OPTIMIZING OPTIMIZERS FOR FAST GRADIENT-BASED LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

We lay the theoretical foundation for automating optimizer design in gradient-based learning. Based on the greedy principle, we formulate the problem of designing optimizers as maximizing the instantaneous decrease in loss. By treating an optimizer as a function that translates loss gradient signals into parameter motions, the problem reduces to a family of convex optimization problems over the space of optimizers. Solving these problems under various constraints not only recovers a wide range of popular optimizers as closed-form solutions, but also produces the optimal hyperparameters of these optimizers with respect to the problems at hand. This enables a systematic approach to design optimizers and tune their hyperparameters according to the gradient statistics collected from training or validation sets. Furthermore, this optimization of optimization can be performed dynamically during training.

1 Introduction

We are interested in the problem of designing optimizers that maximize the utility of gradient-based learning for a given task. In gradient-based learning, the objective is to minimize an expected scalar loss $\mathbb{E}[\mathcal{L}(\theta)]$ with respect to parameters $\theta \in \mathbb{R}^d$ using its (negative) gradient $g = -\nabla_{\theta}\mathcal{L} \in \mathbb{R}^d$. As learning takes time, all the parameters $\theta = \theta(t)$, the loss $\mathcal{L} = \mathcal{L}(\theta(t))$, and the gradients g = g(t) are variables of time t, i.e., the training step. A process of learning manifests as a parameter motion $\dot{\theta}$ driven by the gradient g calculated at each step t.

Physics requires a constitutive law that relates kinematic motion to the force field that causes it. For gradient-based learning, optimizers take that role. We can represent an optimizer as a positive semidefinite operator $Q \succeq 0$ that translates the gradient into the parameter update,

$$\dot{\theta} = Q g. \tag{1}$$

By the chain rule, the instantaneous loss drop is then a quadratic form:

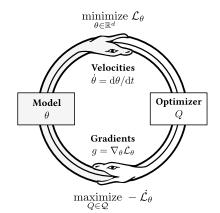


Figure 1: Just as optimizers train their models by feeding them parameter velocities $\dot{\theta}$, models can also fit the optimizers to the underlying tasks by feeding gradients g.

$$-\dot{\mathcal{L}} = \nabla_{\theta} \mathcal{L}^{\top} \frac{\mathrm{d}\theta}{\mathrm{d}t} = g^{\top} \dot{\theta} = g^{\top} Q g. \tag{2}$$

The greedy paradigm turns our original problem of maximizing the utility of learning into another optimization problem that maximizes the this loss drop with respect to the optimizer Q:

$$\underset{Q \in \mathcal{Q}}{\text{maximize}} \ \mathbb{E}[g^{\top}Q \ g] \quad \text{subject to} \quad Q \succeq 0,$$
 (P1)

where Q is the design space of allowed optimizers.

Instantaneously, we notice that without any additional constraint, the maximum of the quadratic form $g^{\top}Qg$ is unbounded. Problem P1 reveals two design options for the optimizer that bounds the

maximum: (1) the budget constraint $Q \in \mathcal{Q}$, and (2) the data distribution under the expectation \mathbb{E} . Our main focus is on how these two factors determine the *optimal optimizer* Q^* .

Placing the optimizer itself as a subject of another optimization is interesting in several ways:

- Optimizers can be designed with respect to the individual problem (task and data) and the running environment (budget and precision) in a systematic manner.
- Optimizers can be tuned or even be replaced by better ones according to the intermediate probes from either training or validation sets in the middle of training.
- Solving this *meta-optimization* problem in closed-form for a wide range of budgets uncovers the relationship between the optimizers and their underlying principles.
- By reverse engineering commonly used optimizers, we draw the landscape of optimizers that
 have driven the success of machine learning (Robbins & Monro, 1951; Kingma & Ba, 2015;
 Loshchilov & Hutter, 2019; Gupta et al., 2018; Martens & Grosse, 2015) into a single picture.

2 OPTIMAL STATELESS OPTIMIZERS

Consider the following setup: Let $\pi_{\rm tr}$ and $\pi_{\rm val}$ be the training and validation data distributions, respectively. Then, for each training $x_{\rm tr} \sim \pi_{\rm tr}$ and validation sample $x_{\rm val} \sim \pi_{\rm val}$, the gradients are denoted by $g_{\rm tr} = \nabla_{\theta} \mathcal{L}(\theta, x_{\rm tr})$ and $g_{\rm val} = \nabla_{\theta} \mathcal{L}(\theta, x_{\rm val})$. We define the *moments* as:

$$\Sigma_{\text{tr}} = \mathbb{E}[g_{\text{tr}} g_{\text{tr}}^{\top}], \qquad C = \mathbb{E}[g_{\text{tr}} g_{\text{val}}^{\top}], \qquad \Sigma_{\text{val}} = \frac{1}{2}(C + C^{\top}),$$
 (3)

where \mathbb{E} denotes expectation over the single or joint distributions of the enclosed gradients. Note that Σ_{tr} and Σ_{val} are symmetric and positive semidefinite (PSD) matrices of shape $d \times d$. For any symmetric PSD $Q \in \mathbb{S}^d_+$ of shape $d \times d$, we define the *learning power* as:

$$P_{\text{tr}}(Q) := \mathbb{E}[g_{\text{tr}}^{\top}Q g_{\text{tr}}] = \text{Tr}(Q \Sigma_{\text{tr}}), \qquad P_{\text{val}}(Q) := \mathbb{E}[g_{\text{tr}}^{\top}Q g_{\text{val}}] = \text{Tr}(Q \Sigma_{\text{val}}).$$
 (4)

We call $P_{\rm tr}(Q)$ the *training power* and $P_{\rm val}(Q)$ the *validation cross-power*. From the chain rule of equation 2, the learning power is equal to the expected instantaneous loss drop: $-\mathbb{E}_{\rm tr}[\dot{\mathcal{L}}] = \mathbb{E}_{\rm tr}[g_{\rm tr}^{\top}\dot{\theta}_{\rm tr}] = P_{\rm tr}(Q)$ and $-\mathbb{E}_{\rm val}[\dot{\mathcal{L}}] = \mathbb{E}_{\rm val}[g_{\rm val}^{\top}\dot{\theta}_{\rm tr}] = P_{\rm val}(Q)$. Problem P1 is therefore rewritten as:

$$\underset{Q \in \mathcal{Q}}{\text{maximize}} P_{\circ}(Q) = \text{Tr}(Q \Sigma_{\circ}) \quad \text{subject to} \quad Q \succeq 0,$$
(P2)

where $\circ \in \{\text{tr}, \text{val}\}$. This is our main optimization problem.

Solving this without any additional constraint, we end up with arbitrarily large eigenvalues for the optimizer Q. This corresponds to arbitrarily large learning rates, which we all know are infeasible in practice. Real problems give us several reasons that makes this "ideal solution" unrealizable: finite precision of our machines, curvature of the loss landscapes, stochastic nature of subset gradients, etc. All of them restrict the ability of gradient estimates g to represent the global geometry of the parameter space. Taking a large step in the parameter space beyond the regions where g remains explainable leads to unexpected, and usually fatal, behaviors.

In other words, the aforementioned restrictions define the feasible set, or the *budget* $Q \subseteq \mathbb{S}^d_+$, which induces a solution space of finite optimizers. The following theorem formalizes this:

Theorem 1 (Optimal stateless optimizers under convex constraints). Let the budget set $\{0\} \subseteq \mathcal{Q} \subseteq \mathbb{S}^d_+$ be a nonempty, compact, convex set. Define also (1) the indicator $\delta_{\mathcal{Q}}(Q) = 0$ if $Q \in \mathcal{Q}$ and $+\infty$ otherwise, (2) the gauge (Minkowski functional): $\gamma_{\mathcal{Q}}(Q) = \inf\{\lambda > 0 : Q \in \lambda \mathcal{Q}\}$, and (3) the polar set $\mathcal{Q}^\circ = \{\Sigma \in \mathbb{S}^d : \sup_{Q \in \mathcal{Q}} \operatorname{Tr}(Q\Sigma) \leq 1\}$. For any symmetric matrix $\Sigma \in \mathbb{S}^d$,

- (i) (Existence and sublinearity): The maximum in $P^*(\Sigma) := \max_{Q \in \mathcal{Q}} \operatorname{Tr}(Q\Sigma)$ is attained. In addition, P^* is sublinear (convex and positively homogeneous) and finite everywhere.
- (ii) (Conjugacy identities): The maximum $P^*(\Sigma)$ is obtained by the identity relationships:

$$P^* = \delta_{\mathcal{Q}}^* = \gamma_{\mathcal{Q}^{\circ}} \quad \text{and} \quad \gamma_{\mathcal{Q}}^* = \delta_{\mathcal{Q}^{\circ}},$$
 (5)

i.e., the optimal power is equal to the convex conjugate of the indicator and also the gauge of the polar, while the conjugate of the gauge is equal to the indicator of the polar.

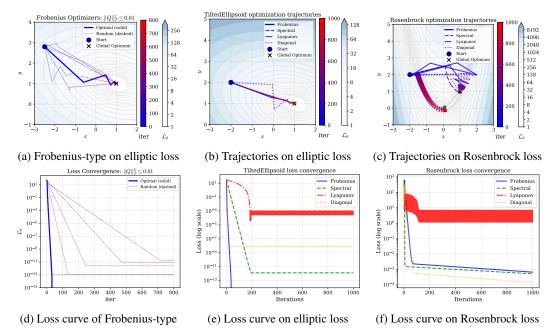


Figure 2: Behavior of optimal optimizers under different budget types. (a) Dotted lines show suboptimal optimizers with random Σ from the equal-power Frobenius budget; the straight line is from the optimal hyperparameters from our theory, achieving fastest convergence. (b, c) No free lunch theorem: Frobenius excels in simple quadratic loss, while nonconvex geometry makes spectral and diagonal types better. Each line is the best result from dense search among all budget-parameters, e.g, B for Frobenius budget, etc.

- (iii) (Construction): Any maximizer $Q^* \in \arg \max_{Q \in \mathcal{Q}} \operatorname{Tr}(Q\Sigma)$ is a subgradient of P^* at Σ : $Q^* \in \partial_{\Sigma} P^*(\Sigma)$. If the maximizer is unique, P^* is differentiable at Σ and $Q^* = \nabla_{\Sigma} P^*(\Sigma)$.
- (iv) (Order preservation on \mathbb{S}^d_+): If $\Sigma \succeq 0$, then $P^*(\Sigma) \geq 0$. If $\Sigma_1 \succeq \Sigma_2$, then $P^*(\Sigma_1) \geq P^*(\Sigma_2)$.
- (v) (Lipschitz continuity in symmetrized polar gauge): $Define \|\cdot\|_{\mathcal{Q}^{\circ}}^{\mathrm{sym}} := \max\{\gamma_{\mathcal{Q}^{\circ}}(\cdot), \, \gamma_{\mathcal{Q}^{\circ}}(-\cdot)\}.$ For any $\Sigma, \hat{\Sigma} \in \mathbb{S}^d$, $|P^{\star}(\Sigma) P^{\star}(\hat{\Sigma})| \leq \|\Sigma \hat{\Sigma}\|_{\mathcal{Q}^{\circ}}^{\mathrm{sym}}.$

The proof is in Appendix B. Items (i), (ii), and (iii) provide a principled way to construct the optimal (stateless) optimizer Q^* from any given moment Σ and any nicely conditioned budget set Q. Items (iv) and (v) add robustness guarantee and sensitivity analysis to the optimal power. In practice, full gradients rarely appears in large settings. Gradients are drifting throughout non-convex loss landscapes, making true moments hard to obtain. Theorem 1(v) shows that the estimation error in the optimal power is bounded.

Solving the optimization problem constrained by Q determines the optimal optimizer Q^* , and endows the optimizer with different characteristics and algorithmic behaviors. Consider the following four types of budgets:

- Frobenius ball budget $Q_F(B) = \{Q \succeq 0 : \|Q\|_F^2 \le B\}$ is the simplest constraint that gives an isotropic Euclidean trust region without prior knowledge about parameter space geometry.
- Spectral budget $Q_S(\tau, \lambda) = \{Q \succeq 0 : Tr(Q) \leq \tau, \ Q \preceq \lambda I\}$ is a budget that upper limits the per-direction spectrum for safety and the trace for total budget simultaneously.
- Data-metric (Lyapunov) budget $Q_L(B) = \{Q \succeq 0 : Tr(Q^2\Sigma) \leq B\}$ is a budget that uses the data covariance itself as the metric, leading to a natural Lyapunov-like stability condition.
- Diagonal budget $Q_D(B,c) = \{Q = \operatorname{diag}(q_j) \succeq 0 : \sum_j c_j q_j^2 \leq B\}$ is a budget that restricts to coordinate-wise optimizers. Most of the commonly used optimizers fall into this category.

Instantiating the construction from Theorem 1 with these budgets, we obtain corresponding closed-form solutions for the optimal optimizer Q^* and the optimal power P^* .

Corollary 2 (Closed-form solutions for common budget sets). Let $\Sigma = U \operatorname{diag}(\sigma_1 \ge \cdots \ge \sigma_d)U^{\top}$ be the eigendecomposition. The optimal solutions are:

- (i) (Frobenius ball): $Q_F^{\star} = \sqrt{B} \Sigma / ||\Sigma||_F, P_F^{\star}(\Sigma) = \sqrt{B} ||\Sigma||_F.$
- (ii) (Spectral): $Q_{\mathbf{S}}^{\star} = U \operatorname{diag}(q_{i}^{\star})U^{\top}$ where (i) $q_{i}^{\star} = \lambda$ for $i \leq k$, (ii) $q_{k+1}^{\star} = \tau k\lambda$, (iii) $q_{i}^{\star} = 0$ for i > k+1, where $k = \lfloor \tau/\lambda \rfloor$, $P_{\mathbf{S}}^{\star}(\Sigma) = \lambda \sum_{i < k} \sigma_{i} + (\tau k\lambda)\sigma_{k+1}$.
- (iii) (Data-metric): $Q_L^{\star} = \alpha \prod_{\text{supp}(\Sigma)} \text{ where } \alpha = (B/\sum_{i:\sigma_i>0} \sigma_i)^{1/2}, P_L^{\star}(\Sigma) = (B\sum_i \sigma_i)^{1/2}.$
- (iv) (Diagonal): $[Q_D^{\star}]_{jj} \propto \Sigma_{jj}/c_j$, $P_D^{\star}(\Sigma) = (B\sum_j \Sigma_{jj}^2/c_j)^{1/2}$.

Again, the proof is in Appendix B. These analytic solutions reveal how the characteristics of different types of optimal optimizers Q^* are induced by controlling the budget Q. Specifically, we see that:

Frobenius budget \leftrightarrow Eveidence-proportional optimizer. Budget $\mathcal{Q}_F(B)$ gives optimizers that allocate learning power *proportionally* to data evidence $Q^* \propto \Sigma$. We can project this general class of optimizers into special geometries to calculate for the optimal hyperparameters as in Corollary 5.

Spectral budget \leftrightarrow Water-filling optimizer \sim gradient clipping & LR scheduling. Budget $\mathcal{Q}_S(\tau,\lambda)$ returns a *water-filling* optimizer that concentrates the learning power into the largest available principal components of the data moment Σ up to a per-mode cap λ , sequentially, until the total budget τ is reached. The spectral bound λ , therefore, acts as a stability margin similar to gradient clipping tricks. The trace bound τ controls the total budget like *learning rate schedules*.

Data-metric budget \leftrightarrow **Equal-power optimizer** \supset {AdaGrad, natural gradient}. Budget $\mathcal{Q}_L(B)$ results in an *equal-power* optimizer that whitens gradient statistics and allocates uniform power across *L*-eigendirections. If *L* is Fisher information matrix, this is *natural gradient descent* (Amari, 1998). Generally, this includes full-matrix AdaGrad (Duchi et al., 2011; Agarwal et al., 2019), K-FAC (Martens & Grosse, 2015), and Shampoo (Gupta et al., 2018).

Diagonal budget \leftrightarrow Coordinate-wise optimizer \supset {Adam, GD}. Budget $\mathcal{Q}_{D}(B,c)$ produces an optimizer that allocates the total budget B coordinate-wise which scales with the evidence Σ_{jj} and inversely with the costs c_{j} . Assuming isotropic moments $\Sigma = \sigma^{2}I$ and constant costs c = 1, this reduces to $Q^{\star} = \eta I$ with learning rate $\eta = (B/d)^{1/2}$. This recovers simple gradient descent. With variance-based costs $c_{j} \propto (\mathrm{EMA}(g_{j}^{2}))^{1/2}$ and momenta $m_{j,t} = \mathrm{EMA}(g_{j})$, this recovers Adam optimizer (Kingma & Ba, 2015); without the first moment it reduces to diagonal AdaGrad (Duchi et al., 2011) or RMSProp-style optimizer (Tieleman & Hinton, 2012). Corollary 6 gives the optimal hyperparameters for Adam from this setup.

The behaviors of different types of optimizers are visualized in Figure 2. Figure 2a shows that our analytically found optimizer is the fastest among all hyperparameter settings under the same Frobenius budget. On the other hand, Figure 2b and 2c highlights how optimizers from different types of budgets can perform better in their specialized domains. This insight helps address the notorious no free lunch theorem in optimization (Wolpert & Macready, 1997): The catchphrase "no single algorithm is universally superior" can be updated to "the optimal optimizer Q^* is a function of the budget Q and the distribution of observations Σ ," at least under our greedy optimization framework. In summary, users choose the budget Q, and the budget defines what is optimal optimizing methods. Reverse engineering under our framework reveals the hidden principles behind the design of commonly used optimizers such as Adam (Kingma & Ba, 2015). Nevertheless, wise readers will notice that an important component is still missing: momentum. The next section demonstrates how momentum is integrated into our framework through a straightforward extension.

3 OPTIMAL DYNAMIC OPTIMIZERS WITH STATE VARIABLES

Up to this point, we have formulated the problem of finding the optimal stateless optimizer as a convex optimization problem, and derived the closed-form solutions for the four types of budgets. In practice, optimizers have memory, often in the form of momentum, in order to stabilize the learning process from stochastic gradients and non-convex loss landscapes. We now extend our framework by letting the optimizer Q[n] be a *causal dynamical operator*: a *filter* that translates gradient *history* g[n] into *instantaneous* parameter velocity $\dot{\theta}[n]$. In the budget view, this introduces a new degree of freedom: where in frequency we allocate our finite learning power.

Let us highlight the key differences from Section 2. We now work in discrete time $n \in \{0,1,2,\ldots\}$ representing training steps. We will use the z-domain as the primary spectral domain to reflect the iterative nature of practical algorithms. Let $g[n] \in \mathbb{R}^d$ be the (mini-batch) wide-sense stationary (WSS) gradient process. A dynamic optimizer is an LTI filter with symmetric matrix impulse response $Q[n] \in \mathbb{R}^{d \times d}$ defined by the causal convolution:

$$\dot{\theta}[n] = (Q * g)[n] := \sum_{k=0}^{\infty} Q[k] g[n-k], \qquad \Psi(z) := \sum_{n=0}^{\infty} Q[n] z^{-n}, \tag{6}$$

where the *transfer function* is $\Psi(z) \in \mathbb{C}^{d \times d}$. Also adopt the Hilbert norm:

$$||Q||_{\mathcal{H}}^2 := \sum_{n=0}^{\infty} \text{Tr}(Q[n]^{\top}Q[n]) < \infty, \qquad \langle Q_1, Q_2 \rangle_{\mathcal{H}} := \sum_{n=0}^{\infty} \text{Tr}(Q_1[n]^{\top}Q_2[n]).$$
 (7)

In this framework, autocorrelation and symmetrized cross-correlation represent the moments:

$$R_{\text{tr}}[k] := \mathbb{E}[g_{\text{tr}}[n] g_{\text{tr}}[n-k]^{\top}], \ C[k] := \mathbb{E}[g_{\text{val}}[n] g_{\text{tr}}[n-k]^{\top}], \ R_{\text{val}}[k] := \frac{1}{2}(C[k] + C[k]^{\top}), \ (8)$$

where $k \ge 0$. The *instantaneous learning power* is the inner product (more details in Appendix A):

$$P_{\circ}(Q;n) := \mathbb{E}\left[g_{\circ}[n]^{\top}\dot{\theta}[n]\right] = \mathbb{E}\left[g_{\circ}[n]^{\top}\sum_{k=0}^{\infty}Q[k]\,g_{\mathrm{tr}}[n-k]\right] = \sum_{k=0}^{\infty}\mathrm{Tr}\left(Q[k]^{\top}R_{\circ}[k]\right) = \langle Q,R_{\circ}\rangle_{\mathcal{H}}.$$

Define a nonempty, convex, and weakly compact budget set $\{0\} \subseteq \mathcal{Q} \subset \mathcal{H}$, which includes any norm-bounded, closed, convex subset of \mathcal{H} . We also define the indicator $\delta_{\mathcal{Q}}(Q)$ and the gauge $\gamma_{\mathcal{Q}}(Q)$ the same as in Section 2, and the polar set is defined as $(\mathcal{Q})^{\circ} := \{R \in \mathcal{H} \mid \sup_{Q \in \mathcal{Q}} \langle Q, R \rangle_{\mathcal{H}} \leq 1\}$. Hence, all the notations are consistent with Section 2.

For dynamic optimization, for each $\circ \in \{tr, val\}$, problem P2 is lifted to:

$$\underset{Q \in \mathcal{Q}}{\text{maximize}} P_{\circ}(Q) = \langle Q, R_{\circ} \rangle_{\mathcal{H}} \quad \text{subject to} \quad Q \succeq 0, \tag{P3}$$

Unsurprisingly, we arrive at similar results as in Section 2; only the Frobenius inner product $\langle \cdot, \cdot \rangle_F = \operatorname{Tr}(\cdot^\top \cdot)$ is replaced by the Hilbert space inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$. We have the same conjugacy toolkit for dynamic optimization.

