 + \beta_t b_t^2$. Absorb the remaining β_t^2 square into $C_1\beta_t^2\delta_t^2 + C_2\beta_t^2$. Then a standard Robbins-Monro argument yields equation 32. □

Theorem 4 (Convergence with log-slowed β_t). *Let $\eta_t = \frac{\eta_0}{(t+t_0)}$ and:*

$$\beta_t = \frac{\beta_0}{(t+t_0)(\log(t+t_0))^q}, \quad q > \frac{1}{2},$$

then $\sum_t \beta_t^2 / \eta_t < \infty$ and $\sum_t \eta_t^2 < \infty$. Then $(\mathbf{w}_t, \lambda_t)$ converges a.s with mean rates:

$$\mathbb{E}[\|\mathbf{w}_t - \mathbf{w}(\lambda_t)\|^2] = O(1/t), \quad \mathbb{E}[|\lambda_t - \lambda_*|^2] = O(\beta_t).$$

1134 *Proof.* Combining Lemmas 5–7 and taking expectations gives

$$1135 \mathbb{E}[V_{t+1}] \leq (1 - \frac{\mu\eta_t}{2})\mathbb{E}[V_t] + \underbrace{C_1 \frac{\beta_t^2}{\eta_t}}_{:=\vartheta_t} \mathbb{E}[V_t] + \frac{C_2\beta_t^2}{\eta_t} \mathbb{E}[e_t^2] + \frac{C_3\beta_t^2}{\eta_t} + \frac{\sigma^2}{2}\eta_t^2 \quad (33)$$

1138 for constants C_1, C_2, C_3 depending on L_ψ, L_w, μ . Choose

$$1141 \eta_t = \frac{\eta_0}{t + t_0}, \quad \beta_t = \frac{\beta_0}{(t + t_0)(\log(t + t_0))^q}, \quad q > \frac{1}{2}.$$

1142 then $\sum_t \vartheta_t = \sum_t C_1 \beta_t^2 / \eta_t < \infty$ and $\sum_t \eta_t^2 < \infty$. Let $M_t := \prod_{s < t} (1 + \vartheta_s)$ (bounded since $\sum \vartheta_s < \infty$). Multiplying equation 33 by M_{t+1} and applying Robbins–Siegmund with $a_t = \mu\eta_t/2$ and Lemma 7 (so $\frac{\beta_t^2}{\eta_t} \mathbb{E}[e_t^2] = O(\beta_t^3/\eta_t)$ is summable) yields:

$$1147 \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}(\lambda_t)\|^2] = O(\eta_t), \quad \mathbb{E}[(\lambda_t - \lambda_\star)^2] = O(\beta_t).$$

1148 Hence $\mathbb{E}[\|\mathbf{w}_t - \mathbf{w}(\lambda_t)\|^2] = O(1/t)$ and $\mathbb{E}\|\lambda_t - \lambda_\star\|^2 = \tilde{O}(1/t)$ (with the log factor from β_t).

1149 In order to get the rate for $\|\mathbf{w}_t - \mathbf{w}_\star\|$, use the decomposition:

$$1151 \|\mathbf{w}_t - \mathbf{w}_\star\| \leq \|\mathbf{w}_t - \mathbf{w}(\lambda_t)\| + \|\mathbf{w}(\lambda_t) - \mathbf{w}(\lambda_\star)\| \leq O(\eta_t) + L_w |\lambda_t - \lambda_\star| \leq O(\eta_t)$$

1152 □

1153 **Why log-slowing is necessary here.** The Lyapunov recursion equation 33 requires:

$$1156 \sum_{t=1}^{\infty} \frac{\beta_t^2}{\eta_t} < \infty$$

1157 With $\eta_t = \eta_0/(t+t_0)$, choosing $\beta_t = \beta_0/((t+t_0)(\log(t+t_0))^q)$ and any $q > \frac{1}{2}$ gives $\sum_t \beta_t^2/\eta_t \sim \sum_t \frac{1}{t(\log t)^{2q}} < \infty$, while preserving Robbins–Monro ($\sum \beta_t = \infty, \sum \beta_t^2 < \infty$).

1160 *The curious reader* can also ponder the possibility of using variance-reduction techniques by increasing the batch size B_t used to estimate ψ_t so that the summability condition is $\sum_t \frac{\beta_t}{\eta_t B_t} < \infty$, which is true when both η_t and β_t are $\Theta(t)$ as long as $B_t > O(t^a)$ for some $a > 0$.

1167 B.8 PROOF OF GROUP SPARSITY CHARACTERIZATION

1168 **Theorem 5** (Practical characterization 2). *If $\exists \mathbf{A} \in \mathcal{S}_{++}^{|\mathbf{b}|}$ such that $\|\nabla_{\mathbf{b}} f(\mathbf{w}(\lambda)) - \mathbf{A} \mathbf{w}_{\mathbf{b}}(\lambda)\| \leq \lambda$ then $\|\mathbf{w}_{\mathbf{b}}\| = 0$. Hence, it follows that for $\psi_{\kappa}^{(\mathbf{A})} : \mathbf{w} \mapsto [\mathbf{u}(\mathbf{w})]_{[\kappa+1]}$ with $\mathbf{u}_j(\mathbf{w}) = \|\nabla_{\mathbf{b}_j} f(\mathbf{w}) - \mathbf{A}_j \mathbf{w}_{\mathbf{b}_j}\|$ and $\mathbf{A} \in \mathcal{S}_{++}^{\mathbf{B}}$, if $\lambda \in \Lambda(\psi_{\kappa}^{(\mathbf{A})})$ then the solution $\mathbf{w}(\lambda)$ has at most κ active blocks, i.e. $\text{nnz}(\{\|\mathbf{w}_{\mathbf{b}_1}(\lambda)\|, \dots, \|\mathbf{w}_{\mathbf{b}_p}(\lambda)\|\}) \leq \kappa$.*

1171 *Proof.* The optimality (KKT) conditions for $\mathbf{w}^* \in \arg \min_{\mathbf{w}} F_{\mathcal{G}}(\cdot; \lambda)$ are

$$1172 \mathbf{0}_n \in \nabla f(\mathbf{w}^*) + \lambda \partial \Omega_{\mathcal{B}}(\mathbf{w}^*). \quad (34)$$

