A Computation and Communication Efficient Method for Distributed Nonconvex Problems in the Partial Participation Setting

Anonymous Author(s)

Affiliation Address email

Abstract

We present a new method that includes three key components of distributed optimization and federated learning: variance reduction of stochastic gradients, partial participation, and compressed communication. We prove that the new method has optimal oracle complexity and state-of-the-art communication complexity in the partial participation setting. Regardless of the communication compression feature, our method successfully combines variance reduction and partial participation: we get the optimal oracle complexity, never need the participation of all nodes, and do not require the bounded gradients (dissimilarity) assumption.

1 Introduction

Federated and distributed learning have become very popular in recent years (Konečný et al., 2016; McMahan et al., 2017). The current optimization tasks require much computational resources and machines. Such requirements emerge in machine learning, where massive datasets and computations are distributed between cluster nodes (Lin et al., 2017; Ramesh et al., 2021). In federated learning, nodes, represented by mobile phones, laptops, and desktops, do not send their data to a server due to privacy and their huge number (Ramaswamy et al., 2019), and the server remotely orchestrates the nodes and communicates with them to solve an optimization problem.

As in classical optimization tasks, one of the main current challenges is to find **computationally**

As in classical optimization tasks, one of the main current challenges is to find **computationally**efficient optimization algorithms. However, the nature of distributed problems induces many other

(Kairouz et al., 2021), including i) **partial participation** of nodes in algorithm steps: due to stragglers

(Li et al., 2020) or communication delays (Vogels et al., 2021), ii) **communication bottleneck**: even

if a node participates, it can be costly to transmit information to a server or other nodes (Alistarh

et al., 2017; Ramesh et al., 2021; Kairouz et al., 2021; Sapio et al., 2019; Narayanan et al., 2019). It

is necessary to develop a method that considers these problems.

24 **Optimization Problem**

Let us consider the nonconvex distributed optimization problem

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

where $f_i: \mathbb{R}^d \to \mathbb{R}$ is a smooth nonconvex function for all $i \in [n] := \{1, \dots, n\}$. The full information about function f_i is stored on i^{th} node. The communication between nodes is maintained in the parameters server fashion (Kairouz et al., 2021): we have a server that receives compressed

- information from nodes, updates a state, and broadcasts an updated model. Since we work in the
- nonconvex world, our goal is to find an ε -solution (ε -stationary point) of (1): a (possibly random) 30
- point $\widehat{x} \in \mathbb{R}^d$, such that $\mathbb{E}[\|\nabla f(\widehat{x})\|^2] \leq \varepsilon$. 31
- We consider three settings: 32
- 1. **Gradient Setting.** The i^{th} node has only access to the gradient $\nabla f_i : \mathbb{R}^d \to \mathbb{R}^d$ of function f_i . 33
- Moreover, the following assumptions for the functions f_i hold. 34
- **Assumption 1.** There exists $f^* \in \mathbb{R}$ such that $f(x) > f^*$ for all $x \in \mathbb{R}$. 35
- **Assumption 2.** The function f is L-smooth, i.e., $\|\nabla f(x) \nabla f(y)\| \le L \|x y\|$ for all $x, y \in \mathbb{R}^d$. 36
- **Assumption 3.** The functions f_i are L_i -smooth for all $i \in [n]$. Let us define $\widehat{L}^2 := \frac{1}{n} \sum_{i=1}^n L_i^2 \cdot \widehat{L}_i^2$
- 2. Finite-Sum Setting. The functions $\{f_i\}_{i=1}^n$ have the finite-sum form

$$f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x), \qquad \forall i \in [n],$$
(2)

- where $f_{ij}: \mathbb{R}^d \to \mathbb{R}$ is a smooth nonconvex function for all $j \in [m]$.
- We assume that Assumptions 1, 2 and 3 hold and the following assumption.
- **Assumption 4.** The function f_{ij} is L_{ij} -smooth for all $i \in [n], j \in [m]$. Let $L_{\max} :=$
- 42 $\max_{i \in [n], j \in [m]} L_{ij}$.
- 3. Stochastic Setting. The function f_i is an expectation of a stochastic function,

$$f_i(x) = \mathcal{E}_{\xi} \left[f_i(x; \xi) \right], \quad \forall i \in [n],$$
 (3)

- where $f_i: \mathbb{R}^d \times \Omega_{\xi} \to \mathbb{R}$. For a fixed $x \in \mathbb{R}$, $f_i(x;\xi)$ is a random variable over some distribution
- \mathcal{D}_i , and, for a fixed $\xi \in \Omega_{\xi}$, $f_i(x;\xi)$ is a smooth nonconvex function. The i^{th} node has only access to a stochastic gradients $\nabla f_i(\cdot;\xi_{ij})$ of the function f_i through the distribution \mathcal{D}_i , where ξ_{ij} is a sample
- from \mathcal{D}_i . We assume that Assumptions 1, 2 and 3 hold and the following assumptions.
- **Assumption 5.** For all $i \in [n]$ and for all $x \in \mathbb{R}^d$, the stochastic gradient $\nabla f_i(x;\xi)$ is unbiased 48
- and has bounded variance, i.e., $\mathbb{E}_{\xi} \left[\nabla f_i(x; \xi) \right] = \nabla f_i(x)$, and $\mathbb{E}_{\xi} \left[\| \nabla f_i(x; \xi) \nabla f_i(x) \|^2 \right] \leq \sigma^2$, 49
- where $\sigma^2 \geq 0$. 50
- **Assumption 6.** For all $i \in [n]$ and for all $x, y \in \mathbb{R}$, the stochastic gradient $\nabla f_i(x; \xi)$ satisfies the
- mean-squared smoothness property, i.e., $\mathrm{E}_{\xi}\left[\left\|\nabla f_{i}(x;\xi) \nabla f_{i}(y;\xi)\right\|^{2}\right] \leq L_{\sigma}^{2}\left\|x y\right\|^{2}$.
- We compare algorithms using the oracle complexity, i.e., the number of (stochastic) gradients that 53
- each node has to calculate to get ε -solution, and the communication complexity, i.e., the number of
- bits that each node has to send to the server to get ε -solution. 55

2.1 Unbiased Compressors

- We use the concept of unbiased compressors to alleviate the communication bottleneck. The unbiased 57 compressors quantize and/or sparsify vectors that the nodes send to the server. 58
- **Definition 1.** A stochastic mapping $\mathcal{C}: \mathbb{R}^d \to \mathbb{R}^d$ is an *unbiased compressor* if there exists $\omega \in \mathbb{R}$ 59
- 60

56

$$E[C(x)] = x, \qquad E[\|C(x) - x\|^2] \le \omega \|x\|^2,$$
 (4)

- for all $x \in \mathbb{R}^d$.
- We denote a set of stochastic mappings that satisfy Definition 1 as $\mathbb{U}(\omega)$. In our methods, the nodes
- make use of unbiased compressors $\{\mathcal{C}_i\}_{i=1}^n$. The community developed a large number of unbiassed

¹Note that this strategy can be used in peer-to-peer communication, assuming that the server is an abstraction and all its algorithmic steps are performed on each node.

²Note that $L \leq \widehat{L}$, $\widehat{L} \leq L_{\max}$, and $\widehat{L} \leq L_{\sigma}$.

Table 1: Summary of methods that solve the problem (1) in the stochastic setting (3). Abbr.: VR (Variance Reduction) = Does a method have the optimal oracle complexity $\mathcal{O}\left(\frac{\sigma^2}{\varepsilon} + \frac{\sigma}{\varepsilon^{3/2}}\right)$? PP (Partial Participation) = Does a method support partial participation from Section 2.2? CC = Does a method have the communication complexity equals to $\mathcal{O}\left(\frac{\omega}{\sqrt{n}\varepsilon}\right)$?

Method	VR	PP	CC	Limitations	
SPIDER, SARAH, PAGE, STORM (Fang et al., 2018; Nguyen et al., 2017) (Li et al., 2021a; Cutkosky and Orabona, 2019)	1	X	X	_	
MARINA (Gorbunov et al., 2021)	✓	X ^(a)	✓ ^(b)	Suboptimal convergence rate (see (Tyurin and Richtárik, 2023)).	
FedPAGE (Zhao et al., 2021b)	X	X ^(a)	X	Suboptimal oracle complexity $\mathcal{O}\left(\frac{\sigma^2}{\varepsilon^2}\right)$.	
FRECON (Zhao et al., 2021a)	Х	1	1	-	
FedAvg (McMahan et al., 2017; Karimireddy et al., 2020b)	Х	1	X	Bounded gradients (dissimilarity) assumption of f_i .	
SCAFFOLD (Karimireddy et al., 2020b)		1	X	Suboptimal convergence rate ^(e) .	
MIME (Veriminal day at al. 2020a) X(d) / X Bounded gradients (dissimilarity)		Calculates full gradient. Bounded gradients (dissimilarity) assumption of f_i . Suboptimal oracle compl. $\mathcal{O}\left(1/\varepsilon^{3/2}\right)$ in the setting (2).			
CE-LSGD (for Partial Participation) ^(c) (Patel et al., 2022) (concurrent work)	1	1	×	Bounded gradients (dissimilarity) assumption of f_i . Suboptimal oracle compl. $\mathcal{O}\left(1/\varepsilon^{3/2}\right)$ in the setting (2).	
DASHA		X	1		
(Tyurin and Richtárik, 2023)	X	or ✓	1	_	
DASHA-PP (new)	✓	✓	1	-	

⁽a) MARINA and FedPAGE, with a small probability, require the participation of all nodes so that they can not support partial participation from Section 2.2. Moreover, these methods provide suboptimal oracle complexities.

- compressors, including Rand K (see Definition 5) (Beznosikov et al., 2020; Stich et al., 2018),
- 65 Adaptive sparsification (Wangni et al., 2018) and Natural compression and dithering (Horváth et al.,
- 2019a). We are aware of correlated compressors by Szlendak et al. (2021) and quantizers by Suresh
- et al. (2022) that help in the homogeneous regimes, but in this work, we are mainly concentrated on
- generic heterogeneous regimes, though, for simplicity, assume the independence of the compressors.
- Assumption 7. $C_i \in \mathbb{U}(\omega)$ for all $i \in [n]$, and the compressors are statistically independent.

2.2 Nodes Partial Participation Assumptions

70

- We now try to formalize the notion of partial participation. Let us assume that we have n events $\{i^{\text{th}} \text{ node is } participating}\}$ with the following properties.
- Assumption 8. The partial participation of nodes has the following distribution: exists constants $p_a \in (0,1]$ and $p_{aa} \in [0,1]$, such that
- Prob (i^{th} node is participating) = $p_a \quad \forall i \in [n]$,
- 76 2. **Prob** $(i^{th} \text{ and } j^{th} \text{ nodes are } participating}) = p_{aa} \quad \forall i \neq j \in [n].$

⁽b) On average, MARINA provides the compressed communication mechanism with complexity $\mathcal{O}\left(\frac{\omega}{\sqrt{n}\varepsilon}\right)$. However, with a small probability, this method sends non-compressed vectors.

⁽c) Note that MIME and CE-LSGD can not be directly compared with DASHA-PP because MIME and CE-LSGD consider the online version of the problem (1), and require more strict assumptions.

⁽d) Although MIME obtains the convergence rate $\mathcal{O}\left(\frac{1}{\varepsilon^{3/2}}\right)$ of a variance reduced method, it requires the calculation of the full (exact) gradients.

⁽e) It can be seen when $\sigma^2=0$. Let us consider the s-nice sampling of the nodes, then SCAFFOLD requires $\mathcal{O}\left(\frac{n^{3/2}}{\varepsilon s^{3/2}}\right)$ communication rounds to get ε -solution, while DASHA-PP requires $\mathcal{O}\left(\frac{\sqrt{n}}{\varepsilon s}\right)$ communication rounds (see Theorem 4 with $\omega=0, b=\frac{p_a}{2-p_a}$, and $p_a=\frac{s}{n}$).

Table 2: Summary of methods that solve the problem (1) in the finite-sum setting (2). Abbr.: VR (Variance Reduction) = Does a method have the optimal oracle complexity $\mathcal{O}\left(m + \frac{\sqrt{m}}{\varepsilon}\right)$? PP and CC are defined in Table 1.

Method	VR	PP	CC	Limitations
SPIDER, PAGE (Fang et al., 2018; Li et al., 2021a)	/	Х	Х	_
MARINA (Gorbunov et al., 2021)	/	X ^(a)	✓ ^(b)	Suboptimal convergence rate (see (Tyurin and Richtárik, 2023)).
ZeroSARAH (Li et al., 2021b)	/	1	Х	Only homogeneous regime, i.e., the functions f_i are equal.
FedPAGE (Zhao et al., 2021b)	Х	X ^(a)	Х	Suboptimal oracle complexity $\mathcal{O}\left(\frac{m}{\varepsilon}\right)$.
DASHA (Tyurin and Richtárik, 2023)	√	X	1	_
DASHA-PP (new)	1	1	1	_

(a), (b) : see Table 1.

77
$$p_{aa} \le p_a^2$$
, (5)

and these events from different communication rounds are independent.

We are not fighting for the full generality and believe that more complex sampling strategies can be considered in the analysis. For simplicity, we settle upon Assumption 8. Standard partial participation strategies, including s-nice sampling, where the server chooses uniformly s nodes without replacement ($p_a = s/n$ and $p_{aa} = \frac{s(s-1)/n(n-1)}{n(n-1)}$), and independent participation, where each node independently participates with probability p_a (due to independence, we have $p_{aa} = p_a^2$), satisfy Assumption 8. In the literature, s-nice sampling is one of the most popular strategies (Zhao et al., 2021a; Richtárik et al., 2021; Reddi et al., 2020; Konečný et al., 2016).

3 Motivation and Related Work

The main goal of our paper is to develop a method for the nonconvex distributed optimization that will include three key features: variance reduction of stochastic gradients, compressed communication, and partial participation. We now provide an overview of the literature (see also Table 1 and Table 2).

1. Variance reduction of stochastic gradients

It is important to consider finite-sum (2) and stochastic (3) settings because, in machine learning tasks, either the number of local functions m is huge or the functions f_i is an expectation of a stochastic function due to the batch normalization (Ioffe and Szegedy, 2015) or random augmentation (Goodfellow et al., 2016), and it is infeasible to calculate the full gradients analytically. Let us recall the results from the nondistributed optimization. In the gradient setting, the optimal oracle complexity is $\mathcal{O}(1/\varepsilon)$, achieved by the vanilla gradient descent (GD) (Carmon et al., 2020; Nesterov, 2018). In the finite-sum setting and stochastic settings, the optimal oracle complexities are $\mathcal{O}\left(m + \frac{\sqrt{m}}{\varepsilon}\right)$ and $\mathcal{O}\left(\frac{\sigma^2}{\varepsilon} + \frac{\sigma}{\varepsilon^{3/2}}\right)$ (Fang et al., 2018; Li et al., 2021a; Arjevani et al., 2019), accordingly, achieved by methods SPIDER, SARAH, PAGE, and STORM from (Fang et al., 2018; Nguyen et al., 2017; Li et al., 2021a; Cutkosky and Orabona, 2019).

2. Compressed communication

In distributed optimization (Ramesh et al., 2021; Xu et al., 2021), lossy communication compression can be a powerful tool to increase the communication speed between the nodes and the server. Different types of compressors are considered in the literature, including unbiased compressors (Alistarh et al., 2017; Beznosikov et al., 2020; Szlendak et al., 2021), contractive (biased) compressors (Richtárik et al., 2021), 3PC compressors (Richtárik et al., 2022). We will focus on unbiased compressors because methods DASHA and MARINA (Tyurin and Richtárik, 2023; Szlendak et al.,

2021; Gorbunov et al., 2021) that employ unbiased compressors provide the current theoretical state-of-the-art (SOTA) communication complexities.

Many methods analyzed optimization methods with the unbiased compressors (Alistarh et al., 2017; Mishchenko et al., 2019; Horváth et al., 2019b; Gorbunov et al., 2021; Tyurin and Richtárik, 2023). In the gradient setting, the methods MARINA and DASHA by Gorbunov et al. (2021) and Tyurin and Richtárik (2023) establish the current SOTA communication complexity, each method needs $\frac{1+\omega/\sqrt{n}}{\varepsilon}$ communication rounds to get an ε -solution. In the finite-sum and stochastic settings, the current SOTA communication complexity is attained by the DASHA method, while maintaining the optimal oracle complexities $\mathcal{O}\left(m+\frac{\sqrt{m}}{\varepsilon\sqrt{n}}\right)$ and $\mathcal{O}\left(\frac{\sigma^2}{\varepsilon n}+\frac{\sigma}{\varepsilon^{3/2}n}\right)$ per node.

3. Partial participation

117

134

136

137

138

147

From the beginning of federated learning era, the partial participation has been considered to be 118 the essential feature of distributed optimization methods (McMahan et al., 2017; Konečný et al., 119 2016; Kairouz et al., 2021). However, previously proposed methods have limitations: i) methods 120 MARINA and FedPAGE from (Gorbunov et al., 2021; Zhao et al., 2021b) still require synchronization 121 of all nodes with a small probability. ii) in the stochastic settings, methods FedAvg, SCAFFOLD, and 122 FRECON with the partial participation mechanism (McMahan et al., 2017; Karimireddy et al., 2020b; Zhao et al., 2021a) provide results without variance reduction techniques from (Fang et al., 2018; Li et al., 2021a; Cutkosky and Orabona, 2019) and, therefore, get suboptimal oracle complexities. Note 125 that FRECON and DASHA reduce the variance only from compressors (in the partial participation and 126 stochastic setting). iii) in the finite-sum setting, the ZeroSARAH method by Li et al. (2021b) focuses 127 on the homogeneous regime only (the functions f_i are equal). iv) The MIME method by Karimireddy 128 et al. (2020a) and the CE-LSGD method (for Partial Participation) by the concurrent paper (Patel 129 et al., 2022) consider the online version of the problem (1). Therefore, MIME and CE-LSGD (for 130 Partial Participation) require stricter assumptions, including the bounded inter-client gradient variance 131 assumption. In the finite-sum setting (2), MIME and CE-LSGD obtain a suboptimal oracle complexity 132 $\mathcal{O}(1/\varepsilon^{3/2})$ while, in the full participation setting, it is possible to get the complexity $\mathcal{O}(1/\varepsilon)$. 133

4 Contributions

135 We propose a new method DASHA-PP for the nonconvex distributed optimization.

- As far as we know, this is the first method that includes three key ingredients of federated learning methods: *variance reduction of stochastic gradients, compressed communication, and partial participation.*
- Moreover, this is the first method that combines *variance reduction of stochastic gradients and* partial participation flawlessly: i) it gets the optimal oracle complexity ii) does not require the participation of all nodes iii) does not require the bounded gradients assumption of the functions f_i .
- We prove convergence rates and show that this method has *the optimal oracle complexity and the*state-of-the-art communication complexity in the partial participation setting. Moreover, in our work,
 we observe a nontrivial side-effect from mixing the variance reduction of stochastic gradients and
 partial participation. It is a general problem not related to our methods or analysis that we discuss in
 Section C.

5 Algorithm Description and Main Challenges Towards Partial Participation

We now present DASHA-PP (see Algorithm 1), a family of methods to solve the optimization problem (1). When we started investigating the problem, we took DASHA as a baseline method for two reasons: the family of algorithms DASHA provides the current state-of-the-art communication complexities in the *non-partial participation* setting, and, unlike MARINA, it does not send non-compressed gradients and does not synchronize all nodes. Let us briefly discuss the main idea of DASHA, its problem in the *partial participation* setting, and why the refinement of DASHA is not an exercise.

In fact, the original DASHA method supports the partial participation of nodes *in the gradient setting*.

