



HAL
open science

DADAO: Decoupled Accelerated Decentralized Asynchronous Optimization

Adel Nabli, Edouard Oyallon

► **To cite this version:**

Adel Nabli, Edouard Oyallon. DADAO: Decoupled Accelerated Decentralized Asynchronous Optimization. 2023. hal-03737694v2

HAL Id: hal-03737694

<https://hal.science/hal-03737694v2>

Preprint submitted on 14 Feb 2023 (v2), last revised 15 Nov 2023 (v3)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

DADAO: Decoupled Accelerated Decentralized Asynchronous Optimization

Adel Nabli

ISIR, CNRS, Sorbonne University
MILA, Concordia University
adel.nabli@sorbonne-universite.fr

Edouard Oyallon

ISIR, CNRS, Sorbonne University

February 14, 2023

Abstract

This work introduces DADAO: the first decentralized, accelerated, asynchronous, primal, first-order algorithm to minimize a sum of L -smooth and μ -strongly convex functions distributed over a given network of size n . Our key insight is based on modeling the local gradient updates and gossip communication procedures with separate independent Poisson Point Processes. This allows us to decouple the computation and communication steps, which can be run in parallel, while making the whole approach completely asynchronous, leading to communication acceleration compared to synchronous approaches. Our new method employs primal gradients and does not use a multi-consensus inner loop nor other ad-hoc mechanisms such as Error Feedback, Gradient Tracking, or a Proximal operator. By relating the inverse of the smallest positive eigenvalue of the Laplacian matrix χ_1 and the maximal resistance $\chi_2 \leq \chi_1$ of the graph to a sufficient minimal communication rate between the nodes of the network, we show that our algorithm requires $\mathcal{O}(n\sqrt{\frac{L}{\mu}} \log(\frac{1}{\epsilon}))$ local gradients and only $\mathcal{O}(n\sqrt{\chi_1\chi_2}\sqrt{\frac{L}{\mu}} \log(\frac{1}{\epsilon}))$ communications to reach a precision ϵ , up to logarithmic terms. Thus, we simultaneously obtain an accelerated rate for both computations and communications, leading to an improvement over state-of-the-art works, our simulations further validating the strength of our relatively unconstrained method. We also propose a SDP relaxation to find the optimal gossip rate of each edge minimizing the total number of communications for a given graph, resulting in faster convergence compared to standard approaches relying on uniform communication weights. Our source code is released on a public repository.

1 Introduction

In recent years, the increased amount of available data as well as the proliferation of highly-parallelizable and connected hardware have brought significant changes in the way we process data. These developments have led to the need for efficient and scalable methods for distributed optimization, particularly in the context of machine learning. Indeed, in scenarios where data is distributed across multiple nodes, such as in edge computing or distributed sensor networks, leveraging the local resources of each device is a topic of significant interest. In other settings, such as clusters, spreading the compute load is ideally done to obtain a linear speedup in the number of nodes. In a typical distributed training framework, the goal is to minimize a sum of functions $(f_i)_{i \leq n}$ split across n nodes of a computer network. A corresponding optimization procedure involves alternating local computations on the nodes and communications along the edges \mathcal{E} of the network. In the decentralized setting, there is no central machine aggregating the information sent by the workers: nodes are only allowed to communicate with their neighbors in the network. This work addresses simultaneously multiple limitations of existing decentralized algorithms while guaranteeing fast convergence rates.

Synchronous lock. Optimal methods [42, 27] have been derived for synchronous first-order algorithms, whose executions are blocked until all nodes have reached a predefined state (*e.g.*, they must all finish computing local gradients before the round of communication begins), which limits their efficiency in practice as they can heavily be impacted by a few slow nodes or edges in the graph (the *straggler problem*). To tackle the synchronous lock, we rely on the continuized framework [12], itself derived from the randomized gossip model [5]. In randomized gossip, nodes update their local values at random times using pairwise communication updates named gossip. Thus, iterates are randomized, labeled with a continuous-time index that only need local clocks to be synchronized at the beginning of the procedure (in opposition to a *global* iteration count that has to be known by all at all time) and performed locally with no regards to a specific global ordering of events. While based on discrete events and thus readily implementable, the continuized framework simplifies the analysis by leveraging continuous proof tools.

Coupled lock. However, in [12], gradient and gossip operations are coupled: each communication along an edge first requires the computation of the gradients of the two functions locally stored on the corresponding nodes. As more communication steps than gradient computations are necessary to reach an ϵ precision, even in an optimal framework [27, 42], the coupling leads to an overload in terms of gradient steps. Moreover, coupling computations and communications implies they must be performed *sequentially*, decoupling them allows both tasks to be done in *parallel*, allowing an additional speedup.

To our knowledge, our work is the first primal method to tackle those locks simultaneously while obtaining accelerated rates for both computations and communications. We propose a novel algorithm (DADAO: Decoupled Accelerated Decentralized Asynchronous Optimization) based on a combination of similar formulations to [27, 13, 18] in the continuized framework of [12]. We study:

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

where each $f_i : \mathbb{R}^d \rightarrow \mathbb{R}$ is a μ -strongly convex and L -smooth function computed in one of the n nodes of a network. We derive a first-order optimization algorithm that only uses primal gradients and relies on Point-wise Poisson Processes (P.P.P.s [31]) modeling of the communication and gradient occurrences, leading to accelerated communication and computation rates. Furthermore, our communication bounds rely on the maximal resistance of a graph rather than the largest eigenvalue of a Laplacian, leading to an additional acceleration compared to works which rely on synchrony. Our framework is based on a simple fixed-point iteration and kept minimal: it only involves primal computations with an additional momentum term. Thus, we do not add other cumbersome designs such as the Error-Feedback or Forward-Backward used in [27] (whose adaptation to asynchronous settings is for now unclear). While we do not consider the delays bound to appear in practice (we assume instantaneous communications and computations), we remove the coupling lock by performing gradient and gossip steps in parallel. Tab. 1 compares DADAO with other approaches and shows it is the only work to achieve accelerated rates both in number of communication and gradients. Then, we assume that each edge fires at a rate that can be adjusted *a priori*, yet the network has to verify a set of physical constraints, such as a maximal bandwidth condition. This contrasts with common gossip algorithms, which either leave the weights of the gossip matrix to be defined by the user or propose simple appropriate heuristics, such as the use of Metropolis weights, see *e.g.* [48, 43, 37, 17]. We cast a relaxation of the minimization of the communication rate as a novel convex optimization problem, which can be efficiently solved with classical solvers. This allows for outperforming communication rates obtained from uniform or standardized Laplacian weights, which are often suggested in the literature.

Contributions. (1) We propose a primal algorithm with provable guarantees in the context of asynchronous decentralized learning. (2) This algorithm is the first to reach accelerated rates for both communications and computations while not requiring ad-hoc mechanisms obtained from an inner loop. (3) We propose a simple theoretical framework compared to concurrent works, we show that our rates are better than previous works, and (4) we illustrate this theoretical comparison numerically. (5) We propose a SDP procedure to optimize the communication rate of our algorithm for a user-specified graph.

Structure of the paper. In Sec. 3.1, we describe our work hypothesis and our model of a decentralized environment, while Sec. 3.2 describes our dynamic. Sec. 3.3 states our convergence guarantees and highlights that the communication and computational rates of our method are better compared to its competitors. Next, Sec. 4 explains how to optimize the edge weights of a given graph topology to minimize DADAO’s overall communication rate. Sec. 5.1 explains our implementation of this algorithm, and finally, Sec. 5.2 verifies our claims numerically. All our experiments are reproducible, using PyTorch [39], our code being online <https://github.com/AdelNabli/DADAO/>.¹

Notations: $f = \mathcal{O}(g)$ means there is a constant $C > 0$ such that $|f| \leq C|g|$, $\{e_i\}_{i \leq d}$ is the canonical basis of \mathbb{R}^d , $d \in \mathbb{N}$, $\mathbf{1}$ is the vector of 1, \mathbf{I} the identity, A^+ is the pseudo-inverse of A . We further write $e_i \triangleq e_i \otimes \mathbf{I}$.

2 Related Work

Table 1: This table shows the strength of DADAO compared to concurrent works for obtaining ϵ -precision. n is the number of nodes, $|\mathcal{E}|$ the number of edges, $\frac{1}{\chi_1}$ the smallest positive eigenvalue of a fixed weighted Laplacian \mathcal{L} , ρ the eigengap and $\chi_2 \leq \chi_1$ the effective resistance. Note that under reasonable assumptions $\sqrt{\chi_1 \chi_2} n = \mathcal{O}(|\mathcal{E}| \sqrt{\rho})$ (see Lemma 3.1). Async., Comm., Grad., M.-C. and Prox. stand respectively for Asynchrony, Communication steps, Gradient steps., Multi-consensus and Proximal operator. As suggested in their respective papers, all the algorithms are run with $\frac{1}{\|\mathcal{L}\|} \mathcal{L}$. For AGT and OGT, the mixing matrix used is stochastic, a more precise comparison is given by Prop. 3.3.

Method	Async.	Decoupled	No Inner Loop (M.-C. or Prox.)	Primal Oracle	Total # Comm.	Total # Grad.
MSDA [42]	\times	\times	\times	\times	$\sqrt{\rho} \mathcal{E} \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}$	$n \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}$
DVR [21]	\times	\times	\times	\checkmark	$\sqrt{\rho} \mathcal{E} \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}$	$n \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}$
ADOM+ [27]	\times	\times	\times	\checkmark	$\rho \mathcal{E} \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}$	$n \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}$
TVR [18]	\times	\checkmark	\times	\checkmark	$\rho \mathcal{E} \frac{L}{\mu} \log \frac{1}{\epsilon}$	$n \frac{L}{\mu} \log \frac{1}{\epsilon}$
AGT [34]	\times	\times	\times	\checkmark	$\sqrt{\rho} \mathcal{E} \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}$	$n \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}$
OGT [44]	\times	\checkmark	\checkmark	\checkmark	$\sqrt{\rho} \mathcal{E} \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}$	$n \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}$
ESDACD [19]	\checkmark	\times	\checkmark	\times	$\sqrt{\chi_1 \chi_2} n \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}$	$\sqrt{\chi_1 \chi_2} n \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}$
Continuized [12]	\checkmark	\times	\checkmark	\times	$\sqrt{\chi_1 \chi_2} n \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}$	$\sqrt{\chi_1 \chi_2} n \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}$
DADAO (ours)	\checkmark	\checkmark	\checkmark	\checkmark	$\sqrt{\chi_1 \chi_2} n \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}$	$n \sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}$

Continuized and asynchronous algorithms. We highly rely on the elegant continuized framework [12], which allows obtaining simpler proofs and brings the flexibility of asynchronous algorithms. We reemphasize that identically to

¹**Erratum:** An earlier version of this preprint wrongly claimed that the method was suitable for time-varying settings whereas further analysis is required for this to be properly shown.

[12], the result of this paper is *a stochastic discrete algorithm with a continuous proof*: our proof framework is not based on the discretization of an Ordinary Differential Equation (ODE) but rather studies *the evolution of a Stochastic Differential Equation (SDE) with jumps*. However, by contrast to [12], in our work, we significantly reduce the necessary amount of gradient steps compared to [12] while maintaining the same amount of activated edges. Another type of asynchronous algorithm can also be found in [32], yet it fails to obtain Nesterov’s accelerated rates for lack of momentum. We note that [33] studies the robustness to delays yet requires a shared memory and thus applies to a different context than decentralized optimization. [18] is a promising approach for modeling random communication on graphs yet fails to obtain acceleration in a neat framework without inner loops.

Decentralized algorithms with fixed topology. [42] is the first work to derive an accelerated algorithm for decentralized optimization, and it links the convergence speed to the Laplacian eigengap. The corresponding algorithm uses a dual formulation and a Chebychev acceleration (synchronous and only for fixed topology). Yet, as stated in Tab. 2, it still requires many edges to be activated. Furthermore, under a relatively flexible condition on the intensity of our P.P.P.s, we show that our work improves over bounds that depend on the spectral gap. An emerging line of work following this formulation employs the continuized framework [14, 12, 13], but unfortunately do not use a primal oracle, as they rely on the gradients of the Fenchel conjugate. Finally, we note that the work of [13] incorporates delays in their model, using the same technique as our work, yet transferring this robustness to another setting remains unclear. Reducing the number of communication has been studied in [36], but without obtaining accelerated rates. [22] allows for fast communication and gossip rates yet requires a proximal step and synchrony between nodes to apply a momentum variable.

Finite sum acceleration. If each local function f_i is a sum of elementary functions $\sum_{j=1}^m f_{i,j}$ with a favorable conditioning, an additional acceleration is possible, as observed by [22, 21, 18], which is a different problem from Eq. 1. This can be viewed as a cluster of nodes with infinite connectivity and an efficient decentralized algorithm should automatically adapt to such structure. Thus, we focused on the setting $m = 1$, which allows a fair comparison with these works.

Error feedback/Gradient tracking. A major lock for asynchrony is the use of Gradient Tracking [26, 37, 34] or Error Feedback [45, 28]. Indeed, gradient operations are locally tracked by a running-mean variable which must be synchronously updated at each gradient update, making it incompatible with an asynchronous framework. Furthermore, acceleration requires an undesirable multi-consensus inner loop. We emphasize that [44] allows to decouple the gradient updates from communication, yet the framework is still synchronous,

leading to synchronous communication rates, that asynchrony can improve (see Tab. 1).

Decoupling procedures. Decoupling subsequent steps of optimization procedures traditionally leads to speed-ups [22, 18, 3, 4]. This contrasts with methods which couple gradient and gossip updates, so that they happen in a predefined order, i.e., simultaneously [12] or sequentially [27, 25]. In decoupled optimization procedures, inner-loops are not desirable as they require an external procedure that can be potentially slow and need a block-barrier instruction during the algorithm’s execution (e.g., [22]). It means in particular that it is preferable to avoid approaches such as Catalyst [35], multi-consensus steps [30] or Tchebychev acceleration of consensus [42].

Resistance of a graph. The maximal resistance of a graph is a widely studied quantity, particularly in physics [24, 46, 23], as it is a refined geometric invariant of graphs. The resistance of a graph corresponds to the commute time of a Markov Chain [7]. However, beyond the Continuized framework [12] or acceleration of consensus problems [6, 2, 16], we are unaware of generic, asynchronous, accelerated decentralized optimization procedures that rely on this quantity. Also, it has a more physical interpretation than the Laplacian’s norm, as it can be computed via Ohm’s and Kirchoff’s Circuit Laws (see [7]).

3 Accelerated Asynchronous Algorithm

3.1 Gossip Framework

We consider the problem defined by Eq. 1 in a distributed environment constituted by n nodes whose dynamic is indexed by a continuous time index $t \in \mathbb{R}^+$. Each node has a local memory and can compute a local gradient ∇f_i , as well as elementary operations, in an instantaneous manner. As said above, having no delay is less realistic, yet adding them also leads to significantly more difficult proofs whose adaptation to our framework remains largely unclear. Next, we will assume that our computations and gossip result from independent homogeneous P.P.P. with no delay. For the sake of simplicity, we assume that all nodes can compute a gradient at the same rate:

Assumption 3.1 (Homogeneous gradient computations). *The gradient computations are normalized to fire independently at a rate of 1 computation per time unit. For the i -th worker, we write $N_i(t)$ the corresponding P.P.P. of rate 1, as well as $\mathbf{N}(t) = (N_i(t))_{i \leq n}$.*

Remark 3.1. *The P.P.P $N_i(t)$ on node i means that taking gradient steps at i are discrete events, but the time intervals between two events is a random variable following an exponential law of parameter 1. Thus, the expected waiting time between two gradient steps on the i -th worker is 1 time unit.*

Next, we model the bandwidth of each connection. For an edge (i, j) belonging to \mathcal{E} , the set of edges of a graph we assume connected (3.2), we write $M_{ij}(t)$ the P.P.P. with rate $0 < \lambda_{ij} < \infty$. When this P.P.P. fires, both nodes share and update their local memories. The rate λ_{ij} is adjustable locally by machine i while λ_{ji} is controlled by machine j . While λ_{ij} and λ_{ji} may be different, we highlight that the communication process is symmetric, *i.e.* both nodes update their local memories when either one of the two corresponding P.P.P.s fires. Thus, in the corresponding undirected graph \mathcal{E} , each edge (i, j) will fire at a rate of $\lambda_{ij} + \lambda_{ji}$. Given our notations, if $(i, j) \notin \mathcal{E}$, then the connection between (i, j) can be thought as a P.P.P. with intensity 0. Taking the λ_{ij} as edge weights, we introduce the subsequent graph Laplacian, which is the expected Laplacian of our graph:

$$\Lambda \triangleq \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} (e_i - e_j)(e_i - e_j)^\top.$$

We write $\mathbf{\Lambda} \triangleq \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} (\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top$ its tensorized counter-part that will be useful for our Lyapunov-based proofs. Following [42], we will compare this quantity to the following projector:

$$\pi \triangleq \mathbf{I} - \frac{1}{n} \mathbf{1}\mathbf{1}^\top = \frac{1}{2n} \sum_{1 \leq i, j \leq n} (\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top.$$

Using the spectral norm, we introduce the algebraic connectivity of our network using Λ , as in [27]: $\chi_1 \triangleq \|\Lambda^+\|$. We might also write $\chi_1[\Lambda]$ to avoid confusion, depending on the context.

