



**HAL**  
open science

# Block Delayed Majorize-Minimize Subspace Algorithm for Large Scale Image Restoration

Mathieu Chalvidal, Emilie Chouzenoux, Jean-Baptiste Fest, Claire Lefort

► **To cite this version:**

Mathieu Chalvidal, Emilie Chouzenoux, Jean-Baptiste Fest, Claire Lefort. Block Delayed Majorize-Minimize Subspace Algorithm for Large Scale Image Restoration. 2022. hal-04281658v2

**HAL Id: hal-04281658**

**<https://hal.science/hal-04281658v2>**

Preprint submitted on 7 Jan 2023 (v2), last revised 13 Nov 2023 (v4)

**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.

# Block Delayed Majorize-Minimize Subspace Algorithm for Large Scale Image Restoration‡

Mathieu Chalvidal<sup>(a)</sup>, Emilie Chouzenoux<sup>(b)</sup>, Jean-Baptiste  
Fest<sup>(b)§</sup>, and Claire Lefort<sup>(c)</sup>

(a) ANITI - CerCo UMR CNRS 5549, Toulouse France - Brown University, USA

(b) CVN, CentraleSupélec, Inria, Université Paris-Saclay, France.

(c) XLIM Research Institute, UMR CNRS 7252, Université de Limoges, France.

**Abstract.** In this work, we propose an asynchronous Majorization-Minimization (MM) algorithm for solving large scale differentiable non-convex optimization problems. The proposed algorithm runs efficient MM memory gradient updates on block of coordinates, in a parallel and possibly asynchronous manner. We establish the convergence of the resulting sequence of iterates under mild assumptions. The performance of the algorithm is illustrated on the restoration of 3D images degraded by depth-variant 3D blur, arising in multiphoton microscopy. Significant computational time reduction, scalability and robustness are observed on synthetic data, when compared to state-of-the-art methods. Experiments on the restoration of real acquisitions of a muscle structure illustrate the qualitative performance of our approach and its practical applicability.

**Keywords** Distributed optimization, asynchronous implementation, block alternating method, majorization-minimization, non-convex optimization, image restoration, depth-variant blur, multiphoton microscopy.

Submitted to: XXX

‡ This research work received funding support from the European Research Council Starting Grant MAJORIS ERC-2019-STG-850925.

§ Corresponding author. Contact: jean-baptiste.fest@centralesupelec.fr, 9 rue Joliot Curie, 91190 Gif-sur-Yvette, France.

## 1. Introduction

Large-scale optimization algorithms, benefiting from fast convergence, capable of utilizing modern computing infrastructures, and dealing with distributed datasets are becoming compulsory for solving inverse problems in modern imaging [23]. The ever-growing need for fast processing solutions that can operate on high-dimensional problems (i.e implying a huge number of variables) calls for the development of parallel methods harnessing the power of distributed computational architectures. In addition, the expansion of IoT systems and remote highly parallel computing induce new network issues with specific constraints. For instance, instabilities may occur whenever the volume of data dwarfs the memory capacity of a single agent or when the processing power is shared (potentially unevenly) between devices [46]. Several classes of so-called distributed optimization methods, have been investigated under various assumptions on the computing scenario and on the optimization problem itself, that we review hereafter (see also [77, 73]).

Distributed optimization approaches inherit from block alternating methods. In the latter, at each iteration, only a subset of the variables are updated, by minimizing the objective function with respect to only those variables, the others being fixed. The blocks are selected sequentially following a cyclic (or quasi-cyclic) order or a random rule. Exact minimization with respect to a given block of variables is rarely possible in a closed form. It is not even desirable as it may lead to convergence issues [68]. More efficient and stable block alternating schemes rely on a so-called Majorization-Minimization (MM) strategy [43]. It consists in building, at each iteration, a majorizing approximation for the objective function within the active block of variables, whose minimizer has a more tractable form. Many powerful algorithms fall within this framework, such as BSUM [40], PALM [8], multiplicative methods for NMF [47], to name a few. By exploiting the structure of the objective function, block alternating MM methods can reach fast convergence rates [28, 31, 54, 57] while offering theoretical guarantees in non-convex cases [8, 21, 9].

When the problem size increases, as in 3D microscopy imaging [48] and astronomy [58, 61], running block alternating methods gets inefficient. Parallel implementations have been devised, where the block updates are performed simultaneously, allowing to distribute computations on different nodes (or machines) [10, 65, 63]. Implementation on parallel architecture requires to pay attention to communication cost. The latter can be reduced by resorting to an asynchronous parallel implementation, yielding the so-called distributed optimization approach. Each computation node has now its own iteration loop, so local variables are updated without the need to wait for distant variables update. This however raises challenging questions, in terms of convergence analysis, as the communication delays may introduce instabilities. A plethora of recent works have focused on proposing distributed optimization algorithms with assessed convergence, based on stochastic proximal primal [37, 50, 52] or primal-dual [60, 39, 75, 14, 58, 1] techniques. Recent contributions in the field of federated learning are also highly related

[41, 51, 70]. However, as the aforementioned works rely on specific fixed-point analysis tools involving Fenchel-Rockafellar duality [4], the proposed algorithms are limited to convex (sometimes even strongly convex) optimization and often require specific probabilistic assumptions on the block update rule difficult to meet in practice. In the context of MM algorithms, although the need for distributed implementation strategies is crucial (see the discussion in [40] and the specific examples in [69, 30]), theoretical results regarding convergence guarantees of MM technique in a distributed context are rather scarce. Let us first mention the work of [24, 25], that proposes an asynchronous version of PALM, with proven convergence of the iterates in non-convex case, and good practical behaviour [67]. The convergence of distributed MM methods was also explored in the recent works [11, 49]. However, the analysis of [49] is limited to the convex case. In [11], the analysis covers non-convex terms in the objective function, but it only shows the convergence of the sequence of objective function values, and not the convergence of the iterates themselves (thus, the results is weaker than the one of [24]).

In this work, we aim at solving a smooth optimization problem of the form

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad f(\mathbf{x}), \quad (1)$$

where  $f : \mathbb{R}^N \mapsto \mathbb{R}$  is (Fréchet) differentiable but non necessarily convex. In the context of inverse problems in imaging,  $f$  typically reads as the sum of a data fidelity term (e.g., a least-squares term) measuring the discrepancy between an acquired, degraded (e.g., blurry, noisy) image, and its estimate (usually, through a linear observation operator), and a regularization term incorporating prior information on the sought solution [23, 5] (see also our Section 5). We introduce the block delayed MM memory gradient (BD3MG) algorithm for the resolution of Problem (1). BD3MG is a distributed MM algorithm designed for an efficient implementation on a multi-CPU computing architecture, such as a high performance calculation unit. Our contributions are<sup>†</sup>:

- Introduction of the BD3MG algorithm, that implements an advanced distributed asynchronous update rule within the block alternating MM method we recently proposed in [32].
- Proof for the convergence of BD3MG iterates to a stationary point of  $f$  under mild assumptions (in particular, no convexity is assumed), using recent tools of Lyapunov analysis [71].
- Illustration of the performance of BD3MG by means of various experiments on a real inverse problem of 3D image restoration arising in the context of multiphotonic microscopy.

The paper is organized as follows. Section 2 introduces our notations, recalls the principle of MM schemes and finally presents our proposed algorithm. Section 3 states

<sup>†</sup> A preliminary version of this work has been presented in the conference proceedings [13]. The convergence result was weaker, and stated without proof. The experimental validation was limited to a single, simpler, numerical scenario.

our mathematical assumptions for the convergence analysis and presents preliminary technical propositions and lemmas. Section 4 presents our main theoretical contribution, dedicated to the convergence analysis of the proposed BD3MG scheme. Section 5 illustrates the qualitative and computational performance of BD3MG in the applicative context of 3D image deblurring in the presence of a depth-variant 3D blur. Section 6 concludes the paper.

## 2. Proposed algorithm

### 2.1. Notations

Throughout the paper, we consider the euclidean space  $\mathbb{R}^N$  endowed with the usual scalar product  $\langle \cdot | \cdot \rangle$  (or, equivalently,  $\cdot^\top \cdot$ ) and the norm  $\| \cdot \|$ .  $\mathbf{0}_N$  is the vector with null entries of  $\mathbb{R}^N$ .  $\mathbf{I}_N$  is the identity matrix of  $\mathbb{R}^N$ . We use the short notation  $\llbracket 1, N \rrbracket$ , to denote  $\{1, 2, \dots, N\}$ , i.e. the set of integers from 1 to  $N$ .  $\mathbb{S}^N$  denotes the set of symmetric matrices of  $\mathbb{R}^{N \times N}$ , and  $\mathbb{S}_+^N$  (resp.  $\mathbb{S}_{++}^N$ ) the set of positive (resp definite positive) symmetric matrices. Given some  $\mathbf{M} \in \mathbb{S}_{++}^N$ , we denote by  $\| \cdot \|_{\mathbf{M}}$  the induced weighted euclidean norm, such that, for all  $\mathbf{v} \in \mathbb{R}^N$ ,  $\| \mathbf{v} \|_{\mathbf{M}}^2 = \mathbf{v}^\top \mathbf{M} \mathbf{v}$ . We use the Loewner orders symbols  $<$  and  $\leq$ , to compare real symmetric matrices  $(\mathbf{A}, \mathbf{B}) \in (\mathbb{S}^N)^2$  i.e.,  $\mathbf{A} \leq \mathbf{B}$  (resp.  $\mathbf{A} < \mathbf{B}$ ) is verified when difference  $\mathbf{B} - \mathbf{A}$  belongs to  $\mathbb{S}_+^N$  (resp.  $\mathbb{S}_{++}^N$ ).

Let us introduce extra notations, that will be useful to present block coordinate optimization strategy. Most notations hereafter are reminiscent from [32]. Let  $\mathcal{S} \subset \llbracket 1, N \rrbracket$ .

- ▷ We denote by  $\bar{\mathcal{S}}$  its complementary set  $\llbracket 1, N \rrbracket \setminus \mathcal{S}$ ,  $|\mathcal{S}|$  its cardinal and  $(\mathbb{R}^{|\mathcal{S}|}, \langle \cdot, \cdot \rangle)$  the resulting euclidean space (with a slight abuse of notation). Moreover, we also denote by  $\mathbb{S}^{|\mathcal{S}|}, \mathbb{S}_+^{|\mathcal{S}|}, \mathbb{S}_{++}^{|\mathcal{S}|}$  respectively the set of symmetric, symmetric positive, and symmetric definite positive matrices of  $\mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ .
- ▷ Let  $\mathbf{x} = (x_n)_{n \in \llbracket 1, N \rrbracket} \in \mathbb{R}^N$ . We denote  $\mathbf{x}_{(\mathcal{S})} = (x_i)_{i \in \mathcal{S}} \in \mathbb{R}^{|\mathcal{S}|}$  the vector gathering the entries of  $\mathbf{x}$  with indexes within the set  $\mathcal{S}$  of coordinates.
- ▷ Let  $\mathbf{x} \in \mathbb{R}^N$ .  $\nabla f(\mathbf{x})$  is the gradient of  $f$  evaluated at  $\mathbf{x}$ . Moreover,  $\nabla_{(\mathcal{S})} f(\mathbf{x}) = \left( [\nabla f(\mathbf{x})]_i \right)_{i \in \mathcal{S}} \in \mathbb{R}^{|\mathcal{S}|}$  denotes the partial gradient of  $f$  with respect to the coordinates with indexes in  $\mathcal{S}$ , evaluated at  $\mathbf{x}$ .
- ▷ Let  $\mathbf{M} \in \mathbb{S}^N$ . We denote the (symmetric) sub-matrix  $\mathbf{M}_{(\mathcal{S})} = (M_{i,j})_{(i,j) \in \mathcal{S}^2} \in \mathbb{S}^{|\mathcal{S}|}$ . If  $\mathbf{M}_{(\mathcal{S})} \in \mathbb{S}_+^{|\mathcal{S}|}$ , we define the induced weighted Euclidean norm as  $\| \cdot \|_{\mathbf{M}_{(\mathcal{S})}}$ .
- ▷ For any  $\mathbf{x} \in \mathbb{R}^N$ , we introduce the restriction of  $f$  to the block  $\mathcal{S}$  and vector  $\mathbf{x}$  as the function  $f_{(\mathcal{S})}(\cdot, \mathbf{x}) : \mathbf{v} \in \mathbb{R}^{|\mathcal{S}|} \mapsto f(\mathbf{u})$  where  $\mathbf{u}$  is related to  $(\mathbf{v}, \mathbf{x})$  through the relations  $\mathbf{u}_{(\mathcal{S})} = \mathbf{v}$  and  $\mathbf{u}_{(\bar{\mathcal{S}})} = \mathbf{x}_{(\bar{\mathcal{S}})}$ .

## 2.2. Block MM principle

MM approach to the resolution of Problem (1) is a generic iterative procedure where each iteration amounts to minimizing (exactly or not) a surrogate for  $f$  satisfying a majorizing property [66, 76, 43, 42]. The theoretical and practical properties of an MM algorithm greatly depend on (i) the family of considered surrogates, (ii) the procedure to minimize it. In this work, we focus on quadratic MM techniques, where  $f$  is such that it can be upper bounded by quadratic functions (typically,  $f$  is Lipschitz differentiable). In such context, the inner step of an MM algorithm amounts to minimizing a quadratic function on  $\mathbb{R}^N$  or, otherwise stating, to invert an  $N \times N$  system. In the large scale context, this is not desirable and various approaches have been proposed to cope with the curse of dimensionality in MM quadratic methods [16, 17, 20, 40, 10, 65, 32]. In particular, to limit the dependence of the MM algorithm on the dimension of the problem, block alternating approaches have been developed. In these schemes, at each iteration only a subset of the variables are updated [40], giving rise to so-called block MM algorithms, that we describe hereafter.

Define a partition  $\mathbb{T}$  of  $\llbracket 1, N \rrbracket$ . The block MM approach requires to build a majorizing surrogate for the restriction  $f_{(\mathcal{S})}(\cdot, \mathbf{x})$  for any  $\mathbf{x} \in \mathbb{R}^N$  and block index  $\mathcal{S} \in \mathbb{T}$ . Let us assume the existence of a mapping  $\mathbf{A} : \mathbf{x} \in \mathbb{R}^N \mapsto \mathbf{A}(\mathbf{x}) \in \mathbb{S}_{++}^N$  such that for all  $\mathcal{S} \in \mathbb{T}$ ,  $\mathbf{v} \in \mathbb{R}^{|\mathcal{S}|}$  and  $\mathbf{x} \in \mathbb{R}^N$

$$Q_{\mathcal{S}}(\mathbf{v}, \mathbf{x}) = f(\mathbf{x}) + \langle \nabla_{(\mathcal{S})} f(\mathbf{x}), \mathbf{v} - \mathbf{x}_{(\mathcal{S})} \rangle + \frac{1}{2} \|\mathbf{v} - \mathbf{x}_{(\mathcal{S})}\|_{\mathbf{A}_{(\mathcal{S})}(\mathbf{x})}^2, \quad (2)$$

fulfills the majorizing condition

$$Q_{\mathcal{S}}(\mathbf{v}, \mathbf{x}) \geq f_{(\mathcal{S})}(\mathbf{v}, \mathbf{x}). \quad (3)$$

Note that, by (2),

$$Q_{\mathcal{S}}(\mathbf{x}_{(\mathcal{S})}, \mathbf{x}) = f_{(\mathcal{S})}(\mathbf{x}_{(\mathcal{S})}, \mathbf{x}). \quad (4)$$

The existence of such mapping  $\mathbf{A}$  can be ensured under mild assumptions. For instance, it is satisfied as soon as  $f$  is Lipschitz differentiable. Moreover, [21, Remark 2.4] shows that, as soon as the above mapping holds for  $\mathcal{S} = \mathbb{R}^N$ , it stays valid for any block subset  $\mathcal{S} \subset \llbracket 1, N \rrbracket$ . Examples of constructions of majorant mappings have been extensively discussed in [66, 17, 65] for optimization problems arising in the fields of inverse problems, image processing and telecommunication.

Once the block majorant approximations (2) satisfying (3) are built, the block MM (B2M) algorithm reads [43] (also called BSUM in [40]):

$$(\forall k \in \mathbb{N}) \quad \begin{cases} \text{Choose } \mathcal{S}^k \in \mathbb{T}, \\ \mathbf{x}_{(\mathcal{S}^k)}^{k+1} \in \arg \min_{\mathbf{v} \in \mathbb{R}^{|\mathcal{S}^k|}} Q_{\mathcal{S}^k}(\mathbf{v}, \mathbf{x}^k), \\ \mathbf{x}_{(\mathcal{S}^k)}^{k+1} = \mathbf{x}_{(\mathcal{S}^k)}^k. \end{cases} \quad (\text{B2M})$$

Hereabove,  $(\mathcal{S}^k)_{k \in \mathbb{N}}$  is a sequence of subsets (i.e., blocks) chosen in the predefined partition  $\mathbb{T}$ . The most current strategy is to adopt a cyclic rule, where each element of  $\mathbb{T}$  is selected sequentially until the end of the partition list, and then the loop is repeated until convergence of the algorithm. A more flexible option is to adopt a so-called quasi-cyclic (or acyclic) rule where each  $\mathcal{S} \in \mathbb{T}$  must be updated at least once per  $K$  iterations period.

The interest of scheme (B2M) and more generally block coordinate methods notably lies in the large scale context involving a very huge  $N$ , for which dealing with all the coordinates of the current iterate may be too high time consuming and even infeasible due to memory limitations. However, block MM methods require a sequential update of the blocks and thus, by construction, might require many iterations to reach convergence. To limit this issue, (block) diagonal mappings have been considered for instance in [65, 10]. The underlying idea is to choose the mapping so that the inner minimization problem in (B2M) is separable, and thus can be performed in parallel over the entries of the selected block. This yields the so-called block parallel MM schemes that take advantage of recent technological advances in parallel computing on multicore architectures. In particular, these methods can tailor the number of available processors to the computational load. However, such block diagonal structure may be detrimental to the approximation quality of the surrogates, and thus reduce again the practical convergence rate. In the present work, we opt for not making any extra structural assumption on the majorant mapping, thanks to the introduction of two catalyzers into (B2M), namely (i) a subspace acceleration approach, (ii) a distributed asynchronous update strategy, that we describe hereafter.

### 2.3. Subspace acceleration

Our first catalyst is to introduce a subspace acceleration [66], in (B2M). This strategy has been initially introduced for full-batch MM algorithms (i.e., without any block coordinate strategy) in [16]. Convergence analysis can be found in [17, 15, 20, 19, 18] under various situations. We recently extended this strategy to cope with block coordinate updates with the form of (B2M) [32], leading to the B2MS (Block MM Subspace) scheme that we present hereafter.

