

---

# Meta-optimization for Deep Learning via Nonstochastic Control

---

Anonymous Authors<sup>1</sup>

## Abstract

Hyperparameter tuning in mathematical optimization is a notoriously difficult problem. Recent tools from online control give rise to a provable methodology for hyperparameter tuning in convex optimization called meta-optimization. In this work, we extend this methodology to nonconvex optimization and the training of deep neural networks. We present an algorithm for nonconvex meta-optimization that leverages the reduction from nonconvex optimization to convex optimization, and investigate its applicability for deep learning tasks on academic-scale datasets.

## 1. Introduction

Hyperparameter tuning for deep learning is notoriously difficult and resource-consuming. It is therefore an extremely well-studied problem, with numerous approaches including Bayesian optimization (Snoek et al., 2012), bandit algorithms (Li et al., 2018), meta-gradient methods (Baydin et al., 2017), and spectral methods (Hazan et al., 2018).

A recent paradigm based on online control (Chen & Hazan, 2023) gives the first provable guarantees for hyperparameter tuning for smooth convex optimization. There are several novel aspects to this approach: (1) Instead of learning the hyperparameters in one shot, the goal is to repeatedly solve the given optimization problem, learning from experience and converging to the performance of the best hyperparameters. (2) Parameters are automatically tuned by a feedback control algorithm based on novel techniques from online control.

Convexity is important for meta optimization since it is based on regret minimization, which is in general intractable for nonconvex problems. Regret minimization for the convex objectives of meta-optimization allows convergence to the best method in hindsight, and it is not immediately clear

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

what the nonconvex analogue is.

In this work, we extend the online control based framework of meta-optimization to nonconvex optimization. Given the framework’s guarantees for smooth convex functions, we leverage a reduction from nonconvex to convex optimization inspired by (Agarwal et al., 2019). We propose an algorithm that can learn to adapt to the problem over many episodes and eventually reach an approximate stationary point. As the number of episodes increases, it converges at a rate that is determined by the performance of the best optimization algorithm from a class of methods.

We conduct experiments on academic-scale workloads, including image classification and machine translation. These initial experiments demonstrate the applicability of the algorithm. In addition, we ablate over several design choices and empirically verify our assumptions.

### 1.1. Related work

**Parameter-free optimization** Parameter-free optimization methods are adaptations of first-order methods that remove the need to tune certain hyperparameters, such as the learning rate. Methods in this space include coin-betting (Orabona & Pál, 2016), adaptation to the diameter of the decision set (Defazio & Mishchenko, 2023), and many more (Ivgy et al., 2023), (Cutkosky et al., 2024), (Lu et al., 2022). The meta-optimization approach is more general in two aspects: (1) it attempts to learn the best algorithm for specific objective functions, rather than a class of functions (for example smooth functions with a specific smoothness parameter) (2) for quadratic problems, it has guarantees over a larger class of methods, including preconditioned methods with a fixed preconditioner. However, it is not parameter-free since there are tunable parameters in the method.

**Control for optimization** Control and optimization are closely related fields, starting from Lyapunov’s work and its application to the design and analysis of optimization algorithms, see (Chen & Hazan, 2023) for more background. (Lessard et al., 2016) apply control theory to the analysis of optimization algorithms on a single problem instance, giving a general framework for obtaining convergence guarantees for a variety of gradient-based methods. (Casgrain & Kratsios, 2021) study the characterization of the regret-

optimal algorithm given an objective function, using a value function-based approach motivated by optimal control. Our work builds upon (Chen & Hazan, 2023), which studies an online control based methodology for regret minimization and apply it to convex optimization.

## 2. Preliminaries

**Meta-optimization** In meta-optimization (Chen & Hazan, 2023), we are given a **sequence of optimization problems**, called episodes. The goal is to design an algorithm that, over many episodes, can perform as well as the best algorithm in a benchmark algorithm class. We denote the number of episodes as  $N$ , and in each episode, we perform  $T$  steps of optimization. Since the given problems are solved individually, we assume that at the beginning of an episode, the iterate is reinitialized to an arbitrary starting point  $x_{i,1}$ . (Chen & Hazan, 2023) proposes a method based on online control, and guarantees that over  $N$  episodes, the method converges to the performance of the best first-order gradient method from a general class of methods.

**Extension to nonconvex stochastic optimization** We extend the meta-optimization framework to nonconvex optimization in the finite-sum setting. Denote the sequence of  $N$  objective functions as  $\{f_i\}_{i=1}^N$ , in this setting, each function is a finite sum of  $n$  nonconvex functions:

$$f_i(x) = \frac{1}{n} \sum_{j=1}^n f_{i,j}(x).$$

In each episode, we are given an objective function  $f_i$ , and at each time step, we have access to a mini-batch of  $f_{i,j}$ 's. Notably, this setting formalizes the problem of neural network training, where the objective function is the average loss on training examples.

The goal of stochastic nonconvex optimization is to obtain an  $\varepsilon$ -stationary point in expectation for each objective function: an  $x_i$  such that  $\mathbb{E} [\|\nabla f_i(x_i)\|] \leq \varepsilon$  for every  $i \in [N]$ . The expectation is taken over the randomness of the mini-batches and possibly the optimization algorithm. As is standard in the literature, we assume each  $f_{i,j}$  is smooth and has bounded function value.

**Definition 1.** We say a function is  $\beta$ -smooth if for every  $x, y$ ,  $\|\nabla f(x) - \nabla f(y)\| \leq \beta\|x - y\|$ .

**Assumption 1.** For all  $i, j$ ,  $0 \leq f_{i,j}(x) \leq M$  for all  $x$ , and  $f_{i,j}$  is  $\beta$ -smooth.

In episode  $i$ , at time  $t$  (denoted as step  $(i, t)$ ), an optimization algorithm  $\mathcal{A}$  chooses a point  $x_{i,t} \in \mathbb{R}^d$ . Then it receives a mini-batch of examples  $B_{i,t}$  of size  $b$ , and suffers the nonconvex cost  $f_{i,t}(x_{i,t}) = \frac{1}{b} \sum_{j \in B_{i,t}} f_{i,j}(x_{i,t})$ . The protocol of this setting is formally defined in Algorithm 2 in the

appendix. Our goal is to design an algorithm  $\mathcal{A}$  whose convergence rate for finding an approximate stationary point approaches that of the best algorithm in a benchmark class.

**Reduction from nonconvex to convex optimization** It is known that deterministic non-convex optimization can be reduced to solving a sequence of strongly convex problems (Agarwal et al., 2019; Paquette et al., 2018; Wang & Srebro, 2019). Each of the sub-problems in this sequence regularizes the original nonconvex function with the  $\ell_2$  regularizer, so that the sub-problem is strongly convex. The reduction is stated in Algorithm 3 in the appendix.

## 3. Algorithm and guarantees

We leverage the reduction from nonconvex to convex optimization to design our algorithm. The reduction shows that if one can minimize a sequence of strongly convex problems, then an approximate stationary point can be found among these minimizers. Using this framework, we can apply convex meta-optimization to the sequence of strongly convex functions and obtain a method whose convergence is characterized by the performance of the best method on that sequence of strongly convex functions. This guarantee is different from the meta-optimization guarantee one can achieve in convex optimization through regret minimization. In general, regret minimization for nonconvex functions is computationally intractable, and alternative notions of regret were introduced in (Hazan et al., 2017).

---

### Algorithm 1 Nonconvex stochastic meta-optimization

---

**Require:** episode number  $N$ , epoch number  $K$ , inner step number  $S$  such that  $KS = T$ , smoothness parameter  $\beta$ , initial points  $\{x_{i,1}\}_i$ , convex meta-optimization algorithm  $\mathcal{A}$ .

- 1: **for**  $i = 1, \dots, N$  **do**
  - 2:     Re-initialize iterate to  $x_{i,1}$ .
  - 3:     **for**  $k = 1, \dots, K$  **do**
  - 4:         Set  $x_{i,k,1} = x_{i,k}$  and denote  $f_{i,k}(x) = f_i(x) + \beta\|x - x_{i,k}\|^2$  as the regularized objective.
  - 5:         **for**  $s = 1, \dots, S$  **do**
  - 6:             Play  $x_{i,k,s}$  and receive a batch of examples  $B_{i,k,s}$  of size  $b$ .
  - 7:             Define  $f_{i,k,s}(x) = \frac{1}{b} \sum_{j \in B_{i,k,s}} f_{i,j}(x) + \beta\|x - x_{i,k}\|^2$ .
  - 8:             Update  $x_{i,k,s+1} = \mathcal{A}(f_{1,1,1}, \dots, f_{i,k,s})$ .
  - 9:             Update  $x_{i,k+1} = \frac{1}{S} \sum_{s=1}^S x_{i,k,s}$ .
- 

We present the guarantee of our main algorithm, Algorithm 1, in the theorem below. The class of benchmark algorithms  $\pi$  contains methods whose updates are linear functions of past gradients. If the losses are quadratic, it can simulate first-order gradient-based methods such as gradient descent,

momentum, and preconditioned methods. Due to limited space, we defer additional assumptions and technical details to Appendix B.