Assumption 3.2 (Connected graph). *The set of edges \mathcal{E} defines a connected graph such that $\chi_1[\Lambda] < \infty$.*

Next, the maximal effective resistance of the network, as in [12, 11], is:

$$\chi_2 \triangleq \frac{1}{2} \sup_{(i,j) \in \mathcal{E}} (e_i - e_j)^\top \Lambda^+ (e_i - e_j).$$

A standard quantity [42], which is used to control the number of synchronous gossips steps, is the spectral gap, given by $\rho \triangleq \|\Lambda\| \|\Lambda^+\|$. We also introduce the value κ , the ratio of communication frequency between the fastest and slowest edges:

$$\kappa \triangleq \frac{\sup_{(i,j) \in \mathcal{E}} \lambda_{ij} + \lambda_{ji}}{\inf_{(i,j) \in \mathcal{E}} \lambda_{ij} + \lambda_{ji}}.$$

This ratio is typically bounded, for instance, in the case of a graph with constant edge weights or for $\lambda_{ij} = \frac{1}{d_i}$ with d_i the degree of the i -th node in a bounded degree graph or a regular graph. We prove the following Lemma (proved in Appendix A), which is useful to control χ_1, χ_2 and compare our bounds with works that rely on the spectral gap of a graph:

Lemma 3.1 (Effective resistance). *The spectrum of Λ is non-negative. Also, we have $\chi_1 = +\infty$ iff $\bar{\mathcal{E}}$ is not a connected graph. Also, if the graph is connected, then:*

$$\frac{n-1}{\text{Tr } \Lambda} \leq \chi_2 \leq \chi_1.$$

Furthermore, we also have the following:

$$\sqrt{\chi_1 \chi_2} \text{Tr } \Lambda \leq \sqrt{\rho} \sqrt{\kappa n |\bar{\mathcal{E}}|}.$$

The last part of this Lemma indicates that our method requires less communications than methods depending on the spectral gap when no degenerated behavior on the graph's connectivity happens, *i.e.* when κ is adequately bounded, which is a standard assumption [20].

Remark 3.2. *For synchronous frameworks [27, 42], the spectral quantity extracted from their gossip matrix is a given measure of the connectedness of their graphs that is used afterwards to deduce the right number of synchronous communication rounds between two gradient steps. In our framework, Λ directly contains the information of both the topology \mathcal{E} and the edge communication rates λ_{ij} , thus $\chi_1[\Lambda]$ must rather be understood as an indicator of how well the graphs connectedness and the chosen communication strategy interact. In fact, we will see later that there is a condition on χ_1, χ_2 for DADAO to converge, see Appendix F for further discussion.*

3.2 Dynamic to optimum

Next, we follow a standard approach [29, 27, 41, 18] for solving Eq. 1 (see Appendix B for details), leading to studying, for $0 < \nu < \mu$, the following Lagrangian:

$$\inf_{x \in \mathbb{R}^{n \times d}} \sup_{\substack{y \in \mathbb{R}^{n \times d} \\ z \in \mathbb{R}^{n \times d}}} \sum_{i=1}^n f_i(x_i) - \frac{\nu}{2} \|x\|^2 - \langle x, y \rangle - \frac{1}{2\nu} \|\pi z + y\|^2.$$

For $f(x) = \sum_{i=1}^n f_i(x_i)$, the saddle points (x^*, y^*, z^*) of the above Lagrangian are given by:

$$\begin{cases} \nabla f(x^*) - \nu x^* - y^* & = 0 \\ y^* + \pi z^* + \nu x^* & = 0 \\ \pi z^* + \pi y^* & = 0. \end{cases} \quad (2)$$

Our algorithm is based on a fixed-point algorithm to obtain those saddle points, which is a similar idea to [28], which is only restricted to a setting without communication acceleration. Furthermore, contrary to [27], we do not employ a Forward-Backward algorithm, which requires both an extra-inversion step and additional regularity on the considered proximal operator. Not only does this condition not hold in this particular case, but this is not desirable in a continuous framework where iterates are not ordered in a predefined sequence and

require a local descent at each instant. Another major difference is that no Error-feedback is required by our approach, which allows unlocking asynchrony while simplifying the proofs and decreasing the required number of communications. Instead, we show it is enough to incorporate a standard fixed point algorithm, *without any specific preconditioning* (see [8]). We consider the following dynamic:

$$\begin{aligned}
dx_t &= \eta(\tilde{x}_t - x_t)dt - \gamma(\nabla f(x_t) - \nu x_t - \tilde{y}_t) d\mathbf{N}(t) \\
d\tilde{x}_t &= \tilde{\eta}(x_t - \tilde{x}_t)dt - \tilde{\gamma}(\nabla f(x_t) - \nu x_t - \tilde{y}_t) d\mathbf{N}(t) \\
d\tilde{y}_t &= -\theta(y_t + z_t + \nu \tilde{x}_t)dt + (\delta + \tilde{\delta})(\nabla f(x_t) - \nu x_t - \tilde{y}_t)d\mathbf{N}(t) \\
dy_t &= \alpha(\tilde{y}_t - y_t)dt \\
dz_t &= \alpha(\tilde{z}_t - z_t)dt - \beta \sum_{(i,j) \in \mathcal{E}} (\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top (y_t + z_t) dM_{ij}(t) \\
d\tilde{z}_t &= \tilde{\alpha}(z_t - \tilde{z}_t)dt - \tilde{\beta} \sum_{(i,j) \in \mathcal{E}} (\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top (y_t + z_t) dM_{ij}(t),
\end{aligned} \tag{3}$$

where $\nu, \tilde{\eta}, \eta, \gamma, \alpha, \tilde{\alpha}, \theta, \delta, \tilde{\delta}, \beta, \tilde{\beta}$ are undetermined real-valued parameters. As in [38], variables are paired to obtain a Nesterov acceleration. The variables (x, y) allow decoupling the gossip steps from the gradient steps using independent P.P.s. Furthermore, the Lebesgue integrable path of \tilde{y}_t does not correspond to a standard momentum, as in a continuized framework [12]; however, it turns out to be a crucial component of our method. Compared to [27], no extra multi-consensus step needs to be integrated. Our formulation of an asynchronous gossip step is similar to [12], which introduces a stochastic variable on edges; however, contrary to this work, our gossip and gradient computations are decoupled. We emphasize that while the dynamic 3 is formulated using SDEs [1], which brings the power of the continuous-time analysis toolbox, it is still *event-based* and thus discrete in nature. Hence, the dynamic can be efficiently implemented in practice as explained in Sec. 5.1 and Appendix H.

3.3 Theoretical guarantees

We prove the following in Appendix C.

Theorem 3.2. *Assume each f_i is μ -strongly convex and L -smooth. Assume 3.1, and that $2\chi_1\chi_2 \leq 1$. Then there exists some parameters for the dynamic Eq. (3) (given in Appendix C.1), such that for any initialization $x_0 \in \ker(\pi)$, and $\tilde{x}_0 = x_0, y_0 = \tilde{y}_0 = \nabla f(x_0) - \frac{\mu}{2}x_0, z_0 = \tilde{z}_0 = -\pi y_0$, we get for $t \in \mathbb{R}^+$:*

$$\mathbb{E}[\|x_t - x^*\|^2] \leq \left(\frac{1}{2} + \frac{23L}{8\mu} + 2\frac{L^2}{\mu^2}\right)\|x_0 - x^*\|^2 e^{-\frac{t}{8\sqrt{2}}\sqrt{\frac{\mu}{L}}}$$

Also, the expected number of oracle gradient calls is nt and the expected number of edges activated is: $\frac{t}{2}\text{Tr } \Lambda$.

The condition $2\chi_1\chi_2 \leq 1$ can simply be understood as whether or not the chosen communication strategy suits sufficiently well the graph topology \mathcal{E} for

DADAO to converge using the expected rates of communication $(\lambda_{ij})_{(i,j) \in \mathcal{E}}$ compared to the expected rate of one gradient step per time unit on each node we assumed in Assumption 3.1 (see Appendix F for a more complete discussion). Thus, given a graph topology, it is straightforward to obtain a Laplacian verifying $2\chi_1\chi_2 \leq 1$. Indeed, we have the following:

Remark 3.3. *For any graph topology \mathcal{E} and any choice of corresponding graph Laplacian \mathcal{L} verifying $\chi_1[\mathcal{L}]$ bounded, there is a communication strategy $(\lambda_{ij})_{(i,j) \in \mathcal{E}}$ such that Λ , the Laplacian of \mathcal{E} with edges $(i, j) \in \mathcal{E}$ weighted with their expected communication rate λ_{ij} , verifies $2\chi_1[\Lambda]\chi_2[\Lambda] \leq 1$. One such example is given by:*

$$\Lambda = \sqrt{2\chi_1[\mathcal{L}]\chi_2[\mathcal{L}]} \mathcal{L}.$$

This property allows us to use in DADAO the Laplacians introduced in previous work, to which we can now compare. It leads to the following proposition (proved in Appendix D), which shows that DADAO obtains better complexities than concurrent works while starting from the same family of Laplacians:

Proposition 3.3 (Comparison with concurrent work).

- If $\kappa = \mathcal{O}(\frac{|\mathcal{E}|}{n})$, then DADAO obtains better communication rate than MSDA [42], and requires strictly less communications for the complete graph.
- For any fixed Laplacian valid for ADOM+ [27], DADAO obtains a better communication rate than ADOM+, and requires strictly less communications for the complete graph.
- If $\kappa = \mathcal{O}(\frac{|\mathcal{E}|}{n})$, for a valid Laplacian using the Gossip matrix of OGT [44] or AGT [34], DADAO obtains a better communication rate than both Gradient Tracking methods, and requires strictly less communications for the complete graph.
- DADAO requires fewer gradient computations than the Continuzied framework [13], and has a strictly better computation rate for the cycle graph.

We highlight that [42] claimed that their algorithm is optimal because they study the number of computations and *synchronized* gossips on a worst-case graph; our claim is, *by nature* different, as we are interested in the number of edges fired rather than the number of synchronized gossip rounds. Indeed, in an asynchronous framework, there is no notion of *round of* communication, which allows this framework to enjoy the advantageous rates of randomized procedures. Moreover, not only this measure of complexity is standard in asynchronous frameworks (*e.g.*, see [5]), but also, if we are aiming for the most frugal procedure, minimizing both the total number of computations and communications is of interest. Tab. 2 predicts the behavior of our algorithm on various classes of graphs encoded via a normalized Laplacian. It shows that systematically, our algorithm leads to the best complexities. For example, in the case of

a complete graph, one synchronized gossip round requires a total of $|\mathcal{E}| = \mathcal{O}(n^2)$ communications whereas we show that a rate of only $\mathcal{O}(n)$ communications per time unit suffices for DADAO in this case. We note that the graph class depicted in Tab. 2 was used as worst-case examples for proving the optimality of [42] in a synchronous context.

Table 2: Complexity for various graphs using $\frac{1}{\|\mathcal{E}\|}\mathcal{L}$ with \mathcal{L} the standard Laplacian with unit edge weights. In this case, $\rho = \chi_1$. The complexities are reported per time unit so that all algorithms reach ϵ -precision at the same time. We have, respectively, for a star/line or cyclic/complete graph and the d -dimensional grid: $\chi_1 = \mathcal{O}(n)$, $\chi_2 = \mathcal{O}(n)$, $\text{Tr} = \mathcal{O}(1)$ / $\chi_1 = \mathcal{O}(n^2)$, $\chi_2 = \mathcal{O}(1)$, $\text{Tr} = \mathcal{O}(n)$ / $\chi_1 = \mathcal{O}(1)$, $\chi_2 = \mathcal{O}(1)$, $\text{Tr} = \mathcal{O}(n)$ / $\chi_1 = \mathcal{O}(n^{2/d})$, $\chi_2 = \mathcal{O}(1)$, $\text{Tr} = \mathcal{O}(n)$.

Method	# edges activated per time unit				# gradients computed per time unit			
	Star	Line	Complete	d -grid	Star	Line	Complete	d -grid
[27] ADOM+	n^2	n^3	n^2	$n^{1+2/d}$	n	n	n	n
[42] MSDA	$n^{3/2}$	n^2	n^2	$n^{1+1/d}$	n	n	n	n
[12] Continuzied	n	n^2	n	$n^{1+1/d}$	n	n^2	n	$n^{1+1/d}$
Centralized	n	-	-	-	n	-	-	-
DADAO (ours)	n	n^2	n	$n^{1+1/d}$	n	n	n	n

4 Minimizing the communication rate as a SDP

For a given connected graph topology \mathcal{E} , we introduce $\mathbb{R}_+^{\mathcal{E}} \triangleq \{\lambda, \lambda_{ij} \geq 0 \text{ and } \lambda_{ij} = 0, \text{ if } (i, j) \notin \mathcal{E}\}$. We want to find a set of edge firing rates that would minimize the overall number of communications. Given a set of weights $\lambda \in \mathbb{R}_+^{\mathcal{E}}$ defining a Laplacian $\Lambda(\lambda)$, it is sufficient to use the edge firing rates given by the normalized Laplacian $\Lambda(\lambda)\sqrt{\chi_1[\Lambda(\lambda)]\chi_2[\Lambda(\lambda)]}$ to guarantee convergence. Thus, the quantity we want to minimize is the trace of this normalized Laplacian, *i.e.* $\text{Tr} \Lambda(\lambda)\sqrt{2\chi_1[\Lambda(\lambda)]\chi_2[\Lambda(\lambda)]}$. However, except in a simplistic scenario, it is difficult to estimate optimal edge weights for a given graph. This has been discussed in [16] for the sum of the effective resistances of a graph, and in [47, 5] for $\chi_1[\Lambda(\lambda)]$. Similarly to [16], we propose an SDP to minimize the corresponding communication rate, which relies on a slightly technical convexity lemma (see Appendix E.1). However, $\lambda \rightarrow \chi_2[\Lambda(\lambda)]$ is not convex. Indeed, the supremum is taken over a set of edge *dynamically* defined through λ , $\mathcal{E}(\lambda) \triangleq \{(i, j), \lambda_{ij} > 0\}$ which makes it non-convex: an optimal solution might lead to removing some edges from the graph and consider the maximum effective resistance of this subgraph. This forces us to slightly relax this condition by taking the supremum over a predefined *fixed* set of edges, and we introduce for any $\lambda \in \mathbb{R}_+^{\mathcal{E}}$:

$$\chi_2^{\mathcal{E}}(\lambda) \triangleq \frac{1}{2} \max_{(i,j) \in \mathcal{E}} (e_i - e_j)^\top \Lambda^+(\lambda) (e_i - e_j).$$

This function is convex, and we note that by definition $\chi_2^{\mathcal{E}}(\lambda) \geq \chi_2[\Lambda(\lambda)]$. We now propose several structural constraints on the weights of a Laplacian:

- (\mathcal{C}_1) **Global bandwidth:** $\sum_{(ij) \in \mathcal{E}} \lambda_{ij} \leq \lambda$, where $\lambda > 0$ is a total bandwidth constraint,
- (\mathcal{C}_2) **Local connectivity:** $\sum_{j, (ij) \in \mathcal{E}} \lambda_{ij} \leq \lambda_i$, where $\lambda_i > 0$ is a local constraint on node connectivity,
- (\mathcal{C}_3) **Local bandwidth:** $\sum_{j, (ij) \in \mathcal{E}} \lambda_{ji} + \lambda_{ij} \leq \lambda_i$, where $\lambda_i > 0$ is a local constraint on edge connectivity to node i .

In the following, we will generically call $\mathcal{C} \in \{\mathcal{C}_1, \mathcal{C}_2, \mathcal{C}_3\}$ one of these linear constraints. Getting back to our initial problem, we want to consider:

$$\begin{aligned} & \text{minimize} && \text{Tr } \Lambda(\lambda) \sqrt{2\chi_1[\Lambda(\lambda)]\chi_2[\Lambda(\lambda)]} \\ & \text{subject to} && \lambda \in \mathcal{C} \cap \mathbb{R}_+^{\mathcal{E}} \end{aligned} \quad (4)$$

which allows connected subgraphs of \mathcal{E} to be considered. However, in order to frame this minimization problem as a SDP, we first need to relax it, see Appendix E.2 for a complete discussion. This leads to the following proposition, proved in Appendix E.3:

Proposition 4.1 (Minimizing communications as a SDP). *Given a fixed topology \mathcal{E} , the problem:*

$$\begin{aligned} & \text{minimize} && \sqrt{2\chi_1[\Lambda(\lambda)]\chi_2^{\mathcal{E}}(\lambda)} \\ & \text{subject to} && \lambda \in \mathcal{C} \cap \mathbb{R}_+^{\mathcal{E}} \end{aligned} \quad (5)$$

is equivalent to the following SDP:

$$\begin{aligned} & \text{minimize} && t_1 + t_2 \\ & \text{subject to} && \lambda \in u\mathcal{C}, u \geq 1, \lambda_{ij} \geq 0 \forall (i, j) \in \mathcal{E} \\ & && \Lambda = \sum_{ij} \lambda_{ij} (e_i - e_j)(e_i - e_j)^{\top} \\ & && \begin{pmatrix} \Lambda & u(e_i - e_j) \\ u(e_i - e_j)^{\top} & t_1 \end{pmatrix} \succcurlyeq 0 \quad \forall (i, j) \in \mathcal{E} \\ & && \begin{pmatrix} \Lambda & \mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^{\top} \\ \mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^{\top} & t_2\mathbf{I} \end{pmatrix} \succcurlyeq 0 \end{aligned}$$

In Appendix E.3, we prove that, running the accelerated gossip procedures on the Barbell graph using uniform weights is $\mathcal{O}(n^{3/2})$ worse than with weights solution to the relaxed problem (5). This is illustrated in Fig. 1 (a) and confirmed by Fig. 1 (b), which displays the impact of choosing sub-optimal gossip matrices for the convergence rates of gossip-based methods in practice. For more numerical results using our SDP solvers, see Appendix E.5.