Starting with the (B2M) iteration, the subspace acceleration subspace acceleration subspace acceleration consists in performing the minimization of the majorant function within the current block  $\mathcal{S}^k$  in a constrained vectorial subspace spanned by a small

number  $M_k \geq 1$  of search directions. This reads:

$$(\forall k \in \mathbb{N}) \quad \begin{cases} \text{Choose } \mathcal{S}^k \in \mathbb{T}, \\ \text{Choose } \mathbf{D}^k \in \mathbb{R}^{M_k \times |\mathcal{S}^k|}, \\ \mathbf{v}^k \in \arg \min_{\mathbf{v} \in \mathbb{R}^{M_k}} Q_{\mathcal{S}^k}(\mathbf{x}_{(\mathcal{S}^k)}^k + \mathbf{D}^k \mathbf{v}, \mathbf{x}^k), \\ \mathbf{x}_{(\mathcal{S}^k)}^{k+1} = \mathbf{x}_{(\mathcal{S}^k)}^k + \mathbf{D}^k \mathbf{v}^k, \\ \mathbf{x}_{(\overline{\mathcal{S}^k})}^{k+1} = \mathbf{x}_{(\overline{\mathcal{S}^k})}^k, \end{cases} \quad (\text{B2MS})$$

Hereabove, for every  $k \in \mathbb{N}$ ,  $\mathbf{D}^k \in \mathbb{R}^{M_k \times |\mathcal{S}^k|}$  is the so-called subspace matrix. It stacks, row-wise,  $M_k \geq 1$  vectors of dimension  $|\mathcal{S}^k|$ , spanning a vectorial subspace within which we seek for a minimizer of the majorant function  $Q_{\mathcal{S}^k}(\cdot, \mathbf{x}^k)$  (i.e., our next iterate). The advantage is to reduce again the dimensionality of the inner MM problems, without jeopardizing the convergence rate [19]. Several choices for the subspace matrix are discussed in [16, 20, 18]. Intensive comparisons in the fields of inverse problems, image processing and machine learning (e.g., [33, 15]), have shown the superiority of the so-called memory gradient subspace which seems to reach the best compromise between simplicity and efficiency. In the context of (B2MS), this amounts to defining, for every  $k \in \mathbb{N}$ , the memory gradient matrix  $\mathbf{D}^k = \left[ \nabla_{(\mathcal{S}^k)} f(\mathbf{x}^k), \mathbf{x}_{(\mathcal{S}^k)}^k - \mathbf{x}_{(\mathcal{S}^k)}^{k-1} \right]$  (with the convention  $\mathbf{x}_{-1} = \mathbf{0}_N$ ), so that  $M_k = 2$ . When combined with a block diagonal majorant mapping, (B2MS) becomes equivalent to the BP3MG method considered in [13] for 3D image deblurring. The convergence properties of (B2MS) have recently been studied in [32].

#### 2.4. Block Delayed Majorize-Minimize Memory Gradient (BD3MG)

The second catalyst we introduce is the main contribution of this paper, namely the introduction of a distributed asynchronous update rule within (B2MS). Our motivation is to make the algorithm well suited to an implementation on a multi-core / multi-processor architecture, while not being endangered by potential communication delays within the computing units. Let us consider a computing architecture with  $C$  units (or cores), each of them being able to communicate (i.e., send or receive) information to a master node. The architecture thus considered is forming a *star* graph as presented in Figure (1c). The two other graph topologies are discarded from this present study (see, for example, [1] for an efficient distributed method running on a generic hypergraph topology).

The proposed method BD3MG is presented in Algorithms 1-2, describing the iterations of the master (i.e., node 0) and a given worker/node  $c \in \llbracket 1, C \rrbracket$ , respectively. Let us describe these two algorithms. Each computation node  $c \in \llbracket 1, C \rrbracket$  updates (independently from the other nodes) a subset of coordinates  $\mathcal{S}_c \in \mathbb{T}$  (which can change over the process) by applying an MM iteration including a memory gradient acceleration and thus “books” its running block  $\mathcal{S}_c$  so that no other worker overwrites the associated



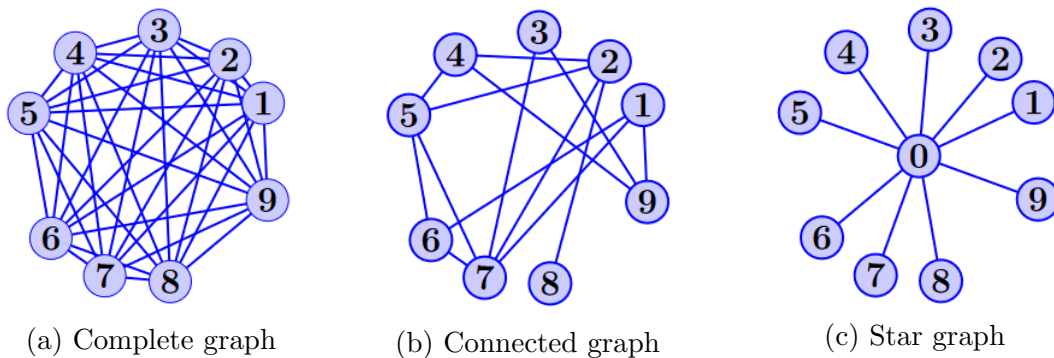


Figure 1: Examples of graph topologies. The graph in (c) is encompassed by our framework.

coordinates. Conversely, any other  $\mathcal{S} \in \mathbb{T} \setminus \{\mathcal{S}_c\}$  remains free to be updated by other workers. Communication steps are performed in order to maintain convergence to a minimizer of the globally shared objective function  $f$  and to control the propagation of errors. Basically, even if the other workers are still busy on their tasks, every time a worker  $c \in \llbracket 1, C \rrbracket$  ends one MM iteration on its running block  $\mathcal{S}_c$ , it sends a feedback to the master. As a response, the latter updates it with most recent available information, and assigns it a new task.

We denote  $(\mathbf{x}^k)_{k \in \mathbb{N}}$  the sequence of iterates gathered by the master loop. For any given node index  $c \in \llbracket 1, C \rrbracket$  and  $k \in \mathbb{N}$ ,  $\mathcal{S}_c^k$  denotes the block of coordinates processed by worker  $c$  during step  $k$ . We impose, by construction, that two nodes do not update the same block of coordinates at the same time, so that we ensure the no-overlap condition

$$(\forall k \in \mathbb{N}) \quad (\forall (c, c') \in \llbracket 1, C \rrbracket^2, c \neq c') \quad \mathcal{S}_c^k \cap \mathcal{S}_{c'}^k = \emptyset. \quad (5)$$

At iteration  $k \in \mathbb{N}$ , worker  $c^k \in \llbracket 1, C \rrbracket$ , updates the block of coordinates  $\mathcal{S}_{c^k}^k$  and sends to the master a vector  $\mathbf{d}_k$  of size  $|\mathcal{S}_{c^k}^k|$ . The corresponding indexes of variable  $\mathbf{x}^k$  within block  $\mathcal{S}_{c^k}^k$  are then incremented with  $\mathbf{d}_k$  while the others remain unchanged, thus defining  $\mathbf{x}^{k+1}$ . The master then defines a new set of coordinates  $\mathcal{S}_{c^k}^{k+1}$  to be treated by worker  $c^k$ , so as to satisfy the non-overlap rule. The master informs worker  $c^k$  of this new running set of coordinates, and sends him the most recent information  $\mathbf{x}^{k+1}$  and the difference  $(\mathbf{x}^{k+1} - \mathbf{x}^k)_{(\mathcal{S}_{c^k}^{k+1})}$ . Meanwhile, the other workers keep processing their allocated indexes. The master then waits until a new worker (possibly the same one) sends a new increment.

Let us now make a focus on the worker loop described in Algorithm 2. Remark that, even if worker  $c$  has access to some properties of function  $f$  (i.e., the expression for its gradient and for its majorizing approximation  $(Q_{\mathcal{S}})_{\mathcal{S} \in \mathbb{S}}$ ), it is not informed about the work done by the master or those of the other workers. It can only rely on the data the master sends to it to perform its local task. From the viewpoint of the worker, a triplet set  $(\mathbf{x}, \mathcal{S}, \mathbf{d}) \in \mathbb{R}^N \times \mathbb{T} \times \mathbb{R}^{|\mathcal{S}|}$  is received from the master and must be used to perform its MM update with memory gradient acceleration. The worker is in charge of

**Algorithm 1.** *BD3MG algorithm - Master loop**Initialization.*

- (a) Set  $k = 0$  and  $\mathbf{x}^0 \in \mathbb{R}^N$ .
- (b) Set  $\mathcal{S}_1^0, \dots, \mathcal{S}_C^0 \in \mathbb{T}$  such that  $\forall (c, c') \in \llbracket 1, C \rrbracket^2$ ,  $\mathcal{S}_c^0 \cap \mathcal{S}_{c'}^0 = \emptyset$ .
- (c) 0-th transmission: For every  $c \in \llbracket 1, C \rrbracket$ , **send**  $(\mathbf{x}^0, \mathcal{S}_c^0, 0_{|\mathcal{S}_c^0|})$  to worker  $c$

*While a stopping criterion is not met:*

(Wait for a feedback from a worker)

- (a)  $(k + 1)$ -th reception: **Receive**  $\mathbf{d}_k$  from a worker  $c^k$ .
- (b) Update  $\mathbf{x}_{\binom{\mathcal{S}_c^k}{c^k}}^{k+1} = \mathbf{x}_{\binom{\mathcal{S}_c^k}{c^k}}^k + \mathbf{d}_k$  and  $\mathbf{x}_{\binom{\mathcal{S}_c^k}{c^k}}^{k+1} = \mathbf{x}_{\binom{\mathcal{S}_c^k}{c^k}}^k$ .
- (c) Set  $\mathcal{S}_1^{k+1}, \dots, \mathcal{S}_C^{k+1} \in \mathbb{T}$  such that  $\mathcal{S}_{c^k}^{k+1} \in \mathbb{T} \setminus \{\mathcal{S}_c^k\}_{c \neq c^k}$  and,  $(\forall c \in \llbracket 1, C \rrbracket \setminus \{c^k\})$ ,  $\mathcal{S}_c^{k+1} = \mathcal{S}_c^k$ .
- (d)  $(k + 1)$ -th transmission: **Send**  $\left( \mathbf{x}^{k+1}, \mathcal{S}_{c^k}^{k+1}, (\mathbf{x}^{k+1} - \mathbf{x}^k)_{\binom{\mathcal{S}_{c^k}^{k+1}}{c^k}} \right)$  to worker  $c^k$ .
- (e)  $k = k + 1$

*End While**Output.* Vector  $\mathbf{x}^k$ .

first building the new memory gradient matrix

$$\mathbf{D}(\mathbf{x}, \mathcal{S}, \mathbf{d}) = [-\nabla_{(\mathcal{S})} f(\mathbf{x}) \mid \mathbf{d}] \in \mathbb{R}^{|\mathcal{S}| \times 2}. \quad (6)$$

and then compute the MM increment  $\mathbf{d}' \in \mathbb{R}^{|\mathcal{S}|}$  defined as

$$\mathbf{d}' = \mathbf{D}(\mathbf{x}, \mathcal{S}, \mathbf{d}) \mathbf{u} \quad (7)$$

$$\text{with } \mathbf{u} \in \arg \min_{\mathbf{v} \in \mathbb{R}^2} Q_{\mathcal{S}}(\mathbf{x} + \mathbf{D}(\mathbf{x}, \mathcal{S}, \mathbf{d}) \mathbf{v}, \mathbf{x}). \quad (8)$$

Note that the uniqueness of the solution for problem (7)-(8) is not guaranteed in general. To overcome such an obstacle, we follow the strategy in [16], and retain the pseudo-inverse solution given by

$$\mathbf{u} = - \left( \mathbf{D}(\mathbf{x}, \mathcal{S}, \mathbf{d})^\top \mathbf{A}_{(\mathcal{S})}(\mathbf{x}) \mathbf{D}(\mathbf{x}, \mathcal{S}, \mathbf{d}) \right)^\dagger \mathbf{D}(\mathbf{x}, \mathcal{S}, \mathbf{d}) \nabla_{(\mathcal{S})} f(\mathbf{x}), \quad (9)$$

where  $\dagger$  refers to the Moore-Penrose pseudo-inverse. Such solution notably verifies the normal equation

$$\langle \nabla_{(\mathcal{S})} f(\mathbf{x}), \mathbf{D}(\mathbf{x}, \mathcal{S}, \mathbf{d}) \mathbf{u} \rangle = -\|\mathbf{D}(\mathbf{x}, \mathcal{S}, \mathbf{d}) \mathbf{u}\|_{\mathbf{A}_{(\mathcal{S})}(\mathbf{x})}^2. \quad (10)$$

**Algorithm 2.** *BD3MG algorithm - Worker loop*

- (Wait for a feedback from the master)
- (a) *Receive*  $(\mathbf{x}, \mathcal{S}, \mathbf{d})$  from Master.
  - (b)  $\mathbf{D}(\mathbf{x}, \mathcal{S}, \mathbf{d}) = [-\nabla_{(\mathcal{S})}f(\mathbf{x}) \mid \mathbf{d}]$ .
  - (c) Compute  $\nabla_{(\mathcal{S})}f(\mathbf{x})$  and  $\mathbf{A}_{(\mathcal{S})}(\mathbf{x})$ .
  - (d)  $\mathbf{u} = -(\mathbf{D}(\mathbf{x}, \mathcal{S}, \mathbf{d})^\top \mathbf{A}_{(\mathcal{S})}(\mathbf{x}) \mathbf{D}(\mathbf{x}, \mathcal{S}, \mathbf{d}))^\dagger \mathbf{D}(\mathbf{x}, \mathcal{S}, \mathbf{d}) \nabla_{(\mathcal{S})}f(\mathbf{x})$ .
  - (e) *Send*  $\mathbf{d}' = \mathbf{D}(\mathbf{x}, \mathcal{S}, \mathbf{d})\mathbf{u}$  to the Master.

*2.5. Distributed structure of BD3MG*

We first have to make an important remark, regarding the communication load in terms of memory, in between the master and the workers. Consider a worker associated to the block index  $\mathcal{S}$ . According to Algorithm 2, the worker receives three quantities, namely  $\mathbf{x}$  of  $N$  real values, the set of integer indexes  $\mathcal{S}$  with cardinality  $|\mathcal{S}|$  and the vector  $\mathbf{d}$  of  $|\mathcal{S}|$  real values. The sent vector  $\mathbf{d}'$  is again made of  $|\mathcal{S}|$  real values. Clearly, the main memory load is related to the reception of vector  $\mathbf{x}$ . One should however notice that the worker only uses  $\mathbf{x}$  to compute  $\nabla_{(\mathcal{S})}f(\mathbf{x})$  and  $\mathbf{A}_{(\mathcal{S})}(\mathbf{x})$ . In most situations encountered in inverse problems of imaging,  $f$  shows some inherent separable structure, so that both of these quantities only depend on a subset of entries of vector  $\mathbf{x}$  that can be of small cardinality compared to  $N$ . The practical implementation of Algorithm 2 should account for this specific situation, in order to avoid memory saturation and important communication delays. We give a detailed analysis for this aspect, in the case of our experimental example, in the Section 5.1.4.

The proposed distributed structure of BDM3G follows a star graph. Practically, it means that one of the computing unit has a higher load, in terms of memory, since it must process the full vector  $\mathbf{x}$  of size  $N$ , while the memory load of the workers is limited, as we discussed hereabove. This can be viewed as a limitation for the proposed method. The extension of our analysis to the case of a hypergraph distributed framework would require to be more specific about the structure of function  $f$  (in the line of the study of [1]), which might reduce the versatility of the algorithm. Up to our knowledge, this analysis is not straightforward and is thus left as future work.

*2.6. Equivalent form for BD3MG*

The way we introduced our scheme BD3MG in the previous subsection was “implementation-oriented”. In order to study its convergence behaviour, we must exhibit an equivalent form of it, mimicking the one of its non distributed counterpart, (B2MS). To do so, it is necessary to formalize the information gap between the master and the workers during the iterative process.

As we have already mentioned, all the information available to a worker (except those on  $f$  and  $(Q_S)_{S \in \mathbb{T}}$ ) is sent to it by the master only after it produces a feedback. For a given  $k \in \mathbb{N}$ , worker  $c^k$  does not receive any information between the  $(k+1)$ -reception and the previous one it made. During this time, its counterparts  $c \in \llbracket 1, C \rrbracket \setminus \{c^k\}$  may have performed additional updates to the master without  $c^k$  being informed. This results in an information mismatch, that we propose to formalize through a vector  $\mathbf{x}^{\iota_k}$  where

$$(\forall k \in \mathbb{N}) \quad \iota_k = \begin{cases} 0 & \text{if } k = 0, \\ \max \left( \{ \ell \in \llbracket 1, k \rrbracket \mid c^{\ell-1} = c^k \} \cup \{0\} \right), & \text{otherwise.} \end{cases} \quad (11)$$

This vector corresponds to the iteration index of the working variable of worker  $c^k$ , which does not necessarily matches with the vector  $\mathbf{x}^k$  manipulated by the master.

Let us list herebelow some situations of interest given the value of  $\iota_k$  at some iteration  $k \in \mathbb{N}$ :

- If  $\iota_k = 0$ , and  $k > 0$ , it means that  $\{ \ell \in \llbracket 1, k \rrbracket \mid c^{\ell-1} = c^k \}$  is an empty set. Hence, the worker  $c^k$  never returned any feedback to the master before the iteration  $k$ . Note that  $\iota_0 = 0$  by construction.
- If  $\iota_k = k$ , we thus have  $c^{k-1} = c^k$ . Hence, worker  $c^k$  was in charge of the two most recent updates, namely the  $(k+1)$ -th and the  $k$ -th ones. As a consequence, to prepare the  $(k+1)$ -th update, worker  $c^k$  received vector  $\mathbf{x}^k$  from the master.
- More generally, if  $\iota_k > 0$ , it follows that worker  $c^k$  at least returned one feedback to the master before iteration  $k$ . And we have the relation  $c^{\iota_k-1} = c^k$ .

Moreover, the non-overlap rule translates into

$$(\forall k \in \mathbb{N}) \quad \mathbf{x}_{(S_{c^k}^{\iota_k})} = \mathbf{x}_{(S_{c^k}^k)}. \quad (12)$$

For instance, if  $\iota_k = k-1$  for some  $k > 1$ , this indicates that  $c^{k-2} = c^k$  and  $c^{k-1} \neq c^k$ . The worker  $c^k$  thus proceeded to the  $(k-1)$ -th and  $(k+1)$ -th reception of the master while the  $k$ -th was made by another  $\tilde{c}^k$  who received the vector  $\mathbf{x}^k$  (from the master). However, since worker  $c^k$  was still processing block  $S_{c^k}^{k-1}$ , the master was not able to update the associated coordinate for computing  $\mathbf{x}^k$  from  $\mathbf{x}^{k-1}$  for worker  $c^k$ , i.e  $\mathbf{x}_{(S_{c^k}^{\iota_k})}^{k-1} = \mathbf{x}_{(S_{c^k}^k)}$ , which is typical from an asynchronous scheme.

More generally, when it comes to dealing with asynchronous algorithms, the use of a specific indexes with similar roles than our  $\iota_k$  ( $k \in \mathbb{N}$ ) is often necessary to build a theoretical delay model and thus to formulate an equivalent scheme being more compact and easier to analyse [24].

