
Think out of the “Box”: Generically-Constrained Asynchronous Composite Optimization and Hedging

Pooria Joulani*
DeepMind, UK
pjoulani@google.com

András György
DeepMind, UK
agyorgy@google.com

Csaba Szepesvári
DeepMind, UK
szepi@google.com

Abstract

We present two new algorithms, ASYNCADA and HEDGEHOG, for asynchronous sparse online and stochastic optimization. ASYNCADA is, to our knowledge, the first asynchronous stochastic optimization algorithm with finite-time data-dependent convergence guarantees for generic convex constraints. In addition, ASYNCADA: (a) allows for proximal (i.e., composite-objective) updates and adaptive step-sizes; (b) enjoys any-time convergence guarantees without requiring an exact global clock; and (c) when the data is sufficiently sparse, its convergence rate for (non-)smooth, (non-)strongly-convex, and even a limited class of non-convex objectives matches the corresponding serial rate, implying a theoretical “linear speed-up”. The second algorithm, HEDGEHOG, is an asynchronous parallel version of the Exponentiated Gradient (EG) algorithm for optimization over the probability simplex (a.k.a. Hedge in online learning), and, to our knowledge, the first asynchronous algorithm enjoying linear speed-ups under sparsity with non-SGD-style updates. Unlike previous work, ASYNCADA and HEDGEHOG and their convergence and speed-up analyses are not limited to individual coordinate-wise (i.e., “box-shaped”) constraints or smooth and strongly-convex objectives. Underlying both results is a generic analysis framework that is of independent interest, and further applicable to distributed and delayed feedback optimization.

1 Introduction

Many modern machine learning methods are based on iteratively optimizing a regularized objective. Given a convex, non-empty set of feasible model parameters $\mathcal{X} \subset \mathbb{R}^d$, a differentiable *loss* function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, and a convex (possibly non-differentiable) *regularizer* function $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$, these methods seek the parameter vector $x^* \in \mathcal{X}$ that minimizes $f + \phi$ (assuming a minimizer exists):

$$x^* = \arg \min_{x \in \mathcal{X}} f(x) + \phi(x). \quad (1)$$

In particular, empirical risk minimization (ERM) methods such as (regularized) least-squares, logistic regression, LASSO, and support vector machines solve optimization problems of the form (1). In these cases, $f(x) = \frac{1}{m} \sum_{i=1}^m F(x, \xi_i)$ is the average of the loss $F(x, \xi_i)$ of the model parameter x on the given training data $\xi_1, \xi_2, \dots, \xi_m$ and $\phi(x)$ is a norm (or a combination of norms) on \mathbb{R}^d (e.g., $F(x, \xi) = \log(1 + \exp(x^\top \xi))$ and $\phi(x) = \frac{1}{2} \|x\|_2^2$ in linear logistic regression [13]).

To bring the power of modern parallel computing architectures to such optimization problems, several papers in the past decade have studied parallel variants of the stochastic optimization algorithms applied to these problems. Here one of the main questions is to quantify the cost of parallelization, that is, how much extra work is needed by a parallel algorithm to achieve the same accuracy as its serial variant. Ideally, a parallel algorithm is required to do no more work than the serial version, but

*Work partially done when the author was at the University of Alberta, Edmonton, AB, Canada.

this is very hard to achieve in our case. Instead, a somewhat weaker goal is to ensure that the price of parallelism is at most a constant factor: that is, the parallel variant needs at most constant-times more updates (or work). In other words, using τ parallel process requires a wall-clock running time that is only $O(1/\tau)$ -times that of the serial variant. In this case we say that the parallel algorithm achieves a *linear speed-up*. Of particular interest are asynchronous lock-free algorithms, where Recht et al. [30] demonstrated first that linear speed-ups are possible: They showed that if τ processes run stochastic gradient descent (SGD) and apply their updates to the same shared iterate without locking, then the overall algorithm (called Hogwild!) converges after the same amount of work as serial SGD, up to a multiplicative factor that increases with the number of concurrent processes and decreases with the sparsity of the problem. Thus, if the problem is sparse enough, this penalty can be considered a constant, and the algorithm achieves linear speed-up. Several follow-up work (see e.g., [20, 18, 17, 27, 24, 10, 29, 7, 11, 4, 2, 3, 19, 31, 33, 32, 35, 36, 12, 6, 28] and the references therein) have demonstrated linear speed-ups for methods based on (block-)coordinate descent (BCD), as well as other variants of SGD such as SVRG [15], SAGA [8], ADAGRAD [22, 9], and SGD with a time-decaying step-size. Despite the great advances, however, several problems remain open.²

First, the existing convergence guarantees concern SGD when the constraint set \mathcal{X} is box-shaped, that is, a Cartesian product of (block-)coordinatewise constraints $\mathcal{X} = \times_{i=1}^d \mathcal{X}_i$. This leaves it unclear whether existing techniques apply to stochastic optimization algorithms that operate on non-box-shaped constraints (e.g., on the ℓ_2 ball), or algorithms that use a non-Euclidean regularizer, such as the exponentiated gradient (EG) algorithm used on the probability simplex (see, e.g., [34, 14]).

Second, with the exception of the works of Duchi et al. [10] and Pan et al. [26] (which still require box-shaped constraints), and De Sa et al. [7] (which only bounds the probability of “failure”, i.e., of producing no iterates in the ϵ -ball around x^*), the existing analyses demonstrating linear speed-ups are limited to strongly-convex (or Polyak-Łojasiewicz) objectives. Thus, so far it has remained unclear whether a similar speed-up analysis is possible if the objective is simply convex or smooth [20], or if we are in the closely-related online-learning setting with the objective changing over time.

Third, with the exception of the work of Pedregosa et al. [27] (which still requires box-shaped constraints, block-separable ϕ and strongly-convex f), the existing analyses do not take advantage of the structure of problem (1). In particular, when ϕ is “simple to optimize” over \mathcal{X} (formally defined as having access to a proximal operator oracle, as we make precise in what follows), serial algorithms such as Proximal-SGD take advantage of this property to achieve considerably faster convergence rates. Asynchronous variants of the Proximal-SGD algorithm with such faster rates have so far been unavailable for non-strongly-convex objectives and non-box constraints.

1.1 Contributions

In this paper we address the aforementioned problems and present algorithms that are applicable to general convex constraint sets, not just box-shaped \mathcal{X} , but still achieve linear speed-ups (under sparsity) for non-smooth and non-strongly-convex (as well as smooth or strongly convex) objectives, and even for a specific class of non-convex problems. This is achieved through our new asynchronous optimization algorithm, ASYNCADA, which generalizes the ASYNC-ADAGRAD (and ASYNC-DA) algorithm of Duchi et al. [10] to proximal updates and its data-dependent bound to arbitrary constraint sets. Instantiations of ASYNCADA under different settings are given in Table 1. Indeed, the results are obtained by a more general analysis framework, built on the work of Duchi et al. [10], that yields data-dependent convergence guarantees for a generic class of adaptive, composite-objective online optimization algorithms undergoing perturbations to their “state”. We further use this framework to derive the first asynchronous online and stochastic optimization algorithm with non-box constraints that uses non-Euclidean regularizers. In particular, we present and analyze HEDGEHOG, the parallel asynchronous variant of the EG algorithm, also known as Hedge in online linear optimization [34, 14],

² In this paper, we do not further consider BCD-based methods, for two main reasons: a) in general, a BCD update may unnecessarily slow down the convergence of the algorithm by focusing only on a single coordinate of the gradient information, especially in the sparse-data problems we consider in this paper (see, e.g., Pedregosa et al. [27, Appendix F]); and b) BCD algorithms typically apply only to box-shaped constraints, which is what our algorithms are designed to be able to avoid. We would like to note, however, that our stochastic gradient oracle set-up (Section 2) does allow for building an unbiased gradient estimate using only one randomly-selected (block-)coordinate, as done in BCD methods. Nevertheless, the literature on parallel asynchronous BCD algorithms is vast, including especially algorithms for proximal, non-strongly-convex, and non-convex optimization; see, e.g., [29, 11, 4, 2, 3, 19, 31, 33, 32, 35, 36, 12, 6, 28] and the references therein.

Algorithm	\mathcal{X}	Nonsmooth	Smooth f	Strongly-convex	Smooth f + Strongly-convex
SGD (DA)	\mathbb{R}^d	[10, 26] ✓	[26] ✓	[26] ✓	[30, 7, 20, 17, 24, 26] ✓
SGD (MD)	\square	[10, 26]	[26]	[26]	[30, 7, 20, 17, 24, 26]
DA	\circ	✓	✓	✓	✓
AG / DA	\square	[10, 26] ✓	[26] ✓	[26] ✓	[26] ✓
AG / DA	\circ	✓	✓	✓	✓
Prox-MD	\square	-	-	-	[27]
Prox-DA	\circ	✓	✓	✓	✓
Prox-AG	\circ	✓	✓	✓	✓
Hedge/EG	\triangle	✓	✓	✓	✓

Table 1: (Star-)convex optimization settings under which sufficient sparsity results in linear speed-up. Previous work are cited under the settings they address. A ✓ indicates a setting covered by the results in this paper. The symbols \square , \triangle , and \circ indicate, respectively, the case when the constraint set is box-shaped, the probability simplex, or any convex constraint set with a projection oracle. AG, DA, and MD stand, respectively, for ADAGRAD, Dual-Averaging, and Mirror Descent, while Prox-AG, Prox-DA, and Prox-MD denote their proximal variants (using the proximal operator of ϕ).

and show that it enjoys similar parallel speed-up regimes as ASYNCADA. The results are derived for the more general setting of noisy online optimization, and the generic framework is of independent interest, in particular in the related settings of distributed and delayed-feedback learning.

The rest of the paper is organized as follows: The optimization problem and its solution with serial algorithms are described in Section 2 and Section 3, respectively. The generic perturbed-iterate framework is given in Section 4. Our main algorithms, ASYNCADA and HEDGEHOG are presented and analyzed in Section 5 and Section 6, respectively. Conclusions are drawn and some open problems are discussed in Section 7, while omitted technical details are given in the appendices.

1.2 Notation and definitions

We use $[n]$ to denote the set $\{1, 2, \dots, n\}$, $\mathbb{I}\{\mathcal{E}\}$ for the indicator of an event \mathcal{E} , and $\sigma(\mathcal{H})$ to denote the sigma-field generated by a set \mathcal{H} of random variables. The j -th coordinate of a vector $a \in \mathbb{R}^d$ is denoted $a^{(j)}$. For $\alpha \in \mathbb{R}^d$ with positive entries, $\|\cdot\|_\alpha$ denotes the α -weighted Euclidean norm, given by $\|x\|_\alpha^2 = \frac{1}{2} \sum_{j=1}^d \alpha^{(j)} (x^{(j)})^2$, and $\|\cdot\|_{\alpha,*}$ its dual. We use $(a_t)_{t=i}^j$ to denote a sequence a_i, a_{i+1}, \dots, a_j and define $a_{i:j} := \sum_{t=i}^j a_t$, with $a_{i:j} := 0$ if $i > j$. Given a differentiable function $h : \mathbb{R}^d \rightarrow \mathbb{R}$, the *Bregman divergence* of $y \in \mathbb{R}^d$ from $x \in \mathbb{R}^d$ with respect to (w.r.t.) h is given by $\mathcal{B}_h(y, x) := h(y) - h(x) - \langle \nabla h(x), y - x \rangle$. It can be shown that a differentiable function is convex if and only if $\mathcal{B}_h(x, y) \geq 0$ for all $x, y \in \mathbb{R}^d$. The function $h : \mathbb{R}^d \rightarrow \mathbb{R}$ is μ -strongly convex w.r.t. a norm $\|\cdot\|$ on \mathbb{R}^d if and only if for all $x, y \in \mathbb{R}^d$ $\mathcal{B}_h(x, y) \geq \frac{\mu}{2} \|x - y\|^2$, and smooth w.r.t. a norm $\|\cdot\|$ if and only if for all $x, y \in \mathbb{R}^d$, $|\mathcal{B}_h(x, y)| \leq \frac{1}{2} \|x - y\|^2$. A differentiable function f is *star-convex* if and only if there exists a global minimizer x^* of f such that for all $x \in \mathbb{R}^d$, $\mathcal{B}_f(x^*, x) \geq 0$.

2 Problem setting: noisy online optimization

We consider a generic iterative optimization setting that enables us to study both online learning and stochastic composite optimization. The problem is defined by a (known) constraint set \mathcal{X} and a (known) convex (possibly non-differentiable) function ϕ , as well as differentiable functions f_1, f_2, \dots about which an algorithm learns iteratively. At each iteration $t = 1, 2, \dots$, the algorithm picks an iterate $x_t \in \mathcal{X}$, and observes an unbiased estimate $g_t \in \mathbb{R}^d$ of the gradient $\nabla f_t(x_t)$, that is, $\mathbb{E}\{g_t | x_t\} = \nabla f_t(x_t)$. The goal is to minimize the composite-objective online regret after T iterations, given by

$$R_T^{(f+\phi)} = \sum_{t=1}^T (f_t(x_t) + \phi(x_t) - f_t(x_T^*) - \phi(x_T^*)),$$

where $x_T^* = \arg \min_{x \in \mathcal{X}} \left\{ \sum_{t=1}^T (f_t(x) + \phi(x)) \right\}$. In the absence of noise (i.e., when $g_t = \nabla f_t(x_t)$), this reduces to the (composite-objective) online (convex) optimization setting [34, 14].

Stochastic optimization, online regret, and iterate averaging. If $f_t = f$ for all $t = 1, 2, \dots$, we recover the stochastic optimization setting, with the algorithm aiming to minimize the composite objective $f + \phi$ over \mathcal{X} while receiving noisy estimates of ∇f at points $(x_t)_{t=1}^T$. The algorithm's online regret can then be used to control the optimization risk: Since $f_t \equiv f$, we have $x_T^* = x^* = \arg \min_{x \in \mathcal{X}} \{f(x) + \phi(x)\}$, and by Jensen's inequality, if f is convex and $\bar{x}_T = \frac{1}{T}x_{1:T}$ is the average iterate,

$$f(\bar{x}_T) + \phi(\bar{x}_T) - f(x^*) - \phi(x^*) \leq \frac{1}{T}R_T^{(f+\phi)}.$$

In addition, if f is non-convex but \bar{x}_T is selected uniformly at random from x_1, \dots, x_T , then the above bound holds in expectation. As such, in the rest of the paper we study the optimization risk through the lens of online regret.

Stochastic first-order oracle. Throughout the paper, we assume that at time t , the noisy gradient estimate g_t is given by a randomized first-order oracle³ $g_t : \mathbb{R}^d \times \Xi \rightarrow \mathbb{R}^d$, where Ξ is some space of random variables, and there exists a sequence $(\xi_t)_{t=1}^T$ of independent elements from Ξ , with distribution \mathbb{P}_Ξ , such that $\int_\Xi g_t(x, \xi) d\mathbb{P}_\Xi(\xi) = \nabla f_t(x)$ for all $x \in \mathcal{X}$.