1173 These specialize block-wise as follows:

$$1174 \begin{cases} \|\nabla_{\mathbf{b}} f(\mathbf{w}^*)\|_2 \leq \lambda, & \text{if } \mathbf{w}_{\mathbf{b}}^* = 0, \\ \nabla_{\mathbf{b}} f(\mathbf{w}^*) = -\lambda \frac{\mathbf{w}_{\mathbf{b}}^*}{\|\mathbf{w}_{\mathbf{b}}^*\|_2}, & \text{if } \mathbf{w}_{\mathbf{b}}^* \neq 0. \end{cases} \quad (35)$$

1175 Fix any $\mathbf{A} \in \mathcal{S}_{++}^{\mathbf{B}}$ and define the per-block nonnegative scores:

$$1176 u_j(x) := \|\nabla_{\mathbf{b}_j} f(\mathbf{w}) - \mathbf{A}_j \mathbf{w}_{\mathbf{b}_j}\|_2, \quad j = 1, \dots, p, \quad (36)$$

1177 and the order-statistic functional

$$1178 \psi_{\kappa}^{(\mathbf{A})}(\mathbf{w}) := [u_1(\mathbf{w}), \dots, u_p(\mathbf{w})]_{[\kappa+1]} \in \mathbb{R}_+. \quad (37)$$

Fix a block $\mathbf{b} = \mathbf{b}_j$. If $\mathbf{w}_\mathbf{b}^* = 0$, there is nothing to prove. Suppose instead that $\mathbf{w}_\mathbf{b}^* \neq 0$. By the block-wise KKT condition equation 35,

$$\nabla_{\mathbf{b}} f(\mathbf{w}^*) = -\lambda \frac{\mathbf{w}_\mathbf{b}^*}{\|\mathbf{w}_\mathbf{b}^*\|_2}.$$

Write $\mathbf{w}_\mathbf{b}^* = r\mathbf{v}$ with $r := \|\mathbf{w}_\mathbf{b}^*\|_2 > 0$ and $\|\mathbf{v}\|_2 = 1$. Then

$$\nabla_{\mathbf{b}} f(\mathbf{w}^*) - \mathbf{A}_j \mathbf{w}_\mathbf{b}^* = -\lambda \mathbf{v} - r \mathbf{A}_j \mathbf{v}.$$

Let $\gamma := \mathbf{v}^\top \mathbf{A}_j \mathbf{v} > 0$ (since $\mathbf{A}_i \succ 0$) and decompose $\mathbf{A}_j \mathbf{v} = \gamma \mathbf{v} + \mathbf{n}$ with $\mathbf{v} \perp \mathbf{n}$. Then

$$\|-\lambda \mathbf{v} - r \mathbf{A}_j \mathbf{v}\|_2^2 = \|(-\lambda - r\gamma)\mathbf{v} - r\mathbf{n}\|_2^2 = (\lambda + r\gamma)^2 + r^2 \|\mathbf{n}\|_2^2 \geq (\lambda + r\gamma)^2 > \lambda^2,$$

because $r\gamma > 0$. Hence $\|\nabla_{\mathbf{b}} f(\mathbf{w}^*) - \mathbf{A}_j \mathbf{w}_\mathbf{b}^*\|_2 > \lambda$. Take the contrapositive to yield the claim. \square

C EXTENDED NOTES

C.1 RELAXING THE STRONG CONVEXITY ASSUMPTION

Theorem C.1. Fix $\lambda > 0$ and let $\gamma > 0$. Define $\mathbf{w}_\gamma := \arg \min_{\mathbf{w}} F_\gamma(\mathbf{w})$ where the γ -strongly convex Tikhonov-perturbed objective is:

$$F_\gamma(\mathbf{w}) := F(\mathbf{w}, \lambda) + \frac{\gamma}{2} \|\mathbf{w}\|^2,$$

Let $\mathbf{w}^\dagger = \arg \min_{\mathbf{w} \in \mathbf{w}(\lambda)} \|\mathbf{w}\|$ be the minimum-Euclidean-norm solution. Then, for every $\gamma > 0$ the following holds:

- **Norm shrinkage.** $\|\mathbf{w}_\gamma\| \leq \|\mathbf{w}^\dagger\|$.
- **Minimal-norm selection.** $\mathbf{w}_\gamma \rightarrow \mathbf{w}^\dagger$ as $\gamma \rightarrow 0^+$.
- **Exact value-gap identity and bound.**

$$0 \leq F(\mathbf{w}_\gamma, \lambda) - F(\mathbf{w}^\dagger, \lambda) \leq \frac{\gamma}{2} (\|\mathbf{w}^\dagger\|^2 - \|\mathbf{w}_\gamma\|^2) \quad (38)$$

Proof. We have the optimality conditions:

$$\mathbf{0}_n \in \partial F(\mathbf{w}_\gamma, \lambda) + \gamma \mathbf{w}_\gamma, \quad \text{and} \quad \mathbf{0}_n \in \partial F(\mathbf{w}^\dagger, \lambda). \quad (39)$$

Pick $g_\gamma \in \partial F(\mathbf{w}_\gamma, \lambda)$ and $g^* \in \partial F(\mathbf{w}^\dagger, \lambda)$ with $g_\gamma = -\gamma \mathbf{w}_\gamma$ and $g^* = 0$. Maximal monotonicity from Equation (23) yields:

$$\langle g_\gamma - g^*, \mathbf{w}_\gamma - \mathbf{w}^\dagger \rangle \geq 0 \implies -\gamma \langle \mathbf{w}_\gamma, \mathbf{w}_\gamma - \mathbf{w}^\dagger \rangle \geq 0.$$

Hence, $\langle \mathbf{w}_\gamma, \mathbf{w}^\dagger \rangle \geq \|\mathbf{w}_\gamma\|^2$, and by Cauchy–Schwarz $\|\mathbf{w}_\gamma\| \leq \|\mathbf{w}^\dagger\|$.