With this aim in mind, let us introduce the shorter notations

$$(\forall k \in \mathbb{N}) \quad \begin{cases} \mathcal{B}^k = S_{c^k}^k, \\ \mathcal{D}^k = D \left( \mathbf{x}^{\iota_k}, \mathcal{B}^k, (\mathbf{x}^{\iota_k} - \mathbf{x}^{\iota_k-1})_{(\mathcal{B}^k)} \right), \end{cases} \quad (13)$$

and  $\mathbf{D}^k = \mathbf{D} \left( \mathbf{x}^{\iota_k}, \mathcal{B}^k, (\mathbf{x}^{\iota_k} - \mathbf{x}^{\iota_k-1})_{(\mathcal{B}^k)} \right)$  with convention  $\mathbf{x}^{-1} = \mathbf{0}_N$ . Then, the master/worker BD3MG loops from Algorithms 1-2 can be rewritten equivalently in a single compact scheme as:

$$(\forall k \in \mathbb{N}) \quad \begin{cases} \text{Let } c^k \in \llbracket 1, C \rrbracket, \\ \mathbf{u}^k = - \left( (\mathbf{D}^k)^\top \mathbf{A}_{(\mathcal{B}^k)} (\mathbf{x}^{\iota_k}) \mathbf{D}^k \right)^\dagger (\mathbf{D}^k)^\top \nabla_{(\mathcal{B}^k)} f(\mathbf{x}^{\iota_k}), \\ \mathbf{x}_{(\mathcal{B}^k)}^{k+1} = \mathbf{x}_{(\mathcal{B}^k)}^k + \mathbf{D}^k \mathbf{u}^k, \\ \mathbf{x}_{(\overline{\mathcal{B}^k})}^{k+1} = \mathbf{x}_{(\overline{\mathcal{B}^k})}^k, \end{cases} \quad (14)$$

where we noticed that (12) now reads (using (13))

$$(\forall k \in \mathbb{N}) \quad \mathbf{x}_{(\mathcal{B}^k)}^{\iota_k} = \mathbf{x}_{(\mathcal{B}^k)}^k. \quad (15)$$

For every  $k \in \mathbb{N}$ , according to (14),  $\mathbf{u}_k$  still reads (9) and thus verifies (10) with  $\mathbf{D}^k = \mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k$ . The optimality equation can be rewritten as:

$$\left( \nabla_{(\mathcal{B}^k)} f(\mathbf{x}^{\iota_k}) \right)^\top \left( \mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k \right) = - \left\| \mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k \right\|_{\mathbf{A}_{(\mathcal{B}^k)}(\mathbf{x}^{\iota_k})}^2. \quad (16)$$

The two next Sections are dedicated to establish the convergence of the iterates produced by (14).

### 2.7. Link with existing works

Let us discuss the links between our proposed scheme BD3MG and existing methods from the literature. When  $\iota_k = k$  for any  $k \in \mathbb{N}$  in BD3MG, the algorithm identifies with our block alternating scheme B2MS [32] where the blocks of variables were updated sequentially in a non parallel (thus, not asynchronous) manner. This present paper can thus be viewed as an extension of the framework and of the convergence analysis of [32] to the distributed setting. Other related methods are [24, 11, 49], and our convergence analysis relies on similar tools than the one from [24]. Assuming zero-valued non-smooth terms in [24, 11, 49] (i.e., the objective function is differentiable), these methods identify with particular instances of BD3MG that (i) would not implement any subspace acceleration (i.e.,  $\mathbf{D}_k = \mathbf{I}_N$  in (14)), (ii) would rely on the simple Lipschitz-based majorant metric (i.e.,  $\mathbf{A}_{(\mathcal{B}^k)}(\mathbf{x}^{\iota_k}) = \mathcal{L}\mathbf{I}_{|\mathcal{B}^k|}$  in (14)) in the case of [24]. As a consequence, assuming differentiability of all terms, our convergence analysis presented in the next section thus also covers the schemes of [24, 11, 49]. Up to our knowledge, our work is the first to show convergence of the iterates of a distributed MM algorithm involving generic quadratic surrogates and subspace acceleration, in the non-convex setting. Finally, we would like to point out that the 3MG update performed in Alg. 2 identifies with a non-linear conjugate gradient (NLCG) update, for a specific (and closed form) pair of stepsize and conjugacy parameters (see discussion in [18, Sec. 1]). Therefore, our work can also be understood as the first convergence analysis of a distributed NLCG method in the

non-convex setting. A comparative study will be conducted in our experimental section to illustrate the superiority of BD3MG with respect to the aforementioned existing methods in terms of convergence speed.

### 3. Assumptions and preliminary results

#### 3.1. Assumptions

In order to analyse the asymptotic behaviour of the sequence  $(\mathbf{x}^k)_{k \in \mathbb{N}}$  generated by scheme (14), we introduce technical assumptions both on function  $f$  and on the parameters of BD3MG method.

**Assumption 1.** *Function  $f$  is coercive, continuously differentiable on  $\mathbb{R}^N$ , and has a  $\mathcal{L}$ -Lipschitzian gradient with  $\mathcal{L} > 0$ , i.e.*

$$(\forall (\mathbf{x}, \mathbf{y}) \in (\mathbb{R}^N)^2) \quad \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq \mathcal{L} \|\mathbf{x} - \mathbf{y}\|. \quad (17)$$

Assumption 1 ensures the existence of solutions for Problem (1) (by coercivity). Moreover, (17) in Assumption 1 guarantees the existence of a quadratic function (2) satisfying (3), setting  $\mathbf{A} : \mathbf{x} \mapsto \mathcal{L}\mathbf{I}_N$ . Another direct consequence is

$$(\forall \mathcal{S} \in \mathbb{T})(\forall (\mathbf{x}, \mathbf{y}) \in (\mathbb{R}^N)^2) \quad \|\nabla_{(\mathcal{S})} f(\mathbf{x}) - \nabla_{(\mathcal{S})} f(\mathbf{y})\| \leq \mathcal{L} \|\mathbf{x} - \mathbf{y}\|, \quad (18)$$

since  $\|\nabla_{(\mathcal{S})} f(\mathbf{x}) - \nabla_{(\mathcal{S})} f(\mathbf{y})\| \leq \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|$  for all  $\mathcal{S} \in \mathbb{T}$  and  $(\mathbf{x}, \mathbf{y}) \in (\mathbb{R}^N)^2$ .

**Assumption 2.** *(Boundedness of delay) For every  $k \in \mathbb{N}$ , and every worker  $c^k \in \llbracket 1, N \rrbracket$ , the set  $\mathcal{S}_c^k$  is not empty and there exists  $\tau \in \mathbb{N}^*$  such that*

$$(\forall k \geq \tau) \quad \llbracket 1, N \rrbracket = \bigcup_{j=k-\tau}^{k-1} \mathcal{B}^j, \quad (19)$$

using the notation in (13).

Assumption 2 gives an upper bound on the delay  $\tau$ . Each of the block of variables should be updated within a time frame of at most  $\tau$  iterations and thus the workers must follow a certain regularity. Such a condition follows a similar goal than quasi-cyclic rule frequently assumed in block coordinate methods [32, 40]. From Assumption 2, we deduce the following Proposition, which appears fundamental for the rest of our convergence study. It guarantees that, for a given  $k \in \mathbb{N}$ , the vector treated by worker  $c^k$  before its feedback (i.e the  $(k+1)$ -th master's reception) is not "too old" compared to the iteration index.

**Proposition 3.1.** *Under Assumption 2, for every  $k \geq \tau$ , the index  $\iota_k$  given in (11) belongs to  $\llbracket k - \tau + 1, k \rrbracket$ .*

*Proof.* Let  $k \geq \tau$ , where  $\tau > 0$  is defined in Assumption 2. Inequality  $\iota_k \leq k$  directly comes from Definition (11). We prove the lower bound on  $\iota_k$  by contradiction. Let us suppose that  $\iota_k \leq k - \tau$ . Two situations may arise.

*Case 1:*  $\iota_k = 0$ . By definition,  $c^0, \dots, c^{k-1} \neq c^k$  and an easy induction gives  $\mathcal{S}_{c^k}^0 = \dots = \mathcal{S}_{c^k}^k$ . Non-overlap rule (5) with  $c^0, \dots, c^{k-1} \neq c^k$  yields

$$(\forall j \in \llbracket 0, k-1 \rrbracket) \quad \mathcal{S}_{c^k}^j \cap \mathcal{S}_{c^j}^j = \mathcal{S}_{c^k}^k \cap \mathcal{B}^j \quad (20)$$

$$= \emptyset. \quad (21)$$

Since  $\mathcal{S}_{c^k}^k$  is non empty by Assumption 2, condition (20) ensures the existence of some  $i_k \in \llbracket 1, N \rrbracket$  verifying  $i_k \notin \bigcup_{j=0}^{k-1} \mathcal{B}^j$  contradicting  $\bigcup_{j=k-\tau}^{k-1} \mathcal{B}^j = \llbracket 1, N \rrbracket$ , as  $k \geq \tau$ .

*Case 2:*  $\iota_k > 0$ . We have  $c^{\iota_k-1} = c^k$  and a finite induction leads to

$$(\forall j \in \llbracket \iota_k, k \rrbracket) \quad \mathcal{S}_{c^k}^{\iota_k} = \mathcal{S}_{c^{\iota_k-1}}^{\iota_k} = \mathcal{S}_{c^{\iota_k-1}}^j = \mathcal{S}_{c^k}^j. \quad (22)$$

Majoration  $\iota_k \leq k - \tau$  implies that

$$(\forall j \in \llbracket k - \tau, k \rrbracket) \quad \mathcal{S}_{c^k}^{\iota_k} = \mathcal{S}_{c^k}^j. \quad (23)$$

Non-overlap rule (5) with  $c^{k-\tau}, \dots, c^{k-1} \neq c^k$  then gives

$$(\forall j \in \llbracket k - \tau, k - 1 \rrbracket) \quad \mathcal{S}_{c^k}^j \cap \mathcal{S}_{c^j}^j = \mathcal{S}_{c^k}^{\iota_k} \cap \mathcal{B}^j \quad (24)$$

$$= \emptyset. \quad (25)$$

Since  $\mathcal{S}_{c^k}^{\iota_k}$  is non empty, Condition (24) thus ensures the existence of  $i_k \in \llbracket 1, N \rrbracket$  verifying  $i_k \notin \bigcup_{j=k-\tau}^{k-1} \mathcal{B}^j$  which contradicts  $\bigcup_{j=k-\tau}^{k-1} \mathcal{B}^j = \llbracket 1, N \rrbracket$ . □

**Assumption 3.** (*Curvature of majorizing matrix*)

(i) The mapping  $\mathbf{A} : \mathbf{x} \in \mathbb{R}^N \mapsto \mathbf{A}(\mathbf{x}) \in \mathbb{S}_{++}^N$  is such that (3) holds. Moreover, there exists  $\bar{\nu} > 0$  such that, for all  $\mathcal{S} \in \mathbb{T}$  and  $\mathbf{v} \in \mathbb{R}^{|\mathcal{S}|}$ ,

$$0 < \mathbf{A}_{(\mathcal{S})}(\mathbf{v}) \leq \bar{\nu} \mathbf{I}_{|\mathcal{S}|}. \quad (26)$$

(ii) There exists  $\underline{\nu} > 0$  such that, for all  $k \in \mathbb{N}$ ,

$$\mathbf{\Gamma}_c^k = \mathbf{A}_{(\mathcal{B}^k)}(\mathbf{x}^{\iota_k}) - \frac{1}{2} \mathbf{A}_{(\mathcal{B}^k)}(\mathbf{x}^k) \geq \left( \frac{\mathcal{L} \sqrt{\tau} (1 + \tau)}{2} + \underline{\nu} \right) \mathbf{I}_{|\mathcal{B}^k|}. \quad (27)$$

Assumption 3(i) is standard in optimization literature dealing with MM methods involving quadratic surrogates [16]. Assumption 3(ii) assumes that the spectrum of the difference between delayed and exact majorizing matrices of the partial quadratic majoring functions is strictly greater than a certain constant. This hypothesis controls the length of the MM increments performed by each worker. It aims at ensuring consistency between the asynchronous updates, by directly relating the worst-case curvature of the function  $f$  (parameterized by the Lipschitz constant  $\mathcal{L}$ ) and the worst-case communication delay (parameterized by the constant  $\tau$ ). Condition (27) is key to ensure a condition descent for the general process generated by BD3MG scheme (see subsection 4.1). Assumption 3(ii) becomes redundant with Assumption 3(i) in the case when no delay occurs (i.e.,  $\tau = 0$ ). A detailed constructive example on how to meet Assumption 3(ii) will be provided in our experimental Section 5.

## 3.2. Technical lemmas

We conclude this section by presenting some preliminary results that will be useful for our convergence analysis.

**Lemma 3.1.** *Under Assumption 2, for every  $k \geq \tau$ ,*

$$\|\mathbf{x}^k - \mathbf{x}^{\iota_k}\|^2 \leq \tau \sum_{j=k-\tau+1}^k \|\mathbf{x}^j - \mathbf{x}^{j-1}\|^2. \quad (28)$$

*Proof.* Let  $k \in \mathbb{N}$ . If  $\iota_k = k$ , inequality (28) is trivial. For the rest of the proof we thus suppose  $\iota_k \leq k - 1$ . According to the definition of the euclidean norm we have

$$\|\mathbf{x}^k - \mathbf{x}^{\iota_k}\|^2 = \sum_{n=1}^N (x_n^k - x_n^{\iota_k})^2. \quad (29)$$

Then, for all  $n \in \llbracket 1, N \rrbracket$ , the Jensen's inequality leads to

$$(x_n^k - x_n^{\iota_k})^2 = \left( \sum_{j=\iota_k+1}^k (x_n^j - x_n^{j-1}) \right)^2 \leq (k - \iota_k) \sum_{j=\iota_k+1}^k (x_n^j - x_n^{j-1})^2. \quad (30)$$

Moreover, Proposition 3.1 ensures that  $\iota_k$  belongs to  $\llbracket k - \tau + 1, k \rrbracket$ . As a consequence

$$(\forall n \in \llbracket 1, N \rrbracket) \quad (x_n^k - x_n^{\iota_k})^2 \leq \tau \sum_{j=k-\tau+1}^k (x_n^j - x_n^{j-1})^2. \quad (31)$$

We then replace (31) in (29), which yields

$$\|\mathbf{x}^k - \mathbf{x}^{\iota_k}\|^2 \leq \tau \sum_{j=k-\tau+1}^k \sum_{n=1}^N (x_n^j - x_n^{j-1})^2. \quad (32)$$

Relation (28) directly comes from the identification of the inner sum of (32) as  $\|\mathbf{x}^j - \mathbf{x}^{j-1}\|^2$  for all  $j \in \llbracket k - \tau + 1, k \rrbracket$ .  $\square$

Lemma 3.1 provides a bound on the residual between  $\mathbf{x}^k$  and the delayed vector  $\mathbf{x}^{\iota_k}$  updated by worker  $c^k$  at iteration  $k \in \mathbb{N}$ . The right term in (28) can be understood as the extra information available to the master, when compared to the one available to worker  $c^k$ . This Lemma will allow to establish a descent condition on the BD3MG process in the next Section.

**Lemma 3.2.** *Under Assumptions 1 and 3(i), for every  $k \in \mathbb{N}$ ,*

$$\|\nabla_{(\mathcal{B}^k)} f(\mathbf{x}^{\iota_k})\|^2 \leq \bar{\nu}^2 \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2. \quad (33)$$

*Proof.* Let  $k \in \mathbb{N}$ . Let us analyse the quantity  $f(\mathbf{x}^{\iota_k}) - Q_{\mathcal{B}^k}(\mathbf{x}_{(\mathcal{B}^k)}^{k+1}, \mathbf{x}^{\iota_k})$ .



On the one hand, function  $\Psi_k : \alpha \in \mathbb{R} \mapsto Q_{\mathcal{B}^k}(\mathbf{x}_{(\mathcal{B}^k)}^k - \mathbf{D}^k \alpha \mathbf{e}, \mathbf{x}^{\iota_k})$  with  $\mathbf{e} = (1, 0)^\top$  is a second degree convex polynomial with a unique minimizer that reads

$$\hat{\alpha}_k = \frac{\|\nabla_{(\mathcal{B}^k)} f(\mathbf{x}^{\iota_k})\|^2}{\|\nabla_{(\mathcal{B}^k)} f(\mathbf{x}^{\iota_k})\|_{\mathbf{A}_{(\mathcal{B}^k)}(\mathbf{x}^{\iota_k})}^2}. \quad (34)$$

Since  $\mathbf{u}^k$  is a minimizer of  $Q_{\mathcal{B}^k}(\mathbf{x}_{(\mathcal{B}^k)}^{\iota_k} + \mathbf{D}^k \cdot, \mathbf{x}^{\iota_k}) = Q_{\mathcal{B}^k}(\mathbf{x}_{(\mathcal{B}^k)}^k + \mathbf{D}^k \cdot, \mathbf{x}^{\iota_k})$ , with  $\mathbf{x}_{(\mathcal{B}^k)}^{\iota_k} = \mathbf{x}_{(\mathcal{B}^k)}^k$  by Equation (11), we deduce that

$$Q_{\mathcal{B}^k}(\mathbf{x}_{(\mathcal{B}^k)}^{k+1}, \mathbf{x}^{\iota_k}) \leq \Psi_k(\hat{\alpha}_k) = f(\mathbf{x}^{\iota_k}) - \frac{1}{2} \hat{\alpha}_k \|\nabla_{(\mathcal{B}^k)} f(\mathbf{x}^{\iota_k})\|^2. \quad (35)$$

From Assumption 3(i),  $\hat{\alpha}_k$  verifies  $\hat{\alpha}_k \geq \bar{\nu}^{-1}$ . Equation (35) can thus be rewritten as

$$f(\mathbf{x}^{\iota_k}) - Q_{\mathcal{B}^k}(\mathbf{x}_{(\mathcal{B}^k)}^{k+1}, \mathbf{x}^{\iota_k}) \geq \frac{1}{2\bar{\nu}} \|\nabla_{(\mathcal{B}^k)} f(\mathbf{x}^{\iota_k})\|^2. \quad (36)$$

On the other hand, using (15) from Definition (2), and Equation (16) yield

$$\begin{aligned} & f(\mathbf{x}^{\iota_k}) - Q_{\mathcal{B}^k}(\mathbf{x}_{(\mathcal{B}^k)}^{k+1}, \mathbf{x}^{\iota_k}) \\ &= \left\langle \nabla_{(\mathcal{B}^k)} f(\mathbf{x}^{\iota_k}), \mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^{\iota_k} \right\rangle + \frac{1}{2} \left\| \mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^{\iota_k} \right\|_{\mathbf{A}_{(\mathcal{B}^k)}(\mathbf{x}^{\iota_k})}^2 \\ &= \left\langle \nabla_{(\mathcal{B}^k)} f(\mathbf{x}^{\iota_k}), \mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k \right\rangle + \frac{1}{2} \left\| \mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k \right\|_{\mathbf{A}_{(\mathcal{B}^k)}(\mathbf{x}^k)}^2 \\ &= \frac{1}{2} \left\| \mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k \right\|_{\mathbf{A}_{(\mathcal{B}^k)}(\mathbf{x}^{\iota_k})}^2. \end{aligned} \quad (37)$$

The combination of (36) and (37) leads to

$$\|\nabla_{(\mathcal{B}^k)} f(\mathbf{x}^{\iota_k})\|^2 \leq \bar{\nu} \left\| \mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k \right\|_{\mathbf{A}_{(\mathcal{B}^k)}(\mathbf{x}^{\iota_k})}^2. \quad (38)$$

Finally, Equation (33) comes using Assumption 3(i), and in particular,

$$\left\| \mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k \right\|_{\mathbf{A}_{(\mathcal{B}^k)}(\mathbf{x}^{\iota_k})}^2 \leq \bar{\nu} \left\| \mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k \right\|^2 \quad (39)$$

$$= \bar{\nu} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2. \quad (40)$$

□

Lemma 3.2 generalizes the decreasing behavior of standard MM schemes [17, 32] to the asynchronous context. It is not directly invoked in our main convergence proof but serves as an intermediary to show the following technical result.