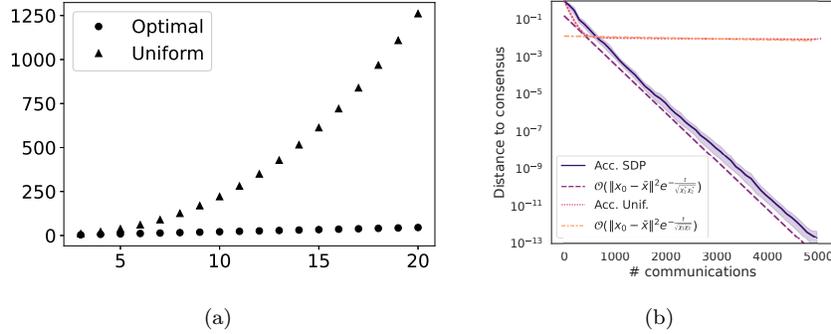


Figure 1: (a) Optimal value of $\sqrt{2\chi_1\chi_2}$ for problem (5) under constraint \mathcal{C}_1 compared to the one obtained with uniform weights $1/|\mathcal{E}|$, as a function of n for the barbell graph. (b) Convergence rate for the accelerated gossip procedure on a barbell graph $K_{50} - K_{50}$ using both the edge weights given by our SDP and the uniform ones. For reference, we plot the theoretical rates, denoting by χ_1^*, χ_2^* the values obtained using the SDP weights, and by χ_1, χ_2 the ones obtained with $1/|\mathcal{E}|$. Averaged over 20 runs.

5 Practical implementation

5.1 Algorithm

To study the trajectories of $X_t \triangleq (x_t, \tilde{x}_t, \tilde{y}_t), Y_t \triangleq (y_t, z_t, \tilde{z}_t)$, we use the following, equivalent to (3):

$$\begin{aligned} dX_t &= a_1(X_t, Y_t)dt + b_1(X_t)d\mathbf{N}(t) \\ dY_t &= a_2(X_t, Y_t)dt + \sum_{(i,j) \in \mathcal{E}} b_2^{ij}(Y_t)dM_{ij}(t), \end{aligned} \quad (6)$$

where $a_1, a_2, b_1 = (b_1^i)_i, (b_2^{ij})_{ij}$ are smooth functions. We now describe the algorithm used to implement the dynamics of Eq. (3). Let us write $T_1^{(i)} < T_2^{(i)} < \dots < T_k^{(i)} < \dots$ the time of the k -th event on the i -th node, which is either an edge activation, either a gradient update. We remind that the spiking times of a specific event correspond to random variables with independent exponential increments and can thus be generated at the beginning of our simulation. They can also be generated on the fly and locally to stress the locality and asynchronicity of our algorithm. We write $X_t = (X_t^{(i)})_i$ and $Y_t = (Y_t^{(i)})_i$, then on the i -th node, at the k -th iteration, we integrate on $[T_k^{(i)}; T_{k+1}^{(i)}]$ the ODE

$$\begin{aligned} dX_t &= a_1(X_t, Y_t)dt \\ dY_t &= a_2(X_t, Y_t)dt \end{aligned}$$

to define the values right before the spike. One can easily find the matrix $\mathcal{A} \in \mathbb{R}^{6 \times 6}$ such that:

$$\begin{pmatrix} X_{T_{k+1}^{(i)-}}^{(i)} \\ Y_{T_{k+1}^{(i)-}}^{(i)} \end{pmatrix} = \exp\left((T_{k+1}^{(i)} - T_k^{(i)})\mathcal{A}\right) \begin{pmatrix} X_{T_k^{(i)}}^{(i)} \\ Y_{T_k^{(i)}}^{(i)} \end{pmatrix}. \quad (7)$$

Next, if one has a gradient update, then:

$$X_{T_{k+1}^{(i)}}^{(i)} = X_{T_{k+1}^{(i)-}}^{(i)} + b_1 \left(X_{T_{k+1}^{(i)-}}^{(i)} \right).$$

Otherwise, if the edge (i, j) or (j, i) is activated, a communication bridge is created between both nodes i and j . In this case, the local update on i writes:

$$Y_{T_{k+1}^{(i)}}^{(i)} = Y_{T_{k+1}^{(i)-}}^{(i)} + b_2 \left(Y_{T_{k+1}^{(i)-}}^{(i)}, Y_{T_{k+1}^{(i)-}}^{(j)} \right).$$

Note that, even if this event takes place along an edge (i, j) , we can write it separately for nodes i and j by making sure they both have the events $T_{k_i}^{(i)} = T_{k_j}^{(j)}$, for some $k_i, k_j \in \mathbb{N}$ corresponding to this communication. As advocated, all those operations are local, and we summarize in the Alg. 1 the algorithmic block corresponding to our implementation. See Appendix H for more details.

Algorithm 1: This algorithm block describes our implementation on each local machine. The *ODE* routine is described by Eq. 7.

Input: On each machine $i \in \{1, \dots, n\}$, gradient oracle ∇f_i , parameters

μ, L, χ_1, t_{\max} .

- 1 **Initialize** on each machine $i \in \{1, \dots, n\}$:
 - 2 | Set $X^{(i)}, Y^{(i)}, T^{(i)}$ to 0 and set \mathcal{A} via Eq. (88);
 - 3 **Synchronize** the clocks of all machines ;
 - 4 **In parallel on workers** $i \in \{1, \dots, n\}$, **while** $t < t_{\max}$, **continuously do:**
 - 5 | $t \leftarrow \text{clock}()$;
 - 6 | **if** there is an event at time t **then**
 - 7 | | $(X^{(i)}, Y^{(i)}) \leftarrow \text{ODE}(\mathcal{A}, t - T^{(i)}, X^{(i)}, Y^{(i)})$;
 - 8 | | **if** the event is to take a gradient step **then**
 - 9 | | | $X^{(i)} \leftarrow X^{(i)} + b_1(X^{(i)})$;
 - 10 | | **else if** the event is to communicate with j **then**
 - 11 | | | $Y^{(i)} \leftarrow Y^{(i)} + b_2(Y^{(i)}, Y^{(j)})$; // Happens at j simultaneously.
 - 12 | | $T^{(i)} \leftarrow t$;
 - 13 **return** $(x_{t_{\max}}^{(i)})_{1 \leq i \leq n}$.
-

5.2 Numerical results

In this section, we study the behavior of our method in a standard experimental setting (*e.g.*, see [27, 12]). In order to compare to methods using the gradient of the Fenchel conjugate [12] in our experiments, we restrict ourselves to a situation where it is easily computable. Thus, we perform the empirical risk minimization for the decentralized linear regression task given by:

$$f_i(x) = \frac{1}{m} \sum_{j=1}^m \|a_{ij}^\top x - c_{ij}\|^2, \quad (8)$$

where $a_{ij} \in \mathbb{R}^d$, and $c_{ij} \in \mathbb{R}$ correspond to m local data points stored at node i . We follow a protocol similar to [27]: we generate n independent synthetic datasets with the `make_regression` functions of scikit-learn [40], each worker storing $m = 100$ data points. We recall that the metrics of interest are the total number of local gradient steps and the total number of individual messages exchanged (*i.e.*, *number of edges that fired*) to reach an ϵ -precision. We systematically used the proposed hyper-parameters of each reference paper for our implementation without any specific fine-tuning.

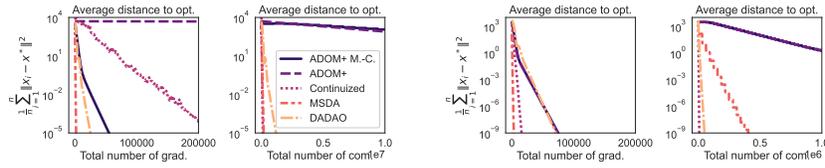


Figure 2: Comparison between ADOM+ [27], the continuized framework [12], MSDA [42] and DADAO, using the same data for the linear regression task, and the same graphs (*from left to right: line with $n = 150$, complete with $n = 250$*).

Comparison between all methods. We fix the Laplacian matrix via Eq. (3.3) to compare simultaneously to the continuized framework [12], MSDA [42] and ADOM+ [27]. We compare ourselves to both versions of ADOM+: with and without the Multi-Consensus (M.-C.). We report in Fig. 2 results corresponding to the complete graph with $n = 250$ nodes and the line graph of size $n = 150$. While sharing the same asymptotic rate (see Fig. 7 for experimental confirmation), we note that the Continuized framework [12] and MSDA [42] have better absolute constants than DADAO, giving them an advantage both in terms of the number of communication and gradient steps. However, in the continuized framework, the gradient and communication steps being coupled, the number of gradient computations can potentially be orders of magnitude worse than our algorithm, which is reflected by Fig. 2 for the line graph. As for MSDA, Tab. 2 showed they do not have the best communication rates on certain

classes of graphs, as confirmed to the right in Fig. 2 for MSDA and in Fig. 7. Thanks to its M.-C. procedure, ADOM+ can significantly reduce the number of necessary gradient steps. Yet, consistently with our analysis in Prop. 3.3, our method is systematically better in all settings in terms of communications.

In conclusion, while several methods can share similar convergence rates, ours is the only one to perform at least as well as its competitors in every setting for different graph’s topology, as predicted by Tab. 1.

6 Conclusion

In this work, we have proposed a novel stochastic algorithm for the decentralized optimization of a sum of smooth and strongly convex functions. We have demonstrated, theoretically and empirically, that this algorithm leads systematically to a substantial acceleration compared to state-of-the-art works. Furthermore, our algorithm is asynchronous, decoupled, primal, and does not rely on an extra inner loop: each of these properties makes it suitable for real applications. We also proposed a novel approach for minimizing the accelerated communication rate of asynchronous gossip algorithms. In future work, we would like to explore the robustness of such algorithms to more challenging variabilities occurring in real-life applications such as time-varying networks and to extend our work to less regular functions and stochastic analysis.

Acknowledgements

This work was supported by Project ANR-21-CE23-0030 ADONIS and EMERG-ADONIS from Alliance SU. In addition, the authors would like to thank Raphaël Berthier, Mathieu Even, Hadrien Hendrikx, and Dmitry Kovalev for their helpful discussions.

References

- [1] Ludwig Arnold. Stochastic differential equations. *New York*, 1974.
- [2] Necdet Serhat Aybat and Mert Gürbüzbalaban. Decentralized computation of effective resistances and acceleration of consensus algorithms. In *2017 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*, pages 538–542, 2017.
- [3] Eugene Belilovsky, Michael Eickenberg, and Edouard Oyallon. Decoupled greedy learning of CNNs. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 736–745. PMLR, 13–18 Jul 2020.

- [4] Eugene Belilovsky, Louis Leconte, Lucas Caccia, Michael Eickenberg, and Edouard Oyallon. Decoupled greedy learning of cnns for synchronous and asynchronous distributed learning. *arXiv preprint arXiv:2106.06401*, 2021.
- [5] S. Boyd, A. Ghosh, B. Prabhakar, and D. Shah. Randomized gossip algorithms. *IEEE Transactions on Information Theory*, 52(6):2508–2530, 2006.
- [6] Bugra Can, Saeed Soori, Necdet Serhat Aybat, Maryam Mehri Dehnavi, and Mert Gürbüzbalaban. Randomized gossiping with effective resistance weights: Performance guarantees and applications. *IEEE Transactions on Control of Network Systems*, 9(2):524–536, 2022.
- [7] Ashok K Chandra, Prabhakar Raghavan, Walter L Ruzzo, Roman Smolensky, and Prason Tiwari. The electrical resistance of a graph captures its commute and cover times. *computational complexity*, 6(4):312–340, 1996.
- [8] Laurent Condat, Daichi Kitahara, Andrés Contreras, and Akira Hirabayashi. Proximal splitting algorithms for convex optimization: A tour of recent advances, with new twists, 2019.
- [9] Easley David and Kleinberg Jon. *Networks, Crowds, and Markets: Reasoning About a Highly Connected World*. Cambridge University Press, USA, 2010.
- [10] Steven Diamond and Stephen Boyd. CVXPY: A Python-embedded modeling language for convex optimization. *Journal of Machine Learning Research*, 17(83):1–5, 2016.
- [11] Wendy Ellens, Floske M Spieksma, Piet Van Mieghem, Almerima Jamakovic, and Robert E Kooij. Effective graph resistance. *Linear algebra and its applications*, 435(10):2491–2506, 2011.
- [12] Mathieu Even, Raphaël Berthier, Francis Bach, Nicolas Flammarion, Hadrien Hendrikx, Pierre Gaillard, Laurent Massoulié, and Adrien Taylor. A continuized view on nesterov acceleration for stochastic gradient descent and randomized gossip. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.
- [13] Mathieu Even, Hadrien Hendrikx, and Laurent Massoulié. Decentralized optimization with heterogeneous delays: a continuous-time approach. *arXiv preprint arXiv:2106.03585*, 2021.
- [14] Mathieu Even, Hadrien Hendrikx, and Laurent Massoulié. Asynchrony and acceleration in gossip algorithms, 2020.
- [15] Marguerite Frank. The braess paradox. *Mathematical Programming*, 20(1):283–302, 1981.

- [16] Arpita Ghosh, Stephen Boyd, and Amin Saberi. Minimizing effective resistance of a graph. *SIAM review*, 50(1):37–66, 2008.
- [17] Eduard Gorbunov, Alexander Rogozin, Aleksandr Beznosikov, Darina Dvinskikh, and Alexander Gasnikov. *Recent Theoretical Advances in Decentralized Distributed Convex Optimization*, pages 253–325. Springer International Publishing, Cham, 2022.
- [18] Hadrien Hendrikx. A principled framework for the design and analysis of token algorithms, 2022.
- [19] Hadrien Hendrikx, Francis Bach, and Laurent Massoulié. Accelerated decentralized optimization with local updates for smooth and strongly convex objectives. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 897–906. PMLR, 2019.
- [20] Hadrien Hendrikx, Francis Bach, and Laurent Massoulié. An accelerated decentralized stochastic proximal algorithm for finite sums. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [21] Hadrien Hendrikx, Francis Bach, and Laurent Massoulié. Dual-free stochastic decentralized optimization with variance reduction. *Advances in neural information processing systems*, 33:19455–19466, 2020.
- [22] Hadrien Hendrikx, Francis Bach, and Laurent Massoulié. An optimal algorithm for decentralized finite-sum optimization. *SIAM Journal on Optimization*, 31(4):2753–2783, 2021.
- [23] Douglas J Klein. Resistance-distance sum rules. *Croatica chemica acta*, 75(2):633–649, 2002.
- [24] Douglas J Klein and Milan Randić. Resistance distance. *Journal of mathematical chemistry*, 12(1):81–95, 1993.
- [25] Anastasia Koloskova, Nicolas Loizou, Sadra Boreiri, Martin Jaggi, and Sebastian Stich. A unified theory of decentralized sgd with changing topology and local updates. In *International Conference on Machine Learning*, pages 5381–5393. PMLR, 2020.
- [26] Anastasiia Koloskova, Tao Lin, and Sebastian U Stich. An improved analysis of gradient tracking for decentralized machine learning. *Advances in Neural Information Processing Systems*, 34:11422–11435, 2021.
- [27] Dmitry Kovalev, Elnur Gasanov, Alexander Gasnikov, and Peter Richtárik. Lower bounds and optimal algorithms for smooth and strongly convex decentralized optimization over time-varying networks. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.

- [28] Dmitry Kovalev, Alexander Gasnikov, and Peter Richtárik. Accelerated primal-dual gradient method for smooth and convex-concave saddle-point problems with bilinear coupling. *arXiv preprint arXiv:2112.15199*, 2021.
- [29] Dmitry Kovalev, Adil Salim, and Peter Richtarik. Optimal and practical algorithms for smooth and strongly convex decentralized optimization. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 18342–18352. Curran Associates, Inc., 2020.
- [30] Dmitry Kovalev, Egor Shulgin, Peter Richtárik, Alexander V Rogozin, and Alexander Gasnikov. Adom: accelerated decentralized optimization method for time-varying networks. In *International Conference on Machine Learning*, pages 5784–5793. PMLR, 2021.
- [31] Günter Last and Mathew Penrose. *Lectures on the Poisson process*, volume 7. Cambridge University Press, 2017.
- [32] Jonas Latz. Analysis of stochastic gradient descent in continuous time. *Statistics and Computing*, 31(4):1–25, 2021.
- [33] Rémi Leblond, Fabian Pedregosa, and Simon Lacoste-Julien. Improved asynchronous parallel optimization analysis for stochastic incremental methods. *arXiv preprint arXiv:1801.03749*, 2018.
- [34] Huan Li and Zhouchen Lin. Accelerated gradient tracking over time-varying graphs for decentralized optimization. *arXiv preprint arXiv:2104.02596*, 2021.
- [35] Hongzhou Lin, Julien Mairal, and Zaid Harchaoui. A universal catalyst for first-order optimization. *Advances in neural information processing systems*, 28, 2015.
- [36] Konstantin Mishchenko, Grigory Malinovsky, Sebastian Stich, and Peter Richtárik. Proxskip: Yes! local gradient steps provably lead to communication acceleration! finally! *arXiv preprint arXiv:2202.09357*, 2022.
- [37] Angelia Nedic, Alex Olshevsky, and Wei Shi. Achieving geometric convergence for distributed optimization over time-varying graphs. *SIAM Journal on Optimization*, 27(4):2597–2633, 2017.
- [38] Yurii Nesterov. *Introductory lectures on convex optimization: A basic course*, volume 87. Springer Science & Business Media, 2003.
- [39] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle,

- A. Beygelzimer, F. d Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- [40] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [41] Adil Salim, Laurent Condat, Dmitry Kovalev, and Peter Richtárik. An optimal algorithm for strongly convex minimization under affine constraints. In *International Conference on Artificial Intelligence and Statistics*, pages 4482–4498. PMLR, 2022.
- [42] Kevin Scaman, Francis Bach, Sébastien Bubeck, Yin Tat Lee, and Laurent Massoulié. Optimal algorithms for smooth and strongly convex distributed optimization in networks. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 3027–3036. PMLR, 06–11 Aug 2017.
- [43] Wei Shi, Qing Ling, Gang Wu, and Wotao Yin. Extra: An exact first-order algorithm for decentralized consensus optimization. *SIAM Journal on Optimization*, 25(2):944–966, 2015.
- [44] Zhuoqing Song, Lei Shi, Shi Pu, and Ming Yan. Optimal gradient tracking for decentralized optimization. *arXiv preprint arXiv:2110.05282*, 2021.
- [45] Sebastian U. Stich and Sai Praneeth Karimireddy. The error-feedback framework: Sgd with delayed gradients. *Journal of Machine Learning Research*, 21(237):1–36, 2020.
- [46] Vaya Sapobi Samui Vos. Methods for determining the effective resistance. *Mestrado, Mathematisch Instituut Universiteit Leiden*, 2016.
- [47] Lin Xiao and S. Boyd. Fast linear iterations for distributed averaging. In *42nd IEEE International Conference on Decision and Control (IEEE Cat. No.03CH37475)*, volume 5, pages 4997–5002 Vol.5, 2003.
- [48] Lin Xiao, Stephen Boyd, and Seung-Jean Kim. Distributed average consensus with least-mean-square deviation. *Journal of Parallel and Distributed Computing*, 67(1):33–46, 2007.

