
ASAGA: Asynchronous Parallel SAGA

Rémi Leblond

INRIA - Sierra Project-team
École normale supérieure, Paris

Fabian Pedregosa

INRIA - Sierra Project-team
École normale supérieure, Paris

Simon Lacoste-Julien

Department of CS & OR (DIRO)
Université de Montréal, Montréal

Abstract

We describe ASAGA, an asynchronous parallel version of the incremental gradient algorithm SAGA that enjoys fast linear convergence rates. Through a novel perspective, we revisit and clarify a subtle but important technical issue present in a large fraction of the recent convergence rate proofs for asynchronous parallel optimization algorithms, and propose a simplification of the recently introduced “perturbed iterate” framework that resolves it. We thereby prove that ASAGA can obtain a theoretical linear speedup on multi-core systems even without sparsity assumptions. We present results of an implementation on a 40-core architecture illustrating the practical speedup as well as the hardware overhead.

1 Introduction

We consider the unconstrained optimization problem of minimizing a *finite sum* of smooth convex functions:

$$\min_{x \in \mathbb{R}^d} f(x), \quad f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x), \quad (1)$$

where each f_i is assumed to be convex with L -Lipschitz continuous gradient, f is μ -strongly convex and n is large (for example, the number of data points in a regularized empirical risk minimization setting). We define a condition number for this problem as $\kappa := L/\mu$. A flurry of randomized incremental algorithms (which at each iteration select i at random and process only one gradient f'_i) have recently been proposed to solve (1) with a fast¹ linear convergence rate, such as SAG (Le

¹Their complexity in terms of gradient evaluations to reach an accuracy of ϵ is $O((n + \kappa) \log(1/\epsilon))$, in contrast to $O(n\kappa \log(1/\epsilon))$ for batch gradient descent in the worst case.

Roux et al., 2012), SDCA (Shalev-Shwartz and Zhang, 2013), SVRG (Johnson and Zhang, 2013) and SAGA (Defazio et al., 2014). These algorithms can be interpreted as variance reduced versions of the stochastic gradient descent (SGD) algorithm, and they have demonstrated both theoretical and practical improvements over SGD (for the *finite sum* optimization problem (1)).

In order to take advantage of the multi-core architecture of modern computers, the aforementioned optimization algorithms need to be adapted to the asynchronous parallel setting, where multiple threads work concurrently. Much work has been devoted recently in proposing and analyzing asynchronous parallel variants of algorithms such as SGD (Niu et al., 2011), SDCA (Hsieh et al., 2015) and SVRG (Reddi et al., 2015; Mania et al., 2015; Zhao and Li, 2016). Among the incremental gradient algorithms with fast linear convergence rates that can optimize (1) in its general form, only SVRG has had an asynchronous parallel version proposed.² No such adaptation has been attempted yet for SAGA, even though one could argue that it is a more natural candidate as, contrarily to SVRG, it is not epoch-based and thus has no synchronization barriers at all.

Contributions. In Section 2, we present a novel sparse variant of SAGA that is more adapted to the parallel setting than the original SAGA algorithm. In Section 3, we present ASAGA, a lock-free asynchronous parallel version of Sparse SAGA that does not require consistent reads. We propose a simplification of the “perturbed iterate” framework from Mania et al. (2015) as a basis for our convergence analysis. At the same time, through a novel perspective, we revisit and clarify a technical problem present in a large fraction of the literature on randomized asynchronous parallel algorithms (with the exception of Mania et al. (2015), which also highlights this issue): namely, they all assume unbiased gradient estimates, an assumption that is inconsistent with their proof technique without fur-

²We note that SDCA requires the knowledge of an explicit μ -strongly convex regularizer in (1), whereas SAG / SAGA are adaptive to any local strong convexity of f (Schmidt et al., 2016; Defazio et al., 2014). This is also true for a variant of SVRG (Hofmann et al., 2015).

ther synchronization assumptions. In Section 3.3, we present a tailored convergence analysis for ASAGA. Our main result states that ASAGA obtains the same geometric convergence rate per update as SAGA when the overlap bound τ (which scales with the number of cores) satisfies $\tau \leq \mathcal{O}(n)$ and $\tau \leq \mathcal{O}(\frac{1}{\sqrt{\Delta}} \max\{1, \frac{n}{\kappa}\})$, where $\Delta \leq 1$ is a measure of the sparsity of the problem, notably implying that a linear speedup is theoretically possible even without sparsity in the well-conditioned regime where $n \gg \kappa$. In Section 4, we provide a practical implementation of ASAGA and illustrate its performance on a 40-core architecture, showing improvements compared to asynchronous variants of SVRG and SGD.

Related Work. The seminal textbook of Bertsekas and Tsitsiklis (1989) provides most of the foundational work for parallel and distributed optimization algorithms. An asynchronous variant of SGD with constant step size called HOGWILD was presented by Niu et al. (2011); part of their framework of analysis was re-used and inspired most of the recent literature on asynchronous parallel optimization algorithms with convergence rates, including asynchronous variants of coordinate descent (Liu et al., 2015), SDCA (Hsieh et al., 2015), SGD for non-convex problems (De Sa et al., 2015; Lian et al., 2015), SGD for stochastic optimization (Duchi et al., 2015) and SVRG (Reddi et al., 2015; Zhao and Li, 2016). These papers make use of an unbiased gradient assumption that is not consistent with the proof technique, and thus suffers from technical problems³ that we highlight in Section 3.2.

The “perturbed iterate” framework presented in Mania et al. (2015) is to the best of our knowledge the only one that does not suffer from this problem, and our convergence analysis builds heavily from their approach, while simplifying it. In particular, the authors assumed that f was both strongly convex and had a bound on the gradient, two *inconsistent* assumptions in the unconstrained setting that they analyzed. We overcome these difficulties by using tighter inequalities that remove the requirement of a bound on the gradient. We also propose a more convenient way to label the iterates (see Section 3.2). The sparse version of SAGA that we propose is also inspired from the sparse version of SVRG proposed by Mania et al. (2015). Reddi et al. (2015) presents a hybrid algorithm called HSAG that includes SAGA and SVRG as special cases. Their asynchronous analysis is epoch-based though, and thus does not handle a fully asynchronous version of SAGA as we do. Moreover, they require consistent reads and do not propose an efficient sparse implementation for SAGA, in contrast to ASAGA.

³Except Duchi et al. (2015) that can be easily fixed by incrementing their global counter *before* sampling.

Notation. We denote by \mathbb{E} a full expectation with respect to all the randomness, and by \mathbf{E} the *conditional* expectation of a random i (the index of the factor f_i chosen in SGD-like algorithms), conditioned on all the past, where “past” will be clear from the context. $[x]_v$ is the coordinate v of the vector $x \in \mathbb{R}^d$. x^+ represents the updated parameter vector after one algorithm iteration.

2 Sparse SAGA

Borrowing our notation from Hofmann et al. (2015), we first present the original SAGA algorithm and then describe a novel sparse variant that is more appropriate for a parallel implementation.

Original SAGA Algorithm. The standard SAGA algorithm (Defazio et al., 2014) maintains two moving quantities to optimize (1): the current iterate x and a table (memory) of historical gradients $(\alpha_i)_{i=1}^n$.⁴ At every iteration, the SAGA algorithm samples uniformly at random an index $i \in \{1, \dots, n\}$, and then executes the following update on x and α (for the unconstrained optimization version):

$$x^+ = x - \gamma(f'_i(x) - \alpha_i + \bar{\alpha}); \quad \alpha_i^+ = f'_i(x), \quad (2)$$

where γ is the step size and $\bar{\alpha} := 1/n \sum_{i=1}^n \alpha_i$ can be updated efficiently in an online fashion. Crucially, $\mathbf{E}\alpha_i = \bar{\alpha}$ and thus the update direction is unbiased ($\mathbf{E}x^+ = x - \gamma f'(x)$). Furthermore, it can be proven (see Defazio et al. (2014)) that under a reasonable condition on γ , the update has vanishing variance, which enables the algorithm to converge linearly with a constant step size.

Motivation for a Variant. In its current form, every SAGA update is dense even if the individual gradients are sparse due to the historical gradient ($\bar{\alpha}$) term. Schmidt et al. (2016) introduced a special implementation with lagged updates where every iteration has a cost proportional to the size of the support of $f'_i(x)$. However, this subtle technique is not easily adaptable to the parallel setting (see App. F.2). We therefore introduce Sparse SAGA, a novel variant which explicitly takes sparsity into account and is easily parallelizable.

Sparse SAGA Algorithm. As in the Sparse SVRG algorithm proposed in Mania et al. (2015), we obtain Sparse SAGA by a simple modification of the parameter update rule in (2) where $\bar{\alpha}$ is replaced by a sparse version equivalent in expectation:

$$x^+ = x - \gamma(f'_i(x) - \alpha_i + D_i \bar{\alpha}), \quad (3)$$

where D_i is a diagonal matrix that makes a weighted projection on the support of f'_i . More precisely, let S_i

⁴For linear predictor models, the memory α_i^0 can be stored as a scalar. Following Hofmann et al. (2015), α_i^0 can be initialized to any convenient value (typically 0), unlike the prescribed $f'_i(x_0)$ analyzed in Defazio et al. (2014).

be the support of the gradient f'_i function (i.e., the set of coordinates where f'_i can be nonzero). Let D be a $d \times d$ diagonal reweighting matrix, with coefficients $1/p_v$ on the diagonal, where p_v is the probability that dimension v belongs to S_i when i is sampled uniformly at random in $\{1, \dots, n\}$. We then define $D_i := P_{S_i} D$, where P_{S_i} is the projection onto S_i . The normalization from D ensures that $\mathbf{E} D_i \bar{\alpha} = \bar{\alpha}$, and thus that the update is still unbiased despite the projection.

Convergence Result for (Serial) Sparse SAGA.

For clarity of exposition, we model our convergence result after the simple form of Hofmann et al. (2015, Corollary 3) (note that the rate for Sparse SAGA is the same as SAGA). The proof is given in Appendix B.

Theorem 1. *Let $\gamma = \frac{a}{5L}$ for any $a \leq 1$. Then Sparse SAGA converges geometrically in expectation with a rate factor of at least $\rho(a) = \frac{1}{5} \min\{\frac{1}{n}, a\frac{1}{\kappa}\}$, i.e., for x_t obtained after t updates, we have $\mathbb{E}\|x_t - x^*\|^2 \leq (1 - \rho)^t C_0$, where $C_0 := \|x_0 - x^*\|^2 + \frac{1}{5L^2} \sum_{i=1}^n \|\alpha_i^0 - f'_i(x^*)\|^2$.*

Comparison with Lagged Updates. The lagged updates technique in SAGA is based on the observation that the updates for component $[x]_v$ can be delayed until this coefficient is next accessed. Interestingly, the expected number of iterations between two steps where a given dimension v is involved in the partial gradient is p_v^{-1} , where p_v is the probability that v is involved. p_v^{-1} is precisely the term which we use to multiply the update to $[x]_v$ in Sparse SAGA. Therefore one may view the Sparse SAGA updates as *anticipated* SAGA updates, whereas those in the Schmidt et al. (2016) implementation are *lagged*.

Although Sparse SAGA requires the computation of the p_v probabilities, this can be done during a first pass through the data (during which constant step size SGD may be used) at a negligible cost. In our experiments, both Sparse SAGA and SAGA with lagged updates had similar convergence in terms of number of iterations, with the Sparse SAGA scheme being slightly faster in terms of runtime. We refer the reader to Schmidt et al. (2016) and Appendix F for more details.

3 Asynchronous Parallel Sparse SAGA

As most recent parallel optimization contributions, we use a similar hardware model to Niu et al. (2011). We have multiple cores which all have read and write access to a shared memory. They update a central parameter vector in an asynchronous and lock-free fashion. Unlike Niu et al. (2011), we *do not* assume that the vector reads are consistent: multiple cores can read and write different coordinates of the shared vector at the same time. This means that a full vector read for a core might not correspond to any consistent

state in the shared memory at any specific point in time.

3.1 Perturbed Iterate Framework

We first review the “perturbed iterate” framework recently introduced by Mania et al. (2015) which will form the basis of our analysis. In the sequential setting, stochastic gradient descent and its variants can be characterized by the following update rule:

$$x_{t+1} = x_t - \gamma g(x_t, i_t), \quad (4)$$

where i_t is a random variable independent from x_t and we have the unbiasedness condition $\mathbf{E}g(x_t, i_t) = f'(x_t)$ (recall that \mathbf{E} is the relevant-past conditional expectation with respect to i_t).

Unfortunately, in the parallel setting, we manipulate stale, inconsistent reads of shared parameters and thus we do not have such a straightforward relationship. Instead, Mania et al. (2015) proposed to separate \hat{x}_t , the actual value read by a core to compute an update, with x_t , a “virtual iterate” that we can analyze and is *defined* by the update equation: $x_{t+1} := x_t - \gamma g(\hat{x}_t, i_t)$. We can thus interpret \hat{x}_t as a noisy (perturbed) version of x_t due to the effect of asynchrony. In the specific case of (Sparse) SAGA, we have to add the additional read memory argument $\hat{\alpha}^t$ to our update:

$$\begin{aligned} x_{t+1} &:= x_t - \gamma g(\hat{x}_t, \hat{\alpha}^t, i_t); \\ g(\hat{x}_t, \hat{\alpha}^t, i_t) &:= f'_{i_t}(\hat{x}_t) - \hat{\alpha}^t_{i_t} + D_{i_t} (1/n \sum_{i=1}^n \hat{\alpha}^t_i). \end{aligned} \quad (5)$$

We formalize the precise meaning of x_t and \hat{x}_t in the next section. We first note that all the papers mentioned in the related work section that analyzed asynchronous parallel randomized algorithms assumed that the following unbiasedness condition holds:

$$\left[\begin{array}{l} \text{unbiasedness} \\ \text{condition} \end{array} \right] \quad \mathbf{E}[g(\hat{x}_t, i_t) | \hat{x}_t] = f'(\hat{x}_t). \quad (6)$$

This condition is at the heart of most convergence proofs for randomized optimization methods.⁵ Mania et al. (2015) correctly pointed out that most of the literature thus made the often implicit assumption that i_t is independent of \hat{x}_t . But as we explain below, this assumption is incompatible with a non-uniform asynchronous model in the analysis approach used in most of the recent literature.

3.2 On the Difficulty of Labeling the Iterates

Formalizing the meaning of x_t and \hat{x}_t highlights a subtle but important difficulty arising when analyzing

⁵A notable exception is SAG (Le Roux et al., 2012) which has biased updates, yielding a significantly more complex convergence proof. Making SAG unbiased leads to SAGA (Defazio et al., 2014) and a much simpler proof.

randomized parallel algorithms: what is the meaning of t ? This is the problem of *labeling* the iterates for the purpose of the analysis, and this labeling can have randomness itself that needs to be taken in consideration when interpreting the meaning of an expression like $\mathbb{E}[x_t]$. In this section, we contrast three different approaches in a unified framework. We notably clarify the dependency issues that the labeling from Mania et al. (2015) resolves and propose a new, simpler labeling which allows for much simpler proof techniques. We consider algorithms that execute in parallel the following four steps, where t is a global labeling that needs to be defined:

1. Read the information in shared memory (\hat{x}_t).
2. Sample i_t .
3. Perform some computations using (\hat{x}_t, i_t) .
4. Write an update to shared memory.

The “After Write” Approach. We call the “after write” approach the standard global labeling scheme used in Niu et al. (2011) and re-used in all the later papers that we mentioned in the related work section, with the notable exceptions of Mania et al. (2015) and Duchi et al. (2015). In this approach, t is a (virtual) global counter recording the number of *successful writes* to the shared memory x (incremented after step 4 in (7)); x_t thus represents the (true) content of the shared memory after t updates. The interpretation of the crucial equation (5) then means that \hat{x}_t represents the (delayed) local copy value of the core that made the $(t+1)^{\text{th}}$ successful update; i_t represents the factor sampled by this core for this update. Notice that in this framework, the value of \hat{x}_t and i_t is unknown at “time t ”; we have to wait to the later time when the next core writes to memory to finally determine that its local variables are the ones labeled by t . We thus see that here \hat{x}_t and i_t are not necessarily independent – they share dependence through the t label assignment. In particular, if some values of i_t yield faster updates than others, it will influence the label assignment defining \hat{x}_t . We illustrate this point with a concrete problematic example in Appendix A that shows that in order to preserve the unbiasedness condition (6), the “after write” framework makes the implicit assumption that the computation time for the algorithm running on a core is independent of the sample i chosen. This assumption seems overly strong in the context of potentially heterogeneous factors f_i ’s, and is thus a fundamental flaw for analyzing non-uniform asynchronous computation.

The “Before Read” Approach. Mania et al. (2015) addresses this issue by proposing instead to increment the global t counter just *before* a new core starts to *read* the shared memory (before step 1 in (7)). In their framework, \hat{x}_t represents the (inconsistent) read that was made by this core in this computational

block, and i_t represents the picked sample. The update rule (5) represents a *definition* of the meaning of x_t , which is now a “virtual iterate” as it does not necessarily correspond to the content of the shared memory at any point. The real quantities manipulated by the algorithm in this approach are the \hat{x}_t ’s, whereas x_t is used only for the analysis – the critical quantity we want to see vanish is $\mathbb{E}\|\hat{x}_t - x^*\|^2$. The independence of i_t with \hat{x}_t can be simply enforced in this approach by making sure that the way the shared memory x is read does not depend on i_t (e.g. by reading all its coordinates in a fixed order). Note that this means that we have to read all of x ’s coordinates, regardless of the size of f_{i_t} ’s support. This is a much weaker condition than the assumption that all the computation in a block does not depend on i_t as required by the “after write” approach, and is thus more reasonable.

A New Global Ordering: the “After Read” Approach. The “before read” approach gives rise to the following complication in the analysis: \hat{x}_t can depend on i_r for $r > t$. This is because t is a global time ordering only on the assignment of computation to a core, not on when \hat{x}_t was finished to be read. This means that we need to consider both the “future” and the “past” when analyzing x_t . To simplify the analysis (which proved crucial for our ASAGA proof), we thus propose a third way to label the iterates: \hat{x}_t represents the $(t+1)^{\text{th}}$ *fully completed read* (t incremented after step 1 in (7)). As in the “before read” approach, we can ensure that i_t is independent of \hat{x}_t by ensuring that how we read does not depend on i_t . But unlike in the “before read” approach, t here now does represent a global ordering on the \hat{x}_t iterates – and thus we have that i_r is independent of \hat{x}_t for $r > t$. Again using (5) as the definition of the virtual iterate x_t as in the perturbed iterate framework, we then have a very simple form for the value of x_t and \hat{x}_t (assuming atomic writes, see Property 3 below):

$$\begin{aligned}
 x_t &= x_0 - \gamma \sum_{u=0}^{t-1} g(\hat{x}_u, \hat{\alpha}^u, i_u); \\
 [\hat{x}_t]_v &= [x_0]_v - \gamma \sum_{\substack{u=0 \\ \text{u s.t. coordinate } v \text{ was written} \\ \text{for } u \text{ before } t}}^{t-1} [g(\hat{x}_u, \hat{\alpha}^u, i_u)]_v.
 \end{aligned} \tag{8}$$

The main idea of the perturbed iterate framework is to use this handle on $\hat{x}_t - x_t$ to analyze the convergence for x_t . In this paper, we can instead give directly the convergence of \hat{x}_t , and so unlike in Mania et al. (2015), we do not require that there exists a T such that x_T lives in shared memory.