**Lemma 3.3.** *Under Assumptions 1 and 3(i), for all  $k \geq 2\tau$ ,*

$$\|\nabla f(\mathbf{x}^k)\| \leq \mathcal{L}\tau \sum_{j=k-2\tau+1}^k \|\mathbf{x}^j - \mathbf{x}^{j-1}\| + \bar{\nu} \sum_{j=k-\tau}^k \|\mathbf{x}^{j+1} - \mathbf{x}^j\|. \quad (41)$$

*Proof.* Let  $k \geq 2\tau$ . Assumption 2 allows us to bound the gradient of  $f$  at  $\mathbf{x}^k$ , as

$$\|\nabla f(\mathbf{x}^k)\|^2 \leq \sum_{\ell=k-\tau}^{k-1} \|\nabla_{(\mathcal{B}^\ell)} f(\mathbf{x}^k)\|^2 \leq \left( \sum_{\ell=k-\tau}^{k-1} \|\nabla_{(\mathcal{B}^\ell)} f(\mathbf{x}^k)\| \right)^2. \quad (42)$$

Let us extract the root of the above terms, and use triangular and gradient-Lipschitz inequalities, leading to

$$\begin{aligned} \|\nabla f(\mathbf{x}^k)\| &\leq \sum_{\ell=k-\tau}^{k-1} \|\nabla_{(\mathcal{B}^\ell)} f(\mathbf{x}^k) - \nabla_{(\mathcal{B}^\ell)} f(\mathbf{x}^{\ell})\| + \sum_{\ell=k-\tau}^{k-1} \|\nabla_{(\mathcal{B}^\ell)} f(\mathbf{x}^{\ell})\| \\ &\leq \mathcal{L} \sum_{\ell=k-\tau}^{k-1} \|\mathbf{x}^k - \mathbf{x}^{\ell}\| + \sum_{j=k-\tau}^{k-1} \|\nabla_{(\mathcal{B}^j)} f(\mathbf{x}^{j})\|. \end{aligned} \quad (43)$$

For all  $\ell \in \llbracket k - \tau, k - 1 \rrbracket$ , by Proposition 4.1,  $\ell \geq \ell - \tau + 1 \geq k - 2\tau$ . Thus,

$$\|\mathbf{x}^k - \mathbf{x}^{\ell}\| \leq \sum_{j=k-2\tau+1}^k \|\mathbf{x}^j - \mathbf{x}^{j-1}\|. \quad (44)$$

The right term of (44) does not depend on index  $\ell$ . Using (44) and inequality (33) finally proves the result.  $\square$

Lemma 3.3 is useful as it provides a bound on the gradient at step  $\mathbf{x}^k$  only depending on the  $2\tau + 1$  past iterates  $\mathbf{x}^k, \dots, \mathbf{x}^{k-2\tau}$ .

**Lemma 3.4.** *Let  $(u^k)_{k \in \mathbb{N}}$  and  $(v^k)_{k \in \mathbb{N}}$  be two sequences of positive real numbers. If there exists  $P \in \mathbb{N}$  and  $k^* \geq P$  such that*

$$(\forall k \geq k^*) \quad u^k \leq r \sum_{i=k-P}^{k-1} u^i + v^{k-1}, \quad (45)$$

with  $r < 1/P$  and  $\sum_{k=0}^{+\infty} v^k < +\infty$ , then  $\sum_{k=0}^{+\infty} u^k < +\infty$ .

*Proof.* Summing (45) from  $k^*$  to  $n \geq k^*$  leads to

$$\sum_{k=k^*}^n u^k \leq r \sum_{k=k^*}^n \sum_{j=k-P}^{k-1} u^j + \sum_{k=k^*}^n v^{k-1}, \quad (46)$$

with

$$\sum_{k=k^*}^n \sum_{j=k-P}^{k-1} u^j = \sum_{k=k^*}^n \sum_{j=1}^P u^{k-j} = \sum_{j=1}^P \sum_{k=k^*-j}^{n-j} u^k \leq \sum_{j=1}^P \sum_{k=0}^n u^k. \quad (47)$$

Plugging (47) into (46), yields

$$\sum_{k=k^*}^n u^k \leq rP \sum_{k=k^*}^n u^k + \left( rP \sum_{k=0}^{k^*-1} u^k + \sum_{k=k^*}^n v^{k-1} \right) \leq rP \sum_{k=k^*}^n u^k + \left( rP \sum_{k=0}^{k^*-1} u^k + \sum_{k=0}^{+\infty} v^k \right), \quad (48)$$

that is  $(1 - rP) \sum_{k=k^*}^n u^k \leq rP \sum_{k=0}^{k^*-1} u^k + \sum_{k=0}^{+\infty} v^k$ . With  $0 < 1 - rP < 1$ , we deduce the summability of  $(u^k)_{k \in \mathbb{N}}$ .  $\square$

Lemma 3.4 is a technical result to ensure the convergence of some real series. Several variants of inequality (45) have been used to prove the finite length of iterative processes and then their convergence [24, 8].

#### 4. Convergence results

Let us now state our main theoretical results, that relate to the convergence properties of BD3MG iterates. Our proof line is reminiscent of [24, 8] and follows similar steps that we summarize hereafter. First, starting from the majoration property (3) and using Lemma 3.1, we will establish a descent inequality. The latter is the key point of the rest of the proof. In particular, it will allow to show convergence of  $(f(\mathbf{x}_k))_{k \in \mathbb{N}}$ . Then, Lemma 3.3 will ensure that  $(\nabla f(\mathbf{x}_k))_{k \in \mathbb{N}}$  converges to  $\mathbf{0}_N$ , and usual topological properties will serve to show that the set of cluster points  $\mathcal{C}(\mathbf{x}^k)_{k \in \mathbb{N}}$  of  $(\mathbf{x}^k)_{k \in \mathbb{N}}$  lies in the set of stationary point of  $f$ . Finally, we will exhibit a Lyapunov function from our descent inequality and will resort to the Kurdyka-Łojasewicz (KL) inequality [8] to prove our main theorem, showing the convergence of the BDM3G iterates and providing a rate of convergence.

##### 4.1. Descent inequality

**Proposition 4.1.** *Under Assumptions 1-2-3, there exists a positive sequence  $(\nu_k)_{k \geq \tau}$  such that*

$$(\forall k \geq \tau) \quad f(\mathbf{x}^{k+1}) + \nu_{k+1} \leq f(\mathbf{x}^k) + \nu_k - \underline{\nu} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2. \quad (49)$$

*Proof.* By definition of the majorant function (3), for every  $k \in \mathbb{N}$ ,

$$f(\mathbf{x}^{k+1}) \leq f(\mathbf{x}^k) + \left\langle \nabla_{(\mathcal{B}^k)} f(\mathbf{x}^k), \mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k \right\rangle + \frac{1}{2} \|\mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k\|_{\mathbf{A}_{(\mathcal{B}^k)}(\mathbf{x}^k)}^2. \quad (50)$$

Decomposing the scalar product term then yields, for every  $k \in \mathbb{N}$ ,

$$f(\mathbf{x}^{k+1}) \leq f(\mathbf{x}^k) + R_k + \left\langle \nabla_{(\mathcal{B}^k)} f(\mathbf{x}^k), \mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k \right\rangle + \frac{1}{2} \|\mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k\|_{\mathbf{A}_{(\mathcal{B}^k)}(\mathbf{x}^k)}^2, \quad (51)$$

$$\text{with } R_k = \left\langle \nabla_{(\mathcal{B}^k)} f(\mathbf{x}^k) - \nabla_{(\mathcal{B}^k)} f(\mathbf{x}^k), \mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k \right\rangle.$$

Let  $\tau$  defined as in Assumption 2. A majoration of  $R_k$  for every  $k \geq \tau$  comes by using successively Cauchy-Schwartz inequality,  $\mathcal{L}$  gradient-Lipschitz inequality from Assumption 1, and Lemma 3.1:

$$\begin{aligned}
(\forall k \geq \tau) \quad R_k &\leq \mathcal{L} \|\mathbf{x}^k - \mathbf{x}^{\iota_k}\| \|\mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k\| \\
&\leq \frac{\mathcal{L}}{2\sqrt{\tau}} \|\mathbf{x}^k - \mathbf{x}^{\iota_k}\|^2 + \frac{\mathcal{L}\sqrt{\tau}}{2} \|\mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k\|^2, \\
&\leq \frac{\mathcal{L}\sqrt{\tau}}{2} \sum_{j=k-\tau+1}^k \|\mathbf{x}^j - \mathbf{x}^{j-1}\|^2 + \frac{\mathcal{L}\sqrt{\tau}}{2} \|\mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k\|^2. \tag{52}
\end{aligned}$$

We then set, for all  $k \geq \tau$ ,  $\nu_k = \frac{\mathcal{L}\sqrt{\tau}}{2} \sum_{j=k-\tau+1}^k (j - k + \tau) \|\mathbf{x}^j - \mathbf{x}^{j-1}\|^2$ . Since  $\mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k = \mathbf{x}^{k+1} - \mathbf{x}^k$ , (52) also reads

$$\begin{aligned}
(\forall k \geq \tau) \quad R_k &\leq \nu_k - \nu_{k+1} + \frac{\mathcal{L}\tau\sqrt{\tau}}{2} \|\mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k\|^2 + \frac{\mathcal{L}\sqrt{\tau}}{2} \|\mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k\|^2, \\
&= \nu_k - \nu_{k+1} + \frac{\mathcal{L}\sqrt{\tau}(1 + \tau)}{2} \|\mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k\|^2. \tag{53}
\end{aligned}$$

Moreover, Equation (16) ensures that

$$(\forall k \geq \tau) \quad \left\langle \nabla_{(\mathcal{B}^k)} f(\mathbf{x}^{\iota_k}), \mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k \right\rangle = -\|\mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k\|_{\mathbf{A}_{(\mathcal{B}^k)}(\mathbf{x}^{\iota_k})}^2. \tag{54}$$

Replacing both (53) and (54) in (51) gives, for all  $k \geq \tau$ ,

$$\begin{aligned}
f(\mathbf{x}^{k+1}) + \nu_{k+1} &\leq f(\mathbf{x}^k) + \nu_k + \frac{\mathcal{L}\sqrt{\tau}(1 + \tau)}{2} \|\mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k\|^2 - \|\mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k\|_{\Gamma_c^k}^2 \\
&= f(\mathbf{x}^k) + \nu_k - \|\mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k\|_{\Gamma_c^k - \frac{\mathcal{L}\sqrt{\tau}(1+\tau)}{2} \mathbf{I}_{|\mathcal{B}^k|}}^2, \tag{55}
\end{aligned}$$

with  $\Gamma_c^k$  defined in Assumption 3(ii). (49) is a direct consequence of Assumption 3(ii) remarking that  $\|\mathbf{x}_{(\mathcal{B}^k)}^{k+1} - \mathbf{x}_{(\mathcal{B}^k)}^k\| = \|\mathbf{x}^{k+1} - \mathbf{x}^k\|$ .  $\square$

#### 4.2. General behaviour

We now state our first convergence Theorem for BD3MG algorithm.

**Theorem 4.1.** *Under Assumptions 1-2-3, sequence  $(f(\mathbf{x}^k))_{k \in \mathbb{N}}$  generated by BD3MG converges to a finite limit  $f^\infty$ . Moreover,  $(\nabla f(\mathbf{x}^k))_{k \in \mathbb{N}}$  converges to  $\mathbf{0}_N$ .*

*Proof.* Coercivity of  $f$  (Assumption 1) and (49) guarantee that  $(f(\mathbf{x}^k) + \nu_k)_{k \in \mathbb{N}}$  is a decreasing and lower-bounded sequence. It is thus converging to a real value  $f^\infty$ . Equation (49) then directly leads to  $\sum_{k=0}^{+\infty} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 < +\infty$ . On the first hand, using the same notation  $(\nu_k)_{k \in \mathbb{N}}$  introduced in our proof of Proposition 4.1, we have

$$(\forall k \geq \tau) \quad \nu_k \leq \frac{\mathcal{L}\tau\sqrt{\tau}}{2} \sum_{j=k-\tau+1}^{+\infty} \|\mathbf{x}^j - \mathbf{x}^{j-1}\|^2. \tag{56}$$

Thus, the sequence  $(\nu_k)_{k \in \mathbb{N}}$  converges to 0 and so, by Proposition 4.1,  $(f(\mathbf{x}^k))_{k \in \mathbb{N}}$  goes to  $f^\infty$ . On the other hand, Lemma 3.3 gives

$$(\forall k \geq 2\tau) \quad \|\nabla f(\mathbf{x}^k)\| \leq \mathcal{L}\tau \sum_{j=k-2\tau+1}^{+\infty} \|\mathbf{x}^j - \mathbf{x}^{j-1}\| + \sum_{j=k-\tau}^{+\infty} \|\mathbf{x}^{j+1} - \mathbf{x}^j\|, \quad (57)$$

which enables to conclude that  $(\nabla f(\mathbf{x}^k))_{k \in \mathbb{N}}$  converges to  $\mathbf{0}_N$ .  $\square$

**Proposition 4.2.** *Under Assumptions 1-2-3,  $\mathcal{C}((\mathbf{x}^k)_{k \in \mathbb{N}})$ , defined as the set of accumulation points of  $(\mathbf{x}^k)_{k \in \mathbb{N}}$ , is non empty, compact and is a subset of the set of stationary points of  $f$ . Moreover,  $f$  takes value  $f^\infty$  on  $\mathcal{C}((\mathbf{x}^k)_{k \in \mathbb{N}})$ .*

*Proof.* Coercivity of  $f$  (by Assumption 1) and convergence of  $(f(\mathbf{x}^k))_{k \in \mathbb{N}}$  to  $f^\infty$  (by Theorem 4.1) show that  $(\mathbf{x}^k)_{k \in \mathbb{N}}$  is a bounded sequence and  $\mathcal{C}((\mathbf{x}^k)_{k \in \mathbb{N}})$  is non empty and compact. Convergence of  $(\nabla f(\mathbf{x}^k))_{k \in \mathbb{N}}$  to  $\mathbf{0}_N$  (by Theorem 4.1) guarantees that every point  $\mathbf{x}^* \in \mathcal{C}((\mathbf{x}^k)_{k \in \mathbb{N}})$  is a stationary point of  $f$ . Moreover, using again  $(f(\mathbf{x}^k))_{k \in \mathbb{N}}$  converging to  $f^\infty$  yields  $f^\infty = f(\mathbf{x}^*)$  for every  $\mathbf{x}^* \in \mathcal{C}((\mathbf{x}^k)_{k \in \mathbb{N}})$  which concludes the proof.  $\square$

### 4.3. Lyapunov-based asymptotical analysis

In order to establish the convergence of the iterates of BD3MG, we will follow an analysis relying on the use of a Lyapunov function. Such proof technique has also been used in [24, 74, 71]. The idea is to exhibit a function, related to the loss function  $f$  but non necessarily equals to it, that decreases monotonically along the iterative process. Given (49), a natural choice is

$$L : \mathbf{Z} = \begin{pmatrix} \mathbf{Z}_0 \\ \vdots \\ \mathbf{Z}_\tau \end{pmatrix} \in \mathbb{R}^{(\tau+1)N} \mapsto f(\mathbf{Z}_0) + \frac{\mathcal{L}\sqrt{\tau}}{2} \sum_{\ell=1}^{\tau} (\tau - \ell + 1) \|\mathbf{Z}_\ell - \mathbf{Z}_{\ell-1}\|^2. \quad (58)$$

Let us denote, for every  $k \geq \tau$ ,  $\mathbf{Z}^k = \begin{pmatrix} \mathbf{x}^k \\ \vdots \\ \mathbf{x}^{k-\tau} \end{pmatrix} \in \mathbb{R}^{(\tau+1)N}$ , with  $\mathbf{x}^k$  the  $k$ -th BD3MG iterate. Then, the descent condition from Proposition 4.1 can be rewritten as

$$(\forall k \geq \tau) \quad L(\mathbf{Z}^{k+1}) \leq L(\mathbf{Z}^k) - \nu \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2. \quad (59)$$

The structure of  $L$  allows to build an upper bound of its gradient norm along the iterates, where the bound depends only on the differences of the past iterates:

**Lemma 4.1.** *There exists  $\rho > 0$  such that*

$$(\forall k \geq \tau) \quad \|\nabla L(\mathbf{Z}^k)\| \leq \rho \sum_{j=k-2\tau+1}^k \|\mathbf{x}^j - \mathbf{x}^{j-1}\|. \quad (60)$$

*Proof.* Function  $L$  is differentiable. The expression of its gradient is

$$\left(\forall \mathbf{Z} \in \mathbb{R}^{(\tau+1)N}\right) \quad \nabla L(\mathbf{Z}) = \mathbf{g}_0 + \mathcal{L}\sqrt{\tau} \sum_{\ell=1}^{\tau} (\tau - \ell + 1) \mathbf{g}_\ell, \quad (61)$$

$$\text{where } \mathbf{g}_0 = \begin{pmatrix} \nabla f(\mathbf{Z}_0) \\ \mathbf{0}_{\tau N} \end{pmatrix} \text{ and } (\forall \ell \in \llbracket 1, \tau \rrbracket) \quad \mathbf{g}_\ell = \begin{pmatrix} \mathbf{0}_{(\ell-1)N} \\ \mathbf{Z}_{\ell-1} - \mathbf{Z}_\ell \\ \mathbf{Z}_\ell - \mathbf{Z}_{\ell-1} \\ \mathbf{0}_{(\tau-\ell)N} \end{pmatrix}. \quad (62)$$

Let us apply twice the Jensen inequality for the square of the norm and then the majoration  $\tau - \ell + 1 \leq \tau$  for  $1 \leq \ell \leq \tau$ . This yields

$$\begin{aligned} \left(\forall \mathbf{Z} \in \mathbb{R}^{(\tau+1)N}\right) \quad \|\nabla L(\mathbf{Z})\|^2 &\leq 2\|\mathbf{g}_0\|^2 + 2(\mathcal{L}\tau)^2 \sum_{\ell=1}^{\tau} (\tau - \ell + 1)^2 \|\mathbf{g}_\ell\|^2 \\ &= 2\|\nabla f(\mathbf{Z}_0)\|^2 + 4(\mathcal{L}\tau)^2 \sum_{\ell=1}^{\tau} (\tau - \ell + 1)^2 \|\mathbf{Z}_\ell - \mathbf{Z}_{\ell-1}\|^2 \\ &\leq 2\|\nabla f(\mathbf{Z}_0)\|^2 + 4\mathcal{L}^2\tau^4 \sum_{\ell=1}^{\tau} \|\mathbf{Z}_\ell - \mathbf{Z}_{\ell-1}\|^2. \end{aligned} \quad (63)$$

Using  $\sqrt{a^2 + b^2} \leq a + b$  for the two quantities at the right of (63) and then standard norm majoration inequalities, we get:

$$\left(\forall \mathbf{Z} \in \mathbb{R}^{(\tau+1)N}\right) \quad \|\nabla L(\mathbf{Z})\| \leq \sqrt{2}\|\nabla f(\mathbf{Z}_0)\| + 2\mathcal{L}\tau^2 \sum_{\ell=1}^{\tau} \|\mathbf{Z}_\ell - \mathbf{Z}_{\ell-1}\|. \quad (64)$$

The application of (64) to sequence  $(\mathbf{Z}^k)_{k \in \mathbb{N}}$  leads to

$$(\forall k \geq \tau) \quad \|\nabla L(\mathbf{Z}^k)\| \leq \sqrt{2}\|\nabla f(\mathbf{x}_k)\| + 2\mathcal{L}\tau^2 \sum_{j=k-\tau+1}^k \|\mathbf{x}^j - \mathbf{x}^{j-1}\|. \quad (65)$$

By Lemma 3.3 and (65), we finally deduce that

$$\begin{aligned} (\forall k \geq 2\tau) \quad \|\nabla L(\mathbf{Z}^k)\| &\leq \sqrt{2}\mathcal{L}\tau \sum_{j=k-2\tau+1}^k \|\mathbf{x}^j - \mathbf{x}^{j-1}\| + \sqrt{2\bar{\nu}} \sum_{j=k-\tau}^k \|\mathbf{x}^{j+1} - \mathbf{x}^j\| \\ &\quad + 2\mathcal{L}\tau^2 \sum_{j=k-\tau+1}^k \|\mathbf{x}^j - \mathbf{x}^{j-1}\| \\ &\leq \left(\sqrt{2}L\tau + \sqrt{2\bar{\nu}} + 2L\tau^2\right) \sum_{j=k-2\tau+1}^k \|\mathbf{x}^j - \mathbf{x}^{j-1}\|, \end{aligned} \quad (66)$$

which concludes the proof taking  $\rho = \sqrt{2}L\tau + \sqrt{2\bar{\nu}} + 2L\tau^2$ .  $\square$

The following analysis makes use of recent theoretical results around the KL inequality [3, 8] that we recall hereafter. For every  $\zeta > 0$ , we denote by  $\Phi_\zeta$  the set of concave functions  $\varphi: [0, \zeta) \mapsto \mathbb{R}_+$  verifying :

- $\varphi(0) = 0$ .
- $\varphi \in C^1((0, \zeta))$  and is continuous in 0.
- $\forall s \in (0, \zeta), \varphi'(s) > 0$ .