Appendix

Table of Contents

A	Communication bounds	21
B	Saddle Point Reformulation	23
C	Proof of the main Theorem (Th. 3.2)	23
	C.1 Proof of the Lemma C.1	26
D	Comparison of complexities with related work	31
E	Optimizing the communication rates with a SDP	34
	E.1 Convexity lemmas	35
	E.2 Problem relaxation	36
	E.3 SDP formulation	37
	E.4 Optimal weights for the Barbell Graph	38
	E.5 Results from our SDP solver	42
F	Physical interpretation of the sufficient condition	43
G	DADAO vs MSDA on the star graph	43
H	Practical Implementation	44
	H.1 Simulating the Poisson Point Processes	44
	H.2 Pseudo Code	45

A Communication bounds

The following properties will be used all along the proofs of the Lemmas and Theorems and are related to the communication of our nodes.

Lemma A.1. *Under the assumptions of Theorem 3.2, if $z_0, \tilde{z}_0 \in \text{span}(\pi)$, then $z_t, \tilde{z}_t \in \text{span}(\pi)$ almost surely.*

Proof. It's clear that for any i, j , we get:

$$\pi(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top = (\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top.$$

Thus, the variations of (z_t, \tilde{z}_t) belong to $\text{span}(\pi)$, and so is the trajectory. \square

Lemma A.2 (Effective resistance contraction). *For $(i, j) \in \mathcal{E}$ and any $x \in \mathbb{R}^d$, we have:*

$$\|(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top x\|_{\Lambda^+}^2 \leq \chi_2 \|(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top x\|.$$

Proof. Indeed, we note that:

$$\|(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top x\|_{\Lambda^+}^2 = x^\top (\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top \Lambda^+ (\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top x \quad (9)$$

$$\leq 2\chi_2 x^\top (\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top x \quad (10)$$

$$= \chi_2 \|(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top x\|^2 \quad (11)$$

□

Lemma A.3 (Bound on the resistance).

For $(i, j) \in \mathcal{E}$, $(e_i - e_j)^\top \Lambda^+ (e_i - e_j) \leq \frac{1}{\lambda_{ij} + \lambda_{ji}}$.

Proof. For a graph such that $|\bar{\mathcal{E}}| = 1$, the inequality is trivial. Now, we assume that there are other edges than (i, j) or (j, i) . Thus, we have:

$$\Lambda \succcurlyeq (\lambda_{ij} + \lambda_{ji})(e_i - e_j)(e_i - e_j)^\top + \tilde{\lambda}\tilde{\pi}$$

where $\tilde{\pi}(e_i - e_j) = 0$, $\tilde{\pi}\mathbf{1} = 0$, $\text{rank}(\tilde{\pi}) = n - 1$, $\tilde{\pi}$ orthogonal projector and $\tilde{\lambda} > 0$ small enough. In this case: $\left((\lambda_{ij} + \lambda_{ji})(e_i - e_j)(e_i - e_j)^\top + \tilde{\lambda}\tilde{\pi}\right)^+ \succcurlyeq \Lambda^+$, but this implies that $\frac{1}{\lambda_{ij} + \lambda_{ji}} \geq (e_i - e_j)^\top \Lambda^+ (e_i - e_j)$. □

Proof of Lemma 3.1. First, we note that Λ is symmetric and has a non-negative spectrum, as:

$$x^\top \Lambda x = \sum_{(ij) \in \mathcal{E}} \lambda_{ij} \|x_i - x_j\|^2.$$

From this, we also clearly see that $\chi_1 = +\infty$ iff the graph is disconnected. Next, assuming that the graph is connected, 0 is an eigenvalue of Λ with multiplicity 1 and by definition of χ_1 , we have $\chi_1 \geq \chi_2$. As we also have the following:

$$\sum_{(i,j) \in \mathcal{E}} \lambda_{ij} (e_i - e_j)^\top \Lambda^+ (e_i - e_j) = \text{Tr}(\Lambda^+ \Lambda) = n - 1,$$

we can write:

$$n - 1 \leq 2\chi_2 \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} = \chi_2 \text{Tr} \Lambda$$

and get $\frac{n-1}{\text{Tr} \Lambda} \leq \chi_2$. Finally, note that: $\text{Tr}(\Lambda) = 2 \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} \leq 2|\bar{\mathcal{E}}| \sup_{(i,j) \in \mathcal{E}} (\lambda_{ij} + \lambda_{ji})$ and $\text{Tr}(\Lambda) \leq (n-1)\|\Lambda\|$, so that, using Lemma A.3:

$$\sqrt{\chi_2} \text{Tr}(\Lambda) \leq \frac{1}{\sqrt{2 \inf_{(i,j) \in \mathcal{E}} (\lambda_{ij} + \lambda_{ji})}} \sqrt{\text{Tr} \Lambda} \sqrt{\text{Tr} \Lambda} \quad (12)$$

$$\leq \sqrt{\frac{\text{Tr} \Lambda}{2 \inf_{(i,j) \in \mathcal{E}} (\lambda_{ij} + \lambda_{ji})}} \sqrt{2|\bar{\mathcal{E}}| \sup_{(i,j) \in \mathcal{E}} (\lambda_{ij} + \lambda_{ji})} \quad (13)$$

$$\leq \sqrt{\kappa} \sqrt{|\bar{\mathcal{E}}|} \sqrt{(n-1) \|\Lambda\|} \quad (14)$$

□

B Saddle Point Reformulation

With $0 < \nu < \mu$ and introducing an extra dual variable \hat{x} , we get:

$$\begin{aligned} \inf_{x \in \mathbb{R}^d} \sum_{i=1}^n f_i(x) &= \inf_{\substack{x, \hat{x} \in \mathbb{R}^{n \times d} \\ x = \hat{x}, \pi \hat{x} = 0}} \sum_{i=1}^n f_i(x_i) - \frac{\nu}{2} \|x\|^2 + \frac{\nu}{2} \|\hat{x}\|^2 \\ &= \inf_{x, \hat{x} \in \mathbb{R}^{n \times d}} \sup_{y, z \in \mathbb{R}^{n \times d}} \sum_{i=1}^n f_i(x_i) - \frac{\nu}{2} \|x\|^2 + \frac{\nu}{2} \|\hat{x}\|^2 + \langle y, \hat{x} - x \rangle + \langle z, \pi \hat{x} \rangle \\ &= \inf_{x \in \mathbb{R}^{n \times d}} \sup_{y, z \in \mathbb{R}^{n \times d}} \inf_{\hat{x} \in \mathbb{R}^{n \times d}} \sum_{i=1}^n f_i(x_i) - \frac{\nu}{2} \|x\|^2 + \frac{\nu}{2} \|\hat{x}\|^2 + \langle y, \hat{x} - x \rangle + \langle z, \pi \hat{x} \rangle \\ &= \inf_{x \in \mathbb{R}^{n \times d}} \sup_{y, z \in \mathbb{R}^{n \times d}} \sum_{i=1}^n f_i(x_i) - \frac{\nu}{2} \|x\|^2 - \langle x, y \rangle - \frac{1}{2\nu} \|\pi z + y\|^2. \end{aligned}$$

C Proof of the main Theorem (Th. 3.2)

Basic reminds about the Bregman divergence For a smooth convex function F ,

$$d_F(x, y) \triangleq F(x) - F(y) - \langle \nabla F(y), x - y \rangle$$

is its Bregman divergence. Next, we set $\nu = \frac{\mu}{2}$ such that, for $F(x) = \sum_{i=1}^n f_i(x_i) - \frac{\nu}{2} \|x\|^2$, then F is L -smooth, ν -strongly convex, and we get:

$$\frac{1}{2L} \|\nabla F(x) - \nabla F(y)\|^2 \leq d_F(x, y) \leq \frac{L}{2} \|x - y\|^2,$$

and

$$\frac{\nu}{2} \|x - y\|^2 \leq d_F(x, y) \leq \frac{1}{2\nu} \|\nabla F(x) - \nabla F(y)\|^2,$$

Proof of Theorem 3.2. We introduce for a positive semi-definite matrix A , $\|x\|_A \triangleq x^\top A x$. For our proof, we rely on the notation of Eq. 6, and we introduce

$X \triangleq (x, \tilde{x}, \tilde{y}), Y \triangleq (y, z, \tilde{z})$ and the following Lyapunov potential:

$$\begin{aligned} \Phi(t, X, Y) \triangleq & A_t d_F(x, x^*) + \tilde{A}_t \|\tilde{x} - x^*\|^2 + B_t \|y - y^*\|^2 + \tilde{B}_t \|\tilde{y} - y^*\|^2 \\ & + C_t \|z + y - z^* - y^*\|^2 + \tilde{C}_t \|\tilde{z} - z^*\|_{\Lambda^+}^2, \end{aligned}$$

where $A_t, \tilde{A}_t, B_t, \tilde{B}_t, C_t, \tilde{C}_t$ are non-negative functions to be defined. We will use this potential to control the trajectories.

Because Φ is smooth, the SDE is a smooth trajectory, we get via Ito's formula [31] applied to the semi-martingale (X_t, Y_t) on any intervals $[0, T]$:

$$\begin{aligned} \Phi(T, X_T, Y_T) = & \Phi(0, X_0, Y_0) + \int_0^T \langle \nabla \Phi(t, X_t, Y_t), \begin{pmatrix} 1 \\ a_1(X_t, Y_t) \\ a_2(X_t, Y_t) \end{pmatrix} \rangle dt \\ & + \sum_{i=1}^n \int_0^T (\Phi(t, X_t + b_1^i(X_t), Y_t) - \Phi(t, X_t, Y_t)) dt \\ & + \sum_{(i,j) \in \mathcal{E}} \int_0^T (\Phi(t, X_t, Y_t + b_2^{ij}(Y_t)) - \Phi(t, X_t, Y_t)) \lambda_{ij} dt \\ & + \Theta_T, \end{aligned}$$

where the following quantity is a Martingale:

$$\begin{aligned} \Theta_u \triangleq & \sum_{i=1}^n \int_0^u (\Phi(t, X_{t-}, Y_{t-} + b_1^i(X_{t-})) - \Phi(t, X_{t-}, Y_{t-})) (dN_i(t) - dt) \\ & + \sum_{(i,j) \in \mathcal{E}} \int_0^u (\Phi(t, X_{t-}, Y_{t-} + b_2^{ij}(X_{t-}, Y_{t-})) - \Phi(t, X_{t-}, Y_{t-})) (dM_{ij}(t) - \lambda_{ij} dt). \end{aligned}$$

Taking the expectation, we get that, as the initialization is deterministic:

$$\begin{aligned} \mathbb{E}[\Phi(T, X_T, Y_T)] = & \Phi(0, X_0, Y_0) + \int_0^T \langle \nabla \Phi(t, X_t, Y_t), \begin{pmatrix} 1 \\ a_1(X_t, Y_t) \\ a_2(X_t, Y_t) \end{pmatrix} \rangle dt \\ & + \sum_{i=1}^n \int_0^T (\Phi(t, X_t + b_1^i(X_t), Y_t) - \Phi(t, X_t, Y_t)) dt \\ & + \sum_{(i,j) \in \mathcal{E}} \int_0^T (\Phi(t, X_t, Y_t + b_2^{ij}(Y_t)) - \Phi(t, X_t, Y_t)) \lambda_{ij} dt, \end{aligned}$$

To show that the integrand term is negative, we will use the following technical Lemma, which is also difficult to prove and whose proof is deferred to Appendix C.1:

Lemma C.1. *If:*

$$\begin{aligned} \eta &= \frac{1}{8}\sqrt{\frac{\nu}{L}} & \gamma &= \frac{1}{4L} & \delta &= \frac{1}{4}\sqrt{\frac{\nu}{L}} & \alpha &= \frac{1}{4}\sqrt{\frac{\nu}{L}} & \beta &= \frac{1}{2} & \theta &= \frac{1}{2}\sqrt{\frac{L}{\nu}} \\ \tilde{\eta} &= \frac{1}{8}\sqrt{\frac{\nu}{L}} & \tilde{\gamma} &= \frac{1}{4\sqrt{\nu L}} & \tilde{\delta} &= 1 & \tilde{\alpha} &= \frac{1}{8}\sqrt{\frac{\nu}{L}} & \tilde{\beta} &= 2\chi_1[\Lambda]\sqrt{\frac{L}{\nu}} & \nu &= \frac{\mu}{2} \end{aligned}$$

and

$$\begin{aligned} A_t &= e^{\frac{t}{8\sqrt{2}}\sqrt{\frac{\mu}{L}}} & \tilde{A}_t &= \frac{\mu}{8}e^{\frac{t}{8\sqrt{2}}\sqrt{\frac{\mu}{L}}} & \tilde{B}_t &= \frac{1}{8L}e^{\frac{t}{8\sqrt{2}}\sqrt{\frac{\mu}{L}}} \\ B_t &= \frac{1}{16L}e^{\frac{t}{8\sqrt{2}}\sqrt{\frac{\mu}{L}}} & C_t &= \frac{1}{2\mu}e^{\frac{t}{8\sqrt{2}}\sqrt{\frac{\mu}{L}}} & \tilde{C}_t &= \frac{1}{32\chi_1 L}e^{\frac{t}{8\sqrt{2}}\sqrt{\frac{\mu}{L}}}. \end{aligned}$$

then:

$$\begin{aligned} &\langle \nabla\Phi(t, X_t, Y_t), \begin{pmatrix} 1 \\ a_1(X_t, Y_t) \\ a_2(X_t, Y_t) \end{pmatrix} \rangle + (\Phi(t, X_t + b_1(X_t), Y_t) - \Phi(t, X_t, Y_t)) \\ &+ \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} (\Phi(t, X_t, Y_t + b_2^{ij}(Y_t)) - \Phi(t, X_t, Y_t)) \leq 0 \text{ a.s. } . \end{aligned}$$

Now, we remark that if we have $F(x) = f(x) - \frac{\mu}{4}\|x\|^2$, and initialize with $\tilde{x}_0 = x_0$, $y_0 = \tilde{y}_0 = \nabla F(x_0)$ and $z_0 = \tilde{z}_0 = -\pi \nabla F(x_0)$, then, given the linear relation between $A_t, \tilde{A}_t, B_t, \tilde{B}_t, C_t, \tilde{C}_t$, the L smoothness and the fact π is an orthogonal projection, we get:

$$\begin{aligned} \Phi(0, X_0, Y_0) &\leq d_F(x_0, x^*) + \frac{\mu}{8}\|x_0 - x^*\|^2 + \frac{1}{16L}\|\nabla F(x_0) - \nabla F(x^*)\|^2 \\ &+ \frac{1}{8L}\|\nabla F(x_0) - \nabla F(x^*)\|^2 + \frac{1}{2\mu}\|(\mathbf{I} - \pi)(\nabla F(x_0) - \nabla F(x^*))\|^2 \\ &+ \frac{1}{32}\frac{\chi_1(0)}{L\chi_1}\|\pi(\nabla F(x_0) - \nabla F(x^*))\|^2 \\ &\leq \left(\frac{L}{2} + \frac{\mu}{8} + \frac{L}{16} + \frac{L}{8} + \frac{L^2}{2\mu} + \frac{L}{32}\right)\|x_0 - x^*\|^2 \end{aligned}$$

In particular, as $\frac{\mu}{4}\|x - x^*\|^2 \leq d_F(x, x^*)$ it implies that:

$$\mathbb{E}[\|x_t - x^*\|^2] \leq \left(\frac{1}{2} + \frac{23L}{8\mu} + 2\frac{L^2}{\mu^2}\right)\|x_0 - x^*\|^2 e^{-\frac{t}{8\sqrt{2}}\sqrt{\frac{\mu}{L}}}$$

Finally, we note that the expected number of gradients between $[0, T]$ is given by:

$$\mathbb{E}\left[\sum_{i=1}^n N_i(T)\right] = nT,$$

and similarly, the number of edges activated is given by:

$$\mathbb{E}\left[\sum_{1 \leq i, j \leq n} M_{ij}(T)\right] = \sum_{1 \leq i, j \leq n} \int_0^T \lambda_{ij} dt = \frac{T}{2} \text{Tr } \Lambda.$$

□

C.1 Proof of the Lemma C.1

We first state a couple of inequalities that we will combine to obtain a bound on our Lyapunov function.