3.3 Analysis setup

We describe ASAGA, a sparse asynchronous parallel implementation of Sparse SAGA, in Algorithm 1 in the

Algorithm 1 ASAGA (analyzed algorithm)

```

1: Initialize shared variables  $x$  and  $(\alpha_i)_{i=1}^n$ 
2: keep doing in parallel
3:  $\hat{x} =$  inconsistent read of  $x$ 
4:  $\forall j, \hat{\alpha}_j =$  inconsistent read of  $\alpha_j$ 
5: Sample  $i$  uniformly at random in  $\{1, \dots, n\}$ 
6: Let  $S_i$  be  $f_i$ 's support
7:  $[\bar{\alpha}]_{S_i} = 1/n \sum_{k=1}^n [\hat{\alpha}_k]_{S_i}$ 
8:  $[\delta x]_{S_i} = -\gamma(f'_i(\hat{x}) - \hat{\alpha}_i + D_i[\bar{\alpha}]_{S_i})$ 
9:
10: for  $v$  in  $S_i$  do
11:    $[x]_v \leftarrow [x]_v + [\delta x]_v$  // atomic
12:    $[\alpha_i]_v \leftarrow [f'_i(\hat{x})]_v$  // atomic
13:   // (' $\leftarrow$ ' denotes a shared memory update.)
14: end for
15: end parallel loop
    
```

theoretical form that we analyze, and in Algorithm 2 as its practical implementation. Before stating its convergence, we highlight some properties of Algorithm 1 and make one central assumption.

Property 1 (independence). *Given the “after read” global ordering, i_r is independent of $\hat{x}_t \forall r \geq t$.*

We enforce the independence for $r = t$ in Algorithm 1 by having the core read all the shared data parameters and historical gradients before starting their iterations. Although this is too expensive to be practical if the data is sparse, this is required by the theoretical Algorithm 1 that we can analyze. As Mania et al. (2015) stress, this independence property is assumed in most of the parallel optimization literature. The independence for $r > t$ is a consequence of using the “after read” global ordering instead of the “before read” one.

Property 2 (Unbiased estimator). *The update, $g_t := g(\hat{x}_t, \hat{\alpha}^t, i_t)$, is an unbiased estimator of the true gradient at \hat{x}_t (i.e. (5) yields (6) in conditional expectation).*

This property is crucial for the analysis, as in most related literature. It follows by the independence of i_t with \hat{x}_t and from the computation of $\bar{\alpha}$ on line 7 of Algorithm 1, which ensures that $\mathbb{E}\hat{\alpha}_i = 1/n \sum_{k=1}^n [\hat{\alpha}_k]_{S_i} = [\bar{\alpha}]_{S_i}$, making the update unbiased. In practice, recomputing $\bar{\alpha}$ is not optimal, but storing it instead introduces potential bias issues in the proof (as detailed in Appendix G.3).

Property 3 (atomicity). *The shared parameter coordinate update of $[x]_v$ on line 11 is atomic.*

Since our updates are additions, this means that there are no overwrites, even when several cores compete for the same resources. In practice, this is enforced by using *compare-and-swap* semantics, which are heavily optimized at the processor level and have minimal overhead. Our experiments with non-thread safe algorithms (i.e. where this property is not verified, see Figure 6 of Appendix G) show that compare-and-swap

Algorithm 2 ASAGA (implementation)

```

1: Initialize shared variables  $x$ ,  $(\alpha_i)_{i=1}^n$  and  $\bar{\alpha}$ 
2: keep doing in parallel
3: Sample  $i$  uniformly at random in  $\{1, \dots, n\}$ 
4: Let  $S_i$  be  $f_i$ 's support
5:  $[\hat{x}]_{S_i} =$  inconsistent read of  $x$  on  $S_i$ 
6:  $\hat{\alpha}_i =$  inconsistent read of  $\alpha_i$ 
7:  $[\bar{\alpha}]_{S_i} =$  inconsistent read of  $\bar{\alpha}$  on  $S_i$ 
8:  $[\delta\alpha]_{S_i} = f'_i([\hat{x}]_{S_i}) - \hat{\alpha}_i$ 
9:  $[\delta x]_{S_i} = -\gamma([\delta\alpha]_{S_i} + D_i[\bar{\alpha}]_{S_i})$ 
10: for  $v$  in  $S_i$  do
11:    $[x]_v \leftarrow [x]_v + [\delta x]_v$  // atomic
12:    $[\alpha_i]_v \leftarrow [\alpha_i]_v + [\delta\alpha]_v$  // atomic
13:    $[\bar{\alpha}]_v \leftarrow [\bar{\alpha}]_v + 1/n[\delta\alpha]_v$  // atomic
14: end for
15: end parallel loop
    
```

is necessary to optimize to high accuracy.

Finally, as is standard in the literature, we make an assumption on the maximum delay that asynchrony can cause – this is the *partially asynchronous* setting as defined in Bertsekas and Tsitsiklis (1989):

Assumption 1 (bounded overlaps). *We assume that there exists a uniform bound, called τ , on the maximum number of iterations that can overlap together. We say that iterations r and t overlap if at some point they are processed concurrently. One iteration is being processed from the start of the reading of the shared parameters to the end of the writing of its update. The bound τ means that iterations r cannot overlap with iteration t for $r \geq t + \tau + 1$, and thus that every coordinate update from iteration t is successfully written to memory before the iteration $t + \tau + 1$ starts.*

Our result will give us conditions on τ subject to which we have linear speedups. τ is usually seen as a proxy for p , the number of cores (which lowerbounds it). However, though τ appears to depend linearly on p , it actually depends on several other factors (notably the data sparsity distribution) and can be orders of magnitude bigger than p in real-life experiments. We can upper bound τ by $(p-1)R$, where R is the ratio of the maximum over the minimum iteration time (which encompasses theoretical aspects as well as hardware overhead). More details can be found in Appendix E.

Explicit effect of asynchrony. By using the overlap Assumption 1 in the expression (8) for the iterates, we obtain the following explicit effect of asynchrony that is crucially used in our proof:

$$\hat{x}_t - x_t = \gamma \sum_{u=(t-\tau)_+}^{t-1} G_u^t g(\hat{x}_u, \hat{\alpha}^u, i_u), \quad (9)$$

where G_u^t are $d \times d$ diagonal matrices with terms in $\{0, +1\}$. We know from our definition of t and x_t that

every update in \hat{x}_t is already in x_t – this is the 0 case. Conversely, some updates might be late: this is the +1 case. \hat{x}_t may be lacking some updates from the “past” in some sense, whereas given our global ordering definition, it cannot contain updates from the “future”.

3.4 Convergence and speedup results

We now state our main theoretical results. We give an outline of the proof in Section 3.5 and its full details in Appendix C. We first define a notion of problem sparsity, as it will appear in our results.

Definition 1 (Sparsity). *As in Niu et al. (2011), we introduce $\Delta_r := \max_{v=1..d} |\{i : v \in S_i\}|$. Δ_r is the maximum right-degree in the bipartite graph of the factors and the dimensions, i.e., the maximum number of data points with a specific feature. For succinctness, we also define $\Delta := \Delta_r/n$. We have $1 \leq \Delta_r \leq n$, and hence $1/n \leq \Delta \leq 1$.*

Theorem 2 (Convergence guarantee and rate of ASAGA). *Suppose $\tau < n/10$.⁶ Let*

$$a^*(\tau) := \frac{1}{32 \left(1 + \tau\sqrt{\Delta}\right) \xi(\kappa, \Delta, \tau)} \quad (10)$$

where $\xi(\kappa, \Delta, \tau) := \sqrt{1 + \frac{1}{8\kappa} \min\left\{\frac{1}{\sqrt{\Delta}}, \tau\right\}}$

(note that $\xi(\kappa, \Delta, \tau) \approx 1$ unless $\kappa < 1/\sqrt{\Delta}$ ($\leq \sqrt{n}$)).

For any step size $\gamma = \frac{a}{L}$ with $a \leq a^*(\tau)$, the inconsistent read iterates of Algorithm 1 converge in expectation at a geometric rate of at least: $\rho(a) = \frac{1}{5} \min\left\{\frac{1}{n}, \frac{1}{\kappa}\right\}$, i.e., $\mathbb{E}f(\hat{x}_t) - f(x^*) \leq (1 - \rho)^t \tilde{C}_0$, where \tilde{C}_0 is a constant independent of t ($\approx \frac{n}{\gamma} C_0$ with C_0 as defined in Theorem 1).

This result is very close to SAGA’s original convergence theorem, but with the maximum step size divided by an extra $1 + \tau\sqrt{\Delta}$ factor. Referring to Hofmann et al. (2015) and our own Theorem 1, the rate factor for SAGA is $\min\{1/n, 1/\kappa\}$ up to a constant factor. Comparing this rate with Theorem 2 and inferring the conditions on the maximum step size $a^*(\tau)$, we get the following conditions on the overlap τ for ASAGA to have the same rate as SAGA (comparing upper bounds).

Corollary 3 (Speedup condition). *Suppose $\tau \leq \mathcal{O}(n)$ and $\tau \leq \mathcal{O}\left(\frac{1}{\sqrt{\Delta}} \max\{1, \frac{n}{\kappa}\}\right)$. Then using the step size $\gamma = a^*(\tau)/L$ from (10), ASAGA converges geometrically with rate factor $\Omega\left(\min\left\{\frac{1}{n}, \frac{1}{\kappa}\right\}\right)$ (similar to SAGA), and is thus linearly faster than its sequential counterpart up to a constant factor. Moreover, if $\tau \leq \mathcal{O}\left(\frac{1}{\sqrt{\Delta}}\right)$, then a universal step size of $\Theta\left(\frac{1}{L}\right)$ can be used for ASAGA to be adaptive to local strong convexity with a similar rate to SAGA (i.e., knowledge of κ is not required).*

⁶ASAGA can actually converge for any τ , but the maximum step size then has a term of $\exp(\tau/n)$ in the denominator with much worse constants. See Appendix C.8.

Interestingly, in the well-conditioned regime ($n > \kappa$, where SAGA enjoys a range of stepsizes which all give the same contraction ratio), ASAGA can get the same rate as SAGA even in the non-sparse regime ($\Delta = 1$) for $\tau < \mathcal{O}(n/\kappa)$. This is in contrast to the previous work on asynchronous incremental gradient methods which required some kind of sparsity to get a theoretical linear speedup over their sequential counterpart (Niu et al., 2011; Mania et al., 2015). In the ill-conditioned regime ($\kappa > n$), sparsity is required for a linear speedup, with a bound on τ of $\mathcal{O}(\sqrt{n})$ in the best-case (though degenerate) scenario where $\Delta = 1/n$.

Comparison to related work.

- We give the first convergence analysis for an asynchronous parallel version of SAGA (note that Reddi et al. (2015) only covers an epoch based version of SAGA with random stopping times, a fairly different algorithm).
- Theorem 2 can be directly extended to a parallel extension of the SVRG version from Hofmann et al. (2015), which is adaptive to the local strong convexity with similar rates (see Appendix C.2).
- In contrast to the parallel SVRG analysis from Reddi et al. (2015, Thm. 2), our proof technique handles inconsistent reads and a non-uniform processing speed across f_i ’s. Our bounds are similar (noting that Δ is equivalent to theirs), except for the adaptivity to local strong convexity: ASAGA does not need to know κ for optimal performance, contrary to parallel SVRG (see App. C.2 for more details).
- In contrast to the SVRG analysis from Mania et al. (2015, Thm. 14), we obtain a better dependence on the condition number in our rate ($1/\kappa$ vs. $1/\kappa^2$ for them) and on the sparsity (they get $\tau \leq \mathcal{O}(\Delta^{-1/3})$), while we remove their gradient bound assumption. We also give our convergence guarantee on \hat{x}_t during the algorithm, whereas they only bound the error for the “last” iterate x_T .

3.5 Proof outline

We give here the outline of our proof. Its full details can be found in Appendix C.

Let $g_t := g(\hat{x}_t, \hat{\alpha}^t, i_t)$. By expanding the update equation (5) defining the virtual iterate x_{t+1} and introducing \hat{x}_t in the inner product term, we get:

$$\begin{aligned} \|x_{t+1} - x^*\|^2 &= \|x_t - x^*\|^2 - 2\gamma \langle \hat{x}_t - x^*, g_t \rangle \\ &\quad + 2\gamma \langle \hat{x}_t - x_t, g_t \rangle + \gamma^2 \|g_t\|^2. \end{aligned} \quad (11)$$

In the sequential setting, we require i_t to be independent of x_t to get unbiasedness. In the perturbed iterate framework, we instead require that i_t is independent of \hat{x}_t (see Property 1). This crucial property enables us to use the unbiasedness condition (6) to write: $\mathbb{E}\langle \hat{x}_t - x^*, g_t \rangle = \mathbb{E}\langle \hat{x}_t - x^*, f'(\hat{x}_t) \rangle$. We thus

take the expectation of (11) that allows us to use the μ -strong convexity of f :⁷

$$\langle \hat{x}_t - x^*, f'(\hat{x}_t) \rangle \geq f(\hat{x}_t) - f(x^*) + \frac{\mu}{2} \|\hat{x}_t - x^*\|^2.$$

With further manipulations on the expectation of (11), including the use of the standard inequality $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$ (see Section C.3), we obtain our basic recursive contraction inequality:

$$a_{t+1} \leq \left(1 - \frac{\gamma\mu}{2}\right) a_t + \gamma^2 \mathbb{E} \|g_t\|^2 - 2\gamma e_t + \underbrace{\gamma\mu \mathbb{E} \|\hat{x}_t - x_t\|^2 + 2\gamma \mathbb{E} \langle \hat{x}_t - x_t, g_t \rangle}_{\text{additional asynchrony terms}}, \quad (12)$$

where $a_t := \mathbb{E} \|x_t - x^*\|^2$ and $e_t := \mathbb{E} f(\hat{x}_t) - f(x^*)$.

In the sequential setting, one crucially uses the negative suboptimality term $-2\gamma e_t$ to cancel the variance term $\gamma^2 \mathbb{E} \|g_t\|^2$ (thus deriving a condition on γ). Here, we need to bound the additional asynchrony terms using the same negative suboptimality in order to prove convergence and speedup for our parallel algorithm – thus getting stronger constraints on the maximum step size.

The rest of the proof then proceeds as follows:

- Lemma 1: we first bound the additional asynchrony terms in (12) in terms of past updates $(\mathbb{E} \|g_u\|^2, u \leq t)$. We achieve this by crucially using the expansion (9) for $x_t - \hat{x}_t$, together with the sparsity inequality (44) (which is derived from Cauchy-Schwartz, see Appendix C.4).
- Lemma 2: we then bound the updates $\mathbb{E} \|g_u\|^2$ with respect to past suboptimalities $(e_v)_{v \leq u}$. From our analysis of Sparse SAGA in the sequential case:

$$\mathbb{E} \|g_t\|^2 \leq 2\mathbb{E} \|f'_{i_t}(\hat{x}_t) - f'_{i_t}(x^*)\|^2 + 2\mathbb{E} \|\hat{\alpha}_{i_t}^t - f'_{i_t}(x^*)\|^2$$

We bound the first term by $4Le_t$ using Hofmann et al. (2015, Equation (8)). To express the second term in terms of past suboptimalities, we note that it can be seen as an expectation of past first terms with an adequate probability distribution which we derive and bound.

- By substituting Lemma 2 into Lemma 1, we get a master contraction inequality (28) in terms of a_{t+1} , a_t and $e_u, u \leq t$.
- We define a novel Lyapunov function $\mathcal{L}_t = \sum_{u=0}^t (1 - \rho)^{t-u} a_u$ and manipulate the master inequality to show that \mathcal{L}_t is bounded by a contraction, subject to a maximum step size condition on γ (given in Lemma 3, see Appendix C.1).
- Finally, we unroll the Lyapunov inequality to get the convergence Theorem 2.

⁷Here is our departure point with Mania et al. (2015) who replaced the $f(\hat{x}_t) - f(x^*)$ term with the lower bound $\frac{\mu}{2} \|\hat{x}_t - x^*\|^2$ in this relationship (see their Equation (2.4)), yielding an inequality too loose to get fast rates for SVRG.

4 Empirical results

We now present the main results of our empirical comparison of asynchronous SAGA, SVRG and HOGWILD. Additional results, including convergence and speedup figures with respect to the number of iteration and measures on the τ constant are available in the appendix.

4.1 Experimental setup

Models. Although ASAGA can be applied more broadly, we focus on logistic regression, a model of particular practical importance. The associated objective function takes the following form: $\frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-b_i a_i^\top x)) + \frac{\lambda}{2} \|x\|^2$, where $a_i \in \mathbb{R}^p$ and $b_i \in \{-1, +1\}$ are the data samples.

Datasets. We consider two sparse datasets: RCV1 (Lewis et al., 2004) and URL (Ma et al., 2009); and a dense one, Covtype (Collobert et al., 2002), with statistics listed in the table below. As in Le Roux et al. (2012), Covtype is standardized, thus 100% dense. Δ is $\mathcal{O}(1)$ in all datasets, hence not very insightful when relating it to our theoretical results. Deriving a less coarse sparsity bound remains an open problem.

	n	d	density	L
RCV1	697,641	47,236	0.15%	0.25
URL	2,396,130	3,231,961	0.004%	128.4
Covtype	581,012	54	100%	48428

Hardware and software. Experiments were run on a 40-core machine with 384GB of memory. All algorithms were implemented in Scala. We chose this high-level language despite its typical 20x slowdown compared to C (when using standard libraries, see Appendix G.2) because our primary concern was that the code may easily be reused and extended for research purposes (to this end, we have made all our code available at <https://github.com/RemiLeblond/ASAGA>).

4.2 Implementation details

Exact regularization. Following Schmidt et al. (2016), the amount of regularization used was set to $\lambda = 1/n$. In each update, we project the gradient of the regularization term (we multiply it by D_i as we also do with the vector $\bar{\alpha}$) to preserve the sparsity pattern while maintaining an unbiased estimate of the gradient. For squared ℓ_2 , the Sparse SAGA updates becomes: $x^+ = x - \gamma(f'_{i_t}(x) - \alpha_i + D_i \bar{\alpha} + \lambda D_i x)$.

Comparison with the theoretical algorithm. The algorithm we used in the experiments is fully detailed in Algorithm 2. There are two differences with Algorithm 1. First, in the implementation we pick i_t at random *before* we read data. This enables us to only read the necessary data for a given iteration (i.e. $[\hat{x}_t]_{S_i}, [\hat{\alpha}_{i_t}^t], [\bar{\alpha}^t]_{S_i}$). Although this violates

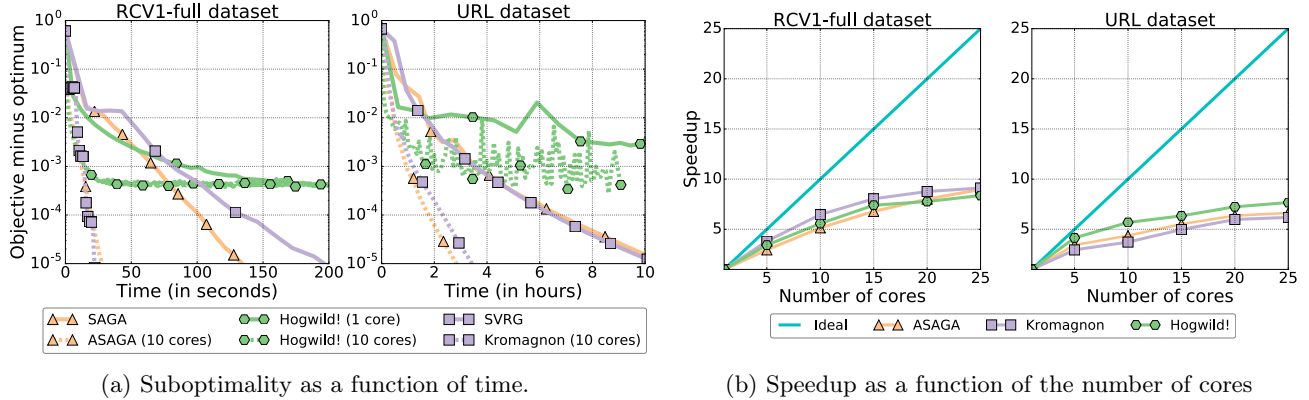


Figure 1: **Convergence and speedup for asynchronous stochastic gradient descent methods.** We display results for RCV1 and URL. Results for Covtype can be found in Appendix D.2.