We are then ready to introduce the so-called KL property. [3, 8]

**Definition 4.1.** [KL property] A differentiable function  $g : \mathbb{R}^d \rightarrow \mathbb{R}$ , with  $d \geq 1$ , satisfies the Kurdyka-Lojasiewicz (KL) property on  $E \subset \mathbb{R}^d$  if, for every  $\mathbf{z} \in E$  and every bounded neighborhood  $V$  of  $\mathbf{z}$ , there exist  $\zeta > 0$  and  $\varphi \in \Phi_\zeta$  such that every  $\mathbf{x} \in E \cap \{\mathbf{x} \text{ s.t. } |g(\mathbf{x}) - g(\mathbf{z})| < \zeta\}$ ,

$$\|\nabla g(\mathbf{x})\| \varphi'(|g(\mathbf{x}) - g(\mathbf{z})|) \geq 1. \quad (67)$$

We also recall the following Lemma:

**Lemma 4.2.** [Uniform KL property [8, Lemma 6]] Let  $C$  a compact set of  $\mathbb{R}^d$  and  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  a differentiable function satisfying KL property on  $C$  and constant on the latter. Then, there exist  $\epsilon, \zeta > 0$  and  $\varphi \in \Phi_\zeta$  such that every  $\bar{\mathbf{x}} \in C$  and all  $\mathbf{x} \in \mathbb{R}^d$  satisfying both  $d(\mathbf{x}, C) < \epsilon$ ,  $0 < g(\mathbf{x}) - g(\bar{\mathbf{x}}) < \zeta$ , we have

$$\|\nabla g(\mathbf{x})\| \varphi'(|g(\mathbf{x}) - g(\bar{\mathbf{x}})|) \geq 1. \quad (68)$$

**Proposition 4.3.** Under Assumptions 1-2-3, if  $L$  defined in (58) fulfills the KL property on  $\mathbb{R}^{(\tau+1)N}$  then, considering  $g = L$ ,  $C = \mathcal{C}((\mathbf{Z}^k)_{k \in \mathbb{N}})$  with  $L(C) = \{f^\infty\}$ , there exists  $\epsilon^L, \zeta^L$  and  $\phi^L \in \Phi_{\zeta^L}$  such that  $L$  satisfies (68).

*Proof.* This is a direct consequence of Lemma 4.2. Continuity of  $L$  is clear. We still have to verify the compactness of  $\mathcal{C}((\mathbf{Z}^k)_{k \in \mathbb{N}})$  and that  $L$  is constant valued on that set.  $\mathcal{C}((\mathbf{Z}^k)_{k \in \mathbb{N}})$  is closed. Moreover, it is straightforward to show that this set is included in the Cartesian product  $\left[\mathcal{C}((\mathbf{x}^k)_{k \in \mathbb{N}})\right]^{\tau+1}$ , where  $\mathcal{C}((\mathbf{x}^k)_{k \in \mathbb{N}})$  is compact.  $\mathcal{C}((\mathbf{Z}^k)_{k \in \mathbb{N}})$  is thus bounded and, finally, it is compact.

Let  $\mathbf{Z} \in \mathcal{C}((\mathbf{Z}^k)_{k \in \mathbb{N}})$ . We have  $L(\mathbf{Z}^k) = f(\mathbf{x}^k) + \nu_k$  for all  $k \in \mathbb{N}$ . From our proof of Theorem 4.1, it follows that sequence  $(L(\mathbf{Z}^k))_{k \in \mathbb{N}}$  converges to  $f^\infty$ . Continuity of  $L$  finally ensures that  $f^\infty = L(\mathbf{Z})$ . This proves that  $f$  is constant valued on  $\mathcal{C}((\mathbf{Z}^k)_{k \in \mathbb{N}})$  (and equals to  $f^\infty$ ).  $\square$

#### 4.4. Convergence of the iterates

We are now ready to state our second convergence Theorem for BD3MG algorithm, characterizing the convergence of  $(\mathbf{x}_k)_{k \in \mathbb{N}}$ .

**Theorem 4.2.** Let assume that Assumptions 1-2-3 hold. Assume furthermore that the Lyapunov function  $L$  in (58) satisfies the KL property on  $\mathbb{R}^{(\tau+1)N}$ . Then, sequence  $(\mathbf{x}^k)_{k \in \mathbb{N}}$  is of finite length, i.e :

$$\sum_{k=0}^{+\infty} \|\mathbf{x}^{k+1} - \mathbf{x}^k\| < +\infty, \quad (69)$$

and converges to a stationary point of  $f$ .

*Proof.* Let us start considering the case when there exists some  $k_0 \in \mathbb{N}$  where  $L(\mathbf{Z}^{k_0}) = f^\infty$ . Since  $(L(\mathbf{Z}^k))_{k \in \mathbb{N}}$  is decreasing sequence converging to  $f^\infty$  (see proof of Proposition 4.3), it follows that  $L(\mathbf{Z}^k) = f^\infty$  for all  $k \geq k_0$ . (59) then gives

$$(\forall k \geq k_0) \quad \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 \leq \underline{\nu}^{-1} \left( L(\mathbf{Z}^k) - L(\mathbf{Z}^{k+1}) \right) = 0, \quad (70)$$

ensuring that  $(\mathbf{x}^k)_{k \in \mathbb{N}}$  has a finite length and  $\mathbf{x}^k$ ,  $k \geq 0$ , is a stationary point of  $f$ .

We now suppose that, for all  $k \in \mathbb{N}$ ,  $L(\mathbf{Z}^{k_0}) \neq f^\infty$ . We aim at exhibiting a uniform KL inequality on sequence  $(L(\mathbf{Z}^k))_{k \in \mathbb{N}}$ . To do so, let us peruse the quantities  $\epsilon^L, \eta^L, \varphi^L$  arising from Proposition 4.3.

On the one hand, the decrease of  $(L(\mathbf{Z}^k))_{k \in \mathbb{N}}$  implies that, for all  $k \in \mathbb{N}$ ,  $L(\mathbf{Z}^k) > f^\infty$ . The set  $\mathcal{C}((\mathbf{x}^k)_{k \in \mathbb{N}})$  is non empty (see proof of Proposition 4.2), so is the set  $\mathcal{C}((\mathbf{Z}^k)_{k \in \mathbb{N}})$ . Let  $\mathbf{Z} \in \mathcal{C}((\mathbf{Z}^k)_{k \in \mathbb{N}})$  an element of such set i.e., a cluster point of  $(\mathbf{Z}^k)_{k \in \mathbb{N}}$ . From Proposition 4.3,  $L(\mathbf{Z}) = f^\infty$ . Hence,  $L(\mathbf{Z}^k) - L(\mathbf{Z}) > 0$  for all  $k \in \mathbb{N}$ .

On the other hand,  $(L(\mathbf{Z}^k))_{k \in \mathbb{N}}$  converges to  $f^\infty = L(\mathbf{Z})$ . The boundedness of  $(\mathbf{x}^k)_{k \in \mathbb{N}}$  also ensures this of  $(\mathbf{Z}^k)_{k \in \mathbb{N}}$ .

We deduce the existence of some  $k_1 \geq 2\tau$  such that

$$(\forall k \geq k_1) \quad 0 < L(\mathbf{Z}^k) - L(\mathbf{Z}) < \eta^L, \quad d\left(\mathbf{Z}^k, \mathcal{C}((\mathbf{Z}^k)_{k \in \mathbb{N}})\right) < \epsilon^L. \quad (71)$$

From Proposition 4.3, the uniform KL property on  $L$  holds i.e.,

$$(\forall k \geq k_1) \quad \|\nabla L(\mathbf{Z}^k)\| \left(\varphi^L\right)' \left(L(\mathbf{Z}^k) - L(\mathbf{Z})\right) \geq 1. \quad (72)$$

Moreover, setting  $\Delta^k = \varphi^L(L(\mathbf{Z}^k) - L(\mathbf{Z})) - \varphi^L(L(\mathbf{Z}^{k+1}) - L(\mathbf{Z}))$  for all  $k \in \mathbb{N}$ , concavity of  $\varphi^L$  and (59) ensure that

$$\begin{aligned} (\forall k \geq k_1) \quad \Delta^k &\geq \left(\varphi^L\right)' \left(L(\mathbf{Z}^k) - L(\mathbf{Z})\right) \left(L(\mathbf{Z}^k) - L(\mathbf{Z}^{k+1})\right) \\ &\geq \bar{\nu} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 \left(\varphi^L\right)' \left(L(\mathbf{Z}^k) - L(\mathbf{Z})\right). \end{aligned} \quad (73)$$

The combination of the latter with (72) leads to

$$(\forall k \geq k_1) \quad \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 \leq \bar{\nu}^{-1} \Delta^k \|\nabla L(\mathbf{Z}^k)\|. \quad (74)$$

By Lemma 4.1, we can upper bound the gradient term in (74). This gives

$$(\forall k \geq k_1) \quad \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 \leq \rho \bar{\nu}^{-1} \Delta^k \sum_{j=k-2\tau+1}^k \|\mathbf{x}^j - \mathbf{x}^{j-1}\|. \quad (75)$$

Passing to the root and using the classical identity  $\sqrt{ab} \leq a/c + bc/4$ , with  $a = \sum_{j=k-2\tau+1}^k \|\mathbf{x}^j - \mathbf{x}^{j-1}\|$  for all  $k \geq k_1$ ,  $b = \Delta^k$ , both positive for all  $k \geq k_1$  and some



$c > 0$  is generic, leads to

$$(\forall k \geq k_1) \quad \|\mathbf{x}^{k+1} - \mathbf{x}^k\| \leq \frac{\sqrt{\rho\nu}^{-1/2}}{c} \sum_{j=k-2\tau+1}^k \|\mathbf{x}^j - \mathbf{x}^{j-1}\| + \frac{c\sqrt{\rho\nu}^{-1/2}}{4} \Delta^k. \quad (76)$$

Since  $(\Delta^k)_{k \in \mathbb{N}}$  is summable (as a telescopic sequence), we can apply Lemma 3.4 with some  $c > 2\tau\sqrt{\rho\nu}^{-1/2}$  so that  $2\tau\frac{\sqrt{\rho\nu}^{-1/2}}{c} < 1$ . This shows that sequence  $(\mathbf{x}^k)_{k \in \mathbb{N}}$  has a finite length.

This finite length property entails that  $(\mathbf{x}^k)_{k \in \mathbb{N}}$  is a Cauchy sequence and thus a converging one. The final conclusion directly comes from Proposition 3.1, ensuring that every accumulation point of  $(\mathbf{x}^k)_{k \in \mathbb{N}}$  is a stationary point of  $f$ .  $\square$

#### 4.5. Discussion

Under the KL condition for the Lyapunov function  $L$  defined in (58), we were able to demonstrate the convergence of sequence  $(\mathbf{x}^k)_{k \in \mathbb{N}}$  to a stationary point of  $f$ . Let us notice that  $f$  satisfying KL property does not necessary imply that  $L$  does. Still, our assumption on  $L$  can be verified in practice for a wide class of functions  $f$ . For instance, following the discussion in [24, section 6], if  $f$  is semi-algebraic [7, 8], then the required condition on  $L$  in Theorem 4.2 is satisfied, with function  $\varphi^L = \kappa(\cdot)^{1-\theta}$  for a some  $(\kappa, \theta) \in \mathbb{R}_+^* \times (0, 1)$ . Such situation will be met in our experimental settings in Section 5. Extending Theorem 4.2 to any KL function  $f$  would be an interesting avenue for future work but up to our knowledge, it does not seem straightforward.

## 5. Application to 3D image restoration

### 5.1. Problem statement

*5.1.1. Observation model.* We focus on the inverse problem of restoring a vectorized 3D volume  $\bar{\mathbf{x}}$  of size  $N = N_X \times N_Y \times N_Z$  given blurry and noisy observation  $\mathbf{y} \in \mathbb{R}^N$ . We consider a depth-variant 3D blur operator  $\mathbf{H} \in \mathbb{R}^{N \times N}$  associated to kernels with support size  $M = M_X \times M_Y \times M_Z$ , and additive i.i.d. Gaussian noise with standard deviation  $\sigma > 0$ , so that the observed volume is related to  $\bar{\mathbf{x}}$  through,

$$\mathbf{y} = \mathbf{H}\bar{\mathbf{x}} + \mathbf{b}, \quad (77)$$

with vector  $\mathbf{b} \in \mathbb{R}^N$  accounting for the noise. The goal is to solve the inverse problem of estimating  $\bar{\mathbf{x}}$  given  $\mathbf{y}$  and  $\mathbf{H}$ . Depth-variant blurs are commonly encountered in 3D microscopy [62, 38, 45, 44], due to optical aberrations. They are particular cases of spatially-variant blurs [12, 53]. The degradation operator  $\mathbf{H}$  raises specific challenges due to its high computational cost. Several strategies have been investigated in the case of 2D spatially variant blur maps encountered for instance in astronomical imaging [27, 26, 29]. The extension to 3D maps of these methods is however not covered up to our knowledge. This motivates the use of a distributed optimization approach for solving the inverse problem (77).

*5.1.2. Objective function* We adopt a variational strategy, which consists in seeking for an estimate of  $\bar{\mathbf{x}}$  that minimizes a penalized least squares criterion  $f$ . A hybrid regularization term is employed incorporating prior knowledge on the smoothness and the range of the sought solution. The objective function reads:

$$(\forall \mathbf{x} \in \mathbb{R}^N) \quad f(\mathbf{x}) = \sum_{s=1}^S f_s(\mathbf{L}_s \mathbf{x}), \quad (78)$$

where, for every  $s \in \{1, \dots, S\}$ ,  $\mathbf{L}_s \in \mathbb{R}^{P_s \times N}$ ,  $P_s \in \mathbb{N}^*$ , and  $f_s$  is a function from  $\mathbb{R}^{P_s}$  to  $\mathbb{R}$ .  $f_1 \circ \mathbf{L}_1$  represents the data fidelity term while the other terms are regularization terms. Here, we set  $S = 4$  and

- $P_1 = N$ ,  $\mathbf{L}_1 = \mathbf{H}$ ,  $f_1 = \frac{1}{2} \|\cdot - \mathbf{y}\|^2$ ,
- $P_2 = N$ ,  $\mathbf{L}_2 = \mathbf{I}_N$ ,  $f_2 = \eta d_{[x_{\min}, x_{\max}]}^2$ ,
- $P_3 = 2N$ ,  $\mathbf{L}_3 = [(\mathbf{V}^X)^\top (\mathbf{V}^Y)^\top]^\top$ ,  $f_3 = \lambda \sum_{n=1}^N \sqrt{[\cdot]_n^2 + [\cdot]_{N+n}^2 + \delta^2}$ ,
- $P_4 = N$ ,  $\mathbf{L}_4 = \mathbf{V}^Z$ ,  $f_4 = \kappa \|\cdot\|^2$ .

Hereabove,  $(\eta, \lambda, \delta, \kappa) \in (0, +\infty)^4$  are hyper-parameters. The linear operators  $\mathbf{V}^X, \mathbf{V}^Y, \mathbf{V}^Z \in \mathbb{R}^{N \times N}$  are discrete gradient operators along X (horizontal), Y (vertical), and Z (longitudinal) directions of the 3D volume. Function  $d_{[x_{\min}, x_{\max}]}^2$  states for the squared distance to set  $[x_{\min}, x_{\max}]^N \subset \mathbb{R}^N$ , with  $(x_{\min}, x_{\max}) \in \mathbb{R}^2$  minimal and maximal bounds on the sought intensity values. The later term can be viewed as an exterior penalty function [18]. ‡

*5.1.3. Majorant mapping.* In order to implement BD3MG, we must build a majorant mapping ensuring the majorization condition (3). First, let us notice that, according to (78), the gradient of  $f$  reads

$$(\forall \mathbf{x} \in \mathbb{R}^N) \quad \nabla f(\mathbf{x}) = \sum_{s=1}^S \mathbf{L}_s^\top \varphi_s(\mathbf{L}_s \mathbf{x}), \quad (79)$$

with, for every  $s \in \{1, \dots, S\}$ ,  $\varphi_s : \mathbb{R}^{P_s} \rightarrow \mathbb{R}^{P_s}$  the gradient operator of  $f_s$ . Then, function  $f$  fits within the class of half-quadratic majorizing constructions initially introduced in [34, 35] and later analysed in [55, 2, 16]. A general structure for the majorant mapping of (78) is thus

$$(\forall \mathbf{x} \in \mathbb{R}^N) \quad \mathbf{A}(\mathbf{x}) = \sum_{s=1}^S \mathbf{L}_s^\top \text{Diag} \left\{ [\omega_s(\mathbf{L}_s \mathbf{x})]_p \right\} \mathbf{L}_s, \quad (80)$$

‡ Function  $\mathbf{x} \in \mathbb{R}^N \mapsto d_E^2(\mathbf{x})$  is 2-Lipschitz differentiable as soon as  $E$  is non-empty closed and convex set [4]. Denoting by  $p_E$  the orthogonal projection operator, its gradient then corresponds to  $\mathbf{x} \in \mathbb{R}^N \mapsto 2(\mathbf{x} - P_E(\mathbf{x}))$ .

where, for every  $s \in \{1, \dots, S\}$ ,  $\omega_s : \mathbb{R}^{P_s} \rightarrow ]0, +\infty[^{P_s}$  is a majorizing potential that depends on the properties of  $(f_s)_{1 \leq s \leq S}$  [20, Tab. I]. In our case, for every  $s \in \{1, \dots, 4\}$ , each of these terms is  $f_s$  is  $\beta_s$ -Lipschitz differentiable with

$$\begin{cases} \beta_1 = 1, \\ \beta_2 = 4\eta, \\ \beta_3 = \lambda\delta^{-1}, \\ \beta_4 = 2\kappa. \end{cases} \quad (81)$$

Then, from descent lemma [6], a valid choice is  $\omega_s(\cdot) = \alpha\beta_s\mathbf{1}_{P_s}$  with some  $\alpha \geq 1$  [16]. We adopt this simple strategy for functions  $f_1$ ,  $f_2$  and  $f_4$ , which yields

$$\begin{cases} \omega_1(\cdot) = \alpha\mathbf{1}_N \\ \omega_2(\cdot) = 4\alpha\eta\mathbf{1}_N \\ \omega_4(\cdot) = 2\alpha\kappa\mathbf{1}_N. \end{cases} \quad (82)$$

Regarding function  $f_3$ , a more sophisticated majorization is adopted, inherited from half-quadratic strategies [2, 55]:

$$(\forall \mathbf{v} \in \mathbb{R}^{2N}) \quad \omega_3(\mathbf{v}) = \lambda \left[ \begin{array}{c} \left( 1/\sqrt{v_n^2 + v_{N+n}^2 + \delta^2} \right)_{1 \leq n \leq N} \\ \left( 1/\sqrt{v_n^2 + v_{N+n}^2 + \delta^2} \right)_{1 \leq n \leq N} \end{array} \right]. \quad (83)$$

A quadratic majorant function satisfying (3) for a given block  $\mathcal{S} \in \mathbb{T}$  can then be obtained using (2) with (79) and (80).