**Theorem 2.** Let  $x_{i,k}^*$  be a minimizer of  $f_{i,k}$ , and let  $g_k^2 = \frac{1}{K} \sum_{k=1}^K \|f_i(x_{i,k})\|^2$  be the average squared gradient norm. Under Assumptions 1, 2, and 4, using the bandit meta-optimization algorithm (Algorithm 4 in appendix) as  $\mathcal{A}$  in Algorithm 1 yields the following guarantee:

$$\mathbb{E} \left[ \frac{1}{N} \sum_{j=1}^N g_k^2 \right] \leq O\left(\frac{1}{K}\right) + \tilde{O}((NKS)^{-\frac{1}{4}}) + \frac{6\beta}{NKS} \min_{\pi \in \Pi} \mathbb{E} \left[ \sum_{i=1}^N \sum_{k=1}^K \sum_{s=1}^S (f_{i,k,s}(x_{i,k,s}^\pi) - f_{i,k}(x_{i,k}^*)) \right].$$

Note that the last term in the guarantee is at least of order  $O(\frac{1}{S})$ , since the iterates are reset for every  $k$  and  $i$ . As  $N$  grows large, the convergence rate is thus dominated by  $O(\frac{1}{K}) + O(\frac{1}{S})$ , and since  $T = KS$ , we set  $K, S = O(\sqrt{T})$ . Under this choice of  $K$  and  $S$ , the convergence rate is consistent with the rate of SGD. However, it is possible to refine the bound for certain cases, as we show in the appendix.

## 4. Experiments

We run experiments on three deep learning workloads of increasing scale: MNIST handwritten digit recognition, CIFAR-10 image classification, and WMT-17 English-to-German machine translation. On each dataset, we apply a fundamental deep learning architecture (MLP, CNN, and transformer, respectively) and compare our method against the common deep learning optimization techniques. We train in two optimization settings: gradient descent on a fixed batch (i.e.  $f_{i,j} \equiv f_i$  for a fixed  $f_i$ ) and stochastic gradient descent (i.e. a minibatch of  $f_{i,j}$ 's is sampled i.i.d. each step). For the methods labeled "ours", we make practical modifications to Algorithm 1 with the considerations mentioned in Appendix E.1 to arrive at Algorithm 6; this is done in Jax and the code may be found at <https://anonymous.4open.science/t/meta-opt-8916/>. We compare against the following fully-tuned baselines: vanilla gradient descent, momentum, Adam with weight decay (AdamW), hypergradient descent (Baydin et al., 2017), Distance-over-Gradients (Ivgi et al., 2023), D-Adaptation (Defazio & Mishchenko, 2023), and the Mechanic algorithm (Cutkosky et al., 2024). Each of the latter 4 optimizers is built on top of either vanilla gradient descent or vanilla Adam; to reduce clutter, we only plot these final 4 optimizers when they match or outperform their tuned vanilla counterparts. We refer to AdamW, D-Adaptation, and Mechanic collectively as the "adaptive methods", which are not captured theoretically by our benchmark algorithm class.

See Appendix E.2 for more information about baselines and experimental hyperparameters such as batch size and number of iterations.

## 5. Results

We present our main experimental results below, for additional ablations and empirical verification of our assumptions, see Appendix E.3.

**Deterministic optimization** The training loss in the deterministic setting (i.e. where each training step is over the same subset of data) allows us to inspect the optimization performance of our algorithm without the effect of noise. In Figure 1, the training losses of the various optimization algorithms are plotted across episodes. We see that our method improves over time – during the first episode it is worse than gradient descent, but after a handful of deterministic episodes it matches or outperforms many other baselines (note that on WMT there is a separation between adaptive methods and non-adaptive methods). Though complexity and the required number of iterations vary between tasks, we find compelling evidence that even on large deep learning workloads our method finds consistent improvement over episodes.

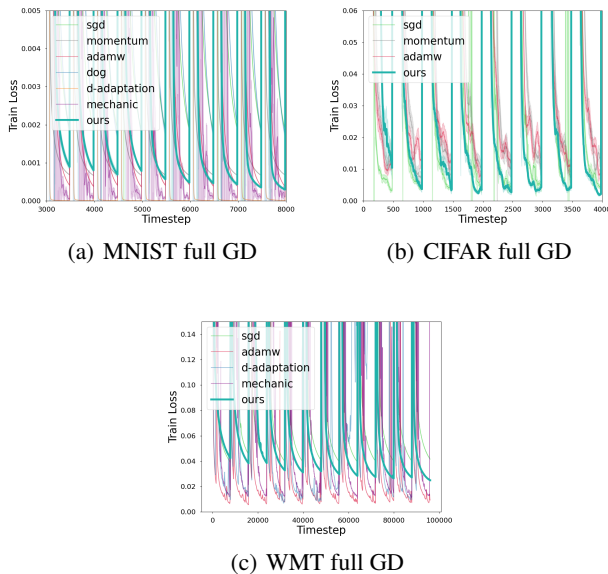


Figure 1. Full gradient descent training with a fixed batch on the three workloads. MNIST and CIFAR are averaged over 5 trials. Losses smoothed with a mean filter.

Furthermore, since meta-optimization is a convex relaxation of the learning to learn problem, we expect that the meta-optimization process is well-behaved in the deterministic setting. This is indeed seen experimentally, as the improve-

ment across episodes is monotonic. Moreover, for each fixed workload, the resulting optimal controller is independent of the initial learning rate  $\eta$  or the hyperparameters of the nonstochastic control algorithm. This stability allows our algorithm to converge properly every time we run it; by contrast, none of the self-tuning baselines converged on the CIFAR workload and only the adaptive methods converged on WMT.

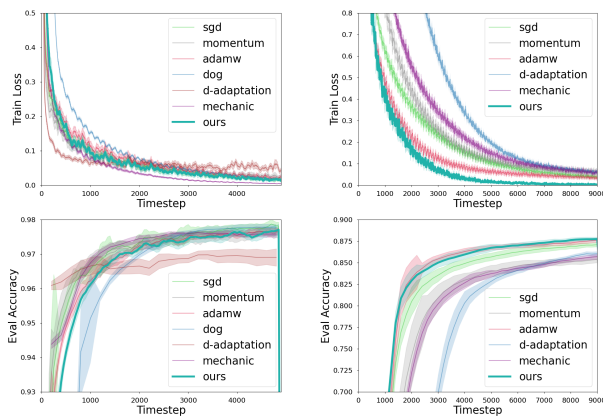
**Stochastic optimization** We also test meta-optimization in the stochastic deep learning setting. Experimentally, we found that the meta-optimization algorithm in the stochastic setting was not as stable on large workloads (see Appendix E.3 for a short explanation). To mitigate this, we take the controllers learned in the deterministic setting and deploy them with frozen parameters to the stochastic setting. Figure 2 shows the performance with this approach. Our method is able to outperform the non-adaptive baselines in terms of training loss, demonstrating that the optimal controllers transfer from deterministic to the stochastic setting. For evaluation metrics, we see that on MNIST and CIFAR our method is able to generalize as well as the baselines. On the WMT workload, however, we once again see a qualitative separation between the adaptive methods and non-adaptive methods (those roughly captured by our benchmark algorithm class). We are investigating this generalization gap in our ongoing work (see Appendix F for a discussion of current questions and future work). When we compare our method to the self-tuning baselines, we see that meta-optimization is consistently competitive in training while methods like DoG, D-Adaptation, and Mechanic can unpredictably suffer on certain workloads.

## 6. Conclusion

In this work, we show that the meta-optimization framework is a promising direction towards the automation of optimization methods in deep learning. On the theoretical side, we extend meta-optimization to the nonconvex setting by leveraging a nonconvex to convex reduction. We give the accompanying convergence guarantee of our method, which depends on the performance of the best algorithm in a class of algorithms. Experimentally, we demonstrated our algorithm’s ability to improve an optimizer across episodes to become competitive against tuned baselines on deep learning workloads. We hope that this initial exploration of meta-optimization in deep learning inspires investigation into the generalization properties, scaling or transfer behavior, and efficient parallelization of the algorithm.

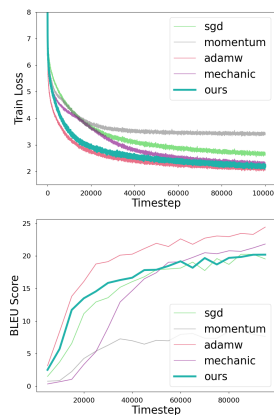
## References

Agarwal, N., Bullins, B., Chen, X., Hazan, E., Singh, K., Zhang, C., and Zhang, Y. Efficient full-matrix adaptive



(a) MNIST stochastic

(b) CIFAR stochastic



(c) WMT stochastic

Figure 2. Stochastic minibatch gradient descent on the three workloads. MNIST and CIFAR are averaged over 5 trials. Losses smoothed with a mean filter.

regularization. In Chaudhuri, K. and Salakhutdinov, R. (eds.), *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pp. 102–110. PMLR, 09–15 Jun 2019.

Anava, O., Hazan, E., and Mannor, S. Online convex optimization against adversaries with memory and application to statistical arbitrage, 2014.

Baydin, A. G., Cornish, R., Rubio, D. M., Schmidt, M., and Wood, F. Online learning rate adaptation with hypergradient descent. *arXiv preprint arXiv:1703.04782*, 2017.