Since the nodes only do the following steps (see full algorithm in Algorithm 6):

$$g_i^{t+1} = g_i^t + C_i \left(\nabla f_i(x^{t+1}) - (1-a) \nabla f_i(x^t) - a g_i^t \right).$$

Algorithm 1 DASHA-PP

```
1: Input: starting point x^0 \in \mathbb{R}^d, stepsize \gamma > 0, momentum a \in (0,1], momentum b \in
                (0,1], probability p_{\text{page}} \in (0,1] (only in DASHA-PP-PAGE), batch size B (only in DASHA-PP-
                PAGE, DASHA-PP-FINITE-MVR and DASHA-PP-MVR), probability p_a \in (0,1] that a node is
  PAGE, DASHA-FF-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-INITE-IN
                           for i = 1, \ldots, n in parallel do
   7:
                                     if i^{th} node is participating (a) then
   8:
                                             Calculate k_i^{t+1} using Algorithm 2, 3, 4 or 5 h_i^{t+1} = h_i^t + \frac{1}{p_a} k_i^{t+1}
   9:
                                            m_i^{t+1} = \mathcal{C}_i \left( \frac{1}{p_\mathrm{a}} k_i^{t+1} - \frac{a}{p_\mathrm{a}} \left( g_i^t - h_i^t \right) \right)
g_i^{t+1} = g_i^t + m_i^{t+1}
Send m_i^{t+1} to the server
13:
                                  else h_{ij}^{t+1} = h_{ij}^t \text{ (only in DASHA-PP-FINITE-MVR)}  h_{i}^{t+1} = h_i^t, \quad g_i^{t+1} = g_i^t, \quad m_i^{t+1} = 0
16:
17:
                        end for g^{t+1} = g^t + \tfrac{1}{n} \textstyle \sum_{i=1}^n m_i^{t+1}
18:
21: Output: \hat{x}^T chosen uniformly at random from \{x^t\}_{k=0}^{T-1}
                (a): For the formal description see Section 2.2.
```

Algorithm 2 Calculate k_i^{t+1} for DASHA-PP in the gradient setting. See line 9 in Alg. 1

```
1: k_i^{t+1} = \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t))
```

Algorithm 3 Calculate k_i^{t+1} for DASHA-PP-PAGE in the finite-sum setting. See line 9 in Alg. 1

```
1: \mbox{ Generate a random set } I_i^t \mbox{ of size } B \mbox{ from } [m] \mbox{ with replacement } \\ 2: \mbox{ } k_i^{t+1} = \begin{cases} \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - \frac{b}{p_{\rm page}} \left( h_i^t - \nabla f_i(x^t) \right), \\ \mbox{ with probability } p_{\rm page} \mbox{ on all } participating \mbox{ nodes,} \\ \frac{1}{B} \sum_{j \in I_i^t} \left( \nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^t) \right), \\ \mbox{ with probability } 1 - p_{\rm page} \mbox{ on all } participating \mbox{ nodes} \end{cases}
```

Algorithm 4 Calc. k_i^{t+1} for DASHA-PP-FINITE-MVR in the finite-sum setting. See line 9 in Alg. 1

```
1: Generate a random set I_i^t of size B from [m] without replacement 

2: k_{ij}^{t+1} = \begin{cases} \frac{m}{B} \left( \nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^t) - b \left( h_{ij}^t - \nabla f_{ij}(x^t) \right) \right), & j \in I_i^t, \\ 0, & j \not\in I_i^t, \end{cases}
3: h_{ij}^{t+1} = h_{ij}^t + \frac{1}{p_a} k_{ij}^{t+1}
4: k_i^{t+1} = \frac{1}{m} \sum_{j=1}^m k_{ij}^{t+1}
```

Algorithm 5 Calculate k_i^{t+1} for DASHA-PP-MVR in the stochastic setting. See line 9 in Alg. 1

```
1: Generate i.i.d. samples \{\xi_{ij}^{t+1}\}_{j=1}^{B} of size B from \mathcal{D}_{i}.

2: k_{i}^{t+1} = \frac{1}{B} \sum_{j=1}^{B} \nabla f_{i}(x^{t+1}; \xi_{ij}^{t+1}) - \frac{1}{B} \sum_{j=1}^{B} \nabla f_{i}(x^{t}; \xi_{ij}^{t+1}) - b \left( h_{i}^{t} - \frac{1}{B} \sum_{j=1}^{B} \nabla f_{i}(x^{t}; \xi_{ij}^{t+1}) \right)
```

The partial participation mechanism (independent participation from Section 2.2) can be easily implemented here if we temporally redefine the compressor and use another one instead:

$$\mathcal{C}_i^p := \begin{cases} \frac{1}{p}\mathcal{C}_i, & \text{with pr. } p_{\mathbf{a}}, \\ 0, & \text{with pr. } 1-p_{\mathbf{a}}. \end{cases} \Rightarrow g_i^{t+1} = \begin{cases} g_i^t + \frac{1}{p_{\mathbf{a}}}\mathcal{C}_i\left(\nabla f_i(x^{t+1}) - (1-a)\nabla f_i(x^t) - ag_i^t\right), p_{\mathbf{a}} \\ g_i^t, & \text{with pr. } 1-p_{\mathbf{a}}. \end{cases}$$

With probability 1-p, a node does not update g_i and does not send anything to the server. The main observation is that we can do this trick since g_i^{t+1} depends only on the vectors x^{t+1} , x^t , and g_i^t .

However, we focus our attention on partial participation *in the finite-sum and stochastic settings*. Consider the nodes' steps in DASHA-MVR (see Algorithm 7) that is designed for the stochastic setting:

$$h_i^{t+1} = \nabla f_i(x^{t+1}; \xi_i^{t+1}) + (1-b)(h_i^t - \nabla f_i(x^t; \xi_i^{t+1})), \tag{6}$$

$$g_i^{t+1} = g_i^t + C_i \left(h_i^{t+1} - h_i^t - a \left(g_i^t - h_i^t \right) \right). \tag{7}$$

Even if we use the same trick for (7), we still have to update (6) in every iteration of the algorithm since g_i^{t+1} additionally depends on h_i^{t+1} and h_i^t . In other words, if a node does not update g_i and does not send anything to the server, it still has to update h_i , what is impossible without the points x^{t+1} and x^t . One of the main challenges was to "guess" how to generalize (6) and (7) to the partial participation setting. We now provide a solution (DASHA-PP-MVR with the batch size B=1):

$$\begin{split} h_i^{t+1} &= h_i^t + \frac{1}{p_{\rm a}} k_i^{t+1}, \ k_i^{t+1} = \nabla f_i(x^{t+1}; \xi_i^{t+1}) - \nabla f_i(x^t; \xi_i^{t+1}) - b \left(h_i^t - \nabla f_i(x^t; \xi_i^{t+1}) \right), \\ g_i^{t+1} &= g_i^t + \mathcal{C}_i \left(\frac{1}{p_{\rm a}} k_i^{t+1} - \frac{a}{p_{\rm a}} \left(g_i^t - h_i^t \right) \right) \text{ with } p_{\rm a}, \text{ and } h_i^{t+1} = h_i^t, \ g_i^{t+1} = g_i^t \text{ with } 1 - p_{\rm a}. \end{split} \tag{8}$$

Now both control variables g_i^t and h_i^t do not change with the probability $1-p_{\rm a}$. However, when the $i^{\rm th}$ node participates, the update rule of g_i^{t+1} and h_i^{t+1} was modified to make the proof work. When $p_{\rm a}=1$ (no partial participation), the updates rule from (8) reduce to (6) and (7).

The theoretical analysis of the new algorithm became more complicated: unlike (6) and (7), the control variables h_i^{t+1} and g_i^{t+1} in (8) (see also main Algorithm 1) are coupled by the randomness from the partial participation. Going deeper into details, one can compare Lemma I.2 from (Tyurin and Richtárik, 2023) and Lemma 5. The former lemma does not use the knowledge about the 171 172 173 update rule of h_i^{t+1} , works with one expectation $E_{\mathcal{C}}[\cdot]$, uses only (4), (12), and (13). The latter 174 lemma additionally requires and uses the structure of the update rule of h_i^{t+1} , surgically copes with the expectations $E_{\mathcal{C}}\left[\cdot\right]$ and $E_{p_a}\left[\cdot\right]$ (for instance, it is not trivial in each order one should apply the 175 176 expectations), and uses the sampling lemma (Lemma 1). The same reasoning applies to other parts of 177 the analysis and the finite-sum setting: the generalization of the previous algorithm and the additional 178 randomness from the partial participation required us to rethink the previous proofs. 179

At the first reading of the proofs, we suggest the reader follow the proof of Theorem 2 in the gradient setting, which takes a small part of the paper. Although the appendix seems to be dense and large, the size is justified by the fact that we consider four different settings and PŁ-condition $(4 \times 2 \text{ tracks})$ of proofs. The theory is designed so that the proofs do not repeat steps of each other and use one framework).

6 Theorems

185

186

187

188

189

190

191

192

We now present the convergence rates theorems of DASHA-PP in different settings. We will compare the theorems with the results of the current state-of-the-art methods, MARINA and DASHA, that work in the full participation setting. Suppose that MARINA or DASHA converges to ε -solution after T communication rounds. Then, ideally, we would expect the convergence of the new algorithms to ε -solution after up to T/p_a communication rounds due to the partial participation constraints³. The detailed analysis of the algorithms under Polyak-Łojasiewicz condition we provide in Section F. Let us define $\Delta_0 := f(x^0) - f^*$.

³We check this numerically in Section A.

6.1 Gradient Setting 193

Theorem 2. Suppose that Assumptions 1, 2, 3, 7 and 8 hold. Let us take $a = \frac{p_a}{2\omega+1}$, $b = \frac{p_a}{2-p_b}$, 194

$$\gamma \leq \left(L + \left[\frac{48\omega\left(2\omega + 1\right)}{np_{\mathrm{a}}^2} + \frac{16}{np_{\mathrm{a}}^2}\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)\right]^{1/2}\widehat{L}\right)^{-1},$$

- and $g_i^0 = h_i^0 = \nabla f_i(x^0)$ for all $i \in [n]$ in Algorithm 1 (DASHA-PP), then $\mathbb{E}\left[\left\|\nabla f(\widehat{x}^T)\right\|^2\right] \leq \frac{2\Delta_0}{\gamma T}$. 195
- Let us recall the convergence rate of MARINA or DASHA, the number of communication rounds to get 196
- ε -solution equals $\mathcal{O}\left(\frac{\Delta_0}{\varepsilon}\left[L+\frac{\omega}{\sqrt{n}}\widehat{L}\right]\right)$, while the rate of DASHA-PP equals $\mathcal{O}\left(\frac{\Delta_0}{\varepsilon}\left[L+\frac{\omega+1}{p_a\sqrt{n}}\widehat{L}\right]\right)$. Up to Lipschitz constants factors, we get the degeneration up to $^{1}\!/p_a$ factor due to the partial 197
- 198
- participation. This is the expected result since each worker sends useful information only with 199
- the probability p_a . 200

6.2 Finite-Sum Setting

Theorem 3. Suppose that Assumptions 1, 2, 3, 4, 7, and 8 hold. Let us take $a = \frac{p_a}{2\omega+1}$, $b = \frac{p_{page}p_a}{2-p_a}$, 202 probability $p_{page} \in (0,1]$, 203

$$\gamma \leq \left(L + \left[\frac{48\omega(2\omega + 1)}{np_{\rm a}^2}\left(\hat{L}^2 + \frac{(1 - p_{page})L_{\rm max}^2}{B}\right) + \frac{16}{np_{\rm a}^2p_{page}}\left(\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)\hat{L}^2 + \frac{(1 - p_{page})L_{\rm max}^2}{B}\right)\right]^{1/2}\right)^{-1}$$

- and $g_i^0 = h_i^0 = \nabla f_i(x^0)$ for all $i \in [n]$ in Algorithm 1 (DASHA-PP-PAGE) then $\mathbb{E}\left[\left\|\nabla f(\widehat{x}^T)\right\|^2\right] \leq$

201

- We now choose p_{page} to balance heavy full gradient and light mini-batch calculations. Let us define 206
- $\mathbb{1}_{p_a} := \sqrt{1 \frac{p_{aa}}{p_a}} \in [0, 1]$. Note that if $p_a = 1$ then $p_{aa} = 1$ and $\mathbb{1}_{p_a} = 0$. 207
- **Corollary 1.** Let the assumptions from Theorem 3 hold and $p_{page} = B/(m+B)$. Then DASHA-PP-PAGE 208 needs 209

$$T := \mathcal{O}\left(\frac{\Delta_0}{\varepsilon} \left[L + \frac{\omega}{p_a \sqrt{n}} \left(\widehat{L} + \frac{L_{\text{max}}}{\sqrt{B}} \right) + \frac{1}{p_a} \sqrt{\frac{m}{n}} \left(\frac{\mathbb{1}_{p_a} \widehat{L}}{\sqrt{B}} + \frac{L_{\text{max}}}{B} \right) \right] \right) \tag{9}$$

- communication rounds to get an ε -solution and the expected number of gradient calculations per 210 node equals $\mathcal{O}(m+BT)$. 211
- The convergence rate the rate of the current state-of-the-art method DASHA-PAGE without partial 212
- participation equals $\mathcal{O}\left(\frac{\Delta_0}{\varepsilon}\left[L+\frac{\omega}{\sqrt{n}}\left(\widehat{L}+\frac{L_{\max}}{\sqrt{B}}\right)+\sqrt{\frac{m}{n}}\frac{L_{\max}}{B}\right]\right)$. Let us closer compare it with (9). As expected, we see that the second term w.r.t. ω degenerates up to $1/p_a$. Surprisingly, the third term 213
- 214
- w.r.t. $\sqrt{m/n}$ can degenerate up to \sqrt{B}/p_a when $\widehat{L}\approx L_{\rm max}$. Hence, in order to keep degeneration up to 215
- $1/p_a$, one should take the batch size $B=\mathcal{O}\left(L_{\max}^2/\widehat{L}^2\right)$. This interesting effect we analyze separately 216
- in Section C. The fact that the degeneration is up to $1/p_a$ we check numerically in Section A. 217
- In the following corollary, we consider Rand K compressors (see Definition 5) and show that with 218
- the particular choice of parameters, up to the Lipschitz constants factors, DASHA-PP-PAGE gets the 219
- optimal oracle complexity and SOTA communication complexity. Indeed, comparing the following 220
- result with (Tyurin and Richtárik, 2023)[Corollary 6.6], one can see that we get the degeneration 221
- up to $1/p_a$ factor, which is expected in the partial participation setting. Note that the complexities 222
- improve with the number of workers n.
- **Corollary 2.** Suppose that assumptions of Corollary 1 hold, $B \leq \min \left\{ \frac{1}{p_a} \sqrt{\frac{m}{n}}, \frac{L_{\max}^2}{1^2 p_a \hat{L}^2} \right\}^5$, and we
- use the unbiased compressor RandK with $K = \Theta(Bd/\sqrt{m})$. Then the communication complexity of 225

If
$$\mathbbm{1}_{p_a} = 0$$
, then $\frac{L_{\sigma}^2}{\mathbbm{1}_{p_a}^2 \hat{L}^2} = +\infty$

⁴The choice of the compressor is driven by simplicity, and the following analysis can be used for other unbiased compressors

Algorithm 1 is $\mathcal{O}\left(d + \frac{L_{\max}\Delta_0 d}{p_a \varepsilon \sqrt{n}}\right)$, and the expected number of gradient calculations per node equals $\mathcal{O}\left(m + \frac{L_{\max}\Delta_0 \sqrt{m}}{p_a \varepsilon \sqrt{n}}\right)$.

227
$$\mathcal{O}\left(m + \frac{L_{\max}\Delta_0\sqrt{m}}{p_{\mathrm{a}}\varepsilon\sqrt{n}}\right)$$

The convergence rate of DASHA-PP-FINITE-MVR is provided in Section E.5. 228

6.3 Stochastic Setting

- We define $h^t := \frac{1}{n} \sum_{i=1}^n h_i^t$. 230

Theorem 4. Suppose that Assumptions 1, 2, 3, 5, 6, 7 and 8 hold. Let us take
$$a = \frac{p_a}{2\omega + 1}$$
, 232 $b \in \left(0, \frac{p_a}{2-p_a}\right]$, $\gamma \le \left(L + \left[\frac{48\omega(2\omega + 1)}{np_a^2}\left(\hat{L}^2 + \frac{(1-b)^2L_\sigma^2}{B}\right) + \frac{12}{np_ab}\left(\left(1 - \frac{p_{aa}}{p_a}\right)\hat{L}^2 + \frac{(1-b)^2L_\sigma^2}{B}\right)\right]^{1/2}\right)^{-1}$, and 233 $g_i^0 = h_i^0$ for all $i \in [n]$ in Algorithm 1 (DASHA-PP-MVR). Then

$$\begin{split} & \mathbf{E}\left[\left\|\nabla f(\hat{x}^{T})\right\|^{2}\right] \leq \frac{1}{T}\left[\frac{2\Delta_{0}}{\gamma} + \frac{2}{b}\left\|h^{0} - \nabla f(x^{0})\right\|^{2} + \left(\frac{32b\omega(2\omega+1)}{np_{\mathrm{a}}^{2}} + \frac{4\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)}{np_{\mathrm{a}}}\right)\left(\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{0} - \nabla f_{i}(x^{0})\right\|^{2}\right)\right] \\ & + \left(\frac{48b^{2}\omega(2\omega+1)}{p_{i}^{2}} + \frac{12b}{p_{\mathrm{a}}}\right)\frac{\sigma^{2}}{nB}. \end{split}$$

- In the next corollary, we choose momentum b and initialize vectors h_i^0 to get ε -solution.
- Corollary 3. Suppose that assumptions from Theorem 4 hold, momentum b 235
- $\Theta\left(\min\left\{\frac{p_{\mathrm{a}}}{\omega}\sqrt{\frac{n\varepsilon B}{\sigma^{2}}},\frac{p_{\mathrm{a}}n\varepsilon B}{\sigma^{2}}\right\}\right),\ \frac{\sigma^{2}}{n\varepsilon B}\ \geq\ 1,\ \textit{and}\ \ h_{i}^{0}\ =\ \frac{1}{B_{\mathrm{init}}}\sum_{k=1}^{B_{\mathrm{init}}}\nabla f_{i}(x^{0};\xi_{ik}^{0})\ \textit{for all}\ i\ \in\ [n],$
- and batch size $B_{\text{init}} = \Theta\left(\frac{\sqrt[c]{p_a}B}{b}\right)$, then Algorithm I (DASHA-PP-MVR) needs

$$T := \mathcal{O} \Bigg(\frac{\Delta_0}{\varepsilon} \Bigg[L + \frac{\omega}{p_{\rm a} \sqrt{n}} \left(\widehat{L} + \frac{L_\sigma}{\sqrt{B}} \right) + \frac{\sigma}{p_{\rm a} \sqrt{\varepsilon} n} \left(\frac{\mathbbm{1}_{p_{\rm a}} \widehat{L}}{\sqrt{B}} + \frac{L_\sigma}{B} \right) \Bigg] + \frac{\sigma^2}{\sqrt{p_{\rm a}} n \varepsilon B} \Bigg)$$

- communication rounds to get an ε -solution and the number of stochastic gradient calculations per 238 node equals $\mathcal{O}(B_{\text{init}} + BT)$. 239
- The convergence rate of the DASHA-SYNC-MVR, the state-of-the-art method without partial participa-
- tion, equals $\mathcal{O}\left(\frac{\Delta_0}{\varepsilon}\left[L + \frac{\omega}{\sqrt{n}}\left(\widehat{L} + \frac{L_{\sigma}}{\sqrt{B}}\right) + \frac{\sigma}{\sqrt{\varepsilon n}}\frac{L_{\sigma}}{B}\right] + \frac{\sigma^2}{n\varepsilon B}\right)$. Similar to Section 6.2, we see that in 241
- the regimes when $\widehat{L} \approx L_{\sigma}$ the third term w.r.t. $1/\varepsilon^{3/2}$ can degenerate up to \sqrt{B}/p_a . However, if we take
- $B = \mathcal{O}\left(L_{\sigma}^{2}/\hat{L}^{2}\right)$, then the degeneration of the third term will be up to $1/p_{a}$. This effect we analyze in 243
- Section C. The fact that the degeneration is up to $1/p_a$ we check numerically in Section A.
- In the following corollary, we consider Rand K compressors (see Definition 5) and show that with 245
- the particular choice of parameters, up to the Lipschitz constants factors, DASHA-PP-MVR gets the 246
- optimal oracle complexity and SOTA communication complexity of DASHA-SYNC-MVR method. 247
- Indeed, comparing the following result with (Tyurin and Richtárik, 2023)[Corollary 6.9], one can see 248
- that we get the degeneration up to $1/p_a$ factor, which is expected in the partial participation setting. 249
- Note that the complexities improve with the number of workers n.
- **Corollary 4.** Suppose that assumptions of Corollary 3 hold, batch size $B \leq \min \left\{ \frac{\sigma}{p_a \sqrt{\varepsilon} n}, \frac{L_{\sigma}^2}{1\frac{\sigma}{2} \hat{L}^2} \right\}$, 251
- we take RandK compressors with $K=\Theta\left(\frac{Bd\sqrt{\varepsilon n}}{\sigma}\right)$. Then the communication complexity equals
- $\mathcal{O}\left(\frac{d\sigma}{\sqrt{p_a}\sqrt{n\varepsilon}} + \frac{L_\sigma\Delta_0d}{p_a\sqrt{n\varepsilon}}\right)$, and the expected number of stochastic gradient calculations per node equals $\mathcal{O}\left(\frac{\sigma^2}{\sqrt{p_a}n\varepsilon} + \frac{L_\sigma\Delta_0\sigma}{p_e\varepsilon^{3/2}n}\right)$.
- We are aware that the initial batch size $B_{\rm init}$ can be suboptimal w.r.t. ω in DASHA-PP-MVR in some 255
- regimes (see also (Tyurin and Richtárik, 2023)). This is a side effect of mixing the variance reduction 256
- of stochastic gradients and compression. However, Corollary 4 reveals that we can escape these 257
- regimes by choosing the parameter K of RandK compressors in a particular way. To get the complete 258
- picture, we analyze the same phenomenon under PŁ condition (see Section F) and provide a new
- method DASHA-PP-SYNC-MVR (see Section G).