*5.1.4. Distributed implementation* We implement BD3MG algorithm as presented in Section 2.4. Our code is available at <https://github.com/mathieuuchal/BD3MG>. We split the 3D volume into 2D slices along the depth axis  $z \in \{1, \dots, N_z\}$ , and consider each 2D slice as an individual block upon which workers can compute an update. Assuming a lexicographic ordering of the voxels, this means that the following partition is adopted:

$$\mathbb{T} = \left\{ \left[ (i-1)N_x N_y + 1, iN_x N_y \right] \mid 1 \leq i \leq N_z \right\}. \quad (84)$$

BD3MG is implemented on a star graph of workers with a specific master aggregating the current solution. For a given number of active cores  $C_{\text{tot}} = C + 1$  of the computer (or of the cluster), one is used as the master process to manage the computation split between the workers while all the  $C (= C_{\text{tot}} - 1)$  others, are computing updates asynchronously on planar blocks (i.e., Algorithm 2). We initially set, for every  $c \in \{1, \dots, C\}$ ,  $\mathcal{S}_c^0$  corresponding to the index set of the  $((c-1)\lfloor \frac{N_z}{C} \rfloor + 1)$ -th 2D slice in the volume. Then, at each iteration  $k$ , the master requires worker  $c^k$  to process the 2D slice with index set  $\mathcal{S}_{c^k}^{k+1}$ , by applying a first-in, first-out basis. The worker  $c^k$  hence computes the update

for the 2D slice that has been modified the longest time ago, assuming it is available (i.e., not processed in the same time by another worker). A cyclic block update is used as default choice, if several blocks are available (this typically arises in the first iterations). Furthermore, the master controls that each slice has been updated at least once every  $\tau$  iterations. Regarding data exchange, as emphasized in Section 2.5, in practice, it is not necessary to share the full vector  $\mathbf{x}$  with all the workers. Consider a worker update associated to the block  $\mathcal{S} \in \mathbb{T}$ . The worker has to compute  $\nabla_{(\mathcal{S})}f(\mathbf{x})$  and  $\mathbf{A}_{(\mathcal{S})}(\mathbf{x})$ . Because of the structure of (78), these quantities actually only depend on a subpart of vector  $\mathbf{x}$ , defined by  $(x_n)_{n \in \mathbb{V}_{\mathcal{S}}}$ , with  $\mathbb{V}_{\mathcal{S}} \subset \llbracket 1, N \rrbracket$  a set which has low cardinality compared to the full volume size  $N$ . Let us explicit this set for our practical example. The key ingredients to account for are (i) the presence of null entries in the linear operators  $(\mathbf{L}_s)_{1 \leq s \leq S}$ , (ii) the (almost) separability of operators  $(\varphi_s, \omega_s)_{1 \leq s \leq S}$ . We introduce the following sets, for every  $s \in \{1, \dots, S\}$ ,

$$(\forall n \in \{1, \dots, N\}) \quad \text{col}_{n,s} = \{p \in \{1, \dots, P_s\} \text{ s.t. } (\mathbf{L}_s)_{p,n} \neq 0\}, \quad (85)$$

$$(\forall p \in \{1, \dots, P_s\}) \quad \text{row}_{p,s} = \{n \in \{1, \dots, N\} \text{ s.t. } (\mathbf{L}_s)_{p,n} \neq 0\}, \quad (86)$$

Moreover the separable structures of  $(\varphi_s, \omega_s)_{1 \leq s \leq S}$  ensure that for every  $s \in \{1, \dots, S\}$  and  $p \in \{1, \dots, P_s\}$ , there exists a subset  $\mathcal{V}_{s,p} \subset \llbracket 1, P_s \rrbracket$  as well as two functions  $\tilde{\varphi}_{s,p} : \mathbb{R}^{|\mathcal{V}_{s,p}|} \rightarrow \mathbb{R}$  and  $\tilde{\omega}_{s,p} : \mathbb{R}^{|\mathcal{V}_{s,p}|} \rightarrow (0, +\infty)$  such that

$$(\forall \mathbf{v} \in \mathbb{R}^{P_s}) \quad \begin{cases} [\varphi_s(\mathbf{v})]_p = \tilde{\varphi}_{s,p}(\mathbf{v}_{\mathcal{V}_{s,p}}), \\ [\omega_s(\mathbf{v})]_p = \tilde{\omega}_{s,p}(\mathbf{v}_{\mathcal{V}_{s,p}}). \end{cases} \quad (87)$$

Considering this, we can now rewrite the expressions  $\nabla_{(\mathcal{S})}f(\mathbf{x})$  and  $\mathbf{A}_{(\mathcal{S})}(\mathbf{x})$  as

$$(\forall \mathbf{x} \in \mathbb{R}^N) \quad \nabla_{(\mathcal{S})}f(\mathbf{x}) = \left( [\nabla f(\mathbf{x})]_i \right)_{i \in \mathcal{S}}, \quad (88)$$

with, for every  $i \in \mathcal{S}$ ,

$$[\nabla f(\mathbf{x})]_i = \sum_{s=1}^S \left[ \mathbf{L}_s^\top \varphi_s(\mathbf{L}_s \mathbf{x}) \right]_i, \quad (89)$$

$$= \sum_{s=1}^S \sum_{p=1}^{P_s} (\mathbf{L}_s)_{p,i} [\varphi_s(\mathbf{L}_s \mathbf{x})]_p, \quad (90)$$

$$= \sum_{s=1}^S \sum_{p \in \text{col}_{i,s}} (\mathbf{L}_s)_{p,i} [\varphi_s(\mathbf{L}_s \mathbf{x})]_p, \quad (91)$$

$$= \sum_{s=1}^S \sum_{p \in \text{col}_{i,s}} (\mathbf{L}_s)_{p,i} \tilde{\varphi}_{s,p}([\mathbf{L}_s \mathbf{x}]_{(\mathcal{B}_{s,p})}), \quad (92)$$

$$= \sum_{s=1}^S \sum_{p \in \text{col}_{i,s}} (\mathbf{L}_s)_{p,i} \tilde{\varphi}_{s,p} \left( \left[ \sum_{n=1}^N (\mathbf{L}_s)_{\ell,n} x_n \right]_{\ell \in \mathcal{V}_{s,p}} \right), \quad (93)$$

$$= \sum_{s=1}^S \sum_{p \in \text{col}_{i,s}} (\mathbf{L}_s)_{p,i} \tilde{\varphi}_{s,p} \left( \left[ \sum_{n \in \text{row}_{\ell,s}} (\mathbf{L}_s)_{\ell,n} x_n \right]_{\ell \in \mathcal{V}_{s,p}} \right), \quad (94)$$

Similar computation shows that, for every  $(i, j) \in \mathcal{S}^2$ ,

$$[\mathbf{A}(\mathbf{x})]_{i,j} = \sum_{s=1}^S \sum_{p \in \text{col}_{i,s} \cap \text{col}_{j,s}} (\mathbf{L}_s)_{p,i} (\mathbf{L}_s)_{p,j} \tilde{\omega}_{s,p} \left( \left[ \sum_{n \in \text{row}_{\ell,s}} (\mathbf{L}_s)_{\ell,n} x_n \right]_{\ell \in \mathcal{V}_{s,p}} \right). \quad (95)$$

Hence, (94)-(95) reflect the fact that the only coordinates of the vector  $\mathbf{x}$  that are manipulated to compute the gradient and majorant mapping related to block  $\mathcal{S}$ , belong to  $\mathbb{V}_{\mathcal{S}}$  where

$$\mathbb{V}_{\mathcal{S}} = \bigcup_{i \in \mathcal{S}} \bigcup_{s \in \{1, \dots, S\}} \bigcup_{p \in \text{col}_{i,s}} \bigcup_{\ell \in \mathcal{V}_{s,p}} \text{row}_{\ell,s}. \quad (96)$$

Since matrices  $(\mathbf{L}_s)_{1 \leq s \leq S}$  are very sparse and functions  $(\varphi_s, \omega_s)_{1 \leq s \leq S}$  close to separable ones, the cardinality of the involved sets in (96) is small so that the memory load for each communication in between master and worker is limited §.

*5.1.5. Validity of Assumptions.* Let us discuss the validity of Assumptions 1, 2 and 3 for the considered problem and implementation.

§ Let us point out that it is not necessary for the worker to have access to the full operators  $(\mathbf{L}_s)_{1 \leq s \leq S}$  but only to the entries involved in expressions (94)-(95). In particular, in our implementation, only the kernel blurs related to the depths in a neighborhood of the 2D slice are shared, reducing again the memory load (see for instance `apply_PSFvar3Dz_block.py` in our shared code folder).

*Assumption 1.* Function  $f$  in (78) is differentiable. Moreover, it has a  $\mathcal{L}$ -Lipschitzian gradient with  $\mathcal{L} = \sum_{s=1}^S \beta_s \|\mathbf{L}_s\|^2$ , where  $\|\cdot\|$  denotes the spectral norm over matrices and  $(\beta_s)_{1 \leq s \leq S}$  were given in the previous subsection. According to [64, Prop. 2.5], a sufficient condition for  $f$  to be coercive is  $\text{Ker}(\mathbf{H}) = \{\mathbf{0}_N\}$ . This latter is verified in our experiments, since  $\mathbf{H}$  is a full-rank operator. Thus, Assumption 1 holds.

*Assumption 2.* This assumption relates to the practical implementation of BD3MG and requires every subset of variables to be updated within a finite number of iterations. In practice, we introduced a safety check in the master loop, that introduces an idle time if a slice has not been updated within the last  $\tau$  iterations with  $\tau$  a predefined value. In our implementation, each worker is in average in charge of  $\frac{N_z}{C}$  2D slices, of the volume. We thus set  $\tau = 2 \left\lceil \frac{N_z}{C} \right\rceil$ , that is each worker is allowed to spend, in average twice more time to update one slice than another. Given our block selection rule, with balanced computational load per slide, and relying on first-in, first-out, this situation could only arise if a worker experienced a major delay, which never occurred in our experiments.

*Assumption 3.* This assumption relates to the majorant mapping. To check this assumption, we proceed in three steps. On the one hand, we have,

$$(\forall \mathbf{x} \in \mathbb{R}^N) \quad \mathbf{A}(\mathbf{x}) \geq \mathbf{L}_2^\top \text{Diag} \left\{ \omega_2(\mathbf{L}_2 \mathbf{x}) \right\} \mathbf{L}_2 \geq \alpha \eta \mathbf{I}_N. \quad (97)$$

On the other hand, according to definition (80) and those of  $\omega_1, \dots, \omega_4$

$$(\forall \mathbf{x} \in \mathbb{R}^N) \quad \mathbf{A}(\mathbf{x}) \leq \left( \sum_{s=1}^S \|\mathbf{L}_s\|^2 \max_{1 \leq p \leq P_s} [\omega_s(\mathbf{L}_s \mathbf{x})]_p \right) \mathbf{I}_N \leq \bar{\nu} \mathbf{I}_N, \quad (98)$$

with

$$\bar{\nu} = \alpha \left( \sum_{s=1}^S \beta_s \|\mathbf{L}_s\|^2 \right). \quad (99)$$

Considering (97), (98)-(99) and the fact that any sub-matrix  $\mathcal{M}_{(\mathcal{S})}$  ( $\mathcal{S} \subset \llbracket 1, N \rrbracket$ ) of a (symmetric) positive matrix  $\mathbf{M}$  remains positive, the chosen mapping  $\mathbf{A}$  thus respects

conditions imposed by Assumption 3(i). Moreover, for all  $(\mathbf{x}, \mathbf{y}) \in (\mathbb{R}^N)^2$ ,

$$\begin{aligned}
& \mathbf{A}(\mathbf{x}) - \frac{1}{2}\mathbf{A}(\mathbf{y}) \\
&= \sum_{s=1}^S (\mathbf{L}_s)^\top \text{Diag}_{1 \leq p \leq P_s} \left\{ \left( [\omega_s(\mathbf{L}_s \mathbf{x})]_p - \frac{1}{2} [\omega_s(\mathbf{L}_s \mathbf{y})]_p \right) \right\} \mathbf{L}_s \\
&= \frac{\alpha}{2} \sum_{s \in \{1, 2, 4\}} (\mathbf{L}_s)^\top \mathbf{L}_s + (\mathbf{L}_3)^\top \text{Diag}_{1 \leq p \leq P_3} \left\{ \left( [\omega_3(\mathbf{L}_3 \mathbf{x})]_p - \frac{1}{2} [\omega_3(\mathbf{L}_3 \mathbf{y})]_p \right) \right\} \mathbf{L}_3 \\
&\geq \frac{\alpha}{2} \mathbf{L}_2^\top \mathbf{L}_2 + (\mathbf{L}_3)^\top \text{Diag}_{1 \leq p \leq P_3} \left\{ \left( [\omega_3(\mathbf{L}_3 \mathbf{x})]_p - \frac{1}{2} [\omega_3(\mathbf{L}_3 \mathbf{y})]_p \right) \right\} \mathbf{L}_3 \\
&\geq \eta(\alpha) \mathbf{I}_N \text{ with } \eta(\alpha) = \frac{\alpha}{2} - \frac{8\lambda}{2\delta},
\end{aligned} \tag{100}$$

as  $\|\mathbf{L}_3\|^2 = 8$ . Under the same previous remark on the block positivity preservation, Assumption 3(ii) is verified considering  $\alpha$  large enough (i.e so as for  $\eta(\alpha)$  to strictly exceed bound  $\frac{\mathcal{L}\sqrt{\tau}(1+\tau)}{2}$ ). In practice, we opt for  $\alpha = 1.1 \times \left( \mathcal{L}\sqrt{\tau}(1+\tau) + \frac{8\lambda}{\delta} \right)$ . The associated  $\underline{\nu}$  in (27) is  $\underline{\nu} = \eta(\alpha) - \frac{\mathcal{L}\sqrt{\tau}(1+\tau)}{2}$ .

*Convergence result* In a nutshell, Assumptions 1-2-3 are fulfilled in our experiments, so that Theorem 4.1 holds. Moreover, function  $f$  is semi-algebraic, hence so is the Lyapunov function  $L$  (see discussion in Sec. 4.5). Thus, Theorem 4.2 holds.

## 5.2. Comparative analysis on a controlled scenario

We first set  $\bar{\mathbf{x}}$  as the 3D microscopic image `FlyBrain` with size  $N = N_X \times N_Y \times N_Z = 256 \times 256 \times 57$ . The linear operator  $\mathbf{H}$  models a 3D depth-varying Gaussian blur. For each depth  $z \in \{1, \dots, N_Z\}$ , the blur kernel is characterized by different variance and rotation parameters  $(\sigma_X(z), \sigma_Y(z), \sigma_Z(z), \varphi_Y(z), \varphi_Z(z))$ , following the model from [72]. In practice, the values of these five parameters are chosen randomly through a uniform distribution over  $[0, 3] \times [0, 3] \times [0, 4] \times [0, 2\pi] \times [0, 2\pi]$ , sampled independently for every  $z$ . The support of the blur is then truncated to reach a kernel size of  $M = M_X \times M_Y \times M_Z = 5 \times 5 \times 11$ , which appears large enough to avoid spurious ringing effects.

A zero-mean white Gaussian noise with standard deviation  $\sigma = 4 \times 10^{-2}$  is then added to the blurred volume. The regularization parameters  $(\lambda, \delta, \kappa, \eta) = (1, 1, 10^{-1}, 10^{-3})$  are chosen empirically so as to maximize the Signal-to-Noise Ratio (SNR) of the restored volume. Moreover, we set  $(x_{\min}, x_{\max}) = (0, 1)$ , equal to the range of the ground truth image. In order to illustrate the acceleration induced by the proposed BD3MG, we run a comparative analysis between different versions of the optimization scheme, in the spirit of an ablation study. Namely, we propose to compare BD3MG with three methods listed hereafter.

|| <https://imagej.nih.gov/ij/images/>

- The 3MG algorithm [16, 17] is considered as the baseline. At each iteration, this algorithm builds the majorant mapping as in Sec. 5.1.3 and computes memory gradient updates on the full volume, without any parallelization.
- The Asynchronous Block Gradient Descent (ABGD) algorithm implements the method from [56]. It performs parallel asynchronous gradient descent updates over the slices of the volume. We adopt here the same parallelization settings as for our BD3MG. Updates correspond to the standard gradient descent on the selected planar blocks, using a fixed step-size  $\mu$  ensuring convergence of the iterative scheme, namely  $\mu = 0.99/(1 + \kappa + 2\lambda/\delta + 2\kappa)$ .
- The BP3MG algorithm from [10, 32] runs a synchronous version of BD3MG algorithm. The master process carries out the main loop of [10, Alg 4.3]. At each iteration  $k \in \mathbb{N}$ , it selects  $C$  block indices (following a cyclic rule) and sends to each worker  $c \in \llbracket 1, C \rrbracket$  the required data allowing it to update  $\mathcal{S}_c^k$ , the associated block. Workers process their block in parallel, wait for each other to finish their tasks, combine their respective updates into a unique vector  $(\mathbf{x}^j)_{j \in \mathcal{S}_1^k \cup \dots \cup \mathcal{S}_C^k}$  and finally send the latter to the master. The majorant mapping is set as a block diagonal matrix, allowing synchronous parallel updates, as described in [10]. This approach could be interpreted as a special case of BD3MG with a single worker (potentially composed of several subworkers) sending its update (potentially composed of several sub-updates) to the central process  $\mathcal{S}^k = \{\mathcal{S}_c^k\}_{c \in C}$ . Thanks to the specific structure of the majorant mapping in BP3MG, there is no mismatch in information between central process and workers in this synchronous version, the delay vector  $i_k$  always equals  $k$ . Nonetheless, the block diagonal form of the majorant mapping of BP3MG is at the price of a lower quality of approximation of the cost function, which might result in slower convergence.

All methods are implemented in *Python* using the built-in *Multiprocessing* library as well as *Numpy* and *Scipy* for both data manipulation and scientific computing. The experiments of this section are conducted on an Intel® Xeon(R) W-2135 CPU with  $C_{\text{tot}} = 12$  cores clocked at 3.70GHz. All the versions were initialized with  $\mathbf{x}^0 = \mathbf{0}_N$  leading to an initial value  $f(\mathbf{x}^0) = 91292.92$ . For every iteration  $k \in \mathbb{N}^*$ , we monitor the cost function  $f(\mathbf{x}^k)$ , the normalized increment  $\|\mathbf{x}^{k+1} - \mathbf{x}^k\|/\|\mathbf{x}^k\|$ , the signal to noise ratio (SNR, in dB) defined as

$$\text{SNR} = 20 \log_{10} \left( \frac{\|\bar{\mathbf{x}}\|}{\|\bar{\mathbf{x}} - \mathbf{x}^k\|} \right), \quad (102)$$

and the reconstruction error  $\|\bar{\mathbf{x}} - \mathbf{x}^k\|$ . The evolution of these metrics along time for the tested algorithms is displayed on Figure 2. We then set a stopping criterion  $\|\mathbf{x}^{k+1} - \mathbf{x}^k\| \leq \tilde{\varepsilon} \|\mathbf{x}^k\|$ . The obtained solution is denoted as  $\mathbf{x}^f$ . We display in Table 1 the metrics for the stopping criterion threshold  $\tilde{\varepsilon} = 10^{-3}$ . Table 1 and Figure 2 show that BD3MG exhibits a faster practical convergence than its competitors. Both BD3MG and ABGD are asynchronous distributed schemes, and the former implements an accelerated



version of the gradient descent involved in the latter. The MM metric and the subspace scheme in BD3MG act as catalyzers, improving the convergence rate compared to ABGD which relies on a simple steepest descent with fixed stepsize. BP3MG and BD3MG are based on the same inherent optimization scheme 3MG. However, BP3MG uses a simplified block diagonal majorant mapping, and imposes synchronous updates, which might yield idle times. These differences can explain why BD3MG converges faster than BP3MG. Finally, 3MG does not exploit the multicore structure of the computing architecture, and thus shows higher computational time.