Casgrain, P. and Kratsios, A. Optimizing optimizers: Regret-



- 220 optimal gradient descent algorithms. In *Proceedings of*  
 221 *Thirty Fourth Conference on Learning Theory*, pp. 883–  
 222 926. PMLR, 2021.
- 223 Chen, X. and Hazan, E. Online control for meta-  
 224 optimization. In *Thirty-seventh Conference on Neural*  
 225 *Information Processing Systems*, 2023.
- 227 Cutkosky, A., Defazio, A., and Mehta, H. Mechanic: A  
 228 learning rate tuner. *Advances in Neural Information Pro-*  
 229 *cessing Systems*, 36, 2024.
- 230 Defazio, A. and Mishchenko, K. Learning-rate-free learning  
 231 by d-adaptation. In *International Conference on Machine*  
 232 *Learning*, pp. 7449–7479. PMLR, 2023.
- 234 Flaxman, A. D., Kalai, A. T., and McMahan, H. B. Online  
 235 convex optimization in the bandit setting: gradient de-  
 236 scent without a gradient. In *Proceedings of the Sixteenth*  
 237 *Annual ACM-SIAM Symposium on Discrete Algorithms*,  
 238 SODA '05, pp. 385–394, USA, 2005. Society for Indus-  
 239 trial and Applied Mathematics. ISBN 0898715857.
- 240 Gradu, P., Hallman, J., and Hazan, E. Non-stochastic control  
 241 with bandit feedback. In Larochelle, H., Ranzato, M.,  
 242 Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in*  
 243 *Neural Information Processing Systems*, volume 33, pp.  
 244 10764–10774. Curran Associates, Inc., 2020.
- 246 Gradu, P., Hazan, E., and Minasyan, E. Adaptive regret for  
 247 control of time-varying dynamics. In Matni, N., Morari,  
 248 M., and Pappas, G. J. (eds.), *Proceedings of The 5th An-*  
 249 *annual Learning for Dynamics and Control Conference*, vol-  
 250 *ume 211 of Proceedings of Machine Learning Research*,  
 251 pp. 560–572. PMLR, 15–16 Jun 2023.
- 252 Hazan, E. and Singh, K. Introduction to online nonstochastic  
 253 control, 2023.
- 255 Hazan, E., Singh, K., and Zhang, C. Efficient regret mini-  
 256 mization in non-convex games. In Precup, D. and Teh,  
 257 Y. W. (eds.), *Proceedings of the 34th International Con-*  
 258 *ference on Machine Learning*, volume 70 of *Proceedings*  
 259 *of Machine Learning Research*, pp. 1433–1441. PMLR,  
 260 06–11 Aug 2017.
- 262 Hazan, E., Klivans, A., and Yuan, Y. Hyperparameter opti-  
 263 mization: a spectral approach. In *International Confer-*  
 264 *ence on Learning Representations*, 2018.
- 265 Ivgi, M., Hinder, O., and Carmon, Y. Dog is sgd’s best  
 266 friend: A parameter-free dynamic step size schedule.  
 267 In *International Conference on Machine Learning*, pp.  
 268 14465–14499. PMLR, 2023.
- 270 Lessard, L., Recht, B., and Packard, A. Analysis and de-  
 271 sign of optimization algorithms via integral quadratic  
 272 constraints. *SIAM Journal on Optimization*, 26(1):57–95,  
 273 2016.
- 274 Li, L., Jamieson, K., DeSalvo, G., Rostamizadeh, A., and  
 Talwalkar, A. Hyperband: A novel bandit-based approach  
 to hyperparameter optimization. *Journal of Machine*  
*Learning Research*, 18(185):1–52, 2018.
- Lu, Z., Xia, W., Arora, S., and Hazan, E. Adaptive gra-  
 dient methods with local guarantees. *arXiv preprint*  
*arXiv:2203.01400*, 2022.
- Orabona, F. and Pál, D. Coin betting and parameter-free on-  
 line learning. *Advances in Neural Information Processing*  
*Systems*, 29, 2016.
- Paquette, C., Lin, H., Drusvyatskiy, D., Mairal, J., and  
 Harchaoui, Z. Catalyst for gradient-based nonconvex  
 optimization. In Storkey, A. and Perez-Cruz, F. (eds.),  
*Proceedings of the Twenty-First International Conference*  
*on Artificial Intelligence and Statistics*, volume 84 of  
*Proceedings of Machine Learning Research*, pp. 613–622.  
 PMLR, 09–11 Apr 2018.
- Snoek, J., Larochelle, H., and Adams, R. P. Practical  
 bayesian optimization of machine learning algorithms.  
*Advances in neural information processing systems*, 25,  
 2012.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones,  
 L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention  
 is all you need, 2023.
- Wang, W. and Srebro, N. Stochastic nonconvex optimiza-  
 tion with large minibatches. In Garivier, A. and Kale, S.  
 (eds.), *Proceedings of the 30th International Conference*  
*on Algorithmic Learning Theory*, volume 98 of *Proceed-*  
*ings of Machine Learning Research*, pp. 857–882. PMLR,  
 22–24 Mar 2019.

## A. Additional preliminaries

### A.1. Additional algorithms

The protocol for meta-optimization and the reduction from nonconvex to convex optimization are presented in the following algorithms.

---

#### Algorithm 2 Meta-optimization

---

**Require:** number of episodes  $N$ , number of steps per episode  $T$ , algorithm  $\mathcal{A}$ , initializers  $\{x_{i,1}\}_{i=1}^N$ .

- 1: **for**  $i = 1, \dots, N$  **do**
  - 2:     **for**  $t = 1, \dots, T$  **do**
  - 3:         Play  $x_{i,t} = \mathcal{A}(f_{1,1}, x_{1,1}, \dots, f_{i,t-1}, x_{i,t-1}) \in \mathbb{R}^d$  if  $t > 1$ ; else play  $x_{i,1}$ .
  - 4:         Receive a mini-batch of examples  $B_{i,t}$  of size  $b$ .
  - 5:         Obtain the loss function  $f_{i,t} = \frac{1}{b} \sum_{j \in B_{i,t}} f_{i,j}$ .
  - 6:         Suffer loss  $f_{i,t}(x_{i,t})$  and compute the stochastic gradient  $\tilde{\nabla} f_{i,t} = \frac{1}{b} \sum_{j \in B_{i,t}} \nabla f_{i,j}(x_{i,t})$ .
- 

---

#### Algorithm 3 Reduction from nonconvex to convex optimization

---

**Require:** epoch number  $K$ , number of inner steps  $S$  such that  $T = KS$ , convex optimization algorithm  $\mathcal{A}$ , smoothness parameter  $\beta$ , initial point  $x_1$ , nonconvex function  $f$ .

- 1: **for**  $k = 1, \dots, K$  **do**
  - 2:     Consider the function  $f_k(x) = f(x) + \beta \|x - x_k\|^2$ .
  - 3:     Starting from  $x_k$ , apply  $\mathcal{A}$  for  $S$  steps on  $f_k$  with mini-batch access, obtain  $x_{k+1}$ .
  - 4: **return**  $x^* = \arg \min_{\{x_k\}_{k=1}^K} \|\nabla f(x_k)\|$ .
- 

## B. Algorithm details

The framework of meta-optimization applies online control methods to a particular dynamical system that describes the optimization process. We apply meta-optimization to the stochastic,  $\ell_2$ -regularized strongly convex functions  $f_{i,k,s}$ , and describe the dynamical system below. The dynamical system is similar to the one for smooth convex optimization put forth in (Chen & Hazan, 2023).

**The dynamical system** For each episode  $i$ , denote  $H_{i,k,s}$  to be the matrix that satisfies

$$\nabla f_{i,k,s}(x_{i,k,s}) = H_{i,k,s}(x_{i,k,s} - x_{i,k,s-1}) + \nabla f_{i,k,s}(x_{i,k,s-1}). \quad (1)$$

If each  $f_{i,j}$  is quadratic, then  $H_{i,k,s}$  has the following explicit form,

$$H_{i,k,s} = \frac{1}{b} \sum_{j \in B_{i,k,s}} \nabla^2 f_{i,j} + 2\beta.$$

For general smooth functions that are twice differentiable, this matrix exists and each row contains certain second-order information of  $f_{i,k,s}$ . To see this, observe that we can apply the mean value theorem to each coordinate of  $\nabla f_{i,k,s}$  to obtain  $H_{i,k,s}$ .

The linear dynamical system we consider is

$$\begin{bmatrix} x_{i,k,s+1} \\ x_{i,k,s} \\ \nabla f_{i,k,s}(x_{i,k,s}) \end{bmatrix} = \begin{bmatrix} (1-\delta)I & 0 & -\eta I \\ I & 0 & 0 \\ H_{i,k,s} & -H_{i,k,s} & 0 \end{bmatrix} \times \begin{bmatrix} x_{i,k,s} \\ x_{i,k,s-1} \\ \nabla f_{i,k,s-1}(x_{i,k,s-1}) \end{bmatrix} \quad (2)$$

$$+ \begin{bmatrix} I & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \times u_{i,k,s} + \begin{bmatrix} 0 \\ 0 \\ \nabla f_{i,k,s}(x_{i,k,s-1}) \end{bmatrix}. \quad (3)$$

For notational convenience, we write the dynamical system as

$$z_{i,k,s+1} = A_{i,k,s} z_{i,k,s} + B u_{i,k,s} + w_{i,k,s},$$

where  $A_{i,k,s}$  is the system dynamics,  $B$  is a constant control-input matrix, and  $w_{i,k,s}$  is the non-stochastic disturbance that contains the gradient.

The system above is linear time-varying (LTV), and we introduce the following notion of stability for LTV systems standard in the non-stochastic control literature (Gradu et al., 2023; Chen & Hazan, 2023).

**Definition 3** (Sequentially stable). A time-varying linear dynamical system with dynamics  $A_1, \dots, A_T$  is  $(\kappa, \gamma)$  sequentially stable if for all intervals  $I = [r, s] \subseteq [T]$ ,  $\|\prod_{t=s}^r A_t\| \leq \kappa^2(1 - \gamma)^{|I|}$ .

**Assumption 2.** We assume that the dynamical system (2) is  $(\kappa, \gamma)$  sequentially stable with  $\kappa \geq 1$ .

We in addition make the following two assumptions on the system dynamics and the iterates, following the meta-optimization framework. In meta-optimization, the iterates are re-initialized to starting points with bounded norm in each episode; in our algorithm, we essentially have  $NK$  episodes, and the iterates are re-initialized with  $x_{i,k}$  at the beginning of each episode. We note that since each  $f_{i,k,s}$  is strongly convex, the iterates effectively stay within a bounded region, potentially justifying the latter assumption below.