For example, in the finite-sum stochastic optimization case when $f = \sum_i f_i$, selecting one f_i uniformly at random to estimate the gradient corresponds to \mathbb{P}_Ξ being the uniform distribution on $\Xi = \{1, 2, \dots, N\}$ and $g_t(x, \xi_t) = \nabla f_{\xi_t}(x)$, whereas selecting a mini-batch of f_i 's corresponds to Ξ being the set of subsets (of a fixed or varying size) of $\{1, 2, \dots, N\}$ and $g_t(x, \xi_t) = \frac{1}{|\xi_t|} \sum_{i \in \xi_t} \nabla f_i(x)$. This also covers variance-reduced gradient estimates as formed, e.g., by SAGA and SVRG, in which case g_t is built using information from the previous rounds.⁴

3 Preliminaries: analysis in the serial setting

First, we recall the analysis of a generic *serial* dual-averaging algorithm, known as Adaptive Follow-the-Regularized-Leader (ADA-FTRL) [21, 25, 16], that generalizes regularized dual-averaging [37] and captures the dual-averaging variants of SGD, Ada-Grad, Proximal-SGD and EG as special case.

Serial ADA-FTRL. The serial ADA-FTRL algorithm uses a sequence of regularizer functions r_0, r_1, r_2, \dots . At time $t = 1, 2, \dots$, given the previous feedback $g_s \in \mathbb{R}^d$, $s \in [t-1]$, ADA-FTRL selects the next point x_t such that

$$x_t \in \arg \min_{x \in \mathcal{X}} \langle z_{t-1}, x \rangle + t\phi(x) + r_{0:t-1}(x), \quad (2)$$

where $z_{t-1} = g_{1:t-1}$ is the sum of the past feedback. We refer to $(z_t, t, r_{0:t})$ as the *state* of the algorithm at time t , noting that apart from tie-breaking in (2), this state determines x_t .

It is straightforward to verify that with $\phi = 0$, $\mathcal{X} = \mathbb{R}^d$, and $r_{0:t-1} = \frac{\eta}{2} \|\cdot\|^2$ for some $\eta > 0$, we get the SGD update $x_t = -\frac{1}{\eta}g_{1:t-1}$. In addition, using $r_{0:t-1} = \frac{1}{2} \|\cdot\|_{\eta_t}^2$ where $\eta_t^{(i)}, i \in [d]$ are positive step-sizes (possibly adaptively tuned [22, 9]), ADA-FTRL reduces to $x_t = \mathbf{prox}(t\phi, -z_{t-1}, \eta_t)$, where \mathbf{prox} is the generalized *proximal operator* oracle⁵ over \mathcal{X} that, given a function ψ and vectors z and η , returns⁶

$$\mathbf{prox}(\psi, z, \eta) := \arg \min_{x \in \mathcal{X}} \psi(x) + \frac{1}{2} \|x - \eta^{-1} \odot z\|_\eta^2. \quad (3)$$

³With a slight abuse of notation, $g_t(x, \xi)$ (with arguments x, ξ) is from now on used to denote the oracle at time t evaluated at x, ξ , where as g_t (without arguments) denotes the observed noisy gradient $g_t(x_t, \xi_t)$.

⁴Note that in this case ξ_t remains an independent sequence, even though g_t changes with the history.

⁵Serial proximal DA [37] and ADA-FTRL call \mathbf{prox} with $\psi \leftarrow t\phi$, whereas the conventional Proximal-SGD algorithm (based on Mirror-Descent) invokes the proximal operator with $\psi \leftarrow \phi$ irrespective of the iteration; see the paper of Xiao [37, Sections 5 and 6] for a detailed discussion of this phenomenon.

⁶Here η^{-1} denotes the elementwise inverse of η and \odot denotes elementwise multiplication.

When η is the same for all coordinates (in which case we simply treat it as a scalar), this reduces to $\text{prox}(\psi, z, \eta) = \arg \min_{x \in \mathcal{X}} \psi(x) + \frac{\eta}{2} \|x - z/\eta\|^2$, which is the standard proximal operator; the generalized version (3) makes it possible to use coordinatewise step-sizes as in ADAGRAD [22, 9]. Finally, when $\phi = 0$ and \mathcal{X} is the probability simplex, ADA-FTRL with the negentropy regularizer $r_{0:t-1}(x) = r_0(x) = \eta \sum_{i=1}^d x_i \log(x_i)$ for some $\eta > 0$, recovers the update $x_t^{(i)} = C_t \exp(-z_{t-1}^{(i)}/\eta)$ of the EG algorithm, where $C_t = 1/\sum_{j=1}^d \exp(-z_{t-1}^{(j)}/\eta)$ is the constant normalizing x_t to lie in \mathcal{X} . Other choices of r_t recover algorithms such as the p -norm update; we refer to Shalev-Shwartz [34], Hazan [14], McMahan [21], and Orabona et al. [25] for further examples.

Analysis of ADA-FTRL ADA-FTRL and its special cases have been extensively studied in the literature [5, 34, 14, 21, 25, 16]. In particular, it has been shown that under specific conditions on r_t and ϕ , which we discuss in detail in Appendix F, ADA-FTRL enjoys the following bound on the linearized regret [25, 16]:

Theorem 1 (Regret of ADA-FTRL). *For any $x^* \in \mathcal{X}$ and any sequence of vectors $(g_t)_{t=1}^T$ in \mathbb{R}^d , using any sequence of regularizers r_0, r_1, \dots, r_T that are admissible w.r.t. a sequence of norms $\|\cdot\|_{(t)}$ (see Definition 2 in Appendix F), the iterates $(x_t)_{t=1}^T$ generated by ADA-FTRL satisfy*

$$\sum_{t=1}^T (\langle g_t, x_t - x^* \rangle + \phi(x_t) - \phi(x^*)) \leq r_{0:T}(x^*) - \sum_{t=0}^T r_t(x_{t+1}) + \sum_{t=1}^T \frac{1}{2} \|g_t\|_{(t,*)}^2. \quad (4)$$

Importantly, this bound holds for *any* feedback sequence g_t irrespective of the way it is generated, and serves as a solid basis to derive bounds under different assumptions on f , ϕ , and r_t [25, 16].

4 Relaxing the serial analysis: algorithms with perturbed state

In this section, we show that Theorem 1 can be used to analyze ADA-FTRL when its state undergoes specific perturbations. This relaxation of the generic serial analysis framework underlies our analysis of parallel asynchronous algorithms, since parallel algorithms like ASYNCADA and HEDGEHOG can be viewed as *serial* ADA-FTRL algorithms with perturbed states, as we show in Sections 5 and 6.

Perturbed ADA-FTRL. Next, we show that Theorem 1 also provides the basis to analyze ADA-FTRL with perturbed states. Specifically, suppose that instead of (2), the iterate x_t is given by

$$x_t \in \arg \min_{x \in \mathcal{X}} \langle \hat{z}_{t-1}, x \rangle + \hat{t}_t \phi(x) + \hat{r}_{0:t-1}(x), \quad t = 1, 2, \dots, \quad (5)$$

where \hat{z}_{t-1} denotes a *perturbed* version of the dual vector z_{t-1} , \hat{t}_t denotes a perturbed version of ADA-FTRL's iteration counter t , and $\hat{r}_{0:t-1}$ denotes a perturbed version of the regularizer $r_{0:t-1}$. Then, we can analyze the regret of the Perturbed-ADA-FTRL update (5) by comparing x_t to the “ideal” iterate \tilde{x}_t , given by

$$\tilde{x}_t := \arg \min_{x \in \mathcal{X}} \langle z_{t-1}, x \rangle + t\phi(x) + r_{0:t-1}(x), \quad t = 1, 2, \dots \quad (6)$$

Since $(\tilde{x}_t)_{t=1}^T$ is given by a non-perturbed ADA-FTRL update, it enjoys the bound of Theorem 1. The crucial observation of Duchi et al. [10] (who studied the special case of (5) with $\phi = 0$, box-shaped \mathcal{X} , and $\hat{r}_t = r_t$) was that the regret of Perturbed-ADA-FTRL is related to the linearized regret of \tilde{x}_t . When ϕ may be non-zero, we capture this relation by the next lemma, proved in Appendix A:

Lemma 1 (Perturbation penalty of ADA-FTRL). *Consider any sequences $(x_t)_{t=1}^T$ and $(\tilde{x}_t)_{t=1}^T$ in \mathcal{X} , and any sequence $(g_t)_{t=1}^T$ in \mathbb{R}^d . Then, the regret $R_T^{(f+\phi)}$ of the sequence $(x_t)_{t=1}^T$ satisfies*

$$R_T^{(f+\phi)} = \sum_{t=1}^T (\langle g_t, \tilde{x}_t - x^* \rangle + \phi(\tilde{x}_t) - \phi(x^*)) + \tilde{\epsilon}_{1:T} + \delta_{1:T} - B_{1:T}, \quad (7)$$

where $\tilde{\epsilon}_t = \langle g_t, x_t - \tilde{x}_t \rangle + \phi(x_t) - \phi(\tilde{x}_t)$, $\delta_t = \langle \nabla f_t(x_t) - g_t, x_t - x^* \rangle$ and $B_t = \mathcal{B}_{f_t}(x^*, x_t)$.

Since g_t is an unbiased estimate of $\nabla f_t(x_t)$ (conditionally given x_t), $\delta_{1:T}$ is zero in expectation 0, and for \tilde{x}_t given by (6), the first summation is bounded by Theorem 1. Also note that when the f_t are (star-)convex, $-B_{1:T} \leq 0$. Thus, to bound the regret of Perturbed-ADA-FTRL, it only

remains to control the ‘‘perturbation penalty’’ terms $\tilde{\epsilon}_t$ capturing the difference in the composite linear loss $\langle g_t, \cdot \rangle + \phi$ between x_t and \tilde{x}_t . In Appendix A, we use the stability of ADA-FTRL algorithms (Lemma 3) to control $\tilde{\epsilon}_{1:T}$, under a specific perturbation structure (coming from delayed updates to \hat{z}_t) that captures the evolution of the state of asynchronous dual-averaging algorithms like ASYNCADA and HEDGEHOG. Unlike Duchi et al. [10], our derivation applies to *any* convex constraint set \mathcal{X} and, crucially, to ADA-FTRL updates incorporating non-zero ϕ and a perturbed counter \hat{t}_t . The following (informal) theorem, whose formal version is given in Appendix A, captures the result.

Theorem 4 (informal). Under appropriate independence, regularity, and structural assumptions on the regularizers and the perturbations, the Perturbed-ADA-FTRL update (5) satisfies

$$\mathbb{E} \left\{ R_T^{(f+\phi)} \right\} \leq \mathbb{E} \left\{ r_{0:T}(x^*) + \sum_{t=1}^T \left(\frac{1 + p_* \nu_t + \sum_{s:t \in O_s} \frac{\tau_s}{\nu_s}}{2} \|g_t\|_{(t,*)}^2 + \frac{\Delta_t}{\nu_t} \right) - B_{1:T} \right\},$$

where p_* , ν_t , τ_t and Δ_t measure, respectively, the sparsity of the gradient estimates g_t , the difference $\hat{t}_t - t$, and the amount of perturbations in \hat{z}_{t-1} , and $\hat{r}_{0:t-1}$, while O_s is the set of time steps whose attributed perturbations affect iteration s (i.e., their updates are delayed beyond s).

As we show next, we can control the effect of p_* , τ_t and Δ_t in the bound by appropriately tuning \hat{t}_t , resulting in linear speed-ups for ASYNCADA and HEDGEHOG.

5 ASYNCADA: Asynchronous Composite Adaptive Dual Averaging

In this section, we introduce and analyze ASYNCADA for asynchronous noisy online optimization. ASYNCADA consists of τ processes running in parallel (e.g., threads on the same physical machine or computing nodes distributed over a network accessing a shared data store). The processes can access a shared memory, consisting of a *dual* vector $z \in \mathbb{R}^d$ to store the sum of observed gradient estimates g_t , a *step-size* vector $\eta \in \mathbb{R}^d$, and an integer t , referred to as the *clock*, to track the number of iterations completed at each point in time. The processes run copies of Algorithm 1 concurrently.

Algorithm 1: ASYNCADA: Asynchronous Composite Adaptive Dual Averaging

```

1 repeat
2    $\hat{\eta} \leftarrow$  a full (lock-free) read of the shared step-sizes  $\eta$ 
3    $\hat{z} \leftarrow$  a full (lock-free) read of the shared dual vector  $z$ 
4    $t \leftarrow t + 1$  // atomic read-increment
5    $\hat{t} \leftarrow t + \gamma$  // denote  $\hat{z}_{t-1} = \hat{z}$ ,  $\hat{\eta}_t = \hat{\eta}$ ,  $\hat{t}_t = \hat{t}$ 
6   Receive  $\xi_t$ 
7   Compute the next iterate:  $x_t \leftarrow \text{prox}(\hat{t}_t \phi, -\hat{z}_{t-1}, \hat{\eta}_t)$  // prox defined in (3)
8   Obtain the noisy gradient estimate:  $g_t \leftarrow g_t(x_t, \xi_t)$ 
9   for  $j$  such that  $g_t^{(j)} \neq 0$  do  $z^{(j)} \leftarrow z^{(j)} + g_t^{(j)}$  // atomic update
10  Update the shared step-size vector  $\eta$ 
11 until terminated

```

Inconsistent reads. The processes access the shared memory without necessarily acquiring a lock: as in previous Hogwild!-style algorithms [30, 20, 18, 17, 27], we only assume that operations on single coordinates of z and η , as well as on t' , are atomic. This in particular means that the values of \hat{z} or $\hat{\eta}$ read by a process may not correspond to an actual state of z or η at any given point in time, as different processes can modify the coordinates in parallel while the read is taking place. A process π is in write-conflict with another process π' (equivalently, π' is in read-conflict with π) if π' reads parts of the memory which should have been updated by π before. To limit the effects of asynchrony, we assume that a process can be in write- and read conflicts with at most $\tau_c - 1$ processes, respectively.

The role of γ . ASYNCADA uses an over-estimate \hat{t}_t of the current global clock t by an additional γ . This over-estimation enables us to better handle the effect of asynchrony when composite objectives are involved, in particular ensuring the appropriate tuning of ν_t in Theorem 4; see Appendix C. ASYNCADA can nevertheless be run without γ (i.e., with $\gamma = 0$).⁷

⁷ In Theorems 2, 5 and 6, we set γ based on $\tau_* := \max\{\tau_c, \tau\}$. The analysis is still possible, and straightforward, with $\gamma = 0$, but results in a worst constant factor in the rate, as well as an extra additive term of order $\mathcal{O}(\tau_*^2 \Phi)$ where $\Phi = \sup_{x, y \in \mathcal{X}} \{\phi(x) - \phi(y)\}$ is the diameter of \mathcal{X} w.r.t. ϕ . This term does not diminish with p_* and may be unnecessarily large, affecting convergence in early stages of the optimization process.