For any \mathbf{w} ,

$$F_\gamma(\mathbf{w}) \geq F_\gamma(\mathbf{w}_\gamma).$$

With $\mathbf{w} = \mathbf{w}^\dagger$:

$$F(\mathbf{w}_\gamma, \lambda) - F(\mathbf{w}^\dagger, \lambda) \leq \frac{\gamma}{2} (\|\mathbf{w}^\dagger\|^2 - \|\mathbf{w}_\gamma\|^2) \quad (40)$$

we also have:

$$F(\mathbf{w}_\gamma, \lambda) - F(\mathbf{w}^\dagger, \lambda) \geq 0 \quad (41)$$

From Equations (40) and (41) we get Equation (38).

Using the strict convexity of F_γ , an equivalent polarized form yields:

$$F(\mathbf{w}_\gamma, \lambda) - F(\mathbf{w}^\dagger, \lambda) \leq \gamma \langle \mathbf{w}_\gamma, \mathbf{w}^\dagger - \mathbf{w}_\gamma \rangle,$$

Let $\gamma_k \downarrow 0$ and $\mathbf{w}_{\gamma_k} \rightarrow \bar{\mathbf{w}}$ along a convergent subsequence (existence by compactness). Choose $g_{\gamma_k} \in \partial F(\mathbf{w}_{\gamma_k}, \lambda)$ with $g_{\gamma_k} = -\gamma_k \mathbf{w}_{\gamma_k} \rightarrow 0$. Closedness of the graph of ∂F yields $0 \in \partial F(\bar{\mathbf{w}})$, so $\bar{\mathbf{w}} \in \mathbf{w}\lambda$. Moreover, the polarized inequality gives:

$$\langle \mathbf{w}_{\gamma_k}, \mathbf{w}^\dagger - \mathbf{w}_{\gamma_k} \rangle \geq 0 \implies \langle \bar{\mathbf{w}}, \mathbf{w}^\dagger - \bar{\mathbf{w}} \rangle \geq 0.$$

This characterizes the Euclidean projection of $\mathbf{0}_n$ onto the closed convex set $\mathbf{w}(\lambda)$; hence $\bar{\mathbf{w}} = P_{\mathbf{w}\lambda}(0) = \mathbf{w}^\dagger$. As every subsequence has the same limit \mathbf{w}^\dagger , the full sequence satisfies $\mathbf{w}_\gamma \rightarrow \mathbf{w}^\dagger$ as $\gamma \downarrow 0$. \square

C.2 DETERMINISTIC LINEAR RATE VIA INTERMITTENT SCHEME

Theorem 6 (Linear Convergence Rate (Deterministic)). *Let T be a fixed interval length, and assume Ψ is contractive in a neighborhood of λ_* , i.e., there exists a constant $\gamma \in [0, 1)$ such that $|\Psi(\lambda) - \Psi(\lambda_*)| \leq \gamma|\lambda - \lambda_*|$ for all λ in this neighborhood. Consider the sequence $(\lambda_k)_k$ generated by the intermittent schedule $\beta \in (0, 1)$ and let $\lambda_k = \lambda_{kT}$. Then the sequence $(\lambda_k)_k$ converges linearly to a neighborhood of λ_* with a radius of $O(\rho^T)$. Specifically, the following bound holds:*

$$|\lambda_{k+1} - \lambda_*| \leq [1 - \beta(1 - \gamma)]|\lambda_k - \lambda_*| + O\left(\frac{\rho^T}{1 - \sigma}\right). \quad (42)$$

Proof. We analyze the intermittent schedule with constant β . Let $\lambda_k = \lambda_{kT}$ and $\mathbf{w}_k = \mathbf{w}_{kT}$.

In the interval $t \in [kT, (k+1)T - 1]$, $\lambda = \lambda_k$ is fixed. The \mathbf{w} -updates converge geometrically towards $\mathbf{w}(\lambda_k)$: $\|\mathbf{w}_{(k+1)T} - \mathbf{w}(\lambda_k)\| \leq \rho^T \|\mathbf{w}_{kT} - \mathbf{w}(\lambda_k)\|$. Let D be an upper bound on $\|\mathbf{w}_t - \mathbf{w}(\lambda_k)\|$.

The error term in the λ -update at step $k+1$ (using $\mathbf{w}_{(k+1)T}$) is $\epsilon_k = \psi_{\kappa}(\mathbf{w}_{(k+1)T}) - \Psi(\lambda_k)$. We have $|\epsilon_k| = |\psi_{\kappa}(\mathbf{w}_{(k+1)T}) - \psi_{\kappa}(\mathbf{w}(\lambda_k))| \leq L_{\psi} \|\mathbf{w}_{(k+1)T} - \mathbf{w}(\lambda_k)\| \leq L_{\psi} \rho^T \|\mathbf{w}_{kT} - \mathbf{w}(\lambda_k)\| \leq L_{\psi} D \rho^T$. This error can be made small by choosing T large.

For the λ -update: $\lambda_{k+1} = \lambda_k + \beta(\Psi(\lambda_k) - \lambda_k + \epsilon_k)$. We have:

$$\begin{aligned} \Delta_{k+1} &= \Delta_k + \beta(\Psi(\lambda_k) - \Psi(\lambda_*)) - \beta(\lambda_k - \lambda_*) + \beta\epsilon_k \\ &= (1 - \beta)\Delta_k + \beta(\Psi(\lambda_k) - \Psi(\lambda_*)) + \beta\epsilon_k \end{aligned}$$

Using the contractivity assumption, $|\Psi(\lambda_k) - \Psi(\lambda_*)| \leq \gamma|\Delta_k|$:

$$\begin{aligned} |\Delta_{k+1}| &\leq |1 - \beta||\Delta_k| + \beta\gamma|\Delta_k| + \beta|\epsilon_k| \\ &\leq \underbrace{(|1 - \beta| + \beta\gamma)}_{\sigma} |\Delta_k| + \beta L_{\psi} D \rho^T \end{aligned}$$

For $0 < \beta \leq 1$, this requires $\sigma = 1 - \beta + \beta\gamma < 1$, which simplifies to $\beta(1 - \gamma) > 0$ (true since $\beta > 0, \gamma < 1$). Thus, $\sigma = 1 - \beta(1 - \gamma) < 1$.