**Assumption 3.** For all  $i \in [N]$ ,  $k \in [K]$ ,  $s \in [S]$ ,  $\rho(H_{i,k,s}) \leq \beta$ .

**Assumption 4.** For all  $i \in [N]$ ,  $k \in [K]$ ,  $\|x_{i,k}\| \leq R$ .

**The benchmark algorithm class** The benchmark algorithm class we consider consists of methods whose updates are linear functions of past gradients. Since we view optimization from the perspective of online control, this class of methods correspond to a general class of controllers that often appear in the online control literature. This class of controllers is called Disturbance-feedback controllers (DFCs), and for linear time-invariant systems, they can approximate any stabilizing state-feedback controllers. Consequently, if our objective functions are quadratic with uniformly bounded Hessians, this benchmark class of algorithms can simulate gradient descent, momentum, and preconditioned methods with a fixed preconditioner. Since these algorithms have to be stabilizing on the dynamical system described above, only certain values of learning rate, momentum, and preconditioners are allowed. For example, the learning rate can be at most  $O(1/\beta)$ , and the preconditioner  $P$  needs to satisfy  $\rho(PH) < 1/8$ , where  $H$  is the Hessian. For more detailed specifications of the range of parameters, see the full version of (Chen & Hazan, 2023).

### C. Proofs for Section 3

*Proof of Theorem 2.* Since all  $f_{i,k,s}(x)$  are convex and smooth, we can use the bandit meta-optimization algorithm (Algorithm 4) as the black-box optimizer  $\mathcal{A}$ . For any policy  $\pi \in \Pi$ , let  $x_{i,k,s}^\pi$  be its updates, and by Corollary 6 we have the following guarantee:

$$\min_{\pi \in \Pi} \mathbb{E} \left[ \sum_{i=1}^N \sum_{k=1}^K \sum_{s=1}^S (f_{i,k,s}(x_{i,k,s}) - f_{i,k,s}(x_{i,k,s}^\pi)) \right] = \tilde{O}((NKS)^{\frac{3}{4}}). \quad (4)$$

Denote  $\pi^* = \operatorname{argmin}_{\pi \in \Pi} \mathbb{E} \left[ \sum_{i=1}^N \sum_{k=1}^K \sum_{s=1}^S f_{i,k,s}(x_{i,k,s}^\pi) \right]$  as the optimal algorithm in  $\Pi$ .

By Theorem A.3 in (Agarwal et al., 2019), since each  $f_{i,k}$  is smooth and strongly convex, for any  $i, k$ ,

$$f_{i,k}(x_{i,k}) - \min_x f_{i,k}(x) \geq \frac{\|\nabla f_{i,k}(x_{i,k})\|^2}{6\beta}.$$

In addition, we have the following decomposition

$$\begin{aligned}
 f_i(x_{i,k}) - f_i(x_{i,k+1}) &\geq f_{i,k}(x_{i,k}) - \min_x f_{i,k}(x) - (f_{i,k}(x_{i,k+1}) - \min_x f_{i,k}(x)) \\
 &= f_{i,k}(x_{i,k}) - \min_x f_{i,k}(x) - \left( \frac{1}{S} \sum_{s=1}^S f_{i,k}(x_{i,k,s}^{\pi^*}) - \min_x f_{i,k}(x) \right) \\
 &\quad + \left( \frac{1}{S} \sum_{s=1}^S f_{i,k}(x_{i,k,s}^{\pi^*}) - f_{i,k}(x_{i,k+1}) \right) \\
 &\geq f_{i,k}(x_{i,k}) - \min_x f_{i,k}(x) - \left( \frac{1}{S} \sum_{s=1}^S f_{i,k}(x_{i,k,s}^{\pi^*}) - \min_x f_{i,k}(x) \right) \\
 &\quad + \left( \frac{1}{S} \sum_{s=1}^S f_{i,k}(x_{i,k,s}^{\pi^*}) - \frac{1}{S} \sum_{s=1}^S f_{i,k}(x_{i,k,s}) \right).
 \end{aligned}$$

Rearranging the terms, we have

$$f_{i,k}(x_{i,k}) - \min_x f_{i,k} \leq f_i(x_{i,k}) - f_i(x_{i,k+1}) + \frac{1}{S} \sum_{s=1}^S f_{i,k}(x_{i,k,s}^{\pi^*}) - \min_x f_{i,k}(x) \quad (5)$$

$$+ \left( \frac{1}{S} \sum_{s=1}^S f_{i,k}(x_{i,k,s}) - \frac{1}{S} \sum_{s=1}^S f_{i,k}(x_{i,k,s}^{\pi^*}) \right). \quad (6)$$

Using the lower bound on the left hand side and summing up over  $k = 1, 2, \dots, K$ ,

$$\begin{aligned}
 \sum_{k=1}^K \frac{\|\nabla f_{i,k}(x_{i,k})\|^2}{6\beta} &\leq M + \frac{1}{S} \sum_{s,k} \left( f_{i,k}(x_{i,k,s}^{\pi^*}) - \min_x f_{i,k}(x) \right) \\
 &\quad + \left( \frac{1}{S} \sum_{s,k} f_{i,k}(x_{i,k,s}) - \frac{1}{S} \sum_{s,k} f_{i,k}(x_{i,k,s}^{\pi^*}) \right).
 \end{aligned}$$

Summing over  $i$  and taking an average,

$$\begin{aligned}
 \frac{1}{NK} \sum_{i,k} \|\nabla f_i(x_{i,k})\|^2 &\leq O\left(\frac{1}{K}\right) + \frac{1}{NKS} \left( \sum_{i,s,k} f_{i,k}(x_{i,k,s}^{\pi^*}) - \min_x f_{i,k}(x) \right) \\
 &\quad + \frac{1}{NKS} \left( \sum_{i,s,k} f_{i,k}(x_{i,k,s}) - \sum_{i,s,k} f_{i,k}(x_{i,k,s}^{\pi^*}) \right).
 \end{aligned}$$

The theorem follows by taking an expectation over the randomness of the batches, and using the guarantee (4).  $\square$

**Refinement of Theorem 2** We use an approach inspired by (Agarwal et al., 2019) to refine our guarantee. Define  $\lambda_{i,k}$  as the ratio

$$\frac{1}{S} \sum_{s=1}^S f_{i,k}(x_{i,k,s}^{\pi^*}) - \min_x f_{i,k}(x) \leq \lambda_{i,k} \sqrt{\frac{f_{i,k}(x_{i,k}) - \min_x f_{i,k}(x)}{\beta S}}. \quad (7)$$

Then in the inequality (5) above, we have

$$\begin{aligned}
 f_{i,k}(x_{i,k}) - \min_x f_{i,k} - \lambda_{i,k} \sqrt{\frac{f_{i,k}(x_{i,k}) - \min_x f_{i,k}(x)}{\beta S}} &\leq f_i(x_{i,k}) - f_i(x_{i,k+1}) \\
 &\quad + \left( \frac{1}{S} \sum_{s=1}^S f_{i,k}(x_{i,k,s}) - \frac{1}{S} \sum_{s=1}^S f_{i,k}(x_{i,k,s}^{\pi^*}) \right).
 \end{aligned}$$



Observe that when a variable  $y$  satisfies  $y^2 - ay \leq b$ , we can complete the squares and obtain  $(y - \frac{a}{2})^2 \leq b + \frac{a^2}{4}$ . Taking a square root, we have  $y \leq \sqrt{b} + a$ , and squaring both sides, we arrive at  $y^2 \leq 2b + 2a^2$ . Using this result, we have

$$\frac{\|\nabla f_i(x_{i,k})\|^2}{6\beta} \leq \frac{\lambda_{i,k}^2}{\beta S} + 2 \left( f_i(x_{i,k}) - f_i(x_{i,k+1}) + \frac{1}{S} \sum_{s=1}^S f_{i,k}(x_{i,k,s}) - \frac{1}{S} \sum_{s=1}^S f_{i,k}(x_{i,k,s}^{\pi^*}) \right).$$

Summing over  $i, k$ , and taking an average,

$$\frac{1}{NK} \sum_{i,k} \|\nabla f_i(x_{i,k})\|^2 \leq O\left(\frac{1}{K}\right) + O\left(\frac{\sum_{i,k} \lambda_{i,k}^2}{NKS}\right) + \frac{6\beta}{NSK} \sum_{i,k} f_{i,k}(x_{i,k,s}) - f_{i,k}(x_{i,k,s}^{\pi^*}).$$

In particular, we have that

$$\mathbb{E} \left[ \frac{1}{N} \sum_{i=1}^N \min_k \|\nabla f_i(x_{i,k})\|^2 \right] \leq O\left(\frac{1}{K} + \frac{\sum_{i,k} \lambda_{i,k}^2}{NKS} + \tilde{O}((NKS)^{-\frac{1}{4}})\right).$$

Let  $\lambda_0$  denote an upper bound of  $\lambda_{i,k}$ . As  $N$  grows large, the right hand side is dominated by the first two terms, and therefore in this regime we can write

$$\mathbb{E} \left[ \frac{1}{N} \sum_{i=1}^N \min_k \|\nabla f_i(x_{i,k})\|^2 \right] \leq O\left(\frac{1}{K} + \frac{\lambda_0^2}{S}\right).$$

As defined in Equation (7),  $\lambda_{i,k}$  scales with the function value optimality gap of the average iterate under  $\pi^*$ , and  $\sqrt{\frac{1}{\beta S}}$ . By the online to batch reduction, SGD on strongly convex functions converges at a rate of  $\tilde{O}(\frac{1}{S})$ . The learning rate of SGD that attains this rate depends on the strong convexity parameter and the gradient upper bound of the loss function. If all the  $f_i$ 's have the same smoothness parameter, and under the assumption of bounded domain, taking  $\pi^*$  to be SGD with the optimal learning rate,  $\lambda_0$  can be as small as  $\tilde{O}\left(\frac{1}{\sqrt{S}}\right)$ .