Exact vs estimated clock. ASYNCADA as given in Algorithm 1 maintains the exact global clock t . However, this option may not be desirable (or available) in certain asynchronous computing scenarios. For example, if the processes are distributed over a network, then maintaining an exact global clock amounts to changing the pattern of asynchrony and delaying the computations by repeated calls over a network. To mitigate this requirement, in Appendix B we provide ASYNCADA(ρ), a version of ASYNCADA in which the processes update the global clock only every ρ iterations. ASYNCADA as presented in Algorithm 1 is equivalent to ASYNCADA(ρ) with $\rho = 1$, and both algorithms enjoy the same rate of convergence and linear speed-up. Obviously, when $\phi \equiv 0$ and t is not used for setting the step-sizes η either, there is no need to maintain t physically, and Line 4 can be omitted.

Updating the step-sizes η : In Line 10 of Algorithm 1, the step-size η has to be updated based on the information received. The exact way this is done depends on the specific step-size schedule. In particular, we consider two situations: First, when the step-size is either constant or a simple function of t (or \hat{t}_t in case of ASYNCADA(ρ)), and second, when diagonal ADA-GRAD step-sizes are used. In the first case, the vector η need not be kept in the shared memory explicitly, and Lines 2 and 10 can be omitted. In the second case, following [10], we store the sum of squared gradients in the shared η , i.e., Line 10 is implemented as follows:

```
10* for  $j$  such that  $g_t^{(j)} \neq 0$  do  $(\eta^{(j)})^2 \leftarrow (\eta^{(j)})^2 + \alpha^2 (g_t^{(j)})^2$  // atomic update
```

for a fixed hyper-parameter $\alpha > 0$. In this case, we are storing the square of η in the shared memory, so a square root operation needs to be applied after reading the shared memory in Line 2 to retrieve η .

Forming the output \bar{x}_T for stochastic optimization: For stochastic optimization, the algorithm needs to output the average (or randomized) iterate \bar{x}_T at the end. However, this needs no further coordination between the processes. To form the average iterate, it suffices for each process to keep a local running sum of the iterates it produces and the number of updates it makes. At the end, \bar{x}_T is built from these sums and the total number of updates. Alternatively, we can return a random iterate as \bar{x}_T by terminating the algorithm, with probability $1/T$, after calculating x in Line 7.

5.1 Analysis of ASYNCADA

The analysis of ASYNCADA is based on treating it as a special case of Perturbed-ADA-FTRL. In order to be able to use Theorem 4, we start with the following independence assumption on ξ_t :

Assumption 1 (Independence of ξ_t). For all $t = 1, 2, \dots, T$, the t -th sample ξ_t is independent of the history $\hat{\mathcal{H}}_t := \{(\xi_s, \hat{z}_s, \hat{\eta}_{s+1})_{s=1}^{t-1}\}$.

This, in turn, implies that ξ_t is independent of x_t as well as x_s and ξ_s for all $s < t$.

For general (non-box-shaped) \mathcal{X} , Assumption 1 is plausible, as ASYNCADA *needs* to read z (and η) completely and independently of ξ_t . If \mathcal{X} is box-shaped and ϕ is coordinate-separable, however, the values of $x_t^{(j)}$ for different coordinates j can be calculated independently. In this case, the algorithm may first sample ξ_t , and then only read the relevant coordinates j from z (and η) for which g_t may be non-zero, as calculating other values of $x_t^{(j)}$ is unnecessary for calculating g_t . As mentioned by Mania et al. [20], this violates Assumption 1. This is because multiple other processes are updating z and η , and the updates that are included the value read for \hat{z}_{t-1} (and $\hat{\eta}_t$) would then depend on ξ_t . Previous papers either assume that this independence holds in their analysis, e.g., by enforcing a full read of z and η , [20, 18, 17, 27], or rely on the smoothness of the objective to bound the effect of the possible change in the read values [20, Appendix A]. It seems possible to adapt the argument of Mania et al. [20, Appendix A] to ASYNCADA for box-shaped \mathcal{X} , by comparing x_t to the iterate that would have been created based on the content of the shared memory right before the start of the execution of the t -th iteration. This makes the analysis more complicated, and is not necessary when \mathcal{X} is not box-shaped; hence, we do not further pursue this construction in this paper.

Sparsity of the gradient estimates. For $t \in [T]$ and $j \in [d]$, let $p_{t,j}$ to denote the probability that the j -th coordinate of g_t is non-zero given the history $\hat{\mathcal{H}}_t$, that is, $p_{t,j} = \mathbb{P}\{g_t^{(j)} \neq 0 | \hat{\mathcal{H}}_t\}$. Let p_* denote an upper-bound on $\max_{t \in [T], j \in [d]} p_{t,j}$. We use p_* as a measure of the sparsity of the problem.⁸

⁸ In stochastic optimization with a finite-sum objective $f = \sum_{i=1}^m f_i$, where $g_t = \nabla f_{\xi_t}(x_t)$ and $\xi_t \in [m]$ is an index at time t sampled uniformly at random and independently of the history, one could measure the

Non-adaptive and time-decaying step-sizes. We first study the case when η_t is either a constant, or varies only as a function of the estimated iteration count \hat{t}_t . Recall that each concurrent iteration of the algorithms can be in read- and write-conflict with at most $\tau_c - 1$ other iterations, respectively, and that the algorithm uses τ parallel processes. Define $\tau_* = \max\{\tau_c, \tau\}$. The next theorem gives bounds on the regret of ASYNCADA under various scenarios. It is proved in Appendix C, where a similar result is also given for ASYNCADA(ρ) (Theorem 5).

Theorem 2. *Suppose that either all $f_t, t \in [T]$ are convex, or $\phi \equiv 0$ and $f_t \equiv f$ for some star-convex function f . Consider ASYNCADA running under Assumption 1 for $T > \tau_*^2$ updates, using $\gamma = 2\tau_*^2$. Let $\eta_0 > 0$. Then:*

(i) *If $\mathbb{E}\{\|g_t\|_2^2\} \leq G_*^2$ for all $t \in [T]$, then using a fixed $\eta_t = \eta_0\sqrt{T}$ or a time-varying $\eta_t = \eta_0\sqrt{\hat{t}_t}$,*

$$\frac{1}{T}\mathbb{E}\{R_T^{(f+\phi)}\} \leq \frac{1}{\sqrt{T}} \left(\eta_0 \|x^*\|_2^2 + \frac{2(1+p_*\tau_*^2)}{\eta_0} G_*^2 \right). \quad (8)$$

(ii) *If $f_t = f = \mathbb{E}_{\xi \sim \mathbb{P}_\Xi}\{F(x, \xi)\}$, $\sigma_*^2 := \mathbb{E}\{\|\nabla F(x^*, \cdot)\|_2^2\}$, and for all $\xi \in \Xi$, $F(\cdot, \xi)$ is convex and 1-smooth w.r.t. the norm $\|\cdot\|_l$ for some $l \in \mathbb{R}^d$ with positive entries, then given a constant $c_0 > 8(1+p_*\tau_*^2)$ and using a fixed $\eta_{t,i} = c_0 l_i + \eta_0\sqrt{T}$ or a time-varying $\eta_{t,i} = c_0 l_i + \eta_0\sqrt{\hat{t}_t}$,*

$$\frac{1}{T}\mathbb{E}\{R_T^{(f+\phi)}\} \leq \frac{c_0 \|x^*\|_l^2}{T} + \frac{2}{\sqrt{T}} \left(\eta_0 \|x^*\|_2^2 + \frac{4(1+p_*\tau_*^2)}{\eta_0} \sigma_*^2 \right). \quad (9)$$

(iii) *If ϕ is μ -strongly-convex and $\mathbb{E}\{\|g_t\|_2^2\} \leq G_*^2$ for all $t \in [T]$, then using $\eta_t \equiv 0$ or, equivalently, $\text{prox}(\hat{t}_t \phi, -z, 0) := \arg \min_{x \in \mathcal{X}} \hat{t}_t \phi(x) + \langle z, x \rangle = \nabla \phi^*(-z/\hat{t}_t)$,*

$$\frac{1}{T}\mathbb{E}\{R_T^{(f+\phi)}\} \leq \frac{(1+p_*\tau_*^2)G_*^2(1+\log(T))}{\mu T}. \quad (10)$$

Remark 1. If $c = p_*\tau_*^2$ is constant, the bounds match the corresponding serial bounds [16] up to constant factors, implying a linear speed-up. This also extends the analysis of ASYNC-DA [10] to non-box-shaped \mathcal{X} , non-zero ϕ , time-varying step sizes, and smooth and strongly-convex objectives.⁹

Remark 2. Note that (10) holds for all time steps, and converges to zero as T grows, without the knowledge of T or epoch-based updates. In case of ASYNCADA(ρ), the algorithm does not maintain an exact clock either. To our knowledge, this makes ASYNCADA(ρ) the first Hogwild!-style algorithm with an any-time guarantee without maintaining a global clock.

Remark 3. Since strongly convex functions have unbounded gradients on unbounded domains, it is not possible to impose a uniform bound on the gradient of $f + \phi$ in part (iii) for unconstrained optimization (i.e., when $\mathcal{X} = \mathbb{R}^d$). However, we only require the gradients of f , the non-strongly-convex part of the objective, to be bounded, which is a feasible assumption. Similarly, Nguyen et al. [24] analyzed strongly-convex optimization with unconstrained Hogwild! while avoiding the aforementioned uniform boundedness assumption, using a global clock. ASYNCADA(ρ) achieves the same result, but applies to arbitrary convex \mathcal{X} and ϕ , without requiring a global clock.

Adaptive step-sizes. Due to space constraints, we relegate the analysis of ASYNCADA(ρ) with AdaGrad step-sizes given by Line 10* to Appendix D.

6 HEDGEHOG: Hogwild-Style Hedge

Next, we present HEDGEHOG, which is, to our knowledge the first asynchronous version of the EG algorithm. The parallelization scheme is very similar to ASYNCADA, the difference being that EG uses multiplicative updates rather than additive SGD-style updates. We focus only on the case of $\phi \equiv 0$. Each process runs Lines 3–10 of Algorithm 2 concurrently with the other processes, sharing the dual vector z .

sparsity of the problem through a “conflict graph” [30, 20, 17, 27], which is a bi-partite graph with $f_i, i \in [m]$ on the left and coordinates $j \in [d]$ on the right, and an edge between f_i and coordinate j if $\nabla f_i(x)^{(j)}$ can be non-zero for some $x \in \mathcal{X}$. In this graph, let δ_j denote the degree of the node corresponding to coordinate j and Δ_r be the largest $\delta_j, j \in [d]$. Then, it is straightforward to see that $p_{t,j} \leq \delta_j/m$. Thus, $p_* = \Delta_r/m$ is a valid upper-bound, and gives the sparsity measure used, e.g., by Leblond et al. [17] and Pedregosa et al. [27].

⁹Note that under the conditions considered in [10], which include that \mathcal{X} is box-shaped and $\phi = 0$, ASYNC-DA requires a less restrictive sparsity regime of $p_*\tau_* \leq c$ for linear speed-up.

Algorithm 2: HEDGEHOG!: Asynchronous Stochastic Exponentiated Gradient.

Input: Step size η

1 Initialization

2 | Let $z \leftarrow 0$ be the shared sum of observed gradient estimates

3 repeat in parallel by each process

4 | $\hat{z} \leftarrow$ a full lock-free read of the shared dual vector z // $t \leftarrow t + 1$, denote $\hat{z}_{t-1} = \hat{z}$

5 | Receive ξ_t

6 | Compute the next iterate: $w_t^{(i)} \leftarrow \exp\left(-\hat{z}_{t-1}^{(i)}/\eta\right)$, $i = 1, 2, \dots, d$

7 | Normalize: $x_t \leftarrow w_t/\|w_t\|_1$

8 | Obtain the noisy gradient estimate: $g_t \leftarrow g_t(x_t, \xi_t)$

9 | **for** j such that $g_t^{(j)} \neq 0$ **do** $z^{(j)} \leftarrow z^{(j)} + g_t^{(j)}$ // atomic update

10 until terminated

As in ASYNCADA(ρ), we index the iterations by the time they finish the reading of z in Line 4 of HEDGEHOG (“after-read” labeling [18]). Similarly, we use $\hat{\mathcal{H}}_t = \{(\xi_s, \hat{z}_s)_{s=1}^{t-1}\}$ to denote the history of HEDGEHOG at time t , and use $\hat{\mathcal{H}}_t$ to define the sparsity measure p_* as in Section 5.1. Then, we have the following regret bound for HEDGEHOG.

Theorem 3. *Let \mathcal{X} be the probability simplex $\mathcal{X} = \{x | x^{(j)} > 0, \|x\|_1 = 1\}$, and suppose that either f_t are all convex, or $f_t \equiv f$ for a star-convex f . Assume that for all $t \in [T]$, the sampling of ξ_t in Line 5 of HEDGEHOG is independent of the history $\hat{\mathcal{H}}_t$. Then, after T updates, HEDGEHOG satisfies*

$$\mathbb{E}\{R_T^{(f)}\} \leq \eta \log(d) + \sum_{t=1}^T \mathbb{E}\left\{\frac{1 + \sqrt{p_*} \tau_*}{2\eta} \|g_t\|_\infty^2\right\}.$$

Remark 4. As in the case of ASYNCADA, as long as $\sqrt{p_*} \tau_*$ is a constant, the rate above matches the worst-case rate of serial EG up to constant factors, implying a linear speed-up. In particular, given an upper-bound G_* on $\mathbb{E}\{\|g_t\|_\infty\}$ and setting $\eta = G_*/\sqrt{T \log(d)}$, we recover the well-known $\mathcal{O}(G_* \sqrt{T \log(d)})$ rate for EG [14], but in the parallel asynchronous setting.

7 Conclusion, limitations, and future work

We presented and analyzed ASYNCADA, a parallel asynchronous online optimization algorithm with composite, adaptive updates, and global convergence rates under generic convex constraints and convex composite objectives which can be smooth, non-smooth, or non-strongly-convex. We also showed a similar global convergence for the so-called “star-convex” class of non-convex functions. Under all of the aforementioned settings, we showed that ASYNCADA enjoys linear speed-ups when the data is sparse. We also derived and analyzed HEDGEHOG, to our knowledge the first Hogwild-style asynchronous variant of the Exponentiated Gradient algorithm working on the probability simplex, and showed that HEDGEHOG enjoyed similar linear speed-ups.

To derive and analyze ASYNCADA and HEDGEHOG, we showed that the idea of perturbed iterates, used previously in the analysis of asynchronous SGD algorithms, naturally extends to generic dual-averaging algorithms, in the form of a perturbation in the “state” of the algorithm. Then, building on the work of Duchi et al. [10], we studied a unified framework for analyzing generic adaptive dual-averaging algorithms for composite-objective noisy online optimization (including ASYNCADA and HEDGEHOG as special cases). Possible directions for future research include applying the analysis to other problem settings, such as multi-armed bandits. In addition, it remains an open problem whether such an analysis is obtainable for constrained adaptive Mirror Descent without further restrictions on the regularizers (e.g., smoothness of the regularizer seems to help). Finally, the derivation of such data-dependent bounds for the final (rather than the average) iterate in stochastic optimization, without the usual strong-convexity and smoothness assumptions, remains an interesting open problem.