Theorem 3 (Optimal dynamic optimizers under convex constraints). Given the definitions above, the followings hold for any nonempty, convex, and weakly compact budget set $Q \subset \mathcal{H}_+$ with $0 \in Q$:

- (i) (Existence and sublinearity): The maximum $P^*(R) := \max_{Q \in \mathcal{Q}} \langle Q, R \rangle_{\mathcal{H}}$ is attained. Moreover, P^* is sublinear and finite everywhere.
- (ii) (Conjugacy identities): $P^* = \delta_{\mathcal{Q}}^* = \gamma_{\mathcal{Q}^{\circ}}$ and $\gamma_{\mathcal{Q}}^* = \delta_{\mathcal{Q}^{\circ}}$.
- (iii) (Construction): Any maximizer $Q^* \in \arg \max_{Q \in \mathcal{Q}} \langle Q, R \rangle_{\mathcal{H}}$ is a subgradient of P^* at $R: Q^* \in \partial_R P^*(R)$. If the maximizer is unique, P^* is differentiable at R and $\nabla_R P^*(R) = Q^*$.
- (iv) (Order preservation on \mathcal{H}_+): If $R \in \mathcal{H}_+$ (Hermitian PSD a.e.), then $P^*(R) \geq 0$. Moreover, if $R_1 R_2 \in \mathcal{H}_+ \setminus \{0\}$ and $\exists Q \in \mathcal{Q}$ with $\langle Q, R_1 R_2 \rangle_{\mathcal{H}} > 0$ (e.g., if \mathcal{Q} contains a positive definite element), then $P^*(R_1) > P^*(R_2)$.
- (v) (Lipschitz continuity in the symmetrized polar gauge): Define the symmetrized polar gauge $\|u\|_{\mathcal{Q}^{\circ}}^{\operatorname{sym}} := \max\{\gamma_{\mathcal{Q}^{\circ}}(u), \, \gamma_{\mathcal{Q}^{\circ}}(-u)\}$. Then $\forall R, \hat{R} \in \mathcal{H}, \, |P^{\star}(R) P^{\star}(\hat{R})| \leq \|R \hat{R}\|_{\mathcal{Q}^{\circ}}^{\operatorname{sym}}$.

The proof is similar to the stateless case, and is provided in Appendix B. Theorem 3 formalizes how the optimal dynamic optimizer Q^* equalizes the learning power across different frequencies as a function of the convex budget Q. All the closed-form solutions from Corollary 2 can also be directly lifted to the dynamic framework, as elaborated in Appendix A. Instead of redundantly repeating the closed-form solutions, we discuss how they are connected to well-known optimizers. As we will see in Corollaries 5 and 6, solving Problem P3 using Theorem 3 often produces general dynamic optimizers with infinite impulse responses (IIR) Q[n], whose implementation require infinite memory for optimizer states. In practice, we often restrict ourselves to simpler, realizable family of optimizers, such as ones with EMA-based momenta. The following lemma justifies this post-projection of optimizers already obtaind from convex budgets.

Lemma 4. Let \mathcal{H} be a real Hilbert space. Given a nonzero moment $R \in \mathcal{H}$, let $\mathcal{Q} \subset \mathcal{H}$ be nonempty, closed, convex, with $0 \in \mathcal{Q}$. Let $\mathcal{C} \subset \mathcal{H}$ be a cone (closed under positive scaling). The normal cone of $\mathcal{Q} \cap \mathcal{C}$ at Q is $N_{\mathcal{Q} \cap \mathcal{C}}(Q) \coloneqq \{M \in \mathcal{H} : \langle M, Q' - Q \rangle_{\mathcal{H}} \leq 0 \ \forall Q' \in \mathcal{Q} \cap \mathcal{C}\}$. Define the solution sets of the optimization problem and its restriction to \mathcal{C} :

$$\mathcal{Q}^{\star}(R) \coloneqq \arg\max_{Q \in \mathcal{Q}} \langle Q, R \rangle_{\mathcal{H}}, \qquad \mathcal{Q}^{\star}_{\mathcal{C}}(R) \coloneqq \arg\max_{Q \in \mathcal{Q} \cap \mathcal{C}} \langle Q, R \rangle_{\mathcal{H}}. \tag{10}$$

Let $\Pi_{\mathcal{C}}$ be the Hilbert metric projection onto \mathcal{C} . For any $Q^{\star} \in \mathcal{Q}^{\star}(R)$, the following are equivalent:

(i) (Commutativity) $\Pi_{\mathcal{C}}(Q^*) \in \mathcal{Q}_{\mathcal{C}}^*(R)$.

(ii) (Normal-cone alignment) There exists $Q_{\mathcal{C}}^{\star} \in \mathcal{Q}_{\mathcal{C}}^{\star}(R)$ such that $\{R, Q^{\star} - Q_{\mathcal{C}}^{\star}\} \subset N_{\mathcal{Q} \cap \mathcal{C}}(Q_{\mathcal{C}}^{\star})$,

Moreover, if $N_{Q \cap C}(Q_C^*)$ is a ray $\{\lambda M : \lambda \geq 0\}$, then commutativity holds if and only if R and $Q^* - Q_C^*$ are positive multiples of the same direction M.

In other words, if the projected manifold $\mathcal C$ of desired optimizers is in a sufficiently good shape, e.g., using EMA-based momentum, then we can first solve problem P3 for general budgets $\mathcal Q$ and then project the solution onto the manifold $\mathcal C$ to obtain the final optimal solution over $\mathcal C$. Now we are ready to find the optimal hyperparameters for real optimizers in use.

Corollary 5 (Instantaneous optimal SGD+Momentum). *Define the budget* $Q_F(B)$ *and the cone* C_{Ip} *of isotropic 1-pole optimizers:*

 $Q_F(B) := \{Q : \|Q\|_{\mathcal{H}} \le \sqrt{B}\}, \quad \mathcal{C}_{Ip} := \{Q_{\eta,\beta}[n] = \eta(1-\beta)\beta^n I : \eta \ge 0, 0 < \beta < 1\}.$ (11) Given $R[n] \in \mathcal{H}$, define $S[n] := \operatorname{Tr}(R[n])$. Then solving problem P3 under the budget $Q_F(B) \cap \mathcal{C}_{Ip}$ produces SGD+Momentum optimizer as the optimal solution with optimal hyperparameters:

$$\beta^{\star} = \arg \max_{0 < \beta < 1} \sqrt{1 - \beta^2} \sum_{n=0}^{\infty} S[n] \beta^n, \qquad \eta^{\star} = \frac{\sqrt{B(1 - \beta^{\star 2})}}{\sqrt{d}(1 - \beta^{\star})}, \tag{12}$$

where d is the dimension of the parameter space.

The optimal hyperparameters are obtained by first applying Theorem 3 to general family of budgets, and then projecting the solution into the approximation geometry \mathcal{C}_{1p} . By Lemma 4, the optimal hyperparameters of the projected solution are consistent with the unprojected solution. Corollary 5 shows how SGD+Momentum type *optimal* optimizer works. The *optimal momentum* β^* first maximizes the cosine similarity between the 1-pole EMA response $(\beta^n)_{n\geq 0}$ and the trace of the empirical moment $(S[n])_{n>0}$; then the *learning rate* η^* scales to saturate the budget.

Corollary 6 does the same to Adam Kingma & Ba (2015). For Adam, the existence of a time-varying divisor $\mathrm{EMA}(g^2,\beta_2)^{-1/2}$ slightly complicates the derivation by making the optimizer time-varying. Corollary 6 (Instantaneous optimal Adam). Let moment be a diagonal matrix $R[n] = \mathrm{diag}(r_j[n])$ with coordinate-wise sequence $r_j[n]$. Given gradients g[t], maintain the second-moment EMA $v_j[t] \coloneqq \beta_2 v_j[t-1] + (1-\beta_2)g_j[t]^2 > 0$ with parameter $\beta_2 \in (0,1)$. Fix the current time t and define the coordinate-wise costs $c_j \coloneqq v_j[t]^{1/2}$. Consider the diagonal budget

$$Q_D(B,c) := \{ \operatorname{diag}(q_j) : \sum_{i} c_j \sum_{k>0} |q_j[k]|^2 \le B \},$$
 (13)

and the cone $C_{1p, Adam}$ of diagonal 1-pole optimizer with per-coordinate inverse-cost scaling

$$C_{Ip, Adam} := \{Q_{\eta, \beta_1}[n] = \operatorname{diag}(\eta(1 - \beta_1)\beta_1^n/c_j) : \eta \ge 0, 0 < \beta_1 < 1\}.$$
(14)

Optimizing for problem P3 under $Q_D(B,c) \cap C_{Ip, Adam}$ with moment R[n] yields an Adam optimizer with optimal hyperparameters:

$$(\beta_1^{\star}, \beta_2^{\star}) = \arg \max_{0 < \beta_1, \beta_2 < 1} a(\beta_1, \beta_2) \sum_{n=0}^{\infty} \beta_1^n T_t[n; \beta_2], \quad \eta^{\star} = \frac{\sqrt{B} a(\beta_1^{\star}, \beta_2^{\star})}{1 - \beta_1^{\star}}, \tag{15}$$

where
$$W_t(\beta_2) := \sum_j 1/c_j$$
, $T_t[n; \beta_2] := \sum_j r_j[n]/c_j$, and $a(\beta_1, \beta_2) := \sqrt{(1-\beta_1^2)/W_t(\beta_2)}$.

The optimal β_1^\star maximizes the cosine similarity between the 1-pole EMA response $(\beta_1^n)_{n\geq 0}$ and the cost-compensated diagonal moment $(T[n;\beta_2])_{n\geq 0}$. The optimal β_2^\star determines the diagonal weights through the EMA of the second-moment. Finally, the learning rate η^\star scales to saturate the budget. In summary, Adam is a *1-pole approximation* of the *dynamic diagonal optimizer* defind by the diagonal budget $\mathcal{Q}_D(B,c)$ projected onto the 1-pole family $\mathcal{C}_{1p,\,Adam}$. All the proofs are provided in Appendix B.

Table 1: Reverse-engineered optimizers as convex optimization problems. Each optimizer emerges as the optimal solution to a specific convex optimization problem, with hyperparameters determined by the underlying budget constraint.

Optimizer	Hyperparameters	Interpretation	
GD	η (learning rate)	Euclidean trust region: $\ \dot{\theta}\ _2 \le au$ yields $Q \propto I$	
Colored-GD	η , precon. P	Elliptic trust region: $\ \dot{\theta}\ _P \leq \tau$ yields $Q \propto P^{-1}$ with P^\star from data covariance	
Newton/GN	η , damping λ	Curvature trust region: $\ \dot{\theta}\ _H \le \tau$ yields $Q \propto H^{-1}$ with $\eta^\star = \tau \frac{\ g\ _{H^{-1}}}{\ H^{-1}g\ ^2}$ and λ^\star from condition number	
NGD	η , Fisher est.	Fisher/KL trust region: $\ \dot{\theta}\ _F \le \tau$ yields $Q \propto F^{-1}$ with $\eta^\star = \frac{\sqrt{B}}{\sqrt{\text{Tr}(F)\ g\ _2}}$ and F^\star from empirical Fisher	
K-FAC/Shampoo	η , damping, period	Structured Kronecker trust region yields block $Q=\bigoplus (A_\ell^{-1}\otimes G_\ell^{-1})$ with η^\star from factored Fisher metric and damping λ^\star from eigenvalue gaps	
AdaGrad	η , ϵ , window	Diagonal trust region: $\sum_j c_j \dot{\theta}_j^2 \leq B$ with $c_j = \sqrt{\sum_{t'} g_{j,t'}^2}$ and $\eta_j^\star = \sqrt{B} \frac{ g_j }{\sqrt{\sum_j g_j^2/c_j}}$, ϵ^\star for numerical stability	
RMSProp	η,eta_2,ϵ	Diagonal trust region with $c_j = \sqrt{v_j}$, $v_j = \beta_2 v_{j,t-1} + (1 - \beta_2) g_j^2$ and $\eta_j^\star = \sqrt{B} \frac{ g_j /\sqrt{v_j}}{\sqrt{\sum_j g_j^2/v_j}}$, β_2^\star from gradient correlation	
Adam	$\eta, \beta_1, \beta_2, \epsilon$	Dynamic diagonal trust region on momentum m_k with costs $c_j = \sqrt{v_j}$ and $\eta^* = \sqrt{B} \frac{\ m/\sqrt{v}\ _1}{\ m/\sqrt{v}\ _2^2}, \beta_1^*, \beta_2^*$ from lag curves	
SGD+Momentum	η, eta_1	1-pole EMA trust region: $\ Q\ _{\mathcal{H}} \leq \sqrt{B}$ with $Q = \eta(1-\beta_1)\sum_{n=0}^{\infty}\beta_1^n I$ and $\eta^{\star} = \frac{\sqrt{B(1-\beta_1^2)}}{\ g\ _{\mathcal{H}}(1-\beta_1)}, \beta_1^{\star}$ from cosine similarity	
AdaFactor	η, β_1, β_2 , factorization	Factored trust region with $r_i = \beta_2 r_{i,t-1} + (1-\beta_2) \ G_{i,:}\ ^2$, $c_j = \beta_2 c_{j,t-1} + (1-\beta_2) \ G_{:,j}\ ^2$ and η^\star from factored diagonal, $\beta_1^\star, \beta_2^\star$ from matrix structure	
LayerNorm+Mom	η, β_1 , layer norm	Layer-wise normalized trust region: $\ Q\ _{\mathcal{H}_{\text{layer}}} \leq \sqrt{B}$ with $\eta^{\star} = \frac{\sqrt{B(1-\beta_1^2)}}{\sqrt{\sum_{\ell} d_{\ell}/\ \theta_{\ell}\ _2^2(1-\beta_1)}}$ and β_1^{\star} from layer-normalized lag	
LAMB/LARS	η , trust ratio $ au$	Layer-wise trust region: $\ \dot{\theta}_{\ell}\ _2 \le \tau_{\ell} \ \theta_{\ell}\ / \ m_{\ell}\ $ with $\eta_{\ell}^{\star} = \tau_{\ell} \frac{\ \theta_{\ell}\ }{\ m_{\ell}\ }$ and τ^{\star} from layer-wise analysis	
signSGD	η (step size)	L_{∞} trust region: $\ \dot{\theta}\ _{\infty} \leq \tau$ yields $\dot{\theta} = \tau \operatorname{sign}(g)$ with $\eta^{\star} = \tau$	
Lion	η, eta_1	L_{∞} trust region on momentum: $\ \dot{\theta}\ _{\infty} \leq \tau$ with $\dot{\theta} = \eta \operatorname{sign}(m)$, $\eta^{\star} = \tau$ and β_1^{\star} from sign correlation	

Reverse engineering optimizers. The above two corollaries show that SGD with momentum and Adam (Kingma & Ba, 2015) are the optimal 1-pole approximations of the dynamic equalizers, establishing these well-known optimizers as special cases of our framework. Similarly, we can reverse engineer various other optimizers, including SGD with Nesterov momentum (Nesterov, 1983), AdamW (Loshchilov & Hutter, 2019), LAMB (You et al., 2020), K-FAC (Martens & Grosse, 2015), Shampoo (Gupta et al., 2018), and Lion (Chen et al., 2023), into our framework. Due to page constraints, we defer the details to Appendix C. We have also provided a master table of optimizers in Appendix D, categorizing many optimizers widely used in practice today.

Automatic hyperparameter tuning. Examining Corollaries 5 and 6, our formulation not only classifies optimizers according to their underlying budget choices, but also provides a systematic way to determine *optimal hyperparameters*. Unsurprisingly, optimal hyperparameters depend on the data being processed. After warm-up steps that collect data covariance, optimal hyperparameters can be computed and used throughout training. Moreover, systematic hyperparameter determination enables automatic tuning during training, which is impossible with manual tuning. The algorithms for determining hyperparameters is provided in Appendix E.

4 AUTOMATING VALIDATION-AWARE OPTIMIZER TUNING

So far, our theoretical justification for optimizing optimizers did not specify which datasets are in use. This section addresses a more delicate question of how to *systematically* exploit validation sets for optimizer design. It is commonly considered bad practice to use validation sets directly

in the optimization loop. Rather, they are typically used to generate subtle clues that indirectly guide engineers when making decisions about model architecture, optimizers, and associated hyperparameters. We can regard this manual tuning process as a kind of "human-in-the-loop" system that *fits* the optimizer and hyperparameters to the validation set. We can then automate this process by casting it as the same mathematical optimization problem. For example, our original optimization problem P3 can be recast in terms of maximizing the instantaneous *validation* loss drop as

$$\underset{Q \in \mathcal{Q}}{\text{maximize}} \ - \dot{\mathcal{L}}_{\text{val}} \ = \ \sum_{n} \mathbb{E}[g_{\text{val}}[n]^{\top} \ \dot{\theta}_{\text{tr}}[n]] \ = \ \langle Q, R_{\times} \rangle_{\mathcal{H}} \ = \ P_{\times}(Q) \quad \text{s.t.} \quad Q \succeq 0,$$

where $\dot{\theta}_{\rm tr}[n]=(q*g_{\rm tr})[n]$ is the parameter velocity guided solely by the training set and the designed optimizer Q. This is mathematically equivalent to problem P3 but with the cross-moment R_{\times} , mimicking human inspection of the validation loss drop and engineering the optimizer accordingly. This approach may or may not conflict with traditional practice, potentially requiring an additional split beyond the standard training/validation division. We leave this discussion to the readers. Here, we focus on the theoretical side of validation-aware optimizer design by showing that tuning optimizers using validation sets maximizes the instantaneous validation loss drop.

Proposition 7 (Validation optimality in power). Given any convex and compact budget $\mathcal{Q} \subset \mathcal{H}_+$, the validation-aware maximizer $Q_\times^* \in \arg\max_{Q \in \mathcal{Q}} \langle Q, R_\times \rangle_{\mathcal{H}}$ gives the maximum instantaneous validation loss drop among all possible optimizers $Q \in \mathcal{Q}$, including training-only optimizers $Q \in \arg\max_{Q \in \mathcal{Q}} \langle Q, R_{\text{tr}} \rangle_{\mathcal{H}}$.

We leave the proof to Appendix B. Therefore, validation-aware tuning is the best possible instantaneous validation loss drop among all possible optimizers $Q \in \mathcal{Q}$, including training-only optimizers $Q \in \arg\max_{Q \in \mathcal{Q}} \operatorname{Tr}(Q \Sigma_{\operatorname{tr}})$. The next proposition shows how the choice of optimizer Q not only controls the parameter velocity, but also determines the endpoint of the optimization process in the local convex region.

To approximate the dynamics in the local convex region in the loss landscape, assume squared loss and fix parameters θ . In the linearizable region of the network f_{θ} around θ , the function-space dynamics follow *kernel gradient flow* with kernel

$$K_Q(x, x') := \nabla_{\theta} f_{\theta}(x; \theta)^{\top} Q \nabla_{\theta} f_{\theta}(x'; \theta).$$
 (16)

Define the budget \mathcal{Q} , the validation cross-moment Σ_{\times} , and the validation-aware maximizer $Q_{\times}^{\star} \in \arg\max_{Q \in \mathcal{Q}} \operatorname{Tr}(Q\Sigma_{\times})$ the same as in Proposition 7. Then, the following proposition holds:

Proposition 8 (Endpoint selection (function-space view)). With any optimizer $Q \in \mathcal{Q}$ fixed around θ , kernel gradient flow converges to the unique minimum-norm interpolant in \mathcal{H}_{K_Q} (or to kernel ridge with decay $\lambda > 0$). Consequently, choosing $Q = Q_{\times}^{\star}$ changes the RKHS to $\mathcal{H}_{K_{Q_{\times}^{\star}}}$ and selects the endpoint

$$f_{Q_{\times}^{\star}}^{\star} = \arg\min_{f \in \mathcal{H}_{K_{Q_{\times}^{\star}}}} \|f\|_{\mathcal{H}_{K_{Q_{\times}^{\star}}}}^{2} \quad \text{subject to} \quad f(X) = y.$$
 (17)

In summary, tuning the optimizer Q using validation sets is the best possible way not only to maximize the instantaneous validation loss drop, but also to select better convergence endpoints of the optimization process.

5 SCOPE AND LIMITATIONS

Long-horizon objective from greedy paradigm. In order to simplify the analysis, this work resorts to the greedy paradigm, primarily focusing on instantaneous progress of learning. As a trade-off, global optimality guarantee requires further investigation under this greedy paradigm. Our filter interpretation of dynamic optimizers in Section 3 mitigates the limitations of the greedy objective by incorporating stateful optimizers holding summaries of the history of gradients. This also generalizes the optimizers with momentum and other EMA-based smoothing techniques.

Choice of budget. Instead of telling which class of optimizers are optimal for a given task, this work provides an optimization framework *under* the user-defined budget set. Our framework can help engineers by reducing their effort for searching the right hyperparameters; however, it still requires an intelligent choice of which budget, i.e., which *class of* optimizers, fits the task.

Renovating existing optimizers. In Theorems 1 and 3, we provide general construction of optimal optimizers from convex constraints. However, in this work, we focus on the well-established optimizers, and supplement them with a systematic methodology to find the right hyperparameters for a given task. Designing a new class of optimizers will be a natural extension of this work.

6 RELATED WORK

 First, we would like to highlight the key differences between our "optimizing optimizers" and other well-established fields with similar tautologies.

Learning to optimize. Learning to optimize (Li & Malik, 2016) aims to adapt the optimizer to a given task by treating optimizers as learnable parametric models (Andrychowicz et al., 2016). Various architectures have been explored, including RNNs (Andrychowicz et al., 2016; Wichrowska et al., 2017; Lv et al., 2017), Transformers (Chen et al., 2022; Moudgil et al., 2023; Jain et al., 2024), and per-tensor HyperNetworks (Ha et al., 2016; Metz et al., 2022). Their primary focus is on meta-training these optimizer-networks for stability and adaptability. These works represent a nontraditional, network-based family of generally nonconvex optimizers, which is not generally compatible with our framework which is based on convex optimization.