## D. Bandit meta-optimization

We give the details of the bandit meta-optimization algorithm in this section. For any set  $\mathcal{M}$  and  $\delta_M > 0$ , define the Minkowski subset  $\mathcal{M}_{\delta_M} = \{x : \frac{1}{1-\delta_M}x \in \mathcal{M}\}$ , and let  $\mathbb{S}_1^d$  be the  $d$ -dimensional unit sphere.

**Theorem 4** (Theorem 5.1 in (Gradu et al., 2020), Theorem 3.3 in (Chen & Hazan, 2023)). *Under Assumptions 1, 2, 3, 4, Algorithm 4 with  $\eta \leq 1$ ,  $L = \Theta(\log NKS)$ , and setting  $\eta_{i,k,s}^M = \Theta((N(i-1) + K(k-1) + s)^{-3/4} L^{-3/2} G^{-2/3})$ , and perturbation constant  $\delta_M = \Theta((NKS)^{-1/4} L^{-1/2})$  gives the guarantee*

$$\mathbb{E} \left[ \sum_{i,k,s} f_{i,k,s}(x_{i,k,s}) \right] - \min_{\mathcal{A} \in \Pi} \sum_{i,k,s} f_{i,k,s}(x_{i,k,s}^{\mathcal{A}}) \leq \tilde{O}((NKS)^{3/4}),$$

where  $\tilde{O}$ ,  $\Theta$  contain polynomial factors in  $\gamma^{-1}, \beta, \kappa, R, b, d, M$ , and  $\tilde{O}$  in addition contains logarithmic factors in  $K, S, N$ . The benchmark algorithm class  $\Pi$  is the class of DFCs discussed in Appendix B.

The theorem above guarantees the performance of Algorithm 4 under any adversarially chosen functions  $f_{i,k,s}$ . However, for our setting of nonconvex stochastic optimization, it is more useful to derive a guarantee in expectation for randomly chosen functions. We show that such extension is possible, and we start from the bandit convex optimization with memory (BCOWM) problem (Gradu et al., 2020). The guarantee for bandit online control and hence bandit meta-optimization can be derived as corollaries of the BCOWM guarantee.

Consider the basic online learning with memory problem under bandit feedback, where the loss functions  $f_t$  are  $\beta$ -smooth,  $G$ -Lipschitz, and  $M$ -bounded. The domain of decisions  $\mathcal{K}$  has diameter  $D$ , and the  $f_t$ 's are random functions determined by previous decisions of the player. Same as before, let  $\mathcal{K}_\delta$  be the Minkowski set of  $\mathcal{K}$ , and let  $\mathbb{S}_1^d$  be the  $d$ -dimensional unit sphere. The algorithm below, proposed by (Gradu et al., 2020), is an application of zeroth-order method (Flaxman et al., 2005) to the Online Convex Optimization with Memory (OCOWM) setting (Anava et al., 2014).

**Algorithm 4** Bandit meta-optimization

**Require:** episode number  $K$ , system parameters  $\eta, \delta, \kappa, \gamma$ , learning rates  $\{\eta_{i,k,s}^M\}$ , history length  $L$ ,  $\delta_M$ , starting points

$\{x_{i,1}\}_{i=1}^N$ .

- 1: Set:  $\mathcal{M} = \{M = \{M^1, \dots, M^L\} : \|M^l\| \leq \kappa^3(1 - \gamma)^l\}$ .
- 2: Initialize any  $M_{1,1} = \dots = M_{L,1} \in \mathcal{M}_{\delta_M}$ ,  $z_{i,1,1} = [x_{i,1}^\top \ x_{i,1}^\top \ 0]^\top$ .
- 3: Sample  $\epsilon_{1,1}, \dots, \epsilon_{L,1} \in_{\mathbb{R}} \mathbb{S}_1^{L \times 3d \times 3d}$ , set  $\widetilde{M}_{l,1} = M_{l,1} + \delta_M \epsilon_{l,1}$  for  $l = 1, \dots, L$ .
- 4: **for**  $i = 1, \dots, N$  **do**
- 5:     **If**  $i > 1$ , set  $z_{i,1,1} = z_{i-1,K,S+1}$ ,  $M_{i,1,1} = M_{i-1,K,S+1}$ .
- 6:     **for**  $k = 1, \dots, K$  **do**
- 7:         **If**  $k > 1$ , set  $z_{i,k,1} = z_{i,k-1,S+1}$ ,  $M_{i,k,1} = M_{i,k-1,S+1}$ .
- 8:         **for**  $s = 1, \dots, S$  **do**
- 9:             Choose  $u_{i,k,s} = \sum_{l=1}^L \widetilde{M}_{l,k,s}^l w_{i,k,s-1}$ .
- 10:            Receive  $f_{i,k,s}$ , compute  $w_{i,k,s} = \nabla f_{i,k,s}(x_{i,k,s-1})$ . **If**  $s = S$ , compute  $x_{i,k+1} = \frac{1}{S} \sum_{s=1}^S x_{i,k,s}$ , and

$$w_{i,k,S} = \begin{bmatrix} x_{i,k+1} - ((1 - \delta)x_{i,k,S} - \eta \nabla f_{i,k,S-1}(x_{i,k,S-1}) + \bar{u}_{i,k,S}) \\ x_{i,k+1} - x_{i,k,S} \\ \nabla f_{i,k,S}(x_{i,k,S-1}) - \nabla f_{i,k,S}(x_{i,k,S}) \end{bmatrix}, \quad (8)$$

where  $\bar{u}_{i,k,S}$  is the first  $d$  coordinates of the control signal  $u_{i,k,S}$ .

- 11:     Suffer control cost  $f_{i,k,S}(x_{i,k,S})$ .
- 12:     Store the gradient estimator  $g_{i,k,s} = \frac{9d^2 L}{\delta_M} f_{i,k,s}(x_{i,k,s}) \sum_{l=1}^L \epsilon_{i,k,s-l}$  **if**  $s \geq L$ , **else** 0.
- 13:     Perform gradient update on the controller parameters:

$$M_{i,k,s+1} = \Pi_{\mathcal{M}_{\delta_M}}(M_{i,k,s} - \eta_{i,k,s}^M \cdot g_{i,k,s-L}).$$

- 14:     Sample  $\epsilon_{i,k,s+1} \in_{\mathbb{R}} \mathbb{S}_1^{L \times 3d \times 3d}$ , set  $\widetilde{M}_{l,k,s+1} = M_{l,k,s+1} + \delta_M \epsilon_{i,k,s+1}$ .
- 15:     **If**  $k = K$ , compute  $w_{i,K,S}$  similar to (8), so the next state evolves to  $z_{i,K,S+1} = [x_{i,1,1}^\top \ x_{i,1,1}^\top \ 0]^\top$ .

**Algorithm 5** BCO with Memory

**Require:** Decision set  $\mathcal{K}$ , time horizon  $T$ , history length  $L$ , learning rates  $\{\eta_t\}$  and noise magnitude  $\delta$ .

- 1: Initialize  $x_1 = \dots = x_L \in \mathcal{K}_\delta$  arbitrarily, and sample noise  $u_1, \dots, u_L \in \mathbb{S}_1^d$ .
- 2: Set  $y_i = x_i + \delta u_i$  for  $i = 1, \dots, L$ ,  $g_i = 0$  for  $i = 1, \dots, L - 1$ .
- 3: Predict  $y_i$  for  $i = 1, \dots, L - 1$ .
- 4: **for**  $t = L, \dots, T$  **do** Play  $y_t$ , and suffer loss  $f_t(y_{t-L+1:t})$ .
- 5:     Store gradient estimate  $g_t = \frac{d}{\delta} f_t(y_{t-L+1:t}) \sum_{i=0}^{L-1} u_{t-i}$ .
- 6:     Set  $x_{t+1} = \Pi_{\mathcal{K}_\delta} [x_t - \eta_t \cdot g_{t-L+1}]$ .
- 7:     Sample  $u_{t+1} \in \mathbb{S}_1^d$ , set  $y_{t+1} = x_{t+1} + \delta u_{t+1}$ .

**Theorem 5.** Suppose the loss functions  $f_t$  are random functions determined by previous iterates  $y_1, \dots, y_{t-1}$ . Let  $O$  denote polynomial dependence on  $D, d, M, L, G, \beta$ . Taking  $\eta_t = O(t^{-3/4}), \delta = O(T^{-1/4})$ , Algorithm 5 produces  $y_t$ 's that satisfy

$$\mathbb{E} \left[ \sum_{t=L}^T f_t(y_{t-L+1:t}) \right] - \min_{x \in \mathcal{K}} \mathbb{E} \left[ \sum_{t=L}^T f_t(x, \dots, x) \right] \leq O(T^{3/4}).$$

*Proof.* We largely follow the proof of Theorem 3.1 in (Gradu et al., 2020). Let  $x^*$  be any comparator in  $\mathcal{K}$ , and  $x_\delta^*$  be the projection of  $x^*$  in the Minkowski set. Let  $\tilde{f}(x) = f(x, \dots, x)$  be the shorthand notation.

$$\mathbb{E} \left[ \sum_{t=L}^T f_t(y_{t-L+1:t}) - \sum_{t=L}^T \tilde{f}_t(x^*) \right] = \mathbb{E} \left[ \sum_{t=L}^T (f_t(y_{t-L+1:t}) - \tilde{f}_t(x^*)) \right] - \mathbb{E} \left[ \sum_{t=L}^T \tilde{f}_{t-L+1}(x_t) - \tilde{f}_{t-L+1}(x_\delta^*) \right] \quad (9)$$

$$+ \mathbb{E} \left[ \sum_{t=L}^T \tilde{f}_{t-L+1}(x_t) - \tilde{f}_{t-L+1}(x_\delta^*) \right] \quad (10)$$