Property 1, it still performs well in practice.

Second, we maintain $\bar{\alpha}^t$ in memory. This saves the cost of recomputing it at every iteration (which we can no longer do since we only read a subset data). Again, in practice the implemented algorithm enjoys good performance. But this design choice raises a subtle point: the update is not guaranteed to be unbiased in this setup (see Appendix G.3 for more details).

4.3 Results

We first compare three different asynchronous variants of stochastic gradient methods on the aforementioned datasets: ASAGA, presented in this work, KROMAGNON, the asynchronous sparse SVRG method described in Mania et al. (2015) and HOGWILD (Niu et al., 2011). Each method had its step size chosen so as to give the fastest convergence (up to 10^{-3} in the special case of HOGWILD). The results can be seen in Figure 1a: for each method we consider its asynchronous version with both one (hence sequential) and ten processors. This figure reveals that the asynchronous version offers a significant speedup over its sequential counterpart.

We then examine the speedup relative to the increase in the number of cores. The speedup is measured as time to achieve a suboptimality of 10^{-5} (10^{-3} for HOGWILD) with one core divided by time to achieve the same suboptimality with several cores, averaged over 3 runs. Again, we choose step size leading to fastest convergence (see Appendix G.2 for information about the step sizes). Results are displayed in Figure 1b.

As predicted by our theory, we observe linear “theoretical” speedups (i.e. in terms of number of iterations, see Appendix D.2). However, with respect to running time, the speedups seem to taper off after 20 cores. This phenomenon can be explained by the fact that our hardware model is by necessity a simplification of reality. As noted in Duchi et al. (2015), in a modern machine there is no such thing as *shared memory*. Each

core has its own levels of cache (L1, L2, L3) in addition to RAM. The more cores are used, the lower in the memory stack information goes and the slower it gets. More experimentation is needed to quantify that effect and potentially increase performance.

5 Conclusions and future work

We have described ASAGA, a novel sparse and fully asynchronous variant of the incremental gradient algorithm SAGA. Building on the recently proposed “perturbed iterate” framework, we have introduced a novel analysis of the algorithm and proven that under mild conditions ASAGA is linearly faster than SAGA. Our empirical benchmarks confirm speedups up to 10x.

Our proof technique accommodates more realistic settings than is usually the case in the literature (e.g. inconsistent reads/writes and an unbounded gradient); we obtain tighter conditions than in previous work. In particular, we show that sparsity is not always necessary to get linear speedups. Further, we have proposed a novel perspective to clarify an important technical issue present in most of the recent convergence rate proofs for asynchronous parallel optimization algorithms.

Schmidt et al. (2016) have shown that SAG enjoys much improved performance when combined with non-uniform sampling and line-search. We have also noticed that our Δ_r constant (being essentially a maximum) sometimes fails to accurately represent the full sparsity distribution of our datasets. Finally, while our algorithm can be directly ported to a distributed master-worker architecture, its communication pattern would have to be optimized to avoid prohibitive costs. Limiting communications can be interpreted as artificially increasing the delay, yielding an interesting trade-off between delay influence and communication costs.

A final interesting direction for future analysis is the further exploration of the τ term, which we have shown encompasses more complexity than previously thought.

Acknowledgments

We would like to thank Xinghao Pan for sharing with us their implementation of KROMAGNON, as well as Alberto Chiappa for spotting a typo in the proof. This work was partially supported by a Google Research Award and the MSR-Inria Joint Center. FP acknowledges financial support from the chaire *Économie des nouvelles données* with the *data science* joint research initiative with the *fonds AXA pour la recherche*.

References

- D. P. Bertsekas and J. N. Tsitsiklis. *Parallel and Distributed Computation: Numerical Methods*. Prentice Hall, 1989.
- R. Collobert, S. Bengio, and Y. Bengio. A parallel mixture of svms for very large scale problems. *Neural Comput.*, 14:1105–1114, 2002.
- C. De Sa, C. Zhang, K. Olukotun, and C. Ré. Taming the wild: a unified analysis of Hogwild!-style algorithms. In *NIPS*, 2015.
- A. Defazio, F. Bach, and S. Lacoste-Julien. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. In *NIPS*, 2014.
- J. C. Duchi, S. Chaturapruek, and C. Ré. Asynchronous stochastic convex optimization. In *NIPS*, 2015.
- T. Hofmann, A. Lucchi, S. Lacoste-Julien, and B. McWilliams. Variance reduced stochastic gradient descent with neighbors. In *NIPS*, 2015.
- C.-J. Hsieh, H.-F. Yu, and I. Dhillon. PASSCoDe: Parallel ASynchronous Stochastic dual Co-ordinate Descent. In *ICML*, 2015.
- R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *NIPS*, 2013.
- J. Konecny and P. Richtarik. Semi-stochastic gradient descent methods. *arXiv:1312.1666*, 2013.
- N. Le Roux, M. Schmidt, and F. Bach. A stochastic gradient method with an exponential convergence rate for finite training sets. In *NIPS*, 2012.
- D. D. Lewis, Y. Yang, T. G. Rose, and F. Li. RCV1: A new benchmark collection for text categorization research. *JMLR*, 5:361–397, 2004.
- X. Lian, Y. Huang, Y. Li, and J. Liu. Asynchronous parallel stochastic gradient for nonconvex optimization. In *NIPS*, 2015.
- J. Liu, S. J. Wright, C. Ré, V. Bittorf, and S. Sridhar. An asynchronous parallel stochastic coordinate descent algorithm. *JMLR*, 16:285–322, 2015.
- C. Ma, V. Smith, M. Jaggi, M. I. Jordan, P. Richtarik, and M. Takac. Adding vs. averaging in distributed primal-dual optimization. In *ICML*, 2015.
- J. Ma, L. K. Saul, S. Savage, and G. M. Voelker. Identifying suspicious URLs: an application of large-scale online learning. In *ICML*, 2009.
- H. Mania, X. Pan, D. Papailiopoulos, B. Recht, K. Ramchandran, and M. I. Jordan. Perturbed iterate analysis for asynchronous stochastic optimization. *arXiv:1507.06970v2*, 2015.
- F. Niu, B. Recht, C. Re, and S. Wright. Hogwild: a lock-free approach to parallelizing stochastic gradient descent. In *NIPS*, 2011.
- S. J. Reddi, A. Hefny, S. Sra, B. Póczos, and A. Smola. On variance reduction in stochastic gradient descent and its asynchronous variants. In *NIPS*, 2015.
- M. Schmidt, N. Le Roux, and F. Bach. Minimizing finite sums with the stochastic average gradient. *F. Math. Program.*, 2016.
- S. Shalev-Shwartz and T. Zhang. Stochastic dual coordinate ascent methods for regularized loss. *JMLR*, 14:567–599, 2013.
- S.-Y. Zhao and W.-J. Li. Fast asynchronous parallel stochastic gradient descent. In *AAAI*, 2016.

Appendix Outline:

- In Appendix A, we give a simple example illustrating why the “After Write” approach can break the crucial unbiasedness condition (6) needed for standard convergence proofs.
- In Appendix B, we adapt the proof from Hofmann et al. (2015) to prove Theorem 1, our convergence result for serial Sparse SAGA.
- In Appendix C, we first give a detailed outline and then the complete details for the proof of convergence for ASAGA (Theorem 2) as well as its linear speedup regimes (Corollary 3).
- In Appendix D, we analyze additional experimental results, including a comparison of serial SAGA algorithms and a look at “theoretical speedups” for ASAGA.
- In Appendix E, we take a closer look at the τ constant. We argue that it encompasses more complexity than is usually implied in the literature, as additional results that we present indicate.
- In Appendix F, we compare the lagged updates implementation of SAGA with our sparse algorithm, and explain why adapting the former to the asynchronous setting is difficult.
- In Appendix G, we give additional details about the datasets and our implementation.

A Problematic Example for the “After Write” Approach

We provide a concrete example to illustrate the non-independence issue arising from the “after write” approach. Suppose that we have two cores and that f has two factors: f_1 which has support on only one variable, and f_2 which has support on 10^6 variables and thus yields a gradient step that is significantly more expensive to compute. In the “after write” approach, x_0 is the initial content of the memory, and we do not officially know yet whether \hat{x}_0 is the local copy read by the first core or the second core, but we are sure that $\hat{x}_0 = x_0$ as no update can occur in shared memory without incrementing the counter. There are four possibilities for the next step defining x_1 depending on which index i was sampled on each core. If any core samples $i = 1$, we know that $x_1 = x_0 - \gamma f'_1(x_0)$ as it will be the first (much faster update) to complete. This happens in 3 out of 4 possibilities; we thus have that $\mathbb{E}x_1 = x_0 - \gamma(\frac{3}{4}f'_1(x_0) + \frac{1}{4}f'_2(x_0))$ – we see that this analysis scheme *does not* satisfy the crucial unbiasedness condition (6).

To understand this subtle point better, note that in this very simple example, i_0 and i_1 are not independent. We can show that $P(i_1 = 2 \mid i_0 = 2) = 1$. They share dependency through the labeling assignment.

The only way we can think to resolve this issue and ensure unbiasedness in the “after write” framework is to assume that the computation time for the algorithm running on a core is independent of the sample i chosen. This assumption seems overly strong in the context of potentially heterogeneous factors f_i ’s, and is thus a fundamental flaw in the “after write” framework that has mostly been ignored in the recent asynchronous optimization literature.

We note that Bertsekas and Tsitsiklis (1989) briefly discussed this issue in Section 7.8.3 of their book, stressing that their analysis for SGD required that the scheduling of computation was independent from the randomness from SGD, but they did not offer any solution if this assumption was not satisfied. Both the “before read” labeling from Mania et al. (2015) and our proposed “after read” labeling resolve this issue.

B Proof of Theorem 1

Proof sketch for Hofmann et al. (2015). As we will heavily reuse the proof technique from Hofmann et al. (2015), we start by giving its sketch.

First, the authors combine classical strong convexity and Lipschitz inequalities to derive the inequality Hofmann et al. (2015, Lemma 1):

$$\mathbf{E}\|x^+ - x^*\|^2 \leq (1 - \gamma\mu)\|x - x^*\|^2 + 2\gamma^2\mathbf{E}\|\alpha_i - f'_i(x^*)\|^2 + (4\gamma^2L - 2\gamma)(f(x) - f(x^*)). \quad (13)$$

This gives a contraction term, as well as two additional terms; $2\gamma^2\mathbf{E}\|\alpha_i - f'_i(x^*)\|^2$ is a positive variance term, but $(4\gamma^2L - 2\gamma)(f(x) - f(x^*))$ is a negative suboptimality term (provided γ is small enough). The suboptimality term can then be used to cancel the variance one.

Second, the authors use a classical smoothness upper bound to control the variance term and relate it to the suboptimality. However, since the α_i are partial gradients computed at previous time steps, the upper bounds of the variance involve suboptimality at previous time steps, which are not directly relatable to the current suboptimality.

Third, to circumvent this issue, a Lyapunov function is defined to encompass both current and past terms. To finish the proof, Hofmann et al. (2015) show that the Lyapunov function is a contraction.

Proof outline. Fortunately, we can reuse most of the proof from Hofmann et al. (2015) to show that Sparse SAGA converges at the same rate as regular SAGA. In fact, once we establish that Hofmann et al. (2015, Lemma 1) is still verified we are done.

To prove this, we derive close variants of equations (6) and (9) in their paper, which we remind the reader of here:

$$\mathbf{E}\|f'_i(x) - \bar{\alpha}_i\|^2 \leq 2\mathbf{E}\|f'_i(x) - f'_i(x^*)\|^2 + 2\mathbf{E}\|\bar{\alpha}_i - f'_i(x^*)\|^2, \quad \text{Hofmann et al. (2015, Eq.(6))}$$

$$\mathbf{E}\|\bar{\alpha}_i - f'_i(x^*)\|^2 \leq \mathbf{E}\|\alpha_i - f'_i(x^*)\|^2. \quad \text{Hofmann et al. (2015, Eq.(9))}$$

Deriving Hofmann et al. (2015, Equation (6)). We first show that the update estimator is unbiased. The estimator is unbiased if:

$$\mathbf{E}D_i\bar{\alpha} = \mathbf{E}\alpha_i = \frac{1}{n} \sum_{i=1}^n \alpha_i. \quad (14)$$

We have:

$$\mathbf{E}D_i\bar{\alpha} = \frac{1}{n} \sum_{i=1}^n D_i\bar{\alpha} = \frac{1}{n} \sum_{i=1}^n P_{S_i}D\bar{\alpha} = \frac{1}{n} \sum_{i=1}^n \sum_{v \in S_i} \frac{[\bar{\alpha}]_v e_v}{p_v} = \sum_{v=1}^d \left(\sum_{i|v \in S_i} 1 \right) \frac{[\bar{\alpha}]_v e_v}{np_v},$$

where e_v is the vector whose only nonzero component is the v component which is equal to 1.

By definition, $\sum_{i|v \in S_i} 1 = np_v$, which gives us Equation (14).

We define $\bar{\alpha}_i := \alpha_i - D_i\bar{\alpha}$ (contrary to Hofmann et al. (2015) where the authors define $\bar{\alpha}_i := \alpha_i - \bar{\alpha}$ since they do not concern themselves with sparsity). Using the inequality $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$, we get:

$$\mathbf{E}\|f'_i(x) - \bar{\alpha}_i\|^2 \leq 2\mathbf{E}\|f'_i(x) - f'_i(x^*)\|^2 + 2\mathbf{E}\|\bar{\alpha}_i - f'_i(x^*)\|^2, \quad (15)$$

which is our equivalent to Hofmann et al. (2015, Eq.(6)), where only our definition of $\bar{\alpha}_i$ differs.

Deriving [Hofmann et al. \(2015, Equation \(9\)\)](#). We want to prove [Hofmann et al. \(2015, Eq.\(9\)\)](#):

$$\mathbf{E}\|\bar{\alpha}_i - f'_i(x^*)\|^2 \leq \mathbf{E}\|\alpha_i - f'_i(x^*)\|^2. \quad (16)$$

We have:

$$\mathbf{E}\|\bar{\alpha}_i - f'_i(x^*)\|^2 = \mathbf{E}\|\alpha_i - f'_i(x^*)\|^2 - 2\mathbf{E}\langle \alpha_i - f'_i(x^*), D_i \bar{\alpha} \rangle + \mathbf{E}\|D_i \bar{\alpha}\|^2. \quad (17)$$

Let $D_{-i} := P_{S_i^c} D$; we then have the orthogonal decomposition $D\alpha = D_i\alpha + D_{-i}\alpha$ with $D_i\alpha \perp D_{-i}\alpha$, as they have disjoint support. We now use the orthogonality of $D_{-i}\alpha$ with any vector with support in S_i to simplify the expression (17) as follows:

$$\begin{aligned} \mathbf{E}\langle \alpha_i - f'_i(x^*), D_i \bar{\alpha} \rangle &= \mathbf{E}\langle \alpha_i - f'_i(x^*), D_i \bar{\alpha} + D_{-i} \bar{\alpha} \rangle && (\alpha_i - f'_i(x^*) \perp D_{-i} \alpha) \\ &= \mathbf{E}\langle \alpha_i - f'_i(x^*), D \bar{\alpha} \rangle \\ &= \langle \mathbf{E}(\alpha_i - f'_i(x^*)), D \bar{\alpha} \rangle \\ &= \langle \mathbf{E} \alpha_i, D \bar{\alpha} \rangle && (f'(x^*) = 0) \\ &= \bar{\alpha}^\top D \bar{\alpha}. && (18) \end{aligned}$$

Similarly,

$$\begin{aligned} \mathbf{E}\|D_i \bar{\alpha}\|^2 &= \mathbf{E}\langle D_i \bar{\alpha}, D_i \bar{\alpha} \rangle \\ &= \mathbf{E}\langle D_i \bar{\alpha}, D \bar{\alpha} \rangle && (D_i \alpha \perp D_{-i} \alpha) \\ &= \langle \mathbf{E} D_i \bar{\alpha}, D \bar{\alpha} \rangle \\ &= \bar{\alpha}^\top D \bar{\alpha}. && (19) \end{aligned}$$

Putting it all together,

$$\mathbf{E}\|\bar{\alpha}_i - f'_i(x^*)\|^2 = \mathbf{E}\|\alpha_i - f'_i(x^*)\|^2 - \bar{\alpha}^\top D \bar{\alpha} \leq \mathbf{E}\|\alpha_i - f'_i(x^*)\|^2. \quad (20)$$

This is our version of [Hofmann et al. \(2015, Equation \(9\)\)](#), which finishes the proof of [Hofmann et al. \(2015, Lemma 1\)](#). The rest of the proof from [Hofmann et al. \(2015\)](#) can then be reused without modification to obtain [Theorem 1](#). \square

C Proof of [Theorem 2](#) and [Corollary 3](#)

C.1 Detailed outline

We first give a detailed outline of the proof. The complete proof is given in the rest of [Appendix C](#).

Initial recursive inequality. Let $g_t := g(\hat{x}_t, \hat{\alpha}^t, i_t)$. From the update equation (5) defining the virtual iterate x_{t+1} , the perturbed iterate framework ([Mania et al., 2015](#)) gives:

$$\begin{aligned} \|x_{t+1} - x^*\|^2 &= \|x_t - \gamma g_t - x^*\|^2 \\ &= \|x_t - x^*\|^2 + \gamma^2 \|g_t\|^2 - 2\gamma \langle x_t - x^*, g_t \rangle \\ &= \|x_t - x^*\|^2 + \gamma^2 \|g_t\|^2 - 2\gamma \langle \hat{x}_t - x^*, g_t \rangle + 2\gamma \langle \hat{x}_t - x_t, g_t \rangle. \end{aligned} \quad (21)$$

Note that we have introduced \hat{x}_t in the inner product because g_t is a function of \hat{x}_t , not x_t .

In the sequential setting, we require i_t to be independent of x_t to get unbiasedness. In the perturbed iterate framework, we instead require that i_t is independent of \hat{x}_t (see [Property 1](#)). This crucial property

enables us to use the unbiasedness condition (6) to write: $\mathbb{E}\langle \hat{x}_t - x^*, g_t \rangle = \mathbb{E}\langle \hat{x}_t - x^*, f'(\hat{x}_t) \rangle$. We thus take the expectation of (21) that allows us to use the μ -strong convexity of f :⁸

$$\langle \hat{x}_t - x^*, f'(\hat{x}_t) \rangle \geq f(\hat{x}_t) - f(x^*) + \frac{\mu}{2} \|\hat{x}_t - x^*\|^2. \quad (22)$$

With further manipulations on the expectation of (21), including the use of the standard inequality $\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2$ (see Section C.3), we obtain our basic recursive contraction inequality:

$$a_{t+1} \leq \left(1 - \frac{\gamma\mu}{2}\right)a_t + \gamma^2\mathbb{E}\|g_t\|^2 + \underbrace{\gamma\mu\mathbb{E}\|\hat{x}_t - x^*\|^2 + 2\gamma\mathbb{E}\langle \hat{x}_t - x_t, g_t \rangle}_{\text{additional asynchrony terms}} - 2\gamma e_t, \quad (23)$$

where $a_t := \mathbb{E}\|x_t - x^*\|^2$ and $e_t := \mathbb{E}f(\hat{x}_t) - f(x^*)$.