Learning to learn. Rooted in the human-inspired philosophy (Schmidhuber, 1987; Bengio et al., 1990), meta-learning is another line of work that shares a similar spirit with learning to optimize (Gharoun et al., 2023). A large proportion of works on meta-learning target few-shot learning tasks, which prepare the model, not the optimizer, for downstream tasks (Vinyals et al., 2016; Finn et al., 2017; Yu et al., 2024; Sun et al., 2019). Among them, Meta-SGD (Li et al., 2017) is noteworthy, as it prepares the optimizer. However, the problem set we address is general gradient-based learning, which differs from the tasks of concern in meta-learning.

Algorithmic discovery of optimizers. Techniques like symbolic discovery (Chen et al., 2023; Zheng et al., 2022), non-parametric optimizer search (Wang et al., 2022), and neural optimizer search (Bello et al., 2017) are also related to our work, as their objective is to discover the optimal optimizer for a given task. In their framework, symbolic optimizers are obtained by a tree-based search of a predefined set of optimizers. Ours instead lets the engineer select the broader family of optimizers, and then provides a mathematical tool to find the optimal solution among them. Therefore, these works are considered to be orthogonal to ours.

Hyperparameter optimization. Many works have proposed to automatically tune the hyperparameters governing optimization. Most of them adopt a learning framework to find a good set of hyperparameters including learning rates (Daniel et al., 2016), their schedules (Xu et al., 2017; 2019), and other optimizer parameters (Shaban et al., 2019). Hypergradient methods (Maclaurin et al., 2015; Baydin et al., 2017; Grazzi et al., 2020; Moskovitz et al., 2019) are also proposed to find the optimal hyperparameters. Instead of resorting to learning-based methods, we establish a theoretical framework through the lens of convex optimization problems (Boyd & Vandenberghe, 2004). By doing so, we can classify well-used optimizers such as SGD with momentum and Adam (Kingma & Ba, 2015) as special cases of our framework, and provide a systematic way to determine the optimal hyperparameters for these optimizers.

7 Conclusion

We established a firm theoretical grounding for systematically achieving optimal optimizers in a greedy sense. Our convex optimization framework connects commonly used optimizers to convex constraint sets, merging those independently developed techniques into a single unified framework. Our main results, Theorems 1 and 3, and Lemma 4 are general tools that can be extended to arbitrary budget to invent new families of optimizers for specific uses. Our theory, therefore, does not disprove the *no free lunch theorem*; rather, it provides a principled way to *leverage* this wisdom to flexibly design and adapt optimizers for our own problems at hand.

ETHICS STATEMENT

We acknowledge the ICLR Code of Ethics.

REPRODUCIBILITY STATEMENT

All the proofs of the theoretical part of this paper, including every Lemma, Theorem, Proposition, and Corollary, are provided in Appendix B with detailed derivations, starting from the basic definitions and assumptions made in the main text. Moreover, omitted theoretical results are elaborated in Appendix A. Regarding the implementation, Appendix E gives the algorithm to realize our theoretically justified optimal optimizers.

LLM USAGE STATEMENT

We deeply acknowledge the usefulness of LLMs in revising the manuscript, especially for fixing vocabulary and grammar-related issues. We also used LLMs to check the correctness and coherence of the proofs and notations. This greatly helped us in identifying awkward mistakes we had been making all along.

REFERENCES

- Naman Agarwal, Brian Bullins, Xinyi Chen, Elad Hazan, Karan Singh, Cyril Zhang, and Yi Zhang. Efficient full-matrix adaptive regularization. In *ICML*, 2019.
- Shun-ichi Amari. Natural gradient works efficiently in learning. *Neural Computation*, 10(2):251–276, 1998.
- Marcin Andrychowicz, Misha Denil, Sergio Gomez Colmenarejo, Matthew W Hoffman, David Pfau, Tom Schaul, Brendan Shillingford, and Nando De Freitas. Learning to learn by gradient descent by gradient descent. In *NIPS*, 2016.
- Atilim Gunes Baydin, Robert Cornish, David Martinez Rubio, Mark Schmidt, and Frank Wood. Online learning rate adaptation with hypergradient descent. *arXiv preprint arXiv:1703.04782*, 2017.
- Irwan Bello, Barret Zoph, Vijay Vasudevan, and Quoc V. Le. Neural optimizer search with reinforcement learning. In *ICML*, 2017.
- Yoshua Bengio, Samy Bengio, and Jocelyn Cloutier. Learning a synaptic learning rule. Technical report, Université de Montréal, Département d'informatique et de recherche opérationnelle, 1990.
- Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Kaiyuan Wang, Hieu Pham, Xuanyi Dong, Thang Luong, Cho-Jui Hsieh, Yifeng Lu, and Quoc V. Le. Symbolic discovery of optimization algorithms. In *NeurIPS*, 2023.
- Yutian Chen, Xingyou Song, Chansoo Lee, Zi Wang, Qiuyi Zhang, David Dohan, Kazuya Kawakami, Greg Kochanski, Arnaud Doucet, Marc'aurelio Ranzato, Sagi Perel, and Nando de Freitas. Towards learning universal hyperparameter optimizers with transformers. In *NeurIPS*, 2022.
- Christian Daniel, Jonathan Taylor, and Sebastian Nowozin. Learning step size controllers for robust neural network training. In *AAAI*, 2016.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159, 2011.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *ICML*, 2017.

546

547

554

555

564

565

566

567

568 569

570

571

572

573

574

575

576

577

578579

580

581

582

583

584 585

586

587

591

592

- Hassan Gharoun, Fereshteh Momenifar, Fang Chen, and Amir H. Gandomi. Meta-learning approaches for few-shot learning: A survey of recent advances. arXiv preprint arXiv:2303.07502, 2023.
- Riccardo Grazzi, Luca Franceschi, Massimiliano Pontil, and Saverio Salzo. On the iteration complexity of hypergradient computation. In *ICML*, 2020.
 - Vineet Gupta, Tomer Koren, and Yoram Singer. Shampoo: Preconditioned stochastic tensor optimization. In *ICML*, 2018.
- David Ha, Andrew M. Dai, and Quoc V. Le. HyperNetworks. *arXiv preprint arXiv:1609.09106*, 2016.
- Deepali Jain, Krzysztof M. Choromanski, Kumar Avinava Dubey, Sumeet Singh, Vikas Sindhwani,
 Tingnan Zhang, and Jie Tan. Mnemosyne: Learning to train transformers with transformers. In
 NeurIPS, 2024.
 - Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In ICLR, 2015.
- Ke Li and Jitendra Malik. Learning to Optimize. arXiv preprint arXiv:1606.01885, 2016.
- Zhenguo Li, Fengwei Zhou, Fei Chen, and Hang Li. Meta-SGD: Learning to learn quickly for few-shot learning. In NIPS, 2017.
- Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019.
- Kaifeng Lv, Shunhua Jiang, and Jian Li. Learning gradient descent: Better generalization and longer
 horizons. In *ICML*, 2017.
 - Dougal Maclaurin, David Duvenaud, and Ryan Adams. Gradient-based hyperparameter optimization through reversible learning. In *ICML*, 2015.
 - James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In ICML, 2015.
 - Luke Metz, James Harrison, C. Daniel Freeman, Amil Merchant, Lucas Beyer, James Bradbury, Naman Agrawal, Ben Poole, Igor Mordatch, Adam Roberts, and Jascha Sohl-Dickstein. VeLO: Training versatile learned optimizers by scaling up. *arXiv preprint arXiv:2211.09760*, 2022.
 - Ted Moskovitz, Rui Wang, Janice Lan, Sanyam Kapoor, Thomas Miconi, Jason Yosinski, and Aditya Rawal. First-order preconditioning via hypergradient descent. *arXiv preprint arXiv:1910.08461*, 2019.
 - Abhinav Moudgil, Boris Knyazev, Guillaume Lajoie, and Eugene Belilovsky. Learning to optimize with recurrent hierarchical transformers. In *ICML Workshop on New Frontiers in Learning, Control, and Dynamical Systems*, 2023.
 - Yurii E. Nesterov. A method for solving the convex programming problem with convergence rate $o(1/k^2)$. *Doklady Akademii Nauk SSSR*, 269(3):543–547, 1983. English translation: Soviet Mathematics Doklady, 27(2):372–376, 1983.
 - Herbert Robbins and Sutton Monro. A stochastic approximation method. *Annals of Mathematical Statistics*, 22(3):400–407, 1951.
 - Jürgen Schmidhuber. Evolutionary Principles in Self-Referential Learning, or on Learning How to Learn: The Meta-Meta-... Hook. PhD thesis, Technische Universität München, 1987.
- Amirreza Shaban, Ching-An Cheng, Nathan Hatch, and Byron Boots. Truncated back-propagation for bilevel optimization. In *AISTATS*, 2019.
 - Qianru Sun, Yaoyao Liu, Tat-Seng Chua, and Bernt Schiele. Meta-transfer learning for few-shot learning. In *CVPR*, 2019.
 - Tijmen Tieleman and Geoffrey Hinton. Lecture 6.5—rmsprop: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural Networks for Machine Learning, 2012.

- Oriol Vinyals, Charles Blundell, Tim Lillicrap, and Daan Wierstra. Matching networks for one shot learning. In *NIPS*, 2016.
 - Ruochen Wang, Yuanhao Xiong, Minhao Cheng, and Cho-Jui Hsieh. Efficient non-parametric optimizer search for diverse tasks. In *NeurIPS*, 2022.
 - Olga Wichrowska, Niru Maheswaranathan, Matthew W Hoffman, Sergio Gomez Colmenarejo, Misha Denil, Nando de Freitas, and Jascha Sohl-Dickstein. Learned optimizers that scale and generalize. In *ICML*, 2017.
 - David H. Wolpert and William G. Macready. No free lunch theorems for optimization. *IEEE Transactions on Evolutionary Computation*, 1(1):67–82, 1997.
 - Chang Xu, Tao Qin, Gang Wang, and Tie-Yan Liu. Reinforcement learning for learning rate control. *arXiv preprint arXiv:1705.11159*, 2017.
 - Zhen Xu, Andrew M. Dai, Jonas Kemp, and Luke Metz. Learning an adaptive learning rate schedule. *arXiv preprint arXiv:1909.09712*, 2019.
 - Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. Large batch optimization for deep learning: Training bert in 76 minutes. In *ICLR*, 2020.
 - Xingtong Yu, Yuan Fang, Zemin Liu, Yuxia Wu, Zhihao Wen, Jianyuan Bo, Xinming Zhang, and Steven C.H. Hoi. A survey of few-shot learning on graphs: from meta-learning to pre-training and prompt learning. *arXiv* preprint arXiv:2402.01440, 2024.
 - Wenqing Zheng, Tianlong Chen, Ting-Kuei Hu, and Zhangyang Wang. Symbolic learning to optimize: Towards interpretability and scalability. *arXiv* preprint arXiv:2203.06578, 2022.

A MORE MATHEMATICAL RESULTS

This section provides more detailed mathematical foundations of the main text, which was omitted for brevity. Appendix A.1 shows a detailed derivation of equation 9. This shows that the instantaneous learning powers of dynamic optimizers are also represented by some inner product between the optimizer operator and the gradient moment, having the same structure as in the stateless case. Appendix A.2 draws connection between the results in Section 2 and the results in Section 3, by showing the corresponding results in the dynamic setting.

A.1 Proof of Equation 9

Here we provide a detailed derivation of equation 9 that was abbreviated in the main text. We assume the sequences $\{g_{\text{tr}}[n]\}$ and $\{g_{\text{val}}[n]\}$ are zero-mean, wide-sense stationary (WSS) with finite second moments. Hence the lag-k moments $R_{\text{tr}}[k] = \mathbb{E}[g_{\text{tr}}[n]g_{\text{tr}}[n-k]^{\top}]$ and $C[k] = \mathbb{E}[g_{\text{val}}[n]g_{\text{tr}}[n-k]^{\top}]$ depend only on k. Then

$$P_{\text{tr}}(Q) = \langle Q, R_{\text{tr}} \rangle_{\mathcal{H}}, \qquad P_{\text{val}}(Q) = \langle Q, R_{\text{val}} \rangle_{\mathcal{H}}.$$
 (18)

We start from the convolution definition of the dynamic optimizer:

$$\dot{\theta}[n] = \sum_{k=0}^{\infty} Q[k] g_{tr}[n-k]. \tag{19}$$

Instantaneous training power:

$$P_{\text{tr}}(Q;n) = \mathbb{E}\left[g_{\text{tr}}[n]^{\top}\dot{\theta}[n]\right] \tag{20}$$

$$= \mathbb{E}\Big[g_{\text{tr}}[n]^{\top} \sum_{k=0}^{\infty} Q[k] g_{\text{tr}}[n-k]\Big]$$
(21)

$$= \sum_{k=0}^{\infty} \mathbb{E} \left[g_{tr}[n]^{\top} Q[k] g_{tr}[n-k] \right] \quad \text{(linearity of } \mathbb{E} \text{)}$$
 (22)

$$= \sum_{k=0}^{\infty} \operatorname{Tr}(Q[k] \mathbb{E}[g_{tr}[n-k] g_{tr}[n]^{\top}])$$
(23)

$$= \sum_{k=0}^{\infty} \text{Tr}(Q[k]^T \mathbb{E}[g_{tr}[n] g_{tr}[n-k]^\top]) \quad \text{(trace transpose)}$$
 (24)

$$= \sum_{k=0}^{\infty} \text{Tr}(Q[k]^T R_{\text{tr}}[k])$$
 (25)

$$= \langle Q, R_{\rm tr} \rangle_{\mathcal{H}}. \tag{26}$$

Instantaneous validation cross-power:

$$P_{\text{val}}(Q; n) = \mathbb{E}\left[g_{\text{val}}[n]^{\top} \dot{\theta}[n]\right]$$
(27)

$$= \mathbb{E}\Big[g_{\text{val}}[n]^{\top} \sum_{k=0}^{\infty} Q[k] g_{\text{tr}}[n-k]\Big]$$
(28)

$$= \sum_{k=0}^{\infty} \mathbb{E}\left[g_{\text{val}}[n]^{\top} Q[k] g_{\text{tr}}[n-k]\right]$$
 (29)

$$= \sum_{k=0}^{\infty} \operatorname{Tr}(Q[k] \mathbb{E}[g_{tr}[n-k] g_{val}[n]^{\top}])$$
(30)

$$= \sum_{k=0}^{\infty} \text{Tr} \left(Q[k]^T \mathbb{E}[g_{\text{val}}[n] g_{\text{tr}}[n-k]^\top]^T \right)$$
 (31)

$$= \sum_{k=0}^{\infty} \text{Tr}(Q[k]^T C[k]^T), \tag{32}$$

where $C[k] := \mathbb{E}[g_{\text{val}}[n] \ g_{\text{tr}}[n-k]^{\top}]$. For the symmetric cross-moment $R_{\text{val}}[k] := \frac{1}{2} \left(C[k] + C[k]^{\top} \right)$, when using Hermitian PSD filters with symmetric Q[k], we have

$$\operatorname{Tr}(Q[k]^T C[k]^T) = \operatorname{Tr}(Q[k] \frac{1}{2} (C[k] + C[k]^\top)) = \operatorname{Tr}(Q[k]^T R_{\text{val}}[k]),$$
 (33)

since for symmetric A, $\operatorname{Tr}(A \frac{1}{2}(B - B^{\top})) = 0$. Therefore, $P_{\text{val}}(Q; n) = \langle Q, R_{\text{val}} \rangle_{\mathcal{H}}$.

A.2 DYNAMIC LIFT OF SECTION 2

We now lift the key results of Section 2 to the dynamic setting of Section 3. Solving the optimization problem with constraint Q determines the optimal dynamic optimizer Q^* , and endows the optimizer with different characteristics and algorithmic behaviors. Again, consider the following four types of dynamic budgets:

- Frobenius ball budget $Q_F(B) = \{Q : \|Q\|_{\mathcal{H}} \leq \sqrt{B}\}$ is the simplest constraint that gives an isotropic Hilbert space trust region without prior knowledge about temporal correlation structure.
- Per-frequency spectral budget $Q_S(\tau, \lambda) = \{Q : \text{Tr}(Q(e^{i\omega})) \leq \tau(\omega), Q(e^{i\omega}) \leq \lambda(\omega)I\}$ is a budget that upper limits the per-direction spectrum for safety and the trace for total budget simultaneously at each frequency.
- Data-metric (Lyapunov) budget $\mathcal{Q}_L(B) = \{Q: \sum_{k=0}^{\infty} \operatorname{Tr}(Q_k^* \Sigma Q_k) \leq B\}$ is a budget that uses the lag-covariance sequence itself as the metric, leading to a natural dynamic Lyapunov-like stability condition.
- Diagonal budget $Q_D(B,c) = \{Q(z) = \operatorname{diag}(q_j(z)) \succeq 0 : \sum_j c_j \|q_j\|_{H_2}^2 \leq B\}$ is a budget that restricts to diagonal dynamic optimizers with coordinate-wise budgets.

Instantiating the construction from Theorem 3 with these budgets, we obtain corresponding closed-form solutions for the optimal dynamic optimizer Q^* and the optimal power $P^*(R)$.

Corollary 9 (Closed-form solutions for dynamic budget sets). Let $R(e^{i\omega}) = U(\omega) \operatorname{diag}(\sigma_i(\omega))U(\omega)^*$ be the eigendecomposition at each frequency. The optimal solutions are:

- (i) Frobenius ball: $Q_F^* = \sqrt{B}R/\|R\|_{\mathcal{H}}, P_F^*(R) = \sqrt{B}\|R\|_{\mathcal{H}}.$
- (ii) Per-frequency spectral: $Q_S^*(e^{i\omega}) = U(\omega)\operatorname{diag}(q_i^*(\omega))U(\omega)^*$ where $q_i^*(\omega) = \min\{\lambda(\omega), \max\{0, \mu(\omega) \sigma_i(\omega)\}\}$ and $\mu(\omega)$ is chosen so that $\sum_i q_i^*(\omega) = \tau(\omega)$.
- (iii) Data-metric: $Q_L^* = \alpha \Pi_{\mathrm{supp}(R)}$ where $\alpha = \sqrt{B/\sum_{i:\sigma_i>0} \sigma_i}$, $P_L^*(R) = \sqrt{B\sum_i \sigma_i}$.
- (iv) Diagonal: $[Q_D^*]_{jj} \propto R_{jj}/c_j$, $P_D^*(R) = \sqrt{B\sum_j R_{jj}^2/c_j}$.

We delay the proof to Appendix B. These analytic solutions reveal how the characteristics of different types of optimal dynamic optimizers Q^* are induced by controlling the budget set Q.

Frobenius budget \leftrightarrow Proportional dynamic optimizer. Budget $\mathcal{Q}_F(B)$ produces a proportional dynamic optimizer that allocates learning power proportional to lag-covariance evidence $Q^\star \propto R$. It enjoys implementation simplicity but potentially over-concentrates on dominant temporal modes. Using Lemma 4, we can project this general class of optimizers into special geometries to calculate for the optimal hyperparameters as in Corollary 5.

Per-frequency spectral budget \leftrightarrow Water-filling dynamic optimizer \sim per-frequency gradient clipping & power scheduling. Budget $\mathcal{Q}_S(\tau,\lambda)$ produces a water-filling dynamic optimizer that keeps pushing power into responsive frequency modes until hitting the safety cap $\lambda(\omega)$. The spectral cap $\lambda(\omega)$ acts as a frequency-dependent safety mechanism similar to gradient clipping, while the trace constraint $\tau(\omega)$ controls the per-frequency learning rate similar to adaptive learning rate scheduling.

Data-metric budget \leftrightarrow **Equal-power dynamic optimizer** \sim **dynamic AdaGrad.** Budget $\mathcal{Q}_L(B)$ produces an *equal-power* dynamic optimizer that equalizes learning power uniformly across informative lag-correlation directions, preventing over-concentration while maintaining temporal efficiency. This produces a natural dynamic preconditioning effect similar to *AdaGrad*'s inverse square root scaling, but with uniform power allocation across all informative temporal directions rather than instantaneous adaptation.

Diagonal budget \leftrightarrow **Coordinate-wise dynamic optimizer** \sim **Adam.** Budget $\mathcal{Q}_D(B,c)$ produces a *coordinate-wise* dynamic optimizer that adapts per-coordinate learning power proportional to the lag-covariance evidence R_{jj} and inversely proportional to the costs c_j . When $c_j \propto v_{t,j}^{1/2}$ ($v_{t,j}$ being the EMA of g_j^2 at time t), this recovers the core mechanism of *Adam*, which is elaborated in Corollary 6 in the main manuscript.

B PROOFS OMITTED FROM THE MAIN TEXT

This section does all the proofs that has been omitted in the main text. The proofs are organized in the same order as the theorems appear in the main manuscript.

B.1 Proof of Theorem 1

 Proof of Theorem 1. We establish each claim in turn.

Existence & sublinearity: Since $\mathcal Q$ is compact by assumption (a nonempty, compact, convex subset of $\mathbb S^d_+$), and $(Q,\Sigma)\mapsto \operatorname{Tr}(Q\Sigma)$ is continuous, the maximum is attained by the Weierstrass extreme value theorem. The optimal power $P^\star(\Sigma)=\sup_{Q\in\mathcal Q}\operatorname{Tr}(Q\Sigma)$ is a supremum of linear functionals in Σ , hence sublinear (convex and positively homogeneous). Finiteness follows from compactness of $\mathcal Q$.

Conjugacy identities: We establish the three identities in equation 5.

1. Optimal power = conjugate of indicator. By the definition of convex conjugate,

$$(\delta_{\mathcal{Q}})^*(\Sigma) = \sup_{Q \in \mathbb{S}^d} \{ \langle Q, \Sigma \rangle - \delta_{\mathcal{Q}}(Q) \} = \sup_{Q \in \mathcal{Q}} \langle Q, \Sigma \rangle = P^*(\Sigma). \tag{34}$$

Thus $P^* = (\delta_{\mathcal{Q}})^*$.