We bound (9) and (10) separately. We start with (9), which can be bounded for any sequence of random variables  $u_1, \dots, u_T$ . Fix  $u_1, \dots, u_T$ , we have

$$\begin{aligned} f_t(y_{t-L+1:t}) - \tilde{f}_t(x_{t+L-1}) &= f_t(x_{t-L+1:t} + \delta u_{t-L+1:t}) - \tilde{f}_t(x_{t+L-1}) \\ &\leq f_t(x_{t-L+1:t}) - \tilde{f}_t(x_{t+L-1}) + \delta G \sqrt{L} \\ &\leq G \|x_{t-L+1:t} - (x_{t+L-1}, \dots, x_{t+L-1})\| + \delta G \sqrt{L} \\ &\leq \frac{2dMGL^2 \eta_{t-L+1}}{\delta} + \delta G \sqrt{L}, \end{aligned}$$

where the first and second inequalities hold by the Lipschitz property of  $f_t$ , and the last inequality is due to Lemma 7. Furthermore, the Lipschitz property of  $f_t$  gives

$$|\tilde{f}_t(x_\delta^*) - \tilde{f}_t(x^*)| \leq G \|(x_\delta^*, \dots, x_\delta^*) - (x^*, \dots, x^*)\| \leq \delta GD \sqrt{L}.$$

Putting the two inequalities together, and accounting for the shift in the index of  $\tilde{f}_{t-L+1}(x_\delta^*)$ ,

$$(2) \leq 2\delta GD \sqrt{L} T + \frac{2dMGL^2}{\delta} \sum_{t=1}^T \eta_t + 2LM.$$

The term (10) can be decomposed as follows,

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=L}^T \tilde{f}_{t-L+1}(x_t) - \tilde{f}_{t-L+1}(x_\delta^*) \right] &\leq \mathbb{E} \left[ \sum_{t=L}^T \nabla \tilde{f}_{t-L+1}(x_t)^\top (x_t - x_\delta^*) \right] \\ &= \mathbb{E} \left[ \sum_{t=L}^T (g_{t-L+1} + (\mathbb{E}_{u_{t-2L+2:t-L+1}}[g_{t-L+1}] - g_{t-L+1})) \right] \\ &\quad + \mathbb{E} \left[ (\nabla \tilde{f}_{t-L+1}(x_t) - \mathbb{E}_{u_{t-2L+2:t-L+1}}[g_{t-L+1}])^\top (x_t - x_\delta^*) \right]. \end{aligned}$$

Since  $x_{t+1}$  is a projected gradient descent step from  $x_t$  with the gradient estimator  $g_{t-L+1}$ , we have

$$2g_{t-L+1}^\top (x_t - x_\delta^*) \leq \frac{1}{\eta_t} (\|x_t - x_\delta^*\|^2 - \|x_{t+1} - x_\delta^*\|^2) + \eta_t \|g_{t-L+1}\|^2, \quad (\text{Eq. 3.2 in (Gradu et al., 2020)})$$

and

$$\begin{aligned} \sum_{t=L}^T g_{t-L+1}^\top (x_t - x_\delta^*) &\leq \frac{1}{2} \sum_{t=L}^T \left( \|x_t - x_\delta^*\|^2 \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) + \eta_t \|g_{t-L+1}\|^2 \right) + \frac{\|x_L - x_\delta^*\|^2}{\eta_{L-1}} \\ &\leq \frac{D^2}{2} \left( \frac{1}{\eta_T} + \frac{1}{\eta_{L-1}} \right) + \frac{d^2 M^2 L}{2\delta^2} \sum_{t=L}^T \eta_t, \end{aligned}$$

where the bound on  $\|g_t\|$  is in the proof of Lemma 7.

By Lemma 8, we also have

$$\begin{aligned} & \mathbb{E} \left[ (\nabla \tilde{f}_{t-L+1}(x_t) - \mathbb{E}_{u_{t-2L+2:t-L+1}}[g_{t-L+1}])^\top (x_t - x_\delta^*) \right] \\ & \leq \mathbb{E} \left[ \|\nabla \tilde{f}_{t-L+1}(x_t) - \mathbb{E}_{u_{t-2L+2:t-L+1}}[g_{t-L+1}]\| \|x_t - x_\delta^*\| \right] \\ & \leq D \mathbb{E} \left[ \|\nabla \tilde{f}_{t-L+1}(x_t) - \mathbb{E}_{u_{t-2L+2:t-L+1}}[g_{t-L+1}]\| \right] \\ & \leq 2\eta_{t-L+1} \frac{dML^{5/2}\beta D}{\delta} + \frac{d\delta L^2 D}{2}. \end{aligned}$$

Lastly, by Lemma 9,

$$\mathbb{E} \left[ \sum_{t=L}^T (\mathbb{E}_{u_{t-2L+2:t-L+1}}[g_{t-L+1}] - g_{t-L+1})^\top (x_t - x_\delta^*) \right] \leq \frac{2d^2 M^2 L^2}{\delta^2} \sum_{t=L}^T \eta_{t-L+1}.$$

Summing up the the three inequalities, (3) can be bounded by

$$(3) \leq \frac{D^2}{\eta_T} + \left( \frac{3d^2 M^2 L^2}{\delta^2} + \frac{2dML^{5/2}\beta D}{\delta} \right) \sum_{t=1}^T \eta_t + \frac{d\delta L^2 DT}{2}.$$

Putting everything together, the expected regret can be bounded by

$$\begin{aligned} \mathbb{E} \left[ \sum_{t=L}^T f_t(y_{t-L+1:t}) - \sum_{t=L}^T \tilde{f}_t(x^*) \right] & \leq 2LM + \frac{D^2}{\eta_T} + \left( \frac{3d^2 M^2 L^2}{\delta^2} + \frac{4dGML^{5/2}\beta D}{\delta} \right) \sum_{t=1}^T \eta_t \\ & \quad + \frac{5d\delta GL^2 DT}{2}. \end{aligned}$$

Let  $O$  denote polynomial dependence on  $D, d, M, L, G, \beta$ . Taking  $\eta_t = O(t^{-3/4}), \delta = O(T^{-1/4})$ , we have  $\sum_{t=1}^T \eta_t \leq O(T^{1/4})$ , and

$$\mathbb{E} \left[ \sum_{t=L}^T f_t(y_{t-L+1:t}) - \sum_{t=L}^T \tilde{f}_t(x^*) \right] \leq O(T^{3/4}).$$

□

**Corollary 6.** *Under the same assumptions as Theorem 5, and setting  $\eta_{i,k,s}^M, \delta_M, L$  correctly, Algorithm 4 produces a sequence of controls  $M_{i,k,s}$  that satisfy*

$$\mathbb{E} \left[ \sum_{i,k,s} f_{i,k,s}(x_{i,k,s}) \right] \leq \min_{\pi \in \Pi} \mathbb{E} \left[ \sum_{i,k,s} f_{i,k,s}(x_{i,k,s}^\pi) \right] + \tilde{O}((NKS)^{3/4}).$$

**Lemma 7** (Variant of Lemma A.6 in (Gradu et al., 2020)). *Fixing  $u_1, \dots, u_T$ , Algorithm 5 produces a sequence of  $x_t$  such that*

$$\|x_{t-H+1:t} - (x_{t+H-1}, \dots, x_{t+H-1})\| \leq 2\eta_{t-H+1} \frac{dCH^2}{\delta}.$$

660 *Proof.* Fix  $u_1, \dots, u_T$ .

$$\begin{aligned}
 661 \quad & \|x_{t-L+1:t} - (x_{t+L-1}, \dots, x_{t+L-1})\|^2 = \sum_{i=0}^{L-1} \|x_{t-i} - x_{t+L-1}\|^2 \\
 662 \quad & \\
 663 \quad & \\
 664 \quad & \\
 665 \quad & \leq \sum_{i=0}^{L-1} \left( \sum_{j=1}^{i+L-1} \|x_{t+L-j} - x_{t+L-j-1}\| \right)^2 \\
 666 \quad & \\
 667 \quad & \\
 668 \quad & \leq \sum_{i=1}^{L-1} \left( \sum_{j=1}^{i+L-1} \eta_{t+L-1-j} \|g_{t-j}\| \right)^2 \\
 669 \quad & \\
 670 \quad & \\
 671 \quad & \leq \eta_{t-L+1}^2 \sum_{i=1}^{L-1} \left( \sum_{j=1}^{i+L-1} \|g_{t-j}\| \right)^2 \\
 672 \quad & \\
 673 \quad & \leq 4\eta_{t-L+1}^2 L^3 \frac{d^2 M^2 L}{\delta^2}, \\
 674 \quad & \\
 675 \quad & \\
 676 \quad & 
 \end{aligned}$$

677 where the second-to-last inequality holds since the stepsize is non-increasing, and the last inequality is true because for any  
678  $t$ ,

$$680 \quad \|g_t\| = \frac{d}{\delta} \|f_t(y_{t-L-1:t}) \sum_{i=0}^{L-1} u_{t-i}\| \leq \frac{dM}{\delta} \sqrt{L}.$$

683 The lemma follows by taking a square root on both sides.  $\square$

684 **Lemma 8** (Variant of Lemma A.12 in (Gradu et al., 2020)). *Conditioned on  $u_1, \dots, u_{t-2L+1}$ , for any sequence of*  
685  *$u_{t-2L+2:t-L+1}$  that determines  $x_t$ ,*

$$687 \quad \|\mathbb{E}_{u_{t-2L+2:t-L+1}}[g_{t-L+1}] - \nabla \tilde{f}_{t-L+1}(x_t)\| \leq 2\eta_{t-L+1} \frac{dML^{5/2}\beta}{\delta} + \frac{d\delta L^2}{2}.$$

690 *Proof.* By definition, after fixing  $u_{1:t-2L+1}$ , the following quantities and functions are deterministic:  $g_1, \dots, g_{t-2L+1}$ ,  
691  $x_1, \dots, x_{t-L+1}$ , and  $f_1, \dots, f_{t-L+2}$ . For a function  $f$  that takes in  $L$  inputs, let  $\nabla_i f(x_0, \dots, x_{L-1}) = \frac{\partial f(x_0, \dots, x_{L-1})}{\partial x_i}$   
692 denote the gradient of  $f$  with respect to  $x_i$ .