References

- Alistarh, D., Grubic, D., Li, J., Tomioka, R., and Vojnovic, M. (2017). QSGD: Communicationefficient SGD via gradient quantization and encoding. In *Advances in Neural Information Process*ing Systems (NIPS), pages 1709–1720.
- Arjevani, Y., Carmon, Y., Duchi, J. C., Foster, D. J., Srebro, N., and Woodworth, B. (2019). Lower bounds for non-convex stochastic optimization. *arXiv preprint arXiv:1912.02365*.
- Beznosikov, A., Horváth, S., Richtárik, P., and Safaryan, M. (2020). On biased compression for distributed learning. *arXiv preprint arXiv:2002.12410*.
- Carmon, Y., Duchi, J. C., Hinder, O., and Sidford, A. (2020). Lower bounds for finding stationary
 points i. *Mathematical Programming*, 184(1):71–120.
- Chang, C.-C. and Lin, C.-J. (2011). LIBSVM: a library for support vector machines. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 2(3):1–27.
- Cutkosky, A. and Orabona, F. (2019). Momentum-based variance reduction in non-convex SGD.
 arXiv preprint arXiv:1905.10018.
- Fang, C., Li, C. J., Lin, Z., and Zhang, T. (2018). SPIDER: Near-optimal non-convex optimization via stochastic path integrated differential estimator. In *NeurIPS Information Processing Systems*.
- Goodfellow, I., Bengio, Y., Courville, A., and Bengio, Y. (2016). *Deep learning*, volume 1. MIT Press.
- Gorbunov, E., Burlachenko, K., Li, Z., and Richtárik, P. (2021). MARINA: Faster non-convex distributed learning with compression. In *38th International Conference on Machine Learning*.
- Horváth, S., Ho, C.-Y., Horvath, L., Sahu, A. N., Canini, M., and Richtárik, P. (2019a). Natural compression for distributed deep learning. *arXiv preprint arXiv:1905.10988*.
- Horváth, S., Kovalev, D., Mishchenko, K., Stich, S., and Richtárik, P. (2019b). Stochastic distributed learning with gradient quantization and variance reduction. *arXiv preprint arXiv:1904.05115*.
- Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by
 reducing internal covariate shift. In *International Conference on Machine Learning*, pages 448–456. PMLR.
- Kairouz, P., McMahan, H. B., Avent, B., Bellet, A., Bennis, M., Bhagoji, A. N., Bonawitz, K.,
 Charles, Z., Cormode, G., Cummings, R., et al. (2021). Advances and open problems in federated
 learning. Foundations and Trends® in Machine Learning, 14(1–2):1–210.
- Karimireddy, S. P., Jaggi, M., Kale, S., Mohri, M., Reddi, S. J., Stich, S. U., and Suresh, A. T. (2020a). Mime: Mimicking centralized stochastic algorithms in federated learning. *arXiv preprint arXiv:2008.03606*.
- Karimireddy, S. P., Kale, S., Mohri, M., Reddi, S., Stich, S., and Suresh, A. T. (2020b). Scaffold:
 Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pages 5132–5143. PMLR.
- Konečný, J., McMahan, H. B., Yu, F. X., Richtárik, P., Suresh, A. T., and Bacon, D. (2016). Federated learning: Strategies for improving communication efficiency. *arXiv preprint arXiv:1610.05492*.
- Li, T., Sahu, A. K., Zaheer, M., Sanjabi, M., Talwalkar, A., and Smith, V. (2020). Federated optimization in heterogeneous networks. *Proceedings of Machine Learning and Systems*, 2:429–450.
- Li, Z., Bao, H., Zhang, X., and Richtárik, P. (2021a). PAGE: A simple and optimal probabilistic gradient estimator for nonconvex optimization. In *International Conference on Machine Learning*, pages 6286–6295. PMLR.
- Li, Z., Hanzely, S., and Richtárik, P. (2021b). ZeroSARAH: Efficient nonconvex finite-sum optimization with zero full gradient computation. *arXiv preprint arXiv:2103.01447*.

- Lin, Y., Han, S., Mao, H., Wang, Y., and Dally, W. J. (2017). Deep gradient compression: Reducing the communication bandwidth for distributed training. *arXiv preprint arXiv:1712.01887*.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., and y Arcas, B. A. (2017). Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR.
- Mishchenko, K., Gorbunov, E., Takáč, M., and Richtárik, P. (2019). Distributed learning with compressed gradient differences. *arXiv preprint arXiv:1901.09269*.
- Narayanan, D., Harlap, A., Phanishayee, A., Seshadri, V., Devanur, N. R., Ganger, G. R., Gibbons,
 P. B., and Zaharia, M. (2019). PipeDream: generalized pipeline parallelism for dnn training. In
 Proceedings of the 27th ACM Symposium on Operating Systems Principles, pages 1–15.
- Nesterov, Y. (2018). Lectures on convex optimization, volume 137. Springer.
- Nguyen, L., Liu, J., Scheinberg, K., and Takáč, M. (2017). SARAH: A novel method for machine learning problems using stochastic recursive gradient. In *The 34th International Conference on Machine Learning*.
- Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein,
 N., Antiga, L., et al. (2019). Pytorch: An imperative style, high-performance deep learning library.
 In Advances in Neural Information Processing Systems (NeurIPS).
- Patel, K. K., Wang, L., Woodworth, B., Bullins, B., and Srebro, N. (2022). Towards optimal communication complexity in distributed non-convex optimization. In *Advances in Neural Information Processing Systems*.
- Ramaswamy, S., Mathews, R., Rao, K., and Beaufays, F. (2019). Federated learning for emoji prediction in a mobile keyboard. *arXiv preprint arXiv:1906.04329*.
- Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Radford, A., Chen, M., and Sutskever, I. (2021). Zero-shot text-to-image generation. *arXiv preprint arXiv:2102.12092*.
- Reddi, S., Charles, Z., Zaheer, M., Garrett, Z., Rush, K., Konečný, J., Kumar, S., and McMahan, H. B. (2020). Adaptive federated optimization. *arXiv preprint arXiv:2003.00295*.
- Richtárik, P., Sokolov, I., and Fatkhullin, I. (2021). EF21: A new, simpler, theoretically better, and practically faster error feedback. *In Neural Information Processing Systems*, 2021.
- Richtárik, P., Sokolov, I., Fatkhullin, I., Gasanov, E., Li, Z., and Gorbunov, E. (2022). 3PC: Three point compressors for communication-efficient distributed training and a better theory for lazy aggregation. *arXiv preprint arXiv:2202.00998*.
- Sapio, A., Canini, M., Ho, C.-Y., Nelson, J., Kalnis, P., Kim, C., Krishnamurthy, A., Moshref, M.,
 Ports, D. R., and Richtárik, P. (2019). Scaling distributed machine learning with in-network
 aggregation. *arXiv* preprint arXiv:1903.06701.
- Stich, S. U., Cordonnier, J.-B., and Jaggi, M. (2018). Sparsified SGD with memory. Advances in
 Neural Information Processing Systems, 31.
- Suresh, A. T., Sun, Z., Ro, J. H., and Yu, F. (2022). Correlated quantization for distributed mean estimation and optimization. *arXiv* preprint arXiv:2203.04925.
- Szlendak, R., Tyurin, A., and Richtárik, P. (2021). Permutation compressors for provably faster distributed nonconvex optimization. *arXiv preprint arXiv:2110.03300*.
- Tyurin, A. and Richtárik, P. (2023). DASHA: Distributed nonconvex optimization with communication compression and optimal oracle complexity. *International Conference on Learning Representations (ICLR)*.
- Vogels, T., He, L., Koloskova, A., Karimireddy, S. P., Lin, T., Stich, S. U., and Jaggi, M. (2021). RelaySum for decentralized deep learning on heterogeneous data. *Advances in Neural Information*

352 Processing Systems, 34.

- Wangni, J., Wang, J., Liu, J., and Zhang, T. (2018). Gradient sparsification for communicationefficient distributed optimization. *Advances in Neural Information Processing Systems*, 31.
- Xu, H., Ho, C.-Y., Abdelmoniem, A. M., Dutta, A., Bergou, E. H., Karatsenidis, K., Canini, M., and Kalnis, P. (2021). Grace: A compressed communication framework for distributed machine learning. In 2021 IEEE 41st International Conference on Distributed Computing Systems (ICDCS), pages 561–572. IEEE.
- Zhao, H., Burlachenko, K., Li, Z., and Richtárik, P. (2021a). Faster rates for compressed federated learning with client-variance reduction. *arXiv preprint arXiv:2112.13097*.
- Zhao, H., Li, Z., and Richtárik, P. (2021b). FedPAGE: A fast local stochastic gradient method for communication-efficient federated learning. *arXiv preprint arXiv:2108.04755*.

663 Contents

364	1	Intr	roduction					
365	imization Problem	1						
366		2.1	Unbiased Compressors	2				
367		2.2	Nodes Partial Participation Assumptions	3				
368	3	Mot	ivation and Related Work	4				
369	4	Con	tributions	5				
370	5	Algo	orithm Description and Main Challenges Towards Partial Participation	5				
371	6	The	orems	7				
372		6.1	Gradient Setting	8				
373		6.2	Finite-Sum Setting	8				
374		6.3	Stochastic Setting	9				
375	A	Nun	nerical Verification of Theoretical Dependencies					
376	В	Orig	ginal DASHA and DASHA-MVR Methods 10					
377	C	Prol	blem of Estimating the Mean in the Partial Participation Setting	17				
378	D	Aux	iliary facts	18				
379		D.1	Sampling Lemma	18				
380		D.2	Compressors Facts	20				
381	E	Proc	ofs of Theorems 2					
382		E.1	Standard Lemmas in the Nonconvex Setting	21				
383		E.2	Generic Lemmas	22				
384		E.3	Proof for DASHA-PP	25				
385		E.4	Proof for DASHA-PP-PAGE	29				
386		E.5	Proof for DASHA-PP-FINITE-MVR	40				
387		E.6	Proof for DASHA-PP-MVR	51				
388	F	Ana	lysis of DASHA-PP under Polyak-Łojasiewicz Condition	64				
389		F.1	Gradient Setting	64				
390		F.2	Finite-Sum Setting	64				
391		F.3	Stochastic Setting	65				
392		F.4	Proofs of Theorems	66				
393			F.4.1 Standard Lemma under Polyak-Łojasiewicz Condition	66				
394			F.4.2 Generic Lemma	66				

399	G.1 Proof	for DASHA-PP-SYNC-MVR	84
398	G Description	n of DASHA-PP-SYNC-MVR	82
397	F.4.5	Proof for DASHA-PP-MVR under PŁ-condition	75
396	F.4.4	Proof for DASHA-PP-PAGE under PŁ-condition	70
395	F.4.3	Proof for DASHA-PP under PŁ-condition	68

400 A Numerical Verification of Theoretical Dependencies

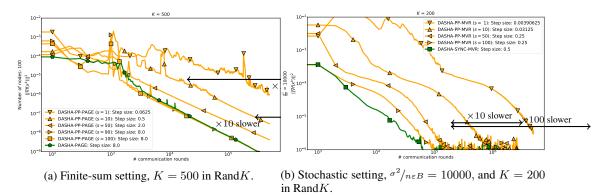


Figure 1: Classification task with the *real-sim* dataset.

Our main goal is to verify the dependeces from the theory. We compare DASHA-PP with DASHA. Clearly, DASHA-PP can not generally perform better than DASHA. In different settings, we verify that the bigger $p_{\rm a}$, the closer DASHA-PP is to DASHA, i.e., DASHA-PP converges no slower than $1/p_{\rm a}$ times.

In all experiments, we take the *real-sim* dataset with dimension d=20,958 and the number of samples equals 72,309 from LIBSVM datasets (Chang and Lin, 2011) (under the 3-clause BSD license), and randomly split the dataset between n=100 nodes equally, ignoring residual samples. In the finite-sum setting, we solve a classification problem with functions

$$f_i(x) := \frac{1}{m} \sum_{j=1}^m \left(1 - \frac{1}{1 + \exp(y_{ij} a_{ij}^\top x)} \right)^2,$$

where $a_{ij} \in \mathbb{R}^d$ is the feature vector of a sample on the i^{th} node, $y_{ij} \in \{-1, 1\}$ is the corresponding label, and m is the number of samples on the i^{th} node for all $i \in [n]$. In the stochastic setting, we consider functions

$$f_i(x_1, x_2) := \mathcal{E}_{j \sim [m]} \left[-\log \left(\frac{\exp \left(a_{ij}^\top x_{y_{ij}} \right)}{\sum_{y \in \{1,2\}} \exp \left(a_{ij}^\top x_y \right)} \right) + \lambda \sum_{y \in \{1,2\}} \sum_{k=1}^d \frac{\{x_y\}_k^2}{1 + \{x_y\}_k^2} \right],$$

where $x_1, x_2 \in \mathbb{R}^d$, $\{\cdot\}_k$ is an indexing operation, $a_{ij} \in \mathbb{R}^d$ is a feature of a sample on the i^{th} node, $y_{ij} \in \{1,2\}$ is a corresponding label, m is the number of samples located on the i^{th} node, constant $\lambda = 0.001$ for all $i \in [n]$.

The code was written in Python 3.6.8 using PyTorch 1.9 (Paszke et al., 2019). A distributed environment was emulated on a machine with Intel(R) Xeon(R) Gold 6226R CPU @ 2.90GHz and 64 cores.

We use the standard setting in experiments 6 where all parameters except step sizes are taken as suggested in theory. Step sizes are finetuned from a set $\{2^i \mid i \in [-10, 10]\}$. We emulate the partial participation setting using s-nice sampling with the number of nodes n=100. We consider the RandK compressor and take the batch size B=1. We plot the relation between communication rounds and values of the norm of gradients at each communication round.

In the finite-sum (Figure 1a) and in the stochastic setting (Figure 1b), we see that the bigger probability $p_a = s/n$ to 1, the closer DASHA-PP to DASHA. Moreover, DASHA-PP with s=10 and s=1 converges approximately $\times 10$ (= $1/p_a$) and $\times 100$ (= $1/p_a$) times slower, accordingly. Our theory predicts such behavior.

⁶Code: https://github.com/mysteryresearcher/dasha-partial-participation

3 B Original DASHA and DASHA-MVR Methods

To simplify the discussion and explanation from the main part, we present the algorithms from (Tyurin and Richtárik, 2023)

Algorithm 6 DASHA

```
1: Input: starting point x^0 \in \mathbb{R}^d, stepsize \gamma > 0, momentum a \in (0,1], number of iterations T \geq 1

2: Initialize g_i^0 \in \mathbb{R}^d on the nodes and g^0 = \frac{1}{n} \sum_{i=1}^n g_i^0 on the server

3: for t = 0, 1, \ldots, T - 1 do

4: x^{t+1} = x^t - \gamma g^t

5: Broadcast x^{t+1}, x^t to all participating^{(a)} nodes

6: for i = 1, \ldots, n in parallel do

7: m_i^{t+1} = \mathcal{C}_i \left( \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - a \left( g_i^t - \nabla f_i(x^t) \right) \right)

8: g_i^{t+1} = g_i^t + m_i^{t+1}

9: Send m_i^{t+1} to the server

10: end for

11: g^{t+1} = g^t + \frac{1}{n} \sum_{i=1}^n m_i^{t+1}

12: end for

13: Output: \hat{x}^T chosen uniformly at random from \{x^t\}_{k=0}^{T-1}
```

Algorithm 7 DASHA-MVR (with batch size B = 1)

```
1: Input: starting point x^0 \in \mathbb{R}^d, stepsize \gamma > 0, momentums a,b \in (0,1], number of iterations T \geq 1

2: Initialize g_i^0 \in \mathbb{R}^d on the nodes and g^0 = \frac{1}{n} \sum_{i=1}^n g_i^0 on the server

3: for t = 0, 1, \ldots, T - 1 do

4: x^{t+1} = x^t - \gamma g^t

5: Broadcast x^{t+1}, x^t to all participating^{(a)} nodes

6: for i = 1, \ldots, n in parallel do

7: h_i^{t+1} = \nabla f_i(x^{t+1}; \xi_i^{t+1}) + (1-b)(h_i^t - \nabla f_i(x^t; \xi_i^{t+1})), \quad \xi_i^{t+1} \sim \mathcal{D}_i

8: m_i^{t+1} = \mathcal{C}_i \left( h_i^{t+1} - h_i^t - a \left( g_i^t - h_i^t \right) \right)

9: g_i^{t+1} = g_i^t + m_i^{t+1}

10: Send m_i^{t+1} to the server

11: end for

12: g^{t+1} = g^t + \frac{1}{n} \sum_{i=1}^n m_i^{t+1}

13: end for

14: Output: \hat{x}^T chosen uniformly at random from \{x^t\}_{k=0}^{T-1}
```

Problem of Estimating the Mean in the Partial Participation Setting

We now provide the example to explain why the only choice of $B = \mathcal{O}\left(\min\left\{\frac{1}{p_a}\sqrt{\frac{m}{n}},\frac{L_{\max}^2}{1\frac{2}{p_a}\widehat{L}^2}\right\}\right)$ and 427

428
$$B = \mathcal{O}\left(\min\left\{\frac{\sigma}{p_a\sqrt{\varepsilon}n}, \frac{L_\sigma^2}{1\frac{p_a}{p_a}\widehat{L}^2}\right\}\right)$$
 in DASHA-PP-PAGE and DASHA-PP-MVR, accordingly, guarantees

the degeneration up to $1/p_a$. This is surprising, because in methods with the variance reduction of 429

430

stochastic gradients (Li et al., 2021a; Tyurin and Richtárik, 2023) we can take the size of batch size $B = \mathcal{O}\left(\sqrt{\frac{m}{n}}\right)$ and $B = \mathcal{O}\left(\frac{\sigma}{\sqrt{\varepsilon n}}\right)$ and guarantee the optimality. Note that the smaller the batch size B, the more the server and the nodes have to communicate to get ε -solution. 431

432

Let us consider the task of estimating the mean of vectors in the distributed setting. Suppose that we 433

have n nodes, and each of them contains m vectors $\{x_{ij}\}_{i=1}^m$, where $x_{ij} \in \mathbb{R}^d$ for all $i \in [n], j \in [m]$. 434

First, let us consider that each node samples a mini-batch I^i of size B with replacement and sends it 435

to the server. Then the server calculates the mean of the mini-batches from nodes. One can easily 436

show that the variance of the estimator is

$$E\left[\left\|\frac{1}{nB}\sum_{i=1}^{n}\sum_{j\in I^{i}}x_{ij} - \frac{1}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}x_{ij}\right\|^{2}\right]$$

$$= \frac{1}{nB}\frac{1}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}\left\|x_{ij} - \frac{1}{m}\sum_{j=1}^{m}x_{ij}\right\|^{2}.$$
(10)

Next, we consider the same task in the partial participation setting with s-nice sampling, i.e., we sample a random set $S \subset [n]$ of $s \in [n]$ nodes without replacement and receive the mini-batches only from the sampled nodes. Such sampling of nodes satisfy Assumption 8 with $p_a = s/n$ and $p_a = \frac{s(s-1)}{n(n-1)}$. In this case, the variance of the estimator (See Lemma 1 with $r_i = 0$ and $s_i = \sum_{j \in I^i} x_{ij}$ is

$$E\left[\left\|\frac{1}{sB}\sum_{i\in S}\sum_{j\in I^{i}}x_{ij} - \frac{1}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}x_{ij}\right\|^{2}\right] \\
= \frac{1}{sB}\underbrace{\frac{1}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}\left\|x_{ij} - \frac{1}{m}\sum_{j=1}^{m}x_{ij}\right\|^{2}}_{\mathcal{L}_{\max}^{2}} \\
+ \frac{n-s}{s(n-1)}\underbrace{\frac{1}{n}\sum_{i=1}^{n}\left\|\frac{1}{m}\sum_{j=1}^{m}x_{ij} - \frac{1}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}x_{ij}\right\|^{2}}_{\widehat{s}_{2}}.$$
(11)

Let us assume that $s \le n/2$. Note that (10) scales with any $B \ge 1$, while (11) only scales when $B = \mathcal{O}\left(\mathcal{L}_{\max}^2/\tilde{\mathcal{L}}^2\right)$. In other words, for large enough B, the variance in (11) does not significantly 444 improves with the growth of B due to the term $\widehat{\mathcal{L}}^2$. In our proof, due to partial participation, the 445 variance from (11) naturally appears, and we get the same effect. As was mentioned in Sections 6.2 and 6.3, it can be seen in our convergence rate bounds.

448 D Auxiliary facts

We list auxiliary facts that we use in our proofs:

1. For all $x, y \in \mathbb{R}^d$, we have

$$||x+y||^2 \le 2 ||x||^2 + 2 ||y||^2$$
 (12)

2. Let us take a *random vector* $\xi \in \mathbb{R}^d$, then

$$E[\|\xi\|^2] = E[\|\xi - E[\xi]\|^2] + \|E[\xi]\|^2.$$
 (13)

452 D.1 Sampling Lemma

This section provides a lemma that we regularly use in our proofs, and it is useful for samplings that satisfy Assumption 8.