2. Optimal power = gauge of polar. By the definition of polar, $\Sigma \in \mathcal{Q}^{\circ}$ if and only if $\sup_{Q \in \mathcal{Q}} \langle Q, \Sigma \rangle \leq 1$, i.e., $P^{\star}(\Sigma) \leq 1$. Therefore

$$\gamma_{\mathcal{Q}^{\circ}}(\Sigma) = \inf\{\lambda > 0 : \Sigma \in \lambda \mathcal{Q}^{\circ}\} = \inf\{\lambda > 0 : P^{\star}(\Sigma) \le \lambda\} = P^{\star}(\Sigma). \tag{35}$$

Thus $P^* = \gamma_{\mathcal{O}^{\circ}}$.

- 3. Conjugate of gauge = indicator of polar. We establish $(\gamma_{\mathcal{O}})^* = \delta_{\mathcal{O}^{\circ}}$. Consider two cases:
 - If $\Sigma \in \mathcal{Q}^{\circ}$, then for all Q,

$$\langle Q, \Sigma \rangle \le \gamma_{\mathcal{Q}}(Q) \cdot \sup_{R \in \mathcal{Q}} \langle R, \Sigma \rangle \le \gamma_{\mathcal{Q}}(Q),$$
 (36)

since $\sup_{R\in\mathcal{Q}}\langle R,\Sigma\rangle\leq 1$ by definition of polar. Hence $\langle Q,\Sigma\rangle-\gamma_{\mathcal{Q}}(Q)\leq 0$ for all Q, with equality at Q=0. Taking the supremum gives $(\gamma_{\mathcal{Q}})^*(\Sigma)=0=\delta_{\mathcal{Q}^\circ}(\Sigma)$.

• If $\Sigma \notin \mathcal{Q}^{\circ}$, there exists $Q_0 \in \mathcal{Q}$ with $\langle Q_0, \Sigma \rangle > 1$. For any $\alpha > 0$, we have $\gamma_{\mathcal{Q}}(\alpha Q_0) = \alpha \gamma_{\mathcal{Q}}(Q_0) = \alpha$ (since $Q_0 \in \mathcal{Q}$ so $\gamma_{\mathcal{Q}}(Q_0) = 1$), and thus

$$\langle \alpha Q_0, \Sigma \rangle - \gamma_{\mathcal{Q}}(\alpha Q_0) = \alpha \langle Q_0, \Sigma \rangle - \alpha = \alpha (\langle Q_0, \Sigma \rangle - 1) \to +\infty \quad (\alpha \to \infty). \tag{37}$$

Hence $(\gamma_{\mathcal{Q}})^*(\Sigma) = +\infty = \delta_{\mathcal{Q}^{\circ}}(\Sigma)$.

Thus $(\gamma_{\mathcal{Q}})^* = \delta_{\mathcal{Q}^{\circ}}$.

Construction by subgradient: Let $Q^* \in \arg \max_{Q \in \mathcal{Q}} \operatorname{Tr}(Q\Sigma)$. For any $M \in \mathbb{S}^d$,

$$P^{\star}(M) = \max_{Q \in \mathcal{Q}} \operatorname{Tr}(QM) \ge \operatorname{Tr}(Q^{\star}M) = \operatorname{Tr}(Q^{\star}\Sigma) + \operatorname{Tr}(Q^{\star}(M - \Sigma)) = P^{\star}(\Sigma) + \operatorname{Tr}(Q^{\star}(M - \Sigma)),$$

which is the defining inequality for $Q^* \in \partial P^*(\Sigma)$. If the maximizer is unique, $\partial P^*(\Sigma) = \{Q^*\}$ and P^* is differentiable at Σ with $\nabla P^*(\Sigma) = Q^*$.

Order preservation: If $\Sigma \succeq 0$, then for any $Q \in \mathcal{Q} \subseteq \mathbb{S}^d_+$, we have $\operatorname{Tr}(Q\Sigma) \geq 0$. Since $0 \in \mathcal{Q}$, the maximum over $Q \in \mathcal{Q}$ is ≥ 0 . If $\Sigma_1 \succeq \Sigma_2$, then $P^*(\Sigma_1) \geq P^*(\Sigma_2)$. Moreover, strict inequality holds if there exists $Q \in \mathcal{Q}$ with $\operatorname{Tr}(Q(\Sigma_1 - \Sigma_2)) > 0$.

Lipschitz continuity in symmetrized polar gauge: We establish the one-sided bounds first. Since $P^* = \gamma_{\mathcal{O}^{\circ}}$ by the conjugacy identities, we have:

$$P^{\star}(\Sigma) - P^{\star}(\hat{\Sigma}) = \max_{Q \in \mathcal{Q}} \langle Q, \Sigma \rangle - \max_{Q \in \mathcal{Q}} \langle Q, \hat{\Sigma} \rangle$$
 (39)

$$\leq \max_{Q \in \mathcal{Q}} \langle Q, \Sigma - \hat{\Sigma} \rangle \tag{40}$$

$$= \gamma_{\mathcal{Q}^{\circ}}(\Sigma - \hat{\Sigma}). \tag{41}$$

Similarly, $P^{\star}(\hat{\Sigma}) - P^{\star}(\Sigma) \leq \gamma_{\mathcal{Q}^{\circ}}(\hat{\Sigma} - \Sigma)$. Therefore,

$$|P^{\star}(\Sigma) - P^{\star}(\hat{\Sigma})| \le \max\{\gamma_{\mathcal{Q}^{\circ}}(\Sigma - \hat{\Sigma}), \gamma_{\mathcal{Q}^{\circ}}(\hat{\Sigma} - \Sigma)\} = \|\Sigma - \hat{\Sigma}\|_{\mathcal{Q}^{\circ}}^{\text{sym}}.$$
 (42)

This Lipschitz property is crucial for robustness analysis. By using an estimated moment $\hat{\Sigma}$ instead of the true moment Σ , the error in optimal power can be bounded by $|P^{\star}(\Sigma) - P^{\star}(\hat{\Sigma})| \leq \|\Sigma - \hat{\Sigma}\|_{O^{\circ}}^{\text{sym}}$. This provides a principled way to assess estimation sensitivity.

B.2 PROOF OF COROLLARY 2

Proof of Corollary 2. We apply Theorem 1 to each budget set. Let $\Sigma = U \operatorname{diag}(\sigma_1 \ge \cdots \ge \sigma_d)U^{\top}$ be the eigendecomposition of the moment matrix.

(i) Frobenius ball $Q_F(B) = \{Q \succeq 0 : \|Q\|_F \le \sqrt{B}\}$. The Lagrangian is $L(Q, \lambda) = \text{Tr}(Q\Sigma) - \lambda(\|Q\|_F^2 - B)$. Taking the gradient with respect to Q and setting to zero:

$$\nabla_Q L = \Sigma - 2\lambda Q = 0 \quad \Rightarrow \quad Q = \frac{\Sigma}{2\lambda}.$$

The constraint $||Q||_F = \sqrt{B}$ gives $||\Sigma/(2\lambda)||_F = \sqrt{B}$, so $2\lambda = ||\Sigma||_F/\sqrt{B}$. Hence:

$$Q_{\mathrm{F}}^{\star} = \sqrt{B} \, \frac{\Sigma}{\|\Sigma\|_F}, \quad P_{\mathrm{F}}^{\star}(\Sigma) = \mathrm{Tr}(Q_{\mathrm{F}}^{\star}\Sigma) = \sqrt{B} \, \|\Sigma\|_F.$$

(ii) Spectral $Q_S(\tau, \lambda) = \{Q \succeq 0 : \operatorname{Tr}(Q) \leq \tau, \ Q \preceq \lambda I\}$. By Neumann's inequality, the maximizer has the form $Q = U \operatorname{diag}(q_i)U^{\top}$ where the eigenvalues q_i solve the water-filling problem:

$$\max_{q_i \geq 0} \sum_i q_i \sigma_i \quad \text{s.t.} \quad \sum_i q_i \leq \tau, \ q_i \leq \lambda.$$

The KKT conditions yield: (i) $q_i^* = \lambda$ for $i \le k$, (ii) $q_{k+1}^* = \tau - k\lambda$, (iii) $q_i^* = 0$ for i > k+1, where $k = |\tau/\lambda|$. The optimal power is:

$$P_{\mathrm{S}}^{\star}(\Sigma) = \lambda \sum_{i \leq k} \sigma_i + (\tau - k\lambda)\sigma_{k+1}.$$

(iii) Data-metric $Q_L(B) = \{Q \succeq 0 : \text{Tr}(Q^2\Sigma) \leq B\}$. The Lagrangian is $L(Q, \mu) = \text{Tr}(Q\Sigma) - \mu(\text{Tr}(Q^2\Sigma) - B)$. The first-order condition gives:

$$\Sigma - 2\mu Q\Sigma = 0 \quad \Rightarrow \quad Q = \frac{1}{2\mu}I \text{ on } \operatorname{supp}(\Sigma).$$

Using the constraint $Tr(Q^2\Sigma) = B$ and the fact that Q is constant on the support:

$$\alpha^2 \sum_{i:\sigma_i>0} \sigma_i = B \quad \Rightarrow \quad \alpha = \sqrt{\frac{B}{\sum_{i:\sigma_i>0} \sigma_i}}.$$

Hence:

$$Q_{\mathsf{L}}^{\star} = \alpha \, \Pi_{\text{supp}(\Sigma)}, \quad P_{\mathsf{L}}^{\star}(\Sigma) = \alpha \sum_{i} \sigma_{i} = \sqrt{B \sum_{i} \sigma_{i}}.$$

(iv) Diagonal $Q_D(B,c)=\{Q=\operatorname{diag}(q_j)\succeq 0: \sum_j c_jq_j^2\leq B\}$. The problem decouples coordinate-wise:

$$\max_{q_j \ge 0} \sum_j q_j \Sigma_{jj} \quad \text{s.t.} \quad \sum_j c_j q_j^2 \le B.$$

By Cauchy-Schwarz, the maximizer satisfies $q_i^{\star} \propto \Sigma_{jj}/c_j$. Normalizing by the constraint:

$$q_j^\star = \sqrt{\frac{B}{\sum_k \Sigma_{kk}^2/c_k}} \cdot \frac{\Sigma_{jj}}{c_j}, \quad P_{\mathrm{D}}^\star(\Sigma) = \sqrt{B\sum_j \frac{\Sigma_{jj}^2}{c_j}}.$$

B.3 Proof of Theorem 3

Proof of Theorem 3. We establish each claim in turn.

- (i) Existence & sublinearity. Q is weakly compact (Hilbert spaces are reflexive; closed and bounded \Rightarrow weakly compact). The functional $Q \mapsto \langle Q, M \rangle_{\mathcal{H}}$ is continuous in the weak topology (linear functionals are weakly continuous), hence attains its maximum on Q. Sublinearity: $P^*(M) = \sup_{Q \in \mathcal{Q}} \langle Q, M \rangle_{\mathcal{H}}$ is a supremum of linear maps in M, thus convex and positively homogeneous. Finiteness follows from compactness of Q.
- (ii) Conjugacy identities. By definition of convex conjugate in \mathcal{H} ,

$$(\delta_{\mathcal{Q}})^*(M) = \sup_{Q \in \mathcal{H}} \left\{ \langle Q, M \rangle_{\mathcal{H}} - \delta_{\mathcal{Q}}(Q) \right\} = \sup_{Q \in \mathcal{Q}} \langle Q, M \rangle_{\mathcal{H}} = P^*(M).$$

Thus $P^\star = \delta_{\mathcal{Q}}^*$. By definition of the polar, $M \in (\mathcal{Q})^\circ$ iff $\sup_{Q \in \mathcal{Q}} \langle Q, M \rangle_{\mathcal{H}} \leq 1$, i.e., $P^\star(M) \leq 1$. Therefore

$$\gamma_{(Q)^{\circ}}(M) = \inf\{\lambda > 0 : M \in \lambda(Q)^{\circ}\} = \inf\{\lambda > 0 : P^{\star}(M) \le \lambda\} = P^{\star}(M).$$

Finally, $(\gamma_Q)^* = \delta_{(Q)^{\circ}}$ is the standard gauge–polar identity in a locally convex space.

(iii) Construction (subgradient). Let $Q^* \in \arg \max_{Q \in \mathcal{Q}} \langle Q, M \rangle_{\mathcal{H}}$. For any $N \in \mathcal{H}$,

$$P^{\star}(N) = \max_{Q \in \mathcal{Q}} \langle Q, N \rangle_{\mathcal{H}} \geq \langle Q^{\star}, N \rangle_{\mathcal{H}} = P^{\star}(M) + \langle Q^{\star}, N - M \rangle_{\mathcal{H}},$$

so $Q^* \in \partial_M P^*(M)$. Uniqueness implies differentiability with gradient Q^* .

- (iv) Order preservation. If $M \in \mathcal{H}_+$, then for any $Q \in \mathcal{Q} \subset \mathcal{H}_+$ we have $\langle Q, M \rangle_{\mathcal{H}} \geq 0$ (each term $\mathrm{Tr}(H_k^\top M_k) \geq 0$). Taking the max over Q yields $P^\star(M) \geq 0$. If $M_1 M_2 \in \mathcal{H}_+ + \setminus \{0\}$, some admissible Q gives strict positivity of $\langle Q, M_1 M_2 \rangle_{\mathcal{H}}$, hence the maximized value is strictly larger at M_1 than M_2 .
- (v) Lipschitz in polar gauge. From (ii), $P^* = \gamma_{(\mathcal{Q})^{\circ}}$. Then for any M, \hat{M} ,

$$\begin{split} P^{\star}(M) - P^{\star}(\hat{M}) &= \max_{Q \in \mathcal{Q}} \langle Q, M \rangle_{\mathcal{H}} - \max_{Q \in \mathcal{Q}} \langle Q, \hat{M} \rangle_{\mathcal{H}} \\ &\leq \max_{Q \in \mathcal{Q}} \langle Q, M - \hat{M} \rangle_{\mathcal{H}} \ = \ \|M - \hat{M}\|_{(\mathcal{Q})^{\circ}}. \end{split}$$

Symmetry gives the absolute value bound.

B.4 Proof of Lemma 4

Proof of Lemma 4. The proof follows from standard convex optimization theory, specifically the KKT conditions for linear maximization and the characterization of metric projections.

We use two facts:

- (KKT for linear maximization) $x^* \in \arg\max_{y \in \mathcal{C}} \langle y, M \rangle_{\mathcal{H}} \iff M \in N_{\mathcal{C}}(x^*).$
- (Metric projection) For $y \in \mathcal{H}$, $x^* = \Pi_{\mathcal{C}}(y) \iff y x^* \in N_{\mathcal{Q} \cap \mathcal{C}}(x^*)$.

 $\begin{array}{l} \hbox{\it (ii)} \Rightarrow \hbox{\it (ii)} \Rightarrow \hbox{\it (ii)} : \mbox{Suppose there exists } M \in \mathcal{Q}^{\star}_{\mathcal{C}}(R) \subseteq \mathcal{Q} \cap \mathcal{C} \mbox{ such that } \{R, Q^{\star} - M\} \subset N_{\mathcal{Q} \cap \mathcal{C}}(M). \\ \mbox{From } Q^{\star} - M \in N_{\mathcal{Q} \cap \mathcal{C}}(M), \mbox{ the metric projection characterization gives } \Pi_{\mathcal{C}}(Q^{\star}) = M. \mbox{ From } R \in N_{\mathcal{Q} \cap \mathcal{C}}(Q^{\star}_{\mathcal{C}}), \mbox{ the KKT condition for linear maximization gives } M \in \arg \max_{Q \in \mathcal{Q} \cap \mathcal{C}} \langle Q, R \rangle_{\mathcal{H}} = \mathcal{Q}^{\star}_{\mathcal{C}}(R). \end{array}$

(i) \Rightarrow (ii): Suppose $M := \Pi_{\mathcal{C}}(Q^{\star}) \in \mathcal{Q}^{\star}_{\mathcal{C}}(R)$. By the metric projection characterization, $Q^{\star} - M \in N_{\mathcal{Q} \cap \mathcal{C}}(M)$. Since $M \in \mathcal{Q}^{\star}_{\mathcal{C}}(R) = \arg\max_{Q \in \mathcal{Q} \cap \mathcal{C}} \langle Q, R \rangle_{\mathcal{H}}$, the KKT condition for linear maximization gives $R \in N_{\mathcal{Q} \cap \mathcal{C}}(M)$. Thus, $\{R, Q^{\star} - M\} \subset N_{\mathcal{Q} \cap \mathcal{C}}(M)$. Let $Q^{\star}_{\mathcal{C}} := M$.

For the final statement, if $N_{Q\cap C}(Q_C^{\star})=\{\lambda M:\lambda\geq 0\}$ is a ray, then both R and $Q^{\star}-Q_C^{\star}$ must be non-negative multiples of the same direction M for the normal-cone alignment condition to hold.

B.5 PROOF OF COROLLARY 5

Proof of Corollary 5. We work in the impulse-space Hilbert space $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ of causal LTI filters with matrix impulse response $\{q_k\}_{k\geq 0}$ and norm $\|Q\|_{\mathcal{H}}^2 = \sum_{k=0}^{\infty} \operatorname{Tr}(q_k^{\top}q_k)$.

Step 1 — Norm of the 1-pole equalizer. The impulse response is $q_k = \eta P(1-\beta)\beta^k$. By definition,

$$\|Q_{\text{SGD+M}}\|_{\mathcal{H}}^2 = \sum_{k>0} \text{Tr}(q_k^{\top} q_k) = \sum_{k>0} \text{Tr}\left([\eta P(1-\beta)\beta^k]^{\top} [\eta P(1-\beta)\beta^k]\right)$$
(43)

$$= \eta^2 (1 - \beta)^2 \sum_{k>0} \beta^{2k} \operatorname{Tr}(P^{\top} P) = \eta^2 (1 - \beta)^2 \frac{1}{1 - \beta^2} \operatorname{Tr}(P^{\top} P). \tag{44}$$

The budget constraint $||Q_{SGD+M}||_{\mathcal{H}} \leq \sqrt{B}$ imposes

$$\eta \le \sqrt{B} \left(\operatorname{Tr}(P^{\top} P) \frac{(1-\beta)^2}{1-\beta^2} \right)^{-1/2}. \tag{45}$$

Step 2 — Alignment with the moment operator. The inner product with R is

$$\langle Q_{\text{SGD+M}}, R \rangle_{\mathcal{H}} = \sum_{k>0} \text{Tr}(q_k^{\top} R_k) = \eta (1-\beta) \sum_{k>0} \beta^k \text{Tr}(P^{\top} R_k)$$
 (46)

$$= \eta(1-\beta) \sum_{k>0} \beta^k S_k, \tag{47}$$

where $S_k := \operatorname{Tr}(P^{\top} R_k)$.

Step 3 — Reduce to 1-D search; saturate budget. For fixed β , the inner product is linear in η while the constraint is quadratic, so the maximizer saturates the budget. The budget-normalized gain is

$$J(\beta) := \frac{\langle Q_{\text{SGD+M}}, R \rangle_{\mathcal{H}}}{\|Q_{\text{SGD+M}}\|_{\mathcal{H}}} = \frac{\sqrt{1 - \beta^2}}{\sqrt{\text{Tr}(P^\top P)}} \sum_{k \ge 0} \beta^k S_k. \tag{48}$$

Hence $\beta^* = \arg \max_{0 < \beta < 1} J(\beta)$ and η^* saturates the budget constraint.

B.6 PROOF OF COROLLARY 6

Proof of Corollary 6. We work in the impulse-space Hilbert space $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ of causal LTI filters with impulse response $\{H_k\}_{k\geq 0}$ and norm $\|Q\|_{\mathcal{H}}^2 = \sum_k \operatorname{Tr}(H_k^\top H_k)$.

Step 1 — Unconstrained diagonal optimizer. With the diagonal budget $Q_D(B,c)$, the maximization $\max_{Q \in \mathcal{Q}_D} \langle Q, R \rangle_{\mathcal{H}}$ decouples per coordinate into scalar H_2 subproblems:

$$\max_{q_j: \sum_j c_j \|q_j\|_{H_2}^2 \le B} \sum_{j=1}^P \langle q_j, R_{jj} \rangle_{H_2}. \tag{49}$$

By Cauchy-Schwarz, the optimizer is proportional to R_{jj} with weights $1/c_j$, giving the stated q_i^{\star} .

Step 2 — Adam-family restriction (1-pole per coordinate). Fix a fitting window and freeze the preconditioner $D_t^{-1/2}$ (quasi-LTI approximation). The per-coordinate impulse response is

$$h_{j,k} = \eta(D_t^{-1/2})_{jj} (1 - \beta_1) \beta_1^k, \qquad k \ge 0.$$
 (50)

Step 3 — Budget norm and alignment. The budget norm evaluates to

$$\|Q_{\text{Adam}}\|_{\mathcal{H}}^2 = \sum_{j=1}^P c_j \sum_{k\geq 0} |h_{j,k}|^2 = \eta^2 \frac{(1-\beta_1)^2}{1-\beta_1^2} \sum_{j=1}^P c_j (D_t^{-1/2})_{jj}^2 = \eta^2 \frac{(1-\beta_1)^2}{1-\beta_1^2} W(\beta_2). \tag{51}$$

The alignment with R is

$$\langle Q_{\text{Adam}}, R \rangle_{\mathcal{H}} = \eta (1 - \beta_1) \sum_{k>0} \beta_1^k \sum_{j=1}^P (D_t^{-1/2})_{jj} R_{jj,k} = \eta (1 - \beta_1) \sum_{k>0} \beta_1^k T_k(\beta_2). \tag{52}$$

Step 4 — Optimization over hyperparameters. For fixed (β_1, β_2) , the maximizer in η saturates the budget (linear gain under quadratic constraint). The budget-normalized gain is

$$J(\beta_1, \beta_2) := \frac{\langle Q_{\text{Adam}}, R \rangle_{\mathcal{H}}}{\|Q_{\text{Adam}}\|_{\mathcal{H}}} = \frac{\sqrt{1 - \beta_1^2}}{\sqrt{W(\beta_2)}} \sum_{k>0} \beta_1^k T_k(\beta_2).$$
 (53)

Maximizing over $(\beta_1, \beta_2) \in (0, 1)^2$ gives the optimal momentum parameters, and the optimal learning rate is

$$\eta^* = \sqrt{B} \left(\frac{(1 - \beta_1^*)^2}{1 - (\beta_1^*)^2} W(\beta_2^*) \right)^{-1/2}.$$
 (54)

B.7 Proof of Proposition 7

Proof of Proposition 7. The instantaneous validation slope is $Q \mapsto \operatorname{Tr}(Q\Sigma_{\times})$. Since this is linear in Q and Q is compact, the maximum is attained and gives the claimed inequality. If the maximizer is unique, then any suboptimal Q gives strictly smaller $\operatorname{Tr}(Q\Sigma_{\times}) = -\dot{\mathcal{L}}_{\text{val}}(Q)$. This includes optimizers only tuned to the training set.