693 By triangle inequality,

$$\begin{aligned}
 694 \quad & \|\mathbb{E}_{u_{t-2L+2:t-L+1}}[g_{t-L+1}] - \nabla \tilde{f}_{t-L+1}(x_t)\| \leq \|\mathbb{E}_{u_{t-2L+2:t-L+1}}[g_{t-L+1}] - \sum_{i=0}^{L-1} \nabla_i f_{t-L+1}(x_{t-2L+2:t-L+1})\| \\
 695 \quad & \\
 696 \quad & + \left\| \sum_{i=0}^{L-1} \nabla_i f_{t-L+1}(x_{t-2L+2:t-L+1}) - \nabla \tilde{f}_{t-L+1}(x_t) \right\| \\
 697 \quad & \\
 698 \quad & \leq \frac{d\delta L^2}{2} + \left\| \sum_{i=0}^{L-1} \nabla_i f_{t-L+1}(x_{t-2L+2:t-L+1}) - \nabla \tilde{f}_{t-L+1}(x_t) \right\|, \\
 699 \quad & \\
 700 \quad & \\
 701 \quad & \\
 702 \quad & \\
 703 \quad & 
 \end{aligned}$$

704 where the second inequality is due to Corollary A.10 in (Gradu et al., 2020). The norm of the sum can be bounded by  
705 smoothness: for any sequence of  $u_{t-2L+2:t-L+1}$ ,

$$\begin{aligned}
 706 \quad & \left\| \sum_{i=0}^{L-1} \nabla_i f_{t-L+1}(x_{t-2L+2:t-L+1}) - \nabla \tilde{f}_{t-L+1}(x_t) \right\|^2 \leq L \sum_{i=0}^{L-1} \|\nabla_i f_{t-L+1}(x_{t-2L+2:t-L+1}) - \nabla_i f_{t-L+1}(x_{t:t})\|^2 \\
 707 \quad & \\
 708 \quad & = L \|\nabla f_{t-L+1}(x_{t-2L+2:t-L+1}) - \nabla f_{t-L+1}(x_{t:t})\|^2 \\
 709 \quad & \\
 710 \quad & \leq L\beta^2 \|x_{t-2L+2:t-L+1} - (x_t, \dots, x_t)\|^2 \\
 711 \quad & \\
 712 \quad & \leq 4\eta_{t-L+1}^2 \frac{d^2 M^2 L^5 \beta^2}{\delta^2}, \\
 713 \quad & \\
 714 \quad & 
 \end{aligned}$$



by Lemma 7. Hence

$$\left\| \sum_{i=0}^{L-1} \nabla_i f_{t-L+1}(x_{t-2L+2:t-L+1}) - \nabla \tilde{f}_{t-L+1}(x_t) \right\| \leq 2\eta_{t-L+1} \frac{dML^{5/2}\beta}{\delta},$$

and

$$\left\| \mathbb{E}_{u_{t-2L+2:t-L+1}}[g_{t-L+1}] - \nabla \tilde{f}_{t-L+1}(x_t) \right\| \leq 2\eta_{t-L+1} \frac{dML^{5/2}\beta}{\delta} + \frac{d\delta L^2}{2}$$

□

**Lemma 9.** *Conditioned on  $u_1, \dots, u_{t-2L+1}$ ,*

$$\mathbb{E}_{u_{t-2L+2:t-L+1}} \left[ \left( \mathbb{E}_{u_{t-2L+2:t-L+1}}[g_{t-L+1}] - g_{t-L+1} \right)^\top (x_t - x_\delta^*) \right] \leq \eta_{t-L+1} \frac{2d^2 M^2 L^2}{\delta^2}.$$

*Proof.* For convenience, let  $\mathbb{E}$  denote the expectation over  $u_{t-2L+2:t-L+1}$ . Note that  $x_t$  is a function of  $g_{t-L+1}$ , which depends on  $u_{t-2L+2:t-L+1}$ . We have

$$\begin{aligned} \mathbb{E} \left[ \left( \mathbb{E}[g_{t-L+1}] - g_{t-L+1} \right)^\top (x_t - x_\delta^*) \right] &= \mathbb{E} \left[ \left( \mathbb{E}[g_{t-L+1}] - g_{t-L+1} \right)^\top (x_{t-L+1} - x_\delta^*) \right] \\ &\quad + \mathbb{E} \left[ \left( \mathbb{E}[g_{t-L+1}] - g_{t-L+1} \right)^\top (x_t - x_{t-L+1}) \right]. \end{aligned}$$

Note that  $x_{t-L+1}$  is fixed conditioned on  $u_1, \dots, u_{t-2L+1}$ , so

$$\mathbb{E} \left[ \left( \mathbb{E}[g_{t-L+1}] - g_{t-L+1} \right)^\top (x_{t-L+1} - x_\delta^*) \right] = 0.$$

The second term satisfies, for any sequence of  $u_{t-2L+2:t-L+1}$ ,

$$\begin{aligned} \left( \mathbb{E}[g_{t-L+1}] - g_{t-L+1} \right)^\top (x_t - x_{t-L+1}) &\leq \left\| \mathbb{E}[g_{t-L+1}] - g_{t-L+1} \right\| \|x_t - x_{t-L+1}\| \\ &\leq 2 \max_{u_{t-2L+1:t-L+1}} \|g_{t-L+1}\| \|x_t - x_{t-L+1}\| \\ &\leq \frac{2dM\sqrt{L}}{\delta} \|x_t - x_{t-L+1}\|. \end{aligned}$$

Similarly to the proof of Lemma 7, we have

$$\|x_t - x_{t-L+1}\| \leq \sum_{i=0}^{L-2} \|x_{t-i} - x_{t-i-1}\| \leq \sum_{i=0}^{L-2} \eta_{t-i-1} \|g_{t-i-L}\| \leq \eta_{t-L+1} \frac{dML^{3/2}}{\delta}.$$

Therefore,

$$\left( \mathbb{E}[g_{t-L+1}] - g_{t-L+1} \right)^\top (x_t - x_{t-L+1}) \leq \eta_{t-L+1} \frac{2d^2 M^2 L^2}{\delta^2},$$

and the lemma follows by summing the two terms. □

## E. Experiments

### E.1. Meta-optimization implementation

The implementation used for the deep learning experiments is the convex stochastic meta-optimization algorithm detailed in Algorithm 2 with the Gradient Perturbation Controller (GPC) as our algorithm  $\mathcal{A}$ . This has two key differences with what is used in our proofs: (1) we do not use the regularized loss functions of the form  $f(x) + \beta\|x - x_k\|^2$  and (2) we use the full GPC algorithm from nonstochastic control (Algorithm 3 in (Hazan & Singh, 2023), as opposed to the bandit version developed in (Gradu et al., 2020)). Note that in the full GPC algorithm, the controller minimizes a surrogate loss that is computed through counterfactual rollouts; therefore, the implemented algorithm must backpropagate through several

training steps in order to perform each meta-update. For more information on nonstochastic control and the counterfactual nature of the GPC algorithm, please see Chapter 7 of (Hazan & Singh, 2023). We learn scalar controller coefficients for computational efficiency (instead of full matrices, though we observed no difference when using diagonal matrices), and we use the Adam optimizer with learning rate  $10^{-4}$  and  $\beta_1 = 0.9, \beta_2 = 0.999$  within the GPC algorithm instead of gradient descent.

An open-source implementation of the algorithm is available at <https://anonymous.4open.science/r/meta-opt-8916/>. For clarity and reproducibility, we also provide a specification of our practical meta-optimization algorithm (in more standard deep learning terminology) as Algorithm 6. In the algorithm below, we set  $H = 32$  (except for WMT, where  $H = 16$  due to memory constraints),  $L = 2$ ,  $\mathcal{A}$  as Adam with learning rate  $10^{-3}$  and  $(\beta_1, \beta_2) = 0.9, 0.999$ , and the initializers selected at random; however, the behavior is quite robust to all these parameters. Note that the memory usage scales with  $H$ , and computation time scales with both  $H$  and  $L$ .

---

**Algorithm 6** Meta-optimization, deep learning implementation

---

**Require:** number of episodes  $N$ , number of steps per episode  $T$ , window size  $H$ , rollout length  $L$ , meta-optimizer  $\mathcal{A}$ , initializers  $\{x_{i,1}\}_{i=1}^N$ , initial learning rate  $\eta$ .