The recurrence $|\Delta_{k+1}| \leq \sigma|\Delta_k| + C\rho^T$ (where $C = \beta L_{\psi} D$) shows linear convergence to a neighborhood of λ_* . Specifically:

$$|\Delta_k| \leq \sigma^k |\Delta_0| + \frac{C\rho^T}{1 - \sigma} \quad (43)$$

The error converges linearly towards a ball of radius $O(\rho^T)$. For convergence to λ_* , T must be large enough relative to the desired precision. The convergence rate is $\sigma = 1 - \beta(1 - \gamma)$ per k (per T steps).

For the sequence \mathbf{w}_t :

$$\begin{aligned} \|\mathbf{w}_{kT} - \mathbf{w}\| &\leq \|\mathbf{w}_{kT} - \mathbf{w}(\lambda_k)\| + \|\mathbf{w}(\lambda_k) - \mathbf{w}(\lambda_*)\| \\ &\leq \rho^T D' + L_h |\lambda_k - \lambda_*| \end{aligned}$$

(where D' bounds $\|\mathbf{w}_{kT} - \mathbf{w}(\lambda_k)\|$, which might depend on k). Both terms decrease towards zero (the first geometrically in T , the second linearly in k). \square

Note on intermittent schedule: If β_t is zero except at $t = kT$, the tracking argument $\|\mathbf{w}_t - \mathbf{w}(\lambda_t)\| \rightarrow 0$ is strengthened. Between λ -updates, \mathbf{w}_t converges geometrically towards the fixed $\mathbf{w}(\lambda_{kT})$. Thus, at times $t = kT$, the error $\|\mathbf{w}_{kT} - \mathbf{w}(\lambda_{(k-1)T})\|$ (or $\|\mathbf{w}_{kT} - \mathbf{w}(\lambda_{kT})\|$ depending on timing) can be made very small by choosing T large, making the noise term ϵ_{kT} in the λ -update small. The overall ODE analysis structure remains applicable.

1296 C.3 SPARSITY ALGEBRA

1297
1298 A novel sparsity algebra is proposed to provide a structured and unified framework for applying di-
1299 verse sparsity constraints to (neural network) parameters. This algebra is composed of three funda-
1300 mental components: the *sparsity block*, the *sparsity group*, and *Group coupling*. These components
1301 can be combined to express a wide range of structured sparsity patterns, from fine-grained N:M
1302 sparsity to coarse-grained channel and layer-level pruning.

1303 C.3.1 FORMULATION

1304
1305 Let a parameter tensor be denoted by $\mathbf{W} \in \mathbb{R}^{d_1 \times \dots \times d_m}$. The proposed sparsity algebra introduces a
1306 hierarchical structure onto this tensor.

1307
1308 **Sparsity Block** The fundamental unit of sparsity is the *sparsity block*. The tensor \mathbf{W} is parti-
1309 tioned into non-overlapping, equally-sized blocks. The block dimensions are defined by a tuple $b =$
1310 (b_1, \dots, b_m) , where each tensor dimension d_i must be divisible by the corresponding block dimen-
1311 sion b_i . This creates a grid of blocks \mathcal{B} with dimensions $B = (\frac{d_1}{b_1}, \dots, \frac{d_m}{b_m})$. Each block $\mathbf{W}_{i_1, \dots, i_m} \in$
1312 $\mathbb{R}^{b_1 \times \dots \times b_m}$ consists of the elements of \mathbf{W} indexed by $[i_1 b_1 : (i_1 + 1)b_1, \dots, i_m b_m : (i_m + 1)b_m]$.
1313 The pruning decision is applied at the block level, meaning all weights within a single block are
1314 pruned or kept together.

1315
1316 **Sparsity Group** A *sparsity group* is a collection of sparsity blocks. The grid of blocks \mathcal{B} is further
1317 partitioned into groups defined by a tuple $g = (g_1, \dots, g_m)$, where for each axis i , the block grid
1318 dimension $B_i = d_i/b_i$ must be divisible by g_i . This results in a grid of groups \mathcal{G} with dimensions
1319 $(\frac{B_1}{g_1}, \dots, \frac{B_m}{g_m})$. Each group contains $|g| := \prod_{i=1}^m g_i$ sparsity blocks. The sparsity level κ is applied
1320 within each group, resulting in a sparsity of $1 - \frac{\kappa}{|g|}$.

1321
1322 **Group Coupling** To enforce sparsity constraints jointly across multiple parameter tensors, the
1323 concept of *Group coupling* is introduced. A set of parameter tensors $\{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(k)}\}$ can be cou-
1324 pled, allowing a global sparsity constraint to be applied across their corresponding groups. This is
1325 particularly useful for maintaining a global level of sparsity or for co-pruning related parameters.
1326 The coupling mechanism aligns the group grids of the different tensors, potentially with permuta-
1327 tions of their axes, and applies a unified sparsity criterion across them.

1328 C.3.2 SPARSITY EXAMPLES

1329
1330 The flexibility of this algebra allows for the formulation of various structured sparsity patterns.

1331
1332 **2:4 Sparsity** A common hardware-accelerated pattern is 2:4 sparsity, where two out of every four
1333 elements must be zero. This can be formulated using our algebra on a weight matrix $\mathbf{W} \in \mathbb{R}^{m \times n}$.
1334 We first define a block size of $b = (1, 1)$ to consider individual weights. We then group these blocks
1335 into 1×4 contiguous sets by defining a group of size $g = (1, 4)$. Within each group, we enforce
1336 that only $\kappa = 2$ elements (blocks) are kept, thereby achieving 2:4 sparsity along the rows.

1337
1338 **κ -Block Row Sparsity** To enforce a fixed number of non-zero blocks in each row of a block-
1339 partitioned matrix, we can define a group that encompasses an entire row of blocks. For a matrix
1340 with block grid dimensions (B_1, B_2) , we set the group size to $g = (1, B_2)$. Within each of these
1341 row-groups, we can then keep κ of blocks per row block (group).