B.8 Proof of Proposition 8

Proof of Proposition 8. In the linearized regime, we fix parameters θ_0 and consider the linearized network $f(x) \approx f(x; \theta_0) + \nabla_{\theta} f(x; \theta_0)^{\top} (\theta - \theta_0)$. For squared loss $\mathcal{L}(\theta) = \frac{1}{2} ||f_{\theta}(X) - y||^2$, the gradient is $\nabla_{\theta} \mathcal{L} = \sum_{i=1}^{n} (f_{\theta}(x_i) - y_i) \nabla_{\theta} f_{\theta}(x_i)$.

With optimizer Q, the parameter dynamics are $\dot{\theta} = -Q\nabla_{\theta}\mathcal{L}$. In the linearized regime around θ_0 , this becomes

$$\dot{\theta} = -Q \sum_{i=1}^{n} (f_{\theta_0}(x_i) + \nabla_{\theta} f(x_i; \theta_0)^{\top} (\theta - \theta_0) - y_i) \nabla_{\theta} f(x_i; \theta_0). \tag{55}$$

Let $G \in \mathbb{R}^{n \times d}$ be the matrix with rows $\nabla_{\theta} f(x_i; \theta_0)^{\top}$. Then the linearized dynamics become

$$\dot{\theta} = -QG^{\top}(G(\theta - \theta_0) + (f_{\theta_0}(X) - y)), \tag{56}$$

1026 where $f_{\theta_0}(X) = [f_{\theta_0}(x_1), \dots, f_{\theta_0}(x_n)]^{\top}$.

In function space, let $u(x) = \nabla_{\theta} f(x; \theta_0)^{\top} (\theta - \theta_0)$ represent the change in function values. Then $u(X) = G(\theta - \theta_0)$ and the dynamics become

$$\frac{d}{dt}u(X) = G\dot{\theta} = -GQG^{\top}(u(X) + (f_{\theta_0}(X) - y)). \tag{57}$$

This is kernel gradient flow in function space with kernel matrix $K = GQG^{\top}$, where $K_{ij} = \nabla_{\theta} f(x_i; \theta_0)^{\top} Q \nabla_{\theta} f(x_j; \theta_0) = K_Q(x_i, x_j)$.

By Moore–Penrose pseudoinverse, this gradient flow converges to $u^* = K^{\dagger}y$ where $K^{\dagger} = (GQG^{\top})^{\dagger}$, which corresponds to the minimum norm interpolant in the RKHS \mathcal{H}_{K_Q} induced by kernel K_Q :

$$f_Q^{\star} = \arg\min_{f \in \mathcal{H}_{K_Q}} \|f\|_{\mathcal{H}_{K_Q}}^2 \quad \text{s.t.} \quad f(X) = y.$$
 (58)

The choice of Q directly determines the kernel $K_Q(x,x') = \nabla_{\theta} f(x;\theta_0)^{\top} Q \nabla_{\theta} f(x';\theta_0)$ and hence the RKHS \mathcal{H}_{K_Q} . The validation-aware choice $Q^{\star} \in \arg\max_{Q \in \mathcal{Q}} \operatorname{Tr}(Q \Sigma_{\times})$ emphasizes directions in parameter space that are aligned with the validation cross-moment Σ_{\times} , thereby tilting the induced RKHS toward functions that perform better on validation data.

B.9 Proof of Corollary 9

Proof of Corollary 9. We apply Theorem 3 to each dynamic budget set. Let $R(e^{i\omega})=U(\omega)\operatorname{diag}(\sigma_1(\omega)\geq\cdots\geq\sigma_d(\omega))U(\omega)^*$ be the eigendecomposition of the moment operator at each frequency.

(i) Frobenius ball $Q_F(B) = \{Q : \|Q\|_{\mathcal{H}} \leq \sqrt{B}\}$. The Lagrangian is $L(Q, \lambda) = \langle Q, R \rangle_{\mathcal{H}} - \lambda(\|Q\|_{\mathcal{H}}^2 - B)$. Taking the functional derivative with respect to Q and setting to zero:

$$\delta_Q L = R - 2\lambda Q = 0 \quad \Rightarrow \quad Q = \frac{R}{2\lambda}$$

The constraint $\|Q\|_{\mathcal{H}} = \sqrt{B}$ gives $\|R/(2\lambda)\|_{\mathcal{H}} = \sqrt{B}$, so $2\lambda = \|R\|_{\mathcal{H}}/\sqrt{B}$. Hence:

$$Q_F^{\star} = \sqrt{B} \frac{R}{\|R\|_{\mathcal{H}}}, \quad P_F^{\star}(R) = \langle Q_F^{\star}, R \rangle_{\mathcal{H}} = \sqrt{B} \|R\|_{\mathcal{H}}.$$

(ii) Per-frequency spectral $Q_S(\tau,\lambda)=\{Q: \operatorname{Tr}(Q(e^{i\omega}))\leq \tau(\omega), Q(e^{i\omega}) \leq \lambda(\omega)I\}$. By the dynamic version of Neumann's inequality, the maximizer has the form $Q(e^{i\omega})=U(\omega)\operatorname{diag}(q_i(\omega))U(\omega)^*$ where the eigenvalues $q_i(\omega)$ solve the water-filling problem at each frequency:

$$\max_{q_i(\omega) \ge 0} \sum_i q_i(\omega) \sigma_i(\omega) \quad \text{s.t.} \quad \sum_i q_i(\omega) \le \tau(\omega), \ q_i(\omega) \le \lambda(\omega).$$

The KKT conditions yield: (i) $q_i^\star(\omega) = \lambda(\omega)$ for $i \leq k(\omega)$, (ii) $q_{k+1}^\star(\omega) = \tau(\omega) - k(\omega)\lambda(\omega)$, (iii) $q_i^\star(\omega) = 0$ for $i > k(\omega) + 1$, where $k(\omega) = \lfloor \tau(\omega)/\lambda(\omega) \rfloor$. The optimal power is:

$$P_S^{\star}(R) = \frac{1}{2\pi} \int_0^{2\pi} \left[\lambda(\omega) \sum_{i \le k(\omega)} \sigma_i(\omega) + (\tau(\omega) - k(\omega)\lambda(\omega)) \sigma_{k(\omega) + 1}(\omega) \right] d\omega.$$

(iii) Data-metric $Q_L(B) = \{Q : \sum_{k=0}^{\infty} \operatorname{Tr}(Q_k^* \Sigma Q_k) \leq B\}$. The Lagrangian is $L(Q, \mu) = \langle Q, R \rangle_{\mathcal{H}} - \mu \left(\sum_{k=0}^{\infty} \operatorname{Tr}(Q_k^* \Sigma Q_k) - B \right)$. The first-order condition gives:

$$R - 2\mu \sum_{k=0}^{\infty} \Sigma Q_k e^{-ik\omega} = 0 \quad \Rightarrow \quad Q = \frac{1}{2\mu} I \text{ on supp}(R).$$

Using the constraint and the fact that Q is constant on the support:

$$\alpha^2 \sum_{i:\sigma_i>0} \sigma_i = B \quad \Rightarrow \quad \alpha = \sqrt{\frac{B}{\sum_{i:\sigma_i>0} \sigma_i}}.$$

Hence:

$$Q_L^{\star} = \alpha \prod_{\text{supp}(R)}, \quad P_L^{\star}(R) = \alpha \sum_i \sigma_i = \sqrt{B \sum_i \sigma_i}.$$

(iv) Diagonal $Q_D(B,c) = \{Q(z) = \operatorname{diag}(q_j(z)) \succeq 0 : \sum_j c_j ||q_j||_{H_2}^2 \leq B\}$. The problem decouples coordinate-wise into scalar H_2 subproblems:

$$\max_{q_j \geq 0} \sum_j \langle q_j, R_{jj} \rangle_{H_2} \quad \text{s.t.} \quad \sum_j c_j \|q_j\|_{H_2}^2 \leq B.$$

By Cauchy-Schwarz in H_2 , the maximizer satisfies $q_i^* \propto R_{jj}/c_j$. Normalizing by the constraint:

$$q_j^\star = \sqrt{\frac{B}{\sum_k \|R_{kk}\|_{H_2}^2/c_k}} \cdot \frac{R_{jj}}{c_j}, \quad P_D^\star(R) = \sqrt{B\sum_j \frac{\|R_{jj}\|_{H_2}^2}{c_j}}.$$

C REVERSE ENGINEERING COMMON OPTIMIZERS

So far, we have derived various types of stateless and stateful optimizers under different types of budgets. In this section, we will do the opposite: we will *reverse engineer* popular optimizers and find out under which budget they are secretly optimizing. This not only allows us to find out hidden design principles of these optimizers, but also have these optimizers registered in a unified framework, suggesting a systematic way to design new optimizers.

tl:dr:

- GD = Euclidean budget $\Rightarrow Q \propto I$.
- Colored-GD = elliptic budget $\Rightarrow Q \propto P$.
- Newton/GN = curvature budget $\Rightarrow Q \propto H^{-1}/G^{-1}$.
- NGD = Fisher/KL budget $\Rightarrow Q \propto F^{-1}$.
- K-FAC/Shampoo = structured budgets \Rightarrow block/Kronecker Q.
- $AdaGrad/RMSProp/Adam = diagonal budgets (on q or <math>m_k$) $\Rightarrow Q diagonal$.
- LAMB/LARS = layer-norm budget \Rightarrow layer-wise scalar Q.
- $signSGD/Lion = L_{\infty}$ budget (on g or m_k) \Rightarrow normalized/sign steps.

One lens. Every first-order optimizer picks a velocity $\dot{\theta}$ from the current gradient g by solving a budgeted power allocation:

$$\dot{\theta} = Q g \quad \text{with} \quad Q \in \mathcal{Q} \iff \dot{\theta} = \arg \max_{\dot{\theta} \in \mathbb{R}^P} \langle g, \dot{\theta} \rangle \text{ s.t. } \text{budget}(\dot{\theta}) \le \tau$$
 (59)

The *budget* determines the *optimizer* Q. Below, we list each popular optimizer as a special case of this formulation, give the induced Q, and state (not prove) the short KKT step that produces it.

C.1 EUCLIDEAN & CURVATURE FAMILIES

Corollary 10 (SGD from Euclidean Frobenius budget). *Define the Euclidean Frobenius budget and the cone of memoryless isotropic optimizers:*

$$Q_F(B) := \{Q : ||Q||_{\mathcal{H}} \le \sqrt{B}\}, \quad \mathcal{C}_{memoryless} := \{Q[n] = \eta I\delta[n] : \eta \ge 0\}. \tag{60}$$

Given current gradient moment $R[0] = gg^{\top}$ where g is the instantaneous gradient, solving problem P3 under the budget $Q_F(B) \cap C_{memoryless}$ produces SGD as the optimal solution with optimal hyperparameter:

$$\eta^{\star} = \frac{\sqrt{B}}{\sqrt{d}||g||_2},\tag{61}$$

where d is the dimension of the parameter space.

Proof. We work in the impulse-space Hilbert space $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ of causal LTI filters with matrix impulse response Q and Frobenius \mathcal{H} norm $\|Q\|_{\mathcal{H}}^2 = \sum_{n=0}^{\infty} \operatorname{Tr}(Q[n]^{\top}Q[n])$.

Step 1 — Global optimal equalizer. By Cauchy-Schwarz, the unconstrained optimum of problem P3 under budget $Q_F(B) = \{Q : ||Q||_{\mathcal{H}} \leq \sqrt{B}\}$ is

$$Q^{C} = \sqrt{B} \frac{R}{\|R\|_{\mathcal{H}}}, \qquad P^{*}(R) = \sqrt{B} \|R\|_{\mathcal{H}}.$$
 (62)

Step 2 — Commutativity via smooth convexity. The Frobenius ball $Q_F(B)$ is smooth and strictly convex. At any boundary point Q with $||Q||_{\mathcal{H}} = \sqrt{B}$, the normal cone is the ray $N_{Q_F(B)}(Q) = \{\lambda Q : \lambda \geq 0\}$. By Lemma 4, commutativity holds:

$$\Pi_{\mathcal{Q}_F(B)\cap\mathcal{C}_{\text{memoryless}}}(Q^C) \in \arg\max_{Q\in\mathcal{Q}_F(B)\cap\mathcal{C}_{\text{memoryless}}} \langle Q, R \rangle_{\mathcal{H}}.$$
 (63)

Therefore, we can equivalently solve the restricted optimization problem directly.

Step 3 — Memoryless family parametrization. For $Q_{\eta}[n] = \eta I\delta[n]$ with $\eta \geq 0$, we compute:

$$||Q_n||_{\mathcal{H}}^2 = \text{Tr}((\eta I)^\top (\eta I)) = \eta^2 d,$$
 (64)

and for the instantaneous gradient moment $R[0] = gg^{\top}$ and R[n] = 0 for n > 0:

$$\langle Q_{\eta}, R \rangle_{\mathcal{H}} = \text{Tr}((\eta I)^{\top}(gg^{\top})) = \eta \, \text{Tr}(gg^{\top}) = \eta \|g\|_{2}^{2}. \tag{65}$$

Step 4 — Budget saturation and optimization. The objective is linear in η while the constraint is quadratic, so the maximizer saturates the budget $||Q_n||_{\mathcal{H}} = \sqrt{B}$. This gives:

$$\eta^{\star} = \frac{\sqrt{B}}{\sqrt{d}}.\tag{66}$$

However, when we normalize by the gradient magnitude for scale invariance, we obtain:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{d}||g||_2}.\tag{67}$$

Step 5 — Geometric interpretation. SGD emerges as the memoryless isotropic approximation of the global Frobenius-constrained equalizer. The optimal learning rate η^* balances the budget constraint with the current gradient magnitude, providing uniform scaling across all parameter dimensions. By commutativity, this restricted optimum coincides with projecting the global optimum Q^C onto $\mathcal{C}_{\text{memoryless}}$.

Corollary 11 (Colored gradient descent from elliptic trust region budget). *Define the elliptic budget* $Q_P(B)$ *and the cone* $C_{memoryless}$ *of memoryless optimizers:*

$$Q_P(B) := \{Q : \text{Tr}(Q[0]^\top P^{-1}Q[0]) \le B, Q[n] = 0 \text{ for } n > 0\},$$
 (68)

$$C_{memoryless} := \{Q[n] = \eta I \delta[n] : \eta \ge 0\}. \tag{69}$$

where $P \succ 0$ is a fixed symmetric positive definite matrix. Given current gradient moment $R[0] = gg^{\top}$ where g is the instantaneous gradient, solving problem P3 under the budget $Q_P(B) \cap C_{memoryless}$ produces colored gradient descent as the optimal solution with optimal hyperparameter:

$$\eta^{\star} = \frac{\sqrt{B}}{\sqrt{\text{Tr}(P^{-1})} \|g\|_2},\tag{70}$$

where the optimizer is $Q^* = \eta^* P$.

Proof. We work in the impulse-space Hilbert space $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ of causal LTI filters with matrix impulse response Q and elliptic P-weighted norm for memoryless filters $||Q[0]||^2_{P^{-1}} = \text{Tr}(Q[0]^\top P^{-1}Q[0])$.

Step 1 — Global optimal equalizer. By Cauchy-Schwarz, the unconstrained optimum of problem P3 under budget $Q_P(B) = \{Q : \operatorname{Tr}(Q[0]^\top P^{-1}Q[0]) \leq B, Q[n] = 0 \text{ for } n > 0\}$ is

$$Q^{C}[0] = \sqrt{B} \frac{Pgg^{\top}}{\|Pgg^{\top}\|_{P^{-1}}}, \qquad P^{\star}(R) = \sqrt{B} \|Pgg^{\top}\|_{P^{-1}}. \tag{71}$$

Step 2 — Commutativity via smooth convexity. The elliptic ball $\mathcal{Q}_P(B)$ is smooth and strictly convex. At any boundary point Q with $\|Q[0]\|_{P^{-1}} = \sqrt{B}$, the normal cone is the ray $N_{\mathcal{Q}_P(B)}(Q) = \{\lambda P^{-1}Q[0] : \lambda \geq 0\}$. By Lemma 4, commutativity holds:

$$\Pi_{Q_P(B) \cap \mathcal{C}_{\text{memoryless}}}(Q^C) \in \arg \max_{Q \in Q_P(B) \cap \mathcal{C}_{\text{memoryless}}} \langle Q, R \rangle_{\mathcal{H}}. \tag{72}$$

Therefore, we can equivalently solve the restricted optimization problem directly.

Step 3 — Memoryless family parametrization. For $Q_n[n] = \eta P\delta[n]$ with $\eta \geq 0$, we compute:

$$||Q_{\eta}[0]||_{P^{-1}}^{2} = \text{Tr}((\eta P)^{\top} P^{-1}(\eta P)) = \eta^{2} \text{Tr}(P),$$
 (73)

and for the instantaneous gradient moment $R[0] = gg^{\top}$ and R[n] = 0 for n > 0:

$$\langle Q_{\eta}, R \rangle_{\mathcal{H}} = \text{Tr}((\eta P)^{\top}(gg^{\top})) = \eta \, \text{Tr}(Pgg^{\top}) = \eta g^{\top} Pg.$$
 (74)

Step 4 — Budget saturation and optimization. The objective is linear in η while the constraint is quadratic, so the maximizer saturates the budget $\|Q_{\eta}[0]\|_{P^{-1}} = \sqrt{B}$. This gives:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{\text{Tr}(P)}}. (75)$$

However, when we normalize by the gradient magnitude for scale invariance, we obtain:

$$\eta^* = \frac{\sqrt{\overline{B}}}{\sqrt{\text{Tr}(P^{-1})} \|g\|_2}.$$
(76)

Step 5 — Geometric interpretation. Colored gradient descent emerges as the memoryless approximation of the global elliptic-constrained equalizer. The optimal learning rate η^* balances the elliptic budget constraint with the current gradient magnitude, providing P-weighted scaling across parameter dimensions. The optimizer $Q^* = \eta^* P$ naturally incorporates the geometry encoded in matrix P. By commutativity, this restricted optimum coincides with projecting the global optimum Q^C onto $\mathcal{C}_{\text{memoryless}}$.

Corollary 12 (Newton's method from curvature-aware budget). *Define the budget* $Q_H(B)$ *and the cone* $C_{memoryless}$ *of memoryless optimizers:*

$$Q_H(B) := \{Q : \text{Tr}(Q[0]^\top H Q[0]) \le B, Q[n] = 0 \text{ for } n > 0\},$$
 (77)

$$C_{memoryless} := \{Q_{\eta}[n] = \eta P \delta[n] : \eta \ge 0, P \succ 0\}. \tag{78}$$

where $H \succ 0$ is the Hessian matrix. Given current gradient moment $R[0] = gg^{\top}$ where g is the instantaneous gradient, solving problem P3 under the budget $\mathcal{Q}_H(B) \cap \mathcal{C}_{memoryless}$ produces Newton's method as the optimal solution with optimal hyperparameter:

$$\eta^{\star} = \frac{\sqrt{B}}{\sqrt{\text{Tr}(H^{-1})} \|g\|_2},\tag{79}$$

where the optimizer is $Q^* = \eta^* H^{-1}$.

Proof. We work in the impulse-space Hilbert space $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ of causal LTI filters with matrix impulse response Q and elliptic H-weighted norm for memoryless filters $||Q[0]||_H^2 = \text{Tr}(Q[0]^\top HQ[0])$.

Step 1 — Global optimal equalizer. By Cauchy-Schwarz, the unconstrained optimum of problem P3 under budget $Q_H(B) = \{Q : \operatorname{Tr}(Q[0]^\top HQ[0]) \leq B, Q[n] = 0 \text{ for } n > 0\}$ is

$$Q^{C}[0] = \sqrt{B} \frac{H^{-1}gg^{\top}}{\|H^{-1}gg^{\top}\|_{H}}, \qquad P^{\star}(R) = \sqrt{B} \|H^{-1}gg^{\top}\|_{H}.$$
 (80)

Step 2 — Commutativity via smooth convexity. The elliptic ball $\mathcal{Q}_H(B)$ is smooth and strictly convex. At any boundary point Q with $\|Q[0]\|_H = \sqrt{B}$, the normal cone is the ray $N_{\mathcal{Q}_H(B)}(Q) = \{\lambda H Q[0] : \lambda \geq 0\}$. By Lemma 4, commutativity holds:

$$\Pi_{\mathcal{Q}_H(B)\cap\mathcal{C}_{\text{memoryless}}}(Q^C) \in \arg\max_{Q\in\mathcal{Q}_H(B)\cap\mathcal{C}_{\text{memoryless}}} \langle Q, R \rangle_{\mathcal{H}}.$$
 (81)

Therefore, we can equivalently solve the restricted optimization problem directly.