References

- [1] Heinz H Bauschke and Patrick L Combettes. *Convex analysis and monotone operator theory in Hilbert spaces*. Springer Science & Business Media, 2011.
- [2] Loris Cannelli et al. “Asynchronous Parallel Algorithms for Nonconvex Big-Data Optimization. Part I: Model and Convergence”. In: *arXiv preprint arXiv:1607.04818* (2017).
- [3] Loris Cannelli et al. “Asynchronous Parallel Algorithms for Nonconvex Big-Data Optimization. Part II: Complexity and Numerical Results”. In: *arXiv preprint arXiv:1701.04900* (2017).
- [4] Loris Cannelli et al. “Asynchronous parallel algorithms for nonconvex optimization”. In: *arXiv preprint arXiv:1607.04818* (2016).
- [5] Nicolò Cesa-Bianchi and Gábor Lugosi. *Prediction, Learning, and Games*. New York, NY, USA: Cambridge University Press, 2006.
- [6] Damek Davis, Brent Edmunds, and Madeleine Udell. “The sound of apalm clapping: Faster nonsmooth nonconvex optimization with stochastic asynchronous palm”. In: *Advances in Neural Information Processing Systems*. 2016, pp. 226–234.
- [7] Christopher De Sa et al. “Taming the Wild: A Unified Analysis of Hogwild!-Style Algorithms”. In: *arXiv preprint arXiv:1506.06438* (2015).
- [8] Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. “SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives”. In: *Advances in neural information processing systems*. 2014, pp. 1646–1654.
- [9] John Duchi, Elad Hazan, and Yoram Singer. “Adaptive Subgradient Methods for Online Learning and Stochastic Optimization”. In: *Journal of Machine Learning Research* 12 (July 2011), pp. 2121–2159.
- [10] John Duchi, Michael I Jordan, and Brendan McMahan. “Estimation, optimization, and parallelism when data is sparse”. In: *Advances in Neural Information Processing Systems*. 2013, pp. 2832–2840.
- [11] Francisco Facchinei, Gesualdo Scutari, and Simone Sagratella. “Parallel selective algorithms for nonconvex big data optimization”. In: *IEEE Transactions on Signal Processing* 63.7 (2015), pp. 1874–1889.
- [12] Olivier Fercoq and Peter Richtárik. “Optimization in high dimensions via accelerated, parallel, and proximal coordinate descent”. In: *SIAM Review* 58.4 (2016), pp. 739–771.
- [13] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*. Vol. 1. Springer series in statistics Springer, Berlin, 2001.
- [14] Elad Hazan. “Introduction to online convex optimization”. In: *Foundations and Trends in Optimization* 2.3-4 (2016), pp. 157–325.
- [15] Rie Johnson and Tong Zhang. “Accelerating stochastic gradient descent using predictive variance reduction”. In: *Advances in Neural Information Processing Systems*. 2013, pp. 315–323.
- [16] Pooria Joulani, András György, and Csaba Szepesvári. “A Modular Analysis of Adaptive (Non-) Convex Optimization: Optimism, Composite Objectives, and Variational Bounds”. In: *Proceedings of Machine Learning Research (Algorithmic Learning Theory 2017)*. 2017, pp. 681–720.
- [17] Rémi Leblond, Fabian Pedregosa, and Simon Lacoste-Julien. “Improved asynchronous parallel optimization analysis for stochastic incremental methods”. In: *arXiv preprint arXiv:1801.03749* (2018).
- [18] Rémi Leblond, Fabian Pedregosa, and Simon Lacoste-Julien. “ASAGA: asynchronous parallel SAGA”. In: *arXiv preprint arXiv:1606.04809* (2016).
- [19] Ji Liu et al. “An asynchronous parallel stochastic coordinate descent algorithm”. In: *arXiv preprint arXiv:1311.1873* (2013).
- [20] H. Mania et al. “Perturbed Iterate Analysis for Asynchronous Stochastic Optimization”. In: *ArXiv e-prints* (July 2015). arXiv: 1507.06970 [stat.ML].
- [21] H. Brendan McMahan. “A survey of Algorithms and Analysis for Adaptive Online Learning”. In: *Journal of Machine Learning Research* 18.90 (2017), pp. 1–50.
- [22] H. Brendan McMahan and Matthew Streeter. “Adaptive bound optimization for online convex optimization”. In: *Proceedings of the 23rd Conference on Learning Theory*. 2010.

- [23] Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*. Vol. 87. Springer Science & Business Media, 2013.
- [24] Lam M Nguyen et al. “SGD and Hogwild! convergence without the bounded gradients assumption”. In: *arXiv preprint arXiv:1802.03801* (2018).
- [25] Francesco Orabona, Koby Crammer, and Nicolò Cesa-Bianchi. “A generalized online mirror descent with applications to classification and regression”. English. In: *Machine Learning* 99.3 (2015), pp. 411–435.
- [26] Xinghao Pan et al. “Cyclades: Conflict-free asynchronous machine learning”. In: *Advances in Neural Information Processing Systems*. 2016, pp. 2568–2576.
- [27] Fabian Pedregosa, Rémi Leblond, and Simon Lacoste-Julien. “Breaking the Nonsmooth Barrier: A Scalable Parallel Method for Composite Optimization”. In: *Advances in Neural Information Processing Systems*. 2017, pp. 55–64.
- [28] Zhimin Peng et al. “Arock: an algorithmic framework for asynchronous parallel coordinate updates”. In: *SIAM Journal on Scientific Computing* 38.5 (2016), A2851–A2879.
- [29] Meisam Razaviyayn et al. “Parallel successive convex approximation for nonsmooth nonconvex optimization”. In: *Advances in Neural Information Processing Systems*. 2014, pp. 1440–1448.
- [30] Benjamin Recht et al. “Hogwild: A Lock-Free Approach to Parallelizing Stochastic Gradient Descent”. In: *Advances in Neural Information Processing Systems 24*. Ed. by J. Shawe-Taylor et al. Curran Associates, Inc., 2011, pp. 693–701.
- [31] Gesualdo Scutari, Francisco Facchinei, and Lorenzo Lampariello. “Parallel and distributed methods for constrained nonconvex optimization—Part I: Theory”. In: *IEEE Transactions on Signal Processing* 65.8 (2016), pp. 1929–1944.
- [32] Gesualdo Scutari and Ying Sun. “Parallel and distributed successive convex approximation methods for big-data optimization”. In: *Multi-agent Optimization*. Springer, 2018, pp. 141–308.
- [33] Gesualdo Scutari et al. “Parallel and distributed methods for constrained nonconvex optimization-part ii: Applications in communications and machine learning”. In: *IEEE Transactions on Signal Processing* 65.8 (2016), pp. 1945–1960.
- [34] Shai Shalev-Shwartz. “Online learning and online convex optimization”. In: *Foundations and Trends in Machine Learning* 4.2 (2011), pp. 107–194.
- [35] Tao Sun, Robert Hannah, and Wotao Yin. “Asynchronous coordinate descent under more realistic assumptions”. In: *Advances in Neural Information Processing Systems*. 2017, pp. 6182–6190.
- [36] Yu-Xiang Wang et al. “Parallel and distributed block-coordinate Frank-Wolfe algorithms”. In: *International Conference on Machine Learning*. 2016, pp. 1548–1557.
- [37] Lin Xiao. “Dual averaging method for regularized stochastic learning and online optimization”. In: *Advances in Neural Information Processing Systems*. 2009, pp. 2116–2124. (Visited on 02/05/2015).

A Proofs for the generic framework

Proof of Lemma 1. The proof follows in the same way as in the serial setting [16]. For $t \in [T]$,

$$\begin{aligned} f_t(x_t) - f_t(x^*) &= \langle \nabla f_t(x_t), x_t - x^* \rangle - \mathcal{B}_{f_t}(x^*, x_t) \\ &= \langle g_t, x_t - x^* \rangle + \langle \nabla f_t(x_t) - g_t, x_t - x^* \rangle - \mathcal{B}_{f_t}(x^*, x_t) \\ &= \langle g_t, \tilde{x}_t - x^* \rangle + \langle g_t, x_t - \tilde{x}_t \rangle + \delta_t - B_t \\ &= \langle g_t, \tilde{x}_t - x^* \rangle + \phi(\tilde{x}_t) - \phi(x_t) + \tilde{\epsilon}_t + \delta_t - B_t \end{aligned}$$

Adding $\phi(x_t) - \phi(x^*)$ to both sides and summing over t completes the proof. \square

Perturbation structure. We assume that the difference of \hat{z}_{t-1} and z_{t-1} is that zero or more coordinates $g_s^{(j)}$ from the past feedback vectors $g_s, s \in [t-1]$, can be missing from (i.e., not added in) the perturbed dual vector \hat{z}_{t-1} . Formally, for all $t \in [T]$ and $j \in [d]$,

$$\hat{z}_{t-1}^{(j)} = g_{1:t-1}^{(j)} - \sum_{s \in O_{t,j}} g_s^{(j)}, \quad (11)$$

where $O_{t,j}$ is the subset of the past indices $[t-1]$ corresponding to the missing updates at the j -th coordinate. Written in a more compact form,

$$\hat{z}_{t-1} = g_{1:t-1} - \sum_{s \in O_t} I_{t,s} g_s, \quad (12)$$

where $O_t = \cup_j O_{t,j}$ is the set of all time steps with missing information at time t (that is, the set of iterations with which iteration t is in read-conflict), and $I_{t,s}, s \in [t-1]$, are diagonal $d \times d$ matrices with $I_{t,s}^{(j,j)} = 1$ if $g_s^{(j)}$ is missing from \hat{z}_{t-1} and 0 otherwise. We define $\tau_{t,j} = |O_{t,j}|$ and $\tau_t = |O_t|$ to denote, respectively, the total number of missing updates to the j -th coordinate of \hat{z}_{t-1} , and to the whole vector \hat{z}_{t-1} . Similarly, we assume that the time-counter \hat{t}_t may not be equal to t , and the cumulative regularizers $r_{0:t}$ and $\hat{r}_{0:t}$, can be different, with the latter using only some of the past updates made to $r_{0:t}$. However, the exact perturbation in \hat{t}_t and $\hat{r}_{0:t}$ depends on the specifics of the algorithm. Our analysis isolates these perturbations in individual terms, which we can subsequently study on a case-by-case basis. We make the following assumption on \hat{t}_t and the sequence of actual regularizers $(\hat{r}_t)_{t=0}^T$ and ideal regularizers $(r_t)_{t=0}^T$.

Assumption 2. The regularizers $r_t, \hat{r}_t, t = 0, 1, \dots, T$, are admissible ADA-FTRL regularizers (Definition 2) with the same sequence of norms $\|\cdot\|_{(t)}$, and the sequence of norms is non-decreasing: $\|\cdot\|_{(t)} \geq \|\cdot\|_{(t-1)}$ for all $t = 1, 2, \dots, T$. Finally, $r_t \geq 0, t = 0, 1, 2, \dots, T$, and $\hat{t}_t > t, t = 1, 2, \dots, T$.

Intuitively, Assumption 2 states that the regularizers \hat{r}_t are not fundamentally different from the regularizers r_t as far as the basic properties of ADA-FTRL are concerned. In particular, the assumption is satisfied if $(r_t)_{t=0}^T$ is admissible with a non-decreasing sequence of norms and the perturbation increases the curvature, that is, $\hat{r}_{0:t-1} - r_{0:t-1}$ is convex. Finally, the assumption $\hat{t}_t > t$ helps us in providing bounds for composite-objective learning, as will become clear later.

Independence assumption. Similarly to the standard serial setting, we will assume that the outcome ξ_t at time t is independent of the history that determines x_t . In the case of perturbed ADA-FTRL, we define the history to depend on the actual states the *perturbed* ADA-FTRL algorithm has gone through:

Definition 1 (History of the perturbed game). For $t = 1, 2, \dots, T$, the *history of the perturbed game* up to time t is defined as

$$\hat{\mathcal{H}}_t = \left\{ (\xi_s, \hat{z}_s, \hat{t}_s, \hat{r}_{0:s})_{s=1}^{t-1} \right\},$$

where $\hat{z}_s, \hat{r}_{0:s}, \hat{t}_s$ are the dual vector, regularizer and time-counter used by the $(s+1)$ -th perturbed ADA-FTRL update.

We assume that the stochastic outcomes are independent of the history:

Assumption 3 (Independence of ξ_t). For all $t = 1, 2, \dots, T$, the t -th sample ξ_t is independent of the history $\hat{\mathcal{H}}_t$.

This in turn means that ξ_t is independent of x_t as well as x_s and ξ_s for all $s < t$.

We call a norm $\|\cdot\|$ a *weighted q -norm* if there exists $q > 0$ and $a_j, j \in [d]$ such that for all $x \in \mathbb{R}^d$,

$$\|x\| = \left(\sum_{j=1}^d a_j |x^{(j)}|^q \right)^{1/q}. \quad (13)$$

The next theorem describes a generic data-dependent bound on the regret of perturbed ADA-FTRL.

Theorem 4. *Suppose that Perturbed-ADA-FTRL is run under Assumption 3, and Assumption 2 holds such that for each $t \in [T]$, $\|\cdot\|_{(t)}$ is a weighted q -norm with $q = 1$ or $q = 2$. For all $t \in [T]$, define $\Delta_t = r_{0:t-1}(x_t) - r_{0:t-1}(\tilde{x}_t) + \hat{r}_{0:t-1}(\tilde{x}_t) - \hat{r}_{0:t-1}(x_t)$, and $\nu_t = \hat{t}_t - t$ with the \hat{t}_t used in the Perturbed-ADA-FTRL update (5). Then, the regret of Perturbed-ADA-FTRL satisfies*

$$\mathbb{E}\left\{R_T^{(f+\phi)}\right\} \leq \mathbb{E}\left\{r_{0:T}(x^*) + \sum_{t=1}^T \left(\frac{1 + p_* \nu_t + \sum_{s:t \in O_s} \frac{\tau_s}{\nu_s}}{2} \|g_t\|_{(t,*)}^2 + \frac{\Delta_t}{\nu_t} \right) - B_{1:T} \right\},$$

where p_* is a global upper-bound on $\mathbb{P}\{g_t^{(j)} \neq 0 | \hat{\mathcal{H}}_t\}$.

A.1 Proof of Theorem 4

First, we upper-bound $\tilde{\epsilon}_t$ in terms of the difference between \tilde{x}_t and x_t .