Slices of the reconstructed volume are displayed in Figure 3, revealing fine details of the image recovered by the restoration procedure.

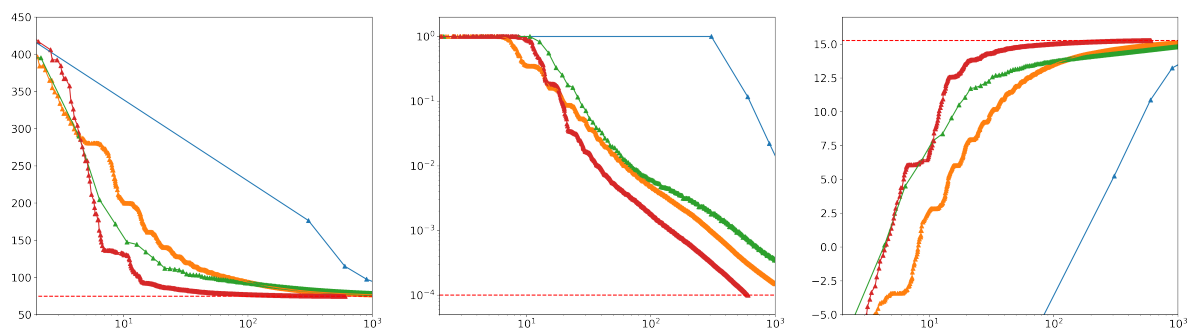


Figure 2: Evolution of quantitative metrics along time (in seconds), for algorithms 3MG (blue), ABGD (orange), BP3MG (green) and BD3MG (red), for FlyBrain restoration. Evolution of reconstruction error  $\|\mathbf{x}^k - \bar{\mathbf{x}}\|$  (left), relative increment  $\|\mathbf{x}^{k+1} - \mathbf{x}^k\|/\|\mathbf{x}^k\|$  (middle), and SNR in dB (right).

Version	//	Asy.	MM	SNR (dB)	$f(\mathbf{x}^f)$	Error	Time ( $\times$ Acc.)
3MG	$\times$	$\times$	$\checkmark$	14.72	1266.04	79.28	1683.79 (-)
ABGD	$\checkmark$	$\checkmark$	$\times$	15.11	1268.80	76.73	305.76 (5.51)
BP3MG	$\checkmark$	$\times$	$\checkmark$	15.21	1264.08	75.33	489.99 (3.44)
BD3MG	$\checkmark$	$\checkmark$	$\checkmark$	<b>15.26</b>	<b>1261.59</b>	<b>74.72</b>	<b>147.16 (11.44)</b>

Table 1: Characteristics and performances of compared algorithms on the Flybrain restoration task, for reaching the stopping criterion with  $\tilde{\epsilon} = 10^{-3}$ . “//” = Parallel, “Asy.” = Asynchronous, “MM” = Majorize-Minimize scheme. Time is in seconds and “ $\times$  Acc.” is the acceleration ratio with respect to 3MG running time.

### 5.3. Effect of an imbalanced computing power

In order to further demonstrate the advantages of BD3MG over its synchronous counterpart BP3MG, we tested the methods under different computing environments by synthetically modeling stochastic delays in the computing loop of workers. More

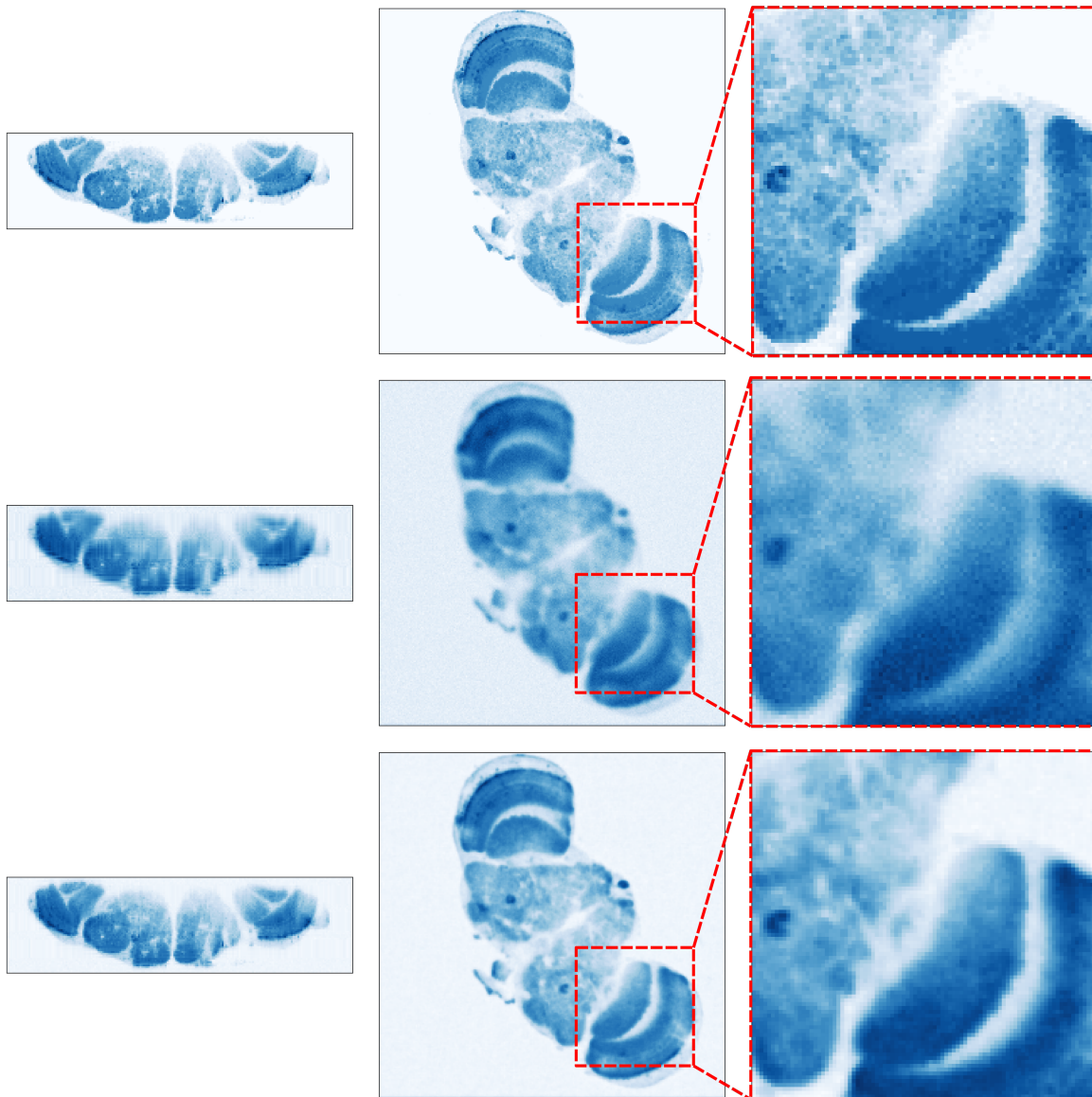


Figure 3: Restoration results of *Flybrain*: ground truth volume (top), degraded version (middle), and results of BD3MG restoration (bottom). Visual comparisons along the  $X - Z$  axis (left) the  $X - Y$  axis (middle) and zoomed details (right). The optimization process recovers fine details of the original volume that were lost in its degraded version.

specifically, the same restoration task and computer characteristics than in the previous section is considered, again with  $C_{\text{tot}} = 12$  active cores. We introduce artificial perturbation in the computing environment by randomly “freezing” some worker processes for a certain amount of time (i.e., delay) following the three scenarios below:

- **Type I:** One of the workers is consistently affected by a delay that follows a uniform distribution  $\mathcal{U}([0, 1])$  (in sec.). The other cores are not affected by any delay.
- **Type II:** Two worker cores are not affected by any delay while the others 9 agents

are delayed in the following fashion:

3 cores hold a delay following a uniform distribution  $\mathcal{U}([0, 1])$ .

3 cores hold a delay following a uniform distribution  $\mathcal{U}([0, 0.5])$ .

3 cores hold a delay following a uniform distribution  $\mathcal{U}([0, .25])$ .

- **Type III:** All worker cores are affected by a delay that follows a uniform distribution  $\mathcal{U}([0, 1])$ .

Method (Scenario)	SNR (dB)	$f(\mathbf{x}^f)$	Time (s.)
3MG (no delay)	18.1	1 247.0	1 683.8
BP3MG (Type I)	17.9	1 247.1	623.1
BD3MG (Type I)	18.7	1 246.0	<b>211.3</b>
BP3MG (Type II)	17.9	1 247.1	707.9
BD3MG (Type II)	18.7	1 246.0	<b>220.7</b>
BP3MG (Type III)	17.9	1 247.1	752.8
BD3MG (Type III)	18.7	1 246.0	<b>219.9</b>

Table 2: Performances of BP3MG and BD3MG under imbalanced computed power, for reaching the stopping criterion with  $\tilde{\varepsilon} = 5 \times 10^{-4}$  for Flybrain restoration. We additionally provide results for the vanilla 3MG algorithm for sake of comparison.

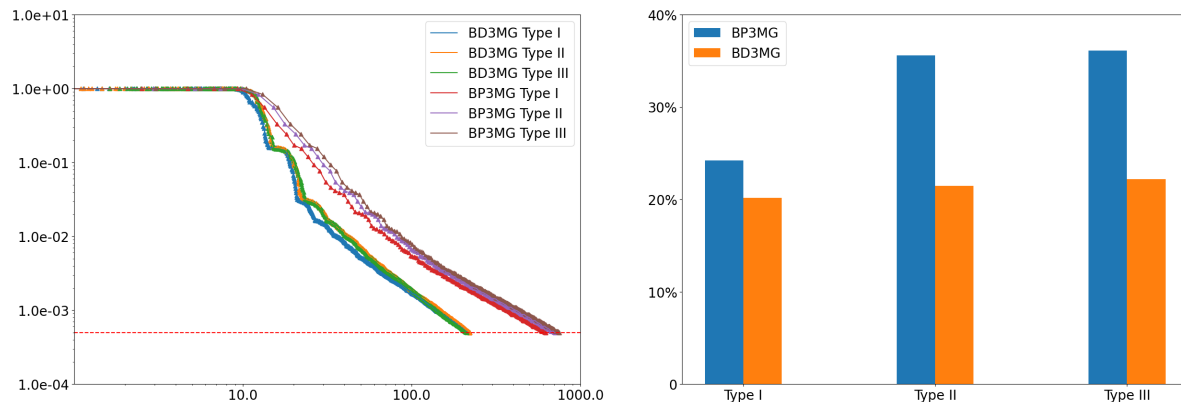


Figure 4: Numerical comparisons between BD3MG and BP3MG for FlyBrain restoration under imbalanced computing power: evolution of the relative increment  $\|\mathbf{x}^{k+1} - \mathbf{x}^k\| / \|\mathbf{x}^k\|$  along time (in sec.) for each of the three experimental settings in log-log scale (left), and averaged ratio of workers CPU idle time over the entire optimization process for each scenario (right).

The results are summarized in Table 2 and Figure 4. We also report the results of a plain, not delayed, 3MG implementation, for the sake of comparison. In all three scenarios, BD3MG outperforms its synchronous version BP3MG, in terms of computation time while reaching similar final criterion value and SNR. The criteria decrease is faster

for BD3MG which can be explained by two main differences with BP3MG. First, the majorant mapping of BD3MG performs a tighter approximation of the cost function than BP3MG, thus leading intrinsically to improved convergence rate. Second, BD3MG is asynchronous by essence and thus it is resilient to communication delays as soon as they are bounded, as shown by our convergence analysis. In contrast, BP3MG simply waits for all workers to finalize their update, to force the synchronicity, which yields slowdown in case of delayed workers. A more efficient and dynamic handling of the workload is performed in BD3MG, as shown in Figure 4 where CPU idle time is consistently lower for BD3MG than for BP3MG. We note that in asymmetric settings such as (Type II) and (Type III), BD3MG proved to be particularly efficient in reducing the synchronicity constraint of BP3MG for “fast” workers. The comparable results for BD3MG on all three scenarios further suggest that the proposed algorithm is robust to an imbalance in the computing power of workers. Moreover, despite the delayed feedbacks of the workers, both BP3MG and BD3MG remain competitive with respect to the vanilla 3MG, which shows the great interest of a parallel friendly algorithmic structure in this context.

#### 5.4. Scalability assessment.

In order to assess the scalability properties of BD3MG, we further analyse the speed-up generated by the number of cores available. We consider the restoration problem of the 3D image **Aneurysm**¶ of size  $N = N_x \times N_y \times N_z = 256 \times 256 \times 154$ , under the same degradation operator and noise level than in the previous example. Figure 5 presents the acceleration ratio between the required computation time for a single active worker versus the computation time of up to 30 active workers in reaching the stopping criterion with  $\tilde{\epsilon} = 10^{-3}$ . The regularization parameters are set empirically to  $(\lambda, \delta, \kappa, \eta) = (1, 1, 10^{-1}, 10^{-3})$  to maximize the final SNR and the same blur kernel than in the previous subsection is used. The computations were performed using HPC resources from the **Oscar** - Ocean State Center for Advanced Resources of the Center for Computation and Visualization, Brown University. The hardware is an Intel Corei9 CPU with up to 48 physical cores at 3.3 GHz and 300G of RAM. Results found in Figure 5 illustrate the great potential of scalability of the proposed algorithm. As the number of core increases, a mild saturation effect is observed (in agreement with Amdahl’s law [59]).

#### 5.5. Application to real data from multiphoton microscopy

We finally illustrate the performance of BD3MG on a restoration task of real multiphoton microscopy data specifically acquired for this experiment. Multiphoton microscopy is an interesting solution for the 3D and submicrometric characterization of biomedical structures, it is label-free and contactless [36]. Such a solution takes

¶ <http://www.mathworks.com/matlabcentral/fileexchange/25987-showvol-isosurface-render/>

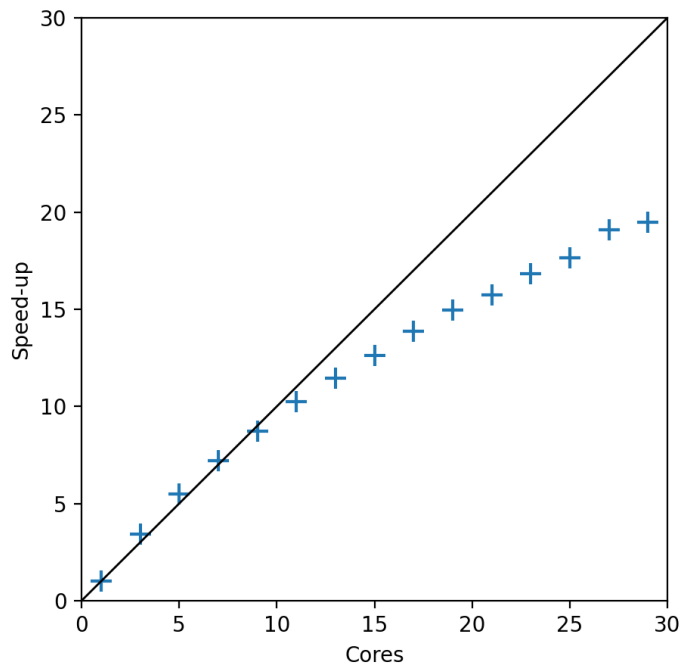


Figure 5: Speed-up ratio of the computation time for 1 to 30 cores for BD3MG for the restoration of **Aneurysm**.

advantage of optical sectioning, an optical property resulting from the nonlinear optical processes involved. 3D images are produced with sub-micrometer resolution without slicing the sample. We use an instrumental solution for acquisition rested on a commercial system from Olympus (BX61WI) coupled with a multiphoton water immersion objective (Olympus XLPLN25XWMP, 25 $\times$ , NA 1.05). A laser system, emitting femtosecond pulses centred at 810 nm with 10 nm of spectral bandwidth, is used for production of the nonlinear phenomena of second harmonic generation (SHG) and two-photon fluorescence (TPF). The biomedical sample is made of a whole mouse muscle, the Extensor digitorum longus (EDL), isolated from tendon to tendon. Sub-micrometric fluorescent microspheres emitting in the green range are included into the EDL and spread homogeneously all along the whole muscle structure. Under such an experimental protocol, the production of two 3D images is obtained. The first channel contains the SHG from the myosin of the muscle and the second channel displayed the TPF of microspheres used for calibrating the instrumental PSF. A hundred of 2D image slices of SHG and TPF are produced, with 0.1  $\mu\text{m}$  resolution along depth axis Z and 0.049 $\mu\text{m}$  $\times$ 0.049 $\mu\text{m}$  resolution over X – Y horizontal-vertical axis. The acquisition recording starts 140  $\mu\text{m}$  under the sample surface for a total sample thickness of 180  $\mu\text{m}$ . For this range of depth, the imaging of biological samples is degraded by scattering effects. Both raw volumes (i.e., SHG and TPF) dimension have 2048  $\times$  2048  $\times$  100 voxels, from which we extract a subpart with size  $N = N_X \times N_Y \times N_Z = 256 \times 256 \times 100$  voxels

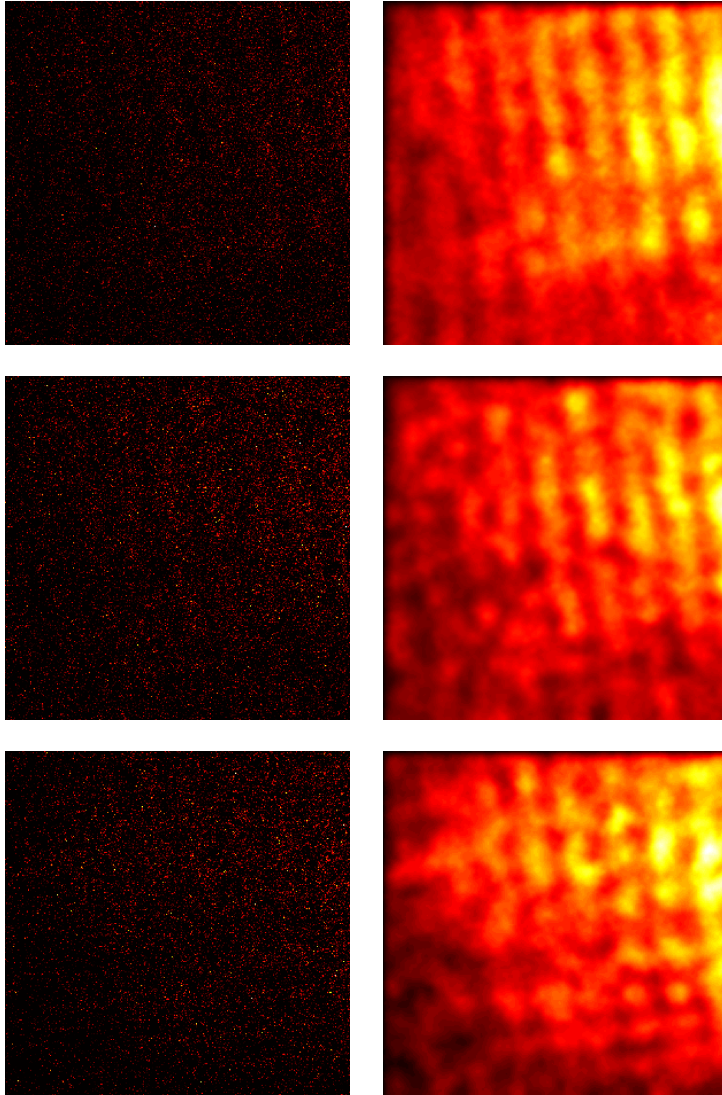


Figure 6: Slices ( $12, 5\mu m \times 12, 5\mu m$ ) for depths  $z = 5, 25$  and  $70$  (from top to bottom) of the original acquisition (left) and after restoration (right). The comparisons show that the definition of the muscular structure has been enhanced by the reconstruction.

for the purpose of our study.