Inequality (23) is a midway point between the one derived in the proof of Lemma 1 in Hofmann et al. (2015) and Equation (2.5) in Mania et al. (2015), because we use the tighter strong convexity bound (22) than in the latter (giving us the important extra term $-2\gamma e_t$).

In the sequential setting, one crucially uses the negative suboptimality term $-2\gamma e_t$ to cancel the variance term $\gamma^2\mathbb{E}\|g_t\|^2$ (thus deriving a condition on γ). In our setting, we need to bound the additional asynchrony terms using the same negative suboptimality in order to prove convergence and speedup for our parallel algorithm – this will give stronger constraints on the maximum step size.

The rest of the proof then proceeds as follows:

1. By using the expansion (9) for $\hat{x}_t - x_t$, we can bound the additional asynchrony terms in (23) in terms of the past updates ($\mathbb{E}\|g_u\|^2, u \leq t$). This gives Lemma 1 below.
2. We then bound the updates $\mathbb{E}\|g_t\|^2$ in terms of past suboptimalities $(e_u)_{u \leq v}$ by using standard SAGA inequalities and carefully analyzing the update rule for α_i^+ (2) in expectation. This gives Lemma 2 below.
3. By substituting Lemma 2 into Lemma 1, we get a master contraction inequality (28) in terms of a_{t+1} , a_t and $e_u, u \leq t$.
4. We define a novel Lyapunov function $\mathcal{L}_t = \sum_{u=0}^t (1-\rho)^{t-u} a_u$ and manipulate the master inequality to show that \mathcal{L}_t is bounded by a contraction, subject to a maximum step size condition on γ (given in Lemma 3 below).
5. Finally, we unroll the Lyapunov inequality to get the convergence Theorem 2.

We list the key lemmas below with their proof sketch, and give the detailed proof in the later sections of Appendix C.

Lemma 1 (Inequality in terms of $g_t := g(\hat{x}_t, \hat{\alpha}^t, i_t)$). *For all $t \geq 0$:*

$$a_{t+1} \leq \left(1 - \frac{\gamma\mu}{2}\right)a_t + \gamma^2 C_1 \mathbb{E}\|g_t\|^2 + \gamma^2 C_2 \sum_{u=(t-\tau)_+}^{t-1} \mathbb{E}\|g_u\|^2 - 2\gamma e_t, \quad (24)$$

where $C_1 := 1 + \sqrt{\Delta}\tau$ and $C_2 := \sqrt{\Delta} + \gamma\mu C_1$.

To prove this lemma we need to bound both $\mathbb{E}\|\hat{x}_t - x^*\|^2$ and $\mathbb{E}\langle \hat{x}_t - x_t, g_t \rangle$ with respect to $(g_u, u \leq t)$. We achieve this by crucially using Equation (9), together with the following proposition, which we derive by a combination of Cauchy-Schwartz and our sparsity definition (see Section C.4).

$$\mathbb{E}\langle G_u^t g_u, g_t \rangle \leq \frac{\sqrt{\Delta}}{2} (\mathbb{E}\|g_u\|^2 + \mathbb{E}\|g_t\|^2). \quad (25)$$

⁸Note that here is our departure point with Mania et al. (2015) who replaced the $f(\hat{x}_t) - f(x^*)$ term with the lower bound $\frac{\mu}{2}\|\hat{x}_t - x^*\|^2$ in this relationship (see their Equation (2.4)), thus yielding an inequality too loose afterwards to get the fast rates for SVRG.

Lemma 2 (Suboptimality bound on $\mathbb{E}\|g_t\|^2$). For all $t \geq 0$,

$$\mathbb{E}\|g_t\|^2 \leq 4Le_t + \frac{4L}{n} \sum_{u=1}^{t-1} \left(1 - \frac{1}{n}\right)^{(t-2\tau-u-1)_+} e_u + 4L\left(1 - \frac{1}{n}\right)^{(t-\tau)_+} \tilde{e}_0. \quad (26)$$

where $\tilde{e}_0 := \frac{1}{2L} \mathbb{E}\|\alpha_i^0 - f'_i(x^*)\|^2$.⁹

From our Sparse SAGA proof we know that (see Appendix B):

$$\mathbb{E}\|g_t\|^2 \leq 2\mathbb{E}\|f'_i(\hat{x}_t) - f'_i(x^*)\|^2 + 2\mathbb{E}\|\hat{\alpha}_{i_t}^t - f'_i(x^*)\|^2. \quad (27)$$

We can handle the first term by taking the expectation over a Lipschitz inequality (Hofmann et al. (2015, Equations (7) and (8))). All that remains to prove the lemma is to express the $\mathbb{E}\|\hat{\alpha}_{i_t}^t - f'_i(x^*)\|^2$ term in terms of past suboptimality. We note that it can be seen as an expectation of past first terms with an adequate probability distribution which we derive and bound.

From our algorithm, we know that each dimension of the memory vector $[\hat{\alpha}_i]_v$ contains a partial gradient computed at some point in the past $[f'_i(\hat{x}_{u^t_{i,v}})]_v$ ¹⁰ (unless $u = 0$, in which case we replace the partial gradient with α_i^0). We then derive bounds on $P(u^t_{i,v} = u)$ and sum on all possible u . Together with clever conditioning, we obtain Lemma 2 (see Section C.5).

Master inequality. Let H_t be defined as $H_t := \sum_{u=1}^{t-1} \left(1 - \frac{1}{n}\right)^{(t-2\tau-u-1)_+} e_u$. Then, by setting (26) into Lemma 1, we get (see Section C.6):

$$\begin{aligned} a_{t+1} \leq & \left(1 - \frac{\gamma\mu}{2}\right)a_t - 2\gamma e_t + 4L\gamma^2 C_1 \left(e_t + \left(1 - \frac{1}{n}\right)^{(t-\tau)_+} \tilde{e}_0\right) + \frac{4L\gamma^2 C_1}{n} H_t \\ & + 4L\gamma^2 C_2 \sum_{u=(t-\tau)_+}^{t-1} \left(e_u + \left(1 - \frac{1}{n}\right)^{(u-\tau)_+} \tilde{e}_0\right) + \frac{4L\gamma^2 C_2}{n} \sum_{u=(t-\tau)_+}^{t-1} H_u. \end{aligned} \quad (28)$$

Lyapunov function and associated recursive inequality. We now have the beginning of a contraction with additional positive terms which all converge to 0 as we near the optimum, as well as our classical negative suboptimality term. This is not unusual in the variance reduction literature. One successful approach in the sequential case is then to define a Lyapunov function which encompasses all terms and is a true contraction (see Defazio et al. (2014); Hofmann et al. (2015)). We emulate this solution here. However, while all terms in the sequential case only depend on the current iterate, t , in the parallel case we have terms “from the past” in our inequality. To resolve this issue, we define a more involved Lyapunov function which also encompasses past iterates:

$$\mathcal{L}_t = \sum_{u=0}^t (1 - \rho)^{t-u} a_u, \quad 0 < \rho < 1, \quad (29)$$

where ρ is a target contraction rate that we define later.

Using the master inequality (28), we get (see Appendix C.7):

$$\begin{aligned} \mathcal{L}_{t+1} &= (1 - \rho)^{t+1} a_0 + \sum_{u=0}^t (1 - \rho)^{t-u} a_{u+1} \\ &\leq (1 - \rho)^{t+1} a_0 + \left(1 - \frac{\gamma\mu}{2}\right)\mathcal{L}_t + \sum_{u=1}^t r_u^t e_u + r_0^t \tilde{e}_0. \end{aligned} \quad (30)$$

⁹We introduce this quantity instead of e_0 so as to be able to handle the arbitrary initialization of the α_i^0 .

¹⁰More precisely: $\forall t, i, v \exists u^t_{i,v} < t$ s.t. $[\hat{\alpha}_i^t]_v = [f'_i(\hat{x}_{u^t_{i,v}})]_v$.

The aim is to prove that \mathcal{L}_t is bounded by a contraction. We have two promising terms at the beginning of the inequality, and then we need to handle the last term. Basically, we can rearrange the sums in (28) to expose a simple sum of e_u multiplied by factors r_u^t .

Under specific conditions on ρ and γ , we can prove that r_u^t is negative for all $u \geq 1$, which coupled with the fact that each e_u is positive means that we can safely drop the sum term from the inequality. The r_0^t term is a bit trickier and is handled separately.

In order to have a bound on e_t directly rather than on $\mathbb{E}\|\hat{x}_t - x^*\|^2$, we then introduce an additional γe_t term on both sides of (30). The bound on γ under which the modified $r_t^t + \gamma$ is negative is then twice as small (we could have used any multiplier between 0 and 2γ , but chose γ for simplicity's sake). This condition is given in the following Lemma.

Lemma 3 (Sufficient condition for convergence). *Suppose $\tau < n/10$ and $\rho \leq 1/4n$. If*

$$\gamma \leq \gamma^* = \frac{1}{32L(1 + \sqrt{\Delta}\tau)\sqrt{1 + \frac{1}{8\kappa} \min(\tau, \frac{1}{\sqrt{\Delta}})}} \quad (31)$$

then for all $u \geq 1$, the r_u^t from (30) verify:

$$r_u^t \leq 0; \quad r_t^t + \gamma \leq 0, \quad (32)$$

and thus we have:

$$\gamma e_t + \mathcal{L}_{t+1} \leq (1 - \rho)^{t+1} a_0 + (1 - \frac{\gamma\mu}{2})\mathcal{L}_t + r_0^t \tilde{e}_0. \quad (33)$$

We obtain this result after carefully deriving the r_u^t terms. We find a second-order polynomial inequality in γ , which we simplify down to (31) (see Appendix C.8).

We can then finish the argument to bound the suboptimality error e_t . We have:

$$\mathcal{L}_{t+1} \leq \gamma e_t + \mathcal{L}_{t+1} \leq (1 - \frac{\gamma\mu}{2})\mathcal{L}_t + (1 - \rho)^{t+1}(a_0 + A\tilde{e}_0). \quad (34)$$

We have two linearly contracting terms. The sum contracts linearly with the worst rate between the two (the smallest geometric rate factor). If we define $\rho^* := \nu \min(\rho, \gamma\mu/2)$, with $0 < \nu < 1$,¹¹ then we get:

$$\gamma e_t + \mathcal{L}_{t+1} \leq (1 - \frac{\gamma\mu}{2})^{t+1} \mathcal{L}_0 + (1 - \rho^*)^{t+1} \frac{a_0 + A\tilde{e}_0}{1 - \eta} \quad (35)$$

$$\gamma e_t \leq (1 - \rho^*)^{t+1} \left(\mathcal{L}_0 + \frac{1}{1 - \eta} (a_0 + A\tilde{e}_0) \right), \quad (36)$$

where $\eta := \frac{1-M}{1-\rho^*}$ with $M := \max(\rho, \gamma\mu/2)$. Our geometric rate factor is thus ρ^* (see Appendix C.9).

C.2 Extension to SVRG

Our proof can easily be adapted to accommodate the SVRG variant introduced in Hofmann et al. (2015), which is closer to SAGA than the initial SVRG algorithm and which is adaptive to local strong convexity (it does not require the inner loop epoch size $m = \Omega(\kappa)$ as a hyperparameter). In this variant, instead of computing a full gradient every iterations, a random binary variable U with probability $P(U = 1) = 1/n$ is sampled at the beginning of every iteration to determine whether a full gradient is computed or a normal SVRG step is made. If $U = 1$, then a full gradient is computed. Otherwise the algorithm takes a normal inner SVRG step.¹²

¹¹ ν is introduced to circumvent the problematic case where ρ and $\gamma\mu/2$ are too close together.

¹²Note that the parallel implementation is not very straightforward, as it requires a way to communicate to cores when they should start computing a batch gradient instead of inner steps.

To prove convergence, all one has to do is to modify Lemma 2 very slightly (the only difference is that the $(t - 2\tau - u - 1)_+$ exponent is replaced by $(t - u)$ and the rest of the proof can be used as is). The justification for this small tweak is that the batch steps in SVRG are fully synchronized. More details can be found in Section C.5 (see footnote 16).

By using our “after read” labeling, we were also able to derive a convergence and speedup proof for the original SVRG algorithm, but the proof technique diverges after Lemma 1. This is beyond the scope of this paper, so we omit it here. Using the “after read” labeling and a different proof technique from Mania et al. (2015), we obtain an epoch size in $\mathcal{O}(\kappa)$ instead of $\mathcal{O}(\kappa^2)$ and a dependency in our overlap bound in $\mathcal{O}(\Delta^{-1/2})$ instead of $\mathcal{O}(\Delta^{-1/3})$.

ASAGA vs. asynchronous SVRG. There are several scenarios in which ASAGA can be practically advantageous over its closely related cousin, asynchronous SVRG (note though that “asynchronous” SVRG still requires a synchronization step to compute the full gradients).

First, while SAGA trades memory for less computation, in the case of generalized linear models the memory cost can be reduced to $\mathcal{O}(n)$, which is the same as for SVRG. This is of course also true for their asynchronous counterparts.

Second, as ASAGA does not require any synchronization steps, it is better suited to heterogeneous computing environments (where cores have different clock speeds or are shared with other applications).

Finally, ASAGA does not require knowing the condition number κ for optimal convergence in the sparse regime. It is thus adaptive to local strong convexity, whereas SVRG is not. Indeed, SVRG and its asynchronous variant require setting an additional hyper-parameter – the epoch size m – which needs to be at least $\Omega(\kappa)$ for convergence but yields a slower effective convergence rate than ASAGA if it is set much bigger than κ . SVRG thus requires tuning this additional hyper-parameter or running the risk of either slower convergence (if the epoch size chosen is much bigger than the condition number) or even not converging at all (if m is chosen to be much smaller than κ).¹³

C.3 Initial recursive inequality derivation

We start by proving Equation (23). Let $g_t := g(\hat{x}_t, \hat{\alpha}^t, i_t)$. From (5), we get:

$$\begin{aligned} \|x_{t+1} - x^*\|^2 &= \|x_t - \gamma g_t - x^*\|^2 = \|x_t - x^*\|^2 + \gamma^2 \|g_t\|^2 - 2\gamma \langle x_t - x^*, g_t \rangle \\ &= \|x_t - x^*\|^2 + \gamma^2 \|g_t\|^2 - 2\gamma \langle \hat{x}_t - x^*, g_t \rangle + 2\gamma \langle \hat{x}_t - x_t, g_t \rangle. \end{aligned}$$

In order to prove Equation (23), we need to bound the $-2\gamma \langle \hat{x}_t - x^*, g_t \rangle$ term. Thanks to Property 1, we can write:

$$\mathbb{E} \langle \hat{x}_t - x^*, g_t \rangle = \mathbb{E} \langle \hat{x}_t - x^*, \mathbf{E} g_t \rangle = \mathbb{E} \langle \hat{x}_t - x^*, f'(\hat{x}_t) \rangle.$$

We can now use a classical strong convexity bound as well as a squared triangle inequality to get:

$$\begin{aligned} -\langle \hat{x}_t - x^*, f'(\hat{x}_t) \rangle &\leq -(f(\hat{x}_t) - f(x^*)) - \frac{\mu}{2} \|\hat{x}_t - x^*\|^2 && \text{(Strong convexity bound)} \\ -\|\hat{x}_t - x^*\|^2 &\leq \|\hat{x}_t - x_t\|^2 - \frac{1}{2} \|x_t - x^*\|^2 && (\|a + b\|^2 \leq 2\|a\|^2 + 2\|b\|^2) \\ -2\gamma \mathbb{E} \langle \hat{x}_t - x^*, g_t \rangle &\leq -\frac{\gamma\mu}{2} \mathbb{E} \|x_t - x^*\|^2 + \gamma\mu \mathbb{E} \|\hat{x}_t - x_t\|^2 - 2\gamma (\mathbb{E} f(\hat{x}_t) - f(x^*)). && (37) \end{aligned}$$

Putting it all together, we get the initial recursive inequality (23), rewritten here explicitly:

$$a_{t+1} \leq \left(1 - \frac{\gamma\mu}{2}\right) a_t + \gamma^2 \mathbb{E} \|g_t\|^2 + \gamma\mu \mathbb{E} \|\hat{x}_t - x_t\|^2 + 2\gamma \mathbb{E} \langle \hat{x}_t - x_t, g_t \rangle - 2\gamma e_t, \quad (38)$$

where $a_t := \mathbb{E} \|x_t - x^*\|^2$ and $e_t := \mathbb{E} f(\hat{x}_t) - f(x^*)$.

¹³Note that as SAGA (and contrary to the original SVRG), the SVRG variant from Hofmann et al. (2015) does not require knowledge of κ and is thus adaptive to local strong convexity, which carries over to its asynchronous adaptation.

C.4 Proof of Lemma 1

To prove Lemma 1, we now bound both $\mathbb{E}\|\hat{x}_t - x_t\|^2$ and $\mathbb{E}\langle \hat{x}_t - x_t, g_t \rangle$ with respect to $\mathbb{E}\|g_u\|^2, u \leq t$.

We start by proving a relevant property of Δ , which enables us to derive an essential inequality for both these terms, given in Proposition 1 below. We reuse the sparsity constant introduced in Reddi et al. (2015) and relate it to the one we have defined earlier, Δ_r :

Remark 1. Let D be the smallest constant such that:

$$\mathbb{E}\|x\|_i^2 = \frac{1}{n} \sum_{i=1}^n \|x\|_i^2 \leq D\|x\|^2 \quad \forall x \in \mathbb{R}^d, \quad (39)$$

where $\|\cdot\|_i$ is defined to be the ℓ_2 -norm restricted to the support S_i of f_i . We have:

$$D = \frac{\Delta_r}{n} = \Delta. \quad (40)$$

Proof. We have:

$$\mathbb{E}\|x\|_i^2 = \frac{1}{n} \sum_{i=1}^n \|x\|_i^2 = \frac{1}{n} \sum_{i=1}^n \sum_{v \in S_i} [x]_v^2 = \frac{1}{n} \sum_{v=1}^d \sum_{i|v \in S_i} [x]_v^2 = \frac{1}{n} \sum_{v=1}^d \delta_v [x]_v^2, \quad (41)$$

where $\delta_v := \mathbf{card}(i \mid v \in S_i)$.

This implies:

$$D \geq \frac{1}{n} \sum_{v=1}^d \delta_v \frac{[x]_v^2}{\|x\|^2}. \quad (42)$$

Since D is the minimum constant satisfying this inequality, we have:

$$D = \max_{x \in \mathbb{R}^d} \frac{1}{n} \sum_{v=1}^d \delta_v \frac{[x]_v^2}{\|x\|^2}. \quad (43)$$

We need to find x such that it maximizes the right-hand side term. Note that the vector $([x]_v^2/\|x\|^2)_{v=1..d}$ is in the unit probability simplex, which means that an equivalent problem is the maximization over all convex combinations of $(\delta_v)_{v=1..d}$. This maximum is found by putting all the weight on the maximum δ_v , which is Δ_r by definition.