Lemma 2. *Consider Perturbed-ADA-FTRL under the conditions of Theorem 4. Let $\beta_t \in \mathbb{R}^d$ be given by $\beta_t^{(j)} = \mathbb{I}\{g_t^{(j)} \neq 0\}$, and use \odot to denote elementwise vector multiplication. Then,*

- For any positive real number c_t and any norm $\|\cdot\|$, we have

$$\tilde{\epsilon}_t + \phi(\tilde{x}_t) - \phi(x_t) \leq \frac{c_t}{2} \|g_t\|_*^2 + \frac{1}{2c_t} \|\beta_t \odot (x_t - \tilde{x}_t)\|^2,$$

- In the stochastic setting under Assumption 3, for any $c_t > 0$ and any norm $\|\cdot\|$,

$$\mathbb{E}\{\tilde{\epsilon}_t + \phi(\tilde{x}_t) - \phi(x_t)\} \leq \mathbb{E}\left\{\frac{c_t}{2} \|\nabla f_t(x_t)\|_*^2 + \frac{1}{2c_t} \|x_t - \tilde{x}_t\|^2\right\}.$$

- Under Assumption 3, for any $q \geq 1$, any weighted q -norm $\|\cdot\|$ determined by the history $\hat{\mathcal{H}}_t$, and any positive scalar $c_t \in \sigma(\hat{\mathcal{H}}_t)$,

$$\mathbb{E}\{\tilde{\epsilon}_t + \phi(\tilde{x}_t) - \phi(x_t)\} \leq \mathbb{E}\left\{\frac{c_t}{2} \|g_t\|_*^2\right\} + p_*^{(1/q)} \mathbb{E}\left\{\frac{1}{2c_t} \|(x_t - \tilde{x}_t)\|^2\right\},$$

where p_* is a global upper-bound on $\mathbb{P}\{g_t^{(j)} \neq 0 | \hat{\mathcal{H}}_t\}$. In case of $q = 2$, the bound still holds if $p_*^{1/2}$ is replaced with p_* .

Proof of Lemma 2. To get the first inequality, note that $g_t = \beta_t \odot g_t$ by definition. The bound then follows by the Fenchel-Young inequality.

To get the second bound, note that $x_t, \tilde{x}_t \in \sigma(\hat{\mathcal{H}}_t)$ by construction, so by Assumption 3,

$$\mathbb{E}\{\langle g_t - \nabla f_t(x_t), x_t - \tilde{x}_t \rangle\} = \mathbb{E}\left\{\langle \mathbb{E}\{g_t - \nabla f_t(x_t) | \hat{\mathcal{H}}_t\}, x_t - \tilde{x}_t \rangle\right\} = 0.$$

Thus, $\mathbb{E}\{\tilde{\epsilon}_t + \phi(\tilde{x}_t) - \phi(x_t)\} = \mathbb{E}\{\langle \nabla f_t(x_t), x_t - \tilde{x}_t \rangle\}$, and the result follows by the Fenchel-Young inequality.

To get the third bound, we first start with the simpler case of $q = 2$, using $a \in \sigma(\hat{\mathcal{H}}_t)$ to denote the associated weighting vector, then apply the first inequality and take expectation of the terms $\|\beta_t \odot (x_t - \tilde{x}_t)\|^2$. Note that by construction, $x_t, \tilde{x}_t \in \sigma(\hat{\mathcal{H}}_t)$. Furthermore, by assumption, $c_t, a \in \sigma(\hat{\mathcal{H}}_t)$. Hence,

$$\begin{aligned}
\mathbb{E}\left\{\frac{1}{2c_t}\|\beta_t \odot (x_t - \tilde{x}_t)\|^2\right\} &= \mathbb{E}\left\{\mathbb{E}\left\{\frac{1}{2c_t}\|\beta_t \odot (x_t - \tilde{x}_t)\|^2 \mid \hat{\mathcal{H}}_t\right\}\right\} \\
&= \mathbb{E}\left\{\sum_{j=1}^d \mathbb{E}\left\{\frac{1}{2c_t}a^{(j)}\beta_t^{(j)}(x_t^{(j)} - \tilde{x}_t^{(j)})^2 \mid \hat{\mathcal{H}}_t\right\}\right\} \\
&= \mathbb{E}\left\{\sum_{j=1}^d \mathbb{E}\left\{\mathbb{I}\{g_t^{(j)} \neq 0\} \mid \hat{\mathcal{H}}_t\right\} \frac{1}{2c_t}a^{(j)}(x_t^{(j)} - \tilde{x}_t^{(j)})^2\right\} \\
&= \mathbb{E}\left\{\sum_{j=1}^d p_{t,j} \frac{1}{2c_t}a^{(j)}(x_t^{(j)} - \tilde{x}_t^{(j)})^2\right\} \\
&\leq \left(\max_{j \in [d]} p_{t,j}\right) \mathbb{E}\left\{\frac{1}{2c_t} \sum_{j=1}^d a^{(j)}(x_t^{(j)} - \tilde{x}_t^{(j)})^2\right\},
\end{aligned}$$

completing the proof.

To get the bound for any $q \geq 1$, first note that when $q \in [1, \infty)$, the function $h : [0, \infty) \rightarrow \mathbb{R}$ given by $h(x) := x^{1/q}$ (with $h(0) := 0$) is concave for all $x > 0$. Thus, by Jensen's inequality, $\mathbb{E}\{h(X)\} \leq h(\mathbb{E}\{X\})$ for any non-negative random variable X . Next, we let the q -norm in question be given by (13), with $a \in \sigma(\hat{\mathcal{H}}_t)$ denoting the associated weighting vector, and continue as in the case of $q = 2$ above:

$$\begin{aligned}
\mathbb{E}\left\{\frac{1}{2c_t}\|\beta_t \odot (x_t - \tilde{x}_t)\|^2\right\} &= \mathbb{E}\left\{\mathbb{E}\left\{\frac{1}{2c_t}\|\beta_t \odot (x_t - \tilde{x}_t)\|^2 \mid \hat{\mathcal{H}}_t\right\}\right\} \\
&= \mathbb{E}\left\{\frac{1}{2c_t} \mathbb{E}\left\{\left(\sum_{j=1}^d a^{(j)}\beta_t^{(j)} |x_t^{(j)} - \tilde{x}_t^{(j)}|^q\right)^{2/q} \mid \hat{\mathcal{H}}_t\right\}\right\} \\
&\leq \mathbb{E}\left\{\frac{1}{2c_t} \left(\mathbb{E}\left\{\left(\sum_{j=1}^d a^{(j)}\beta_t^{(j)} |x_t^{(j)} - \tilde{x}_t^{(j)}|^q\right)^2 \mid \hat{\mathcal{H}}_t\right\}\right)^{1/q}\right\},
\end{aligned}$$

where the last inequality follows since $\mathbb{E}\{h(X) \mid \hat{\mathcal{H}}_t\} \leq h(\mathbb{E}\{X \mid \hat{\mathcal{H}}_t\})$ by the concavity of h as argued above, where $X = \left(\sum_{j=1}^d a^{(j)}\beta_t^{(j)} |x_t^{(j)} - \tilde{x}_t^{(j)}|^q\right)^2$. On the other hand, since h is also increasing, we can bound $h(\mathbb{E}\{X \mid \hat{\mathcal{H}}_t\})$ by first upper-bounding $\mathbb{E}\{X \mid \hat{\mathcal{H}}_t\}$. In particular,

$$\begin{aligned}
\mathbb{E}\{X \mid \hat{\mathcal{H}}_t\} &= \mathbb{E}\left\{\left(\sum_{j=1}^d a^{(j)}\beta_t^{(j)} |x_t^{(j)} - \tilde{x}_t^{(j)}|^q\right)^2 \mid \hat{\mathcal{H}}_t\right\} \\
&= \mathbb{E}\left\{\sum_{j=1}^d a^{(j)}\beta_t^{(j)} |x_t^{(j)} - \tilde{x}_t^{(j)}|^q \left(\sum_{i=1}^d a^{(i)}\beta_t^{(i)} |x_t^{(i)} - \tilde{x}_t^{(i)}|^q\right) \mid \hat{\mathcal{H}}_t\right\} \\
&\leq \mathbb{E}\left\{\sum_{j=1}^d a^{(j)}\beta_t^{(j)} |x_t^{(j)} - \tilde{x}_t^{(j)}|^q \left(\sum_{i=1}^d a^{(i)} |x_t^{(i)} - \tilde{x}_t^{(i)}|^q\right) \mid \hat{\mathcal{H}}_t\right\}
\end{aligned}$$

$$\begin{aligned}
&= \sum_{j=1}^d a^{(j)} \mathbb{E} \left\{ \beta_t^{(j)} \mid \hat{\mathcal{H}}_t \right\} \left| x_t^{(j)} - \tilde{x}_t^{(j)} \right|^q \left(\sum_{i=1}^d a^{(i)} \left| x_t^{(i)} - \tilde{x}_t^{(i)} \right|^q \right) \\
&\leq p_* \sum_{j=1}^d a^{(j)} \left| x_t^{(j)} - \tilde{x}_t^{(j)} \right|^q \left(\sum_{i=1}^d a^{(i)} \left| x_t^{(i)} - \tilde{x}_t^{(i)} \right|^q \right) \\
&= p_* \|x_t - \tilde{x}_t\|^{2q}.
\end{aligned}$$

Thus, $h \left(\mathbb{E} \left\{ X \mid \hat{\mathcal{H}}_t \right\} \right) \leq h \left(p_* \|x_t - \tilde{x}_t\|^{2q} \right) = p_*^{1/q} \|x_t - \tilde{x}_t\|^2$. Thus,

$$\begin{aligned}
\mathbb{E} \left\{ \frac{1}{2c_t} \|\beta_t \odot (x_t - \tilde{x}_t)\|^2 \right\} &\leq \mathbb{E} \left\{ \frac{1}{2c_t} h \left(\mathbb{E} \left\{ X \mid \hat{\mathcal{H}}_t \right\} \right) \right\} \\
&\leq \mathbb{E} \left\{ \frac{1}{2c_t} p_*^{1/q} \|x_t - \tilde{x}_t\|^2 \right\},
\end{aligned}$$

completing the proof of the third bound. \square

Thus, controlling the regret in perturbed optimization reduces to picking a suitable norm $\|\cdot\|$ and applying Lemma 2 at each time step t , and then controlling the differences $x_t - \tilde{x}_t$. To that end, we use the stability of ADA-FTRL updates, that is, that the difference of two ADA-FTRL iterates is controlled by the difference in the two states of the algorithm resulting in the iterates. The following lemma provides this stability bound.

Lemma 3. *Let $(x_t)_{t=1}^T$ and $(\tilde{x}_t)_{t=1}^T$ be given by updates (5) and (6), respectively, and suppose that Assumption 2 holds. Define $\Delta_t = r_{0:t-1}(x_t) - r_{0:t-1}(\tilde{x}_t) + \hat{r}_{0:t-1}(\tilde{x}_t) - \hat{r}_{0:t-1}(x_t)$ for $t = 1, 2, \dots, T$. Then, for all $t = 1, 2, \dots, T+1$,*

$$\frac{1}{2} \|x_t - \tilde{x}_t\|_{(t)}^2 \leq \frac{1}{2} \left\| \sum_{s \in O_t} I_{t,s} g_s \right\|_{(t,*)}^2 + (t - \hat{t}_t) (\phi(x_t) - \phi(\tilde{x}_t)) + \Delta_t. \quad (14)$$

Proof of Lemma 3. Since both $(r_t)_{t=1}^T$ and $(\hat{r}_t)_{t=1}^T$ are admissible, the ADA-FTRL margin lemma [16, Lemma 24 (Appendix F)] applied to the update (5) implies that for all $t = 1, 2, \dots, T$,

$$\langle \hat{z}_{t-1}, \tilde{x}_t - x_t \rangle + \hat{t}_t (\phi(\tilde{x}_t) - \phi(x_t)) + \hat{r}_{0:t-1}(\tilde{x}_t) - \hat{r}_{0:t-1}(x_t) \geq \mathcal{B}_{\hat{t}_t \phi + \hat{r}_{0:t-1}}(\tilde{x}_t, x_t),$$

while for update (6) we have

$$\langle z_{t-1}, x_t - \tilde{x}_t \rangle + t (\phi(x_t) - \phi(\tilde{x}_t)) + r_{0:t-1}(x_t) - r_{0:t-1}(\tilde{x}_t) \geq \mathcal{B}_{t\phi + r_{0:t-1}}(x_t, \tilde{x}_t).$$

By the strong convexity of $t\phi + r_{0:t-1}$ and $\hat{t}_t\phi + \hat{r}_{0:t-1}$ w.r.t. $\|\cdot\|_{(t)}$, convexity of ϕ , and the fact that $\hat{t}_t > t > 0$ (so that $\hat{t}_t\phi + \hat{r}_{0:t-1}$ is also strongly-convex w.r.t. $\|\cdot\|_{(t)}$), we have $\mathcal{B}_{\hat{t}_t\phi + \hat{r}_{0:t-1}}(\tilde{x}_t, x_t) \geq \frac{1}{2} \|x_t - \tilde{x}_t\|_{(t)}^2$ and $\mathcal{B}_{t\phi + r_{0:t-1}}(x_t, \tilde{x}_t) \geq \frac{1}{2} \|x_t - \tilde{x}_t\|_{(t)}^2$. Adding the above,

$$\begin{aligned}
\frac{1}{2} \|x_t - \tilde{x}_t\|_{(t)}^2 &\leq -\frac{1}{2} \|x_t - \tilde{x}_t\|_{(t)}^2 + \langle z_{t-1} - \hat{z}_{t-1}, x_t - \tilde{x}_t \rangle + (t - \hat{t}_t) (\phi(x_t) - \phi(\tilde{x}_t)) \\
&\quad + (r_{0:t-1}(x_t) - \hat{r}_{0:t-1}(x_t)) - (r_{0:t-1}(\tilde{x}_t) - \hat{r}_{0:t-1}(\tilde{x}_t)) \\
&= -\frac{1}{2} \|x_t - \tilde{x}_t\|_{(t)}^2 + \left\langle \sum_{s \in O_t} I_{t,s} g_s, x_t - \tilde{x}_t \right\rangle + (t - \hat{t}_t) (\phi(x_t) - \phi(\tilde{x}_t)) + \Delta_t \\
&\leq \frac{1}{2} \left\| \sum_{s \in O_t} I_{t,s} g_s \right\|_{(t,*)}^2 + (t - \hat{t}_t) (\phi(x_t) - \phi(\tilde{x}_t)) + \Delta_t, \quad (15)
\end{aligned}$$

where in the last step we have used the Fenchel-Young inequality, completing the proof. \square

We can now prove the theorem.

Proof of Theorem 4. For $t = 1, 2, \dots, T$, recall that the imaginary iterate \tilde{x}_t is defined by (6)

$$\tilde{x}_t = \arg \min_{x \in \mathcal{X}} \langle g_{1:t-1}, x \rangle + t\phi(x) + r_{0:t-1}(x),$$

and note that in addition to the difference between $r_{0:t-1}$ and $\hat{r}_{0:t-1}$, the actual iterate x_t and the imaginary iterate \tilde{x}_t have a difference of $\nu_t\phi(x)$ in their regularization.