1342
1343 **Channel-wise Sparsity** For a 2D convolutional layer, the weight tensor is typically of shape $\mathbf{W} \in$
1344 $\mathbb{R}^{c_{out} \times c_{in} \times k_h \times k_w}$, where c_{out} and c_{in} are the number of output and input channels, and k_h, k_w are
1345 the kernel height and width. To achieve channel-wise sparsity where each output channel connects to
1346 only κ input channels, define the blocks to be the filter weights corresponding to a single input-output
1347 channel connection, i.e., $b = (1, 1, k_h, k_w)$. The resulting block grid has dimensions $(c_{out}, c_{in}, 1, 1)$.
1348 We then group all blocks corresponding to a single output channel by setting the group size to
1349 $g = (1, c_{in}, 1, 1)$. For each of these groups, we enforce that only κ blocks (i.e., connections to κ
input channels) are kept.

1350 **Global Thresholding with Coupling** A significant challenge in pruning is determining the appropriate sparsity level for each layer. Global thresholding applies a single criterion across the entire network, imposing a uniform sparsity pressure without requiring per-layer hyperparameter tuning. Our algebra provides a structured mechanism for this via group coupling.

1354 For instance, consider an MLP with two layers, $\mathbf{W}_1 \in \mathbb{R}^{1024 \times 2048}$ and $\mathbf{W}_2 \in \mathbb{R}^{2048 \times 1024}$. We can prune entire neurons in the hidden layer by coupling the columns of \mathbf{W}_1 with the rows of \mathbf{W}_2 . We set a block size of $b^{(1)} = b^{(2)} = (16, 16)$. The block grids are $B^{(1)} = (64, 128)$ and $B^{(2)} = (128, 64)$. To group along the hidden dimension, we set group sizes $g^{(1)} = (64, 1)$ and $g^{(2)} = (1, 64)$. These yields group grids of shape $(1, 128)$ for \mathbf{W}_1 and $(128, 1)$ for \mathbf{W}_2 . By coupling these two group grids (permuting the axes of one to align them), we can jointly select the top κ groups, effectively keeping only $\kappa \times 16$ neurons in the hidden layer.

1362 Alternatively, for global block-level sparsity, we can set the group size equal to the block grid size for each tensor ($g = B$). This creates one large group per tensor. Coupling these groups allows us to keep κ blocks across the entire 2-layer MLP.

1365 C.4 DENSE- 4:16 BLOCK-SPARSE MATMUL KERNEL

```

1367
1368 import torch
1369 import triton
1370 import triton.language as tl
1371
1372 """
1373 Grouped Block Sparse GEMM Benchmark
1374 =====
1375 Pattern: For every GROUP_SIZE (=16) contiguous 16x16 blocks along the K
1376 dimension,
1377 only NNZ_PER_GROUP (2, 4, 8) blocks are non zero. This is a structured
1378 block sparsity
1379 pattern analogous to extended N:M at a 16x16 block granularity.
1380
1381 Storage Layout:
1382 Indices : [num_col_blocks, k_group_count, NNZ_PER_GROUP]
1383 Values  : [num_col_blocks, k_group_count, NNZ_PER_GROUP, B_K, B_N]
1384 Flattened in rowmajor order (last dimension fastest). The kernel
1385 expects both tensors
1386 flattened to 1-D contiguous buffers with the above logical order.
1387
1388 Kernel Mapping:
1389 Each program instance (pid_m, pid_n) computes a tile of C of shape
1390 (BLOCK_SIZE_M, B_N). It iterates over all K groups and, inside each
1391 group,
1392 iterates the NNZ_PER_GROUP non zero 16x16 blocks, performing one dot
1393 per block.
1394
1395 Limitations / Assumptions:
1396 - K must be divisible by (GROUP_SIZE * B_K)
1397 - N must be divisible by B_N
1398 - Block size B_K x B_N = 16 x 16
1399 """
1400
1401 DTYPE = torch.float16
1402
1403 @triton.jit
1404 def grouped_block_sparse_kernel_vec(
1405     a_ptr,
1406     b_values_ptr,
1407     b_indices_ptr,
1408     c_ptr,
1409     M,

```

```

1404     N,
1405     K,
1406     stride_am,
1407     stride_ak,
1408     stride_cm,
1409     stride_cn,
1410     B_K: tl.constexpr,
1411     B_N: tl.constexpr,
1412     GROUP_SIZE: tl.constexpr, # 16
1413     NNZ_PER_GROUP: tl.constexpr, # 2,4,8
1414     BLOCK_SIZE_M: tl.constexpr,
1415     GROUP_M: tl.constexpr,
1416 ):
1417     """
1418     For each group g:
1419     1. Load NNZ_PER_GROUP local block indices.
1420     2. Map to global block indices (add g * GROUP_SIZE).
1421     3. Gather all (NNZ_PER_GROUP * B_K) columns from A into a single
1422     tile.
1423     4. Load contiguous B values segment for the group.
1424     5. Perform one tl.dot and accumulate.
1425     This reduces loop overhead inside the group and enables a larger
1426     fused dot.
1427     """
1428     pid_m = tl.program_id(axis=0)
1429     pid_n = tl.program_id(axis=1)
1430     num_pid_m = tl.cdiv(M, BLOCK_SIZE_M)
1431     num_pid_n = tl.cdiv(N, B_N)
1432     pid_m, pid_n = tl.swizzle2d(pid_m, pid_n, num_pid_m, num_pid_n,
1433     GROUP_M)
1434     m_start = pid_m * BLOCK_SIZE_M
1435     n_start = pid_n * B_N
1436     offs_m = m_start + tl.arange(0, BLOCK_SIZE_M)
1437     offs_n = n_start + tl.arange(0, B_N)
1438     c_ptrs = c_ptr + offs_m[:, None] * stride_cm + offs_n[None, :] *
1439     stride_cn
1440     c_mask = (offs_m[:, None] < M) & (offs_n[None, :] < N)
1441     accumulator = tl.zeros((BLOCK_SIZE_M, B_N), dtype=tl.float32)
1442     k_blocks = K // B_K
1443     groups_per_col = k_blocks // GROUP_SIZE
1444     indices_per_col = groups_per_col * NNZ_PER_GROUP
1445     values_per_col = indices_per_col * B_K * B_N
1446     col_indices_base = b_indices_ptr + pid_n * indices_per_col
1447     col_values_base = b_values_ptr + pid_n * values_per_col
1448     offs_i = tl.arange(0, NNZ_PER_GROUP)
1449     b_offs = tl.arange(0, NNZ_PER_GROUP * B_K * B_N)
1450     group_indices_ptr = col_indices_base
1451     b_group_start = col_values_base
1452     b_ptrs = b_group_start + b_offs
1453     for g in range(groups_per_col):
1454         b_group_flat = tl.load(b_ptrs)
1455         local_block_indices = tl.load(group_indices_ptr + offs_i)
1456         global_block_indices = g * GROUP_SIZE + local_block_indices
1457         k_offsets_scattered = (
1458             global_block_indices[:, None] * B_K + tl.arange(0, B_K)[None,
1459             :]
1460         )
1461         k_offsets_flat = tl.reshape(k_offsets_scattered, (NNZ_PER_GROUP *
1462             B_K,))
1463         a_ptrs = (