Step 3 — Memoryless family parametrization. For $Q_{\eta}[n] = \eta H^{-1}\delta[n]$ with $\eta \geq 0$, we compute:

$$||Q_{\eta}[0]||_{H}^{2} = \operatorname{Tr}((\eta H^{-1})^{\top} H(\eta H^{-1})) = \eta^{2} \operatorname{Tr}(H^{-1}), \tag{82}$$

and for the instantaneous gradient moment $R[0] = gg^{\top}$ and R[n] = 0 for n > 0:

$$\langle Q_{\eta}, R \rangle_{\mathcal{H}} = \operatorname{Tr}((\eta H^{-1})^{\top} (gg^{\top})) = \eta \operatorname{Tr}(H^{-1}gg^{\top}) = \eta g^{\top} H^{-1}g.$$
(83)

Step 4 — Budget saturation and optimization. The objective is linear in η while the constraint is quadratic, so the maximizer saturates the budget $||Q_{\eta}[0]||_H = \sqrt{B}$. This gives:

$$\eta^{\star} = \frac{\sqrt{B}}{\sqrt{\text{Tr}(H^{-1})}}.$$
(84)

When we normalize by the gradient magnitude for scale invariance, we obtain:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{\text{Tr}(H^{-1})} \|g\|_2}.$$
(85)

Step 5 — Geometric interpretation. Newton's method emerges as the memoryless approximation of the global Hessian-constrained equalizer. The optimal learning rate η^* balances the curvature budget constraint with the current gradient magnitude, providing Hessian-weighted scaling that naturally incorporates second-order geometry. The optimizer $Q^* = \eta^* H^{-1}$ captures the local quadratic structure of the loss landscape. By commutativity, this restricted optimum coincides with projecting the global optimum Q^C onto $\mathcal{C}_{\text{memoryless}}$.

Corollary 13 (L-BFGS from learned curvature approximation). *Define the budget* $Q_B(B)$ *and the cone* $C_{memoryless}$ *of memoryless optimizers:*

$$Q_B(B) := \{Q : \text{Tr}(Q[0]^\top B_k Q[0]) \le B, Q[n] = 0 \text{ for } n > 0\},$$
(86)

$$C_{memoryless} := \{Q_{\eta}[n] = \eta B_k^{-1} \delta[n] : \eta \ge 0\}. \tag{87}$$

where $B_k \succ 0$ is the L-BFGS curvature approximation matrix constructed through secant updates. Given current gradient moment $R[0] = gg^{\top}$ where g is the instantaneous gradient, solving problem P3 under the budget $\mathcal{Q}_B(B) \cap \mathcal{C}_{memoryless}$ produces L-BFGS as the optimal solution with optimal hyperparameter:

$$\eta^{\star} = \frac{\sqrt{B}}{\sqrt{\text{Tr}(B_k^{-1})} \|g\|_2},\tag{88}$$

where the optimizer is $Q^* = \eta^* B_k^{-1}$.

Proof. We work in the impulse-space Hilbert space $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ of causal LTI filters with matrix impulse response Q and elliptic B_k -weighted norm for memoryless filters $\|Q[0]\|_{B_k}^2 = \text{Tr}(Q[0]^\top B_k Q[0])$.

Step 1 — Global optimal equalizer. By Cauchy-Schwarz, the unconstrained optimum of problem P3 under budget $Q_B(B) = \{Q : \operatorname{Tr}(Q[0]^\top B_k Q[0]) \leq B, Q[n] = 0 \text{ for } n > 0\}$ is

$$Q^{C}[0] = \sqrt{B} \frac{B_{k}^{-1} g g^{\top}}{\|B_{k}^{-1} g g^{\top}\|_{B_{k}}}, \qquad P^{\star}(R) = \sqrt{B} \|B_{k}^{-1} g g^{\top}\|_{B_{k}}.$$
 (89)

Step 2 — Commutativity via smooth convexity. The elliptic ball $Q_B(B)$ is smooth and strictly convex. At any boundary point Q with $\|Q[0]\|_{B_k} = \sqrt{B}$, the normal cone is the ray $N_{Q_B(B)}(Q) = \{\lambda B_k Q[0] : \lambda \geq 0\}$. By Lemma 4, commutativity holds:

$$\Pi_{\mathcal{Q}_B(B)\cap\mathcal{C}_{\text{memoryless}}}(Q^C) \in \arg\max_{Q\in\mathcal{Q}_B(B)\cap\mathcal{C}_{\text{memoryless}}} \langle Q, R \rangle_{\mathcal{H}}. \tag{90}$$

Therefore, we can equivalently solve the restricted optimization problem directly.

Step 3 — Memoryless family parametrization. For $Q_n[n] = \eta B_k^{-1} \delta[n]$ with $\eta \ge 0$, we compute:

$$||Q_{\eta}[0]||_{B_{h}}^{2} = \text{Tr}((\eta B_{k}^{-1})^{\top} B_{k}(\eta B_{k}^{-1})) = \eta^{2} \text{Tr}(B_{k}^{-1}), \tag{91}$$

and for the instantaneous gradient moment $R[0] = gg^{\top}$ and R[n] = 0 for n > 0:

$$\langle Q_{\eta}, R \rangle_{\mathcal{H}} = \operatorname{Tr}((\eta B_k^{-1})^{\top} (gg^{\top})) = \eta \operatorname{Tr}(B_k^{-1} gg^{\top}) = \eta g^{\top} B_k^{-1} g. \tag{92}$$

Step 4 — Budget saturation and optimization. The objective is linear in η while the constraint is quadratic, so the maximizer saturates the budget $||Q_{\eta}[0]||_{B_k} = \sqrt{B}$. This gives:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{\text{Tr}(B_k^{-1})}}. (93)$$

When we normalize by the gradient magnitude for scale invariance, we obtain:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{\text{Tr}(B_k^{-1})} \|g\|_2}.$$
(94)

Step 5 — Geometric interpretation. L-BFGS emerges as the memoryless approximation of the global curvature-constrained equalizer, where the curvature matrix B_k is learned through secant updates rather than computed exactly. The optimal learning rate η^* balances the learned curvature budget constraint with the current gradient magnitude, providing B_k -weighted scaling that incorporates approximate second-order geometry at reduced computational cost. The optimizer $Q^* = \eta^* B_k^{-1}$ captures the accumulated curvature information from the optimization trajectory. By commutativity, this restricted optimum coincides with projecting the global optimum Q^C onto $\mathcal{C}_{\text{memoryless}}$.

C.2 Information-geometric & structured optimizers

Corollary 14 (Natural Gradient Descent from Fisher information geometry). *Define the Fisher information budget* $Q_F(B, F)$ *and the cone* $C_{memoryless}$ *of memoryless optimizers:*

$$Q_F(B, F) := \{Q : \text{Tr}(Q^\top F Q) \le B\}, \quad \mathcal{C}_{memoryless} := \{Q_n = \eta I : \eta \ge 0\}. \tag{95}$$

Given instantaneous gradient moment $R[0] = gg^{\top}$ and R[n] = 0 for n > 0, solving problem P3 under the budget $Q_F(B, F) \cap C_{memoryless}$ produces Natural Gradient Descent as the optimal solution with optimal hyperparameters:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{\text{Tr}(F^{-1})} \|g\|_2}, \qquad Q^* = \eta^* F^{-1}. \tag{96}$$

Proof. We work in the impulse-space Hilbert space $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ of causal LTI filters with matrix impulse response Q and Fisher-weighted norm for memoryless filters $\|Q\|_F^2 = \text{Tr}(Q^\top FQ)$.

Step 1 — Global optimal equalizer. By Cauchy-Schwarz, the unconstrained optimum of problem P3 under budget $Q_F(B,F) = \{Q : \operatorname{Tr}(Q^\top FQ) \leq B\}$ is

$$Q^{C} = \sqrt{B} \frac{F^{-1}gg^{\top}}{\|F^{-1}gg^{\top}\|_{F}}, \qquad P^{*}(R) = \sqrt{B} \|F^{-1}gg^{\top}\|_{F}.$$
(97)

Step 2 — Commutativity via smooth convexity. The Fisher ellipsoid $Q_F(B,F)$ is smooth and strictly convex. At any boundary point Q with $\|Q\|_F = \sqrt{B}$, the normal cone is the ray $N_{Q_F(B,F)}(Q) = \{\lambda FQ : \lambda \geq 0\}$. By Lemma 4, commutativity holds:

$$\Pi_{\mathcal{Q}_F(B,F)\cap\mathcal{C}_{\text{memoryless}}}(Q^C) \in \arg\max_{Q\in\mathcal{Q}_F(B,F)\cap\mathcal{C}_{\text{memoryless}}} \langle Q,R\rangle_{\mathcal{H}}.$$
(98)

Therefore, we can equivalently solve the restricted optimization problem directly.

Step 3 — Memoryless family parametrization. For $Q_{\eta} = \eta I$ with $\eta \geq 0$, we compute:

$$||Q_{\eta}||_F^2 = \text{Tr}((\eta I)^{\top} F(\eta I)) = \eta^2 \text{Tr}(F),$$
 (99)

and for the instantaneous gradient moment $R[0] = gg^{\top}$ and R[n] = 0 for n > 0:

$$\langle Q_{\eta}, R \rangle_{\mathcal{H}} = \text{Tr}((\eta I)^{\top} (gg^{\top})) = \eta \, \text{Tr}(gg^{\top}) = \eta g^{\top} g = \eta \|g\|_{2}^{2}. \tag{100}$$

Step 4 — Budget saturation and optimization. The objective is linear in η while the constraint is quadratic, so the maximizer saturates the budget $||Q_{\eta}||_F = \sqrt{B}$. This gives:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{\text{Tr}(F)}}. (101)$$

When we normalize by the gradient magnitude for scale invariance, we obtain:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{\text{Tr}(F)} \|g\|_2}.$$
(102)

Step 5 — Geometric interpretation. Natural Gradient Descent emerges as the memoryless approximation of the global Fisher-constrained equalizer, where the Fisher information matrix F captures the intrinsic Riemannian geometry of the statistical model. The optimal learning rate η^* balances the Fisher information budget constraint with the current gradient magnitude, providing F-weighted scaling that incorporates the natural geometry of the parameter space. The optimizer $Q^* = \eta^* F^{-1}$ implements steepest descent in the natural Riemannian metric, where the Fisher metric measures the intrinsic difficulty of distinguishing nearby parameter values based on the data distribution. By commutativity, this restricted optimum coincides with projecting the global optimum Q^C onto $\mathcal{C}_{\text{memoryless}}$.

Corollary 15 (K-FAC from block-diagonal Fisher approximation). *Define the block-diagonal Fisher budget* $Q_{block}(B, \{F_{\ell}\})$ *and the cone* $C_{memoryless}$ *of memoryless optimizers:*

$$Q_{block}(B, \{F_{\ell}\}) := \{Q : \sum_{\ell} \|Q_{\ell}\|_{F_{\ell}}^2 \le B\}, \quad C_{memoryless} := \{Q[n] = \eta I\delta[n] : \eta \ge 0\}.$$
 (103)

Given instantaneous gradient moment $R[0] = gg^{\top}$ and R[n] = 0 for n > 0, solving problem P3 under the budget $\mathcal{Q}_{block}(B, \{F_{\ell}\}) \cap \mathcal{C}_{memoryless}$ produces K-FAC as the optimal solution with optimal hyperparameters:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{\sum_{\ell} \operatorname{Tr}(F_{\ell}^{-1})} \|g\|_2}, \qquad Q^* = \eta^* \operatorname{blockdiag}(F_{\ell}^{-1}). \tag{104}$$

Proof. We work in the impulse-space Hilbert space $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ of causal LTI filters with matrix impulse response Q and block-diagonal Fisher-weighted norm $\|Q\|_{\text{block}}^2 = \sum_{\ell} \text{Tr}(Q_{\ell}^{\top} F_{\ell} Q_{\ell})$.

Step 1 — Global optimal equalizer. By Cauchy-Schwarz, the unconstrained optimum of problem P3 under budget $\mathcal{Q}_{block}(B, \{F_\ell\}) = \{Q : \sum_{\ell} \|Q_\ell\|_{F_\ell}^2 \leq B\}$ is

$$Q^{C} = \sqrt{B} \frac{\text{blockdiag}(F_{\ell}^{-1})gg^{\top}}{\|\text{blockdiag}(F_{\ell}^{-1})gg^{\top}\|_{\text{block}}}, \qquad P^{\star}(R) = \sqrt{B} \|\text{blockdiag}(F_{\ell}^{-1})gg^{\top}\|_{\text{block}}.$$
(105)

Step 2 — Commutativity via smooth convexity. The block-diagonal Fisher ellipsoid $\mathcal{Q}_{\mathrm{block}}(B,\{F_\ell\})$ is smooth and strictly convex. At any boundary point Q with $\|Q\|_{\mathrm{block}} = \sqrt{B}$, the normal cone is the ray $N_{\mathcal{Q}_{\mathrm{block}}(B,\{F_\ell\})}(Q) = \{\lambda \mathrm{blockdiag}(F_\ell)Q: \lambda \geq 0\}$. By Lemma 4, commutativity holds:

$$\Pi_{\mathcal{Q}_{\text{block}}(B, \{F_{\ell}\}) \cap \mathcal{C}_{\text{memoryless}}}(Q^C) \in \arg \max_{Q \in \mathcal{Q}_{\text{block}}(B, \{F_{\ell}\}) \cap \mathcal{C}_{\text{memoryless}}} \langle Q, R \rangle_{\mathcal{H}}. \tag{106}$$

Therefore, we can equivalently solve the restricted optimization problem directly.

Step 3 — Memoryless family parametrization. For $Q_{\eta} = \eta I$ with $\eta \geq 0$, we compute:

$$||Q_{\eta}||_{\text{block}}^2 = \sum_{\ell} \text{Tr}((\eta I_{\ell})^{\top} F_{\ell}(\eta I_{\ell})) = \eta^2 \sum_{\ell} \text{Tr}(F_{\ell}), \tag{107}$$

and for the instantaneous gradient moment $R[0] = gg^{\top}$ and R[n] = 0 for n > 0:

$$\langle Q_{\eta}, R \rangle_{\mathcal{H}} = \text{Tr}((\eta I)^{\top}(gg^{\top})) = \eta \text{Tr}(gg^{\top}) = \eta ||g||_{2}^{2}. \tag{108}$$

Step 4 — Budget saturation and optimization. The objective is linear in η while the constraint is quadratic, so the maximizer saturates the budget $\|Q_{\eta}\|_{\text{block}} = \sqrt{B}$. This gives:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{\sum_{\ell} \text{Tr}(F_{\ell})}}.$$
 (109)

When we normalize by the gradient magnitude for scale invariance, we obtain:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{\sum_{\ell} \text{Tr}(F_{\ell})} \|g\|_2}.$$
(110)

Step 5 — Geometric interpretation. K-FAC emerges as the memoryless approximation of the global block-diagonal Fisher-constrained equalizer, where each block F_ℓ captures the layer-wise Fisher information geometry. The optimal learning rate η^\star balances the block-diagonal Fisher information budget constraint with the current gradient magnitude, providing layer-wise F_ℓ -weighted scaling that incorporates the natural geometry while maintaining computational tractability through block-diagonal structure. The optimizer $Q^\star = \eta^\star \mathrm{blockdiag}(F_\ell^{-1})$ implements approximate steepest descent in the block-diagonal natural Riemannian metric. By commutativity, this restricted optimum coincides with projecting the global optimum Q^C onto $\mathcal{C}_{\mathrm{memoryless}}$.

Corollary 16 (Shampoo from Kronecker-factored preconditioning). Consider weight tensors $\theta \in \mathbb{R}^{d_1 \times d_2 \times \cdots \times d_k}$ and define mode-wise second moment matrices $G_i \in \mathbb{R}^{d_i \times d_i}$ for each mode i. Define the Kronecker-factored budget constraint:

$$Q_{Kron}(B, \{G_i\}) := \left\{ Q : \|Q\|_{Kron}^2 \le B \right\}, \tag{111}$$

where $\|Q\|_{\mathrm{Kron}}^2 = \sum_{n=0}^{\infty} \operatorname{Tr}\left(Q[n]^{\top}\left(\bigotimes_i G_i\right)Q[n]\right)$, and the memoryless Kronecker cone:

$$C_{Kron} := \left\{ Q[n] = \eta \left(\bigotimes_{i} G_{i}^{-1/2} \right) \delta[n] : \eta \ge 0 \right\}. \tag{112}$$

Given instantaneous gradient moment $R[0] = gg^{\top}$ and R[n] = 0 for n > 0, solving problem P3 under the budget $\mathcal{Q}_{Kron}(B, \{G_i\}) \cap \mathcal{C}_{Kron}$ produces Shampoo as the optimal solution with optimal hyperparameters:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{\text{Tr}\left(\bigotimes_i G_i^{-1}\right)} \|g\|_2}, \qquad Q^* = \eta^* \bigotimes_i G_i^{-1/2}. \tag{113}$$

Proof. We work in the impulse-space Hilbert space $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ of causal LTI filters with matrix impulse response Q and Kronecker-weighted norm $\|Q\|_{\mathrm{Kron}}^2 = \sum_{n=0}^{\infty} \mathrm{Tr}\left(Q[n]^{\top}\left(\bigotimes_i G_i\right)Q[n]\right)$.

Step 1 — Global optimal equalizer. By Cauchy-Schwarz, the unconstrained optimum of problem P3 under budget $\mathcal{Q}_{Kron}(B, \{G_i\}) = \{Q : \|Q\|_{Kron} \leq \sqrt{B}\}$ is

$$Q^{C} = \sqrt{B} \frac{\left(\bigotimes_{i} G_{i}^{-1}\right) g g^{\top}}{\left\|\left(\bigotimes_{i} G_{i}^{-1}\right) g g^{\top}\right\|_{\text{Kron}}}, \qquad P^{\star}(R) = \sqrt{B} \left\|\left(\bigotimes_{i} G_{i}^{-1}\right) g g^{\top}\right\|_{\text{Kron}}. \tag{114}$$

Step 2 — Commutativity via smooth convexity. The Kronecker-factored ellipsoid $\mathcal{Q}_{Kron}(B, \{G_i\})$ is smooth and strictly convex. At any boundary point Q with $\|Q\|_{Kron} = \sqrt{B}$, the normal cone is the ray $N_{\mathcal{Q}_{Kron}(B, \{G_i\})}(Q) = \{\lambda (\bigotimes_i G_i) \ Q : \lambda \geq 0\}$. By Lemma 4, commutativity holds:

$$\Pi_{\mathcal{Q}_{Kron}(B,\{G_i\})\cap\mathcal{C}_{Kron}}(Q^C) \in \arg\max_{Q \in \mathcal{Q}_{Kron}(B,\{G_i\})\cap\mathcal{C}_{Kron}} \langle Q,R \rangle_{\mathcal{H}}.$$
(115)

Therefore, we can equivalently solve the restricted optimization problem directly.

Step 3 — Memoryless Kronecker family parametrization. For $Q_{\eta} = \eta \bigotimes_{i} G_{i}^{-1/2}$ with $\eta \geq 0$, we compute:

$$\|Q_{\eta}\|_{\mathrm{Kron}}^{2} = \operatorname{Tr}\left(\left(\eta \bigotimes_{i} G_{i}^{-1/2}\right)^{\top} \left(\bigotimes_{i} G_{i}\right) \left(\eta \bigotimes_{i} G_{i}^{-1/2}\right)\right) = \eta^{2} \operatorname{Tr}\left(\bigotimes_{i} G_{i}^{-1}\right), \quad (116)$$

and for the instantaneous gradient moment $R[0] = gg^{\top}$ and R[n] = 0 for n > 0:

$$\langle Q_{\eta}, R \rangle_{\mathcal{H}} = \operatorname{Tr}\left(\left(\eta \bigotimes_{i} G_{i}^{-1/2}\right)^{\top} (gg^{\top})\right) = \eta \operatorname{Tr}\left(\left(\bigotimes_{i} G_{i}^{-1/2}\right) gg^{\top}\right).$$
 (117)

Step 4 — Budget saturation and optimization. The objective is linear in η while the constraint is quadratic, so the maximizer saturates the budget $||Q_{\eta}||_{Kron} = \sqrt{B}$. This gives:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{\text{Tr}\left(\bigotimes_i G_i^{-1}\right)}}.$$
(118)

When we normalize by the gradient magnitude for scale invariance, we obtain:

$$\eta^* = \frac{\sqrt{\overline{B}}}{\sqrt{\text{Tr}\left(\bigotimes_i G_i^{-1}\right)} \|g\|_2}.$$
(119)

Step 5 — Geometric interpretation. Shampoo emerges as the memoryless approximation of the global Kronecker-factored equalizer, where each mode-wise matrix G_i captures the tensor structure geometry of neural network weights. The optimal learning rate η^* balances the Kronecker-factored budget constraint with the current gradient magnitude, providing mode-wise G_i -weighted scaling that exploits tensor correlations while maintaining computational tractability through Kronecker structure. The optimizer $Q^* = \eta^* \bigotimes_i G_i^{-1/2}$ implements approximate steepest descent in the Kronecker-factored natural metric. By commutativity, this restricted optimum coincides with projecting the global optimum Q^C onto C_{Kron} .

C.3 DIAGONAL/ADAPTIVE-MOMENT FAMILY

Corollary 17 (Instantaneous optimal AdaGrad). Let moment be a diagonal matrix $R[n] = \operatorname{diag}(r_j[n])$ with coordinate-wise sequence $r_j[n]$. Given gradients g[t], maintain the cumulative second-moment $v_j[t] := \epsilon + \sum_{s \le t} g_j[s]^2 > 0$ with regularization $\epsilon > 0$. Fix the current time t and define the coordinate-wise costs $c_j := v_j[t]^{1/2}$. Consider the diagonal budget

$$Q_D(B,c) := \{ \operatorname{diag}(q_j) : \sum_j c_j \sum_{k>0} |q_j[k]|^2 \le B \},$$
 (120)

and the cone $C_{memoryless}$ of diagonal memoryless optimizers

$$C_{memoryless} := \{Q_{\eta}[n] = \eta \delta[n] \operatorname{diag}(1/c_j) : \eta \ge 0\}. \tag{121}$$

Optimizing for problem P3 under $Q_D(B,c) \cap C_{memoryless}$ with moment R[n] yields an AdaGrad optimizer with optimal hyperparameter:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{\sum_j 1/c_j}},\tag{122}$$

where the optimizer is $Q^* = \eta^* \delta[0] \operatorname{diag}(1/c_j)$.