- 1: Initialize buffers of the past  $L + 1$  model parameters and data batches, and a buffer for the past  $H + L$  stochastic gradients.
  - 2: Initialize scalar controller parameters  $\{\eta_{i,t,h}\}_{h=0}^{H-1} \subset \mathbb{R}$  with  $\eta_{i,t,h} = \begin{cases} \eta & h = 0 \\ 0 & h \geq 1 \end{cases}$ .
  - 3: **for**  $i = 1, \dots, N$  **do**
  - 4:     **for**  $t = 1, \dots, T$  **do**
  - 5:         Play  $x_{i,t} = x_{i,t-1} - \sum_{h=0}^{H-1} \eta_{i,t,h} \tilde{\nabla} f_{i,t-h}$  if  $t > 1$ ; else play  $x_{i,1}$  (if  $t - h < 0$ , set  $\tilde{\nabla} f_{i,t-h} = 0$ ).
  - 6:         Receive a mini-batch of examples  $B_{i,t}$  of size  $b$ .
  - 7:         Obtain the loss function  $f_{i,t} = \frac{1}{b} \sum_{j \in B_{i,t}} f_{i,j}$ .
  - 8:         Suffer loss  $f_{i,t}(x_{i,t})$  and compute the stochastic gradient  $\tilde{\nabla} f_{i,t} = \frac{1}{b} \sum_{j \in B_{i,t}} \nabla f_{i,j}(x_{i,t})$ .
  - 9:         If  $t \geq H + L$ , compute the surrogate loss by rolling out for  $L$  training steps starting from  $x_{i,t-L}$  using the **past batches**  $\{B_{i,t-\ell}\}_{\ell=L}^1$  with the **current controller**  $\{\eta_{i,t,h}\}_{h=0}^{H-1}$  and the **past gradients**. Evaluate the loss on  $B_{i,t}$  for the parameters at the end of the rollout, and take the gradient with respect to  $(\eta_{i,t,h})_h$ .
  - 10:         Update the controller using this gradient via  $\mathcal{A}$ , and discard the parameters from the end of the rollout.
  - 11:         Append  $\tilde{\nabla} f_{i,t}$  to the gradient buffer,  $x_{i,t}$  to the parameter buffer, and  $B_{i,t}$  to the data buffer.
- 

## E.2. Experimental setup

**Architectures** We used the following commonplace deep learning architectures for the three workloads, and note that the deterministic setting uses the same batch of data throughout training:

- **MNIST:** a 3-layer multilayer perceptron (MLP) with ReLU and 784, 100, 100, and 10 neurons in the input layer, two hidden layers, and output layer, respectively, totaling 90K parameters. We used a batch size of 512 in both the deterministic and stochastic settings with no preprocessing. For the deterministic setting we trained for  $N = 16$  episodes of  $T = 500$  iterations each, and for the stochastic setting we train for 5,000 iterations.
- **CIFAR:** a VGG-16 architecture with an output layer of 10, totaling 15M parameters. We used a batch size of 512 in both the deterministic and stochastic settings with no preprocessing. For the deterministic setting we trained for  $N = 8$  episodes of  $T = 500$  iterations each, and for the stochastic setting we train for 9,000 iterations.
- **WMT:** a base Transformer architecture (as specified in (Vaswani et al., 2023) and implemented in Flax’s WMT tutorial) totaling 65M parameters. We evaluated on the WMT-14 English-to-German dataset, and we used a batch size of 16 in both the deterministic and stochastic settings. For the deterministic setting we trained for  $N = 12$  episodes of  $T = 8,000$  iterations each, and for the stochastic setting we train for 100,000 iterations.

**Baselines** For each of the above workloads, we tried the following deep learning optimizers:

- **SGD:** Gradient descent with weight decay. To tune this baseline, we used a grid search over the learning rate  $\eta$  and the weight decay parameter  $\delta$  taking values  $\eta \in [0.001, 0.01, 0.1, 0.2, 0.4, 1.0]$  and  $\delta \in [0, 10^{-5}, 10^{-4}, 10^{-3}]$ , respectively.

- **MOMENTUM**: Gradient descent with momentum and weight decay. To tune this baseline, we used a grid search over the learning rate  $\eta$ , the momentum parameter  $\mu$ , and the weight decay parameter  $\delta$  taking values  $\eta \in [0.001, 0.01, 0.1, 0.2, 0.4, 1.0]$ ,  $\mu \in [0.9, 0.95, 0.99]$ , and  $\delta \in [0, 10^{-5}, 10^{-4}, 10^{-3}]$ , respectively.
- **ADAMW**: Adam optimizer, with weight decay. To tune this baseline, we used a grid search and swept the learning rate  $\eta$ , momentum parameters  $\beta_1, \beta_2$ , and weight decay parameter  $\delta$  with the values  $\eta \in [10^{-4}, 4 \cdot 10^{-4}, 10^{-3}]$ ,  $\beta_1 \in [0.9, 0.99]$ ,  $\beta_2 \in [0.9, 0.99, 0.999]$ , and  $\delta \in [0, 10^{-5}, 10^{-4}, 10^{-3}]$ , respectively.
- **HGD**: Hypergradient descent acting on the standard gradient descent algorithm (Algorithm 4 in (Baydin et al., 2017)). To tune this baseline, we set the initial learning rate to be the tuned SGD learning rate and swept the meta-learning rate  $\beta$  with the values  $\beta \in [10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}]$ . Hypergradient descent never performed better than tuned vanilla SGD, so we do not plot it in Figures 1 or 2.
- **DOG**: The Distance-over-Gradients (DoG) algorithm (Ivgi et al., 2023). We run this baseline with the given hyperparameters since it is self-tuning, and we use the `optax.contrib` implementation.
- **D-ADAPTATION**: D-Adaptation algorithm acting on the Adam optimizer (Algorithm 5 in (Defazio & Mishchenko, 2023)). We run this baseline with the given hyperparameters since it is self-tuning, and we use the `optax.contrib` implementation.
- **MECHANIC**: the Mechanic (Cutkosky et al., 2024) algorithm acting on the AdamW optimizer described earlier. We tune this baseline with the same grid search used to tune the AdamW optimizer, and we use the `optax.contrib` implementation.

### E.3. Ablations & other experiments

**Sequential stability** One assumption we make that is nonstandard in the deep learning optimization literature is Assumption 2 – the sequential stability of the LTV dynamical system. We numerically verify this notion of stability in Figure 3 for the dynamical system induced by training a small neural network on MNIST (since the size of these matrices scales quadratically with number of parameters, computing this for larger networks is infeasible).

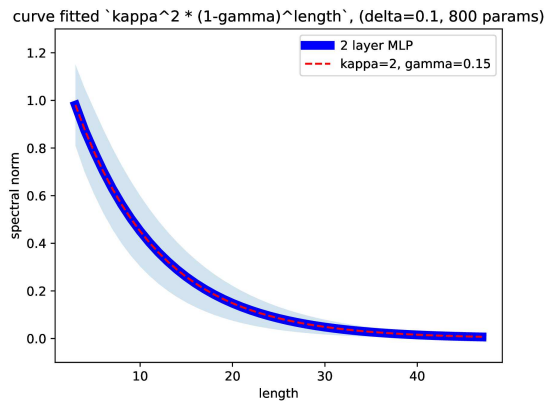


Figure 3. Decay of spectral norm of  $\|\prod_{t=s}^r A_t\|$  as a function of  $|r - s|$  for a small neural network at the beginning of training. Averaged over 10 trials. Assumption 2 is satisfied in this instance with a value of  $\kappa \approx 2.0$ .

**Stochastic meta-optimization** In Figure 4, we show what may happen if the meta-optimization algorithm is run in the stochastic setting; as can be seen, the performance degrades between episodes. This occurs on the more complex datasets, and so we believe it to be a characteristic of how backpropagation through rollouts responds to noise between batches.

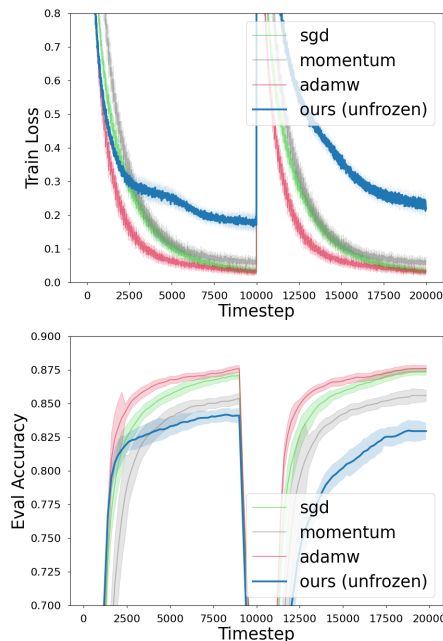


Figure 4. Behavior of the meta-optimization algorithm when training unfrozen in the stochastic regime. The optimizer’s performance degrades over time.

### F. Future work

In this work, we presented an initial exploration into the behavior of our meta-optimization algorithm in deep learning environments. As such, our investigations and design decisions leave much room for improvement and discovery, and we hope that the promising results inspire research to make such methods more practical. We list below several directions of investigation that we think will be fruitful.

**Generalization** As seen in Figure 2, on the WMT workload there is a noticeable generalization gap between adaptive methods (algorithms like Adam and its derivatives) and non-adaptive methods (vanilla gradient descent, momentum, and fixed preconditioners). This mirrors what is seen in many related works, where versions of hyperparameter-tuning algorithms that are built on top of Adam variants perform better than those built on vanilla gradient descent. While training with adaptive methods does not induce a linear dynamical system, we consider it a problem of practical importance to incorporate this adaptivity into the meta-optimization algorithm.

**Scaling** We have demonstrated that the meta-optimization approach is competitive on workloads of different scales. However, for the largest workloads, it would be valuable to understand the transferability of learned optimizers across model scales. Progress in this direction could allow for one to learn a controller on a smaller architecture and transfer it to a larger one, potentially allowing for meta-optimization of large frontier models.

**Optimizer pre-training** The paradigm we proposed for meta-optimization on general deep learning workloads was to learn an optimizer in the deterministic setting on a fixed batch and deploy it in the stochastic setting. However, we are still investigating the effects of the data selection itself on the downstream performance: how do batch size, dataset complexity, and using the same frozen batch affect the optimal controller, and does this impact its transferability to the stochastic setting? Depending on the answers to these questions, there may be more principled or practical ways to learn a robust optimizer that can be deployed in the stochastic minibatch setting.

**Efficient implementation** There is much room for improvement in terms of the implementation and parallelization of the meta-optimization algorithm. At the moment, the optimizer state needs to retain the past  $H$  gradients, which for large models can be a significant memory burden; however, there is ample structure in the gradient buffer and how it is used, and so we expect that an efficient sharding of optimizer state is possible. Furthermore, we anticipate that there are cleverer ways to use Jax’s machinery in order to help with the computational cost of backpropagating through rollouts.