Lemma 1. Suppose that a set S is a random subset of a set [n] such that

457 2.
$$\operatorname{\mathbf{Prob}}(i \in S, j \in S) = p_{\mathrm{aa}}, \quad \forall i \neq j \in [n],$$

458 3.
$$p_{\mathrm{aa}} \leq p_{\mathrm{a}}^2,$$

where $p_a \in (0,1]$ and $p_{aa} \in [0,1]$. Let us take random independent vectors $s_i \in \mathbb{R}^d$ for all $i \in [n]$, nonrandom vector $r_i \in \mathbb{R}^d$ for all $i \in [n]$, and random vectors

$$v_i = \begin{cases} r_i + \frac{1}{p_a} s_i, i \in S, \\ r_i, i \notin S, \end{cases}$$

461 *then*

$$E\left[\left\|\frac{1}{n}\sum_{i=1}^{n}v_{i} - E\left[\frac{1}{n}\sum_{i=1}^{n}v_{i}\right]\right\|^{2}\right] \\
= \frac{1}{n^{2}p_{a}}\sum_{i=1}^{n}E\left[\left\|s_{i} - E\left[s_{i}\right]\right\|^{2}\right] + \frac{p_{a} - p_{aa}}{n^{2}p_{a}^{2}}\sum_{i=1}^{n}\left\|E\left[s_{i}\right]\right\|^{2} + \frac{p_{aa} - p_{a}^{2}}{p_{a}^{2}}\left\|\frac{1}{n}\sum_{i=1}^{n}E\left[s_{i}\right]\right\|^{2} \\
\leq \frac{1}{n^{2}p_{a}}\sum_{i=1}^{n}E\left[\left\|s_{i} - E\left[s_{i}\right]\right\|^{2}\right] + \frac{p_{a} - p_{aa}}{n^{2}p_{a}^{2}}\sum_{i=1}^{n}\left\|E\left[s_{i}\right]\right\|^{2}.$$

462 *Proof.* Let us define additional constants p_{an} and p_{nn} , such that

463 1.
$$\operatorname{\mathbf{Prob}}(i \in S, j \notin S) = p_{\operatorname{an}}, \quad \forall i \neq j \in [n],$$

464 2.
$$\mathbf{Prob} (i \notin S, j \notin S) = p_{nn}, \quad \forall i \neq j \in [n].$$

465 Note, that

$$p_{\rm an} = p_{\rm aa} - p_{\rm a} \tag{14}$$

466 and

$$p_{\rm nn} = 1 - p_{\rm aa} - 2p_{\rm an}. \tag{15}$$

467 Using the law of total expectation and

$$E[v_i] = p_a \left(r_i + E\left[\frac{1}{p_a} s_i\right] \right) + (1 - p_a) r_i = r_i + E[s_i],$$

468 we have

$$\begin{split} & \mathbf{E} \left[\left\| \frac{1}{n} \sum_{i=1}^{n} v_{i} - \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} v_{i} \right] \right\|^{2} \right] \\ & = \frac{1}{n^{2}} \sum_{i=1}^{n} \mathbf{E} \left[\left\| v_{i} - (r_{i} + \mathbf{E} \left[s_{i} \right]) \right\|^{2} \right] \\ & + \frac{1}{n^{2}} \sum_{i \neq j}^{n} \mathbf{E} \left[\left\| v_{i} - (r_{i} + \mathbf{E} \left[s_{i} \right]), v_{j} - (r_{j} + \mathbf{E} \left[s_{j} \right]) \right\rangle \right] \\ & = \frac{p_{a}}{n^{2}} \sum_{i=1}^{n} \mathbf{E} \left[\left\| r_{i} + \frac{1}{p_{a}} s_{i} - (r_{i} + \mathbf{E} \left[s_{i} \right]) \right\|^{2} \right] \\ & + \frac{1 - p_{a}}{n^{2}} \sum_{i \neq j}^{n} \mathbf{E} \left[\left\langle r_{i} + \frac{1}{p_{a}} s_{i} - (r_{i} + \mathbf{E} \left[s_{i} \right]), r_{j} + \frac{1}{p_{a}} s_{j} - (r_{j} + \mathbf{E} \left[s_{j} \right]) \right\rangle \right] \\ & + \frac{2p_{an}}{n^{2}} \sum_{i \neq j}^{n} \mathbf{E} \left[\left\langle r_{i} + \frac{1}{p_{a}} s_{i} - (r_{i} + \mathbf{E} \left[s_{i} \right]), r_{j} - (r_{j} + \mathbf{E} \left[s_{j} \right]) \right\rangle \right] \\ & + \frac{p_{nn}}{n^{2}} \sum_{i \neq j}^{n} \left\langle r_{i} - (r_{i} + \mathbf{E} \left[s_{i} \right]), r_{j} - (r_{j} + \mathbf{E} \left[s_{j} \right]) \right\rangle. \end{split}$$

From the independence of random vectors s_i , we obtain

$$\begin{split} & \mathbf{E} \left[\left\| \frac{1}{n} \sum_{i=1}^{n} v_{i} - \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} v_{i} \right] \right\|^{2} \right] \\ & = \frac{p_{a}}{n^{2}} \sum_{i=1}^{n} \mathbf{E} \left[\left\| \frac{1}{p_{a}} s_{i} - \mathbf{E} \left[s_{i} \right] \right\|^{2} \right] \\ & + \frac{1 - p_{a}}{n^{2}} \sum_{i=1}^{n} \left\| \mathbf{E} \left[s_{i} \right] \right\|^{2} \\ & + \frac{p_{aa} (1 - p_{a})^{2}}{n^{2} p_{a}^{2}} \sum_{i \neq j}^{n} \left\langle \mathbf{E} \left[s_{i} \right], \mathbf{E} \left[s_{j} \right] \right\rangle \\ & + \frac{2 p_{an} (p_{a} - 1)}{n^{2} p_{a}} \sum_{i \neq j}^{n} \left\langle \mathbf{E} \left[s_{i} \right], \mathbf{E} \left[s_{j} \right] \right\rangle \\ & + \frac{p_{nn}}{n^{2}} \sum_{i \neq j}^{n} \left\langle \mathbf{E} \left[s_{i} \right], \mathbf{E} \left[s_{j} \right] \right\rangle. \end{split}$$

470 Using (14) and (15), we have

$$E\left[\left\|\frac{1}{n}\sum_{i=1}^{n}v_{i} - E\left[\frac{1}{n}\sum_{i=1}^{n}v_{i}\right]\right\|^{2}\right]$$

$$= \frac{p_{a}}{n^{2}}\sum_{i=1}^{n}E\left[\left\|\frac{1}{p_{a}}s_{i} - E\left[s_{i}\right]\right\|^{2}\right]$$

$$+ \frac{1 - p_{a}}{n^{2}}\sum_{i=1}^{n}\left\|E\left[s_{i}\right]\right\|^{2}$$

$$\begin{split} &+ \frac{p_{\text{aa}} - p_{\text{a}}^2}{n^2 p_{\text{a}}^2} \sum_{i \neq j}^n \left< \mathbf{E} \left[s_i \right], \mathbf{E} \left[s_j \right] \right> \\ \stackrel{\text{(13)}}{=} \frac{1}{n^2 p_{\text{a}}} \sum_{i=1}^n \mathbf{E} \left[\left\| s_i - \mathbf{E} \left[s_i \right] \right\|^2 \right] \\ &+ \frac{1 - p_{\text{a}}}{n^2 p_{\text{a}}} \sum_{i=1}^n \left\| \mathbf{E} \left[s_i \right] \right\|^2 \\ &+ \frac{p_{\text{aa}} - p_{\text{a}}^2}{n^2 p_{\text{a}}^2} \sum_{i \neq j}^n \left< \mathbf{E} \left[s_i \right], \mathbf{E} \left[s_j \right] \right> \\ &= \frac{1}{n^2 p_{\text{a}}} \sum_{i=1}^n \mathbf{E} \left[\left\| s_i - \mathbf{E} \left[s_i \right] \right\|^2 \right] \\ &+ \frac{p_{\text{a}} - p_{\text{aa}}}{n^2 p_{\text{a}}^2} \sum_{i=1}^n \left\| \mathbf{E} \left[s_i \right] \right\|^2 \\ &+ \frac{p_{\text{aa}} - p_{\text{a}}^2}{p_{\text{a}}^2} \left\| \frac{1}{n} \sum_{i=1}^n \mathbf{E} \left[s_i \right] \right\|. \end{split}$$

Finally, using that $p_{aa} \leq p_a^2$, we have

$$\begin{split} & \mathbf{E}\left[\left\|\frac{1}{n}\sum_{i=1}^{n}v_{i} - \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}v_{i}\right]\right\|^{2}\right] \\ & \leq \frac{1}{n^{2}p_{\mathbf{a}}}\sum_{i=1}^{n}\mathbf{E}\left[\left\|s_{i} - \mathbf{E}\left[s_{i}\right]\right\|^{2}\right] + \frac{p_{\mathbf{a}} - p_{\mathbf{aa}}}{n^{2}p_{\mathbf{a}}^{2}}\sum_{i=1}^{n}\left\|\mathbf{E}\left[s_{i}\right]\right\|^{2}. \end{split}$$

Compressors Facts

472

473

479

We define the Rand K compressor that chooses without replacement K coordinates, scales them by a constant factor to preserve unbiasedness and zero-out other coordinates.

Definition 5. Let us take a random subset S from [d], |S| = K, $K \in [d]$. We say that a stochastic mapping $\mathcal{C}:\mathbb{R}^d o \mathbb{R}^d$ is RandK if

$$C(x) = \frac{d}{K} \sum_{j \in S} x_j e_j,$$

where $\{e_i\}_{i=1}^d$ is the standard unit basis.

Theorem 6. If C is RandK, then $C \in \mathbb{U}\left(\frac{d}{k}-1\right)$. 477

See the proof in (Beznosikov et al., 2020). 478

Proofs of Theorems

There are three different sources of randomness in Algorithm 1: the first one from vectors $\{k_i^{t+1}\}_{i=1}^n$, 480 the second one from compressors $\{C_i\}_{i=1}^n$, and the third one from availability of nodes. We define 481 $\mathbf{E}_{k}\left[\cdot\right]$, $\mathbf{E}_{\mathcal{C}}\left[\cdot\right]$ and $\mathbf{E}_{p_{\mathbf{a}}}\left[\cdot\right]$ to be conditional expectations w.r.t. $\{k_{i}^{t+1}\}_{i=1}^{n},\ \{\mathcal{C}_{i}\}_{i=1}^{n},\$ and availability, accordingly, conditioned on all previous randomness. Moreover, we define $\mathbf{E}_{t+1}\left[\cdot\right]$ to be a conditional 482 483 expectation w.r.t. all randomness in iteration t+1 conditioned on all previous randomness. Note, 484 that $E_{t+1} [\cdot] = E_k [E_{\mathcal{C}} [E_{p_a} [\cdot]]]$. 485 In the case of DASHA-PP-PAGE, there are two different sources of randomness from $\{k_i^{t+1}\}_{i=1}^n$. 486 We define $E_{p_{\text{page}}}[\cdot]$ and $E_{B}[\cdot]$ to be conditional expectations w.r.t. the probabilistic switching and mini-batch indices I_i^t , accordingly, conditioned on all previous randomness. Note, that $E_{t+1}[\cdot] =$ $\mathbf{E}_{B}\left[\mathbf{E}_{\mathcal{C}}\left[\mathbf{E}_{p_{\mathrm{a}}}\left[\mathbf{E}_{p_{\mathrm{page}}}\left[\cdot\right]\right]\right]\right] \text{ and } \mathbf{E}_{t+1}\left[\cdot\right] = \mathbf{E}_{B}\left[\mathbf{E}_{p_{\mathrm{page}}}\left[\mathbf{E}_{\mathcal{C}}\left[\mathbf{E}_{p_{\mathrm{a}}}\left[\cdot\right]\right]\right]\right].$

490 E.1 Standard Lemmas in the Nonconvex Setting

- We start the proof of theorems by providing standard lemmas from the nonconvex optimization.
- **Lemma 2.** Suppose that Assumption 2 holds and let $x^{t+1} = x^t \gamma g^t$. Then for any $g^t \in \mathbb{R}^d$ and
- 493 $\gamma > 0$, we have

$$f(x^{t+1}) \le f(x^t) - \frac{\gamma}{2} \left\| \nabla f(x^t) \right\|^2 - \left(\frac{1}{2\gamma} - \frac{L}{2} \right) \left\| x^{t+1} - x^t \right\|^2 + \frac{\gamma}{2} \left\| g^t - \nabla f(x^t) \right\|^2. \tag{16}$$

494 *Proof.* Using L-smoothness, we have

$$f(x^{t+1}) \le f(x^t) + \left\langle \nabla f(x^t), x^{t+1} - x^t \right\rangle + \frac{L}{2} \|x^{t+1} - x^t\|^2$$

= $f(x^t) - \gamma \left\langle \nabla f(x^t), g^t \right\rangle + \frac{L}{2} \|x^{t+1} - x^t\|^2$.

495 Next, due to $-\langle x, y \rangle = \frac{1}{2} \|x - y\|^2 - \frac{1}{2} \|x\|^2 - \frac{1}{2} \|y\|^2$, we obtain

$$f(x^{t+1}) \le f(x^t) - \frac{\gamma}{2} \left\| \nabla f(x^t) \right\|^2 - \left(\frac{1}{2\gamma} - \frac{L}{2} \right) \left\| x^{t+1} - x^t \right\|^2 + \frac{\gamma}{2} \left\| g^t - \nabla f(x^t) \right\|^2.$$

496

497 **Lemma 3.** Suppose that Assumption 1 holds and

$$\mathrm{E}\left[f(x^{t+1})\right] + \gamma \Psi^{t+1} \le \mathrm{E}\left[f(x^t)\right] - \frac{\gamma}{2} \mathrm{E}\left[\left\|\nabla f(x^t)\right\|^2\right] + \gamma \Psi^t + \gamma C, \tag{17}$$

- where Ψ^t is a sequence of numbers, $\Psi^t \geq 0$ for all $t \in [T]$, constant $C \geq 0$, and constant $\gamma > 0$.
 - $\mathrm{E}\left[\left\|\nabla f(\widehat{x}^T)\right\|^2\right] \le \frac{2\Delta_0}{\gamma T} + \frac{2\Psi^0}{T} + 2C,\tag{18}$

- where a point \hat{x}^T is chosen uniformly from a set of points $\{x^t\}_{t=0}^{T-1}$
- 501 *Proof.* By unrolling (17) for t from 0 to T-1, we obtain

$$\frac{\gamma}{2} \sum_{t=0}^{T-1} \mathrm{E}\left[\left\| \nabla f(x^t) \right\|^2 \right] + \mathrm{E}\left[f(x^T) \right] + \gamma \Psi^T \leq f(x^0) + \gamma \Psi^0 + \gamma T C.$$

We subtract f^* , divide inequality by $\frac{\gamma T}{2}$, and take into account that $f(x) \geq f^*$ for all $x \in \mathbb{R}$, and $\Psi^t \geq 0$ for all $t \in [T]$, to get the following inequality:

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbf{E} \left[\left\| \nabla f(x^t) \right\|^2 \right] \le \frac{2\Delta_0}{\gamma T} + \frac{2\Psi^0}{T} + 2C.$$

It is left to consider the choice of a point \hat{x}^T to complete the proof of the lemma.

Lemma 4. If $0 < \gamma \le (L + \sqrt{A})^{-1}$, L > 0, and $A \ge 0$, then

$$\frac{1}{2\gamma} - \frac{L}{2} - \frac{\gamma A}{2} \ge 0.$$

The lemma can be easily checked with the direct calculation.

506 E.2 Generic Lemmas

Lemma 5. Suppose that Assumptions 7 and 8 hold and let us consider sequences g_i^{t+1} , h_i^{t+1} , and k_i^{t+1} from Algorithm 1, then

$$\mathbb{E}_{\mathcal{C}}\left[\mathbb{E}_{p_{a}}\left[\left\|g^{t+1}-h^{t+1}\right\|^{2}\right]\right] \\
\leq \frac{2\omega}{n^{2}p_{a}}\sum_{i=1}^{n}\left\|k_{i}^{t+1}\right\|^{2}+\frac{a^{2}((2\omega+1)\,p_{a}-p_{aa})}{n^{2}p_{a}^{2}}\sum_{i=1}^{n}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2}+(1-a)^{2}\left\|g^{t}-h^{t}\right\|^{2}, \quad (19)$$

509 and

$$E_{\mathcal{C}}\left[E_{p_{a}}\left[\left\|g_{i}^{t+1}-h_{i}^{t+1}\right\|^{2}\right]\right] \\
\leq \frac{2\omega}{p_{a}}\left\|k_{i}^{t+1}\right\|^{2}+\left(\frac{a^{2}(2\omega+1-p_{a})}{p_{a}}+(1-a)^{2}\right)\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} \quad \forall i \in [n].$$
(20)

First, we estimate $\mathbb{E}_{\mathcal{C}}\left[\mathbb{E}_{p_a}\left[\left\|g^{t+1}-h^{t+1}\right\|^2\right]\right]$:

$$\begin{split} & \mathbb{E}_{\mathcal{C}}\left[\mathbb{E}_{p_{a}}\left[\left\|g^{t+1}-h^{t+1}\right\|^{2}\right]\right] \\ & = \mathbb{E}_{\mathcal{C}}\left[\mathbb{E}_{p_{a}}\left[\left\|g^{t+1}-h^{t+1}-\mathbb{E}_{\mathcal{C}}\left[\mathbb{E}_{p_{a}}\left[g^{t+1}-h^{t+1}\right]\right]\right\|^{2}\right]\right] + \left\|\mathbb{E}_{\mathcal{C}}\left[\mathbb{E}_{p_{a}}\left[g^{t+1}-h^{t+1}\right]\right]\right\|^{2}, \end{split}$$

where we used (13). Due to Assumption 8, we have

$$\begin{split} & \mathbf{E}_{\mathcal{C}} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[g_{i}^{t+1} \right] \right] \\ & = p_{\mathbf{a}} \mathbf{E}_{\mathcal{C}} \left[g_{i}^{t} + \mathcal{C}_{i} \left(\frac{1}{p_{\mathbf{a}}} k_{i}^{t+1} - \frac{a}{p_{\mathbf{a}}} \left(g_{i}^{t} - h_{i}^{t} \right) \right) \right] + (1 - p_{\mathbf{a}}) g_{i}^{t} \\ & = g_{i}^{t} + p_{\mathbf{a}} \mathbf{E}_{\mathcal{C}} \left[\mathcal{C}_{i} \left(\frac{1}{p_{\mathbf{a}}} k_{i}^{t+1} - \frac{a}{p_{\mathbf{a}}} \left(g_{i}^{t} - h_{i}^{t} \right) \right) \right] \\ & = g_{i}^{t} + k_{i}^{t+1} - a \left(g_{i}^{t} - h_{i}^{t} \right), \end{split}$$

512 and

$$E_{\mathcal{C}}\left[E_{p_{a}}\left[h_{i}^{t+1}\right]\right] = p_{a}E_{\mathcal{C}}\left[h_{i}^{t} + \frac{1}{p_{a}}k_{i}^{t+1}\right] + (1 - p_{a})h_{i}^{t} = h_{i}^{t} + k_{i}^{t+1}.$$

513 Thus, we can get

$$\begin{split} & \mathbf{E}_{\mathcal{C}} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| g^{t+1} - h^{t+1} \right\|^{2} \right] \right] \\ & = \mathbf{E}_{\mathcal{C}} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| g^{t+1} - h^{t+1} - \mathbf{E}_{\mathcal{C}} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[g^{t+1} - h^{t+1} \right] \right] \right\|^{2} \right] \right] + (1 - a)^{2} \left\| g^{t} - h^{t} \right\|^{2}. \end{split}$$

Due to the independence of compressors, we can use Lemma 1 with $r_i = g_i^t - h_i^t$ and $s_i =$

515
$$p_{\mathbf{a}}\mathcal{C}_i\left(\frac{1}{p_{\mathbf{a}}}k_i^{t+1} - \frac{a}{p_{\mathbf{a}}}\left(g_i^t - h_i^t\right)\right) - k_i^{t+1}$$
, and obtain
$$\mathbb{E}_{\mathcal{C}}\left[\mathbb{E}_{p_{\mathbf{a}}}\left[\left\|g^{t+1} - h^{t+1}\right\|^2\right]\right]$$

$$\leq \frac{1}{n^{2}p_{a}} \sum_{i=1}^{n} E_{\mathcal{C}} \left[\left\| p_{a} \mathcal{C}_{i} \left(\frac{1}{p_{a}} k_{i}^{t+1} - \frac{a}{p_{a}} \left(g_{i}^{t} - h_{i}^{t} \right) \right) - k_{i}^{t+1} - E_{\mathcal{C}} \left[p_{a} \mathcal{C}_{i} \left(\frac{1}{p_{a}} k_{i}^{t+1} - \frac{a}{p_{a}} \left(g_{i}^{t} - h_{i}^{t} \right) \right) - k_{i}^{t+1} \right] \right\|^{2} \right]$$

$$+ \frac{p_{a} - p_{aa}}{n^{2} p_{a}^{2}} \sum_{i=1}^{n} \left\| E_{\mathcal{C}} \left[p_{a} \mathcal{C}_{i} \left(\frac{1}{p_{a}} k_{i}^{t+1} - \frac{a}{p_{a}} \left(g_{i}^{t} - h_{i}^{t} \right) \right) - k_{i}^{t+1} \right] \right\|^{2}$$

$$\begin{aligned} &+ (1-a)^{2} \left\| g^{t} - h^{t} \right\|^{2} \\ &= \frac{p_{a}}{n^{2}} \sum_{i=1}^{n} E_{\mathcal{C}} \left[\left\| \mathcal{C}_{i} \left(\frac{1}{p_{a}} k_{i}^{t+1} - \frac{a}{p_{a}} \left(g_{i}^{t} - h_{i}^{t} \right) \right) - \left(\frac{1}{p_{a}} k_{i}^{t+1} - \frac{a}{p_{a}} \left(g_{i}^{t} - h_{i}^{t} \right) \right) \right\|^{2} \right] \\ &+ \frac{a^{2} \left(p_{a} - p_{aa} \right)}{n^{2} p_{a}^{2}} \sum_{i=1}^{n} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} + (1-a)^{2} \left\| g^{t} - h^{t} \right\|^{2}. \end{aligned}$$