We follow the computational pipeline FAMOUS previously introduced in [48]. We estimate a depth-variant Gaussian PSF field within the 3D microscopic volume by applying the 3D Gaussian fitting algorithm FIGARO from [22] to volume of interests extracted from the second image channel, displaying fluorescence of calibrated microbeads. Each volume of interest is selected through an automatic search of connected components within a filtered and binarized version of the observed volume. Then, FIGARO method is ran, yielding parameters (i.e., mean, covariance, scaling, shift) of a 3D Gaussian shape. This allows to build, through a simple interpolation strategy, a model for a depth-variant PSF with truncated support of size  $M =$

$M_X \times M_Y \times M_Z = 21 \times 21 \times 21$  (see more details in [48, Sec.2.4]). Since no ground truth is available, the regularization parameters  $(\lambda, \delta, \kappa, \eta) = (10^2, 2, 10, 10^{-3})$ , are selected by retrospective visual inspection. The reconstruction shown in Figure 6 exhibits clear contrasts and sharpness properties. Comparative videos of the original and restored volume are available at <https://github.com/mathieuchal/BD3MG>. The native signal from the raw image was presenting a high level of noise and blur due to the presence of scattering elements all along the 140  $\mu\text{m}$  of sample depth. Thanks to the proposed restoration strategy, the localisation of the myosin in the muscle sample is made possible, and the spatial organization of this protein into the down side of the EDL is revealed. The volume restoration took 305 seconds and  $\sim 2000$  iterations on a  $C_{\text{tot}} = 12$  cores setting, when setting  $\tilde{\varepsilon} = 10^{-3}$ .

## 6. Conclusion

In this paper, we have presented a new block distributed Majorize-Minimize algorithm, BD3MG, devised to tackle large-size differentiable optimization problems met in a wide range of applications. Our main contribution lies in a distributed asynchronous formulation that allows for delays in the current solution computed between workers, while securing convergence guarantees under mild assumptions. Our new algorithm BD3MG has been tested in the context of 3D image restoration with depth-variant blur. Experimental results underlined the speedup potential of this method and its concrete applicability in the field of fluorescence microscopy. Future work will be dedicated to extension to more general distributed graph topologies.

## Declarations

*Funding:* This research work received funding support from the European Research Council Starting Grant MAJORIS ERC-2019-STG-850925.

*Competing interests:* The authors have no relevant financial or non-financial interests to disclose.

## Bibliography

- [1] F. Abboud, M. Stamm, E. Chouzenoux, J.-C. Pesquet, and H. Talbot. Distributed algorithms for proximity operator computation with applications to video processing. *Digital Signal Processing*, 2022. To appear.
- [2] M. Allain, J. Idier, and Y. Goussard. On global and local convergence of half-quadratic algorithms. *IEEE Transactions on Image Processing*, 15(5):1130–1142, May 2006.
- [3] H. Attouch, J. Bolte, P. Redont, and A. Soubeyran. Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the Kurdyka-Lojasiewicz inequality. *Mathematics of Operations Research*, 35(2):438–457, 2010.
- [4] H. Bauschke and P. L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. Springer, 2011.
- [5] M. Bertero. *Introduction to Inverse Problems in Imaging*. CRC Press, 2017.

- [6] D. P. Bertsekas. Nonlinear programming. *Journal of the Operational Research Society*, 48(3):334–334, 1997.
- [7] J. Bochnak, M. Coste, and M.-F. Roy. *Real Algebraic Geometry*, volume 36. Springer Science & Business Media, 2013.
- [8] J. Bolte, S. Sabach, and M. Teboulle. Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Mathematical Programming*, 146(1):459–494, 2014.
- [9] S. Bonettini, M. Prato, and S. Rebegoldi. A block coordinate variable metric linesearch based proximal gradient method. *Computational Optimization and Applications*, 71(1):5–52, 2018.
- [10] S. Cadoni, E. Chouzenoux, J.-C. Pesquet, and C. Chaux. A block parallel majorize-minimize memory gradient algorithm. In *Proceedings of the 23rd IEEE International Conference on Image Processing (ICIP 2016)*, pages 3194–3198, Phoenix, AZ, 25–28 Sep. 2016.
- [11] L. Cannelli, F. Facchinei, G. Scutari, and V. Kungurtsev. Asynchronous optimization over graphs: Linear convergence under error bound conditions. *IEEE Transactions on Automatic Control*, 66(10):4604–4619, 2020.
- [12] A. Chakrabarti, T. Zickler, and W. T. Freeman. Analyzing spatially-varying blur. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2010)*, pages 2512–2519, San Francisco, CA, USA, 13–18 June 2010.
- [13] M. Chalvidal and E. Chouzenoux. Block distributed 3MG algorithm and its application to 3D image restoration. In *Proceedings of the IEEE International Conference on Image Processing (ICIP 2020)*, pages 938–942, Abu Dhabi, United Arab Emirates (virtual), 25–28 Oct. 2020.
- [14] F. Chorobura and I. Necoara. Random coordinate descent methods for nonseparable composite optimization. Technical report, 2022. <https://arxiv.org/abs/2203.14368>.
- [15] E. Chouzenoux and J.-B. Fest. SABRINA: a stochastic subspace majorization-minimization algorithm. Technical report, 2021. <https://optimization-online.org/2021/08/8573/>.
- [16] E. Chouzenoux, J. Idier, and S. Moussaoui. A majorize-minimize strategy for subspace optimization applied to image restoration. *IEEE Transactions on Image Processing*, 20(6):1517–1528, 2010.
- [17] E. Chouzenoux, A. Jezierska, J.-C. Pesquet, and H. Talbot. A majorize-minimize subspace approach for  $\ell_2 - \ell_0$  image regularization. *SIAM Journal on Imaging Sciences*, 6(1):563–591, Jan 2013.
- [18] E. Chouzenoux, S. Martin, and J.-C. Pesquet. A local MM subspace method for solving constrained variational problems in image recovery. *Journal of Mathematical Imaging and Vision*, 2022. to appear.
- [19] E. Chouzenoux and J.-C. Pesquet. Convergence rate analysis of the majorize-minimize subspace algorithm. *IEEE Signal Processing Letters*, 23(9):1284–1288, Sep. 2016.
- [20] E. Chouzenoux and J.-C. Pesquet. A stochastic majorize-minimize subspace algorithm for online penalized least squares estimation. *IEEE Transactions on Signal Processing*, 65(18):4770–4783, 2017.
- [21] E. Chouzenoux, J.-C. Pesquet, and A. Repetti. A block coordinate variable metric forward-backward algorithm. *Journal of Global Optimization*, pages 1–29, Feb. 2016.
- [22] E. Chouzenoux, T. Tsz-Kit Lau, C. Lefort, and J.-C. Pesquet. Optimal multivariate Gaussian fitting with applications to PSF modeling in two-photon microscopy imaging. *Journal of Mathematical Imaging and Vision*, 61(7):1037–1050, Sept. 2019.
- [23] J. Chung, S. Knepper, and J. G. Nagy. *Large-Scale Inverse Problems in Imaging*, pages 47–90. Springer New York, New York, NY, 2015.
- [24] D. Davis. The asynchronous PALM algorithm for nonsmooth nonconvex problems. Technical report, 2016. <https://arxiv.org/abs/1604.00526>.
- [25] D. Davis, M. Udell, and B. Edmunds. The sound of APALM clapping: Faster nonsmooth nonconvex optimization with stochastic asynchronous palm. In *Proceedings of the 30th International Conference on Neural Information Processing Systems (NIPS 2016)*, Barcelona, Spain, 5–10 Dec. 2016.



- [26] L. Denis, E. Thiébaud, F. Soulez, J. M. Becker, and R. Mourya. Fast approximations of shift-variant blur. *International Journal of Computer Vision*, 115:253–278, 2015.
- [27] L. Denis, E. Thiébaud, and F. Soulez. Fast model of space-variant blurring and its application to deconvolution in astronomy. In *2011 18th IEEE International Conference on Image Processing*, pages 2817–2820. IEEE, 2011.
- [28] L. Duval, E. Chouzenoux, A. Repetti, M. Q. Pham, and J.-C. Pesquet. Euclid in a taxicab: Sparse blind deconvolution with smoothed l1/l2 regularization. *IEEE Signal Processing Letters*, 22(5):539–543, May 2015.
- [29] P. Escande and P. Weiss. Sparse wavelet representations of spatially varying blurring operators. *SIAM Journal on Imaging Sciences*, 8:2976–3014, 2015.
- [30] T. Fan and T. Murphey. Majorization minimization methods for distributed pose graph optimization with convergence guarantees. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2020)*, pages 5058–5065, Las Vegas, USA, 2020.
- [31] J. A. Fessler. Grouped coordinate descent algorithms for robust edge-preserving image restoration. In T. J. Schulz, editor, *Image Reconstruction and Restoration II, International Society for Optics and Photonics*, volume 3170, pages 184–194, 1997.
- [32] J.-B. Fest and E. Chouzenoux. Convergence analysis of block majorize-minimize subspace approaches. Technical report, 2021. <https://optimization-online.org/2021/12/8711/>.
- [33] A. Florescu, E. Chouzenoux, J.-C. Pesquet, P. Ciuciu, and S. Ciochina. A majorize-minimize memory gradient method for complex-valued inverse problems. *Signal Processing*, 103:285–295, 2014.
- [34] D. Geman and G. Reynolds. Constrained restoration and the recovery of discontinuities. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 14(3):367–383, 1992.
- [35] D. Geman and C. Yang. Nonlinear image recovery with half-quadratic regularization. *IEEE transactions on Image Processing*, 4(7):932–946, 1995.
- [36] W. Göbel, B. M. Kampa, and F. Helmchen. Imaging cellular network dynamics in three dimensions using fast 3d laser scanning. *Nature methods*, 4(1):73–79, 2007.
- [37] D. Grishchenko, F. Iutzeler, J. Malick, and M.-R. Amini. Asynchronous distributed learning with sparse communications and identification. Technical report, 2018. <https://arxiv.org/abs/1812.03871>.
- [38] S. B. Hadj and L. Blanc-Féraud. Modeling and removing depth variant blur in 3d fluorescence microscopy. In *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 689–692. IEEE, 2012.
- [39] R. Hannah and W. Yin. On unbounded delays in asynchronous parallel fixed-point algorithms. *Journal of Scientific Computing*, 76(1):299–326, Dec 2017.
- [40] M. Hong, M. Razaviyayn, Z. Luo, and J. Pang. A unified algorithmic framework for block-structured optimization involving big data: With applications in machine learning and signal processing. *IEEE Signal Processing Magazine*, 33(1):57–77, Jan 2016.
- [41] S. Horváth, M. Sanjabi, L. Xiao, P. Richtárik, and M. Rabbat. Fedshuffle: Recipes for better use of local work in federated learning. Technical report, 2022. <https://arxiv.org/abs/2204.13169>.
- [42] D. R. Hunter and K. Lange. A tutorial on MM algorithms. *The American Statistician*, 58(1):30–37, 2004.
- [43] M. W. Jacobson and J. A. Fessler. An expanded theoretical treatment of iteration-dependent majorize-minimize algorithms. *IEEE Transactions on Image Processing*, 16(10):2411–2422, Oct. 2007.
- [44] A. Jezierska, H. Talbot, and J.C.Pesquet. Spatially variant psf modeling in confocal macroscopy. In *Proceedings of the 15th IEEE International Symposium on Biomedical Imaging (ISBI 2018)*, pages 489–492. Washington, DC, USA, 4-7 April 2018.
- [45] B. Kim and T. Naemura. Blind depth-variant deconvolution of 3d data in wide-field fluorescence microscopy. *Scientific reports*, 5(1):1–9, 2015.

- [46] J. Konečný, H. B. McMahan, D. Ramage, and P. Richtárik. Federated optimization: Distributed machine learning for on-device intelligence. Technical report, 2016. <https://arxiv.org/abs/1610.02527>.
- [47] D. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorization. In *Proceedings of the 13th International Conference on Neural Information Processing Systems (NIPS 2000)*, page 535541, Denver, Colorado, 2000.
- [48] C. Lefort, M. Chalvidal, A. Parente, V. Blanquet, H. Massias, L. Magnol, and E. Chouzenoux. FAMOUS: a fast instrumental and computational pipeline for multiphoton microscopy applied to 3d imaging of muscle ultrastructure. *Journal of Physics D: Applied Physics*, 54(27):274005, 2021.
- [49] Z. Li, Z. Dong, Z. Liang, and Z. Ding. Surrogate-based distributed optimisation for expensive black-box functions. *Automatica*, 125:109407, 2021.
- [50] X. Lian, Y. Huang, Y. Li, and J. Liu. Asynchronous parallel stochastic gradient for nonconvex optimization. Technical report, 2015. <https://arxiv.org/abs/1506.08272>.
- [51] N. Loizou and P. Richtárik. Revisiting randomized gossip algorithms: General framework, convergence rates and novel block and accelerated protocols. *IEEE Transactions on Information Theory*, 67(12):8300–8324, 2021.
- [52] K. Mishchenko, F. Iutzeler, and J. Malick. A distributed flexible delay-tolerant proximal gradient algorithm. *SIAM Journal on Optimization*, 30(1):933–959, 2020.
- [53] J. Nagy and D. O’Leary. Restoring images degraded by spatially variant blur. *SIAM Journal on Scientific Computing*, 19(4):1063–1082, 1998.
- [54] Y. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *Université catholique de Louvain, Center for Operations Research and Econometrics (CORE), CORE Discussion Papers*, 22, Jan. 2010.
- [55] M. Nikolova and M.-K. Ng. Analysis of half-quadratic minimization methods for signal and image recovery. *SIAM J. Sci. Comput.*, 27:937–966, 2005.
- [56] F. Niu, B. Recht, C. Re, and S. Wright. Hogwild: A lockfree approach to parallelizing stochastic gradient descent. In *Proceedings of the 25th Conference on Advances in Neural Information Processing Systems (NIPS 2011)*, pages 693–701, Granada, Spain, 12-17 Dec. 2011.
- [57] J. Nutini, M. Schmidt, I. H. Laradji, M. Friedlander, and H. Koepke. Coordinate descent converges faster with the Gauss-Southwell rule than random selection. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning (ICML 2015)*, volume 37, page 1632–1641, 2015.
- [58] A. Onose, R. E. Carrillo, A. Repetti, J. D. McEwen, J.-T. Thiran, J.-C. Pesquet, and Y. Wiaux. Scalable splitting algorithms for big-data interferometric imaging in the SKA era. *Monthly Notices of the Royal Astronomical Society*, 462(4):4314–4335, 2016.
- [59] D. A. Patterson, J. L. Hennessy, and D. Goldberg. *Computer architecture: a quantitative approach*, volume 2. Morgan Kaufmann San Mateo, CA, 1990.
- [60] J.-C. Pesquet and A. Repetti. A class of randomized primal-dual algorithms for distributed optimization. *Journal of Nonlinear and Convex Analysis*, 16(12):2353–2490, 2014.
- [61] M. Prato, A. La Camera, S. Bonettini, S. Rebegoldi, M. Bertero, and P. Boccacci. A blind deconvolution method for ground based telescopes and Fizeau interferometers. *New Astronomy*, 40:1–13, 2015.
- [62] C. Preza and J.-A. Conchello. Depth-variant maximum-likelihood restoration for three-dimensional fluorescence microscopy. *JOSA A*, 21(9):1593–1601, 2004.
- [63] P. Richtárik and M. Takáč. Parallel coordinate descent methods for big data optimization. *Mathematical Programming*, 156(1-2):433–484, Apr 2015.
- [64] M. C. Robini and Y. Zhu. Generic half-quadratic optimization for image reconstruction. *SIAM Journal on Imaging Sciences*, 8(3):1752–1797, 2015.
- [65] S. Sotthivirat and J. A. Fessler. Image recovery using partitioned-separable paraboloidal surrogate coordinate ascent algorithms. *IEEE Transactions on Signal Processing*, 11(3):306–317, 2002.

- [66] Y. Sun, P. Babu, and D. P. Palomar. Majorization-minimization algorithms in signal processing, communications, and machine learning. *IEEE Transactions on Signal Processing*, 65(3):794–816, 2016.
- [67] P. Thouvenin, N. Dobigeon, and J.-Y. Tournet. Partially asynchronous distributed unmixing of hyperspectral images. *IEEE Transactions on Geoscience and Remote Sensing*, 57(4), Apr. 2019.
- [68] P. Tseng. Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of Optimization, Theory and Applications*, 109(3):475–494, 2001.
- [69] J. Tuck, D. Hallac, and S. Boyd. Distributed majorization-minimization for Laplacian regularized problems. *IEEE/CAA Journal of Automatica Sinica*, 6(1):45–52, 2019.
- [70] J. Wang, Z. Charles, Z. Xu, G. Joshi, H. B. McMahan, M. Al-Shedivat, G. Andrew, S. Avestimehr, K. Daly, D. Data, et al. A field guide to federated optimization. Technical report, 2021. <https://arxiv.org/abs/2107.06917>.
- [71] A. Wilson. *Lyapunov Arguments in Optimization*. PhD thesis, University of California, Berkeley, 2018.
- [72] O. Wirjadi and T. Breuel. Approximate separable 3D anisotropic Gauss filter. In *Proceedings of the 12nd IEEE International Conference on Image Processing (ICIP 2005)*, volume 2, pages 149–52, Genoa, Italy, 11-14 Sept 2005.
- [73] T. Yang, X. Yi, J. Wu, Y. Yuan, D. Wu, Z. Meng, Y. Hong, H. Wang, Z. Lin, and K. H. Johansson. A survey of distributed optimization. *Annual Reviews in Control*, 47:278–305, 2019.
- [74] S. Zavriev and F. Kostyuk. Heavy-ball method in nonconvex optimization problems. *Computational Mathematics and Modeling*, 4(4):336–341, 1993.
- [75] R. Zhang and J. T. Kwok. Asynchronous distributed ADMM for consensus optimization. In *Proceedings of the 31st International Conference on International Conference on Machine Learning (ICML 2014)*, volume 32, 2014.
- [76] Z. Zhang, J. T. Kwok, and D.-Y. Yeung. Surrogate maximization/minimization algorithms and extensions. *Machine Learning*, 69:1–33, 2007.
- [77] Y. Zheng and Q. Liu. A review of distributed optimization: Problems, models and algorithms. *Neurocomputing*, 483:446–459, 2021.