Starting from the regret decomposition, and using the linear regret of the imaginary iterate \tilde{x}_t , as well as the fact that r_t are non-negative by Assumption 2, we have

$$\begin{aligned} R_T^{(f+\phi)}(x^*) &\leq \sum_{t=1}^T \langle g_t, \tilde{x}_t - x^* \rangle + \tilde{\epsilon}_{1:T} + \delta_{1:T} - B_{1:T} + \sum_{t=1}^T (\phi(\tilde{x}_t) - \phi(x^*)) \\ &\leq r_{0:T}(x^*) - \sum_{t=0}^T r_t(\tilde{x}_{t+1}) + \sum_{t=1}^T \frac{1}{2} \|g_t\|_{(t,*)}^2 + \tilde{\epsilon}_{1:T} + \delta_{1:T} - B_{1:T} \\ &\leq r_{0:T}(x^*) + \sum_{t=1}^T \frac{1}{2} \|g_t\|_{(t,*)}^2 + \tilde{\epsilon}_{1:T} + \delta_{1:T} - B_{1:T}. \end{aligned} \quad (16)$$

In the above, the first inequality follows by Lemma 1. The second inequality follows by bounding the linear regret $\sum_{t=1}^T \langle g_t, \tilde{x}_t - x^* \rangle$ using Theorem 1, and the third by dropping the non-negative terms $r_t(\tilde{x}_{t+1})$.

Next, we bound the penalty terms $\tilde{\epsilon}_{1:T}$. For each $t = 1, 2, \dots, T$, using the fact that $\nu_t > 0$ by Assumption 2 and $\nu_t \in \sigma(\hat{\mathcal{H}}_t)$ by definition, we have

$$\begin{aligned} \mathbb{E}\{\tilde{\epsilon}_t + \phi(\tilde{x}_t) - \phi(x_t)\} &\leq \mathbb{E}\left\{\frac{p_*\nu_t}{2} \|g_t\|_{(t,*)}^2 + \frac{1}{2\nu_t} \|(x_t - \tilde{x}_t)\|_{(t)}^2\right\} \\ &\leq \mathbb{E}\left\{\frac{p_*\nu_t}{2} \|g_t\|_{(t,*)}^2\right\} \\ &\quad + \mathbb{E}\left\{\frac{1}{2\nu_t} \left(\left\|\sum_{s \in O_t} I_{t,s} g_s\right\|_{(t,*)}^2 + 2(\nu_t\phi(\tilde{x}_t) - \nu_t\phi(x_t) + \Delta_t)\right)\right\} \\ &\leq \mathbb{E}\left\{\frac{p_*\nu_t}{2} \|g_t\|_{(t,*)}^2 + \sum_{s \in O_t} \frac{\tau_t}{2\nu_t} \|I_{t,s} g_s\|_{(t,*)}^2 + \frac{\Delta_t}{\nu_t} + \phi(\tilde{x}_t) - \phi(x_t)\right\} \\ &\leq \mathbb{E}\left\{\frac{p_*\nu_t}{2} \|g_t\|_{(t,*)}^2 + \sum_{s \in O_t} \frac{\tau_t}{2\nu_t} \|g_s\|_{(t,*)}^2 + \frac{\Delta_t}{\nu_t} + \phi(\tilde{x}_t) - \phi(x_t)\right\} \\ &\leq \mathbb{E}\left\{\frac{p_*\nu_t}{2} \|g_t\|_{(t,*)}^2 + \sum_{s \in O_t} \frac{\tau_t}{2\nu_t} \|g_s\|_{(s,*)}^2 + \frac{\Delta_t}{\nu_t} + \phi(\tilde{x}_t) - \phi(x_t)\right\}. \end{aligned} \quad (17)$$

The first inequality above uses Lemma 2 with $c_t = p_*\nu_t$ (using the assumption of $\|\cdot\|_{(t)}$ being a weighted q -norm with $q = 1$ or $q = 2$), the second follows by Lemma 3, the third uses the convexity of the norms $\|\cdot\|_{(t,*)}^2$ and Jensen's inequality, the fourth follows because $I_{t,s}$ is a $\{0, 1\}$ -valued diagonal matrix and $\|\cdot\|_{(t)}$ is a weighted q -norm, and hence $\|I_{t,s} g_s\|_{(t,*)} \leq \|g_s\|_{(t,*)}$, and the last line follows because $s \in O_t$ implies $s \leq t$ by construction, and for $s \leq t$, the dual norms satisfy $\|\cdot\|_{(t,*)} \leq \|\cdot\|_{(s,*)}$ by Assumption 2. Summing the second term on the r.h.s. of (17), for $t = 1, 2, \dots, T$, we get

$$\sum_{t=1}^T \sum_{s \in O_t} \frac{\tau_t}{2\nu_t} \|g_s\|_{(s,*)}^2 = \sum_{s=1}^T \left(\sum_{t:s \in O_t} \frac{\tau_t}{2\nu_t}\right) \|g_s\|_{(s,*)}^2, \quad (18)$$

Thus, summing (17) over t , combining with (18), and noting that the terms $\phi(x_t) - \phi(\tilde{x}_t)$ cancel from the sides of the asynchrony penalty bounds (17), we get

$$\mathbb{E}\left\{R_T^{(f+\phi)}(x^*)\right\} \leq \mathbb{E}\left\{r_{0:T}(x^*) + \sum_{t=1}^T \frac{1 + p_*\nu_t + \sum_{s:t \in O_s} \frac{\tau_s}{\nu_s}}{2} \|g_t\|_{(t,*)}^2 + \frac{\Delta_t}{\nu_t} + \delta_{1:T}\right\}.$$

Finally, noting that $x_t \in \sigma(\mathcal{H}_t)$ by definition, by Assumption 3 it follows that $\mathbb{E}\{\delta_t | \mathcal{H}_t\} = 0$ in the stochastic setting. This completes the proof. \square

B ASYNCADA(ρ): ASYNCADA with inexact clock

In this section, we present ASYNCADA(ρ), a more general version of ASYNCADA that maintains the global clock sparsely. In the context of ASYNCADA(ρ), we use t' to denote the *clock* variable in the shared memory, and use t to denote the virtual iteration index as we specify below. The processes run copies of Algorithm 3 concurrently. Each process is also equipped with an internal counter t'' and a function `MaintainClock` to control the updating of the global clock t' .

Similar notes as in ASYNCADA apply regarding the maintenance of the step-size η and the formation of the average iterate. Note, however, that unlike ASYNCADA, step-sizes changing with time need to use \hat{t}_t rather than t , since the latter is not available anymore. As Theorem 5 in Appendix C shows, this has a negligible effect on the convergence guarantees.

Algorithm 3: ASYNCADA(ρ): ASYNCADA with inexact clock

Input: clock update frequency ρ

- 1 Initialize internal local counter $t'' \leftarrow 0$
- 2 **repeat**
- 3 $\hat{\eta} \leftarrow$ a full (lock-free) read of the shared step-sizes η
- 4 $\hat{z} \leftarrow$ a full (lock-free) read of the shared dual vector z
- 5 $\hat{t} \leftarrow$ `MaintainClock()` // $t \leftarrow t + 1$, denote $\hat{z}_{t-1} = \hat{z}$, $\hat{\eta}_t = \hat{\eta}$, $\hat{t}_t = \hat{t}$
- 6 Receive ξ_t
- 7 Compute the next iterate: $x_t \leftarrow \text{prox}(\hat{t}_t \phi, -\hat{z}_{t-1}, \hat{\eta}_t)$ // `prox` defined in (3)
- 8 Obtain the noisy gradient estimate: $g_t \leftarrow g_t(x_t, \xi_t)$
- 9 **for** j such that $g_t^{(j)} \neq 0$ **do** $z^{(j)} \leftarrow z^{(j)} + g_t^{(j)}$ // atomic update
- 10 Update the shared step-size vector η
- 11 **until terminated**

Algorithm 4: Maintaining the local and global iteration counters

- 1 **Function** `MaintainClock()`
- 2 Let $\gamma > t'' \tau_*$
- 3 $t'' \leftarrow t'' + 1$ // count number of iterations by this process
- 4 **if** $t'' \geq \rho$ **then** // Update global clock every ρ local iterations
- 5 $t'' \leftarrow 0$
- 6 $t' \leftarrow t' + \rho$ // atomic read-increment
- 7 **end if**
- 8 **return** $t' + \gamma$ // Use the value of t' read in Line 6 if executed; otherwise read t'
- 9 **end**

Indexing the iterates Unlike ASYNCADA, in ASYNCADA(ρ) the iterates are not physically indexed by the global clock t . As such, at each point in time we define a virtual count of the number of iterations undertaken so far, and then come up with an actual estimate of this virtual global clock. To that end, we use the “after-read” iteration indexing proposed by Leblond et al. [17]: we define the t -th iteration to be the one corresponding to the t -th completion of reading the shared memory (which happens by reading t' in Line 6 or 8 of `MaintainClock`), before the execution of Line 6. This ensures that \hat{z}_{t-1} contains only updates made by processes $s < t$, which proves useful in the analysis.

Estimating the clock. In ASYNCADA(ρ), the processes share an integer t' to estimate the (virtual) iteration count t , which is updated by each process every ρ iterations. In particular, in each iteration a process makes one call to the function `MaintainClock` (Algorithm 4), which increments its local counter t'' of the number of updates made by that process since it last updated the global clock; then, after every ρ local updates, `MaintainClock` increments the shared global clock estimate t' by ρ (and resets t'' for that process). Note that in this way, t' is always an under-estimate of t , and

ASYN-CADA (Algorithm 1) is recovered when $\rho = 1$. Again, when $\phi \equiv 0$ and t is not used for setting the step-sizes η either, there is no need to maintain t' physically, and the call to `MaintainClock` can be omitted in Algorithm 3.

C Proofs for ASYN-CADA and ASYN-CADA(ρ)

We start the analysis by a lemma on the time estimates formed by ASYN-CADA(ρ).

Lemma 4 (Time estimate of ASYN-CADA(ρ)). *Suppose ASYN-CADA(ρ) is run for T iterations with any $\rho \geq 1$, using $\gamma \geq \rho\tau_* + \tau_*^2$. Then, the estimated clock \hat{t}_t is non-decreasing with t , i.e., for all $s, t \in [T]$ with $s < t$, we have $\hat{t}_s \leq \hat{t}_t$. In addition, for all $t \in [T]$, we have $\hat{t}_t > t + \tau_*^2$.*

Proof. Fix $s < t \in [T]$, and note that the value of t' read in `MaintainClock` in the s -th iteration cannot be greater than the value of t' read in the t -th iteration. Specifically, t' can only increase over (physical) time, and the iterations are indexed by the time they make their last reading of the shared memory before the update in Line 7 of ASYN-CADA(ρ), which is the reading (and possibly incrementing) of t' in Line 8 (respectively, Line 6) of `MaintainClock`. Thus, the reading of t' in iteration $s < t$ necessarily has happened before that of t , leading to a smaller value of t' . As all the processes are adding the same fixed value of γ to t' to obtain \hat{t} , this implies $\hat{t}_s \leq \hat{t}_t$.

To see that $\hat{t}_t > t$, fix t and let t' be the value of the global clock estimate at the end of the call to `MaintainClock` in the t -th iteration. Since there have been at most $\rho - 1$ updates in each of the τ processors since the last update of t' by each processor, $t' > t - \rho\tau \geq t - \rho\tau_*$, where the second inequality holds by the definition of τ_* . As such, $\hat{t}_t = t' + \gamma \geq t' + \rho\tau_* + \tau_*^2 > t + \tau_*^2$. \square

Proof of Theorem 2. The theorem follows immediately from the convergence bound of ASYN-CADA(ρ) with $\rho = 1$, given by Theorem 5 below. \square

To analyze ASYN-CADA(ρ), we need to make a slightly modified version of Assumption 1, since \hat{t}_t is now a non-deterministic part of the state of the algorithm:

Assumption 4 (Independence of ξ_t). For all $t = 1, 2, \dots, T$, the t -th sample ξ_t is independent of the history $\hat{\mathcal{H}}_t = \left\{ (\xi_s, \hat{z}_s, \hat{t}_s, \hat{\eta}_{s+1})_{s=1}^{t-1} \right\}$.

Theorem 5. *Suppose that either all $f_t, t \in [T]$ are convex, or $\phi \equiv 0$ and $f_t \equiv f$ for some star-convex function f . Consider ASYN-CADA(ρ) running under Assumption 4 for $T > \tau_*^2$ updates, using $\gamma = 2\tau_*^2$ and any $\rho \leq \tau_*$ in `MaintainClock`. Let $\eta_0 > 0$. Then:*

(i) *If $\mathbb{E}\{\|g_t\|_2^2\} \leq G_*^2$ for all $t \in [T]$, then using a fixed $\eta_t = \eta_0\sqrt{T}$ or a time-varying $\eta_t = \eta_0\sqrt{\hat{t}_t}$,*

$$\frac{1}{T}\mathbb{E}\left\{R_T^{(f+\phi)}\right\} \leq \frac{1}{\sqrt{T}}\left(\eta_0\|x^*\|_2^2 + \frac{2(1+p_*\tau_*^2)}{\eta_0}G_*^2\right). \quad (19)$$

(ii) *If for all $\xi \in \Xi$, $F(\cdot, \xi)$ is convex and 1-smooth w.r.t. a norm $\|\cdot\|_l$, then given a constant $c_0 > 8(1+p_*\tau_*^2)$ and using a fixed $\eta_{t,i} = c_0l_i + \eta_0\sqrt{T}$ or a time-varying $\eta_{t,i} = c_0l_i + \eta_0\sqrt{\hat{t}_t}$,*

$$\frac{1}{T}\mathbb{E}\left\{R_T^{(f+\phi)}\right\} \leq \frac{c_0\|x^*\|_l^2}{T} + \frac{2}{\sqrt{T}}\left(\eta_0\|x^*\|_2^2 + \frac{4(1+p_*\tau_*^2)}{\eta_0}\sigma_*^2\right), \quad (20)$$

where $\sigma_*^2 = \mathbb{E}\{\|g(x^*, \cdot)\|_2^2\}$.