516 From Assumption 7, we have

$$\begin{split} & \mathbf{E}_{\mathcal{C}}\left[\mathbf{E}_{p_{\mathbf{a}}}\left[\left\|g^{t+1}-h^{t+1}\right\|^{2}\right]\right] \\ & \leq \frac{\omega p_{\mathbf{a}}}{n^{2}}\sum_{i=1}^{n}\left\|\frac{1}{p_{\mathbf{a}}}k_{i}^{t+1}-\frac{a}{p_{\mathbf{a}}}\left(g_{i}^{t}-h_{i}^{t}\right)\right\|^{2}+\frac{a^{2}\left(p_{\mathbf{a}}-p_{\mathbf{a}\mathbf{a}}\right)}{n^{2}p_{\mathbf{a}}^{2}}\sum_{i=1}^{n}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2}+\left(1-a\right)^{2}\left\|g^{t}-h^{t}\right\|^{2} \\ & = \frac{\omega}{n^{2}p_{\mathbf{a}}}\sum_{i=1}^{n}\left\|k_{i}^{t+1}-a\left(g_{i}^{t}-h_{i}^{t}\right)\right\|^{2}+\frac{a^{2}\left(p_{\mathbf{a}}-p_{\mathbf{a}\mathbf{a}}\right)}{n^{2}p_{\mathbf{a}}^{2}}\sum_{i=1}^{n}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2}+\left(1-a\right)^{2}\left\|g^{t}-h^{t}\right\|^{2} \\ & \leq \frac{2\omega}{n^{2}p_{\mathbf{a}}}\sum_{i=1}^{n}\left\|k_{i}^{t+1}\right\|^{2}+\frac{a^{2}\left((2\omega+1)p_{\mathbf{a}}-p_{\mathbf{a}\mathbf{a}}\right)}{n^{2}p_{\mathbf{a}}^{2}}\sum_{i=1}^{n}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2}+\left(1-a\right)^{2}\left\|g^{t}-h^{t}\right\|^{2}. \end{split}$$

The second inequality can be proved almost in the same way:

$$\begin{split} & \mathcal{E}_{\mathcal{C}}\left[\mathcal{E}_{p_{a}}\left[\left\|g_{i}^{t+1}-h_{i}^{t+1}\right\|^{2}\right]\right] \\ & = \mathcal{E}_{\mathcal{C}}\left[\mathcal{E}_{p_{a}}\left[\left\|g_{i}^{t+1}-h_{i}^{t+1}-\mathcal{E}_{\mathcal{C}}\left[\mathcal{E}_{p_{a}}\left[g_{i}^{t+1}-h_{i}^{t+1}\right]\right]\right\|^{2}\right]\right] + \left\|\mathcal{E}_{\mathcal{C}}\left[\mathcal{E}_{p_{a}}\left[g_{i}^{t+1}-h_{i}^{t+1}\right]\right]\right\|^{2} \\ & = \mathcal{E}_{\mathcal{C}}\left[\mathcal{E}_{p_{a}}\left[\left\|g_{i}^{t+1}-h_{i}^{t+1}-g_{i}^{t}+a\left(g_{i}^{t}-h_{i}^{t}\right)+h_{i}^{t}\right\|^{2}\right]\right] + (1-a)^{2}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} \\ & = p_{a}\mathcal{E}_{\mathcal{C}}\left[\left\|\mathcal{C}_{i}\left(\frac{1}{p_{a}}k_{i}^{t+1}-\frac{a}{p_{a}}\left(g_{i}^{t}-h_{i}^{t}\right)\right)-\frac{1}{p_{a}}k_{i}^{t+1}+a\left(g_{i}^{t}-h_{i}^{t}\right)\right\|^{2}\right] \\ & + a^{2}(1-p_{a})\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} + (1-a)^{2}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} \\ & = p_{a}\mathcal{E}_{\mathcal{C}}\left[\left\|\mathcal{C}_{i}\left(\frac{1}{p_{a}}k_{i}^{t+1}-\frac{a}{p_{a}}\left(g_{i}^{t}-h_{i}^{t}\right)\right)-\left(\frac{1}{p_{a}}k_{i}^{t+1}-\frac{a}{p_{a}}\left(g_{i}^{t}-h_{i}^{t}\right)\right)\right\|^{2}\right] \\ & + a^{2}\left(1-p_{a}\right)^{2}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} \\ & + a^{2}(1-p_{a})\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} + (1-a)^{2}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} \\ & \leq \frac{\omega}{p_{a}}\left\|k_{i}^{t+1}-a\left(g_{i}^{t}-h_{i}^{t}\right)\right\|^{2} \\ & + \frac{a^{2}(1-p_{a})}{p_{a}}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} + (1-a)^{2}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} \\ & \leq \frac{2\omega}{p_{a}}\left\|k_{i}^{t+1}\right\|^{2} + \frac{a^{2}(2\omega+1-p_{a})}{p_{a}}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} + (1-a)^{2}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2}. \end{split}$$

Lemma 6. Suppose that Assumptions 2, 7, and 8 hold and let us take $a = \frac{p_a}{2\omega + 1}$, then

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t}) - \frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] + \frac{4\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|k_{i}^{t+1}\right\|^{2}\right]. \end{split}$$

Proof. Due to Lemma 2 and the update step from Line 5 in Algorithm 1, we have $\mathrm{E}_{t+1}\left[f(x^{t+1})\right]$

$$\leq \mathbf{E}_{t+1} \left[f(x^{t}) - \frac{\gamma}{2} \left\| \nabla f(x^{t}) \right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2} \right) \left\| x^{t+1} - x^{t} \right\|^{2} + \frac{\gamma}{2} \left\| g^{t} - \nabla f(x^{t}) \right\|^{2} \right]$$

$$= \mathbf{E}_{t+1} \left[f(x^{t}) - \frac{\gamma}{2} \left\| \nabla f(x^{t}) \right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2} \right) \left\| x^{t+1} - x^{t} \right\|^{2} + \frac{\gamma}{2} \left\| g^{t} - h^{t} + h^{t} - \nabla f(x^{t}) \right\|^{2} \right]$$

$$\leq \mathbf{E}_{t+1} \left[f(x^{t}) - \frac{\gamma}{2} \left\| \nabla f(x^{t}) \right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2} \right) \left\| x^{t+1} - x^{t} \right\|^{2} + \gamma \left(\left\| g^{t} - h^{t} \right\|^{2} + \left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \right] \right).$$

- Let us fix some constants $\kappa, \eta \in [0, \infty)$ that we will define later. Combining the last inequality,
- bounds (19), (20) and using the law of total expectation, we get

$$E\left[f(x^{t+1})\right]$$

$$\begin{split} &+\kappa \mathbf{E}\left[\left\|g^{t+1}-h^{t+1}\right\|^{2}\right]+\eta \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1}-h_{i}^{t+1}\right\|^{2}\right]\\ &=\mathbf{E}\left[\mathbf{E}_{t+1}\left[f(x^{t+1})\right]\right]\\ &+\kappa \mathbf{E}\left[\mathbf{E}_{\mathcal{C}}\left[\mathbf{E}_{p_{a}}\left[\left\|g^{t+1}-h^{t+1}\right\|^{2}\right]\right]\right]+\eta \mathbf{E}\left[\mathbf{E}_{\mathcal{C}}\left[\mathbf{E}_{p_{a}}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1}-h_{i}^{t+1}\right\|^{2}\right]\right]\right]\\ &\leq \mathbf{E}\left[f(x^{t})-\frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2}-\left(\frac{1}{2\gamma}-\frac{L}{2}\right)\left\|x^{t+1}-x^{t}\right\|^{2}+\gamma\left(\left\|g^{t}-h^{t}\right\|^{2}+\left\|h^{t}-\nabla f(x^{t})\right\|^{2}\right)\right]\\ &+\kappa \mathbf{E}\left[\frac{2\omega}{n^{2}p_{a}}\sum_{i=1}^{n}\left\|k_{i}^{t+1}\right\|^{2}+\frac{a^{2}((2\omega+1)\,p_{a}-p_{aa})}{n^{2}p_{a}^{2}}\sum_{i=1}^{n}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2}+(1-a)^{2}\left\|g^{t}-h^{t}\right\|^{2}\right]\\ &+\eta \mathbf{E}\left[\frac{2\omega}{np_{a}}\sum_{i=1}^{n}\left\|k_{i}^{t+1}\right\|^{2}+\left(\frac{a^{2}(2\omega+1-p_{a})}{p_{a}}+(1-a)^{2}\right)\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2}\right]\\ &=\mathbf{E}\left[f(x^{t})-\frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2}-\left(\frac{1}{2\gamma}-\frac{L}{2}\right)\left\|x^{t+1}-x^{t}\right\|^{2}+\gamma\left\|h^{t}-\nabla f(x^{t})\right\|^{2}\right]\\ &+\left(\gamma+\kappa\left(1-a\right)^{2}\right)\mathbf{E}\left[\left\|g^{t}-h^{t}\right\|^{2}\right]\\ &+\left(\frac{\kappa a^{2}((2\omega+1)\,p_{a}-p_{aa})}{np_{a}^{2}}+\eta\left(\frac{a^{2}(2\omega+1-p_{a})}{p_{a}}+(1-a)^{2}\right)\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2}\right]\\ &+\left(\frac{2\kappa\omega}{np_{a}}+\frac{2\eta\omega}{p_{a}}\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|k_{i}^{t+1}\right\|^{2}\right]. \end{split}$$

Now, by taking $\kappa = \frac{\gamma}{a}$, we can see that $\gamma + \kappa (1-a)^2 \le \kappa$, and thus

$$\begin{split} & + \frac{\gamma}{a} \mathbf{E} \left[\left\| g^{t+1} - h^{t+1} \right\|^{2} \right] + \eta \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| g_{i}^{t+1} - h_{i}^{t+1} \right\|^{2} \right] \\ & \leq \mathbf{E} \left[f(x^{t}) - \frac{\gamma}{2} \left\| \nabla f(x^{t}) \right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2} \right) \left\| x^{t+1} - x^{t} \right\|^{2} + \gamma \left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \right] \\ & + \frac{\gamma}{a} \mathbf{E} \left[\left\| g^{t} - h^{t} \right\|^{2} \right] \\ & + \left(\frac{\gamma a((2\omega + 1) \, p_{a} - p_{aa})}{n p_{a}^{2}} + \eta \left(\frac{a^{2}(2\omega + 1 - p_{a})}{p_{a}} + (1 - a)^{2} \right) \right) \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \right] \\ & + \left(\frac{2\gamma\omega}{anp_{a}} + \frac{2\eta\omega}{p_{a}} \right) \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| k_{i}^{t+1} \right\|^{2} \right]. \end{split}$$

Next, by taking
$$\eta = \frac{\gamma((2\omega+1)p_a-p_{aa})}{np_a^2}$$
 and considering the choice of a , one can show that
$$\left(\frac{\gamma a((2\omega+1)p_a-p_{aa})}{np_a^2} + \eta\left(\frac{a^2(2\omega+1-p_a)}{p_a} + (1-a)^2\right)\right) \leq \eta.$$
 Thus

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t}) - \frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \left(\frac{2\gamma(2\omega+1)\omega}{np_{\mathbf{a}}^{2}} + \frac{2\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})\omega}{np_{\mathbf{a}}^{3}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|k_{i}^{t+1}\right\|^{2}\right]. \end{split}$$

Considering that $p_{aa} \ge 0$, we can simplify the last term and get

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t}) - \frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \frac{4\gamma(2\omega+1)\omega}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|k_{i}^{t+1}\right\|^{2}\right]. \end{split}$$

528 **E.3 Proof for** DASHA-PP

Lemma 7. Suppose that Assumptions 3 and 8 hold. For h_i^{t+1} and k_i^{t+1} from Algorithm 1 (DASHA-PP)

530 we have

527

I.
$$E_{p_{\mathbf{a}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right]$$

$$\leq \frac{2 \left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}} \right) \widehat{L}^{2}}{n p_{\mathbf{a}}^{2}} \left\| x^{t+1} - x^{t} \right\|^{2} + \frac{2 b^{2} \left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}} \right)}{n^{2} p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} + \left(1 - b \right)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}.$$

2.
$$E_{p_{a}} \left[\left\| h_{i}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right]$$

$$\leq \frac{2(1 - p_{a})}{p_{a}} L_{i}^{2} \left\| x^{t+1} - x^{t} \right\|^{2} + \left(\frac{2b^{2}(1 - p_{a})}{p_{a}} + (1 - b)^{2} \right) \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}, \quad \forall i \in [n].$$

3. $\left\|k_i^{t+1}\right\|^2 \le 2L_i^2 \left\|x^{t+1} - x^t\right\|^2 + 2b^2 \left\|h_i^t - \nabla f_i(x^t)\right\|^2, \quad \forall i \in [n].$

Proof. First, let us proof the bound for $\mathbf{E}_k\left[\mathbf{E}_{p_a}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^2\right]\right]$:

$$\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right]$$

$$= \mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h^{t+1} - \mathbf{E}_{p_{\mathbf{a}}} \left[h^{t+1} \right] \right\|^{2} \right] + \left\| \mathbf{E}_{p_{\mathbf{a}}} \left[h^{t+1} \right] - \nabla f(x^{t+1}) \right\|^{2}.$$

532 Using

$$\mathbb{E}_{p_{\mathbf{a}}} \left[h_i^{t+1} \right] = h_i^t + \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t))$$

and (13), we have

$$\begin{split} & \mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \\ & = \mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h^{t+1} - \mathbf{E}_{p_{\mathbf{a}}} \left[h^{t+1} \right] \right\|^{2} \right] + (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}. \end{split}$$

We can use Lemma 1 with $r_i = h_i^t$ and $s_i = k_i^{t+1}$ to obtain

$$\begin{split} & \mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \\ & \leq \frac{1}{n^{2}p_{\mathbf{a}}} \sum_{i=1}^{n} \left\| k_{i}^{t+1} - k_{i}^{t+1} \right\|^{2} + \frac{p_{\mathbf{a}} - p_{\mathbf{aa}}}{n^{2}p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| k_{i}^{t+1} \right\|^{2} + (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \\ & = \frac{p_{\mathbf{a}} - p_{\mathbf{aa}}}{n^{2}p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right\|^{2} + (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \\ & \leq \frac{2 \left(p_{\mathbf{a}} - p_{\mathbf{aa}} \right)}{n^{2}p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} + \frac{2b^{2} \left(p_{\mathbf{a}} - p_{\mathbf{aa}} \right)}{n^{2}p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} + (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \\ & \leq \frac{2 \left(p_{\mathbf{a}} - p_{\mathbf{aa}} \right) \widehat{L}^{2}}{np_{\mathbf{a}}^{2}} \left\| x^{t+1} - x^{t} \right\|^{2} + \frac{2b^{2} \left(p_{\mathbf{a}} - p_{\mathbf{aa}} \right)}{n^{2}p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} + (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}. \end{split}$$

In the last in inequality, we used Assumption 3. Now, we prove the second inequality:

$$\begin{split} & \mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h_{i}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \\ & = \mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h_{i}^{t+1} - \mathbf{E}_{p_{\mathbf{a}}} \left[h_{i}^{t+1} \right] \right\|^{2} \right] + \left\| \mathbf{E}_{p_{\mathbf{a}}} \left[h_{i}^{t+1} \right] - \nabla f_{i}(x^{t+1}) \right\|^{2} \\ & = \mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h_{i}^{t+1} - \left(h_{i}^{t} + \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b(h_{i}^{t} - \nabla f_{i}(x^{t})) \right) \right\|^{2} \right] + (1 - b)^{2} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \\ & = \frac{(1 - p_{\mathbf{a}})^{2}}{p_{\mathbf{a}}} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b(h_{i}^{t} - \nabla f_{i}(x^{t})) \right\|^{2} \\ & + (1 - p_{\mathbf{a}}) \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b(h_{i}^{t} - \nabla f_{i}(x^{t})) \right\|^{2} + (1 - b)^{2} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \\ & = \frac{(1 - p_{\mathbf{a}})}{p_{\mathbf{a}}} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b(h_{i}^{t} - \nabla f_{i}(x^{t})) \right\|^{2} + (1 - b)^{2} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \\ & \leq \frac{2(1 - p_{\mathbf{a}})}{p_{\mathbf{a}}} L_{i}^{2} \left\| x^{t+1} - x^{t} \right\|^{2} + \left(\frac{2b^{2}(1 - p_{\mathbf{a}})}{p_{\mathbf{a}}} + (1 - b)^{2} \right) \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}. \end{split}$$

Finally, the third inequality of the theorem follows from (12) and Assumption 3.

Theorem 2. Suppose that Assumptions 1, 2, 3, 7 and 8 hold. Let us take $a=\frac{p_a}{2\omega+1}, b=\frac{p_a}{2-p_a}$

$$\gamma \leq \left(L + \left[\frac{48\omega\left(2\omega + 1\right)}{np_{\mathrm{a}}^2} + \frac{16}{np_{\mathrm{a}}^2}\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)\right]^{1/2}\widehat{L}\right)^{-1},$$

and $g_i^0 = h_i^0 = \nabla f_i(x^0)$ for all $i \in [n]$ in Algorithm 1 (DASHA-PP), then $\mathbb{E}\left[\left\|\nabla f(\widehat{x}^T)\right\|^2\right] \leq \frac{2\Delta_0}{\gamma T}$.