Lemma C.2. *First:*

$$\begin{aligned} \phi_A &\triangleq A_t(d_F(x^+, x^*) - d_F(x, x^*)) + \tilde{A}_t(\|\tilde{x}^+ - x^*\|^2 - \|\tilde{x} - x^*\|^2) \\ &\quad + \eta A_t \langle \tilde{x} - x, \nabla F(x) - \nabla F(x^*) \rangle + 2\tilde{\eta} \tilde{A}_t \langle x - \tilde{x}, \tilde{x} - x^* \rangle \end{aligned} \quad (15)$$

$$\begin{aligned} &\leq \|\nabla F(x) - \tilde{y}\|^2 \left(A_t \frac{L\gamma^2}{2} - A_t \gamma + \tilde{A}_t \tilde{\gamma}^2 \right) \\ &\quad + A_t \gamma \langle \nabla F(x) - \tilde{y}, y^* - \tilde{y} \rangle + 2\tilde{\gamma} \tilde{A}_t \langle \tilde{y} - y^*, \tilde{x} - x^* \rangle \\ &\quad - 2\tilde{\gamma} \tilde{A}_t (d_F(\tilde{x}, x^*) + d_F(x^*, x) - d_F(\tilde{x}, x)) \\ &\quad - \eta A_t (d_F(\tilde{x}, x) + d_F(x, x^*) - d_F(\tilde{x}, x^*)) - \tilde{A}_t \tilde{\eta} \|\tilde{x} - x^*\|^2 + \tilde{A}_t \tilde{\eta} \|x - x^*\|^2 \end{aligned} \quad (16)$$

Proof. First, we have to use optimality conditions and smoothness, as well as the separability of F :

$$d_F(x^+, x^*) - d_F(x, x^*) = d_F(x^+, x) - \langle x^+ - x, \nabla F(x^*) - \nabla F(x) \rangle \quad (17)$$

$$\leq \frac{L}{2} \|x^+ - x\|^2 - \langle x^+ - x, \nabla F(x^*) - \nabla F(x) \rangle \quad (18)$$

$$\begin{aligned} &= \frac{L\gamma^2}{2} \|\tilde{y} - \nabla F(x)\|^2 - \gamma \|\nabla F(x) - \tilde{y}\|^2 \\ &\quad + \gamma \langle \nabla F(x) - \tilde{y}, y^* - \tilde{y} \rangle \end{aligned} \quad (19)$$

Next, we note that, again using optimality conditions:

$$\|\tilde{x}^+ - x^*\|^2 - \|\tilde{x} - x^*\|^2 = 2\langle \tilde{x}^+ - \tilde{x}, \tilde{x} - x^* \rangle + \|\tilde{x}^+ - \tilde{x}\|^2 \quad (20)$$

$$= -2\tilde{\gamma} \langle \nabla F(x) - \tilde{y}, \tilde{x} - x^* \rangle + \tilde{\gamma}^2 \|\nabla F(x) - \tilde{y}\|^2 \quad (21)$$

$$= -2\tilde{\gamma} \langle \nabla F(x) - \nabla F(x^*), \tilde{x} - x^* \rangle + 2\tilde{\gamma} \langle \tilde{y} - y^*, \tilde{x} - x^* \rangle + \tilde{\gamma}^2 \|\nabla F(x) - \tilde{y}\|^2 \quad (22)$$

$$= -2\tilde{\gamma} (d_F(\tilde{x}, x^*) + d_F(x^*, x) - d_F(\tilde{x}, x)) + 2\tilde{\gamma} \langle \tilde{y} - y^*, \tilde{x} - x^* \rangle + \tilde{\gamma}^2 \|\nabla F(x) - \tilde{y}\|^2 \quad (23)$$

Momentum in x associated with the term $d_F(x, x^*)$ gives:

$$\eta \langle \tilde{x} - x, \nabla F(x) - \nabla F(x^*) \rangle = -\eta (d_F(\tilde{x}, x) + d_F(x, x^*) - d_F(\tilde{x}, x^*)) \quad (24)$$

and momentum in \tilde{x} associated with $\|\tilde{x} - x^*\|^2$ leads to:

$$2\tilde{\eta} \langle x - \tilde{x}, \tilde{x} - x^* \rangle = -2\tilde{\eta} \|\tilde{x} - x^*\|^2 + 2\tilde{\eta} \langle x - x^*, \tilde{x} - x^* \rangle \leq -\tilde{\eta} \|\tilde{x} - x^*\|^2 + \tilde{\eta} \|x - x^*\|^2 \quad (25)$$

□

Lemma C.3. Next, we show that if $\alpha B_t = \frac{\delta}{2} \tilde{B}_t$:

$$\begin{aligned}
\phi_B &\triangleq B_t(\|y^+ - y^*\|^2 - \|y - y^*\|^2) + \tilde{B}_t(\|\tilde{y}^+ - y^*\|^2 - \|\tilde{y} - y^*\|^2) \\
&\quad + 2\alpha B_t \langle y - y^*, \tilde{y} - y \rangle - 2\theta \tilde{B}_t \langle y + z + \nu \tilde{x}, \tilde{y} - y^* \rangle + 2\alpha C_t \langle \tilde{y} - y, z + y - y^* - z^* \rangle \\
&\hspace{15em} (26) \\
&\leq -\frac{\delta}{2} \tilde{B}_t \|\tilde{y} - y^*\|^2 - \frac{\delta}{2} \tilde{B}_t \|y - y^*\|^2 - 2\tilde{\delta} \tilde{B}_t \langle \nabla F(x) - \tilde{y}, y^* - \tilde{y} \rangle \\
&\quad + \delta \tilde{B}_t \|\nabla F(x) - \nabla F(x^*)\|^2 + \left((\delta + \tilde{\delta})^2 - \delta \right) \tilde{B}_t \|\nabla F(x) - y\|^2 \\
&\quad - 2\theta \tilde{B}_t \langle y + z - y^* - z^*, \tilde{y} - y^* \rangle - 2\theta \nu \tilde{B}_t \langle \tilde{x} - x^*, \tilde{y} - y^* \rangle + 2\alpha C_t \langle \tilde{y} - y, z + y - y^* - z^* \rangle \\
&\hspace{15em} (27)
\end{aligned}$$

Proof. Using optimality conditions:

$$\begin{aligned}
\|\tilde{y}^+ - y^*\|^2 - \|\tilde{y} - y^*\|^2 &= 2\langle \tilde{y} - y^*, \tilde{y}^+ - \tilde{y} \rangle + \|\tilde{y}^+ - \tilde{y}\|^2 \hspace{10em} (28) \\
&= 2\delta \langle \nabla F(x) - \tilde{y}, \tilde{y} - y^* \rangle + 2\tilde{\delta} \langle \nabla F(x) - \tilde{y}, \tilde{y} - y^* \rangle + (\delta + \tilde{\delta})^2 \|\nabla F(x) - \tilde{y}\|^2 \\
&\hspace{15em} (29) \\
&= -2\tilde{\delta} \langle \nabla F(x) - \tilde{y}, y^* - \tilde{y} \rangle + \delta \|\nabla F(x) - \nabla F(x^*)\|^2 - \delta \|\tilde{y} - y^*\|^2 \\
&\quad + \left((\delta + \tilde{\delta})^2 - \delta \right) \|\nabla F(x) - \tilde{y}\|^2 \hspace{10em} (30)
\end{aligned}$$

The momentum in \tilde{y} associated with the term $\|\tilde{y} - y^*\|^2$ gives:

$$-2\theta \tilde{B}_t \langle y + z + \nu \tilde{x}, \tilde{y} - y^* \rangle = -2\theta \tilde{B}_t \langle y + z - y^* - z^*, \tilde{y} - y^* \rangle - 2\theta \nu \tilde{B}_t \langle \tilde{x} - x^*, \tilde{y} - y^* \rangle \hspace{5em} (31)$$

The momentum in y associated with the term $\|y - y^*\|^2$ gives:

$$2\alpha B_t \langle \tilde{y} - y, y - y^* \rangle = -\alpha B_t \|y - y^*\|^2 - \alpha B_t \|\tilde{y} - y\|^2 + \alpha B_t \|\tilde{y} - y^*\|^2 \hspace{5em} (32)$$

and the one associated with $\|y + z - y^* - z^*\|^2$:

$$2\alpha C_t \langle \tilde{y} - y, z + y - y^* - z^* \rangle \hspace{10em} (33)$$

□

Lemma C.4. Finally, assuming $\theta \tilde{B}_t = \tilde{\beta} \tilde{C}_t = \alpha C_t$, letting $1 \geq \tilde{\tau} > 0$,

$$\begin{aligned}
z_{ij}^+ &= z - \beta(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top(y + z) \text{ and } \tilde{z}_{ij}^+ = \tilde{z} - \tilde{\beta}(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top(y + z): \\
\phi_C - 2\theta\tilde{B}_t\langle y + z - y^* - z^*, \tilde{y} - y^* \rangle + 2\alpha C_t\langle \tilde{y} - y, z + y - y^* - z^* \rangle &\triangleq \\
\sum_{ij} \lambda_{ij} C_t \left(\|y + z_{ij}^+ - y^* - z^*\|^2 - \|y + z - y^* - z^*\|^2 \right) & \\
+ \sum_{ij} \lambda_{ij} \tilde{C}_t \left(\|\tilde{z}_{ij}^+ - z^*\|_{\Lambda^+}^2 - \|\tilde{z} - z^*\|_{\Lambda^+}^2 \right) + 2\tilde{\alpha}\tilde{C}_t\langle z - \tilde{z}, \tilde{z} - z^* \rangle_{\Lambda^+} & \quad (34) \\
+ 2\alpha C_t\langle \tilde{z} + \tilde{y} - z - y, z + y - y^* - z^* \rangle - 2\theta\tilde{B}_t\langle y + z - y^* - z^*, \tilde{y} - y^* \rangle & \\
\leq \tilde{\beta}^2 \chi_2 \tilde{C}_t \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} \|(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top(y + z)\|^2 & \\
+ \beta(\beta - 1) C_t \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} \|(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top(y + z)\|^2 & \\
- \alpha C_t \|y + z - y^* - z^*\|^2 + \tilde{\alpha} \chi_1 \tilde{C}_t \|z - z^*\|^2 - \tilde{\alpha} \tilde{C}_t \|\tilde{z} - z^*\|_{\Lambda^+}^2 & \\
- \tilde{\tau} \frac{1}{2} \tilde{\beta} \frac{\nu}{L} \tilde{C}_t \|z - z^*\|^2 + \tilde{\tau} \frac{\nu}{L} \frac{2\alpha\theta}{\delta} B_t \|y - y^*\|^2 & \quad (35)
\end{aligned}$$

Proof. Having in mind that $\pi(y^* + z^*) = 0$ and $\mathbf{\Lambda}^+ \mathbf{\Lambda} = \pi$, we get, using Lemma A.1 and Lemma A.2 on the inequality (39):

$$\Delta_{\tilde{z}} \triangleq \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} (\|\tilde{z}_{ij}^+ - z^*\|_{\Lambda^+}^2 - \|\tilde{z} - z^*\|_{\Lambda^+}^2) \quad (36)$$

$$= \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} 2\langle \tilde{z} - z^*, \tilde{z}_{ij}^+ - \tilde{z} \rangle_{\Lambda^+} + \|\tilde{z}_{ij}^+ - \tilde{z}\|_{\Lambda^+}^2 \quad (37)$$

$$\begin{aligned}
&= -2\tilde{\beta} \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} \langle \tilde{z} - z^*, (\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top(y + z - y^* - z^*) \rangle_{\Lambda^+} \\
&\quad + \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} \tilde{\beta}^2 \|(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top(y + z)\|_{\Lambda^+}^2 \quad (38)
\end{aligned}$$

$$\leq -2\tilde{\beta} \langle \tilde{z} - z^*, \mathbf{\Lambda}^+ \mathbf{\Lambda}(y + z) \rangle + \chi_2 \tilde{\beta}^2 \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} \|(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top(y + z)\|^2 \quad (39)$$

$$= -2\tilde{\beta} \langle \tilde{z} - z^*, \pi(y + z) \rangle + \chi_2 \tilde{\beta}^2 \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} \|(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top(y + z)\|^2 \quad (40)$$

We also have, as $y^+ = y$:

$$\Delta_z \triangleq \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} (\|y^+ + z_{ij}^+ - y^* - z^*\|^2 - \|y + z - y^* - z^*\|^2) \quad (41)$$

$$= 2 \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} \langle y + z_{ij}^+ - y - z, y + z - y^* - z^* \rangle + \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} \|y + z_{ij}^+ - y - z\|^2 \quad (42)$$

$$= -2 \sum_{(i,j) \in \mathcal{E}} \beta \lambda_{ij} \langle (\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top (y + z), y + z - y^* - z^* \rangle + \sum_{(i,j) \in \mathcal{E}} \beta^2 \lambda_{ij} \|(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top (y + z)\|^2 \quad (43)$$

$$= \beta(\beta - 1) \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} \|(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top (y + z)\|^2 \quad (44)$$

The momentum in \tilde{z} associated with $\|\tilde{z} - z^*\|_{\Lambda^+}^2$ gives:

$$2\tilde{\alpha}\tilde{C}_t \langle z - \tilde{z}, \tilde{z} - z^* \rangle_{\Lambda^+} \leq \tilde{\alpha}\chi_1 \tilde{C}_t \|z - z^*\|^2 - \tilde{\alpha}\tilde{C}_t \|\tilde{z} - z^*\|_{\Lambda^+}^2 \quad (45)$$

And the one in z associated with $\|y + z - y^* - z^*\|^2$ gives:

$$2\alpha C_t \langle \tilde{z} - z, z + y - y^* - z^* \rangle \quad (46)$$

Then, assuming that $\theta\tilde{B}_t = \tilde{\beta}\tilde{C}_t = \alpha C_t$, we have:

$$2\alpha C_t \langle \tilde{y} - y, z + y - y^* - z^* \rangle - 2\tilde{\beta}\tilde{C}_t \langle \tilde{z} - z^*, y + z - y^* - z^* \rangle - 2\theta\tilde{B}_t \langle y + z - y^* - z^*, \tilde{y} - y^* \rangle + 2\alpha C_t \langle \tilde{z} - z, z + y - y^* - z^* \rangle \quad (47)$$

$$= -2\alpha C_t \|y + z - y^* - z^*\|^2 \quad (48)$$

At this stage, we split the negative term (48) into two halves, upper-bounding one of the halves by remembering that $\frac{\nu}{L} \leq 1$ and introducing $1 \geq \tilde{\tau} > 0$:

$$-\alpha C_t \|y + z - y^* - z^*\|^2 \leq -\tilde{\tau} \frac{\nu}{L} \alpha C_t \|y + z - y^* - z^*\|^2 \quad (49)$$

$$= -\tilde{\tau} \tilde{\beta} \frac{\nu}{L} \tilde{C}_t \|y + z - y^* - z^*\|^2 \quad (50)$$

$$\leq -\tilde{\tau} \frac{1}{2} \tilde{\beta} \frac{\nu}{L} \tilde{C}_t \|z - z^*\|^2 + \tilde{\tau} \tilde{\beta} \frac{\nu}{L} \tilde{C}_t \|y - y^*\|^2 \quad (51)$$

$$= -\tilde{\tau} \frac{1}{2} \tilde{\beta} \frac{\nu}{L} \tilde{C}_t \|z - z^*\|^2 + \tilde{\tau} \frac{\nu}{L} \frac{2\alpha\theta}{\delta} B_t \|y - y^*\|^2 \quad (52)$$

□

Keeping in mind that $\theta\tilde{B}_t = \tilde{\beta}\tilde{C}_t = \alpha C_t$ and $\frac{\xi}{2}\tilde{B}_t = \alpha B_t$, we put everything together. Defining $\Psi = \frac{\partial\Phi}{\partial t} + \phi_A + \phi_B + \phi_C$, we have:

$$\Psi \leq \|\nabla F(x) - \tilde{y}\|^2 \left(A_t \frac{L\gamma^2}{2} - A_t\gamma + \tilde{A}_t\tilde{\gamma}^2 + \left((\delta + \tilde{\delta})^2 - \delta \right) \tilde{B}_t \right) \quad (53)$$

$$+ \|\tilde{z} - z^*\|_{\mathbf{A}^+}^2 \left(-\tilde{\alpha}\tilde{C}_t + \tilde{C}'_t \right) \quad (54)$$

$$+ \|\tilde{y} - y^*\|^2 \left(\tilde{B}'_t - \frac{\delta}{2}\tilde{B}_t \right) \quad (55)$$

$$+ \|x - x^*\|^2 \left(\tilde{A}_t\tilde{\eta} - \tilde{A}_t \frac{\nu\tilde{\gamma}}{2} \right) \quad (56)$$

$$+ \|\tilde{x} - x^*\|^2 \left(\tilde{A}'_t - \tilde{A}_t\tilde{\eta} \right) \quad (57)$$

$$+ \|\nabla F(x) - \nabla F(x^*)\|^2 \left(\delta\tilde{B}_t - \frac{\tilde{\gamma}}{2L}\tilde{A}_t \right) \quad (58)$$

$$+ \sum_{(i,j) \in \mathcal{E}} \lambda_{ij} \|(\mathbf{e}_i - \mathbf{e}_j)(\mathbf{e}_i - \mathbf{e}_j)^\top (y + z)\|^2 \left(\chi_2 \tilde{\beta}^2 \tilde{C}_t + \beta(\beta - 1)C_t \right) \quad (59)$$

$$+ \|z - z^*\|^2 \left(\chi_1 \tilde{\alpha} - \tilde{\tau} \frac{1}{2} \tilde{\beta} \frac{\nu}{L} \right) \tilde{C}_t \quad (60)$$

$$+ \|y - y^*\| \left(B'_t - \left(1 - \tilde{\tau} \frac{\nu}{L} \frac{2\theta}{\delta} \right) \alpha B_t \right) \quad (61)$$

$$+ \|y + z - y^* - z^*\|^2 \left(C'_t - \alpha C_t \right) \quad (62)$$

$$+ d_F(x, x^*) \left(A'_t - \eta A_t \right) \quad (63)$$

$$+ d_F(\tilde{x}, x) \left(-A_t\eta + 2\tilde{\gamma}\tilde{A}_t \right) \quad (64)$$

$$+ d_F(\tilde{x}, x^*) \left(A_t\eta - 2\tilde{\gamma}\tilde{A}_t \right) \quad (65)$$

$$+ \langle \nabla F(x) - \tilde{y}, y^* - \tilde{y} \rangle \left(-2\tilde{\delta}\tilde{B}_t + \gamma A_t \right) \quad (66)$$

$$+ \langle \tilde{y} - y^*, \tilde{x} - x^* \rangle \left(2\tilde{\gamma}\tilde{A}_t - 2\theta\nu\tilde{B}_t \right) \quad (67)$$