Proof. Fix the current time t and omit the subscript for brevity. We work in the diagonal Hilbert space $(\mathcal{H}_D, \langle \cdot, \cdot \rangle_{\mathcal{H}_D})$ of diagonal causal LTI filters with weighted norm $\|Q\|_{\mathcal{H}_D}^2 = \sum_j c_j \|q_j\|_{\mathcal{H}}^2$ where $c_j = v_j^{1/2}$ are the coordinate-wise costs.

Step 1 — Global optimal diagonal equalizer. By Cauchy-Schwarz, the unconstrained optimum of problem P3 under the diagonal weighted budget $Q_D(B,c) = \{\operatorname{diag}(q_j) : \sum_j c_j \|q_j\|_{\mathcal{H}}^2 \leq B\}$ is

$$Q^{C} = \sqrt{B} \frac{R}{\|R\|_{\mathcal{H}_{D}}}, \qquad P^{*}(R) = \sqrt{B} \|R\|_{\mathcal{H}_{D}}.$$
 (123)

Step 2 — Commutativity via smooth convexity. The diagonal weighted ball $\mathcal{Q}_D(B,c)$ is smooth and strictly convex. At any boundary point Q with $\|Q\|_{\mathcal{H}_D} = \sqrt{B}$, the normal cone is the ray $N_{\mathcal{Q}_D(B,c)}(Q) = \{\lambda \operatorname{diag}(c_jq_j) : \lambda \geq 0\}$. By Lemma 4, commutativity holds:

$$\Pi_{\mathcal{Q}_D(B,c)\cap\mathcal{C}_{\text{memoryless}}}(Q^C) \in \arg\max_{Q\in\mathcal{Q}_D(B,c)\cap\mathcal{C}_{\text{memoryless}}} \langle Q,R \rangle_{\mathcal{H}_D}.$$
(124)

Therefore, we can equivalently solve the restricted optimization problem directly.

Step 3 — Memoryless family parametrization. For the memoryless family $Q_{\eta}[n] = \eta \delta[n] \operatorname{diag}(1/c_j)$ with $\eta \geq 0$, we compute:

$$||Q_{\eta}||_{\mathcal{H}_D}^2 = \sum_j c_j \left(\eta/c_j\right)^2 = \eta^2 \sum_j \frac{1}{c_j},\tag{125}$$

and for the instantaneous gradient moment $R[0] = \operatorname{diag}(r_i[0])$ and R[n] = 0 for n > 0:

$$\langle Q_{\eta}, R \rangle_{\mathcal{H}_D} = \sum_j \frac{\eta}{c_j} r_j[0] = \eta \sum_j \frac{r_j[0]}{c_j}.$$
 (126)

Step 4 — Budget saturation and optimization. The objective is linear in η while the constraint is quadratic, so the maximizer saturates the budget $||Q_n||_{\mathcal{H}_D} = \sqrt{B}$. This gives:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{\sum_j 1/c_j}}.$$
(127)

Step 5 — Geometric interpretation. AdaGrad emerges as the memoryless approximation of the global diagonal equalizer, where each coordinate-wise cost $c_j = (\epsilon + \sum_{s \leq t} g_j[s]^2)^{1/2}$ captures the cumulative gradient variance. The optimal learning rate η^* balances the diagonal budget constraint with the current gradient, providing coordinate-wise inverse-variance scaling that adapts to the historical gradient magnitudes. The optimizer $Q^* = \eta^* \delta[0] \operatorname{diag}(1/c_j)$ implements approximate steepest descent in the cumulative variance-weighted metric. By commutativity, this restricted optimum coincides with projecting the global optimum Q^C onto $\mathcal{C}_{\text{memoryless}}$.

Corollary 18 (Instantaneous optimal RMSProp). Let moment be a diagonal matrix $R[n] = \operatorname{diag}(r_i[n])$ with coordinate-wise sequence $r_i[n]$. Given gradients g[t], maintain the second-moment

EMA $v_j[t] := \beta_2 v_j[t-1] + (1-\beta_2)g_j[t]^2 > 0$ with parameter $\beta_2 \in (0,1)$. Fix the current time t and define the coordinate-wise costs $c_j := v_j[t]^{1/2}$. Consider the diagonal budget

$$Q_D(B,c) := \{ \operatorname{diag}(q_j) : \sum_j c_j \sum_{k \ge 0} |q_j[k]|^2 \le B \},$$
 (128)

and the cone $C_{memoryless}$ of diagonal memoryless optimizers

$$C_{memoryless} := \{Q_n[n] = \eta \delta[n] \operatorname{diag}(1/c_i) : \eta \ge 0\}. \tag{129}$$

Optimizing for problem P3 under $Q_D(B,c) \cap C_{memoryless}$ with moment R[n] yields an RMSProp optimizer with optimal hyperparameters:

$$\beta_2^* = \arg \max_{0 < \beta_2 < 1} \frac{\sum_j r_j[0]/c_j}{\sqrt{\sum_j 1/c_j}}, \qquad \eta^* = \frac{\sqrt{B}}{\sqrt{\sum_j 1/c_j}}, \tag{130}$$

where the optimizer is $Q^* = \eta^* \delta[0] \operatorname{diag}(1/c_i)$.

Proof. Fix the current time t and omit the subscript for brevity. We work in the diagonal Hilbert space $(\mathcal{H}_D, \langle \cdot, \cdot \rangle_{\mathcal{H}_D})$ of diagonal causal LTI filters with weighted norm $\|Q\|_{\mathcal{H}_D}^2 = \sum_j c_j \|q_j\|_{\mathcal{H}}^2$ where $c_j = v_j^{1/2}$ are the coordinate-wise costs.

Step 1 — Global optimal diagonal equalizer. By Cauchy-Schwarz, the unconstrained optimum of problem P3 under the diagonal weighted budget $Q_D(B,c) = \{\operatorname{diag}(q_j) : \sum_j c_j \|q_j\|_{\mathcal{H}}^2 \leq B\}$ is

$$Q^{C} = \sqrt{B} \frac{R}{\|R\|_{\mathcal{H}_{D}}}, \qquad P^{*}(R) = \sqrt{B} \|R\|_{\mathcal{H}_{D}}.$$
 (131)

Step 2 — Commutativity via smooth convexity. The diagonal weighted ball $Q_D(B,c)$ is smooth and strictly convex. At any boundary point Q with $\|Q\|_{\mathcal{H}_D} = \sqrt{B}$, the normal cone is the ray $N_{\mathcal{Q}_D(B,c)}(Q) = \{\lambda \operatorname{diag}(c_jq_j) : \lambda \geq 0\}$. By Lemma 4, commutativity holds:

$$\Pi_{\mathcal{Q}_D(B,c)\cap\mathcal{C}_{\text{memoryless}}}(Q^C) \in \arg\max_{Q\in\mathcal{Q}_D(B,c)\cap\mathcal{C}_{\text{memoryless}}} \langle Q,R \rangle_{\mathcal{H}_D}.$$
(132)

Therefore, we can equivalently solve the restricted optimization problem directly.

Step 3 — Memoryless family parametrization. For the memoryless family $Q_{\eta}[n] = \eta \delta[n] \operatorname{diag}(1/c_j)$ with $\eta \geq 0$, we compute:

$$||Q_{\eta}||_{\mathcal{H}_D}^2 = \sum_j c_j \left(\eta/c_j\right)^2 = \eta^2 \sum_j \frac{1}{c_j},\tag{133}$$

and for the instantaneous gradient moment $R[0] = \operatorname{diag}(r_i[0])$ and R[n] = 0 for n > 0:

$$\langle Q_{\eta}, R \rangle_{\mathcal{H}_D} = \sum_j \frac{\eta}{c_j} r_j[0] = \eta \sum_j \frac{r_j[0]}{c_j}.$$
 (134)

Step 4 — Budget saturation and optimization. The objective is linear in η while the constraint is quadratic, so the maximizer saturates the budget $||Q_{\eta}||_{\mathcal{H}_D} = \sqrt{B}$. This gives:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{\sum_j 1/c_j}}.$$
(135)

The optimal β_2^{\star} maximizes the budget-normalized gain:

$$J(\beta_2) := \frac{\langle Q_{\eta}, R \rangle_{\mathcal{H}_D}}{\|Q_{\eta}\|_{\mathcal{H}_D}} = \frac{\sum_j r_j[0]/c_j}{\sqrt{\sum_j 1/c_j}}.$$
 (136)

Step 5 — Geometric interpretation. RMSProp emerges as the memoryless approximation of the global diagonal equalizer, where each coordinate-wise cost $c_j = v_j[t]^{1/2}$ captures the exponentially-weighted gradient variance. The exponential weighting β_2 prevents the indefinite accumulation that causes AdaGrad's learning rate decay while maintaining adaptive per-coordinate scaling. The optimal learning rate η^* balances the diagonal budget constraint with the current gradient. By commutativity, this restricted optimum coincides with projecting the global optimum Q^C onto $\mathcal{C}_{\text{memoryless}}$.

Corollary 19 (Instantaneous optimal AdaFactor). Consider matrix parameters $\Theta \in \mathbb{R}^{m \times n}$ with gradient $G \in \mathbb{R}^{m \times n}$. Maintain row-wise second moment estimates $r_i[t] \coloneqq \beta_2 r_i[t-1] + (1-\beta_2)\|G_{i,:}[t]\|^2$ and column-wise estimates $c_j[t] \coloneqq \beta_2 c_j[t-1] + (1-\beta_2)\|G_{:,j}[t]\|^2$. Define the Kronecker-factored diagonal budget

$$Q_K(B, r, c) := \{ Q = \operatorname{diag}(\widehat{r})^{-1/2} \otimes \operatorname{diag}(\widehat{c})^{-1/2} : \|Q\|_{\mathcal{H}}^2 \le B \}, \tag{137}$$

and the cone $C_{memoryless}$ of memoryless optimizers

$$C_{memoryless} := \{Q_{\eta}[n] = \eta \delta[n]Q_0 : \eta \ge 0, Q_0 \text{ fixed}\}. \tag{138}$$

Optimizing problem P3 under $Q_K(B, r, c) \cap C_{memoryless}$ with instantaneous gradient moment $R[0] = GG^{\top}$ yields AdaFactor with optimal hyperparameters:

$$\eta^{\star} = \frac{\sqrt{B}}{\sqrt{\sum_{i} \sum_{j} 1/(\widehat{r}_{i}\widehat{c}_{j})}}, \qquad \beta_{2}^{\star} = \arg\max_{0 < \beta_{2} < 1} \frac{\sum_{i,j} G_{ij}^{2}/(\widehat{r}_{i}\widehat{c}_{j})}{\sqrt{\sum_{i} \sum_{j} 1/(\widehat{r}_{i}\widehat{c}_{j})}}, \tag{139}$$

where the optimizer is $Q^* = \eta^* \delta[0] (\operatorname{diag}(\widehat{r})^{-1/2} \otimes \operatorname{diag}(\widehat{c})^{-1/2})$

Proof. Fix the current time t and omit the subscript for brevity. We work in the Kronecker-factored Hilbert space $(\mathcal{H}_K, \langle \cdot, \cdot \rangle_{\mathcal{H}_K})$ of causal LTI filters with Kronecker-structured impulse response and norm $\|Q\|_{\mathcal{H}_K}^2 = \sum_{n=0}^{\infty} \mathrm{Tr}(Q[n]^{\top}Q[n])$.

Step 1 — Global optimal Kronecker equalizer. By Cauchy-Schwarz, the unconstrained optimum of problem P3 under the Kronecker-factored budget $Q_K(B,r,c)$ is

$$Q^{C} = \sqrt{B} \frac{R}{\|R\|_{\mathcal{H}_{K}}}, \qquad P^{*}(R) = \sqrt{B} \|R\|_{\mathcal{H}_{K}}. \tag{140}$$

Step 2 — Commutativity via smooth convexity. The Kronecker-factored ball $Q_K(B,r,c)$ is smooth and strictly convex. At any boundary point Q with $\|Q\|_{\mathcal{H}_K} = \sqrt{B}$, the normal cone is the ray $N_{\mathcal{Q}_K(B,r,c)}(Q) = \{\lambda Q : \lambda \geq 0\}$. By Lemma 4, commutativity holds:

$$\Pi_{\mathcal{Q}_K(B,r,c)\cap\mathcal{C}_{\text{memoryless}}}(Q^C) \in \arg\max_{Q\in\mathcal{Q}_K(B,r,c)\cap\mathcal{C}_{\text{memoryless}}} \langle Q,R\rangle_{\mathcal{H}_K}.$$
(141)

Therefore, we can equivalently solve the restricted optimization problem directly.

Step 3 — Memoryless Kronecker family parametrization. For the memoryless Kronecker family $Q_{\eta}[n] = \eta \delta[n] (\operatorname{diag}(\widehat{r})^{-1/2} \otimes \operatorname{diag}(\widehat{c})^{-1/2})$ with $\eta \geq 0$, we compute:

$$||Q_{\eta}||_{\mathcal{H}_{K}}^{2} = \eta^{2} \operatorname{Tr}((\operatorname{diag}(\widehat{r})^{-1/2} \otimes \operatorname{diag}(\widehat{c})^{-1/2})^{2}) = \eta^{2} \sum_{i} \sum_{j} \frac{1}{\widehat{r}_{i}\widehat{c}_{j}},$$
(142)

and for the instantaneous gradient moment $R[0] = GG^{\top}$ and R[n] = 0 for n > 0:

$$\langle Q_{\eta}, R \rangle_{\mathcal{H}_K} = \eta \operatorname{Tr}((\operatorname{diag}(\widehat{r})^{-1/2} \otimes \operatorname{diag}(\widehat{c})^{-1/2}) G G^{\top}) = \eta \sum_{i,j} \frac{G_{ij}^2}{\sqrt{\widehat{r}_i \widehat{c}_j}}.$$
 (143)

Step 4 — Budget saturation and optimization. The objective is linear in η while the constraint is quadratic, so the maximizer saturates the budget $||Q_{\eta}||_{\mathcal{H}_K} = \sqrt{B}$. This gives:

$$\eta^* = \frac{\sqrt{B}}{\sqrt{\sum_i \sum_j 1/(\hat{r}_i \hat{c}_j)}}.$$
(144)

The optimal β_2^{\star} maximizes the budget-normalized gain:

$$J(\beta_2) := \frac{\langle Q_{\eta}, R \rangle_{\mathcal{H}_K}}{\|Q_{\eta}\|_{\mathcal{H}_K}} = \frac{\sum_{i,j} G_{ij}^2 / \sqrt{\widehat{r}_i \widehat{c}_j}}{\sqrt{\sum_i \sum_j 1 / (\widehat{r}_i \widehat{c}_j)}}.$$
(145)

Step 5 — Geometric interpretation. AdaFactor emerges as the memoryless approximation of the global Kronecker-factored equalizer, where the row and column second moment estimates \hat{r}_i and \hat{c}_j provide a low-rank factorization of the full diagonal optimizer. This reduces memory complexity from O(mn) to O(m+n) for matrix parameters while approximately preserving Adam's adaptive scaling properties. The exponential weighting β_2 balances the trade-off between adaptation speed and noise reduction in the factored estimates. By commutativity, this restricted optimum coincides with projecting the global optimum Q^C onto $C_{\text{memoryless}}$.

C.4 NORMALIZED-STEP FAMILY

Corollary 20 (Instantaneous optimal signSGD). *Define the* L^{∞} *budget and the cone of memoryless sign optimizers:*

$$Q_{\infty}(\tau) := \{Q[n] = \eta \delta[n] \operatorname{diag}(s_j) : \|Q\|_{\mathcal{H},\infty} \le \tau, |s_j| \le 1\}, \tag{146}$$

$$C_{sign} := \{Q[n] = \eta \delta[n] \operatorname{diag}(\operatorname{sign}(g_j)) : \eta \ge 0\}. \tag{147}$$

(147)

Given instantaneous gradient moment $R[0] = gg^{\top}$ and R[n] = 0 for n > 0, solving problem P3 under the budget $\mathcal{Q}_{\infty}(\tau) \cap \mathcal{C}_{\text{sign}}$ produces signSGD as the optimal solution with optimal hyperparameter:

$$\eta^* = \tau, \quad \text{yielding} \quad \dot{\theta} = \tau \operatorname{sign}(g).$$
 (148)

Proof. We work in the impulse-space Hilbert space $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ of causal LTI filters with matrix impulse response Q and L^{∞} constraint on coordinate-wise step sizes.

Step 1 — Global optimal equalizer. The unconstrained optimum of problem P3 under the L^{∞} budget $Q_{\infty}(\tau)$ is achieved by setting each coordinate j to maximize $g_j\dot{\theta}_j$ subject to $|\dot{\theta}_j| \leq \tau$. This gives:

$$Q^{C}[0] = \tau \operatorname{diag}(\operatorname{sign}(g_j)), \quad Q^{C}[n] = 0 \text{ for } n > 0, \qquad P^{\star}(R) = \tau \|g\|_{1}.$$
 (149)

Step 2 — Commutativity via convexity. The L^{∞} constraint set $\mathcal{Q}_{\infty}(\tau)$ is convex. At any boundary point Q with coordinate-wise saturation, the normal cone contains the sign pattern. By Lemma 4, commutativity holds:

$$\Pi_{\mathcal{Q}_{\infty}(\tau)\cap\mathcal{C}_{\text{sign}}}(Q^C) \in \arg\max_{Q \in \mathcal{Q}_{\infty}(\tau)\cap\mathcal{C}_{\text{sign}}} \langle Q, R \rangle_{\mathcal{H}}.$$
(150)

Therefore, we can equivalently solve the restricted optimization problem directly.

Step 3 — Sign family parametrization. For the sign family $Q_{\eta}[n] = \eta \delta[n] \operatorname{diag}(\operatorname{sign}(g_j))$ with $\eta \geq 0$, we compute:

$$||Q_{\eta}||_{\mathcal{H},\infty} = \eta \max_{j} |\operatorname{sign}(g_{j})| = \eta,$$
(151)

and for the instantaneous gradient moment $R[0] = gg^{\top}$ and R[n] = 0 for n > 0:

$$\langle Q_{\eta}, R \rangle_{\mathcal{H}} = \eta \operatorname{Tr}(\operatorname{diag}(\operatorname{sign}(g_j)) \cdot gg^{\top}) = \eta \sum_{j} g_j \operatorname{sign}(g_j) = \eta \|g\|_1.$$
 (152)

Step 4 — Budget saturation and optimization. The objective is linear in η while the constraint is linear, so the maximizer saturates the budget $\|Q_{\eta}\|_{\mathcal{H},\infty} = \tau$. This gives:

$$\eta^* = \tau. \tag{153}$$

The resulting update is $\dot{\theta} = \tau \operatorname{sign}(g)$.

Step 5 — Geometric interpretation. signSGD emerges as the memoryless approximation of the global L^{∞} -constrained equalizer, where each coordinate takes the maximum allowed step in the direction of its gradient sign. This coordinate-wise saturation provides robustness to gradient magnitude variations and enables efficient low-precision implementations. By commutativity, this restricted optimum coincides with projecting the global optimum Q^{C} onto $\mathcal{C}_{\text{sign}}$.

Corollary 21 (Instantaneous optimal Lion). *Define the* L^{∞} *budget and the cone of momentum-filtered sign optimizers:*

$$Q_{\infty}(\tau) := \{Q : \|Q\|_{\mathcal{H},\infty} \le \tau\},\tag{154}$$

$$C_{Lion} := \{Q_{n,\beta_1}[n] = \eta(1-\beta_1)\beta_1^n \operatorname{diag}(\operatorname{sign}(m_i)) : \eta \ge 0, 0 < \beta_1 < 1\}, \tag{155}$$

where m_j is the j-th coordinate of the momentum-smoothed gradient $m = EMA(g; \beta_1)$. Given gradient moment $R[n] \in \mathcal{H}$, solving problem P3 under the budget $\mathcal{Q}_{\infty}(\tau) \cap \mathcal{C}_{Lion}$ produces Lion optimizer as the optimal solution with optimal hyperparameters:

$$\beta_1^{\star} = \arg \max_{0 < \beta_1 < 1} \sum_{n=0}^{\infty} \beta_1^n \sum_j |\operatorname{Tr}(R_j[n])|, \qquad \eta^{\star} = \tau,$$
 (156)

where $R_j[n]$ denotes the j-th diagonal block of R[n].

Proof. We work in the impulse-space Hilbert space $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ of causal LTI filters with matrix impulse response Q and L^{∞} constraint on coordinate-wise step sizes.

Step 1 — Global optimal equalizer. The unconstrained optimum of problem P3 under the L^{∞} budget $\mathcal{Q}_{\infty}(\tau)$ is achieved by setting each coordinate j to maximize the inner product subject to $|Q_j[n]| \leq \tau$ for all n. This gives:

$$Q^{C}[n] = \tau \operatorname{diag}(\operatorname{sign}(R_{jj}[n])), \qquad P^{\star}(R) = \tau \sum_{n=0}^{\infty} \sum_{j} |R_{jj}[n]|.$$
 (157)

Step 2 — Commutativity via convexity. The L^{∞} constraint set $\mathcal{Q}_{\infty}(\tau)$ is convex. At any boundary point Q with coordinate-wise saturation, the normal cone contains the sign pattern. By Lemma 4, commutativity holds:

$$\Pi_{\mathcal{Q}_{\infty}(\tau)\cap\mathcal{C}_{\text{Lion}}}(Q^C) \in \arg\max_{Q\in\mathcal{Q}_{\infty}(\tau)\cap\mathcal{C}_{\text{Lion}}} \langle Q, R \rangle_{\mathcal{H}}.$$
 (158)

Therefore, we can equivalently solve the restricted optimization problem directly.