Proof. Let us fix constants $\nu, \rho \in [0, \infty)$ that we will define later. Considering Lemma 6, Lemma 7, and the law of total expectation, we obtain

$$\mathrm{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathrm{a}}}\mathrm{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathrm{a}} - p_{\mathrm{aa}})}{np_{\mathrm{a}}^{2}}\mathrm{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right]$$

$$\begin{split} &+\nu \mathbf{E}\left[\left\|h^{t+1}-\nabla f(x^{t+1})\right\|^{2}\right]+\rho \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1}-\nabla f_{i}(x^{t+1})\right\|^{2}\right]\\ &=\mathbf{E}\left[f(x^{t+1})\right]+\frac{\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1}-h^{t+1}\right\|^{2}\right]+\frac{\gamma((2\omega+1)\,p_{\mathbf{a}}-p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1}-h_{i}^{t+1}\right\|^{2}\right]\\ &+\nu \mathbf{E}\left[\mathbf{E}_{p_{\mathbf{a}}}\left[\left\|h^{t+1}-\nabla f(x^{t+1})\right\|^{2}\right]\right]+\rho \mathbf{E}\left[\mathbf{E}_{p_{\mathbf{a}}}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1}-\nabla f_{i}(x^{t+1})\right\|^{2}\right]\right]\\ &\leq \mathbf{E}\left[f(x^{t})-\frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2}-\left(\frac{1}{2\gamma}-\frac{L}{2}\right)\left\|x^{t+1}-x^{t}\right\|^{2}+\gamma\left\|h^{t}-\nabla f(x^{t})\right\|^{2}\right]\\ &+\frac{\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t}-h^{t}\right\|^{2}\right]+\frac{\gamma((2\omega+1)\,p_{\mathbf{a}}-p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2}\right]\\ &+\frac{4\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[2\hat{L}^{2}\left\|x^{t+1}-x^{t}\right\|^{2}+2b^{2}\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t}-\nabla f_{i}(x^{t})\right\|^{2}\right]\\ &+\nu \mathbf{E}\left[\frac{2\left(p_{\mathbf{a}}-p_{\mathbf{a}\mathbf{a}}\right)\hat{L}^{2}}{np_{\mathbf{a}}^{2}}\left\|x^{t+1}-x^{t}\right\|^{2}+\frac{2b^{2}\left(p_{\mathbf{a}}-p_{\mathbf{a}\mathbf{a}}\right)}{n^{2}p_{\mathbf{a}}^{2}}\sum_{i=1}^{n}\left\|h_{i}^{t}-\nabla f_{i}(x^{t})\right\|^{2}+\left(1-b\right)^{2}\left\|h^{t}-\nabla f(x^{t})\right\|^{2}\right]\\ &+\rho \mathbf{E}\left[\frac{2(1-p_{\mathbf{a}})}{p_{\mathbf{a}}}\hat{L}^{2}\left\|x^{t+1}-x^{t}\right\|^{2}+\left(\frac{2b^{2}(1-p_{\mathbf{a}})}{p_{\mathbf{a}}}+\left(1-b\right)^{2}\right)\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t}-\nabla f_{i}(x^{t})\right\|^{2}\right]. \end{split}$$

After rearranging the terms, we get

$$\begin{split} & \mathbb{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbb{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \nu\mathbb{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbb{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbb{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbb{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}}\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{8\gamma\omega\left(2\omega+1\right)\,\hat{L}^{2}}{np_{\mathbf{a}}^{2}} - \nu\frac{2\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)\,\hat{L}^{2}}{np_{\mathbf{a}}^{2}} - \rho\frac{2(1-p_{\mathbf{a}})\hat{L}^{2}}{p_{\mathbf{a}}}\right)\mathbb{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \left(\gamma + \nu(1-b)^{2}\right)\mathbb{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \nu\frac{2b^{2}\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}} + \rho\left(\frac{2b^{2}(1-p_{\mathbf{a}})}{p_{\mathbf{a}}} + (1-b)^{2}\right)\right)\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right]. \end{split}$$

542 By taking $\nu=\frac{\gamma}{b},$ one can show that $\left(\gamma+\nu(1-b)^2\right)\leq \nu,$ and

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1) p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho \mathbf{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1) p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \end{split}$$

$$\begin{split} &-\left(\frac{1}{2\gamma}-\frac{L}{2}-\frac{8\gamma\omega\left(2\omega+1\right)\widehat{L}^{2}}{np_{a}^{2}}-\frac{2\gamma\left(p_{a}-p_{aa}\right)\widehat{L}^{2}}{bnp_{a}^{2}}-\rho\frac{2(1-p_{a})\widehat{L}^{2}}{p_{a}}\right)\mathrm{E}\left[\left\|\boldsymbol{x}^{t+1}-\boldsymbol{x}^{t}\right\|^{2}\right]\\ &+\frac{\gamma}{b}\mathrm{E}\left[\left\|\boldsymbol{h}^{t}-\nabla\boldsymbol{f}(\boldsymbol{x}^{t})\right\|^{2}\right]\\ &+\left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{a}^{2}}+\frac{2\gamma b\left(p_{a}-p_{aa}\right)}{np_{a}^{2}}+\rho\left(\frac{2b^{2}(1-p_{a})}{p_{a}}+(1-b)^{2}\right)\right)\mathrm{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|\boldsymbol{h}_{i}^{t}-\nabla\boldsymbol{f}_{i}(\boldsymbol{x}^{t})\right\|^{2}\right]. \end{split}$$

Note that $b = \frac{p_a}{2-p_a}$, thus

$$\begin{split} & \left(\frac{8b^2\gamma\omega(2\omega+1)}{np_{\rm a}^2} + \frac{2\gamma b\left(p_{\rm a} - p_{\rm aa}\right)}{np_{\rm a}^2} + \rho\left(\frac{2b^2(1-p_{\rm a})}{p_{\rm a}} + (1-b)^2\right) \right) \\ & \leq \left(\frac{8b^2\gamma\omega(2\omega+1)}{np_{\rm a}^2} + \frac{2\gamma b\left(p_{\rm a} - p_{\rm aa}\right)}{np_{\rm a}^2} + \rho\left(1-b\right) \right). \end{split}$$

And if we take $ho=rac{8b\gamma\omega(2\omega+1)}{np_{\rm a}^2}+rac{2\gamma(p_{\rm a}-p_{\rm aa})}{np_{\rm a}^2},$ then

$$\left(\frac{8b^2\gamma\omega(2\omega+1)}{np_a^2} + \frac{2\gamma b\left(p_a - p_{aa}\right)}{np_a^2} + \rho\left(1 - b\right)\right) \le \rho,$$

545 and

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{8\gamma\omega\left(2\omega+1\right)\,\hat{L}^{2}}{np_{\mathbf{a}}^{2}} - \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)\,\hat{L}^{2}}{bnp_{\mathbf{a}}^{2}} \right) \mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & - \frac{16b\gamma\omega(2\omega+1)(1-p_{\mathbf{a}})\hat{L}^{2}}{np_{\mathbf{a}}^{3}} - \frac{4\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)(1-p_{\mathbf{a}})\hat{L}^{2}}{np_{\mathbf{a}}^{3}}\right) \mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right]. \end{split}$$

Let us simplify the last inequality. First, note that

$$\frac{16b\gamma\omega(2\omega+1)(1-p_{\mathrm{a}})\widehat{L}^{2}}{np_{\mathrm{a}}^{3}} \leq \frac{16\gamma\omega(2\omega+1)\widehat{L}^{2}}{np_{\mathrm{a}}^{2}},$$

due to $b \leq p_a$. Second,

$$\frac{2\gamma\left(p_{\mathrm{a}}-p_{\mathrm{aa}}\right)\widehat{L}^{2}}{bnp_{\mathrm{a}}^{2}}\leq\frac{4\gamma\left(p_{\mathrm{a}}-p_{\mathrm{aa}}\right)\widehat{L}^{2}}{np_{\mathrm{a}}^{3}},$$

due to $b \ge \frac{p_a}{2}$. All in all, we have

$$E\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{a}} E\left[\|g^{t+1} - h^{t+1}\|^{2}\right] + \frac{\gamma((2\omega+1)p_{a} - p_{aa})}{np_{a}^{2}} E\left[\frac{1}{n}\sum_{i=1}^{n}\|g_{i}^{t+1} - h_{i}^{t+1}\|^{2}\right] + \frac{\gamma}{b} E\left[\|h^{t+1} - \nabla f(x^{t+1})\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{a}^{2}} + \frac{2\gamma(p_{a} - p_{aa})}{np_{a}^{2}}\right) E\left[\frac{1}{n}\sum_{i=1}^{n}\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\|^{2}\right]$$

$$\leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{24\gamma\omega\left(2\omega+1\right)\,\hat{L}^{2}}{np_{\mathbf{a}}^{2}} - \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)\,\hat{L}^{2}}{np_{\mathbf{a}}^{3}}\right)\mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ + \frac{\gamma}{b}\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right].$$

Using Lemma 4 and the assumption about γ , we get

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right]. \end{split}$$

550 It is left to apply Lemma 3 with

$$\begin{split} \Psi^{t} & = & \frac{(2\omega+1)}{p_{\mathrm{a}}} \mathbf{E} \left[\left\| g^{t} - h^{t} \right\|^{2} \right] + \frac{((2\omega+1)\,p_{\mathrm{a}} - p_{\mathrm{aa}})}{np_{\mathrm{a}}^{2}} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \right] \\ & + & \frac{1}{b} \mathbf{E} \left[\left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \right] + \left(\frac{8b\omega(2\omega+1)}{np_{\mathrm{a}}^{2}} + \frac{2\left(p_{\mathrm{a}} - p_{\mathrm{aa}}\right)}{np_{\mathrm{a}}^{2}} \right) \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \right] \end{split}$$

to conclude the proof.

552 E.4 Proof for DASHA-PP-PAGE

553 Let us denote

$$\begin{split} k_{i,1}^{t+1} &:= \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - \frac{b}{p_{\text{page}}} \left(h_i^t - \nabla f_i(x^t) \right), \\ k_{i,2}^{t+1} &:= \frac{1}{B} \sum_{j \in I_i^t} \left(\nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^t) \right), \\ h_{i,1}^{t+1} &:= \begin{cases} h_i^t + \frac{1}{p_{\text{a}}} k_{i,1}^{t+1}, & i^{\text{th}} \text{ node is } \textit{participating}, \\ h_i^t, & \text{otherwise}, \end{cases} \\ h_{i,2}^{t+1} &:= \begin{cases} h_i^t + \frac{1}{p_{\text{a}}} k_{i,2}^{t+1}, & i^{\text{th}} \text{ node is } \textit{participating}, \\ h_i^t, & \text{otherwise}, \end{cases} \end{split}$$

554
$$h_1^{t+1} := \frac{1}{n} \sum_{i=1}^n h_{i,1}^{t+1}$$
, and $h_2^{t+1} := \frac{1}{n} \sum_{i=1}^n h_{i,2}^{t+1}$. Note, that
$$h^{t+1} = \begin{cases} h_1^{t+1}, & \text{with probability } p_{\text{page}}, \\ h_2^{t+1}, & \text{with probability } 1 - p_{\text{page}}. \end{cases}$$

Lemma 8. Suppose that Assumptions 3, 4, and 8 hold. For h_i^{t+1} and k_i^{t+1} from Algorithm 1 (DASHA-PP-PAGE) we have

1.

$$\begin{split} & \mathbf{E}_{B} \left[\mathbf{E}_{p_{a}} \left[\mathbf{E}_{p_{page}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \right] \\ & \leq \left(\frac{2 \left(p_{a} - p_{aa} \right) \widehat{L}^{2}}{n p_{a}^{2}} + \frac{\left(1 - p_{page} \right) L_{\max}^{2}}{n p_{a} B} \right) \left\| x^{t+1} - x^{t} \right\|^{2} \\ & + \frac{2 \left(p_{a} - p_{aa} \right) b^{2}}{n^{2} p_{a}^{2} p_{page}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} + \left(p_{page} \left(1 - \frac{b}{p_{page}} \right)^{2} + \left(1 - p_{page} \right) \right) \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}. \end{split}$$

2

$$\begin{split} &\mathbf{E}_{B}\left[\mathbf{E}_{p_{\mathbf{a}}}\left[\mathbf{E}_{p_{page}}\left[\left\|\boldsymbol{h}_{i}^{t+1} - \nabla f_{i}(\boldsymbol{x}^{t+1})\right\|^{2}\right]\right]\right] \\ &\leq \left(\frac{2\left(1-p_{\mathbf{a}}\right)L_{i}^{2}}{p_{\mathbf{a}}} + \frac{\left(1-p_{page}\right)L_{\max}^{2}}{p_{\mathbf{a}}B}\right)\left\|\boldsymbol{x}^{t+1} - \boldsymbol{x}^{t}\right\|^{2} \\ &+ \left(\frac{2\left(1-p_{\mathbf{a}}\right)b^{2}}{p_{\mathbf{a}}p_{page}} + p_{page}\left(1-\frac{b}{p_{page}}\right)^{2} + \left(1-p_{page}\right)\right)\left\|\boldsymbol{h}_{i}^{t} - \nabla f_{i}(\boldsymbol{x}^{t})\right\|^{2}, \quad \forall i \in [n]. \end{split}$$

3.

$$\begin{split} & \mathbf{E}_{B}\left[\mathbf{E}_{p_{page}}\left[\left\|k_{i}^{t+1}\right\|^{2}\right]\right] \\ & \leq \left(2L_{i}^{2} + \frac{(1-p_{page})L_{\max}^{2}}{B}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \frac{2b^{2}}{p_{page}}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}, \quad \forall i \in [n]. \end{split}$$

557 *Proof.* First, we prove the first inequality of the theorem:

$$\begin{split} & \mathbf{E}_{B} \left[\mathbf{E}_{p_{\text{a}}} \left[\mathbf{E}_{p_{\text{page}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \right] \\ & = p_{\text{page}} \mathbf{E}_{B} \left[\mathbf{E}_{p_{\text{a}}} \left[\left\| h^{t+1}_{1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] + (1 - p_{\text{page}}) \mathbf{E}_{B} \left[\mathbf{E}_{p_{\text{a}}} \left[\left\| h^{t+1}_{2} - \nabla f(x^{t+1}) \right\|^{2} \right] \right]. \end{split}$$

558 Using

$$\begin{split} &\mathbf{E}_{B}\left[\mathbf{E}_{p_{\mathbf{a}}}\left[h_{i,1}^{t+1}\right]\right] = \\ &= p_{\mathbf{a}}h_{i}^{t} + \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - \frac{b}{p_{\mathrm{page}}}\left(h_{i}^{t} - \nabla f_{i}(x^{t})\right) + (1 - p_{\mathbf{a}})h_{i}^{t} \\ &= h_{i}^{t} + \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - \frac{b}{p_{\mathrm{page}}}\left(h_{i}^{t} - \nabla f_{i}(x^{t})\right). \end{split}$$

559 and

$$\begin{split} \mathbf{E}_{B} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[h_{i,2}^{t+1} \right] \right] &= \\ &= p_{\mathbf{a}} h_{i}^{t} + \mathbf{E}_{B} \left[\frac{1}{B} \sum_{j \in I_{i}^{t}} \left(\nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) \right) \right] + (1 - p_{\mathbf{a}}) h_{i}^{t} \\ &= h_{i}^{t} + \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}), \end{split}$$

560 we obtain

$$\begin{split} & \mathbf{E}_{B}\left[\mathbf{E}_{p_{\text{a}}}\left[\mathbf{E}_{p_{\text{page}}}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right]\right]\right] \\ & \stackrel{\text{(13)}}{=} p_{\text{page}}\mathbf{E}_{p_{\text{a}}}\left[\left\|h^{t+1}_{1} - \mathbf{E}_{p_{\text{a}}}\left[h^{t+1}_{1}\right]\right\|^{2}\right] + (1 - p_{\text{page}})\mathbf{E}_{B}\left[\mathbf{E}_{p_{\text{a}}}\left[\left\|h^{t+1}_{2} - \mathbf{E}_{B}\left[\mathbf{E}_{p_{\text{a}}}\left[h^{t+1}_{2}\right]\right]\right\|^{2}\right]\right] \end{split}$$

$$+ p_{\text{page}} \| \mathbf{E}_{p_{a}} \left[h_{1}^{t+1} \right] - \nabla f(x^{t+1}) \|^{2} + (1 - p_{\text{page}}) \| \mathbf{E}_{B} \left[\mathbf{E}_{p_{a}} \left[h_{2}^{t+1} \right] \right] - \nabla f(x^{t+1}) \|^{2}$$

$$= p_{\text{page}} \mathbf{E}_{p_{a}} \left[\| h_{1}^{t+1} - \mathbf{E}_{p_{a}} \left[h_{1}^{t+1} \right] \|^{2} \right] + (1 - p_{\text{page}}) \mathbf{E}_{B} \left[\mathbf{E}_{p_{a}} \left[\| h_{2}^{t+1} - \mathbf{E}_{B} \left[\mathbf{E}_{p_{a}} \left[h_{2}^{t+1} \right] \right] \|^{2} \right] \right]$$

$$+ \left(p_{\text{page}} \left(1 - \frac{b}{p_{\text{page}}} \right)^{2} + (1 - p_{\text{page}}) \right) \| h^{t} - \nabla f(x^{t}) \|^{2}.$$

$$(21)$$

Next, we consider $\mathbf{E}_{p_{\mathbf{a}}}\left[\left\|h_1^{t+1} - \mathbf{E}_{p_{\mathbf{a}}}\left[h_1^{t+1}\right]\right\|^2\right]$. We can use Lemma 1 with $r_i = h_i^t$ and $s_i = k_{i,1}^{t+1}$ to obtain

$$\begin{split} & \mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h_{1}^{t+1} - \mathbf{E}_{p_{\mathbf{a}}} \left[h_{1}^{t+1} \right] \right\|^{2} \right] \\ & \leq \frac{1}{n^{2} p_{\mathbf{a}}} \sum_{i=1}^{n} \left\| k_{i,1}^{t+1} - k_{i,1}^{t+1} \right\|^{2} + \frac{p_{\mathbf{a}} - p_{\mathbf{aa}}}{n^{2} p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| k_{i,1}^{t+1} \right\|^{2} \\ & = \frac{p_{\mathbf{a}} - p_{\mathbf{aa}}}{n^{2} p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - \frac{b}{p_{\mathrm{page}}} \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right\|^{2} \\ & \stackrel{\text{(12)}}{\leq} \frac{2 \left(p_{\mathbf{a}} - p_{\mathbf{aa}} \right)}{n^{2} p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} + \frac{2 \left(p_{\mathbf{a}} - p_{\mathbf{aa}} \right) b^{2}}{n^{2} p_{\mathrm{page}}^{2}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}. \end{split}$$

From Assumption 3, we have

$$\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h_{1}^{t+1} - \mathbf{E}_{p_{\mathbf{a}}} \left[h_{1}^{t+1} \right] \right\|^{2} \right] \\
\leq \frac{2 \left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}} \right) \widehat{L}^{2}}{n p_{\mathbf{a}}^{2}} \left\| x^{t+1} - x^{t} \right\|^{2} + \frac{2 \left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}} \right) b^{2}}{n^{2} p_{\mathbf{a}}^{2} p_{\mathbf{p}age}^{2}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}. \tag{22}$$

Now, we prove the bound for $E_B\left[E_{p_a}\left[\left\|h_2^{t+1}-E_B\left[E_{p_a}\left[h_2^{t+1}\right]\right]\right\|^2\right]\right]$. Considering that minibatches in the algorithm are independent, we can use Lemma 1 with $r_i=h_i^t$ and $s_i=k_{i,2}^{t+1}$ to obtain

$$\begin{split} & \mathbf{E}_{B} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h_{2}^{t+1} - \mathbf{E}_{B} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[h_{2}^{t+1} \right] \right] \right\|^{2} \right] \right] \\ & \leq \frac{1}{n^{2} p_{\mathbf{a}}} \sum_{i=1}^{n} \mathbf{E}_{B} \left[\left\| k_{i,2}^{t+1} - \mathbf{E}_{B} \left[k_{i,2}^{t+1} \right] \right\|^{2} \right] + \frac{p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}}{n^{2} p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| \mathbf{E}_{B} \left[k_{i,2}^{t+1} \right] \right\|^{2} \\ & = \frac{1}{n^{2} p_{\mathbf{a}}} \sum_{i=1}^{n} \mathbf{E}_{B} \left[\left\| \frac{1}{B} \sum_{j \in I_{i}^{t}} \left(\nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) \right) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right] \\ & + \frac{p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}}{n^{2} p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} \\ & = \frac{1}{n^{2} p_{\mathbf{a}} B^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} \\ & = \frac{1}{n^{2} p_{\mathbf{a}} B^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} \\ & = \frac{1}{n^{2} p_{\mathbf{a}} B^{2}} \sum_{i=1}^{n} \sum_{j=1}^{m} \left\| \left(\nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) \right) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \\ & + \frac{p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}}{n^{2} p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} \\ & \leq \frac{1}{n^{2} p_{\mathbf{a}} B^{2}} \sum_{i=1}^{n} \sum_{j=1}^{m} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{ij}(x^{t}) \right\|^{2} + \frac{p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}}{n^{2} p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2}. \end{split}$$

Next, we use Assumptions 3 and 4 to get

$$\mathbb{E}_{B}\left[\mathbb{E}_{p_{a}}\left[\left\|h_{2}^{t+1} - \mathbb{E}_{B}\left[\mathbb{E}_{p_{a}}\left[h_{2}^{t+1}\right]\right]\right\|^{2}\right]\right] \leq \left(\frac{L_{\max}^{2}}{np_{a}B} + \frac{(p_{a} - p_{aa})\widehat{L}^{2}}{np_{a}^{2}}\right)\left\|x^{t+1} - x^{t}\right\|^{2}.$$
(23)

568 Applying (22) and (23) into (21), we get

$$\begin{split} & \mathbf{E}_{B} \left[\mathbf{E}_{p_{\text{a}}} \left[\mathbf{E}_{p_{\text{page}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \right] \\ & \leq p_{\text{page}} \left(\frac{2 \left(p_{\text{a}} - p_{\text{aa}} \right) \hat{L}^{2}}{n p_{\text{a}}^{2}} \left\| x^{t+1} - x^{t} \right\|^{2} + \frac{2 \left(p_{\text{a}} - p_{\text{aa}} \right) b^{2}}{n^{2} p_{\text{a}}^{2} p_{\text{page}}^{2}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \right) + \\ & + \left(1 - p_{\text{page}} \right) \left(\frac{L_{\text{max}}^{2}}{n p_{\text{a}} B} + \frac{\left(p_{\text{a}} - p_{\text{aa}} \right) \hat{L}^{2}}{n p_{\text{a}}^{2}} \right) \left\| x^{t+1} - x^{t} \right\|^{2} \\ & + \left(p_{\text{page}} \left(1 - \frac{b}{p_{\text{page}}} \right)^{2} + \left(1 - p_{\text{page}} \right) \right) \left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \\ & \leq \left(\frac{2 \left(p_{\text{a}} - p_{\text{aa}} \right) \hat{L}^{2}}{n p_{\text{a}}^{2}} + \frac{\left(1 - p_{\text{page}} \right) L_{\text{max}}^{2}}{n p_{\text{a}} B} \right) \left\| x^{t+1} - x^{t} \right\|^{2} \\ & + \frac{2 \left(p_{\text{a}} - p_{\text{aa}} \right) b^{2}}{n^{2} p_{\text{apge}}^{2}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} + \left(p_{\text{page}} \left(1 - \frac{b}{p_{\text{page}}} \right)^{2} + \left(1 - p_{\text{page}} \right) \right) \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}. \end{split}$$