This means that $\Delta = \Delta_r/n$ is indeed the smallest constant satisfying (39). \square

Proposition 1. For any $u \neq t$,

$$\mathbb{E}|\langle g_u, g_t \rangle| \leq \frac{\sqrt{\Delta}}{2} (\mathbb{E}\|g_u\|^2 + \mathbb{E}\|g_t\|^2). \quad (44)$$

Proof. Let $u \neq t$. Without loss of generality, $u < t$.¹⁴ Then:

$$\begin{aligned} \mathbb{E}|\langle g_u, g_t \rangle| &\leq \mathbb{E}\|g_u\|_{i_t} \|g_t\| && \text{(Sparse inner product; support of } g_t \text{ is } S_{i_t}) \\ &\leq \sqrt{\mathbb{E}\|g_u\|_{i_t}^2} \sqrt{\mathbb{E}\|g_t\|^2} && \text{(Cauchy-Schwarz for expectations)} \\ &\leq \sqrt{\Delta \mathbb{E}\|g_u\|^2} \sqrt{\mathbb{E}\|g_t\|^2} && \text{(Remark 1 and } i_t \perp\!\!\!\perp g_u, \forall u < t) \\ &\leq \frac{\sqrt{\Delta}}{2} (\mathbb{E}\|g_u\|^2 + \mathbb{E}\|g_t\|^2). && \text{(AM-GM inequality)} \end{aligned}$$

¹⁴One only has to switch u and t if $u > t$.

All told, we have:

$$\mathbb{E}|\langle g_u, g_t \rangle| \leq \frac{\sqrt{\Delta}}{2}(\mathbb{E}\|g_u\|^2 + \mathbb{E}\|g_t\|^2). \quad (45)$$

□

Bounding $\mathbb{E}\langle \hat{x}_t - x_t, g_t \rangle$ in terms of g_u .

$$\begin{aligned} \frac{1}{\gamma}\mathbb{E}\langle \hat{x}_t - x_t, g_t \rangle &= \sum_{u=(t-\tau)_+}^{t-1} \mathbb{E}\langle G_u^t g_u, g_t \rangle && \text{(by Equation (9))} \\ &\leq \sum_{u=(t-\tau)_+}^{t-1} \mathbb{E}|\langle g_u, g_t \rangle| && (G_u^t \text{ diagonal matrices with terms in } \{0, 1\}) \\ &\leq \sum_{u=(t-\tau)_+}^{t-1} \frac{\sqrt{\Delta}}{2}(\mathbb{E}\|g_u\|^2 + \mathbb{E}\|g_t\|^2) && \text{(by Proposition 1)} \\ &\leq \frac{\sqrt{\Delta}}{2} \sum_{u=(t-\tau)_+}^{t-1} \mathbb{E}\|g_u\|^2 + \frac{\sqrt{\Delta}\tau}{2}\mathbb{E}\|g_t\|^2. && (46) \end{aligned}$$

Bounding $\mathbb{E}\|\hat{x}_t - x_t\|^2$ with respect to g_u Thanks to the expansion for $\hat{x}_t - x_t$ (9), we get:

$$\|\hat{x}_t - x_t\|^2 \leq \gamma^2 \sum_{u,v=(t-\tau)_+}^{t-1} |\langle G_u^t g_u, G_v^t g_v \rangle| \leq \gamma^2 \sum_{u=(t-\tau)_+}^{t-1} \|g_u\|^2 + \gamma^2 \sum_{\substack{u,v=(t-\tau)_+ \\ u \neq v}}^{t-1} |\langle G_u^t g_u, G_v^t g_v \rangle|.$$

Using (44) from Proposition 1, we have that for $u \neq v$:

$$\mathbb{E}|\langle G_u^t g_u, G_v^t g_v \rangle| \leq \mathbb{E}|\langle g_u, g_v \rangle| \leq \frac{\sqrt{\Delta}}{2}(\mathbb{E}\|g_u\|^2 + \mathbb{E}\|g_v\|^2). \quad (47)$$

By taking the expectation and using (47), we get:

$$\begin{aligned} \mathbb{E}\|\hat{x}_t - x_t\|^2 &\leq \gamma^2 \sum_{u=(t-\tau)_+}^{t-1} \mathbb{E}\|g_u\|^2 + \gamma^2 \sqrt{\Delta}(\tau - 1)_+ \sum_{u=(t-\tau)_+}^{t-1} \mathbb{E}\|g_u\|^2 \\ &= \gamma^2 (1 + \sqrt{\Delta}(\tau - 1)_+) \sum_{u=(t-\tau)_+}^{t-1} \mathbb{E}\|g_u\|^2 \\ &\leq \gamma^2 (1 + \sqrt{\Delta}\tau) \sum_{u=(t-\tau)_+}^{t-1} \mathbb{E}\|g_u\|^2. && (48) \end{aligned}$$

We can now rewrite (23) in terms of $\mathbb{E}\|g_t\|^2$, which finishes the proof for Lemma 1 (by introducing C_1

and C_2 as specified in Lemma 1):

$$\begin{aligned}
 a_{t+1} &\leq \left(1 - \frac{\gamma\mu}{2}\right)a_t - 2\gamma e_t + \gamma^2 \mathbb{E}\|g_t\|^2 + \gamma^3 \mu(1 + \sqrt{\Delta}\tau) \sum_{u=(t-\tau)_+}^{t-1} \mathbb{E}\|g_u\|^2 \\
 &\quad + \gamma^2 \sqrt{\Delta} \sum_{u=(t-\tau)_+}^{t-1} \mathbb{E}\|g_u\|^2 + \gamma^2 \sqrt{\Delta}\tau \mathbb{E}\|g_t\|^2 \\
 &\leq \left(1 - \frac{\gamma\mu}{2}\right)a_t - 2\gamma e_t + \gamma^2 C_1 \mathbb{E}\|g_t\|^2 + \gamma^2 C_2 \sum_{u=(t-\tau)_+}^{t-1} \mathbb{E}\|g_u\|^2.
 \end{aligned} \tag{49}$$

□

C.5 Proof of Lemma 2

We now derive our bound on g_t with respect to suboptimality. From Appendix B, we know that:

$$\mathbb{E}\|g_t\|^2 \leq 2\mathbb{E}\|f'_{i_t}(\hat{x}_t) - f'_{i_t}(x^*)\|^2 + 2\mathbb{E}\|\hat{\alpha}_{i_t}^t - f'_{i_t}(x^*)\|^2 \tag{50}$$

$$\mathbb{E}\|f'_{i_t}(\hat{x}_t) - f'_{i_t}(x^*)\|^2 \leq 2L(\mathbb{E}f(\hat{x}_t) - f(x^*)) = 2Le_t. \tag{51}$$

N. B.: In the following, i_t is a random variable picked uniformly at random in $\{1, \dots, n\}$, whereas i is a fixed constant.

We still have to handle the $\mathbb{E}\|\hat{\alpha}_{i_t}^t - f'_{i_t}(x^*)\|^2$ term and express it in terms of past suboptimalities. We know from our definition of t that i_t and \hat{x}_u are independent $\forall u < t$. Given the “after read” global ordering, \mathbf{E} – the expectation on i_t conditioned on \hat{x}_t and all “past” \hat{x}_u and i_u – is well defined, and we can rewrite our quantity as:

$$\begin{aligned}
 \mathbb{E}\|\hat{\alpha}_{i_t}^t - f'_{i_t}(x^*)\|^2 &= \mathbb{E}(\mathbf{E}\|\hat{\alpha}_{i_t}^t - f'_{i_t}(x^*)\|^2) = \mathbb{E}\frac{1}{n} \sum_{i=1}^n \|\hat{\alpha}_i^t - f'_i(x^*)\|^2 \\
 &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}\|\hat{\alpha}_i^t - f'_i(x^*)\|^2.
 \end{aligned}$$

Now, with i fixed, let $u_{i,l}^t$ be the time of the iterate last used to write the $[\hat{\alpha}_i^t]_l$ quantity, i.e. $[\hat{\alpha}_i^t]_l = [f'_i(\hat{x}_{u_{i,l}^t})]_l$. We know¹⁵ that $0 \leq u_{i,l}^t \leq t-1$. To use this information, we first need to split $\hat{\alpha}_i$ along its dimensions to handle the possible inconsistencies among them:

$$\mathbb{E}\|\hat{\alpha}_i^t - f'_i(x^*)\|^2 = \mathbb{E} \sum_{l=1}^d ([\hat{\alpha}_i^t]_l - [f'_i(x^*)]_l)^2 = \sum_{l=1}^d \mathbb{E} \left[([\hat{\alpha}_i^t]_l - [f'_i(x^*)]_l)^2 \right].$$

This gives us:

$$\begin{aligned}
 \mathbb{E}\|\hat{\alpha}_i^t - f'_i(x^*)\|^2 &= \sum_{l=1}^d \mathbb{E} \left[(f'_i(\hat{x}_{u_{i,l}^t})_l - f'_i(x^*)_l)^2 \right] \\
 &= \sum_{l=1}^d \mathbb{E} \left[\sum_{u=0}^{t-1} \mathbb{1}_{\{u_{i,l}^t=u\}} (f'_i(\hat{x}_u)_l - f'_i(x^*)_l)^2 \right] \\
 &= \sum_{u=0}^{t-1} \sum_{l=1}^d \mathbb{E} \left[\mathbb{1}_{\{u_{i,l}^t=u\}} (f'_i(\hat{x}_u)_l - f'_i(x^*)_l)^2 \right].
 \end{aligned} \tag{52}$$

¹⁵In the case where $u = 0$, one would have to replace the partial gradient with α_i^0 . We omit this special case here for clarity of exposition.

We will now rewrite the indicator so as to obtain independent events from the rest of the equality. This will enable us to distribute the expectation. Suppose $u > 0$ ($u = 0$ is a special case which we will handle afterwards). $\{u_{i,l}^t = u\}$ requires two things:

1. at time u , i was picked uniformly at random,
2. (roughly) i was not picked again between u and t .

We need to refine both conditions because we have to account for possible collisions due to asynchrony. We know from our definition of τ that the t^{th} iteration finishes before at $t + \tau + 1$, but it may still be unfinished by time $t + \tau$. This means that we can only be sure that an update selecting i at time v has been written to memory at time t if $v \leq t - \tau - 1$. Later updates may not have been written yet at time t . Similarly, updates before $v = u + \tau + 1$ may be overwritten by the u^{th} update so we cannot infer that they did not select i . From this discussion, we conclude that $u_{i,l}^t = u$ implies that $i_v \neq i$ for all v between $u + \tau + 1$ and $t - \tau - 1$, though it can still happen that $i_v = i$ for v outside this range.

Using the fact that i_u and i_v are independent for $v \neq u$, we can thus upper bound the indicator function appearing in (52) as follows:¹⁶

$$\mathbb{1}_{\{u_{i,l}^t = u\}} \leq \mathbb{1}_{\{i_u = i\}} \mathbb{1}_{\{i_v \neq i \ \forall v \text{ s.t. } u + \tau + 1 \leq v \leq t - \tau - 1\}}. \quad (53)$$

This gives us:

$$\begin{aligned} & \mathbb{E} \left[\mathbb{1}_{\{u_{i,l}^t = u\}} (f'_i(\hat{x}_u)_l - f'_i(x^*)_l)^2 \right] \\ & \leq \mathbb{E} \left[\mathbb{1}_{\{i_u = i\}} \mathbb{1}_{\{i_v \neq i \ \forall v \text{ s.t. } u + \tau + 1 \leq v \leq t - \tau - 1\}} (f'_i(\hat{x}_u)_l - f'_i(x^*)_l)^2 \right] \\ & \leq P\{i_u = i\} P\{i_v \neq i \ \forall v \text{ s.t. } u + \tau + 1 \leq v \leq t - \tau - 1\} \mathbb{E} (f'_i(\hat{x}_u)_l - f'_i(x^*)_l)^2 \quad (i_v \perp\!\!\!\perp \hat{x}_u, \forall v \geq u) \\ & \leq \frac{1}{n} \left(1 - \frac{1}{n}\right)^{(t-2\tau-u-1)_+} \mathbb{E} (f'_i(\hat{x}_u)_l - f'_i(x^*)_l)^2. \end{aligned} \quad (54)$$

Note that the third line used the crucial independence assumption $i_v \perp\!\!\!\perp \hat{x}_u, \forall v \geq u$ arising from our ‘‘After Read’’ ordering. Summing over all dimensions l , we then get:

$$\mathbb{E} \left[\mathbb{1}_{\{u_{i,l}^t = u\}} \|f'_i(\hat{x}_u) - f'_i(x^*)\|^2 \right] \leq \frac{1}{n} \left(1 - \frac{1}{n}\right)^{(t-2\tau-u-1)_+} \mathbb{E} \|f'_i(\hat{x}_u) - f'_i(x^*)\|^2. \quad (55)$$

So now:

$$\begin{aligned} \mathbb{E} \|\hat{\alpha}_{i_t}^t - f'_{i_t}(x^*)\|^2 - \lambda \tilde{e}_0 & \leq \frac{1}{n} \sum_{i=1}^n \sum_{u=1}^{t-1} \frac{1}{n} \left(1 - \frac{1}{n}\right)^{(t-2\tau-u-1)_+} \mathbb{E} \|f'_i(\hat{x}_u) - f'_i(x^*)\|^2 \\ & = \sum_{u=1}^{t-1} \frac{1}{n} \left(1 - \frac{1}{n}\right)^{(t-2\tau-u-1)_+} + \frac{1}{n} \sum_{i=1}^n \mathbb{E} \|f'_i(\hat{x}_u) - f'_i(x^*)\|^2 \\ & = \sum_{u=1}^{t-1} \frac{1}{n} \left(1 - \frac{1}{n}\right)^{(t-2\tau-u-1)_+} + \mathbb{E} \left(\mathbf{E} \|f'_{i_u}(\hat{x}_u) - f'_{i_u}(x^*)\|^2 \right) \quad (i_u \perp\!\!\!\perp \hat{x}_u) \\ & \leq \frac{2L}{n} \sum_{u=1}^{t-1} \left(1 - \frac{1}{n}\right)^{(t-2\tau-u-1)_+} e_u \quad (\text{by Equation (51)}) \\ & = \frac{2L}{n} \sum_{u=1}^{(t-2\tau-1)_+} \left(1 - \frac{1}{n}\right)^{t-2\tau-u-1} e_u + \frac{2L}{n} \sum_{u=\max(1, t-2\tau)}^{t-1} e_u. \end{aligned} \quad (56)$$

¹⁶In the simpler case of the variant of SVRG from Hofmann et al. (2015) as described in C.2, the batch gradient computations are fully synchronized. This means that we can write much the same inequality without having to worry about possible overwrites, thus replacing $\mathbb{1}_{\{i_v \neq i \ \forall v \text{ s.t. } u + \tau + 1 \leq v \leq t - \tau - 1\}}$ by $\mathbb{1}_{\{i_v \neq i \ \forall v \text{ s.t. } u + 1 \leq v \leq t\}}$.

Note that we have excluded \tilde{e}_0 from our formula, using a generic λ multiplier. We need to treat the case $u = 0$ differently to bound $\mathbb{1}_{\{u_{i_t}^t = u\}}$. Because all our initial α_i are initialized to a fixed α_i^0 , $\{u_i^t = 0\}$ just means that i has not been picked between 0 and $t - \tau - 1$, i.e. $\{i_v \neq i \forall v \text{ s.t. } 0 \leq v \leq t - \tau - 1\}$. This means that the $\mathbb{1}_{\{i_u = i\}}$ term in (53) disappears and thus we lose a $\frac{1}{n}$ factor compared to the case where $u > 1$.

Let us now evaluate λ . We have:

$$\begin{aligned} \mathbb{E} \left[\mathbb{1}_{\{u_i^t = 0\}} \|\alpha_i^0 - f'_i(x^*)\|^2 \right] &\leq \mathbb{E} \left[\mathbb{1}_{\{i_v \neq i \forall v \text{ s.t. } 0 \leq v \leq t - \tau - 1\}} \|\alpha_i^0 - f'_i(x^*)\|^2 \right] \\ &\leq P\{i_v \neq i \forall v \text{ s.t. } 0 \leq v \leq t - \tau - 1\} \mathbb{E} \|\alpha_i^0 - f'_i(x^*)\|^2 \\ &\leq \left(1 - \frac{1}{n}\right)^{(t-\tau)_+} \mathbb{E} \|\alpha_i^0 - f'_i(x^*)\|^2. \end{aligned} \quad (57)$$

Plugging (56) and (57) into (50), we get Lemma 2:

$$\mathbb{E} \|g_t\|^2 \leq 4Le_t + \frac{4L}{n} \sum_{u=1}^{t-1} \left(1 - \frac{1}{n}\right)^{(t-2\tau-u-1)_+} e_u + 4L \left(1 - \frac{1}{n}\right)^{(t-\tau)_+} \tilde{e}_0, \quad (58)$$

where we have introduced $\tilde{e}_0 := \frac{1}{2L} \mathbb{E} \|\alpha_i^0 - f'_i(x^*)\|^2$. Note that in the original SAGA algorithm, a batch gradient is computed to set the $\alpha_i^0 = f'_i(x_0)$. In this setting, we can write Lemma 2 using $\tilde{e}_0 \leq e_0$ thanks to (51). In the more general setting where we initialize all α_i^0 to a fixed quantity, we cannot use (51) to bound $\mathbb{E} \|\alpha_i^0 - f'_i(x^*)\|^2$ which means that we have to introduce \tilde{e}_0 .

C.6 Master inequality derivation

Now, if we combine the bound on $\mathbb{E} \|g_t\|^2$ which we just derived (i.e. Lemma 2) with Lemma 1, we get:

$$\begin{aligned} a_{t+1} &\leq \left(1 - \frac{\gamma\mu}{2}\right) a_t - 2\gamma e_t \\ &\quad + 4L\gamma^2 C_1 e_t + \frac{4L\gamma^2 C_1}{n} \sum_{u=1}^{t-1} \left(1 - \frac{1}{n}\right)^{(t-2\tau-u-1)_+} e_u + 4L\gamma^2 C_1 \left(1 - \frac{1}{n}\right)^{(t-\tau)_+} \tilde{e}_0 \\ &\quad + 4L\gamma^2 C_2 \sum_{u=(t-\tau)_+}^{t-1} e_u + 4L\gamma^2 C_2 \sum_{u=(t-\tau)_+}^{t-1} \left(1 - \frac{1}{n}\right)^{(u-\tau)_+} \tilde{e}_0 \\ &\quad + \frac{4L\gamma^2 C_2}{n} \sum_{u=(t-\tau)_+}^{t-1} \sum_{v=1}^{u-1} \left(1 - \frac{1}{n}\right)^{(u-2\tau-v-1)_+} e_v. \end{aligned} \quad (59)$$

If we define $H_t := \sum_{u=1}^{t-1} \left(1 - \frac{1}{n}\right)^{(t-2\tau-u-1)_+} e_u$, then we get:

$$\begin{aligned} a_{t+1} &\leq \left(1 - \frac{\gamma\mu}{2}\right) a_t - 2\gamma e_t \\ &\quad + 4L\gamma^2 C_1 \left(e_t + \left(1 - \frac{1}{n}\right)^{(t-\tau)_+} \tilde{e}_0\right) + \frac{4L\gamma^2 C_1}{n} H_t \\ &\quad + 4L\gamma^2 C_2 \sum_{u=(t-\tau)_+}^{t-1} \left(e_u + \left(1 - \frac{1}{n}\right)^{(u-\tau)_+} \tilde{e}_0\right) + \frac{4L\gamma^2 C_2}{n} \sum_{u=(t-\tau)_+}^{t-1} H_u, \end{aligned} \quad (60)$$

which is the master inequality (28).