```

```
1458         a_ptr
1459         + offs_m[:, None] * stride_am
1460         + k_offsets_flat[None, :] * stride_ak
1461     )
1462     a_tile = tl.load(a_ptrs, mask=(offs_m[:, None] < M), other=0.0)
1463
1464     b_tile = tl.reshape(b_group_flat, (NNZ_PER_GROUP * B_K, B_N))
1465     accumulator += tl.dot(a_tile, b_tile)
1466
1467     # b_group_start += NNZ_PER_GROUP * (B_K * B_N)
1468     group_indices_ptr += NNZ_PER_GROUP
1469     b_ptrs += NNZ_PER_GROUP * (B_K * B_N)
1470     tl.store(c_ptrs, accumulator, mask=c_mask)
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
```

C.5 PROXIMAL ADAM OPTIMIZER

Algorithm 3 Adam Optimizer

Require: Step size η , decay rates $\beta_1, \beta_2 \in [0, 1)$, objective $f(\mathbf{w})$

- 1: **Initialize** $m_0 \leftarrow 0, v_0 \leftarrow 0, t \leftarrow 0$
- 2: **while** $t < T$ **do**
- 3: $t \leftarrow t + 1$
- 4: $g_t \leftarrow \nabla_{\theta} f_t(\mathbf{w}_{t-1})$
- 5: $m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$
- 6: $v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$
- 7: $\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$
- 8: $\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$
- 9: $\mathbf{w}_t \leftarrow \mathbf{w}_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$
- 10: **end while**
- 11: **return** \mathbf{w}_T

Algorithm 4 Group Proximal Adam

Require: Regularization $\Omega_{\mathcal{B}, \mathcal{G}}$ and Adam optimizer.

Initialize Adam(α, β_1, β_2)

while $t < T$ **do**

 Run one Adam step to get $\mathbf{v} = \mathbf{w}_t$

 Define $\mathbf{H} = \text{diag}(\hat{v}_t)$

for each block $b \in \mathcal{B}$ **do**

if $\|\mathbf{H}_{b,b} \mathbf{v}_b\| \leq \eta \lambda$ **then**

$\mathbf{w}_{t,b} \leftarrow \mathbf{0}$.

else

 Solve for μ using Equation (45)

 Set $\mathbf{w}_b \leftarrow (\mathbf{H}_{b,b} + \mu \mathbf{I})^{-1} \mathbf{H}_b \mathbf{v}_b$

end if

end for

end while

return \mathbf{w}_T

To extend the soft-thresholding operator to the class of adaptive gradient methods, we construct a local quadratic approximation of the objective F at iteration t . We employ the norm $\|\cdot\|_{\mathbf{H}_t}$ induced by the symmetric positive-definite preconditioner \mathbf{H}_t , defined such that $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{H}_t} = \langle \mathbf{x}, \mathbf{H}_t \mathbf{y} \rangle$.

The surrogate objective function is given by:

$$F_t(\mathbf{w}) \approx f(\mathbf{w}_t) + \langle \hat{\mathbf{m}}_t, \mathbf{w} - \mathbf{w}_t \rangle + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|_{\mathbf{H}_t}^2 + \Omega_{\mathcal{B}, \mathcal{G}}(\mathbf{w}),$$

where $\hat{\mathbf{m}}_t$ denotes the first moment estimate of the gradient.

The update rule is the minimizer of this regularized quadratic model:

$$\mathbf{w}_{t+1} = \text{prox}_{\eta_t \Omega}^{\mathbf{H}_t}(\mathbf{w}_t - \eta_t \mathbf{H}_t^{-1} \hat{\mathbf{m}}_t), \quad (44)$$

where the proximal operator is defined under the \mathbf{H}_t -norm as:

$$\text{prox}_{\eta_t \Omega}^{\mathbf{H}_t}(\mathbf{u}) \triangleq \arg \min_{\mathbf{w}} \left(\frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{u}\|_{\mathbf{H}_t}^2 + \Omega_{\mathcal{B}, \mathcal{G}}(\mathbf{w}) \right).$$

When $\Omega_{\mathcal{B}, \mathcal{G}}$ is the element-wise ℓ_1 regularization, we obtain the same Proximal Adam expression introduced by Melchior et al. (2020), which is equivalent to using a pseudo element-wise learning rate $\hat{\eta}_t = \frac{\eta_t}{\mathbf{H}_{ii}}$.

Proposition 1 (Separable Adaptive Soft-Thresholding with Adam). *Let $\Omega(\mathbf{w}) = \lambda \|\mathbf{w}\|_1$ and let $\mathbf{H}_t = \text{diag}(\mathbf{h}_t)$ be a diagonal positive-definite matrix. Given the intermediate iterate $\mathbf{u}_t = \mathbf{w}_t - \eta_t \mathbf{H}_t^{-1} \hat{\mathbf{m}}_t$, the update \mathbf{w}_{t+1} is given element-wise by:*

$$w_{t+1,i} = \text{sgn}(\mathbf{u}_{t,i}) \cdot \max\left(0, |u_{t,i}| - \frac{\eta_t \lambda}{\mathbf{h}_{t,i}}\right).$$

Proof. Since \mathbf{H}_t is diagonal and the ℓ_1 norm is separable, the objective $\Phi(\mathbf{w})$ decouples into independent scalar problems for each coordinate i . Dropping the subscript t for brevity:

$$\min_{w_i} \left\{ \frac{h_i}{2\eta} (w_i - u_i)^2 + \lambda |w_i| \right\}.$$

The first-order optimality condition requires zero to be in the subdifferential of the objective with respect to w_i :

$$0 \in \frac{h_i}{\eta} (w_i - u_i) + \lambda \partial |w_i|.$$

1566 Rearranging for u_i :

$$1567 \quad u_i \in w_i + \frac{\eta\lambda}{h_i} \partial|w_i|.$$

1568 which equivalent to the Euclidean soft-thresholding with $\hat{\lambda} = \frac{\lambda}{h_i}$. Hence:

$$1571 \quad w_i = \text{sgn}(u_i) \max\left(0, |u_i| - \frac{\eta\lambda}{h_i}\right).$$

1574 \square

1576 With group ℓ_1 regularization $\Omega_{\mathcal{B},\mathcal{G}}$, there is no closed form for the proximal operator, but we can
1577 compute it efficiently using a root finding method (Becker et al., 2018; Adler et al., 2020).