The proof of the second inequality almost repeats the previous one:

$$E_{B} \left[E_{p_{a}} \left[E_{p_{page}} \left[\left\| h_{i}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \right] \right] \\
= p_{page} E_{B} \left[E_{p_{a}} \left[\left\| h_{i,1}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \right] + (1 - p_{page}) E_{B} \left[E_{p_{a}} \left[\left\| h_{i,2}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \right] \\
\stackrel{\text{(13)}}{=} p_{page} E_{B} \left[E_{p_{a}} \left[\left\| h_{i,1}^{t+1} - E_{B} \left[E_{p_{a}} \left[h_{i,1}^{t+1} \right] \right] \right]^{2} \right] \right] + (1 - p_{page}) E_{B} \left[E_{p_{a}} \left[\left\| h_{i,2}^{t+1} - E_{B} \left[E_{p_{a}} \left[h_{i,2}^{t+1} \right] \right] \right]^{2} \right] \right] \\
+ p_{page} \left\| E_{B} \left[E_{p_{a}} \left[h_{i,1}^{t+1} \right] - \nabla f_{i}(x^{t+1}) \right]^{2} + (1 - p_{page}) \left\| E_{B} \left[E_{p_{a}} \left[h_{i,2}^{t+1} \right] - \nabla f_{i}(x^{t+1}) \right]^{2} \right] \right] \\
= p_{page} E_{B} \left[E_{p_{a}} \left[\left\| h_{i,1}^{t+1} - E_{B} \left[E_{p_{a}} \left[h_{i,1}^{t+1} \right] \right] \right]^{2} \right] + (1 - p_{page}) E_{B} \left[E_{p_{a}} \left[\left\| h_{i,2}^{t+1} - E_{B} \left[E_{p_{a}} \left[h_{i,2}^{t+1} \right] \right] \right]^{2} \right] \right] \\
+ \left(p_{page} \left(1 - \frac{b}{p_{page}} \right)^{2} + (1 - p_{page}) \right) \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} . \tag{24}$$

 $\text{570}\quad \text{Let us consider } \mathbf{E}_{B}\left[\mathbf{E}_{p_{\mathbf{a}}}\left[\left\|h_{i,1}^{t+1}-\mathbf{E}_{B}\left[\mathbf{E}_{p_{\mathbf{a}}}\left[h_{i,1}^{t+1}\right]\right]\right\|^{2}\right]\right]:$

$$\begin{split} & \mathbf{E}_{B} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h_{i,1}^{t+1} - \mathbf{E}_{B} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[h_{i,1}^{t+1} \right] \right] \right\|^{2} \right] \right] \\ & = \mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h_{i,1}^{t+1} - \mathbf{E}_{B} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[h_{i,1}^{t+1} \right] \right] \right\|^{2} \right] \\ & = p_{\mathbf{a}} \left\| h_{i}^{t} + \frac{1}{p_{\mathbf{a}}} k_{i,1}^{t+1} - \left(h_{i}^{t} + \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - \frac{b}{p_{\mathrm{page}}} \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right) \right\|^{2} \\ & + (1 - p_{\mathbf{a}}) \left\| h_{i}^{t} - \left(h_{i}^{t} + \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - \frac{b}{p_{\mathrm{page}}} \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right\|^{2} \\ & = \frac{(1 - p_{\mathbf{a}})^{2}}{p_{\mathbf{a}}} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - \frac{b}{p_{\mathrm{page}}} \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right\|^{2} \\ & + (1 - p_{\mathbf{a}}) \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - \frac{b}{p_{\mathrm{page}}} \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right\|^{2} \\ & = \frac{1 - p_{\mathbf{a}}}{p_{\mathbf{a}}} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - \frac{b}{p_{\mathrm{page}}} \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right\|^{2}. \end{split}$$

571 Considering (12) and Assumption 3, we obtain

$$\mathbf{E}_{B} \left[\mathbf{E}_{p_{a}} \left[\left\| h_{i,1}^{t+1} - \mathbf{E}_{B} \left[\mathbf{E}_{p_{a}} \left[h_{i,1}^{t+1} \right] \right] \right\|^{2} \right] \right] \\
\leq \frac{2 \left(1 - p_{a} \right) L_{i}^{2}}{p_{a}} \left\| x^{t+1} - x^{t} \right\|^{2} + \frac{2 \left(1 - p_{a} \right) b^{2}}{p_{a} p_{\text{page}}^{2}} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}.$$
(25)

Next, we obtain the bound for $\mathrm{E}_{B}\left[\mathrm{E}_{p_{\mathrm{a}}}\left[\left\|h_{i,2}^{t+1}-\mathrm{E}_{B}\left[\mathrm{E}_{p_{\mathrm{a}}}\left[h_{i,2}^{t+1}\right]\right]\right\|^{2}\right]\right]$:

$$E_{B} \left[E_{p_{a}} \left[\| h_{i,2}^{t+1} - E_{B} \left[E_{p_{a}} \left[h_{i,2}^{t+1} \right] \right] \right]^{2} \right] \right] \\
= p_{a} E_{B} \left[\| h_{i}^{t} + \frac{1}{p_{a}} k_{i,2}^{t+1} - \left(h_{i}^{t} + \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right]^{2} \right] \\
+ (1 - p_{a}) E_{B} \left[\| h_{i}^{t} - \left(h_{i}^{t} + \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right]^{2} \right] \\
= p_{a} E_{B} \left[\| \frac{1}{p_{a}} k_{i,2}^{t+1} - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right]^{2} \right] \\
+ (1 - p_{a}) \| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \|^{2} \\
= \frac{1}{p_{a}} E_{B} \left[\| k_{i,2}^{t+1} - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \|^{2} \right] + \frac{(1 - p_{a})^{2}}{p_{a}} \| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \|^{2} \\
+ (1 - p_{a}) \| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \|^{2} \\
= \frac{1}{p_{a}} E_{B} \left[\| k_{i,2}^{t+1} - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \|^{2} \right] + \frac{1 - p_{a}}{p_{a}} \| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \|^{2} \\
\leq \frac{1}{p_{a}} E_{B} \left[\| k_{i,2}^{t+1} - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \|^{2} \right] + \frac{(1 - p_{a}) L_{i}^{2}}{p_{a}} \| x^{t+1} - x^{t} \|^{2}, \tag{26}$$

where we used Assumption 3. By plugging (25) and (26) into (24), we get

$$\begin{split} & \mathbf{E}_{B} \left[\mathbf{E}_{p_{\text{page}}} \left[\left\| h_{i}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \right] \right] \\ & \leq p_{\text{page}} \left(\frac{2 \left(1 - p_{\text{a}} \right) L_{i}^{2}}{p_{\text{a}}} \left\| x^{t+1} - x^{t} \right\|^{2} + \frac{2 \left(1 - p_{\text{a}} \right) b^{2}}{p_{\text{a}} p_{\text{page}}^{2}} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \right) \\ & + \left(1 - p_{\text{page}} \right) \left(\frac{1}{p_{\text{a}}} \mathbf{E}_{B} \left[\left\| k_{i,2}^{t+1} - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right] + \frac{\left(1 - p_{\text{a}} \right) L_{i}^{2}}{p_{\text{a}}} \left\| x^{t+1} - x^{t} \right\|^{2} \right) \\ & + \left(p_{\text{page}} \left(1 - \frac{b}{p_{\text{page}}} \right)^{2} + \left(1 - p_{\text{page}} \right) \right) \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \\ & \leq \frac{2 \left(1 - p_{\text{a}} \right) L_{i}^{2}}{p_{\text{a}}} \left\| x^{t+1} - x^{t} \right\|^{2} + \frac{1 - p_{\text{page}}}{p_{\text{a}}} \mathbf{E}_{B} \left[\left\| k_{i,2}^{t+1} - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right] \\ & + \left(\frac{2 \left(1 - p_{\text{a}} \right) b^{2}}{p_{\text{a}} p_{\text{page}}} + p_{\text{page}} \left(1 - \frac{b}{p_{\text{page}}} \right)^{2} + \left(1 - p_{\text{page}} \right) \right) \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}. \end{split}$$

From the independence of elements in the mini-batch, we obtain

$$\begin{split} & \mathbf{E}_{B} \left[\mathbf{E}_{p_{\text{a}}} \left[\mathbf{E}_{p_{\text{page}}} \left[\left\| h_{i}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \right] \right] \\ & \leq \frac{2 \left(1 - p_{\text{a}} \right) L_{i}^{2}}{p_{\text{a}}} \left\| x^{t+1} - x^{t} \right\|^{2} + \frac{1 - p_{\text{page}}}{p_{\text{a}}} \mathbf{E}_{B} \left[\left\| \frac{1}{B} \sum_{j \in I_{i}^{t}} \left(\nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) \right) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right] \\ & + \left(\frac{2 \left(1 - p_{\text{a}} \right) b^{2}}{p_{\text{a}} p_{\text{page}}} + p_{\text{page}} \left(1 - \frac{b}{p_{\text{page}}} \right)^{2} + \left(1 - p_{\text{page}} \right) \right) \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \end{split}$$

$$\begin{split} &= \frac{2\left(1-p_{\rm a}\right)L_{i}^{2}}{p_{\rm a}}\left\|x^{t+1}-x^{t}\right\|^{2} + \frac{1-p_{\rm page}}{p_{\rm a}B^{2}}\mathbf{E}_{B}\left[\sum_{j\in I_{i}^{t}}\left\|\left(\nabla f_{ij}(x^{t+1})-\nabla f_{ij}(x^{t})\right)-\left(\nabla f_{i}(x^{t+1})-\nabla f_{i}(x^{t})\right)\right\|^{2}\right] \\ &+ \left(\frac{2\left(1-p_{\rm a}\right)b^{2}}{p_{\rm a}p_{\rm page}}+p_{\rm page}\left(1-\frac{b}{p_{\rm page}}\right)^{2}+\left(1-p_{\rm page}\right)\right)\left\|h_{i}^{t}-\nabla f_{i}(x^{t})\right\|^{2} \\ &= \frac{2\left(1-p_{\rm a}\right)L_{i}^{2}}{p_{\rm a}}\left\|x^{t+1}-x^{t}\right\|^{2}+\frac{1-p_{\rm page}}{mp_{\rm a}B}\sum_{j=1}^{m}\left\|\left(\nabla f_{ij}(x^{t+1})-\nabla f_{ij}(x^{t})\right)-\left(\nabla f_{i}(x^{t+1})-\nabla f_{i}(x^{t})\right)\right\|^{2} \\ &+ \left(\frac{2\left(1-p_{\rm a}\right)b^{2}}{p_{\rm a}p_{\rm page}}+p_{\rm page}\left(1-\frac{b}{p_{\rm page}}\right)^{2}+\left(1-p_{\rm page}\right)\right)\left\|h_{i}^{t}-\nabla f_{i}(x^{t})\right\|^{2} \\ &\leq \frac{2\left(1-p_{\rm a}\right)L_{i}^{2}}{p_{\rm a}}\left\|x^{t+1}-x^{t}\right\|^{2}+\frac{1-p_{\rm page}}{mp_{\rm a}B}\sum_{j=1}^{m}\left\|\nabla f_{ij}(x^{t+1})-\nabla f_{ij}(x^{t})\right\|^{2} \\ &+ \left(\frac{2\left(1-p_{\rm a}\right)b^{2}}{p_{\rm a}p_{\rm page}}+p_{\rm page}\left(1-\frac{b}{p_{\rm page}}\right)^{2}+\left(1-p_{\rm page}\right)\right)\left\|h_{i}^{t}-\nabla f_{i}(x^{t})\right\|^{2} \\ &\leq \left(\frac{2\left(1-p_{\rm a}\right)L_{i}^{2}}{p_{\rm a}}+\frac{\left(1-p_{\rm page}\right)L_{\rm max}^{2}}{p_{\rm a}B}\right)\left\|x^{t+1}-x^{t}\right\|^{2} \\ &+ \left(\frac{2\left(1-p_{\rm a}\right)b^{2}}{p_{\rm a}p_{\rm page}}+p_{\rm page}\left(1-\frac{b}{p_{\rm page}}\right)^{2}+\left(1-p_{\rm page}\right)\right)\left\|h_{i}^{t}-\nabla f_{i}(x^{t})\right\|^{2}, \end{split}$$

where we used Assumption 4. Finally, we prove the last inequality:

$$\begin{split} & \mathbf{E}_{B} \left[\mathbf{E}_{p_{\text{page}}} \left[\left\| k_{i}^{t+1} \right\|^{2} \right] \right] \\ & = p_{\text{page}} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - \frac{b}{p_{\text{page}}} \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right\|^{2} \\ & + (1 - p_{\text{page}}) \mathbf{E}_{B} \left[\left\| \frac{1}{B} \sum_{j \in I_{i}^{t}} \left(\nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) \right) \right\|^{2} \right] \\ & \stackrel{\text{(13)}}{=} p_{\text{page}} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - \frac{b}{p_{\text{page}}} \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right\|^{2} \\ & + (1 - p_{\text{page}}) \mathbf{E}_{B} \left[\left\| \frac{1}{B} \sum_{j \in I_{i}^{t}} \left(\nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) \right) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right] \\ & + (1 - p_{\text{page}}) \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} + \frac{2b^{2}}{p_{\text{page}}} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \\ & \leq 2 p_{\text{page}} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} \\ & \leq 2 \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} + \frac{2b^{2}}{p_{\text{page}}} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \\ & \leq 2 \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} + \frac{2b^{2}}{p_{\text{page}}} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \end{split}$$

$$+ (1 - p_{\text{page}}) \mathbf{E}_B \left[\left\| \frac{1}{B} \sum_{j \in I_i^t} \left(\nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^t) \right) - \left(\nabla f_i(x^{t+1}) - \nabla f_i(x^t) \right) \right\|^2 \right].$$

Using the independence of elements in the mini-batch, we have

$$\begin{split} & \mathbf{E}_{B} \left[\mathbf{E}_{p_{\text{page}}} \left[\left\| k_{i}^{t+1} \right\|^{2} \right] \right] \\ & \leq 2 \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} + \frac{2b^{2}}{p_{\text{page}}} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \\ & + \frac{1 - p_{\text{page}}}{B^{2}} \mathbf{E}_{B} \left[\sum_{j \in I_{i}^{t}} \left\| \left(\nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) \right) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right] \\ & = 2 \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} + \frac{2b^{2}}{p_{\text{page}}} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \\ & + \frac{1 - p_{\text{page}}}{Bm} \sum_{j=1}^{m} \left\| \left(\nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) \right) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \\ & \leq 2 \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} + \frac{2b^{2}}{p_{\text{page}}} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \\ & + \frac{1 - p_{\text{page}}}{Bm} \sum_{j=1}^{m} \left\| \nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) \right\|^{2} \end{split}$$

It it left to consider Assumptions 3 and 4 to get

$$\begin{split} & \mathbf{E}_{B}\left[\mathbf{E}_{p_{\text{page}}}\left[\left\|k_{i}^{t+1}\right\|^{2}\right]\right] \\ & \leq \left(2L_{i}^{2} + \frac{(1-p_{\text{page}})L_{\max}^{2}}{B}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \frac{2b^{2}}{p_{\text{page}}}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}. \end{split}$$

Theorem 3. Suppose that Assumptions 1, 2, 3, 4, 7, and 8 hold. Let us take $a = \frac{p_a}{2\omega+1}$, $b = \frac{p_{page}p_a}{2-p_a}$,

580 probability $p_{page} \in (0,1]$,

$$\gamma \leq \left(L + \left[\frac{48\omega(2\omega + 1)}{np_{\rm a}^2}\left(\hat{L}^2 + \frac{(1 - p_{page})L_{\rm max}^2}{B}\right) + \frac{16}{np_{\rm a}^2p_{page}}\left(\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)\hat{L}^2 + \frac{(1 - p_{page})L_{\rm max}^2}{B}\right)\right]^{1/2}\right)^{-1}$$

and $g_i^0 = h_i^0 = \nabla f_i(x^0)$ for all $i \in [n]$ in Algorithm 1 (DASHA-PP-PAGE) then $\mathrm{E}\left[\left\|\nabla f(\widehat{x}^T)\right\|^2\right] \leq 1$

582 $\frac{2\Delta_0}{\gamma T}$.

Proof. Let us fix constants $\nu, \rho \in [0, \infty)$ that we will define later. Considering Lemma 6, Lemma 8, and the law of total expectation, we obtain

$$E\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{a}} E\left[\|g^{t+1} - h^{t+1}\|^{2}\right] + \frac{\gamma((2\omega+1)p_{a} - p_{aa})}{np_{a}^{2}} E\left[\frac{1}{n}\sum_{i=1}^{n}\|g_{i}^{t+1} - h_{i}^{t+1}\|^{2}\right] \\
+ \nu E\left[\|h^{t+1} - \nabla f(x^{t+1})\|^{2}\right] + \rho E\left[\frac{1}{n}\sum_{i=1}^{n}\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\|^{2}\right] \\
\leq E\left[f(x^{t}) - \frac{\gamma}{2}\|\nabla f(x^{t})\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\|x^{t+1} - x^{t}\|^{2} + \gamma\|h^{t} - \nabla f(x^{t})\|^{2}\right] \\
+ \frac{\gamma(2\omega+1)}{p_{a}} E\left[\|g^{t} - h^{t}\|^{2}\right] + \frac{\gamma((2\omega+1)p_{a} - p_{aa})}{np_{a}^{2}} E\left[\frac{1}{n}\sum_{i=1}^{n}\|g_{i}^{t} - h_{i}^{t}\|^{2}\right]$$

$$\begin{split} & + \frac{4\gamma\omega(2\omega+1)}{np_a^2} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^n \left\| k_i^{t+1} \right\|^2 \right] \\ & + \nu \mathbf{E} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^2 \right] + \rho \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^n \left\| h_i^{t+1} - \nabla f_i(x^{t+1}) \right\|^2 \right] \\ & = \mathbf{E} \left[f(x^t) - \frac{\gamma}{2} \left\| \nabla f(x^t) \right\|^2 - \left(\frac{1}{2\gamma} - \frac{L}{2} \right) \left\| x^{t+1} - x^t \right\|^2 + \gamma \left\| h^t - \nabla f(x^t) \right\|^2 \right] \\ & + \frac{\gamma(2\omega+1)}{p_a} \mathbf{E} \left[\left\| g^t - h^t \right\|^2 \right] + \frac{\gamma((2\omega+1)p_a - p_{aa})}{np_a^2} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^n \left\| g_i^t - h_i^t \right\|^2 \right] \\ & + \frac{4\gamma\omega(2\omega+1)}{np_a^2} \mathbf{E} \left[\mathbf{E}_B \left[\mathbf{E}_{p_{\text{page}}} \left[\frac{1}{n} \sum_{i=1}^n \left\| k_i^{t+1} \right\|^2 \right] \right] \right] \\ & + \nu \mathbf{E} \left[\mathbf{E} \mathbf{E} \left[\mathbf{E}_{p_s} \left[\mathbf{E}_{p_{\text{page}}} \left[\frac{1}{n} \sum_{i=1}^n \left\| k_i^{t+1} - \nabla f(x^{t+1}) \right\|^2 \right] \right] \right] \right] \\ & + \rho \mathbf{E} \left[\mathbf{E} \mathbf{E} \left[\mathbf{E}_{p_s} \left[\mathbf{E}_{p_{\text{page}}} \left[\frac{1}{n} \sum_{i=1}^n \left\| h_i^{t+1} - \nabla f_i(x^{t+1}) \right\|^2 \right] \right] \right] \right] \\ & \leq \mathbf{E} \left[f(x^t) - \frac{\gamma}{2} \left\| \nabla f(x^t) \right\|^2 - \left(\frac{1}{2\gamma} - \frac{L}{2} \right) \left\| x^{t+1} - x^t \right\|^2 + \gamma \left\| h^t - \nabla f(x^t) \right\|^2 \right] \\ & + \frac{\gamma(2\omega+1)}{p_a} \mathbf{E} \left[\left\| g^t - h^t \right\|^2 \right] + \frac{\gamma((2\omega+1)p_a-p_{aa})}{np_a^2} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^n \left\| g_i^t - h_i^t \right\|^2 \right] \\ & + \frac{4\gamma\omega(2\omega+1)}{np_a^2} \mathbf{E} \left[\left(2\hat{L}^2 + \frac{(1-p_{\text{page}})L_{\text{max}}^2}{B} \right) \left\| x^{t+1} - x^t \right\|^2 + \frac{2b^2}{p_{\text{page}}} \frac{1}{n} \sum_{i=1}^n \left\| h_i^t - \nabla f_i(x^t) \right\|^2 \right] \\ & + \nu \mathbf{E} \left(\left(\frac{2(p_a - p_{aa})\hat{L}^2}{np_a^2} + \frac{(1-p_{\text{page}})L_{\text{max}}^2}{np_a B} \right) \left\| x^{t+1} - x^t \right\|^2 \\ & + \frac{2(p_a - p_{aa})b^2}{n^2p_{\text{page}}} \sum_{i=1}^n \left\| h_i^t - \nabla f_i(x^t) \right\|^2 + \left(p_{\text{page}} \left(1 - \frac{b}{p_{\text{page}}} \right)^2 + (1-p_{\text{page}}) \right) \left\| h^t - \nabla f(x^t) \right\|^2 \right) \\ & + \rho \mathbf{E} \left(\left(\frac{2(1-p_a)\hat{L}^2}{p_a} + \frac{(1-p_{\text{page}})L_{\text{max}}^2}{p_a B} \right) \left\| x^{t+1} - x^t \right\|^2 \\ & + \left(\frac{2(1-p_a)b^2}{p_a p_{\text{page}}} + p_{\text{page}} \left(1 - \frac{b}{p_{\text{page}}} \right)^2 + (1-p_{\text{page}}) \right) \frac{1}{n} \sum_{i=1}^n \left\| h_i^t - \nabla f_i(x^t) \right\|^2 \right) \end{aligned}$$