C.7 Lyapunov function and associated recursive inequality

We define $\mathcal{L}_t := \sum_{u=0}^t (1-\rho)^{t-u} a_u$ for some target contraction rate $\rho < 1$ to be defined later. We have:

$$\mathcal{L}_{t+1} = (1-\rho)^{t+1} a_0 + \sum_{u=1}^{t+1} (1-\rho)^{t+1-u} a_u = (1-\rho)^{t+1} a_0 + \sum_{u=0}^t (1-\rho)^{t-u} a_{u+1}. \quad (61)$$

We now use our new bound on a_{t+1} , (60):

$$\begin{aligned} \mathcal{L}_{t+1} &\leq (1-\rho)^{t+1} a_0 + \sum_{u=0}^t (1-\rho)^{t-u} \left[\left(1 - \frac{\gamma\mu}{2}\right) a_u - 2\gamma e_u + 4L\gamma^2 C_1 \left(e_u + \left(1 - \frac{1}{n}\right)^{(u-\tau)_+} \tilde{e}_0\right) \right. \\ &\quad \left. + \frac{4L\gamma^2 C_1}{n} H_u + \frac{4L\gamma^2 C_2}{n} \sum_{v=(u-\tau)_+}^{u-1} H_v \right. \\ &\quad \left. + 4L\gamma^2 C_2 \sum_{v=(u-\tau)_+}^{u-1} \left(e_v + \left(1 - \frac{1}{n}\right)^{(v-\tau)_+} \tilde{e}_0\right) \right] \\ &\leq (1-\rho)^{t+1} a_0 + \left(1 - \frac{\gamma\mu}{2}\right) \mathcal{L}_t \\ &\quad + \sum_{u=0}^t (1-\rho)^{t-u} \left[-2\gamma e_u + 4L\gamma^2 C_1 \left(e_u + \left(1 - \frac{1}{n}\right)^{(u-\tau)_+} \tilde{e}_0\right) \right. \\ &\quad \left. + \frac{4L\gamma^2 C_1}{n} H_u + \frac{4L\gamma^2 C_2}{n} \sum_{v=(u-\tau)_+}^{u-1} H_v \right. \\ &\quad \left. + 4L\gamma^2 C_2 \sum_{v=(u-\tau)_+}^{u-1} \left(e_v + \left(1 - \frac{1}{n}\right)^{(v-\tau)_+} \tilde{e}_0\right) \right]. \quad (62) \end{aligned}$$

We can now rearrange the sums to expose a simple sum of e_u multiplied by factors r_u^t :

$$\mathcal{L}_{t+1} \leq (1-\rho)^{t+1} a_0 + \left(1 - \frac{\gamma\mu}{2}\right) \mathcal{L}_t + \sum_{u=1}^t r_u^t e_u + r_0^t \tilde{e}_0. \quad (63)$$

C.8 Proof of Lemma 3

We want to make explicit what conditions on ρ and γ are necessary to ensure that r_u^t is negative for all $u \geq 1$. Since each e_u is positive, we will then be able to safely drop the sum term from the inequality. The r_0^t term is a bit trickier and is handled separately. Indeed, trying to enforce that r_0^t is negative results in a significantly worse condition on γ and eventually a convergence rate smaller by a factor of n than our final result. Instead, we handle this term directly in the Lyapunov function.

Computation of r_u^t . Let's now make the multiplying factor explicit. We assume $u \geq 1$.

We split r_u^t into five parts coming from (62):

- r_1 , the part coming from the $-2\gamma e_u$ terms;
- r_2 , coming from $4L\gamma^2 C_1 e_u$;
- r_3 , coming from $\frac{4L\gamma^2 C_1}{n} H_u$;

- r_4 , coming from $4L\gamma^2 C_2 \sum_{v=(u-\tau)_+}^{u-1} e_v$;
- r_5 , coming from $\frac{4L\gamma^2 C_2}{n} \sum_{v=(u-\tau)_+}^{u-1} H_v$.

r_1 is easy to derive. Each of these terms appears only in one inequality. So for u at time t , the term is:

$$r_1 = -2\gamma(1 - \rho)^{t-u}. \quad (64)$$

For much the same reasons, r_2 is also easy to derive and is:

$$r_2 = 4L\gamma^2 C_1 (1 - \rho)^{t-u}. \quad (65)$$

r_3 is a bit trickier, because for a given $v > 0$ there are several H_u which contain e_v . The key insight is that we can rewrite our double sum in the following manner:

$$\begin{aligned} & \sum_{u=0}^t (1 - \rho)^{t-u} \sum_{v=1}^{u-1} \left(1 - \frac{1}{n}\right)^{(u-2\tau-v-1)_+} e_v \\ &= \sum_{v=1}^{t-1} e_v \sum_{u=v+1}^t (1 - \rho)^{t-u} \left(1 - \frac{1}{n}\right)^{(u-2\tau-v-1)_+} \\ &\leq \sum_{v=1}^{t-1} e_v \left[\sum_{u=v+1}^{\min(t, v+2\tau)} (1 - \rho)^{t-u} + \sum_{u=v+2\tau+1}^t (1 - \rho)^{t-u} \left(1 - \frac{1}{n}\right)^{u-2\tau-v-1} \right] \\ &\leq \sum_{v=1}^{t-1} e_v \left[2\tau(1 - \rho)^{t-v-2\tau} + (1 - \rho)^{t-v-2\tau-1} \sum_{u=v+2\tau+1}^t q^{u-2\tau-v-1} \right] \\ &\leq \sum_{v=1}^{t-1} (1 - \rho)^{t-v} e_v (1 - \rho)^{-2\tau-1} \left[2\tau + \frac{1}{1-q} \right], \end{aligned} \quad (66)$$

where we have defined:

$$q := \frac{1 - 1/n}{1 - \rho}, \quad \text{with the assumption } \rho < \frac{1}{n}. \quad (67)$$

Note that we have bounded the $\min(t, v + 2\tau)$ term by $v + 2\tau$ in the first sub-sum, effectively adding more positive terms.

This gives us that at time t , for u :

$$r_3 \leq \frac{4L\gamma^2 C_1}{n} (1 - \rho)^{t-u} (1 - \rho)^{-2\tau-1} \left[2\tau + \frac{1}{1-q} \right]. \quad (68)$$

For r_4 we use the same trick:

$$\begin{aligned} & \sum_{u=0}^t (1 - \rho)^{t-u} \sum_{v=(u-\tau)_+}^{u-1} e_v = \sum_{v=0}^{t-1} e_v \sum_{u=v+1}^{\min(t, v+\tau)} (1 - \rho)^{t-u} \\ &\leq \sum_{v=0}^{t-1} e_v \sum_{u=v+1}^{v+\tau} (1 - \rho)^{t-u} \leq \sum_{v=0}^{t-1} e_v \tau (1 - \rho)^{t-v-\tau}. \end{aligned} \quad (69)$$

This gives us that at time t , for u :

$$r_4 \leq 4L\gamma^2 C_2 (1 - \rho)^{t-u} \tau (1 - \rho)^{-\tau}. \quad (70)$$

Finally we compute r_5 which is the most complicated term. Indeed, to find the factor of e_w for a given $w > 0$, one has to compute a triple sum, $\sum_{u=0}^t (1-\rho)^{t-u} \sum_{v=(u-\tau)_+}^{u-1} H_v$. We start by computing the factor of e_w in the inner double sum, $\sum_{v=(u-\tau)_+}^{u-1} H_v$.

$$\sum_{v=(u-\tau)_+}^{u-1} \sum_{w=1}^{v-1} \left(1 - \frac{1}{n}\right)^{(v-2\tau-w-1)_+} e_w = \sum_{w=1}^{u-2} e_w \sum_{v=\max(w+1, u-\tau)}^{u-1} \left(1 - \frac{1}{n}\right)^{(v-2\tau-w-1)_+}. \quad (71)$$

Now there are at most τ terms for each e_w . If $w \leq u - 3\tau - 1$, then the exponent is positive in every term and it is always bigger than $u - 3\tau - 1 - w$, which means we can bound the sum by $\tau(1 - \frac{1}{n})^{u-3\tau-1-w}$. Otherwise we can simply bound the sum by τ . We get:

$$\sum_{v=(u-\tau)_+}^{u-1} H_v \leq \sum_{w=1}^{u-2} \left[\mathbb{1}_{\{u-3\tau \leq w \leq u-2\}} \tau + \mathbb{1}_{\{w \leq u-3\tau-1\}} \tau \left(1 - \frac{1}{n}\right)^{u-3\tau-1-w} \right] e_w. \quad (72)$$

This means that for w at time t :

$$\begin{aligned} r_5 &\leq \frac{4L\gamma^2 C_2}{n} \sum_{u=0}^t (1-\rho)^{t-u} \left[\mathbb{1}_{\{u-3\tau \leq w \leq u-2\}} \tau + \mathbb{1}_{\{w \leq u-3\tau-1\}} \tau \left(1 - \frac{1}{n}\right)^{u-3\tau-1-w} \right] \\ &\leq \frac{4L\gamma^2 C_2}{n} \left[\sum_{u=w+2}^{\min(t, w+3\tau)} \tau (1-\rho)^{t-u} + \sum_{u=w+3\tau+1}^t \tau \left(1 - \frac{1}{n}\right)^{u-3\tau-1-w} (1-\rho)^{t-u} \right] \\ &\leq \frac{4L\gamma^2 C_2}{n} \tau \left[(1-\rho)^{t-w} (1-\rho)^{-3\tau} 3\tau \right. \\ &\quad \left. + (1-\rho)^{t-w} (1-\rho)^{-1-3\tau} \sum_{u=w+3\tau+1}^t \left(1 - \frac{1}{n}\right)^{u-3\tau-1-w} (1-\rho)^{-u+3\tau+1+w} \right] \\ &\leq \frac{4L\gamma^2 C_2}{n} \tau (1-\rho)^{t-w} (1-\rho)^{-3\tau-1} \left(3\tau + \frac{1}{1-q} \right). \end{aligned} \quad (73)$$

By combining the five terms together ((64), (65), (68), (70) and (73)), we get that $\forall u$ s.t. $1 \leq u \leq t$:

$$\begin{aligned} r_u^t &\leq (1-\rho)^{t-u} \left[-2\gamma + 4L\gamma^2 C_1 + \frac{4L\gamma^2 C_1}{n} (1-\rho)^{-2\tau-1} \left(2\tau + \frac{1}{1-q} \right) \right. \\ &\quad \left. + 4L\gamma^2 C_2 \tau (1-\rho)^{-\tau} + \frac{4L\gamma^2 C_2}{n} \tau (1-\rho)^{-3\tau-1} \left(3\tau + \frac{1}{1-q} \right) \right]. \end{aligned} \quad (74)$$

Computation of r_0^t . Recall that we treat the \tilde{e}_0 term separately in Section C.5. The initialization of SAGA creates an initial synchronization, which means that the contribution of \tilde{e}_0 in our bound on $\mathbb{E}\|g_t\|^2$ (58) is roughly n times bigger than the contribution of any e_u for $1 < u < t$.¹⁷ In order to safely handle this term in our Lyapunov inequality, we only need to prove that it is bounded by a reasonable constant. Here again, we split r_0^t in five contributions coming from (62):

- r_1 , the part coming from the $-2\gamma e_u$ terms;
- r_2 , coming from $4L\gamma^2 C_1 e_u$;
- r_3 , coming from $4L\gamma^2 C_1 \left(1 - \frac{1}{n}\right)^{(u-\tau)_+} \tilde{e}_0$;

¹⁷This is explained in details right before (57).

- r_4 , coming from $4L\gamma^2 C_2 \sum_{v=(u-\tau)_+}^{u-1} e_v$;
- r_5 , coming from $4L\gamma^2 C_2 \sum_{v=(u-\tau)_+}^{u-1} (1 - \frac{1}{n})^{(v-\tau)_+} \tilde{e}_0$.

Note that there is no \tilde{e}_0 in H_t , which is why we can safely ignore these terms here.

We have $r_1 = -2\gamma(1 - \rho)^t$ and $r_2 = 4L\gamma^2 C_1(1 - \rho)^t$.

Let us compute r_3 .

$$\begin{aligned}
 & \sum_{u=0}^t (1 - \rho)^{t-u} (1 - \frac{1}{n})^{(u-\tau)_+} \\
 &= \sum_{u=0}^{\min(t, \tau)} (1 - \rho)^{t-u} + \sum_{u=\tau+1}^t (1 - \rho)^{t-u} (1 - \frac{1}{n})^{u-\tau} \\
 &\leq (\tau + 1)(1 - \rho)^{t-\tau} + (1 - \rho)^{t-\tau} \sum_{u=\tau+1}^t (1 - \rho)^{\tau-u} (1 - \frac{1}{n})^{u-\tau} \\
 &\leq (1 - \rho)^t (1 - \rho)^{-\tau} \left(\tau + 1 + \frac{1}{1 - q} \right). \tag{75}
 \end{aligned}$$

This gives us:

$$r_3 \leq (1 - \rho)^t 4L\gamma^2 C_1 (1 - \rho)^{-\tau} \left(\tau + 1 + \frac{1}{1 - q} \right). \tag{76}$$

We have already computed r_4 for $u > 0$ and the computation is exactly the same for $u = 0$. $r_4 \leq (1 - \rho)^t 4L\gamma^2 C_2 \tau (1 - \rho)^{-\tau}$.

Finally we compute r_5 .

$$\begin{aligned}
 & \sum_{u=0}^t (1 - \rho)^{t-u} \sum_{v=(u-\tau)_+}^{u-1} (1 - \frac{1}{n})^{(v-\tau)_+} \\
 &= \sum_{v=1}^{t-1} \sum_{u=v+1}^{\min(t, v+\tau)} (1 - \rho)^{t-u} (1 - \frac{1}{n})^{(v-\tau)_+} \\
 &\leq \sum_{v=1}^{\min(t-1, \tau)} \sum_{u=v+1}^{v+\tau} (1 - \rho)^{t-u} + \sum_{v=\tau+1}^{t-1} \sum_{u=v+1}^{\min(t, v+\tau)} (1 - \rho)^{t-u} (1 - \frac{1}{n})^{v-\tau} \\
 &\leq \tau^2 (1 - \rho)^{t-2\tau} + \sum_{v=\tau+1}^{t-1} (1 - \frac{1}{n})^{v-\tau} \tau (1 - \rho)^{t-v-\tau} \\
 &\leq \tau^2 (1 - \rho)^{t-2\tau} + \tau (1 - \rho)^t (1 - \rho)^{-2\tau} \sum_{v=\tau+1}^{t-1} (1 - \frac{1}{n})^{v-\tau} \tau (1 - \rho)^{-v+\tau} \\
 &\leq (1 - \rho)^t (1 - \rho)^{-2\tau} \left(\tau^2 + \tau \frac{1}{1 - q} \right). \tag{77}
 \end{aligned}$$

Which means:

$$r_5 \leq (1 - \rho)^t 4L\gamma^2 C_2 (1 - \rho)^{-2\tau} \left(\tau^2 + \tau \frac{1}{1 - q} \right). \tag{78}$$

Putting it all together, we get that: $\forall t \geq 0$

$$r_0^t \leq (1 - \rho)^t \left[\left(-2\gamma + 4L\gamma^2 C_1 + 4L\gamma^2 C_2 \tau (1 - \rho)^{-\tau} \right) \frac{e_0}{\tilde{e}_0} + 4L\gamma^2 C_1 (1 - \rho)^{-\tau} \left(\tau + 1 + \frac{1}{1 - q} \right) + 4L\gamma^2 C_2 \tau (1 - \rho)^{-2\tau} \left(\tau + \frac{1}{1 - q} \right) \right]. \quad (79)$$

Sufficient condition for convergence. We need all $r_u^t, u \geq 1$ to be negative so we can safely drop them from (63). Note that for every u , this is the same condition. We will reduce that condition to a second-order polynomial sign condition. We also remark that since $\gamma \geq 0$, we can upper bound our terms in γ and γ^2 in this upcoming polynomial, which will give us sufficient conditions for convergence.

Now, as γ is part of C_2 , we need to expand it once more to find our conditions. We have:

$$C_1 = 1 + \sqrt{\Delta} \tau; \quad C_2 = \sqrt{\Delta} + \gamma \mu C_1.$$

Dividing the bracket in (74) by γ and rearranging as a second degree polynomial, we get the condition:

$$4L \left(C_1 + \frac{C_1}{n} (1 - \rho)^{-2\tau-1} \left[2\tau + \frac{1}{1 - q} \right] + \left[\sqrt{\Delta} \tau (1 - \rho)^{-\tau} + \frac{\sqrt{\Delta} \tau}{n} (1 - \rho)^{-3\tau-1} \left(3\tau + \frac{1}{1 - q} \right) \right] \right) \gamma + 8\mu C_1 L \tau \left[(1 - \rho)^{-\tau} + \frac{1}{n} (1 - \rho)^{-3\tau-1} \left(3\tau + \frac{1}{1 - q} \right) \right] \gamma^2 + 2 \leq 0. \quad (80)$$

The discriminant of this polynomial is always positive, so γ needs to be between its two roots. The smallest is negative, so the condition is not relevant to our case (where $\gamma > 0$). By solving analytically for the positive root ϕ , we get an upper bound condition on γ that can be used for any overlap τ and guarantee convergence. Unfortunately, for large τ , the upper bound becomes exponentially small because of the presence of τ in the exponent in (80). More specifically, by using the bound $1/(1 - \rho) \leq \exp(2\rho)$ ¹⁸ and thus $(1 - \rho)^{-\tau} \leq \exp(2\tau\rho)$ in (80), we would obtain factors of the form $\exp(\tau/n)$ in the denominator for the root ϕ (recall that $\rho < 1/n$).

Our Lemma 3 is derived instead under the assumption that $\tau \leq \mathcal{O}(n)$, with the constants chosen in order to make the condition (80) more interpretable and to relate our convergence result with the standard SAGA convergence (see Theorem 1). As explained in Appendix E, the assumption that $\tau \leq \mathcal{O}(n)$ appears reasonable in practice. First, by using Bernoulli's inequality, we have:

$$(1 - \rho)^{k\tau} \geq 1 - k\tau\rho \quad \text{for integers } k\tau \geq 0. \quad (81)$$

To get manageable constants, we make the following slightly more restrictive assumptions on the target rate ρ ¹⁹ and overlap τ :²⁰

$$\rho \leq \frac{1}{4n} \quad (82)$$

$$\tau \leq \frac{n}{10}. \quad (83)$$

¹⁸This bound can be derived from the inequality $(1 - x/2) \geq \exp(-x)$ which is valid for $0 \leq x \leq 1.59$.

¹⁹Note that we already expected $\rho < 1/n$.

²⁰This bound on τ is reasonable in practice, see Appendix E.