1579 **Adam with Block-wise soft-thresholding** With a diagonal \mathbf{H}_t , the problem decouples to a soft
1580 threshold per-block using the group threshold $\lambda := \lambda_i$ for some $i \in \{1 \dots, q\}$. We denote the
1581 intermediate update for block b as:

$$1582 \quad \mathbf{v} = \mathbf{w}_{b,t} - \eta_t \mathbf{H}_{t,b}^{-1} \hat{\mathbf{m}}_{b,t}$$

1584 Without loss of generality, we will drop the block index b and assume we have a single block:

$$1585 \quad \mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{v}\|_{\mathbf{H}_t}^2 + \lambda \|\mathbf{w}\|_2$$

1587 then we need to find \mathbf{w}^* such that:

$$1589 \quad -\frac{1}{\eta_t} \mathbf{H}(\mathbf{w}^* - \mathbf{v}) \in \lambda \partial \|\mathbf{w}^*\|_2$$

1591 It's is straightforward to show that:

$$1592 \quad \|\mathbf{H}\mathbf{v}\| \leq \eta_t \lambda \iff \mathbf{w}^* = \mathbf{0}$$

1594 We notice here that the Euclidean metric ($\mathbf{H} = \mathbf{I}$) is a special case.

1595 If $\|\mathbf{H}\mathbf{v}\| > \eta_t \lambda$ then $\mathbf{w}^* \neq \mathbf{0}$ and $\partial \|\mathbf{w}^*\|_2 = \frac{\mathbf{w}^*}{\|\mathbf{w}^*\|_2}$ and we get:

$$1597 \quad \frac{1}{\eta_t} \mathbf{H}(\mathbf{w}^* - \mathbf{v}) + \lambda \frac{\mathbf{w}^*}{\|\mathbf{w}^*\|_2} = \mathbf{0}$$

1600 Denote $\mu := \frac{\eta_t \lambda}{\|\mathbf{w}^*\|} > 0$, then:

$$1603 \quad \mathbf{H}(\mathbf{w}^* - \mathbf{v}) + \mu \mathbf{w}^* = \mathbf{0}$$

1605 The solution is then:

$$1606 \quad \mathbf{w}^* = (\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{H}\mathbf{v} = \left[\frac{H_{ii}}{H_{ii} + \mu} \right]_{i \in b} \circ \mathbf{v}$$

1609 The scalar $\mu > 0$ is the unique root of the one-dimensional equation for $\mathbf{H} = \text{diag}(H_{ii})$:

$$1611 \quad \eta_t \lambda = \zeta(\mu) = \mu \|\mathbf{H}(\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{H}\mathbf{v}\|_2 = \mu \sqrt{\sum_{i \in b} \left(\frac{H_{ii}}{H_{ii} + \mu} v_i \right)^2} \quad (45)$$

1614 ζ is strictly increasing in μ with $\zeta(0) = 0$ and $\zeta(+\infty) = \|\mathbf{H}\mathbf{v}\|_2 > \eta_t \lambda$, so the equation has a unique
1615 solution. It can be found efficiently by a simple bisection search.

1616 The group soft-thresholding operator for a diagonal Newton matrix is therefore:

$$1618 \quad S_{\eta\lambda, \mathbf{H}}(\mathbf{v}) = \begin{cases} 0, & \text{if } \|\mathbf{H}\mathbf{v}\|_2 \leq \eta\lambda \\ (\mathbf{H} + \mu \mathbf{I})^{-1} \mathbf{H}\mathbf{v}, & \text{otherwise,} \end{cases}$$

where $\mu > 0$ solves Equation (45). In component form:

$$\mathbf{w}_{t+1} = \begin{cases} 0, & i \notin G, \\ \frac{\mathbf{H}_{ii}}{\mathbf{H}_{ii} + \mu} \mathbf{v}_i, & i \in G. \end{cases}$$

When the Hessian is the identity ($H_k = I$), $h_i = 1$ for all i , it reduces to $\eta\lambda = \mu \|\mathbf{v}\|_2 / (1 + \mu)$, yielding the familiar block soft-thresholding $(1 - \frac{\eta\lambda}{\|\mathbf{v}\|_2})_+ \mathbf{v}$.

Let:

$$S := \|\mathbf{H}\mathbf{v}\|_2^2, \quad h_{\min} := \min_i \mathbf{H}_{ii}, \quad m_{\max} := \max_i \mathbf{H}_{ii}.$$

We bound the sum by min/max of h_{ii} , because every term in the sum is positive:

$$\frac{1}{(h_{\max} + \mu)^2} \leq \frac{1}{(\mathbf{H}_{ii} + \mu)^2} \leq \frac{1}{(h_{\min} + \mu)^2},$$

Multiplying by the non-negative constants $\mathbf{H}_{ii}^2 \mathbf{v}_i^2$ and summing gives

$$S \frac{\mu^2}{(h_{\max} + \mu)^2} \leq \lambda^2 \leq S \frac{\mu^2}{(h_{\min} + \mu)^2}.$$

Both sides increase monotonically in μ . Therefore, the solution μ^* is bracketed by the two solutions of the equalities:

$$S \frac{\mu^2}{(h_{\max} + \mu)^2} = \lambda^2, \quad S \frac{\mu^2}{(h_{\min} + \mu)^2} = \lambda^2.$$

Hence, provided that the solution exists (i.e. $\sqrt{S} > \lambda$), the unique positive root is:

$$\mu(h) = \frac{\lambda h}{\sqrt{S} - \lambda}$$

(the denominator is the same for both h_{\min} and h_{\max}).