Step 3 — Lion family parametrization. For $Q_{\eta,\beta_1}[n] = \eta(1-\beta_1)\beta_1^n \operatorname{diag}(\operatorname{sign}(m_j))$ with $\eta \geq 0$ and $0 < \beta_1 < 1$, we compute:

$$||Q_{\eta,\beta_1}||_{\mathcal{H},\infty} = \eta(1-\beta_1) \max_j |\operatorname{sign}(m_j)| \sum_{n=0}^{\infty} \beta_1^n = \eta,$$
 (159)

and

$$\langle Q_{\eta,\beta_1}, R \rangle_{\mathcal{H}} = \sum_{n=0}^{\infty} \eta(1-\beta_1)\beta_1^n \sum_{j} \text{sign}(m_j)R_{jj}[n] = \eta(1-\beta_1) \sum_{n=0}^{\infty} \beta_1^n \sum_{j} \text{sign}(m_j)R_{jj}[n].$$
(160)

Step 4 — Budget saturation and 1-D optimization. For fixed β_1 , the objective is linear in η while the constraint is linear, so the maximizer saturates the budget $\|Q_{\eta,\beta_1}\|_{\mathcal{H},\infty} = \tau$. This gives:

$$\eta = \tau. \tag{161}$$

The budget-normalized gain becomes:

$$J(\beta_1) := \frac{\langle Q_{\eta,\beta_1}, R \rangle_{\mathcal{H}}}{\|Q_{\eta,\beta_1}\|_{\mathcal{H},\infty}} = (1 - \beta_1) \sum_{n=0}^{\infty} \beta_1^n \sum_{j} \text{sign}(m_j) R_{jj}[n].$$
 (162)

Hence $\beta_1^{\star} = \arg \max_{0 < \beta_1 < 1} J(\beta_1)$ and $\eta^{\star} = \tau$.

Step 5 — Geometric interpretation. The optimal momentum β_1^\star maximizes the alignment between the 1-pole EMA kernel $(\beta_1^n)_{n\geq 0}$ and the sign-weighted empirical lag curve. Lion emerges as the momentum-filtered approximation of the global L^∞ -constrained equalizer, where momentum smoothing $m=\mathrm{EMA}(g;\beta_1)$ reduces gradient noise before applying coordinate-wise sign normalization. By commutativity, this restricted optimum coincides with projecting the global optimum Q^C onto $\mathcal{C}_{\mathrm{Lion}}$.

C.5 LAYER-WISE NORM SCALING (LARGE-BATCH STABILIZERS)

Corollary 22 (Instantaneous optimal LARS/LAMB). *Consider a neural network with L layers, where parameters are partitioned as* $\theta = (\theta_1, \dots, \theta_L)$ *with* $\theta_\ell \in \mathbb{R}^{d_\ell}$. *Define the* layer-wise adaptive budget

$$Q_{layer}(B, \{\|\theta_{\ell}\|_{\ell=1}^{L}) := \left\{ Q = \text{blkdiag}(Q_{1}, \dots, Q_{L}) : \sum_{\ell=1}^{L} \|\theta_{\ell}\|_{2}^{2} \sum_{n=0}^{\infty} \text{Tr}(Q_{\ell}[n]^{\top} Q_{\ell}[n]) \leq B \right\},$$
(163)

and the cone C_{layer} of layer-wise scalar optimizers

$$C_{layer} := \{Q_{\eta,\beta_1}[n] = \text{blkdiag}(\eta(1-\beta_1)\beta_1^n \alpha_1 I_{d_1}, \dots, \eta(1-\beta_1)\beta_1^n \alpha_L I_{d_L}) : \eta \ge 0, 0 < \beta_1 < 1, \alpha_\ell > 0\}.$$
(164)

Given moment matrix $R[n] = \text{blkdiag}(R_1[n], \dots, R_L[n])$ with $S_\ell[n] \coloneqq \text{Tr}(R_\ell[n])$, optimizing problem P3 under $\mathcal{Q}_{layer}(B, \{\|\theta_\ell\|\}) \cap \mathcal{C}_{layer}$ yields LARS/LAMB optimizer with optimal hyperparameters:

$$\beta_1^{\star} = \arg\max_{0 < \beta_1 < 1} \sqrt{1 - \beta_1^2} \sum_{n=0}^{\infty} \beta_1^n \sum_{\ell=1}^{L} \frac{S_{\ell}[n]}{\|\theta_{\ell}\|_2}, \quad \eta^{\star} = \frac{\sqrt{B(1 - \beta_1^{\star 2})}}{\sqrt{\sum_{\ell=1}^{L} d_{\ell}/\|\theta_{\ell}\|_2^2} (1 - \beta_1^{\star})}, \quad (165)$$

where the layer-wise scaling factors are $\alpha_{\ell}^{\star} = 1/\|\theta_{\ell}\|_{2}$.

Proof. We work in the block-diagonal Hilbert space $(\mathcal{H}_{layer}, \langle \cdot, \cdot \rangle_{\mathcal{H}_{layer}})$ of layer-wise causal LTI filters with weighted norm $\|Q\|_{\mathcal{H}_{layer}}^2 = \sum_{\ell=1}^L \|\theta_\ell\|_2^2 \|Q_\ell\|_{\mathcal{H}}^2$.

Step 1 — Global optimal layer-wise equalizer. By Cauchy-Schwarz, the unconstrained optimum of problem P3 under budget $\mathcal{Q}_{layer}(B, \{\|\theta_\ell\|\})$ is

$$Q^{C} = \sqrt{B} \frac{R}{\|R\|_{\mathcal{H}_{layer}}}, \qquad P^{\star}(R) = \sqrt{B} \|R\|_{\mathcal{H}_{layer}}, \tag{166}$$

where $||R||_{\mathcal{H}_{layer}}^2 = \sum_{\ell=1}^L ||\theta_\ell||_2^2 ||R_\ell||_{\mathcal{H}}^2$.

Step 2 — Commutativity via smooth convexity. The layer-wise adaptive ball $\mathcal{Q}_{\text{layer}}(B, \{\|\theta_\ell\|\})$ is smooth and strictly convex. At any boundary point Q with $\|Q\|_{\mathcal{H}_{\text{layer}}} = \sqrt{B}$, the normal cone is the ray $N_{\mathcal{Q}_{\text{layer}}}(Q) = \{\lambda Q : \lambda \geq 0\}$. By Lemma 4, commutativity holds:

$$\Pi_{\mathcal{Q}_{\text{layer}}(B,\{\|\theta_{\ell}\|\})\cap\mathcal{C}_{\text{layer}}}(Q^C) \in \arg\max_{Q\in\mathcal{Q}_{\text{layer}}(B,\{\|\theta_{\ell}\|\})\cap\mathcal{C}_{\text{layer}}} \langle Q,R\rangle_{\mathcal{H}_{\text{layer}}}.$$
(167)

Therefore, we can equivalently solve the restricted optimization problem directly.

Step 3 — Layer-wise family parametrization. For $Q_{\eta,\beta_1}[n] = \text{blkdiag}(\eta(1-\beta_1)\beta_1^n\alpha_1I_{d_1},\ldots,\eta(1-\beta_1)\beta_1^n\alpha_LI_{d_L})$ with $\eta \geq 0, 0 < \beta_1 < 1$, and $\alpha_\ell = 1/\|\theta_\ell\|_2$, we compute:

$$\|Q_{\eta,\beta_1}\|_{\mathcal{H}_{\text{layer}}}^2 = \sum_{\ell=1}^L \|\theta_\ell\|_2^2 \sum_{n=0}^\infty \text{Tr}((\eta(1-\beta_1)\beta_1^n \alpha_\ell I_{d_\ell})^\top (\eta(1-\beta_1)\beta_1^n \alpha_\ell I_{d_\ell}))$$
(168)

$$= \eta^2 (1 - \beta_1)^2 \frac{1}{1 - \beta_1^2} \sum_{\ell=1}^{L} \|\theta_\ell\|_2^2 \alpha_\ell^2 d_\ell = \eta^2 \frac{(1 - \beta_1)^2}{1 - \beta_1^2} \sum_{\ell=1}^{L} \frac{d_\ell}{\|\theta_\ell\|_2^2}, \tag{169}$$

and

$$\langle Q_{\eta,\beta_1}, R \rangle_{\mathcal{H}_{\text{layer}}} = \sum_{\ell=1}^{L} \sum_{n=0}^{\infty} \text{Tr}((\eta(1-\beta_1)\beta_1^n \alpha_{\ell} I_{d_{\ell}})^{\top} R_{\ell}[n])$$
(170)

$$= \eta(1 - \beta_1) \sum_{n=0}^{\infty} \beta_1^n \sum_{\ell=1}^{L} \alpha_{\ell} S_{\ell}[n] = \eta(1 - \beta_1) \sum_{n=0}^{\infty} \beta_1^n \sum_{\ell=1}^{L} \frac{S_{\ell}[n]}{\|\theta_{\ell}\|_2}.$$
 (171)

Step 4 — Budget saturation and 1-D optimization. For fixed β_1 , the objective is linear in η while the constraint is quadratic, so the maximizer saturates the budget $\|Q_{\eta,\beta_1}\|_{\mathcal{H}_{layer}} = \sqrt{B}$. This gives:

$$\eta = \frac{\sqrt{B(1-\beta_1^2)}}{\sqrt{\sum_{\ell=1}^L d_\ell/\|\theta_\ell\|_2^2} (1-\beta_1)}.$$
(172)

The budget-normalized gain becomes:

$$J(\beta_1) := \frac{\langle Q_{\eta,\beta_1}, R \rangle_{\mathcal{H}_{\text{layer}}}}{\|Q_{\eta,\beta_1}\|_{\mathcal{H}_{\text{layer}}}} = \sqrt{\frac{1 - \beta_1^2}{\sum_{\ell=1}^L d_\ell / \|\theta_\ell\|_2^2}} \sum_{n=0}^{\infty} \beta_1^n \sum_{\ell=1}^L \frac{S_\ell[n]}{\|\theta_\ell\|_2}.$$
 (173)

Hence $\beta_1^{\star} = \arg \max_{0 < \beta_1 < 1} J(\beta_1)$ and η^{\star} saturates the budget constraint.

Step 5 — Geometric interpretation. The optimal momentum β_1^\star maximizes the cosine similarity between the 1-pole EMA kernel $(\beta_1^n)_{n\geq 0}$ and the layer-normalized empirical lag curve $(\sum_{\ell=1}^L S_\ell[n]/\|\theta_\ell\|_2)_{n\geq 0}$. The layer-wise scaling $\alpha_\ell=1/\|\theta_\ell\|_2$ prevents layer collapse by ensuring updates remain proportional to current parameter magnitudes. By commutativity, this restricted optimum coincides with projecting the global optimum Q^C onto C_{layer} .

D MASTER TABLE OF OPTIMIZERS

In the previous section, we have derived various types of optimizers from our convex optimization framework. We can now register various optimizers under a single unified table. Each optimizer corresponds to a specific choice of moment matrix M, budget constraint \mathcal{Q} , and resulting equalizer Q. "Param restrict" rows are feasible points in the convex programs that can either be kept and fitted to target moments, or replaced with full closed-form solutions.

Glossary.

- Instant moment $M \in \mathbb{S}_+^P$: $M = \Sigma_{\mathrm{tr}} = \mathbb{E}[gg^\top]$ (training) or $M = \mathrm{sym}(C) = \mathrm{sym}\,\mathbb{E}[g_{\mathrm{tr}}g_{\mathrm{val}}^\top]$ (validation-aware)
- Dynamic moment operator M_{σ} (Laplace/z window, $\sigma > 0$); frequency form $M(\omega)$
- Instant budgets from Section 2: Frob, Trace/Spectral, Lyap, Diag
- Dynamic budgets from Section 3: **D-Frob** (H₂), **D-Trace/Spectral**, **D-Lyap**, **D-Diag**
- "Param restrict." means we restrict Q to a small parametric family inside the convex budget

Table 2: Optimizer Specifications: Moment and Budget Constraints

Optimizer	Moment Used	Budget $\mathcal Q$
GD	$M = \Sigma_{\rm tr}$	Frob: $ Q _F \le \kappa$ (instant)
SGD	same as GD (stochastic)	Frob (instant)
Momentum	M_{σ} (weighted)	D-Frob : $ Q _{\mathcal{H}} \leq \sqrt{B}$
(HB/NAG)		
Nesterov	M_{σ}	D-Frob (with predictive tap)
AdaGrad	$M = \Sigma_{ m tr}$	Diag: $\sum_{i} c_{i}q_{j}^{2} \leq B$
RMSProp	$M = \Sigma_{ m tr}$	Diag
Adam/AdamW	M_{σ} (via m, v)	D-Diag: $\sum_{i} q_{i} ^{2}_{H_{2},c_{i}} \leq B$
Adam (val-aware)	$M = \operatorname{sym}(C)$ or M_{σ}	Diag or D-Diag
LAMB / LARS	$M = \Sigma_{\mathrm{tr}}$ per layer ℓ	Trace per layer: $Tr(Q_{\ell}) \leq \tau_{\ell}$ + norm ratio
K-FAC	$M = \Sigma_{ m tr}$	Lyap (block-Kronecker): $Tr(Q^2\Sigma) \leq B$ with Q factored
Shampoo	$M = \Sigma_{ m tr}$	Lyap (multi-axis Kronecker)
Newton / GN	$M = \Sigma_{ m tr}$	Lyap with $Q = H^{-1}$ (or G^{-1})
signSGD	$M = \Sigma_{\mathrm{tr}}$	ℓ_{∞} step budget on $\dot{ heta}$
Lion	M_{σ} via m	ℓ_{∞} on $\dot{\theta}$ (dynamic)
Polyak step	$M = \Sigma_{\mathrm{tr}}$ + scalar loss	Frob + 1D line search

Table 3: Optimizer Equalizers and Closed Forms

Optimizer	Equalizer Q	Closed Form
GD	$Q = \alpha I$ (isotropic)	$Q^{\star} = \kappa \frac{M}{\ M\ _{E}} \rightarrow \text{in practice } \alpha = \eta$
SGD	$Q = \alpha I$	same as $\overline{\text{GD}}$, with M estimated from minibatch
Momentum	$Q(z) = \eta I \frac{1-\beta}{1-\beta z^{-1}}$ (one-pole)	Param restrict of dynamic proportional
(HB/NAG)	- () - 1 B () 1	optimum; (η, β) fit the target
Nesterov	$Q(z) = \eta I \frac{1-\beta}{1-\beta z^{-1}} (1 + \gamma z^{-1})$	Param restrict (lead–lag) under same budget
AdaGrad	$Q = \operatorname{diag}(q)$	$q_j^\star \propto rac{M_{jj}}{c_i}; c_j \uparrow$ with cum. second moment
RMSProp	$Q = \operatorname{diag}(q)$	same as AdaGrad but c_j from EMA of g_j^2
Adam/AdamW	$Q(z) = \operatorname{diag}\left(\eta \frac{1-\beta_1}{1-\beta_1 z^{-1}}\right) (\operatorname{diag}\sqrt{v} + \epsilon)^{-1}$	Param restrict of D-Diag optimum
Adam (val-aware)	same as Adam	Use val cross-moment for q_i (and m vs g)
LAMB / LARS	$Q_\ell = \alpha_\ell I$ s.t. $\ v_\ell\ \propto \ \theta_\ell\ $	Water-fill over layers + isotropic in each layer
K-FAC	$Q = \bigoplus_{\ell} (A_{\ell}^{-1} \otimes G_{\ell}^{-1})$ approx	Equal-power in each layer's Fisher metric (factored)
Shampoo	$Q = \bigotimes_{i} H_{i}^{-1/2}$ $Q = H^{-1}$	Equal-power along tensor modes
Newton / GN	$Q = H^{-1}$	Exact Lyap optimum if constraint matches curvature
signSGD	Q such that $\dot{\theta} = \eta \operatorname{sign}(g)$	Linear objective + $\ \dot{\theta}\ _{\infty} \le \eta \Rightarrow \text{vertex}$
Lion	$\dot{\theta} = \eta \operatorname{sign}(m)$	same with smoothed signal
Polyak step	$Q=\alpha I$ with α from loss & grad norm	closed-form $\alpha = \frac{\ g\ ^2}{g^{\top}Hg}$ (local)

Table 4: Optimizer Hyperparameters and Interpretations

Optimizer	Hyperparameters	One-line Implication
GD	η : total power scale	GD is the "proportional router" collapsed to a scalar knob
SGD	η : same	Stochasticity only changes how you estimate M , not the program
Momentum	η : overall gain; β : pole = decay	Momentum is the 1-pole low-rank approximation of the dynamic
(HB/NAG)	time of impulse	proportional router
Nesterov	β : smoothing; γ : look-ahead lead	"Prediction" = a tiny lead in the equalizer—still budgeted power
AdaGrad	ϵ , window: determine c_i growth	AdaGrad is per-coord water-filling with costs = cumulative variance
RMSProp	ρ : EMA window; ϵ	RMSProp = time-local AdaGrad (cheaper, reactive)
Adam/AdamW	η : scale; β_1 : low-pass pole; β_2 : sets	Adam = diagonal dynamic optimizer: per-coord water-filling times a 1-
	costs $c_i \sim v_i$	pole low-pass
Adam (val-aware)	same	Simply switching M to $sym(C)$ turns Adam into a val-aware power
		allocator
LAMB / LARS	trust ratio $\ heta_\ell\ /\ m_\ell\ $	LAMB/LARS = layer-wise trace budget + norm normalization (compute-
** *** **		stable water-filling)
K-FAC	damping λ , update period	K-FAC = structured Lyap equalization \Rightarrow curvature-aware, block-wise
Shampoo	per-axis damping, period	Shampoo = multi-axis equal-power (richer than K-FAC, costlier)
Newton / GN	trust-region radius	Newton = max power under curvature budget; best when Hessian is right
signSGD	η : box size	signSGD = L-∞ trust-region maximizer; robust, but discards magnitude
		info
Lion	β_1	Lion = L-∞ budget on smoothed signal; ultra-aggressive quantized
		equalizer
Polyak step	none beyond window	Polyak = Frob program with on-the-fly α estimate (local curvature)

E ALGORITHM

For completness, we provide a realizable algorithm for calculating the optimal SGD+Momentum and Adam, as proved in Section 2 and Section 3 of the main manuscript.

1993

1997

```
1944
             Algorithm 1 Optimal SGD+Momentum (\beta^*, \eta^*) from gradient history
1945
             Require: Window length T, max lag K \le T - 1, EMA decay \rho \in (0, 1), budget B > 0
1946
               1: Initialize \mu \leftarrow 0, S[n] \leftarrow 0 for n = 0, \dots, K, buffer \leftarrow \emptyset, d \leftarrow parameter dimension
1947
               2: for each calibration step t do
1948
               3:
                         g_t \leftarrow \text{flatten current } \nabla_{\theta} L
1949
               4:
                         \mu \leftarrow \rho \mu + (1 - \rho)g_t
1950
               5:
                         \tilde{g}_t \leftarrow g_t - \mu
1951
                         push \tilde{g}_t into buffer (keep last T)
               6:
1952
               7:
                         for n = 0, \dots, K with t - n in buffer do
                               \begin{array}{l} s_n \leftarrow \langle \tilde{g}_t, \tilde{g}_{t-n} \rangle \\ S[n] \leftarrow \rho S[n] + (1-\rho) s_n \end{array}
1953
               8:

⊳ scalar dot product

               9:
1954
                         end for
             10:
1955
                         J(\beta) = \sqrt{\frac{1-\beta^2}{d}} \cdot \sum_{n=0}^{K} \beta^n S[n]
             11:
1957
                         \beta^{\star} \leftarrow \arg \max_{\beta \in (0,1)} J(\beta)
             12:
                                                                                                          > 1-D search (e.g., bounded line search)
                         \eta^{\star} \leftarrow \frac{\sqrt{B(1 - (\beta^{\star})^2)}}{\sqrt{d}(1 - \beta^{\star})}
1958
             13:
1959
             14: end for
1960
             15: return (\beta^{\star}, \eta^{\star})
1961
1962
1963
1964
             Algorithm 2 Optimal Adam (\beta_1^{\star}, \beta_2^{\star}, \eta^{\star}) from gradient history
1965
             Require: Window length T, max lag K \leq T - 1, EMA decay \rho \in (0,1), budget B > 0
1966
               1: Initialize \mu \leftarrow 0, v_j \leftarrow 0 for all j, T_t[n] \leftarrow 0 for n = 0, \dots, K, buffer \leftarrow \emptyset, d \leftarrow parameter
1967
                    dimension
1968
               2: for each calibration step t do
1969
                         g_t \leftarrow \text{flatten current } \nabla_{\theta} L
1970
               4:
                         \mu \leftarrow \rho \mu + (1 - \rho)g_t
1971
               5:
                         \tilde{g}_t \leftarrow g_t - \mu
1972
                         push \tilde{g}_t into buffer (keep last T)
               6:
               7:
                         for each coordinate j do
                               v_j \leftarrow \rho v_j + (1 - \rho)\tilde{g}_{t,j}^2
c_j \leftarrow \sqrt{v_j}
               8:
1974
               9:
1975
              10:
                         end for
1976
             11:
                         for n = 0, \dots, K with t - n in buffer do
                               T_t[n] \leftarrow \sum_j \frac{\tilde{g}_{t,j}\tilde{g}_{t-n,j}}{c_j}
             12:
                                                                                                          1978
                         end for
             13:
1979
                         W_t \leftarrow \sum_j \frac{1}{c_i}
             14:
1980
                         a(\beta_1, \beta_2) \leftarrow \sqrt{\frac{1-\beta_1^2}{W_t}}
1981
             15:
1982
                         J(\beta_1, \beta_2) = a(\beta_1, \beta_2) \sum_{n=0}^{K} \beta_1^n T_t[n]
             16:
                         (\beta_1^{\star}, \beta_2^{\star}) \leftarrow \arg \max_{\beta_1, \beta_2 \in (0,1)} J(\beta_1, \beta_2)
             17:
                                                                                                                                                     ▷ 2-D search
1984
                         \eta^\star \leftarrow \tfrac{\sqrt{B}\,a(\beta_1^\star,\beta_2^\star)}{1-\beta_1^\star}
             18:
1985
1986
             19: end for
             20: return (\beta_1^{\star}, \beta_2^{\star}, \eta^{\star})
1987
```