Resolution

Proof of Lemma C.1. Our goal is to put to zero all of the terms appearing next to scalar products and make the factors of positive quantities (norms or divergences) less or equal to zero. Given our relations, we guess that each exponential has the same rate. Thus, with $\tau > 0$, we fix $\frac{\delta}{2} = \tilde{\eta} = \eta = \tilde{\alpha} = \tau\sqrt{\frac{\nu}{L}}$, which leads to $\tilde{\gamma} = \frac{2\tau}{\sqrt{\nu L}}$ using Eq. (56). Also, from Eq. (65):

$$4\tilde{A}_t = \nu A_t.$$

Next, from Eq. (58) and Eq. (67), it's necessary that:

$$2L\delta = \theta\nu,$$

thus $\theta = 4\tau\sqrt{\frac{L}{\nu}}$. From Eq. (67), we get:

$$\tilde{A}_t = 2L\nu\tilde{B}_t.$$

Combining this previous equation with Eq. (66), as $4\tilde{A}_t = \nu A_t$, we have $\tilde{\delta} = 4L\gamma$. Next, Eq. (53) gives, with the equations above:

$$\begin{aligned} A_t\left(\frac{L\gamma^2}{2} - \gamma\right) + \tilde{A}_t\tilde{\gamma}^2 + \left((\delta + \tilde{\delta})^2 - \delta\right)\tilde{B}_t &= A_t\frac{L\gamma^2}{2} - A_t\gamma + \frac{\nu}{4}\tilde{\gamma}^2A_t + \left(\delta^2 + \tilde{\delta}^2 + \delta\right)\frac{A_t}{8L} \\ &= A_t\left(\frac{L\gamma^2}{2} - \gamma + \frac{\nu}{4}\frac{4\tau^2}{\nu L}\right) + A_t\left(2\tau\sqrt{\frac{\nu}{L}} + 4\tau^2\frac{\nu}{L} + 16L^2\gamma^2\right)\frac{1}{8L} \\ &\leq A_t\left(\gamma^2\frac{5}{2}L - \gamma + \frac{5}{4}\frac{\tau^2}{L} + \frac{\sqrt{2}}{8}\frac{\tau}{L}\right) \end{aligned}$$

We thus pick $\gamma = \frac{1}{4L}$ and $\tau = \frac{1}{8}$, so that $\tilde{\delta} = 1$. Via Eq. (61), we fix $\tilde{\tau} = \frac{1}{8} < 1$. With Eq. (60), we then get:

$$\tilde{\beta} = 2\chi_1\sqrt{\frac{L}{\nu}}$$

We also put $\alpha = 2\tau\sqrt{\frac{\nu}{L}}$ and only one last equation, Eq. (59), needs to be satisfied, for which we pick $\beta = \frac{1}{2}$:

$$\chi_2\tilde{\beta}^2\tilde{C}_t + \beta(\beta - 1)C_t = \left(\chi_2\tilde{\beta}\alpha - \frac{1}{4}\right)C_t$$

This implies that $\chi_2\chi_1 \leq \frac{1}{2}$. In summary, we set:

$$\begin{array}{llllllll} \eta = \frac{1}{8}\sqrt{\frac{\nu}{L}} & \gamma = \frac{1}{4L} & \delta = \frac{1}{4}\sqrt{\frac{\nu}{L}} & \alpha = \frac{1}{4}\sqrt{\frac{\nu}{L}} & \beta = \frac{1}{2} & \theta = \frac{1}{2}\sqrt{\frac{L}{\nu}} & \tilde{\tau} = \frac{1}{8} \\ \tilde{\eta} = \frac{1}{8}\sqrt{\frac{\nu}{L}} & \tilde{\gamma} = \frac{1}{4\sqrt{\nu L}} & \tilde{\delta} = 1 & \tilde{\alpha} = \frac{1}{8}\sqrt{\frac{\nu}{L}} & \tilde{\beta} = 2\chi_1\sqrt{\frac{L}{\nu}} & \nu = \frac{\mu}{2} & \tau = \frac{1}{8} \end{array}$$

Now, let's pick: $A_t = e^{t\tau}\sqrt{\frac{\nu}{L}} = e^{\frac{t}{8\sqrt{2}}}\sqrt{\frac{\mu}{L}}$ so that:

$$\begin{aligned} A_t &= e^{\frac{t}{8\sqrt{2}}}\sqrt{\frac{\mu}{L}} & \tilde{A}_t &= \frac{\mu}{8}e^{\frac{t}{8\sqrt{2}}}\sqrt{\frac{\mu}{L}} & \tilde{B}_t &= \frac{1}{8L}e^{\frac{t}{8\sqrt{2}}}\sqrt{\frac{\mu}{L}} \\ B_t &= \frac{1}{16L}e^{\frac{t}{8\sqrt{2}}}\sqrt{\frac{\mu}{L}} & C_t &= \frac{1}{2\mu}e^{\frac{t}{8\sqrt{2}}}\sqrt{\frac{\mu}{L}} & \tilde{C}_t &= \frac{1}{32\chi_1L}e^{\frac{t}{8\sqrt{2}}}\sqrt{\frac{\mu}{L}}. \end{aligned}$$

This implies that $\Psi \leq 0$. □

D Comparison of complexities with related work

Proof of Prop. 3.3. We consider the settings of concurrent works and, given any gossip matrix admissible for them, we show that DADAO has better rates.

Comparison with ADOM+ [27]. The ADOM+ setting is:

- *Gossip matrices:* Laplacians \mathcal{L} with $\|\mathcal{L}x - x\|^2 \leq (1 - \frac{1}{\chi})\|x\|^2$ for $\chi \geq 1$ and $x \in \mathbf{1}^\perp$.
- *Total number of gradients to reach ϵ precision:* $\mathcal{O}(n\sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon})$.

- *Total number of edges activated to reach ϵ precision:* $\mathcal{O}(|\mathcal{E}|\chi\sqrt{\frac{L}{\mu}}\log\frac{1}{\epsilon})$.

If we take an eigenvector of \mathcal{L} for the eigenvalue $\frac{1}{\chi_1[\mathcal{L}]}$, then the Laplacian inequality directly leads to $\left(1 - \frac{1}{\chi_1[\mathcal{L}]}\right)^2 \leq 1 - \frac{1}{\chi}$ and we have:

$$\frac{1}{\chi_1[\mathcal{L}]} \left(\frac{1}{\chi_1[\mathcal{L}]} - 2 \right) \leq -\frac{1}{\chi},$$

leading to:

$$\frac{2}{\chi_1[\mathcal{L}]} \geq \frac{1}{\chi_1[\mathcal{L}]} \left(2 - \frac{1}{\chi_1[\mathcal{L}]} \right) \geq \frac{1}{\chi}.$$

Thus, $\chi_1[\mathcal{L}] \leq 2\chi$. Remind that, by definition, $\chi_2 \leq \chi_1$. Note that in ADOM+, we also have $\|\mathcal{L}\| \leq 2$, so that $\text{Tr}(\mathcal{L}) \leq 2n$. Then, for any Laplacian matrix \mathcal{L} valid for ADOM+, we consider for DADAO a Λ defined as followed:

$$\Lambda = \sqrt{2\chi_1[\mathcal{L}]\chi_2[\mathcal{L}]} \mathcal{L}$$

Then, DADAO has the same gradient complexity as ADOM+, but the number of communications of DADAO is:

$$\frac{T}{2} \text{Tr} \Lambda \leq \frac{1}{2} \sqrt{2\chi_1[\mathcal{L}]\chi_2[\mathcal{L}]} \sqrt{\frac{L}{\mu}} \log\left(\frac{1}{\epsilon}\right) 2n = \mathcal{O}(|\mathcal{E}|\chi\sqrt{\frac{L}{\mu}}\log\frac{1}{\epsilon})$$

Consequently, DADAO is better than ADOM+ for all valid configurations of ADOM+ (in the fixed graph topology setting) and DADAO. We note that for the complete graph, $\chi = \mathcal{O}(1)$ and $|\mathcal{E}| = \mathcal{O}(n^2)$, whereas $n\sqrt{2\chi_1[\mathcal{L}]\chi_2[\mathcal{L}]} = \mathcal{O}(n)$ and DADAO has strictly better communication rates than ADOM+ on this graph.

Comparison with Gradient Tracking methods AGT, OGT [34, 44].

The setting is described by:

- *Gossip matrices:* Any matrix $\mathcal{L} = \mathbf{I} - W$ with W symmetric doubly-stochastic, *i.e.* such that $\sum_i W_{ij} = 1$ and $\sum_j W_{ij} = 1$.
- *Total number of gradients to reach ϵ precision:* $\mathcal{O}(n\sqrt{\frac{L}{\mu}}\log\frac{1}{\epsilon})$.
- *Total number of edges activated to reach ϵ precision:* $\mathcal{O}(|\mathcal{E}|\frac{1}{\sqrt{\theta}}\sqrt{\frac{L}{\mu}}\log\frac{1}{\epsilon})$, where $\theta = 1 - \|W - \frac{1}{n}\mathbf{1}\mathbf{1}^\top\|$.

If $\kappa = \mathcal{O}(\frac{|\mathcal{E}|}{n})$, using Lemma 3.1, we have:

$$\frac{1}{2} \sqrt{2\chi_1[\mathcal{L}]\chi_2[\mathcal{L}]} \text{Tr} \mathcal{L} = \mathcal{O}(|\mathcal{E}|\sqrt{\rho[\mathcal{L}]})$$

Furthermore, as W is stochastic, $\|\mathcal{L}\| \leq 2$ and $\theta \leq \frac{1}{\chi_1[\mathcal{L}]}$, leading to $\rho \leq \frac{2}{\theta}$. Consequently, the communication complexity of DADAO run for $\mathcal{O}(\sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon})$ iterations recovers the rate of GT. Furthermore, for the complete graph, for any Laplacian \mathcal{L} with $\rho[\mathcal{L}] = \mathcal{O}(1)$ (remind that, by definition $\rho \geq 1$), we have $\sqrt{2\chi_1[\mathcal{L}]\chi_2[\mathcal{L}]} \text{Tr } \mathcal{L} = \mathcal{O}(n)$ whereas $|\mathcal{E}| = \mathcal{O}(n^2)$, thus DADAO uses an order of magnitude less communications than GT need to.

Note that for the Star-graph, there is no admissible Laplacian in the framework of [44].

Comparison with MSDA [42]. The setting is described by:

- *Gossip matrices*: any Laplacians admissible for DADAO.
- *Total number of gradients to reach ϵ precision*: $\mathcal{O}(n\sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon})$.
- *Total number of edges activated to reach ϵ precision*: $\mathcal{O}(|\mathcal{E}|\sqrt{\rho}\sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon})$.

If $\kappa = \mathcal{O}(\frac{|\mathcal{E}|}{n})$, using Lemma 3.1, we see that, for any Laplacian \mathcal{L} :

$$\frac{1}{2}\sqrt{2\chi_1[\mathcal{L}]\chi_2[\mathcal{L}]} \text{Tr } \mathcal{L} = \mathcal{O}(|\mathcal{E}|\sqrt{\rho})$$

Consequently, the communication complexity of DADAO run for $\mathcal{O}(\sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon})$ iterations is better than MSDA. Furthermore, for the complete graph, for any Laplacian \mathcal{L} with $\rho[\mathcal{L}] = \mathcal{O}(1)$ (remind that, by definition $\rho \geq 1$), we have $\sqrt{2\chi_1[\mathcal{L}]\chi_2[\mathcal{L}]} \text{Tr } \mathcal{L} = \mathcal{O}(n)$ whereas $|\mathcal{E}| = \mathcal{O}(n^2)$, thus DADAO uses an order of magnitude less communications than MSDA need to.

Comparison with the Continuized framework [12]. In this framework and using their notations, we have:

- *Gossip matrices*: all Laplacian matrices \mathcal{L} verifying $\text{Tr } \mathcal{L} = 2$.
- *Total number of gradients to reach ϵ precision*: $\mathcal{O}(\frac{1}{\theta'_{\text{ARG}}}\sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon})$.
- *Total number of edges activated to reach ϵ precision*: $\mathcal{O}(\frac{1}{\theta'_{\text{ARG}}}\sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon})$.

First, we slightly rephrase one of the proposition of [12] to match our notations:

Lemma D.1. *For a Laplacian \mathcal{L} with $\text{Tr } \mathcal{L} = 2$, the communication and gradient complexity of the continuized framework is given by:*

$$\mathcal{O}(\sqrt{\chi_1[\mathcal{L}]\chi_2[\mathcal{L}]}\sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon}).$$

Proof. Under the notation of [12], $\theta'_{\text{ARG}} = \sqrt{\mu_{\text{gossip}} / \max_{\{v,w\}} \frac{R_{vw}}{P_{vw}}}$, with $\mu_{\text{gossip}} = \frac{1}{\chi_1[\mathcal{L}]}$ in our setting. Moreover, we remind that in [12], $\mathcal{L} = AA^\top$ and that $Ae_{vw} = \sqrt{P_{vw}}(e_v - e_w)$. Thus, by definition:

$$\frac{R_{vw}}{P_{vw}} \triangleq \frac{e_{vw}^\top A^+ Ae_{vw}}{P_{vw}} \quad (68)$$

$$= \frac{e_{vw}^\top A^+(e_v - e_w)}{\sqrt{P_{vw}}} \quad (69)$$

$$= \frac{(A^{+\top} e_{vw})^\top (e_v - e_w)}{\sqrt{P_{vw}}} \quad (70)$$

$$= \frac{((AA^\top)^{+\top} Ae_{vw})^\top (e_v - e_w)}{\sqrt{P_{vw}}} \quad (71)$$

$$= (e_v - e_w)^\top \mathcal{L}^+ (e_v - e_w), \quad (72)$$

and we have $\theta'_{\text{ARG}} = \frac{1}{\sqrt{2\chi_1[\mathcal{L}]\chi_2[\mathcal{L}]}}$. \square

Next, if $\text{Tr } \mathcal{L} = 2$, we proved in Lemma 3.1 that $2\sqrt{\chi_1[\mathcal{L}]\chi_2[\mathcal{L}]} \geq (n-1)$. Thus, we see that the gradient complexity of DADAO in $\mathcal{O}(n\sqrt{\frac{L}{\mu}} \log \frac{1}{\epsilon})$ is always better than the one of the continuized framework. If we write $f(n) = \Omega(g(n))$ the fact that there is a constant $C > 0$ and $n_0 \in \mathbb{N}$ such that $\forall n \geq n_0, Cg(n) \leq f(n)$, then, in the cycle graph, for any \mathcal{L} with $\text{Tr } \mathcal{L} = 2$, $\chi_1[\mathcal{L}] = \Omega(n^3)$ and $\chi_2[\mathcal{L}] = \Omega(n)$. Thus DADAO uses an order of magnitude less gradients than the continuized framework for this graph. \square

E Optimizing the communication rates with a SDP

We remind the following useful Lemma:

Lemma E.1 (Generalized Schur complement). *Let M be a symmetric matrix defined as:*

$$M = \begin{pmatrix} A & B \\ B^\top & C \end{pmatrix}.$$

Then:

- $M \succcurlyeq 0 \Leftrightarrow A \succ 0$, $C - B^\top A^{-1} B \succcurlyeq 0$.
- $M \succcurlyeq 0 \Leftrightarrow A \succcurlyeq 0$, $(\mathbf{I} - AA^+)B = 0$, $C - B^\top A^+ B \succcurlyeq 0$.

E.1 Convexity lemmas

Proposition E.2 (log-convexity lemma). *Fix $\|u\| = 1$, $u \perp \mathbf{1}$ and let $\lambda = (\lambda_{ij}) \in \mathbb{R}_+^{n^2}$ such that $\Lambda(\lambda) = \sum_{1 \leq i, j \leq n} \lambda_{ij} (e_i - e_j)(e_i - e_j)^\top$, and let:*

$$\psi(\lambda, u) = u^\top \Lambda(\lambda) u.$$

Then, $\lambda \rightarrow \psi(\lambda, u)$ is log-convex on $\mathbb{R}_+^{n^2}$. Furthermore, $\lambda \rightarrow \psi(\lambda, u)$ is strictly convex on the convex set $\{\lambda, \chi_1[\Lambda(\lambda)] < \infty\}$ of the connected graphs.