(iii) *If ϕ is μ -strongly-convex and $\mathbb{E}\{\|g_t\|_2^2\} \leq G_*^2$ for all $t \in [T]$, then using $\eta_t \equiv 0$ or, equivalently, $\mathbf{prox}(\hat{t}_t\phi, -z, 0) := \arg \min_{x \in \mathcal{X}} \hat{t}_t\phi(x) + \langle z, x \rangle = \nabla\phi^*(-z/\hat{t}_t)$,*

$$\frac{1}{T}\mathbb{E}\left\{R_T^{(f+\phi)}\right\} \leq \frac{(1+p_*\tau_*^2)G_*^2(1+\log(T))}{\mu T}, \quad (21)$$

Proof. We cast ASYN-CADA(ρ) in the Perturbed-ADA-FTRL framework of Section 4:

- (i) Thanks to the after-read time-indexing discussed above, \hat{z}_t in ASYNCADA(ρ) cannot include any coordinate updates from g_s for $s > t$ since by construction, the reading of z in t has finished before calculating of g_s is started. As such, \hat{z}_{t-1} and z_{t-1} are related to each other by (11) for all $j \in [d]$ and $t \in [T]$.
- (ii) In addition, letting $r_{0:t-1} = \hat{r}_{0:t-1} = \frac{1}{2} \|\cdot\|_{\eta_t}$, it is easy to see that the ASYNCADA(ρ) update $x_t \leftarrow \mathbf{prox}(\hat{t}_t \phi, -\hat{z}_{t-1}, \hat{\eta}_t)$ is equivalent to the perturbed ADA-FTRL update (5).
- (iii) Furthermore, by Lemma 4, \hat{t}_t is non-decreasing with, and greater than, t . Thus, η_t is also non-decreasing with t , and $\hat{t}_t, r_{0:t}$ and $\hat{r}_{0:t}$ satisfy Assumption 2 with norms $\|\cdot\|_t = \|\cdot\|_{\eta_t}$.
- (iv) Finally, Assumption 4 ensures that Assumption 3 holds.

Therefore, letting $\nu_t = \hat{t}_t - t$, applying Theorem 4, and noting that $\Delta_t = 0$ by construction, we get

$$\mathbb{E}\left\{R_T^{(f+\phi)}(x^*)\right\} \leq \mathbb{E}\left\{r_{0:T}(x^*) + \sum_{t=1}^T \frac{1 + p_* \nu_t + \sum_{s:t \in O_s} \frac{\tau_s}{\nu_s}}{2} \|g_t\|_{(t,*)}^2 - B_{1:T}\right\}.$$

Next, the assumption that either f_t is convex, or $f_t \equiv f$ for a star-convex f implies $\mathcal{B}_f(x^*, x_t) \geq 0$; hence, the $B_{1:T}$ terms can be dropped. Also, by construction, $\hat{t}_t \leq t + \gamma$, as t' always under-estimates t . Together with Lemma 4 and since $\gamma = 2\tau_*^2$, this implies $\tau_*^2 < \nu_t \leq 2\tau_*^2$. Furthermore, τ_s is the number of processes iteration s is in read-conflict with, while $\{s : t \in O_s\}$ is the set of iterations t is in write-conflict with; hence $\tau_s \leq \tau_*$ and $|\{s : t \in O_s\}| \leq \tau_*$, implying that the summation above is bounded by 1 and

$$\mathbb{E}\left\{R_T^{(f+\phi)}(x^*)\right\} \leq \mathbb{E}\left\{r_{0:T}(x^*) + \sum_{t=1}^T \frac{1 + 2p_* \tau_*^2 + 1}{2} \|g_t\|_{(t,*)}^2\right\}. \quad (22)$$

Using the definition of G_* , the fact that $r_{0:T} = \frac{\sqrt{\hat{t}_T}}{2} \|\cdot\|^2 \leq \frac{\sqrt{T+\gamma}}{2} \|\cdot\|^2 \leq \frac{\sqrt{3T}}{2} \|\cdot\|^2$, the expansion $\|\cdot\|_{(t,*)}^2 = \frac{1}{\eta_t} \|\cdot\|^2$, and the well-known bound $\sum_{t=1}^T (\sqrt{t})^{-1} \leq 2\sqrt{T}$ [21], we get (19).

To get (20), we continue from (22) but instead upper-bound $\|g_t\|$ as follows:

$$\begin{aligned} \frac{1}{2} \|g_t\|_{(t,*)}^2 &\leq \|\nabla F(x_t, \xi_t) - \nabla F(x^*, \xi_t)\|_{(t,*)}^2 + \|\nabla F(x^*, \xi_t)\|_{(t,*)}^2 \\ &\leq \frac{1}{c_0} \|\nabla F(x_t, \xi_t) - \nabla F(x^*, \xi_t)\|_{l,*}^2 + \frac{1}{\eta_0 \sqrt{t}} \|\nabla F(x^*, \xi_t)\|^2, \end{aligned}$$

where the last step follows by the definition of $\eta_{t,i}$, which in particular implies $\|\cdot\|_{(t)}^2 \geq c_0 \|\cdot\|_l^2$ and $\|\cdot\|_{(t)}^2 \geq \eta_0 \sqrt{\hat{t}_t} \|\cdot\|^2 \geq \eta_0 \sqrt{t} \|\cdot\|^2$ (similarly, in case of a fixed step size, $\|\cdot\|_{(t)}^2 \geq \eta_0 \sqrt{T} \|\cdot\|^2 \geq \eta_0 \sqrt{t} \|\cdot\|^2$). Putting back into (22), we obtain

$$\begin{aligned} \mathbb{E}\left\{R_T^{(f+\phi)}(x^*)\right\} &\leq \mathbb{E}\left\{\frac{c_0}{2} \|x^*\|_l + \frac{\eta_0 \sqrt{3T}}{2} \|x^*\|^2 + \sum_{t=1}^T \frac{2 + 2p_* \tau_*^2}{\eta_0 \sqrt{t}} \|\nabla F(x^*, \xi_t)\|^2\right\} \\ &\quad + \mathbb{E}\left\{\sum_{t=1}^T \frac{2 + 2p_* \tau_*^2}{c_0} \|\nabla F(x_t, \xi_t) - \nabla F(x^*, \xi_t)\|_{l,*}^2\right\} \\ &\leq \mathbb{E}\left\{\frac{c_0}{2} \|x^*\|_l + \frac{\eta_0 \sqrt{3T}}{2} \|x^*\|^2 + \sum_{t=1}^T \frac{2 + 2p_* \tau_*^2}{\eta_0 \sqrt{t}} \sigma_*^2\right\} \\ &\quad + \mathbb{E}\left\{\sum_{t=1}^T \frac{1}{4} \|\nabla F(x_t, \xi_t) - \nabla F(x^*, \xi_t)\|_{l,*}^2\right\} \\ &\leq \mathbb{E}\left\{\frac{c_0}{2} \|x^*\|_l + \frac{\eta_0 \sqrt{3T}}{2} \|x^*\|^2 + \frac{2(2 + 2p_* \tau_*^2)}{\eta_0} \sigma_*^2 \sqrt{T}\right\} \end{aligned}$$

$$\begin{aligned}
& + \mathbb{E} \left\{ \sum_{t=1}^T \frac{1}{2} (F(x_t, \xi_t) - F(x^*, \xi_t) - \langle \nabla F(x^*, \xi_t), x_t - x^* \rangle) \right\} \\
& = \mathbb{E} \left\{ \frac{c_0}{2} \|x^*\|_t + \frac{\eta_0 \sqrt{3T}}{2} \|x^*\|^2 + \frac{4(1 + p_* \tau_*^2)}{\eta_0} \sigma_*^2 \sqrt{T} \right\} \\
& \quad + \mathbb{E} \left\{ \sum_{t=1}^T \frac{1}{2} (f(x_t) - f(x^*) - \langle \nabla f(x^*), x_t - x^* \rangle) \right\} \\
& = \mathbb{E} \left\{ \frac{c_0}{2} \|x^*\|_t + \frac{\eta_0 \sqrt{3T}}{2} \|x^*\|^2 + \frac{4(1 + p_* \tau_*^2)}{\eta_0} \sigma_*^2 \sqrt{T} \right\} \\
& \quad + \mathbb{E} \left\{ \sum_{t=1}^T \frac{1}{2} (f(x_t) - f(x^*) - \langle \phi'(x^*), x^* - x_t \rangle) \right\} \\
& = \mathbb{E} \left\{ \frac{c_0}{2} \|x^*\|_t + \frac{\eta_0 \sqrt{3T}}{2} \|x^*\|^2 + \frac{4(1 + p_* \tau_*^2)}{\eta_0} \sigma_*^2 \sqrt{T} \right\} \\
& \quad + \mathbb{E} \left\{ \sum_{t=1}^T \frac{1}{2} (f(x_t) - f(x^*) + \phi(x_t) - \phi(x^*)) \right\}.
\end{aligned}$$

Here, the first inequality follows by the definition of $\|\cdot\|_{(t,*)}$ and the fact that $\eta_T = \eta_0 \sqrt{\hat{t}_T} \leq \eta_0 \sqrt{T + \gamma} \leq \eta_0 \sqrt{3T}$, the second inequality follows by the definition of c_0 and σ , the third follows by smoothness of F [23] and the bound $\sum_{t=1}^T (\sqrt{t})^{-1} \leq 2\sqrt{T}$ [21], the fourth by the independence of ξ_t from the history, the fifth by the optimality of x^* (where $\phi'(x^*)$ denotes the sub-gradient of ϕ for which $\phi'(x^*) + \nabla f(x^*) = 0$), and the last line follows by convexity of ϕ . Moving the last term to the l.h.s. and multiplying the sides by 2 completes the proof of (20).

To prove (21), note that $t\phi$ is $t\mu$ -strongly-convex by assumption, and thus the sequence of regularizers $r_t = \hat{r}_t = 0$ still satisfy Assumption 2 with the norms $\|\cdot\|_{(t,*)}^2 = \mu t \|\cdot\|^2$. Thus, (22) implies

$$\begin{aligned}
\mathbb{E} \left\{ R_T^{(f+\phi)}(x^*) \right\} & \leq \mathbb{E} \left\{ \sum_{t=1}^T \frac{2(1 + p_* \tau_*^2)}{2\mu t} \|g_t\|^2 \right\} \\
& \leq \frac{2(1 + p_* \tau_*^2)}{2\mu} G_*^2 (1 + \log(T)),
\end{aligned}$$

where in the last step we have used the bound $\sum_{t=1}^T (1/t) \leq 1 + \log(T)$, completing the proof. \square

D Analysis of ASYNCADA(ρ) with ADAGRAD Step-Sizes

In this section, we provide the details of the analysis of ASYNCADA(ρ) with the step-size vector $\hat{\eta}$ tuned adaptively in parallel (e.g., similar to ASYNC-ADAGRAD of Duchi et al. [10]), but with the full generality of ASYNCADA(ρ), including with the use of composite objectives ϕ and an estimated global clock.

First, recall from Section 5 that in this case, Line 10 is replaced with Line 10*: similarly to z , a vector is maintained in the shared memory, storing the sum of squares of the observed gradient values for each coordinate individually. Also, this vector is read and updated in the same way as z (with a square-root applied to each coordinate to recover the vector $\hat{\eta}$ before passing it to the **prox** operator in Line 7). Finally, note that similarly to \hat{z}_{t-1} , $\eta_t^{(j)}$ will lack some of the updates from the iterations happening concurrently with the t -th iteration. The after-read time-indexing defined in Appendix B ensures that these concurrent updates have indices less than t , that is, there exists a set $D_{t,j} \subset [t-1]$ such that for $s \in D_{t,j}$, $\alpha^2 \left(g_s^{(j)}\right)^2$ has not yet been added to the shared memory when $\hat{\eta}_t^{(j)}$ is read.

Formally, for all $j \in [d]$ and all $t \in [T]$, we have

$$\hat{\eta}_t^{(j)} = \alpha \sqrt{\lambda_j + \sum_{s=1}^{t-1} \left(g_s^{(j)}\right)^2 - \sum_{s \in D_{t,j}} \left(g_s^{(j)}\right)^2}, \quad (23)$$

for some hyper-parameter $\lambda_j > 0$ that, as usual in the dual-averaging formulations of ADAGRAD (cf. McMahan [21]), ensures the initial $\hat{\eta}_1$ is well-defined. Then, we have the following theorem:

Theorem 6. *Suppose that either all $f_t, t \in [T]$ are convex, or $\phi \equiv 0$ and $f_t \equiv f$ for some star-convex function f . Consider ASYNCADA(ρ) running under Assumption 4 for $T \geq 1$ updates, using $\gamma = 2\tau_*^2$ and any $\rho \leq \tau_*$ in MaintainClock. For all $j \in [d]$, assume that $|g_t^{(j)}| \leq G_j, t \in [T]$, and $|x^{(j)}| \leq R$ for all $x \in \mathcal{X}$. Then, using ADAGRAD step-sizes given by (23) with $\lambda_j = (\tau_* + 1)G_j^2$ and any $\alpha > 0$, we have*

$$\mathbb{E}\left\{R_T^{(f+\phi)}\right\} \leq \frac{\alpha R^2 \sqrt{\tau_* + 1}}{2} \sum_{j=1}^d G_j + \left(\alpha R^2 + \frac{2 + 2\rho_* \tau_*^2}{\alpha}\right) \sqrt{T} \sum_{j=1}^d G_j \sqrt{p_j}. \quad (24)$$

Remark 5. For a fixed τ_* , the first term in (24) is of lower order. Thus, if $p_* \tau_*^2 \leq c$ for a constant c and we tune α based on R , Theorem 6 implies that the average iterate \bar{x}_T of ASYNCADA(ρ) with ADAGRAD step-sizes satisfies

$$\mathbb{E}\{f(\bar{x}_T) - f(x^*)\} = \mathcal{O}\left(R \sum_{j=1}^d G_j \sqrt{p_j/T}\right),$$

which, up to a constant factor, is the same worst-case convergence rate as that of serial composite-objective ADAGRAD under sparsity (see, e.g., Duchi et al. [9, Corollary 1]). In addition, when $\phi = 0$, this is the best rate attainable by any serial algorithm under sparsity [10, Proposition 1].¹⁰ Thus, Theorem 6 generalizes the linear speed-up result of ASYNC-ADAGRAD to non-box-shaped \mathcal{X} , star-convex objectives, and proximal updates using a convex ϕ with an inexact global clock. It should be noted, however, that the speed-up regime of ASYNC-ADAGRAD with box-shaped constraints is less restrictive than that of Theorem 6, i.e., the condition required by Duchi et al. [10, Eq. (8)] is $p_* \tau_* \leq c$ rather than $p_* \tau_*^2 \leq c$.