We then have:

$$\frac{1}{1-q} \leq \frac{4n}{3} \quad (84)$$

$$\frac{1}{1-\rho} \leq \frac{4}{3} \quad (85)$$

$$k\tau\rho \leq \frac{3}{40} \quad \text{for } 1 \leq k \leq 3 \quad (86)$$

$$(1-\rho)^{-k\tau} \leq \frac{1}{1-k\tau\rho} \leq \frac{40}{37} \quad \text{for } 1 \leq k \leq 3 \text{ and by using (81)}. \quad (87)$$

We can now upper bound loosely the three terms in brackets appearing in (80) as follows:

$$(1-\rho)^{-2\tau-1} \left[2\tau + \frac{1}{1-q} \right] \leq 3n \quad (88)$$

$$\sqrt{\Delta}\tau(1-\rho)^{-\tau} + \frac{\sqrt{\Delta}\tau}{n}(1-\rho)^{-3\tau-1} \left(3\tau + \frac{1}{1-q} \right) \leq 4\sqrt{\Delta}\tau \leq 4C_1 \quad (89)$$

$$(1-\rho)^{-\tau} + \frac{1}{n}(1-\rho)^{-3\tau-1} \left(3\tau + \frac{1}{1-q} \right) \leq 4. \quad (90)$$

By plugging (88)–(90) into (80), we get the simpler sufficient condition on γ :

$$-1 + 16LC_1\gamma + 16LC_1\mu\tau\gamma^2 \leq 0. \quad (91)$$

The positive root ϕ is:

$$\phi = \frac{16LC_1(\sqrt{1 + \frac{\mu\tau}{4LC_1}} - 1)}{32LC_1\mu\tau} = \frac{\sqrt{1 + \frac{\mu\tau}{4LC_1}} - 1}{2\mu\tau}. \quad (92)$$

We simplify it further by using the inequality:²¹

$$\sqrt{x} - 1 \geq \frac{x-1}{2\sqrt{x}} \quad \forall x > 0. \quad (93)$$

Using (93) in (92), and recalling that $\kappa := L/\mu$, we get:

$$\phi \geq \frac{1}{16LC_1\sqrt{1 + \frac{\tau}{4\kappa C_1}}}. \quad (94)$$

Since $\frac{\tau}{C_1} = \frac{\tau}{1+\sqrt{\Delta}\tau} \leq \min(\tau, \frac{1}{\sqrt{\Delta}})$, we get that a sufficient condition on our stepsize is:

$$\gamma \leq \frac{1}{16L(1 + \sqrt{\Delta}\tau)\sqrt{1 + \frac{1}{4\kappa} \min(\tau, \frac{1}{\sqrt{\Delta}})}}. \quad (95)$$

Subject to our conditions on γ , ρ and τ , we then have that: $r_u^t \leq 0$ for all u s.t. $1 \leq u \leq t$. This means we can rewrite (63) as:

$$\mathcal{L}_{t+1} \leq (1-\rho)^{t+1}a_0 + \left(1 - \frac{\gamma\mu}{2}\right)\mathcal{L}_t + r_0^t\tilde{e}_0. \quad (96)$$

²¹This inequality can be derived by using the concavity property $f(y) \leq f(x) + (y-x)f'(x)$ on the differentiable concave function $f(x) = \sqrt{x}$ with $y = 1$.

Now, we could finish the proof from this inequality, but it would only give us a convergence result in terms of $a_t = \mathbb{E}\|x_t - x^*\|^2$. A better result would be in terms of the suboptimality at \hat{x}_t (because \hat{x}_t is a real quantity in the algorithm whereas x_t is virtual). Fortunately, to get such a result, we can easily adapt (96).

We make e_t appear on the left side of (96), by adding γ to r_t^t in (63):²²

$$\gamma e_t + \mathcal{L}_{t+1} \leq (1 - \rho)^{t+1} a_0 + (1 - \frac{\gamma\mu}{2}) \mathcal{L}_t + \sum_{u=1}^{t-1} r_u^t e_u + r_0^t \tilde{e}_0 + (r_t^t + \gamma) e_t. \quad (97)$$

We now require the stronger property that $\gamma + r_t^t \leq 0$, which translates to replacing -2γ with $-\gamma$ in (74):

$$0 \geq \left[-\gamma + 4L\gamma^2 C_1 + \frac{4L\gamma^2 C_1}{n} (1 - \rho)^{-2\tau-1} (2\tau + \frac{1}{1-q}) + 4L\gamma^2 C_2 \tau (1 - \rho)^{-\tau} + \frac{4L\gamma^2 C_2}{n} \tau (1 - \rho)^{-3\tau-1} (3\tau + \frac{1}{1-q}) \right]. \quad (98)$$

We can easily derive a new stronger condition on γ under which we can drop all the $e_u, u > 0$ terms in (97):

$$\gamma \leq \gamma^* = \frac{1}{32L(1 + \sqrt{\Delta}\tau) \sqrt{1 + \frac{1}{8\kappa} \min(\tau, \frac{1}{\sqrt{\Delta}})}}, \quad (99)$$

and thus under which we get:

$$\gamma e_t + \mathcal{L}_{t+1} \leq (1 - \rho)^{t+1} a_0 + (1 - \frac{\gamma\mu}{2}) \mathcal{L}_t + r_0^t \tilde{e}_0. \quad (100)$$

This finishes the proof of Lemma 3. □

C.9 Proof of Theorem 2

End of Lyapunov convergence. We continue with the assumptions of Lemma 3 which gave us (100). Thanks to (79), we can also rewrite $r_0^t \leq (1 - \rho)^{t+1} A$, where A is a constant which depends on n, Δ, γ and L but is finite and crucially does not depend on t . In fact, by reusing similar arguments as in C.8, we can show the bound $A \leq \gamma n$ under the assumptions of Lemma 3 (including $\gamma \leq \gamma^*$).²³ We then have:

$$\begin{aligned} \mathcal{L}_{t+1} &\leq \gamma e_t + \mathcal{L}_{t+1} \leq (1 - \frac{\gamma\mu}{2}) \mathcal{L}_t + (1 - \rho)^{t+1} (a_0 + A \tilde{e}_0) \\ &\leq (1 - \frac{\gamma\mu}{2})^{t+1} \mathcal{L}_0 + (a_0 + A \tilde{e}_0) \sum_{k=0}^{t+1} (1 - \rho)^{t+1-k} (1 - \frac{\gamma\mu}{2})^k. \end{aligned} \quad (101)$$

We have two linearly contracting terms. The sum contracts linearly with the minimum geometric rate factor between $\gamma\mu/2$ and ρ . If we define $m := \min(\rho, \gamma\mu/2)$, $M := \max(\rho, \gamma\mu/2)$ and $\rho^* := \nu m$ with

²²We could use any multiplier from 0 to 2γ , but choose γ for simplicity. For this reason and because our analysis of the r_t^t term was loose, we could derive a tighter bound, but it does not change the leading terms.

²³In particular, note that e_0 does not appear in the definition of A because it turns out that the parenthesis group multiplying e_0 in (79) is negative. Indeed, it contains less positive terms than (74) which we showed to be negative under the assumptions from Lemma 3.

$0 < \nu < 1$,²⁴ we then get:²⁵

$$\begin{aligned}
 \gamma e_t &\leq \gamma e_t + \mathcal{L}_{t+1} \leq \left(1 - \frac{\gamma\mu}{2}\right)^{t+1} \mathcal{L}_0 + (a_0 + A\tilde{e}_0) \sum_{k=0}^{t+1} (1-m)^{t+1-k} (1-M)^k \\
 &\leq \left(1 - \frac{\gamma\mu}{2}\right)^{t+1} \mathcal{L}_0 + (a_0 + A\tilde{e}_0) \sum_{k=0}^{t+1} (1-\rho^*)^{t+1-k} (1-M)^k \\
 &\leq \left(1 - \frac{\gamma\mu}{2}\right)^{t+1} \mathcal{L}_0 + (a_0 + A\tilde{e}_0) (1-\rho^*)^{t+1} \sum_{k=0}^{t+1} (1-\rho^*)^{-k} (1-M)^k \\
 &\leq \left(1 - \frac{\gamma\mu}{2}\right)^{t+1} \mathcal{L}_0 + (1-\rho^*)^{t+1} \frac{1}{1-\eta} (a_0 + A\tilde{e}_0) \\
 &\leq (1-\rho^*)^{t+1} \left(a_0 + \frac{1}{1-\eta} (a_0 + A\tilde{e}_0)\right), \tag{102}
 \end{aligned}$$

where $\eta := \frac{1-M}{1-\rho^*}$. We have $\frac{1}{1-\eta} = \frac{1-\rho^*}{M-\rho^*}$.

By taking $\nu = \frac{4}{5}$ and setting $\rho = \frac{1}{4n}$ – its maximal value allowed by the assumptions of Lemma 3 – we get $M \geq \frac{1}{4n}$ and $\rho^* \leq \frac{1}{5n}$, which means $\frac{1}{1-\eta} \leq 20n$.

All told, using $A \leq \gamma n$, we get:

$$e_t \leq (1-\rho^*)^{t+1} \tilde{C}_0, \tag{103}$$

where

$$\tilde{C}_0 := \frac{21n}{\gamma} \left(\|x_0 - x^*\|^2 + \gamma \frac{n}{2L} \mathbb{E} \|\alpha_i^0 - f'_i(x^*)\|^2 \right). \tag{104}$$

Since we set $\rho = \frac{1}{4n}, \nu = \frac{4}{5}$, we have $\nu\rho = \frac{1}{5n}$. Using a stepsize $\gamma = \frac{a}{L}$ as in Theorem 2, we get $\nu\frac{\gamma\mu}{2} = \frac{2a}{5\kappa}$. We thus obtain a geometric rate of $\rho^* = \min\{\frac{1}{5n}, a\frac{2}{5\kappa}\}$, which we simplified to $\frac{1}{5} \min\{\frac{1}{n}, a\frac{1}{\kappa}\}$ in Theorem 2, finishing the proof. We also observe that $\tilde{C}_0 \leq \frac{60n}{\gamma} C_0$, with C_0 defined in Theorem 1. \square

C.10 Proof of Corollary 3 (speedup regimes)

Referring to Hofmann et al. (2015) and our own Theorem 1, the geometric rate factor of SAGA is $\frac{1}{5} \min\{\frac{1}{n}, \frac{a}{\kappa}\}$ for a stepsize of $\gamma = \frac{a}{5L}$. We start by proving the first part of the corollary which considers the step size $\gamma = \frac{a}{L}$ with $a = a^*(\tau)$. We distinguish between two regimes to study the parallel speedup our algorithm obtains and to derive a condition on τ for which we have a linear speedup.

Big Data. In this regime, $n > \kappa$ and the geometric rate factor of sequential SAGA is $\frac{1}{5n}$. To get a linear speedup (up to a constant factor), we need to enforce $\rho^* = \Omega(\frac{1}{n})$. We recall that $\rho^* = \min\{\frac{1}{5n}, a\frac{1}{5\kappa}\}$.

We already have $\frac{1}{5n} = \Omega(\frac{1}{n})$. This means that we need τ to verify $\frac{a^*(\tau)}{5\kappa} = \Omega(\frac{1}{n})$, where $a^*(\tau) = \frac{1}{32(1+\tau\sqrt{\Delta})\xi(\kappa, \Delta, \tau)}$ according to Theorem 2. Recall that $\xi(\kappa, \Delta, \tau) := \sqrt{1 + \frac{1}{8\kappa} \min\{\frac{1}{\sqrt{\Delta}}, \tau\}}$. Up to a constant factor, this means we can give the following sufficient condition:

$$\frac{1}{\kappa \left(1 + \tau\sqrt{\Delta}\right) \xi(\kappa, \Delta, \tau)} = \Omega\left(\frac{1}{n}\right) \tag{105}$$

²⁴ ν is introduced to circumvent the problematic case where ρ and $\gamma\mu/2$ are too close together, which does not prevent the geometric convergence, but makes the constant $\frac{1}{1-\eta}$ potentially very big (in the case both terms are equal, the sum even becomes an annoying linear term in t).

²⁵Note that if $m \neq \rho$, we can perform the index change $t+1-k \rightarrow k$ to get the sum.

i.e.

$$\left(1 + \tau\sqrt{\Delta}\right) \xi(\kappa, \Delta, \tau) = \mathcal{O}\left(\frac{n}{\kappa}\right). \quad (106)$$

We now consider two alternatives, depending on whether κ is bigger than $\frac{1}{\sqrt{\Delta}}$ or not. If $\kappa \geq \frac{1}{\sqrt{\Delta}}$, then $\xi(\kappa, \Delta, \tau) < 2$ and we can rewrite the sufficient condition (106) as:

$$\tau = \mathcal{O}(1) \frac{n}{\kappa\sqrt{\Delta}}. \quad (107)$$

In the alternative case, $\kappa \leq \frac{1}{\sqrt{\Delta}}$. Since $a^*(\tau)$ is decreasing in τ , we can suppose $\tau \geq \frac{1}{\sqrt{\Delta}}$ without loss of generality and thus $\xi(\kappa, \Delta, \tau) = \sqrt{1 + \frac{1}{8\kappa\sqrt{\Delta}}}$. We can then rewrite the sufficient condition (106) as:

$$\begin{aligned} \frac{\tau\sqrt{\Delta}}{\sqrt{\kappa^4\sqrt{\Delta}}} &= \mathcal{O}\left(\frac{n}{\kappa}\right) \\ \tau &= \mathcal{O}(1) \frac{n}{\sqrt{\kappa^4\sqrt{\Delta}}}. \end{aligned} \quad (108)$$

We observe that since we have supposed that $\kappa \leq \frac{1}{\sqrt{\Delta}}$, we have $\sqrt{\kappa\sqrt{\Delta}} \leq \kappa\sqrt{\Delta} \leq 1$, which means that our initial assumption that $\tau < \frac{n}{10}$ is stronger than condition (108).

We can now combine both cases to get the following sufficient condition for the geometric rate factor of ASAGA to be the same order as sequential SAGA when $n > \kappa$:

$$\tau = \mathcal{O}(1) \frac{n}{\kappa\sqrt{\Delta}}; \quad \tau = \mathcal{O}(n). \quad (109)$$

Ill-conditioned regime. In this regime, $\kappa > n$ and the geometric rate factor of sequential SAGA is $a \frac{1}{\kappa}$. Here, to obtain a linear speedup, we need $\rho^* = \mathcal{O}\left(\frac{1}{\kappa}\right)$. Since $\frac{1}{n} > \frac{1}{\kappa}$, all we require is that $\frac{a^*(\tau)}{\kappa} = \Omega\left(\frac{1}{\kappa}\right)$ where $a^*(\tau) = \frac{1}{32(1+\tau\sqrt{\Delta})\xi(\kappa, \Delta, \tau)}$, which reduces to $a^*(\tau) = \Omega(1)$.

We can give the following sufficient condition:

$$\frac{1}{\left(1 + \tau\sqrt{\Delta}\right) \xi(\kappa, \Delta, \tau)} = \Omega(1) \quad (110)$$

Using that $\frac{1}{n} \leq \Delta \leq 1$ and that $\kappa > n$, we get that $\xi(\kappa, \Delta, \tau) \leq 2$, which means our sufficient condition becomes:

$$\begin{aligned} \tau\sqrt{\Delta} &= \mathcal{O}(1) \\ \tau &= \frac{\mathcal{O}(1)}{\sqrt{\Delta}}. \end{aligned} \quad (111)$$

This finishes the proof for the first part of Corollary 3.

Universal stepsize. If $\tau = \mathcal{O}\left(\frac{1}{\sqrt{\Delta}}\right)$, then $\xi(\kappa, \Delta, \tau) = \mathcal{O}(1)$ and $(1 + \tau\sqrt{\Delta}) = \mathcal{O}(1)$, and thus $a^*(\tau) = \Omega(1)$ (for any n and κ). This means that the universal stepsize $\gamma = \Theta(1/L)$ satisfies $\gamma \leq a^*(\tau)$ for any κ , giving the same rate factor $\Omega(\min\{\frac{1}{n}, \frac{1}{\kappa}\})$ that sequential SAGA has, completing the proof for the second part of Corollary 3. \square

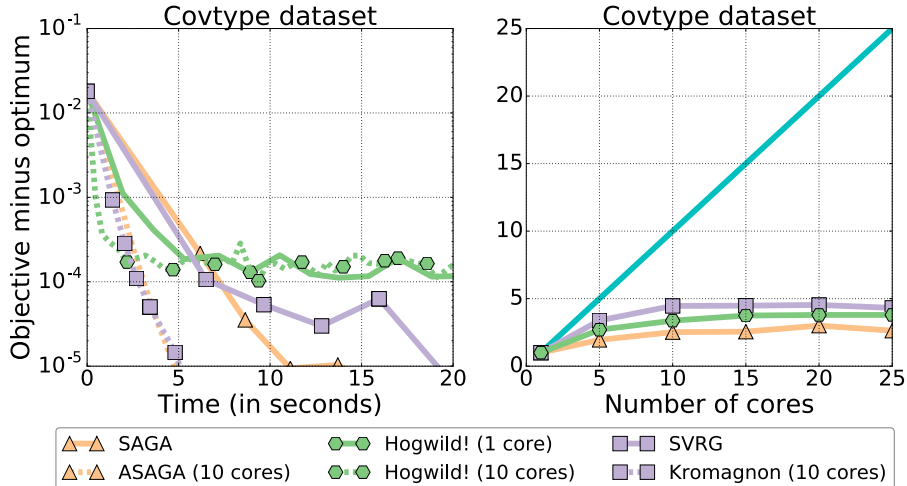


Figure 2: Comparison on the Covtype dataset. Left: suboptimality. Right: speedup. The number of cores in the legend only refers to the left plot.

D Additional experimental results

D.1 Effect of sparsity

Sparsity plays an important role in our theoretical results, where we find that while it is necessary in the “ill-conditioned” regime to get linear speedups, it is not in the “well-conditioned” regime. We confront this to real-life experiments by comparing the convergence and speedup performance of our three asynchronous algorithms on the Covtype dataset, which is fully dense after standardization. The results appear in Figure 2.

While we still see a significant improvement in speed when increasing the number of cores, this improvement is smaller than the one we observe for sparser datasets. The speedups we observe are consequently smaller, and taper off earlier than on our other datasets. However, since the observed “theoretical” speedup is linear (see Section D.2), we can attribute this worse performance to higher hardware overhead. This is expected because each update is fully dense and thus the shared parameters are much more heavily contended for than in our sparse datasets.

One thing we notice when computing the Δ variable for our datasets is that it often fails to capture the full sparsity distribution, being essentially a maximum. This means that Δ can be quite big even for very sparse datasets. Deriving a less coarse bound remains an open problem.

D.2 Theoretical speedups

In the main text of this paper, we show experimental speedup results where suboptimality is a function of the running time. This measure encompasses both theoretical algorithmic properties and hardware overheads (such as contention of shared memory) which are not taken into account in our analysis.

In order to isolate these two effects, we plot our convergence experiments where suboptimality is a function of the number of iterations; thus, we abstract away any potential hardware overhead.²⁶ The experimental results can be seen in Figure 3.

For all three algorithms and all three datasets, the curves for 1 and 10 cores almost coincide, which means that we are indeed in the “theoretical linear speedup” regime. Indeed, when we plotted the

²⁶To do so, we implement a global counter which is sparsely updated (every 100 iterations for example) in order not to modify the asynchrony of the system. This counter is used only for plotting purposes and is not needed otherwise.