Next, we prove our convexity bounds:

Proof of Lemma E.2. We want to show that $\lambda \rightarrow \psi(\lambda, u)$ is log-convex on $\mathbb{R}_+^{n^2}$ and strictly convex on the convex set $\{\lambda, \chi_1(\Lambda(\lambda)) < \infty\}$ of the connected graphs. Having in mind that:

$$\frac{d}{dt} X^+(t) = -X^+(t) \frac{d}{dt} X(t) X^+(t),$$

we note that:

$$\frac{\partial}{\partial \lambda_{ij}} \log \psi(\Lambda) = -\frac{u^\top \Lambda^+(e_i - e_j)(e_i - e_j)^\top \Lambda^+ u}{u^\top \Lambda^+ u} \quad (73)$$

$$= -\frac{(u^\top \Lambda^+(e_i - e_j))^2}{u^\top \Lambda^+ u} \quad (74)$$

and next:

$$\begin{aligned} \frac{\partial^2}{\partial \lambda_{ij} \partial \lambda_{i'j'}} \log \psi(\Lambda) &= \frac{2u^\top \Lambda^+(e_i - e_j) u^\top \Lambda^+(e_{i'} - e_{j'}) (e_{i'} - e_{j'})^\top \Lambda^+(e_i - e_j) u^\top \Lambda^+ u}{(u^\top \Lambda^+ u)^2} \\ &\quad - \frac{(u^\top \Lambda^+(e_i - e_j))^2 (u^\top \Lambda^+(e_{i'} - e_{j'}))^2}{(u^\top \Lambda^+ u)^2} \end{aligned} \quad (75)$$

Now, let the $(\delta_{ij})_{1 \leq i, j \leq n}$ be positive and introduce $\Delta = \sum_{ij} \delta_{ij} (e_i - e_j)(e_i - e_j)^\top$. Then, with $\mathcal{H}(\psi)$ the hessian of $\log \psi$, we observe that:

$$\Delta^\top \mathcal{H}(\psi) \Delta = \sum_{ij, i'j'} \delta_{ij} \delta_{i'j'} \frac{\partial^2}{\partial \lambda_{ij} \partial \lambda_{i'j'}} \psi(\Lambda) \quad (76)$$

$$= \frac{2u^\top \Lambda^+ \Delta \Lambda^+ \Delta \Lambda^+ u (u^\top \Lambda^+ u) - (u^\top \Lambda^+ \Delta \Lambda^+ u)^2}{(u^\top \Lambda^+ u)^2} \quad (77)$$

As $u^\top \Lambda^+ \Delta \Lambda^+ \Delta \Lambda^+ u (u^\top \Lambda^+ u) \geq 0$, $\Delta^\top \mathcal{H}(\psi) \Delta$ would be positive if we show

$$u^\top \Lambda^+ \Delta \Lambda^+ \Delta \Lambda^+ u (u^\top \Lambda^+ u) - (u^\top \Lambda^+ \Delta \Lambda^+ u)^2 \geq 0. \quad (78)$$

As $u \perp \mathbf{1}$, by Schur complement, this would be the case if:

$$\begin{pmatrix} u^\top \Lambda^+ \Delta \Lambda^+ \Delta \Lambda^+ u & u^\top \Lambda^+ \Delta \Lambda^+ u \\ u^\top \Lambda^+ \Delta \Lambda^+ u & u^\top \Lambda^+ u \end{pmatrix} \succcurlyeq 0, \quad (79)$$

but this also writes:

$$\begin{pmatrix} u \\ u \end{pmatrix}^\top \begin{pmatrix} \Lambda^+ \Delta \Lambda^+ \Delta \Lambda^+ & \Lambda^+ \Delta \Lambda^+ \\ \Lambda^+ \Delta \Lambda^+ & \Lambda^+ \end{pmatrix} \begin{pmatrix} u \\ u \end{pmatrix} \succcurlyeq 0. \quad (80)$$

For this to be true for any u , by Schur complement, this would be implied by:

$$\Lambda^+ \succcurlyeq 0 \quad , \quad (\mathbf{I} - \Lambda^+ \Lambda) \Lambda^+ \Delta \Lambda^+ = 0$$

and

$$\Lambda^+ \Delta \Lambda^+ \Delta \Lambda^+ - \Lambda^+ \Delta \Lambda^+ \Lambda \Lambda^+ \Delta \Lambda^+ \succcurlyeq 0.$$

As $\Lambda^+ \Lambda = \pi$, all conditions are met. In particular, the last inequality is an equality, and so is (78). Thus, having:

$$\Delta^\top \mathcal{H}(\psi) \Delta = 0$$

would imply that:

$$u^\top \Lambda^+ \Delta \Lambda^+ \Delta \Lambda^+ u u^\top \Lambda^+ u = 0$$

As $u \neq 0$, if Δ represents a connected graph, this can never be equal to 0. \square

Corollary E.2.1 (Strict convexity of the relaxed communication rate). *Fix a directed graph \mathcal{E} so that $\tilde{\mathcal{E}}$ is connected. Then, $\lambda \rightarrow \sqrt{\chi_1[\Lambda(\lambda)]\chi_2^\mathcal{E}(\lambda)}$ is strictly convex on $\mathbb{R}_+^\mathcal{E}$.*

Proof of Corollary E.2.1. We note that the supremum of logarithmic convex functions is logarithmically convex, thus $\lambda \rightarrow \chi_1[\Lambda(\lambda)]$, $\lambda \rightarrow \chi_2[\Lambda(\lambda)]$ are logarithmically convex and so is their square root. Furthermore, $\lambda \rightarrow \chi_2[\Lambda(\lambda)]$ is a finite supremum (ie, a maximum) of strictly convex function, thus it is strictly convex and so is $\lambda \rightarrow \sqrt{\chi_1[\Lambda(\lambda)]\chi_2[\Lambda(\lambda)]}$. \square

E.2 Problem relaxation

We want to consider the following:

$$\begin{aligned} & \text{minimize} \quad \text{Tr} \Lambda(\lambda) \sqrt{2\chi_1[\Lambda(\lambda)]\chi_2[\Lambda(\lambda)]} \\ & \text{subject to} \quad \lambda \in \mathcal{C} \cap \mathbb{R}_+^\mathcal{E}. \end{aligned} \quad (81)$$

For \mathcal{C}_1 , by homogeneity (the quantity of interest is invariant by rescaling), the optimal solution can be picked as equality in the inequality constraints. For $\mathcal{C}_2, \mathcal{C}_3$, we have an obvious upper bound on the trace, so we can consider:

$$\begin{aligned} & \text{minimize} \quad \sqrt{2\chi_1[\Lambda(\lambda)]\chi_2[\Lambda(\lambda)]} \\ & \text{subject to} \quad \lambda \in \mathcal{C} \cap \mathbb{R}_+^\mathcal{E}. \end{aligned} \quad (82)$$

However, $\lambda \rightarrow \chi_2[\Lambda(\lambda)]$ is not convex as explained above. Hence, we relax the problem by considering:

$$\begin{aligned} & \text{minimize} \quad \sqrt{2\chi_1[\Lambda(\lambda)]\chi_2^\mathcal{E}(\lambda)} \\ & \text{subject to} \quad \lambda \in \mathcal{C} \cap \mathbb{R}_+^\mathcal{E}. \end{aligned} \quad (83)$$

Remark E.1. We note that we have here a Braess paradox [15]: minimizing Eq. 82 compared to minimizing Eq. 83 can lead to a substantially better communication rate. However, finding a solution to problem 82 requires in practice knowing the support of the optimal graph in advance. If some edges are dropped during the optimization procedure (i.e., allocated with a 0 weight), minimizing Eq. 83 with the optimal support would lead to the optimal communication rate.

E.3 SDP formulation

Proof of Prop. 4.1. First, note that $\varphi(u) \triangleq ut_1 + \frac{t_2}{u}$ satisfies $\varphi'(u) = t_1 - \frac{t_2}{u^2}$ so that the minimum is reached in $u^* = \sqrt{\frac{t_2}{t_1}}$, with $\varphi(u^*) = 2\sqrt{t_1 t_2}$. Now, if $t_2 \geq t_1$, $u^* \geq 1$ and we note that $\varphi(u^*) = \inf_{u \geq 1} \varphi(u)$.

Our initial problem writes:

$$\begin{aligned} & \text{minimize} && \sqrt{2\|\Lambda^+\| \frac{1}{2} \sup_{(ij) \in \mathcal{E}} (e_i - e_j)^\top \Lambda^+ (e_i - e_j)} \\ & \text{subject to} && \lambda \in \mathcal{C}, \lambda_{ij} \geq 0, \forall (i, j) \in \mathcal{E} \\ & && \Lambda = \sum_{ij} \lambda_{ij} (e_i - e_j)(e_i - e_j)^\top \end{aligned}$$

Thus, as $\chi_2 \leq \chi_1$, it is equivalent to solving:

$$\begin{aligned} & \text{minimize} && u \sup_{(ij) \in \mathcal{E}} (e_i - e_j)^\top \Lambda^+ (e_i - e_j) + \frac{\|\Lambda^+\|}{u} \\ & \text{subject to} && \lambda \in \mathcal{C}, u \geq 1, \lambda_{ij} \geq 0, \forall (i, j) \in \mathcal{E} \\ & && \Lambda = \sum_{ij} \lambda_{ij} (e_i - e_j)(e_i - e_j)^\top \end{aligned}$$

which also writes:

$$\begin{aligned} & \text{minimize} && ut_1 + \frac{t_2}{u} \\ & \text{subject to} && \lambda \in \mathcal{C}, u \geq 1, \lambda_{ij} \geq 0, \forall (i, j) \in \mathcal{E} \\ & && \Lambda = \sum_{ij} \lambda_{ij} (e_i - e_j)(e_i - e_j)^\top \\ & && t_1 \geq (e_i - e_j)^\top \Lambda^+ (e_i - e_j) \quad \forall (i, j) \in \mathcal{E} \\ & && t_2 \geq \|\Lambda^+\| \end{aligned}$$

equivalent to (with $\tilde{\lambda}_{ij} = u\lambda$ so that $\tilde{\Lambda} = u\Lambda$, $\tilde{\Lambda}^+ = \frac{1}{u}\Lambda^+$ and $\tilde{\lambda} \in u\mathcal{C}$):

$$\begin{aligned}
& \text{minimize} && ut_1 + \frac{t_2}{u} \\
& \text{subject to} && \tilde{\lambda} \in u\mathcal{C}, u \geq 1, \tilde{\lambda}_{ij} \geq 0, \forall (i, j) \in \mathcal{E} \\
& && \tilde{\Lambda} = \sum_{ij} \tilde{\lambda}_{ij} (e_i - e_j)(e_i - e_j)^\top \\
& && t_1 \geq u(e_i - e_j)^\top \tilde{\Lambda}^+ (e_i - e_j) \quad \forall (i, j) \in \mathcal{E} \\
& && t_2 \geq u \|\tilde{\Lambda}^+\|
\end{aligned}$$

equivalent to (with $t'_1 = ut_1, t'_2 = \frac{t_2}{u}$):

$$\begin{aligned}
& \text{minimize} && t'_1 + t'_2 \\
& \text{subject to} && \lambda \in u\mathcal{C}, u \geq 1, \lambda_{ij} \geq 0, \forall (i, j) \in \mathcal{E} \\
& && \Lambda = \sum_{ij} \lambda_{ij} (e_i - e_j)(e_i - e_j)^\top \\
& && t'_1 \geq u^2 (e_i - e_j)^\top \Lambda^+ (e_i - e_j) \quad \forall (i, j) \in \mathcal{E} \\
& && t'_2 \geq \|\Lambda^+\|
\end{aligned}$$

equivalent to, via Schur Complement:

$$\begin{aligned}
& \text{minimize} && t_1 + t_2 \\
& \text{subject to} && \lambda \in u\mathcal{C}, u \geq 1, \lambda_{ij} \geq 0, \forall (i, j) \in \mathcal{E} \\
& && \Lambda = \sum_{ij} \lambda_{ij} (e_i - e_j)(e_i - e_j)^\top \\
& && \begin{pmatrix} \Lambda & u(e_i - e_j) \\ u(e_i - e_j)^\top & t_1 \end{pmatrix} \succcurlyeq 0 \quad \forall (i, j) \in \mathcal{E} \\
& && \begin{pmatrix} \Lambda & \mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top \\ \mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top & t_2\mathbf{I} \end{pmatrix} \succcurlyeq 0
\end{aligned}$$

□

E.4 Optimal weights for the Barbell Graph

In this section, following [16], we make an in-depth study of the communication rates given by our SDP on the Barbell graph $K_n - K_n$ on $2n$ nodes. First, we propose an analytical solution to problem (83) under the global bandwidth constraint \mathcal{C}_1 . In this setting, we show that the optimal communication strategy can be unboundedly better than using uniform communication rates on the edges. Then, we illustrate the Braess paradox on the Barbell graph by showing that removing some edges leads to better convergence rates.

E.4.1 Analytical study

Using the same symmetry argument than [16], we know the problem actually depends on three variables a, b, c as illustrated on Fig. 3. We then have the

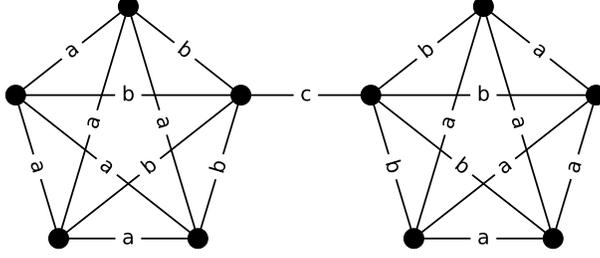


Figure 3: The barbell graph $K_5 - K_5$.

following Lemma, proved in Appendix E.4, which shows that using uniform weights is sub-optimal on the barbell graph.

Lemma E.3 (Barbell graph). *For $K_n - K_n$ with $n \geq 3$, under the constraint $\sum_{(i,j) \in \mathcal{E}} \lambda_{ij} \leq 1$, the optimal parameters for problem (83) are given by: $c^* = \frac{4-n+\sqrt{2n(n-1)}}{2(8+n)} \underset{n \rightarrow \infty}{\sim} \frac{\sqrt{2}-1}{2}$, $b^* = \frac{c^*}{n-2} \left(\sqrt{\frac{2(n-1)}{n}} - \frac{2}{n} \right) \underset{n \rightarrow \infty}{\sim} \frac{\sqrt{2}-1}{n\sqrt{2}}$, $a^* = \frac{1-c^*-2b^*(n-1)}{(n-1)(n-2)} \underset{n \rightarrow \infty}{\sim} \frac{\sqrt{2}-1}{2n^2}$. In this case, $\sqrt{2}\chi_1\chi_2 \underset{n \rightarrow \infty}{\simeq} 2.4n$. Using uniform weights $a = b = c = \frac{1}{|\mathcal{E}|}$ leads to $\sqrt{2}\chi_1\chi_2 \underset{n \rightarrow \infty}{\sim} \frac{n^{5/2}}{\sqrt{2}}$.*

Proof of Lemma E.3. Via a symmetry argument similar to [16], we know that the problem actually depends on three variables a, b, c .

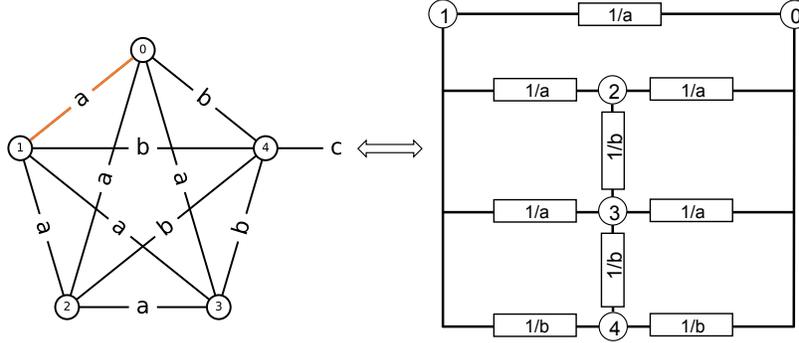


Figure 4: Equivalent resistance network if some current passes between node 0 and node 1 of the Barbell graph $K_5 - K_5$. By symmetry of the resistance network, no current passes through the edges 2 – 3 and 3 – 4 and we have $\frac{1}{2R_a} = a + \frac{2}{\frac{1}{a} + \frac{1}{a}} + \frac{1}{\frac{1}{b} + \frac{1}{b}} = 2a + \frac{b}{2}$.

By symmetry and using the Ohm law, we know that the resistance of the

corresponding edges are given by:

$$2R_c = \frac{1}{c}, \frac{1}{2R_a} = a + \frac{a}{2}(n-3) + \frac{b}{2} = \frac{1}{2}(a(n-1) + b), \frac{1}{2R_b} = b + \frac{1}{\frac{1}{a} + \frac{1}{b}}(n-2),$$

see Fig. 4 for an example of such computation on the Barbell graph $K_5 - K_5$. We can simplify the expressions as follows:

$$R_a = \frac{1}{a(n-1) + b}, \quad R_b = \frac{a+b}{2b}R_a, \quad R_c = \frac{1}{2c}.$$

Directly using the results of [16], the eigenvalues of the graph's Laplacian are:

$$a(n-1) + b, nb, \frac{1}{2}(nb + 2c \pm \sqrt{(2c + nb)^2 - 8bc})$$

and the global bandwidth constraint \mathcal{C}_1 translates into:

$$1 = c + 2(n-1)b + (n-1)(n-2)a.$$

From the drawing, it is natural that $a^* \leq b^* \leq c^*$.

Finding the effective resistance of the graph:

We want to find $\chi_2^* = \max(R_{a^*}, R_{b^*}, R_{c^*})$. We can directly eliminate R_{b^*} from the list as $b^* \geq a^*$, leading to $R_{b^*} \leq R_{a^*}$. From the constraints, we have:

$$\frac{1}{n-1} = \frac{c^*}{n-1} + 2b^* + (n-2)a^* \geq b^* + (n-1)a^*,$$

meaning $R_{a^*} \geq n-1$. Thus, if we assume from the drawing $c^* = \mathcal{O}(1)$, then $R_{c^*} = \mathcal{O}(1)$ and $\chi_2^* = R_{a^*}$. Let us make this assumption for now.

Finding the inverse of the smallest eigenvalue of the Laplacian:

We want to find $\frac{1}{\chi_1^*} = \min\left(a^*(n-1) + b^*, nb^*, \frac{1}{2}(nb^* + 2c^* \pm \sqrt{(2c^* + nb^*)^2 - 8b^*c^*})\right)$. As $b^* \geq a^*$ and for obvious reasons, we directly have:

$$\frac{1}{\chi_1^*} = \min\left(a^*(n-1) + b^*, \frac{1}{2}(nb^* + 2c^* - \sqrt{(2c^* + nb^*)^2 - 8b^*c^*})\right).$$

For now, we will also assume that $\chi_2^* < \chi_1^*$, and as $a^*(n-1) + b^* = \frac{1}{\chi_2^*}$, this means $\frac{1}{\chi_1^*} = \frac{1}{2}(nb^* + 2c^* - \sqrt{(2c^* + nb^*)^2 - 8b^*c^*})$.