With the two extreme values of h we obtain:

$$\mu_{\text{low}} = \frac{\lambda h_{\min}}{\sqrt{S} - \lambda}, \quad \mu_{\text{high}} = \frac{\lambda h_{\max}}{\sqrt{S} - \lambda}$$

and

$$\mu_{\text{low}} \leq \mu^* \leq \mu_{\text{high}}$$

These bounds are inexpensive to compute (one pass to obtain S, h_{\min}, h_{\max}) and provide a tight interval that guarantees convergence of the Newton step for the soft-thresholding operator in a convex optimization problem.

Early exit: if $\eta\lambda \geq \sqrt{S}$ (equivalently $\|H_b \mathbf{v}_b\|_2 \leq \eta\lambda$) we immediately set the block component to zero and skip all further work.

Proximal AdamW To derive the AdamW optimization step within the local linearization framework, we must first understand why standard L_2 regularization fails to yield the desired decoupled weight decay.

Recall the local quadratic approximation for standard Adam at time t , using the norm induced by $\mathbf{H}_t = \text{diag}(\hat{v}_t)^{1/2}$:

$$\min_{\mathbf{w}} \left[f(\mathbf{w}_t) + \langle \mathbf{w} - \mathbf{w}_t, \hat{\mathbf{m}}_t \rangle + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|_{\mathbf{H}_t}^2 \right] \quad (46)$$

Taking the gradient with respect to \mathbf{w} and setting it to zero yields the standard Adam update:

$$\hat{\mathbf{m}}_t + \frac{1}{\eta_t} \mathbf{H}_t (\mathbf{w} - \mathbf{w}_t) = 0 \implies \mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{H}_t^{-1} \hat{\mathbf{m}}_t$$

If we add standard L_2 regularization, $\frac{\rho}{2} \|\mathbf{w}\|_2^2$, to the objective in Eq. equation 46, the gradient condition becomes:

$$\hat{\mathbf{m}}_t + \rho \mathbf{w} + \frac{1}{\eta_t} \mathbf{H}_t (\mathbf{w} - \mathbf{w}_t) = 0$$

Assuming the approximation $\mathbf{w} \approx \mathbf{w}_t$ for the regularization gradient, the update becomes:

$$\mathbf{w}_{t+1} \approx \mathbf{w}_t - \eta_t \mathbf{H}_t^{-1} (\hat{\mathbf{m}}_t + \rho \mathbf{w}_t)$$

Here, the decay term $\rho \mathbf{w}_t$ is scaled by the preconditioner \mathbf{H}_t^{-1} . This couples the regularization with the adaptive learning rates, leading to inconsistent decay across parameters—the issue AdamW aims to solve.

ADAMW FORMULATION

AdamW decouples weight decay from the adaptive gradient step. The target update rule is:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{H}_t^{-1} \hat{\mathbf{m}}_t - \eta_t \rho \mathbf{w}_t$$

To achieve this within the proximal framework, we must adjust the regularization term in the objective function to counteract the geometry of \mathbf{H}_t .

We introduce a **preconditioned linearization** of the regularization term. Instead of the standard Euclidean inner product, we align the penalty with the metric \mathbf{H}_t :

$$\Omega_t(\mathbf{w}) = \rho \langle \mathbf{w} - \mathbf{w}_t, \mathbf{w}_t \rangle_{\mathbf{H}_t}$$

Substituting this into the local approximation objective:

$$f_t^{\text{AdamW}}(\mathbf{w}) = \underbrace{f(\mathbf{w}_t) + \langle \mathbf{w} - \mathbf{w}_t, \hat{\mathbf{m}}_t \rangle}_{\text{Linearized Loss}} + \underbrace{\rho \langle \mathbf{w} - \mathbf{w}_t, \mathbf{w}_t \rangle_{\mathbf{H}_t}}_{\text{Preconditioned Decay}} + \underbrace{\frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|_{\mathbf{H}_t}^2}_{\text{Trust Region}} \quad (47)$$

We minimize $f_t^{\text{AdamW}}(\mathbf{w})$ with respect to \mathbf{w} . Recall that $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathbf{H}_t} = \mathbf{x}^\top \mathbf{H}_t \mathbf{y}$.

$$\nabla_{\mathbf{w}} f_t^{\text{AdamW}}(\mathbf{w}) = \hat{\mathbf{m}}_t + \rho \mathbf{H}_t \mathbf{w}_t + \frac{1}{\eta_t} \mathbf{H}_t (\mathbf{w} - \mathbf{w}_t)$$

Setting the gradient to zero:

$$0 = \hat{\mathbf{m}}_t + \rho \mathbf{H}_t \mathbf{w}_t + \frac{1}{\eta_t} \mathbf{H}_t (\mathbf{w} - \mathbf{w}_t)$$

$$-\frac{1}{\eta_t} \mathbf{H}_t (\mathbf{w} - \mathbf{w}_t) = \hat{\mathbf{m}}_t + \rho \mathbf{H}_t \mathbf{w}_t$$

Multiplying both sides by $-\eta_t \mathbf{H}_t^{-1}$:

$$\mathbf{w} - \mathbf{w}_t = -\eta_t \mathbf{H}_t^{-1} \hat{\mathbf{m}}_t - \eta_t \rho \underbrace{\mathbf{H}_t^{-1} \mathbf{H}_t}_{\mathbf{I}} \mathbf{w}_t$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t (\mathbf{H}_t^{-1} \hat{\mathbf{m}}_t + \rho \mathbf{w}_t)$$

This recovers the exact AdamW update step, where the weight decay $\rho \mathbf{w}_t$ is applied isotropically, independent of the adaptive scaling \mathbf{H}_t^{-1} .

As such, the proximal AdamW is equivalent to proximal Adam if we use the update $\mathbf{v} = \mathbf{w}_t - \eta_t (\mathbf{H}_t^{-1} \hat{\mathbf{m}}_t + \rho \mathbf{w}_t)$