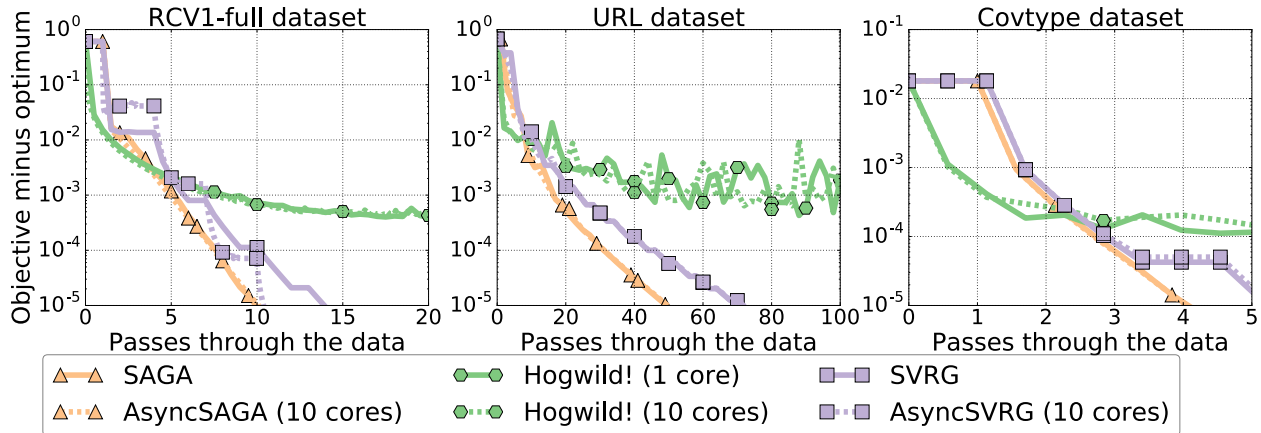


Figure 3: **Theoretical speedups.** Suboptimality with respect to number of iterations for ASAGA SVRG and HOGWILD with 1 and 10 cores. Curves almost coincide, which means the theoretical speedup is almost the number of cores p , hence linear.

amount of iterations required to converge to a given accuracy as a function of the number of cores, we obtained straight horizontal lines for our three algorithms.

The fact that the speedups we observe in running time are less than linear can thus be attributed to various hardware overheads, including shared variable contention – the compare-and-swap operations are more and more expensive as the number of competing requests increases – and cache effects as mentioned in Section 4.3.

E A closer look at the τ constant

E.1 Theory

In the parallel optimization literature, τ is often referred to as a proxy for the number of cores. However, intuitively as well as in practice, it appears that there are a number of other factors that can influence this quantity. We will now attempt to give a few qualitative arguments as to what these other factors might be and how they relate to τ .

Number of cores. The first of these factors is indeed the number of cores.

If we have p cores, $\tau \geq p - 1$. Indeed, in the best-case scenario where all cores have exactly the same execution speed for a single iteration, $\tau = p - 1$.

To get more insight into what τ really encompasses, let us now try to define the worst-case scenario in the preceding example. Consider 2 cores. In the worst case scenario, one core runs while the other is stuck. Then the overlap is t for all t and eventually grows to $+\infty$. If we assume that one core runs twice as fast as the other, then $\tau = 2$. If both run at the same speed, $\tau = 1$.

It appears then that a relevant quantity is R , the ratio between the fastest execution time to the slowest execution time for a single iteration. $\tau \leq (p - 1)R$, which can be arbitrarily bigger than p .

Length of an iteration. There are several factors at play in R itself.

- The first is the speed of execution of the cores themselves (i.e. clock time). The dependency here is quite clear.
- The second is the data matrix itself. If one f_i has support of size n while all the others have support

of size 1, r may eventually become very big.

- The third is the length of the computation itself. The longer our algorithm runs, the more likely it is to explore the potential corner cases of the data matrix.

The overlap is upper bounded by the number of cores times the maximum iteration time over the minimum iteration time (which is linked to the sparsity distribution of the data matrix). This is an upper bound, which means that in some cases it will not really be useful. For example, in the case where one factor has support size 1 and all others have support size d , the probability of the event which corresponds to the upper bound is exponentially small in d . We conjecture that a more useful indicator could be the maximum iteration time over the expected iteration time.

To sum up this preliminary theoretical exploration, the τ term encompasses a lot more complexity than is usually implied in the literature. This is reflected in the experiments we ran, where the constant was orders of magnitude bigger than the number of cores.

E.2 Experimental results

In order to verify our intuition about the τ variable, we ran several experiments on all three datasets, whose characteristics are reminded in Table 1. δ_l^i is the support size of f_i .

Table 1: Density measures including minimum, average and maximum support size δ_l^i of the factors.

	n	d	density	$\max(\delta_l^i)$	$\min(\delta_l^i)$	$\bar{\delta}_l$	$\max(\delta_l^i)/\bar{\delta}_l$
RCV1	697,641	47,236	0.15%	1,224	4	73.2	16.7
URL	2,396,130	3,231,961	0.003%	414	16	115.6	3.58
Covtype	581,012	54	100%	12	8	11.88	1.01

To estimate τ , we compute the average overlap over 100 iterations, which is a lower bound on the actual overlap (which is a maximum, not an average). We then take the maximum observed quantity. We use an average because computing the overlap requires using a global counter, which we do not want to update every iteration since it would make it a heavily contentious quantity susceptible of artificially changing the asynchrony of our algorithm.

The results we observe are order of magnitude bigger than p , indicating that τ can indeed not be dismissed as a mere proxy for the number of cores, but has to be more carefully analyzed.

First, we plot the maximum observed τ as a function of the number of cores (see Figure 4). We observe that the relationship does indeed seem to be roughly linear with respect to the number of cores until 30 cores. After 30 cores, we observe what may be a phase transition where the slope increases significantly.

Second, we measured the maximum observed τ as a function of the number of epochs. We omit the figure since we did not observe any dependency; that is, τ does not seem to depend on the number of epochs. We know that it must depend on the number of iterations (since it cannot be bigger, and is an increasing function with respect to that number for example), but it appears that a stable value is reached quite quickly (before one full epoch is done).

If we allowed the computations to run forever, we would eventually observe an event such that τ would reach the upper bound mentioned in the last section, so it may be that τ is actually a very slowly increasing function of the number of iterations.

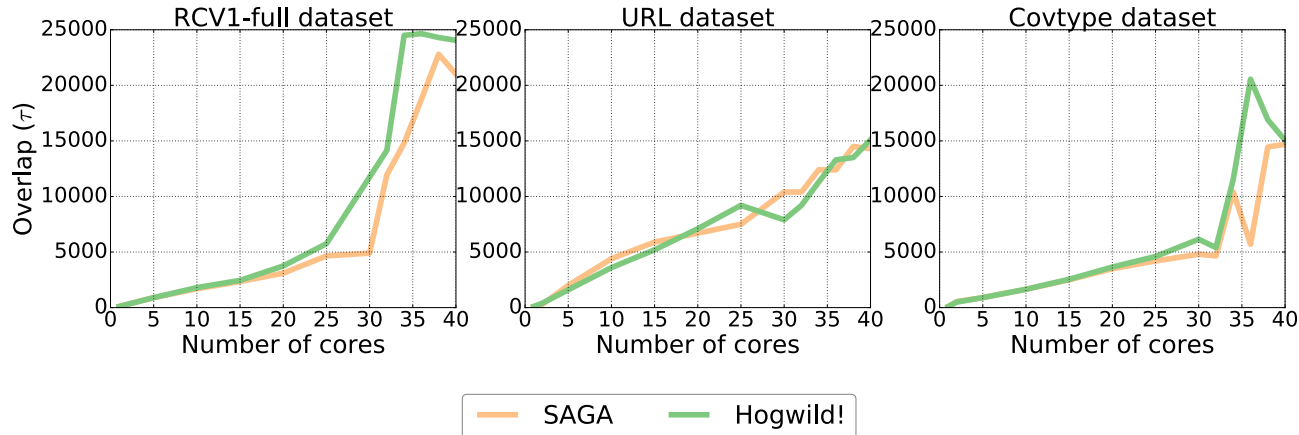


Figure 4: **Overlap**. Overlap as a function of the number of cores for both ASAGA and HOGWILD on all three datasets.

F Lagged updates and Sparsity

F.1 Comparison with Lagged Updates in the sequential case

The lagged updates technique in SAGA is based on the observation that the updates for component $[x]_v$ need not be applied until this coefficient needs to be accessed, that is, until the next iteration t such that $v \in S_{i_t}$. We refer the reader to [Schmidt et al. \(2016\)](#) for more details.

Interestingly, the expected number of iterations between two steps where a given dimension v is involved in the partial gradient is p_v^{-1} , where p_v is the probability that v is involved in a given step. p_v^{-1} is precisely the term which we use to multiply the update to $[x]_v$ in Sparse SAGA. Therefore one may see the updates in Sparse SAGA as *anticipated* updates, whereas those in the [Schmidt et al. \(2016\)](#) implementation are *lagged*. The two algorithms appear to be very close, even though Sparse SAGA uses an expectation to multiply a given update whereas the lazy implementation uses a random variable (with the same expectation). Sparse SAGA therefore uses a slightly more aggressive strategy, which gave faster run-time in our experiments below.

Although Sparse SAGA requires the computation of the p_v probabilities, this can be done during a first pass throughout the data (during which constant step size SGD may be used) at a negligible cost.

In our experiments, we compare the Sparse SAGA variant proposed in Section 2 to two other approaches: the naive (i.e. dense) update scheme and the lagged updates implementation described in [Defazio et al. \(2014\)](#). Note that we use different datasets from the parallel experiments, including a subset of the RCV1 dataset and the realsim dataset. Figure 5 reveals that sparse and lagged updates have a lower cost per iteration, resulting in faster convergence for sparse datasets. Furthermore, while the two approaches had similar convergence in terms of number of iterations, the Sparse SAGA scheme is slightly faster in terms of runtime (and as previously pointed out, sparse updates are better adapted for the asynchronous setting). For the dense dataset (Covtype), the three approaches exhibit a similar performance.

F.2 On the difficulty of parallel lagged updates

In the implementation presented in [Schmidt et al. \(2016\)](#), the dense part ($\bar{\alpha}$) of the updates is deferred. Instead of writing dense updates, counters c_d are kept for each coordinate of the parameter vector – which represent the last time these variables were updated – as well as the average gradient $\bar{\alpha}$ for each coordinate. Then, whenever a component $[\hat{x}]_d$ is needed (in order to compute a new gradient), we subtract $\gamma(t - c_d)[\bar{\alpha}]_d$ from it and c_d is set to t . The reason we can do this without modifying the

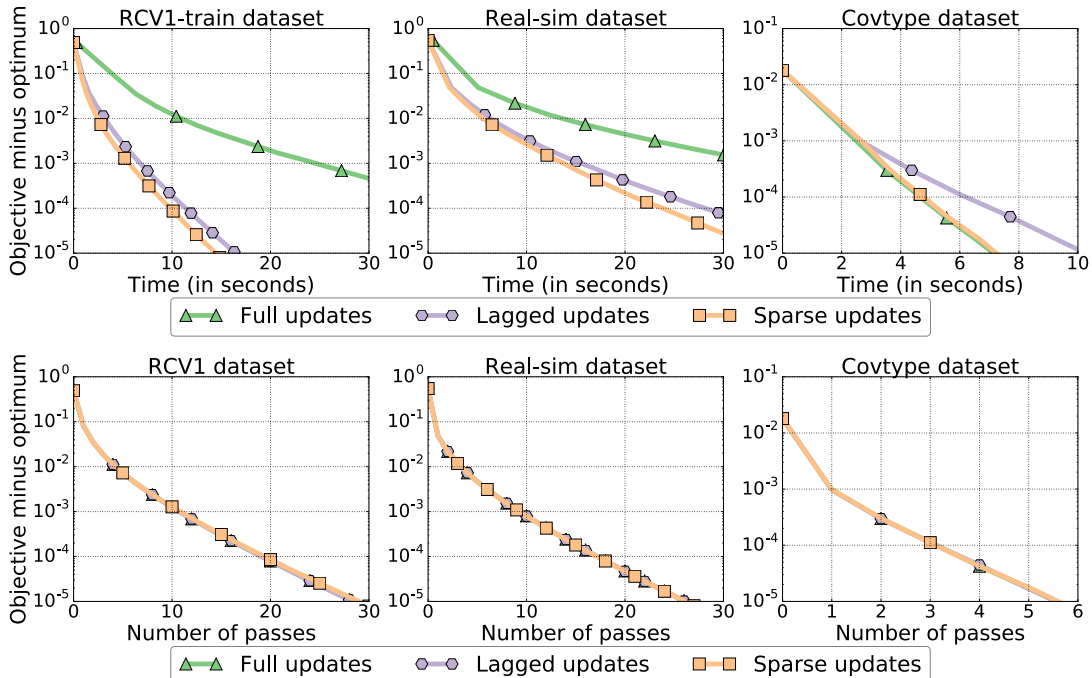


Figure 5: **Lagged vs sparse SAGA updates.** Suboptimality with respect to time for different SAGA update schemes on various datasets. First row: suboptimality as a function of time. Second row: suboptimality as the number of passes over the dataset. For sparse datasets (RCV1 and Real-sim), lagged and sparse updates have a lower cost per iteration which result in faster convergence.

algorithm is that $[\bar{\alpha}]_d$ only changes when $[\hat{x}]_d$ also does.

In the sequential setting, this is strictly the same as doing the updates in a dense way, since the coordinates are only stale when they're not used. Note that at the end of an execution all counters have to be subtracted at once to get the true final parameter vector (and to bring every c_d counter to the final t).

In the parallel setting, several issues arise:

- two cores might be attempting to correct the lag at the same time. In which case since updates are done as additions and not replacements (which is necessary to ensure that there are no overwrites), the lag might be corrected multiple times, i.e. overly corrected.
- we would have to read and write atomically to each $[\hat{x}]_d, c_d, [\bar{\alpha}]_d$ triplet, which is highly impractical.
- we would need to have an explicit global counter, which we do not in ASAGA (our global counter t being used solely for the proof).
- in the dense setting, updates happen coordinate by coordinate. So at time t the number of $\bar{\alpha}$ updates a coordinate has received from a fixed past time c_d is a random variable, which may differ from coordinate to coordinate. Whereas in the lagged implementation, the multiplier is always $(t - c_d)$ which is a constant (conditional to c_d), which means a potentially different \hat{x}_t .

All these points mean both that the implementation of such a scheme in the parallel setting would be impractical, and that it would actually yields a different algorithm than the dense version, which would be even harder to analyze.

G Additional empirical details

G.1 Detailed description of datasets

We run our experiments on four datasets. In every case, we run logistic regression for the purpose of binary classification.

RCV1 ($n = 697,641$, $d = 47,236$). The first is the Reuters Corpus Volume I (RCV1) dataset (Lewis et al., 2004), an archive of over 800,000 manually categorized newswire stories made available by Reuters, Ltd. for research purposes. The associated task is a binary text categorization.

URL ($n = 2,396,130$, $d = 3,231,961$). Our second dataset was first introduced in Ma et al. (2009). Its associated task is a binary malicious url detection. This dataset contains more than 2 million URLs obtained at random from Yahoo’s directory listing (for the “benign” URLs) and from a large Web mail provider (for the “malicious” URLs). The benign to malicious ratio is 2. Features include lexical information as well as metadata. This dataset was obtained from the libsvmtools project.²⁷

Coverttype ($n = 581,012$, $d = 54$). On our third dataset, the associated task is a binary classification problem (down from 7 classes originally, following the pre-treatment of Collobert et al. (2002)). The features are cartographic variables. Contrarily to the first two, this is a dense dataset.

Realsim ($n = 73,218$, $d = 20,958$). We only use our fourth dataset for non-parallel experiments and a specific compare-and-swap test. It constitutes of UseNet articles taken from four discussion groups (simulated auto racing, simulated aviation, real autos, real aviation).

G.2 Implementation details

Hardware. All experiments were run on a Dell PowerEdge 920 machine with 4 Intel Xeon E7-4830v2 processors with 10 2.2GHz cores each and 384GB 1600 Mhz RAM.

Software. All algorithms were implemented in the Scala language and the software stack consisted of a Linux operating system running Scala 2.11.7 and Java 1.6.

We chose this expressive, high level language for our experimentation despite its typical 20x slower performance compared to C because our primary concern was that the code may easily be reused and extended for research purposes (which is harder to achieve with low level, heavily optimized C code; especially for error prone parallel computing).

As a result our timed experiments exhibit sub-optimal running times, e.g. compared to Konecny and Richtarik (2013). This is as we expected. The observed slowdown is both consistent across datasets (roughly 20x) and with other papers that use Scala code (e.g. Mania et al. (2015), Ma et al. (2015, Fig. 2)).

Despite this slowdown, our experiments show state-of-the-art results in convergence per number of iterations. Furthermore, the speed-up patterns that we observe for our implementation of Hogwild and Kromagnon are similar to the ones given in [MN15], Niu et al.[2011] and Reddi et al.[2015] (in various languages).

The code we used to run all the experiments is available at <https://github.com/RemiLeblond/ASAGA>.

Necessity of compare-and-swap operations. Interestingly, we have found necessary to use compare-and-swap instructions in the implementation of ASAGA. In Figure 6, we display suboptimality plots using non-thread safe operations and compare-and-swap (CAS) operations. The non-thread safe version

²⁷<http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/binary.html>

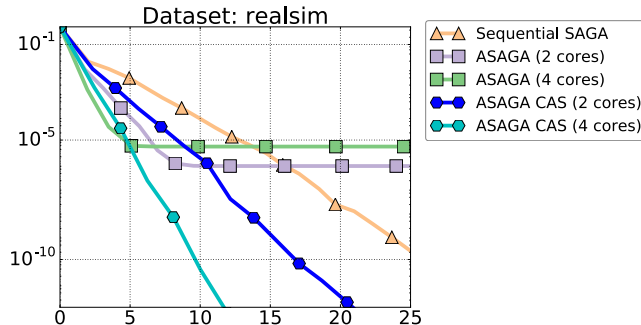


Figure 6: **Compare and swap in the implementation of ASAGA.** Suboptimality as a function of time for ASAGA, both using compare-and-swap (CAS) operations and using standard operations. The graph reveals that CAS is indeed needed in a practical implementation to ensure convergence to a high precision.

starts faster but then fails to converge beyond a specific level of suboptimality, while the compare-and-swap version does converge linearly up to machine precision.

For *compare-and-swap* instructions we used the `AtomicDoubleArray` class from the Google library `Guava`. This class uses an `AtomicLongArray` under the hood (from package `java.util.concurrent.atomic` in the standard Java library), which does indeed benefit from lower-level CPU-optimized instructions.

Efficient storage of the α_i . Storing n gradient may seem like an expensive proposition, but for linear predictor models, one can actually store a single scalar per gradient (as proposed in Schmidt et al. (2016)), which is what we do in our implementation of ASAGA.

Step sizes. For each algorithm, we picked the best step size among 10 equally spaced values in a grid, and made sure that the best step size was never at the boundary of this interval. For Covtype and RCV1, we used the interval $[\frac{1}{10L}, \frac{10}{L}]$, whereas for URL we used the interval $[\frac{1}{L}, \frac{100}{L}]$ as it admitted larger step sizes. It turns out that the best step size was fairly constant for different number of cores for both ASAGA and KROMAGNON, and both algorithms had similar best step sizes.

G.3 Biased update in the implementation

In the implementation detailed in Algorithm 2, $\bar{\alpha}$ is maintained in memory instead of being recomputed for every iteration. This saves both the cost of reading every data point for each iteration and of computing $\bar{\alpha}$ for each iteration.

However, this removes the unbiasedness guarantee. The problem here is the definition of the expectation of $\hat{\alpha}_i$. Since we are sampling uniformly at random, the average of the $\hat{\alpha}_i$ is taken at the precise moment when we read the α_i^t components. Without synchronization, between two reads to a single coordinate in α_i and in $\bar{\alpha}$, new updates might arrive in $\bar{\alpha}$ that are not yet taken into account in α_i . Conversely, writes to a component of α_i might precede the corresponding write in $\bar{\alpha}$ and induce another source of bias.

In order to alleviate this issue, we can use coordinate-level locks on α_i and $\bar{\alpha}$ to make sure they are always synchronized. Such low-level locks are quite inexpensive when d is large, especially when compared to vector-wide locks.

However, as previously noted, experimental results indicate that this fix is not necessary.