Solving the system:

Now, we define $f(a, b, c) = \frac{1}{2\chi_1(a,b,c)\chi_2(a,b,c)}$. We want to find a^*, b^*, c^* solving the following:

$$\begin{aligned} &\text{maximize} && f(a, b, c) = (a(n-1) + b) \left(nb + 2c - \sqrt{(2c + nb)^2 - 8bc} \right) \\ &\text{subject to} && 1 = c + 2(n-1)b + (n-1)(n-2)a \end{aligned}$$

First, from the constraints, we have $a = \frac{2c}{(n-1)(n-2)} - \frac{2b}{n-2}$. Substituting a in $f(a, b, c)$ leads to:

$$f(a, b, c) = \left(\frac{1-c-bn}{n-2} \right) \left(nb+2c - \sqrt{(2c+nb)^2 - 8bc} \right).$$

Differentiating with respect to b and c leads to:

$$\begin{aligned} \frac{\partial f}{\partial b} = 0 &\Leftrightarrow n \left(nb+2c - \sqrt{(2c+nb)^2 - 8bc} \right) = (1-c-bn) \left(n - \frac{n(2c+nb)-4c}{\sqrt{(2c+nb)^2 - 8bc}} \right) \\ \frac{\partial f}{\partial c} = 0 &\Leftrightarrow \left(nb+2c - \sqrt{(2c+nb)^2 - 8bc} \right) = (1-c-bn) \left(2 - \frac{2(2c+nb)-4b}{\sqrt{(2c+nb)^2 - 8bc}} \right) \end{aligned}$$

Subtracting n times the second line to the first leads to the following system:

$$\begin{cases} \left(nb+2c - \sqrt{(2c+nb)^2 - 8bc} \right) &= (1-c-bn) \left(n - \frac{n(2c+nb)-4c}{\sqrt{(2c+nb)^2 - 8bc}} \right) \\ \sqrt{(2c+nb)^2 - 8bc} &= 2c+nb+4 \left(\frac{c}{n} - b \right) \end{cases}$$

Assuming $n > 3$, solving the system leads to:

$$\begin{aligned} c^* &= \frac{4-n+\sqrt{2n(n-1)}}{2(8+n)} \underset{n \rightarrow \infty}{\sim} \frac{\sqrt{2}-1}{2} \\ b^* &= \frac{c^*}{n-2} \left(\sqrt{\frac{2(n-1)}{n}} - \frac{2}{n} \right) \underset{n \rightarrow \infty}{\sim} \frac{\sqrt{2}-1}{n\sqrt{2}} \\ a^* &= \frac{1-c^*-2b^*(n-1)}{(n-1)(n-2)} \underset{n \rightarrow \infty}{\sim} \frac{\sqrt{2}-1}{2n^2}. \end{aligned}$$

We verify that indeed $c^* = \mathcal{O}(1)$ and that $\chi_2^* < \chi_1^*$, validating our assumptions. Finally, we compute:

$$\sqrt{2\chi_1^*\chi_2^*} = \frac{1}{\sqrt{f(a^*, b^*, c^*)}} = 2(\sqrt{2}-1) \sqrt{\frac{2(2+\sqrt{2})}{80-56\sqrt{2}}} n + \mathcal{O}(1).$$

The case for uniform weights:

Now, we assume $a = b = c = \frac{1}{|\mathcal{E}|}$. Then, we necessarily have $\chi_2 = \frac{1}{2c} = \frac{|\mathcal{E}|}{2} = \frac{n(n-1)+1}{2}$. We also evidently have $\frac{1}{\chi_1} = \frac{1}{2}(nb+2c - \sqrt{(2c+nb)^2 - 8bc}) = \frac{c}{2} \left(n+2 - \sqrt{(n+2)^2 - 8} \right)$. Thus, in the end:

$$\sqrt{2\chi_1\chi_2} = (n(n-1)+1) \sqrt{\frac{2}{n+2 - \sqrt{(n+2)^2 - 8}}} = \frac{n^{5/2}}{\sqrt{2}} + \mathcal{O}(\sqrt{n}).$$

□

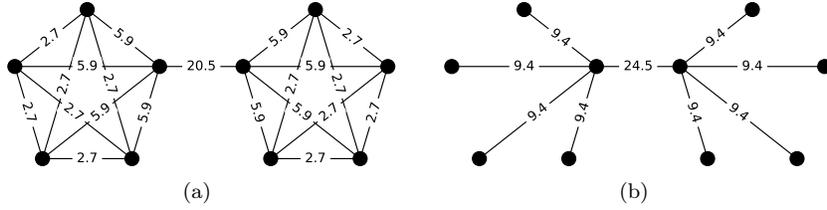


Figure 5: (a) Weights found by our SDP for the Barbell graph $K_5 - K_5$, leading to $\sqrt{\chi_1 \chi_2} \simeq 12.9$. (b) Optimal weights for the subgraph where edges a were removed, leading to a value of $\sqrt{\chi_1 \chi_2} \simeq 10.2$. For the sake of readability, all values were multiplied by a factor of 100.

E.4.2 Effect of the relaxation

As discussed in Rem.E.1, relaxing our problem from Eq. (82) to Eq. (83) prohibits the removal of edges from the given graph. However, as in traffic networks where doing so can sometimes speed up the overall traffic flow [9], considering sub-graphs might lead to better values of $\sqrt{\chi_1 \chi_2}$. As it happens, this is actually the case on the Barbell graph. As shown in Fig. 5, removing edges with weight a leads to a better communication complexity. Thus, the relaxation Eq. (83) cannot find the optimal solution of Eq. (82) in this example. However, it is still orders of magnitude better than putting uniform weights on the edges. In fact, a solution of problem Eq. (83) will always lead to faster rates than using any heuristics putting non 0 weight on the edges of the graph.

E.5 Results from our SDP solver

Efficient SDP solvers exist [10], and we display on Fig. 6 several results obtained thanks to the formulation of Prop. 4.1. Note that the solutions in Fig. 6(a,b)

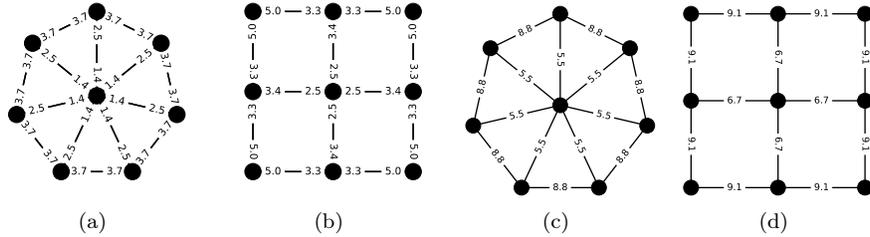


Figure 6: Optimal communication rates for the problem 83 obtained by solving the SDP from Prop. 4.1, under constraints $\forall i, \sum_{j, (i,j) \in \mathcal{E}} \lambda_{ij} \leq 1$ for (a,b) and $\sum_{(i,j) \in \mathcal{E}} \lambda_{ij} \leq 1$ for (c,d). For the sake of readability, we multiplied the values by 10 in (a,b), and by 100 in (c,d).

differ from what the simple heuristic of weighting each directed edge (i, j) with its degree $\frac{1}{d_i}$ would give or what would be prescribed with Metropolis weights, *i.e.* using $\frac{1}{1+\max\{d_i, d_j\}}$ [48, 43, 37, 17]. Also, solutions in Fig. 6(c,d) do not follow the naive approach of putting uniform weights $\frac{1}{|\mathcal{E}|}$ on every edge.

F Physical interpretation of the sufficient condition

Introducing $\lambda \triangleq \sum_{(ij) \in \mathcal{E}} \lambda_{ij}$ and $\mathcal{P}_\Lambda \triangleq \frac{2\Lambda}{\text{Tr } \Lambda}$, we write Λ as the product of these two more interpretable quantities to gain more insight on the condition $2\chi_1[\Lambda]\chi_2[\Lambda] \leq 1$:

$$\Lambda = \lambda \mathcal{P}_\Lambda. \quad (84)$$

In this setting, λ is the expected rate of communication over the *whole graph*, while \mathcal{P}_Λ can be interpreted as the Laplacian of \mathcal{E} with each edge weighted with its probability of having spiked given a communication happened in the graph. We have:

$$\chi_1[\Lambda] = \frac{\chi_1[\mathcal{P}_\Lambda]}{\lambda} \quad ; \quad \chi_2[\Lambda] = \frac{\chi_2[\mathcal{P}_\Lambda]}{\lambda}. \quad (85)$$

\mathcal{P}_Λ being normalized, we could say that the quantities $\chi_1[\mathcal{P}_\Lambda], \chi_2[\mathcal{P}_\Lambda]$ characterize the graph's connectivity while λ is the global rate of communication. Then, using (85), the condition $2\chi_1[\Lambda]\chi_2[\Lambda] \leq 1$ is equivalent to saying

$$\sqrt{2\chi_1[\mathcal{P}_\Lambda]\chi_2[\mathcal{P}_\Lambda]} \leq \lambda, \quad (86)$$

meaning that the global communication rate should be larger than some spectral quantity quantifying the graph's connectivity. If structural constraints on the network makes it impossible to verify this condition, as the notion of time is *only* defined through the rate of gradient steps given by Assumption 3.1 (λ can be interpreted as "*the expected number of communications in the graph between the expected duration separating two subsequent gradient steps on a given node*"), it only means that the gradient steps are happening too fast compared to the ability of the network to communicate, and the rate of gradient steps should be slow-down by a factor of $\sqrt{2\chi_1\chi_2}$.

G DADAO vs MSDA on the star graph

To experimentally confirm that our communication complexity is better than accelerated methods using the spectral gap and abstract away the better absolute constants for MSDA, we ran DADAO and MSDA on the task of distributed linear regression for star graphs of size $n \in \{10, 20, 70, 200, 300, 1000, 2000\}$. We considered the evolution of the average distance to the optimal with the number of gradient steps and communication steps in log scale for each run, and computed the slope of each line. For each graph size, we report in Fig. 7 the

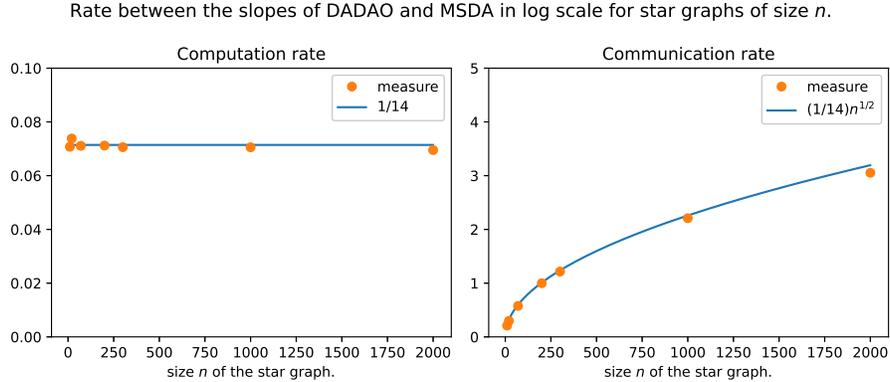


Figure 7: Rate between the slopes in log scale of DADAO and MSDA for star graphs of size $n \in \{10, 20, 70, 200, 300, 1000, 2000\}$.

rate between the slope for DADAO and the slope for MSDA. We remark that the rate between the gradient complexities of DADAO and MSDA is indeed a $\mathcal{O}(1)$ (with a constant value of $\simeq 1/14$) while MSDA is indeed $\mathcal{O}(\sqrt{n})$ worse than DADAO for communications on the star graph, as stated in Tab. 2.

H Practical Implementation

In this section, we describe in more detail the implementation of our algorithm. As we did not physically execute our method on a compute network but carried it out on a single machine, all the asynchronous computations and communications had to be simulated. Thus, we will first discuss the method we followed to simulate our asynchronous framework before detailing the practical steps of our algorithm through a pseudo-code.

H.1 Simulating the Poisson Point Processes

To emulate the asynchronous setting, before running our algorithm, we generate 2 independent sequences of jump times at the graph’s scale: one for the computations and one for the communications. As we considered independent P.P.Ps, the time increments follow a Poisson distribution. At the graph’s scale, each node spiking at a rate of 1, the Poisson parameter for the gradient steps process is n . Following the experimental setting of the Continuized framework [12], we considered that all edges in \mathcal{E} had the same probability of spiking. Thus, given any graph \mathcal{E} and \mathcal{L} its corresponding standard Laplacian with unit edge weights,

we computed the parameter λ of the communication process:

$$\lambda = \sqrt{2\chi_1 \left[\frac{\mathcal{L}}{|\mathcal{E}|} \right] \chi_2 \left[\frac{\mathcal{L}}{|\mathcal{E}|} \right]}. \quad (87)$$

Having generated the 2 sequences of spiking times at the graph's scale, we run our algorithm playing the events in order of appearance, attributing the *location* of the events by sampling uniformly one node if the event is a gradient step or sampling uniformly an edge in \mathcal{E} if it is a communication.

H.2 Pseudo Code

We keep the notations introduced in Eq. (3) and recall the following constant values specified in Appendix C.1:

$$\begin{aligned} \eta &= \frac{1}{8}\sqrt{\frac{\nu}{L}} & \gamma &= \frac{1}{4L} & \delta &= \frac{1}{4}\sqrt{\frac{\nu}{L}} & \alpha &= \frac{1}{4}\sqrt{\frac{\nu}{L}} & \beta &= \frac{1}{2} & \theta &= \frac{1}{2}\sqrt{\frac{L}{\nu}} \\ \tilde{\eta} &= \frac{1}{8}\sqrt{\frac{\nu}{L}} & \tilde{\gamma} &= \frac{1}{4\sqrt{\nu L}} & \tilde{\delta} &= 1 & \tilde{\alpha} &= \frac{1}{8}\sqrt{\frac{\nu}{L}} & \tilde{\beta} &= 2\chi_1[\Lambda]\sqrt{\frac{L}{\nu}} & \nu &= \frac{\mu}{2} \end{aligned}$$

For the sake of completeness, we also specify the matrix \mathcal{A} describing the linear ODE (7):

$$\mathcal{A} = \begin{pmatrix} -\eta & \eta & 0 & 0 & 0 & 0 \\ \tilde{\eta} & -\tilde{\eta} & 0 & 0 & 0 & 0 \\ 0 & 0 & -\alpha & \alpha & 0 & 0 \\ 0 & -\theta\nu & -\theta & 0 & -\theta & 0 \\ 0 & 0 & 0 & 0 & -\alpha & \alpha \\ 0 & 0 & 0 & 0 & \tilde{\alpha} & -\tilde{\alpha} \end{pmatrix} \quad (88)$$

As described in Appendix H.1, we call `PPPspikes` the process mentioned above, returning the ordered sequence of events and time of spikes of the two P.P.Ps. Then, we can write the pseudo-code of our implementation of the DADAO optimizer in Algorithm 2.

Algorithm 2: Pseudo-code of our implementation of DADAO on a single machine.

Input: On each machine $i \in \{1, \dots, n\}$, an oracle able to evaluate ∇f_i ,
Parameters $\mu, L, \chi_1, t_{\max}, n, \lambda$.
The graph \mathcal{E} .

- 1 **Initialize** on each machine $i \in \{1, \dots, n\}$:
- 2 Set $X^{(i)} = (x_i, \tilde{x}_i, \tilde{y}_i)$ and $Y^{(i)} = (y_i, z_i, \tilde{z}_i)$ to 0 ;
- 3 Set constants $\nu, \tilde{\eta}, \eta, \gamma, \alpha, \tilde{\alpha}, \theta, \delta, \tilde{\delta}, \beta, \tilde{\beta}$ using μ, L, χ_1 ;
- 4 Set \mathcal{A} ;
- 5 $T^{(i)} \leftarrow 0$;
- 6 **ListEvents**, **ListTimes** \leftarrow PPPspikes(n, λ, t_{\max}) ;
- 7 $n_{\text{events}} \leftarrow |\text{ListEvents}|$;
- 8 **for** $k \in \llbracket 1, n_{\text{events}} \rrbracket$ **do**
- 9 **if** **ListEvents**[k] *is to take a gradient step* **then**
- 10 $i \sim \mathcal{U}(\llbracket 1, n \rrbracket)$;
- 11 $\begin{pmatrix} X^{(i)} \\ Y^{(i)} \end{pmatrix} \leftarrow \exp((\text{ListTimes}[k] - T^{(i)})\mathcal{A}) \begin{pmatrix} X^{(i)} \\ Y^{(i)} \end{pmatrix}$;
- 12 $g_i \leftarrow (\nabla f_i(x_i) - \nu x_i - \tilde{y}_i)$; // Local gradient computation.
- 13 $x_i \leftarrow x_i - \gamma g_i$;
- 14 $\tilde{x}_i \leftarrow \tilde{x}_i - \tilde{\gamma} g_i$;
- 15 $\tilde{y}_i \leftarrow \tilde{y}_i + (\delta + \tilde{\delta})g_i$;
- 16 $T^{(i)} \leftarrow \text{ListTimes}[k]$;
- 17 **else if** **ListEvents**[k] *is to take a communication step* **then**
- 18 $(i, j) \sim \mathcal{U}(\mathcal{E})$;
- 19 $\begin{pmatrix} X^{(i)} \\ Y^{(i)} \end{pmatrix} \leftarrow \exp((\text{ListTimes}[k] - T^{(i)})\mathcal{A}) \begin{pmatrix} X^{(i)} \\ Y^{(i)} \end{pmatrix}$;
- 20 $\begin{pmatrix} X^{(j)} \\ Y^{(j)} \end{pmatrix} \leftarrow \exp((\text{ListTimes}[k] - T^{(j)})\mathcal{A}) \begin{pmatrix} X^{(j)} \\ Y^{(j)} \end{pmatrix}$;
- 21 $m_{ij} \leftarrow (y_i + z_i - y_j - z_j)$; // Message exchanged.
- 22 $z_i \leftarrow z_i - \beta m_{ij}$;
- 23 $\tilde{z}_i \leftarrow \tilde{z}_i - \tilde{\beta} m_{ij}$;
- 24 $z_j \leftarrow z_j + \beta m_{ij}$;
- 25 $\tilde{z}_j \leftarrow \tilde{z}_j + \tilde{\beta} m_{ij}$;
- 26 $T^{(i)} \leftarrow \text{ListTimes}[k]$;
- 27 $T^{(j)} \leftarrow \text{ListTimes}[k]$;
- 28 **return** $(x_i)_{1 \leq i \leq n}$, the estimate of x^* on each worker i .