Proof. The proof follows as in the case of Theorem 5, in particular noting the implications of the after-read indexing, the fact that Assumption 4 implies Assumption 3, and the equivalence of the **prox** operator to the Perturbed-ADA-FTRL update (5) given a Euclidean regularizer. However, here we use different $r_{0:t}$ and $\hat{r}_{0:t}$, and hence Δ_t does not necessarily vanish anymore. In particular, we let $\tilde{\eta}_t$ to be given by

$$\tilde{\eta}_t^{(j)} = \alpha \sqrt{\lambda_j + \sum_{s=1}^{t-1} \left(g_s^{(j)}\right)^2}$$

and let the idealized $r_{0:t-1}$ be given by

$$r_{0:t-1}(x) = \frac{1}{2} \|x\|_{\tilde{\eta}_t}^2 = \sum_{j=1}^d \frac{\alpha}{2} \sqrt{\lambda_j + \sum_{s=1}^{t-1} \left(g_s^{(j)}\right)^2} \left(x^{(j)}\right)^2,$$

whereas the actual regularizer used by ASYNCADA(ρ) is given by (23), i.e.,

$$\hat{r}_{0:t-1}(x) = \frac{1}{2} \|x\|_{\hat{\eta}_t}^2 = \sum_{j=1}^d \frac{\alpha}{2} \sqrt{\lambda_j + \sum_{s=1}^{t-1} \left(g_s^{(j)}\right)^2 - \sum_{s \in D_{t,j}} \left(g_s^{(j)}\right)^2} \left(x^{(j)}\right)^2.$$

¹⁰Proposition 1 of Duchi et al. [10] establishes, for $\phi = 0$, a lower-bound of $\frac{1}{8} R \sum_{j=1}^d G_j \min\{p_j, \sqrt{p_j/T}\}$. By Theorem 6, ASYNCADA(ρ) with ADAGRAD step-sizes matches this rate (when $\phi = 0$), since in addition to (24), by the sparsity of the gradients and star-convexity of f , we have for any $x \in \mathcal{X}$ and ξ sampled from \mathbb{P}_Ξ that

$$f(x) - f(x^*) \leq \langle \nabla f(x), x - x^* \rangle = \mathbb{E}\{\langle \nabla F(x, \xi), x - x^* \rangle\} \leq \sum_{j \in [d]} G_j p_j R.$$

Note that for all $t \in [T]$, this ensures that $\tilde{\eta}_t^{(j)} \geq \hat{\eta}_t^{(j)} \geq \eta_t^{(j)} := \alpha \sqrt{\sum_{s=1}^t (g_s^{(j)})^2}$ for all $j \in [d]$, since by assumption, $\lambda_j = (\tau_* + 1)G_j^2 \geq \sum_{s \in D_{t,j} \cup \{t\}} (g_s^{(j)})^2 \geq \sum_{s=t-\tau_*}^t (g_s^{(j)})^2$. Hence, both $r_{0:t-1}$ and $\hat{r}_{0:t-1}$ are strongly-convex w.r.t. the norm $\frac{1}{2} \|\cdot\|_{\eta_t}$, and these norms are non-decreasing (since η_t is non-decreasing by definition). Combined with Lemma 4 and the fact that $r_t \geq 0$ as defined above, this implies that Assumption 2 is satisfied. Hence, we can apply Theorem 4 with $\nu_t = \hat{t}_t - t$, recalling that by assumption of (star-)convexity of f_t , we have $B_{1:T} \geq 0$, and obtain

$$\begin{aligned}
\mathbb{E}\left\{R_T^{(f+\phi)}\right\} &\leq \mathbb{E}\left\{\frac{1}{2}\|x^*\|_{\tilde{\eta}_T}^2 + \sum_{t=1}^T \left(\frac{1 + p_*\nu_t + \sum_{s:t \in O_s} \frac{\tau_s}{\nu_s} \|g_t\|_{(t,*)}^2 + \frac{\Delta_t}{\nu_t}\right)\right\} \\
&\leq \mathbb{E}\left\{\frac{1}{2}\|x^*\|_{\tilde{\eta}_T}^2 + \sum_{t=1}^T \left(\frac{2 + 2p_*\tau_*^2}{2} \|g_t\|_{(t,*)}^2 + \frac{\Delta_t}{\nu_t}\right)\right\} \\
&\leq \mathbb{E}\left\{\sum_{j=1}^d \frac{\alpha R^2}{2} \sqrt{\lambda_j + \sum_{t=1}^T (g_t^{(j)})^2} + \sum_{t=1}^T \frac{\Delta_t}{\nu_t}\right\} \\
&\quad + \mathbb{E}\left\{(2 + 2p_*\tau_*^2) \sum_{t=1}^T \sum_{j=1}^d \frac{1}{2\alpha \sqrt{\sum_{s=1}^t (g_s^{(j)})^2}} (g_t^{(j)})^2\right\} \\
&\leq \mathbb{E}\left\{\sum_{j=1}^d \frac{\alpha R^2}{2} \sqrt{\lambda_j + \sum_{t=1}^T (g_t^{(j)})^2} + \sum_{t=1}^T \frac{\Delta_t}{\nu_t}\right\} \\
&\quad + \mathbb{E}\left\{\frac{2 + 2p_*\tau_*^2}{\alpha} \sum_{j=1}^d \sqrt{\sum_{t=1}^T (g_t^{(j)})^2}\right\} \\
&\leq \frac{\alpha R^2}{2} \sum_{j=1}^d \sqrt{\lambda_j} + \sum_{j=1}^d \frac{\alpha R^2}{2} \sqrt{\sum_{t=1}^T \mathbb{E}\left\{(g_t^{(j)})^2\right\}} + \sum_{t=1}^T \mathbb{E}\left\{\frac{\Delta_t}{\nu_t}\right\} \\
&\quad + \frac{2 + 2p_*\tau_*^2}{\alpha} \sum_{j=1}^d \sqrt{\sum_{t=1}^T \mathbb{E}\left\{(g_t^{(j)})^2\right\}}, \tag{25}
\end{aligned}$$

where the second line follows as (22) in the proof of Theorem 5, the third line uses the definition of $\tilde{\eta}_T$ and η_t and the upper-bound R , the fourth line uses the ADAGRAD lemma (McMahan [21, Lemma 4]), and the last line uses concavity of square-root and Jensen's inequality.

Next, we bound the terms Δ_t . Recall that by definition, $\tilde{\eta}_t^{(j)} \geq \hat{\eta}_t^{(j)}$. Hence, $r_{0:t-1} \geq \hat{r}_{0:t-1}$, and

$$\begin{aligned}
\Delta_t &= r_{0:t-1}(x_t) - r_{0:t-1}(\tilde{x}_t) + \hat{r}_{0:t-1}(\tilde{x}_t) - \hat{r}_{0:t-1}(x_t) \leq r_{0:t-1}(x_t) - \hat{r}_{0:t-1}(x_t) \\
&= \sum_{j=1}^d \frac{\tilde{\eta}_t^{(j)} - \hat{\eta}_t^{(j)}}{2} (x_t^{(j)})^2 \\
&\leq \frac{\alpha R^2}{2} \sum_{j=1}^d \left(\sqrt{\lambda_j + \sum_{s=1}^{t-1} (g_s^{(j)})^2} - \sqrt{\lambda_j + \sum_{s=1}^{t-1} (g_s^{(j)})^2 - \sum_{s \in D_{t,j}} (g_s^{(j)})^2} \right) \\
&\leq \sum_{j=1}^d \frac{\alpha R^2 \sum_{s \in D_{t,j}} (g_s^{(j)})^2}{4 \sqrt{\lambda_j + \sum_{s=1}^{t-1} (g_s^{(j)})^2 - \sum_{s \in D_{t,j}} (g_s^{(j)})^2}}
\end{aligned}$$

$$\leq \sum_{j=1}^d \frac{\alpha R^2 \sum_{s \in D_{t,j}} (g_s^{(j)})^2}{4 \sqrt{\sum_{s=1}^t (g_s^{(j)})^2}},$$

where the third line uses the upper-bound R on $|x_t^{(j)}|$, the fourth line uses the inequality $\sqrt{a+b} - \sqrt{a} \leq \frac{b}{2\sqrt{a}}$ which holds for all $a, b > 0$, and the last line uses $\eta_t^{(j)} \leq \hat{\eta}_t^{(j)}$. Noting that by construction, either $\tau_* \geq 1$ (and thus, $\tau_* \leq \tau_*^2 < \nu_t$) or $\Delta_t = 0$, we have

$$\begin{aligned} \sum_{t=1}^T \frac{\Delta_t}{\nu_t} &\leq \frac{\alpha R^2}{4} \sum_{j=1}^d \sum_{t=1}^T \sum_{s \in D_{t,j}} \frac{(g_s^{(j)})^2}{\tau_* \sqrt{\sum_{k=1}^t (g_k^{(j)})^2}} \\ &\leq \frac{\alpha R^2}{4} \sum_{j=1}^d \sum_{t=1}^T \sum_{s \in D_{t,j}} \frac{(g_s^{(j)})^2}{\tau_* \sqrt{\sum_{k=1}^s (g_k^{(j)})^2}} \\ &\leq \frac{\alpha R^2}{4} \sum_{j=1}^d \sum_{s=1}^T \sum_{t: s \in D_{t,j}} \frac{(g_s^{(j)})^2}{\tau_* \sqrt{\sum_{k=1}^s (g_k^{(j)})^2}} \\ &\leq \frac{\alpha R^2}{2} \sum_{j=1}^d \sum_{s=1}^T \frac{(g_s^{(j)})^2}{2 \sqrt{\sum_{k=1}^s (g_k^{(j)})^2}} \\ &\leq \frac{\alpha R^2}{2} \sum_{j=1}^d \sqrt{\sum_{t=1}^T (g_t^{(j)})^2}, \end{aligned}$$

where the second line uses the fact that $s < t$, the third line swaps the summations on t and s , the fourth line uses the fact that $|\{t : s \in D_{t,j}\}| \leq \tau_*$, and the fifth line uses the ADAGRAD lemma (McMahan [21, Lemma 4]). Putting back into (25) and using the concavity of square-root with Jensen's inequality again, we obtain

$$\begin{aligned} \mathbb{E}\{R_T^{(f+\phi)}\} &\leq \frac{\alpha R^2}{2} \sum_{j=1}^d \sqrt{\lambda_j} + \sum_{j=1}^d \alpha R^2 \sqrt{\sum_{t=1}^T \mathbb{E}\{(g_t^{(j)})^2\}} \\ &\quad + \frac{2 + 2p_*\tau_*^2}{\alpha} \sum_{j=1}^d \sqrt{\sum_{t=1}^T \mathbb{E}\{(g_t^{(j)})^2\}}. \end{aligned}$$

Using the fact that $\mathbb{E}\{(g_t^{(j)})^2\} \leq p_{t,j} G_j^2$, the definition of λ_j , and simplifying, we obtain

$$\mathbb{E}\{R_T^{(f+\phi)}\} \leq \frac{\alpha R^2 \sqrt{\tau_* + 1}}{2} \sum_{j=1}^d G_j + \left(\alpha R^2 + \frac{2 + 2p_*\tau_*^2}{\alpha} \right) \sqrt{T} \sum_{j=1}^d G_j \sqrt{p_j},$$

This completes the proof. \square

E Proofs for HEDGEHOG

Proof of Theorem 3. As in the case of ASYNCADEA(ρ), the proof follows by casting HEDGEHOG as Perturbed-ADA-FTRL. In particular, the same relation between \hat{z}_{t-1} and z_{t-1} holds, and it is easy to see that with the after-read time-indexing the HEDGEHOG update corresponds to Perturbed-ADA-FTRL with the regularizer $r_{0:t}(x) = r(x) + \eta \ln(d)$ where $r(x) = \eta \sum_{i=1}^d x^{(i)} \log(x^{(i)})$, which

is 1-strongly-convex w.r.t. the ℓ_1 norm [34]. Note also that we can assume any value for $\hat{t}_t > t$, including $\hat{t}_t = t + \nu_t$ for any $\nu_t > 0$, as we don't use \hat{t}_t in the update and hence don't need to be able to compute it. Then, Assumption 2 is satisfied with $\|\cdot\|_t = \sum_{i=1}^d |x^{(i)}|$ being the ℓ_1 norm, with $\|\cdot\|_{(t,*)} = \|\cdot\|_\infty$. Then, applying Theorem 4 and noting $\phi = 0, \Delta_t = 0, B_{1:T} \geq 0$, we have

$$\mathbb{E}\left\{R_T^{(f)}(x^*)\right\} \leq r(x^*) + \eta \log(d) + \sum_{t=1}^T \mathbb{E}\left\{\frac{1 + p_* \nu_t + \sum_{s:t \in O_s} \frac{\tau_s}{\nu_s} \|g_t\|_{(t,*)}^2}{2}\right\},$$

for any ν_t determined by $\hat{\mathcal{H}}_t$. In particular, letting $\nu_t = \tau_*/\sqrt{p_*}$, recalling that τ_s and $|\{s : t \in O_s\}|$ cannot be larger than τ_* , and noting that $r(x^*) \leq 0$ for any $x^* \in \mathcal{X}$ completes the proof. \square

F Extra details for the analysis of serial ADA-FTRL

A typical proxy for bounding the regret of serial optimization algorithms is *linearizing* the loss, and studying the linearized regret [5, 34, 14]. In particular, we define the linearized forward regret

$$R_T^+(x^*) = \sum_{t=1}^T \langle g_t, x_{t+1} - x^* \rangle,$$

and carry out the analysis in two steps: a decomposition of the regret in terms of the forward regret R_T^+ , followed by a bound on R_T^+ . To that end, we need the following assumption.

Definition 2 (Admissible regularizers.). A sequence of regularizer functions $(r_t)_{t=0}^T$ is “admissible” for ADA-FTRL if and only if all r_t are defined on a common convex domain $S \subset \mathbb{R}^d$, the intersection $\mathcal{X} \cap S$ is non-empty, and there exists a sequence of norms $(\|\cdot\|_{(t)})_{t=1}^T$ such that for all $t = 1, 2, \dots, T$, the cumulative regularizer $t\phi + r_{0:t-1} : S \rightarrow \mathbb{R}$ is lower-semi-continuous and 1-strongly-convex w.r.t. $\|\cdot\|_{(t)}$.

As shown by Lemma 5, admissible regularizers guarantee that the ADA-FTRL updates (2) are well-defined, that is, there exists some $x_{t+1} \in \mathcal{X}$ that satisfies (2), and the associated optimal value is finite.

Lemma 5 (Well-posed ADA-FTRL). *For all $t = 0, 1, \dots, T$, the argmin sets that define x_{t+1} in the ADA-FTRL updates (2) are non-empty, and their optimal values are finite.*

Proof. Fix $t \in [T]$, and consider the extended-value function $h_t = \langle z_{t-1}, \cdot \rangle + r_{0:t-1} + \mathcal{I}_{x \in S \cap \mathcal{X}}$, which is proper, l.s.c. and convex by construction. In addition, since $r_{0:t-1}$ is strongly-convex over S , then h_t is l.s.c. and 1-strong-convex on \mathbb{R}^d . The result then follows by Proposition 17.26 of [1], noting that x_t will be the corresponding minimizer by definition. \square

Then, we have the following bound on the regret of ADA-FTRL.

Theorem 7 (Forward regret of ADA-FTRL, [16]). *For any $x^* \in \mathcal{X}$ and for any sequence of linear losses $\langle g_t, \cdot \rangle, t = 1, 2, \dots, T$, and using any sequence of admissible regularizers r_0, r_1, \dots, r_T , the forward regret of ADA-FTRL satisfies*

$$R_T^+(x^*) \leq r_{0:T}(x^*) - \sum_{t=0}^T r_t(x_{t+1}) + \sum_{t=1}^T (\phi(x^*) - \phi(x_t)) - \sum_{t=1}^T \mathcal{B}_{r_{0:t-1}}(x_{t+1}, x_t). \quad (26)$$

Theorem 1 follows as a direct consequence of the above theorem by using the strong convexity of $r_{0:t-1}$ and the Fenchel-Young inequality; see [16] for further details.