After rearranging the terms, we get

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1) p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \nu \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho \mathbf{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1) p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{4\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} \left(2\widehat{L}^{2} + \frac{(1 - p_{\mathbf{page}})L_{\max}^{2}}{B}\right) \end{split}$$

$$\begin{split} &-\nu\left(\frac{2\left(p_{\mathrm{a}}-p_{\mathrm{aa}}\right)\widehat{L}^{2}}{np_{\mathrm{a}}^{2}}+\frac{(1-p_{\mathrm{page}})L_{\mathrm{max}}^{2}}{np_{\mathrm{a}}B}\right)-\rho\left(\frac{2\left(1-p_{\mathrm{a}}\right)\widehat{L}^{2}}{p_{\mathrm{a}}}+\frac{(1-p_{\mathrm{page}})L_{\mathrm{max}}^{2}}{p_{\mathrm{a}}B}\right)\right)\mathrm{E}\left[\left\|x^{t+1}-x^{t}\right\|^{2}\right]\\ &+\left(\gamma+\nu\left(p_{\mathrm{page}}\left(1-\frac{b}{p_{\mathrm{page}}}\right)^{2}+(1-p_{\mathrm{page}})\right)\right)\mathrm{E}\left[\left\|h^{t}-\nabla f(x^{t})\right\|^{2}\right]\\ &+\left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{\mathrm{a}}^{2}p_{\mathrm{page}}}+\frac{2\nu\left(p_{\mathrm{a}}-p_{\mathrm{aa}}\right)b^{2}}{np_{\mathrm{a}}^{2}p_{\mathrm{page}}}\right.\\ &+\rho\left(\frac{2\left(1-p_{\mathrm{a}}\right)b^{2}}{p_{\mathrm{a}}p_{\mathrm{page}}}+p_{\mathrm{page}}\left(1-\frac{b}{p_{\mathrm{page}}}\right)^{2}+(1-p_{\mathrm{page}})\right)\right)\mathrm{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t}-\nabla f_{i}(x^{t})\right\|^{2}\right]. \end{split}$$

Due to $b = \frac{p_{\text{page}}p_{\text{a}}}{2-p_{\text{a}}} \le p_{\text{page}}$, one can show that $\left(p_{\text{page}}\left(1-\frac{b}{p_{\text{page}}}\right)^2+(1-p_{\text{page}})\right) \le 1-b$. Thus, if we take $\nu=\frac{\gamma}{b}$, then

$$\left(\gamma + \nu \left(p_{\text{page}}\left(1 - \frac{b}{p_{\text{page}}}\right)^2 + (1 - p_{\text{page}})\right)\right) \le \gamma + \nu(1 - b) = \nu,$$

586 therefore

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{4\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}}\left(2\hat{L}^{2} + \frac{(1-p_{\mathsf{page}})L_{\mathsf{max}}^{2}}{B}\right) - \rho\left(\frac{2\left(1-p_{\mathbf{a}}\right)\hat{L}^{2}}{p_{\mathbf{a}}} + \frac{(1-p_{\mathsf{page}})L_{\mathsf{max}}^{2}}{np_{\mathbf{a}}B}\right)\right) \mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}p_{\mathsf{page}}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathsf{a}\mathbf{a}}\right)b}{np_{\mathbf{a}}^{2}p_{\mathsf{page}}} + \rho\left(\frac{2\left(1-p_{\mathbf{a}}\right)b^{2}}{p_{\mathsf{a}}p_{\mathsf{page}}} + p_{\mathsf{page}}\left(1 - \frac{b}{p_{\mathsf{page}}}\right)^{2} + (1-p_{\mathsf{page}})\right)\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right]. \end{split}$$

Next, with the choice of $b=\frac{p_{\mathrm{page}}p_{\mathrm{a}}}{2-p_{\mathrm{a}}},$ we ensure that

$$\left(\frac{2(1-p_{\rm a})b^2}{p_{\rm a}p_{\rm page}} + p_{\rm page}\left(1 - \frac{b}{p_{\rm page}}\right)^2 + (1-p_{\rm page})\right) \le 1 - b.$$

If we take $\rho = \frac{8b\gamma\omega(2\omega+1)}{np_a^2p_{\text{page}}} + \frac{2\gamma(p_a-p_{\text{aa}})}{np_a^2p_{\text{page}}}$, then

$$\left(\frac{8b^2\gamma\omega(2\omega+1)}{np_{\mathrm{a}}^2p_{\mathrm{page}}} + \frac{2\gamma\left(p_{\mathrm{a}} - p_{\mathrm{aa}}\right)b}{np_{\mathrm{a}}^2p_{\mathrm{page}}} + \rho\left(\frac{2\left(1 - p_{\mathrm{a}}\right)b^2}{p_{\mathrm{a}}p_{\mathrm{page}}} + p_{\mathrm{page}}\left(1 - \frac{b}{p_{\mathrm{page}}}\right)^2 + \left(1 - p_{\mathrm{page}}\right)\right)\right) \leq \rho,$$

587 therefore

$$\mathrm{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathrm{a}}}\mathrm{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathrm{a}} - p_{\mathrm{aa}})}{np_{\mathrm{a}}^{2}}\mathrm{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right]$$

$$\begin{split} & + \frac{\gamma}{b} \mathbf{E} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{a}^{2}p_{\mathrm{page}}} + \frac{2\gamma\left(p_{\mathrm{a}} - p_{\mathrm{aa}}\right)}{np_{a}^{2}p_{\mathrm{page}}} \right) \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| h_{i}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \\ & \leq \mathbf{E} \left[f(x^{t}) \right] - \frac{\gamma}{2} \mathbf{E} \left[\left\| \nabla f(x^{t}) \right\|^{2} \right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathrm{a}}} \mathbf{E} \left[\left\| g^{t} - h^{t} \right\|^{2} \right] + \frac{\gamma((2\omega+1)p_{\mathrm{a}} - p_{\mathrm{aa}})}{np_{a}^{2}} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{4\gamma\omega(2\omega+1)}{np_{a}^{2}} \left(2\hat{L}^{2} + \frac{(1-p_{\mathrm{page}})L_{\mathrm{max}}^{2}}{B} \right) \right. \\ & - \frac{\gamma}{bnp_{\mathrm{a}}} \left(2\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}} \right) \hat{L}^{2} + \frac{(1-p_{\mathrm{page}})L_{\mathrm{max}}^{2}}{B} \right) \\ & - \left(\frac{8b\gamma\omega(2\omega+1)}{np_{a}^{3}p_{\mathrm{page}}} + \frac{2\gamma\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}} \right)}{np_{a}^{2}p_{\mathrm{page}}} \right) \left(2\left(1 - p_{\mathrm{a}} \right) \hat{L}^{2} + \frac{(1-p_{\mathrm{page}})L_{\mathrm{max}}^{2}}{B} \right) \right) \mathbf{E} \left[\left\| x^{t+1} - x^{t} \right\|^{2} \right] \\ & + \frac{\gamma}{b} \mathbf{E} \left[\left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{a}^{2}p_{\mathrm{page}}} + \frac{2\gamma\left(p_{\mathrm{a}} - p_{\mathrm{aa}}\right)}{np_{a}^{2}p_{\mathrm{page}}} \right) \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \right]. \end{split}$$

Let us simplify the inequality. First, due to $b \geq \frac{p_{\mathrm{page}}p_{\mathrm{a}}}{2}$, we have

$$\frac{\gamma}{bnp_{\mathrm{a}}}\left(2\left(1-\frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)\widehat{L}^2+\frac{(1-p_{\mathrm{page}})L_{\mathrm{max}}^2}{B}\right)\leq \frac{4\gamma}{np_{\mathrm{a}}^2p_{\mathrm{page}}}\left(\left(1-\frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)\widehat{L}^2+\frac{(1-p_{\mathrm{page}})L_{\mathrm{max}}^2}{B}\right).$$

Second, due to $b \leq p_{\rm a} p_{\rm page}$ and $p_{\rm aa} \leq p_{\rm a}^2$, we get

$$\begin{split} &\left(\frac{8b\gamma\omega(2\omega+1)}{np_{\rm a}^{3}p_{\rm page}} + \frac{2\gamma\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)}{np_{\rm a}^{2}p_{\rm page}}\right) \left(2\left(1 - p_{\rm a}\right)\widehat{L}^{2} + \frac{(1 - p_{\rm page})L_{\rm max}^{2}}{B}\right) \\ &\leq \left(\frac{8\gamma\omega(2\omega+1)}{np_{\rm a}^{2}} + \frac{2\gamma\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)}{np_{\rm a}^{2}p_{\rm page}}\right) \left(2\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)\widehat{L}^{2} + \frac{(1 - p_{\rm page})L_{\rm max}^{2}}{B}\right) \\ &\leq \frac{16\gamma\omega(2\omega+1)}{np_{\rm a}^{2}} \left(\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)\widehat{L}^{2} + \frac{(1 - p_{\rm page})L_{\rm max}^{2}}{B}\right) \\ &+ \frac{4\gamma\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)}{np_{\rm a}^{2}p_{\rm page}} \left(\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)\widehat{L}^{2} + \frac{(1 - p_{\rm page})L_{\rm max}^{2}}{B}\right) \\ &\leq \frac{16\gamma\omega(2\omega+1)}{np_{\rm a}^{2}} \left(\widehat{L}^{2} + \frac{(1 - p_{\rm page})L_{\rm max}^{2}}{B}\right) \\ &+ \frac{4\gamma}{np_{\rm a}^{2}p_{\rm page}} \left(\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)\widehat{L}^{2} + \frac{(1 - p_{\rm page})L_{\rm max}^{2}}{B}\right). \end{split}$$

Combining all bounds together, we obtain the following simplified inequality:

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}p_{\mathrm{page}}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{page}}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \end{split}$$

$$\begin{split} &-\left(\frac{1}{2\gamma}-\frac{L}{2}-\frac{24\gamma\omega(2\omega+1)}{np_{\mathrm{a}}^{2}}\left(\widehat{L}^{2}+\frac{(1-p_{\mathrm{page}})L_{\mathrm{max}}^{2}}{B}\right)\right.\\ &\left.-\frac{8\gamma}{np_{\mathrm{a}}^{2}p_{\mathrm{page}}}\left(\left(1-\frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)\widehat{L}^{2}+\frac{(1-p_{\mathrm{page}})L_{\mathrm{max}}^{2}}{B}\right)\right)\mathrm{E}\left[\left\|\boldsymbol{x}^{t+1}-\boldsymbol{x}^{t}\right\|^{2}\right]\\ &+\frac{\gamma}{b}\mathrm{E}\left[\left\|\boldsymbol{h}^{t}-\nabla f(\boldsymbol{x}^{t})\right\|^{2}\right]+\left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathrm{a}}^{2}p_{\mathrm{page}}}+\frac{2\gamma\left(p_{\mathrm{a}}-p_{\mathrm{aa}}\right)}{np_{\mathrm{a}}^{2}p_{\mathrm{page}}}\right)\mathrm{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|\boldsymbol{h}_{i}^{t}-\nabla f_{i}(\boldsymbol{x}^{t})\right\|^{2}\right]. \end{split}$$

Using Lemma 4 and the assumption about γ , we get

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}p_{\mathrm{page}}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{page}}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}p_{\mathbf{page}}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}p_{\mathbf{page}}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right]. \end{split}$$

It is left to apply Lemma 3 with

$$\begin{split} \Psi^t &= \frac{(2\omega + 1)}{p_{\rm a}} {\rm E} \left[\left\| g^t - h^t \right\|^2 \right] + \frac{((2\omega + 1) \, p_{\rm a} - p_{\rm aa})}{n p_{\rm a}^2} {\rm E} \left[\frac{1}{n} \sum_{i=1}^n \left\| g_i^t - h_i^t \right\|^2 \right] \\ &+ \frac{1}{b} {\rm E} \left[\left\| h^t - \nabla f(x^t) \right\|^2 \right] + \left(\frac{8b\omega (2\omega + 1)}{n p_{\rm a}^2 p_{\rm page}} + \frac{2 \left(p_{\rm a} - p_{\rm aa} \right)}{n p_{\rm a}^2 p_{\rm page}} \right) {\rm E} \left[\frac{1}{n} \sum_{i=1}^n \left\| h_i^t - \nabla f_i(x^t) \right\|^2 \right] \end{split}$$

592 to conclude the proof.

Corollary 1. Let the assumptions from Theorem 3 hold and $p_{page} = B/(m+B)$. Then DASHA-PP-PAGE 593 needs 594

$$T := \mathcal{O}\left(\frac{\Delta_0}{\varepsilon} \left[L + \frac{\omega}{p_a \sqrt{n}} \left(\widehat{L} + \frac{L_{\text{max}}}{\sqrt{B}} \right) + \frac{1}{p_a} \sqrt{\frac{m}{n}} \left(\frac{\mathbb{1}_{p_a} \widehat{L}}{\sqrt{B}} + \frac{L_{\text{max}}}{B} \right) \right] \right)$$
(9)

communication rounds to get an ε -solution and the expected number of gradient calculations per 595 node equals $\mathcal{O}(m+BT)$. 596

Proof. In the view of Theorem 3, it is enough to do

$$T := \mathcal{O}\left(\frac{\Delta_0}{\varepsilon}\left[L + \sqrt{\frac{\omega^2}{np_{\rm a}^2}\left(\widehat{L}^2 + \frac{(1-p_{\rm page})L_{\rm max}^2}{B}\right) + \frac{1}{np_{\rm a}^2p_{\rm page}}\left(\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)\widehat{L}^2 + \frac{(1-p_{\rm page})L_{\rm max}^2}{B}\right)\right]\right)$$

steps to get ε -solution. Using the choice of p_{mega} and the definition of $\mathbb{1}_{p_a}$, we can get (9).

Note that the expected number of gradients calculations at each communication round equals $p_{\mathrm{mega}}m+$ 599

600
$$(1 - p_{\text{mega}})B = \frac{2mB}{m+B} \le 2B$$
.

Corollary 2. Suppose that assumptions of Corollary 1 hold, $B \leq \min \left\{ \frac{1}{p_a} \sqrt{\frac{m}{n}}, \frac{L_{\max}^2}{12 \cdot \hat{L}^2} \right\}^7$, and we

use the unbiased compressor RandK with $K = \Theta(Bd/\sqrt{m})$. Then the communication complexity of

Algorithm 1 is $\mathcal{O}\left(d+\frac{L_{\max}\Delta_0 d}{p_a\varepsilon\sqrt{n}}\right)$, and the expected number of gradient calculations per node equals $\frac{\mathcal{O}\left(m+\frac{L_{\max}\Delta_0\sqrt{m}}{p_a\varepsilon\sqrt{n}}\right)}{^7\text{If }\mathbb{1}_{p_a}=0\text{, then }\frac{L_{\sigma}^2}{1_{p_a}^2L^2}=+\infty}$

604
$$\frac{\mathcal{O}\left(m + \frac{L_{\max}\Delta_0\sqrt{m}}{p_a\varepsilon\sqrt{n}}\right).}{{}^{7}\text{If }\mathbb{1}_{p_a} = 0, \text{ then } \frac{L_{\sigma}^2}{\mathbb{1}_{p_a}^2\tilde{L}^2} = +\infty}$$

605 *Proof.* The communication complexity equals

$$\mathcal{O}\left(d+KT\right) = \mathcal{O}\left(d+\frac{\Delta_0}{\varepsilon}\left[KL+K\frac{\omega}{p_{\rm a}\sqrt{n}}\left(\widehat{L}+\frac{L_{\rm max}}{\sqrt{B}}\right)+K\frac{1}{p_{\rm a}}\sqrt{\frac{m}{n}}\left(\frac{\mathbb{1}_{p_{\rm a}}\widehat{L}}{\sqrt{B}}+\frac{L_{\rm max}}{B}\right)\right]\right).$$

Since $B \leq \frac{L_{\max}^2}{\mathbb{I}_{2}^2 \hat{L}^2}$, we have $\frac{\mathbb{I}_{p_a} \hat{L}}{\sqrt{B}} + \frac{L_{\max}}{B} \leq \frac{2L_{\max}}{B}$ and

$$\mathcal{O}\left(d+KT\right) \ = \ \mathcal{O}\left(d+\frac{\Delta_0}{\varepsilon}\left[KL+K\frac{\omega}{p_{\rm a}\sqrt{n}}\left(\widehat{L}+\frac{L_{\rm max}}{\sqrt{B}}\right)+K\frac{1}{p_{\rm a}}\sqrt{\frac{m}{n}}\frac{L_{\rm max}}{B}\right]\right).$$

Note that $K = \Theta\left(\frac{Bd}{\sqrt{m}}\right) = \mathcal{O}\left(\frac{d}{p_a\sqrt{n}}\right)$ and $\omega + 1 = \frac{d}{K}$ due to Theorem 6, thus

$$\mathcal{O}(d+KT) = \mathcal{O}\left(d + \frac{\Delta_0}{\varepsilon} \left[\frac{d}{p_a \sqrt{n}} L + \frac{d}{p_a \sqrt{n}} \left(\widehat{L} + \frac{L_{\text{max}}}{\sqrt{B}} \right) + \frac{d}{p_a \sqrt{n}} L_{\text{max}} \right] \right)$$

$$= \mathcal{O}\left(d + \frac{L_{\text{max}} \Delta_0 d}{p_a \varepsilon \sqrt{n}} \right).$$

Using the same reasoning, the expected number of gradient calculations per node equals

$$\mathcal{O}(m+BT) = \mathcal{O}\left(m + \frac{\Delta_0}{\varepsilon} \left[BL + B\frac{\omega}{p_{\rm a}\sqrt{n}} \left(\widehat{L} + \frac{L_{\rm max}}{\sqrt{B}}\right) + B\frac{1}{p_{\rm a}}\sqrt{\frac{m}{n}} \left(\frac{\mathbb{1}_{p_{\rm a}}\widehat{L}}{\sqrt{B}} + \frac{L_{\rm max}}{B}\right)\right]\right)$$

$$= \mathcal{O}\left(m + \frac{\Delta_0}{\varepsilon} \left[BL + B\frac{d}{Kp_{\rm a}\sqrt{n}} \left(\widehat{L} + \frac{L_{\rm max}}{\sqrt{B}}\right) + B\frac{1}{p_{\rm a}}\sqrt{\frac{m}{n}}\frac{L_{\rm max}}{B}\right]\right)$$

$$= \mathcal{O}\left(m + \frac{\Delta_0}{\varepsilon} \left[\frac{1}{p_{\rm a}}\sqrt{\frac{m}{n}}L + \frac{\sqrt{m}}{p_{\rm a}\sqrt{n}} \left(\widehat{L} + \frac{L_{\rm max}}{\sqrt{B}}\right) + \frac{1}{p_{\rm a}}\sqrt{\frac{m}{n}}L_{\rm max}\right]\right)$$

$$= \mathcal{O}\left(m + \frac{L_{\rm max}\Delta_0\sqrt{m}}{p_{\rm a}\varepsilon\sqrt{n}}\right).$$

609

610 E.5 Proof for DASHA-PP-FINITE-MVR

Lemma 9. Suppose that Assumptions 3, 4, and 8 hold. For h_i^{t+1} , h_{ij}^{t+1} and k_i^{t+1} from Algorithm 1 (DASHA-PP-FINITE-MVR) we have

1.

$$\begin{split} & \mathbf{E}_{B} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \\ & \leq \left(\frac{2L_{\max}^{2}}{np_{\mathbf{a}}B} + \frac{2\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)\hat{L}^{2}}{np_{\mathbf{a}}^{2}} \right) \left\| x^{t+1} - x^{t} \right\|^{2} \\ & + \frac{2\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)b^{2}}{n^{2}p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} + \frac{2b^{2}}{n^{2}p_{\mathbf{a}}Bm} \sum_{i=1}^{n} \sum_{j=1}^{m} \left\| h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right\|^{2} \\ & + (1 - b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}. \end{split}$$

2.

$$\begin{split} & \mathbf{E}_{B}\left[\mathbf{E}_{p_{\mathbf{a}}}\left[\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right]\right] \\ & \leq \left(\frac{2L_{\max}^{2}}{p_{\mathbf{a}}B} + \frac{2(1-p_{\mathbf{a}})L_{i}^{2}}{p_{\mathbf{a}}}\right)\left\|x^{t+1} - x^{t}\right\|^{2} \\ & + \frac{2b^{2}}{p_{\mathbf{a}}Bm}\sum_{j=1}^{m}\left\|h_{ij}^{t} - \nabla f_{ij}(x^{t})\right\|^{2} + \left(\frac{2\left(1-p_{\mathbf{a}}\right)b^{2}}{p_{\mathbf{a}}} + (1-b)^{2}\right)\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}, \quad \forall i \in [n]. \end{split}$$

3.
$$E_{B}\left[E_{p_{a}}\left[\left\|h_{ij}^{t+1} - \nabla f_{ij}(x^{t+1})\right\|^{2}\right]\right]$$

$$\leq \frac{2\left(1 - \frac{p_{a}B}{m}\right)L_{\max}^{2}}{\frac{p_{a}B}{m}}\left\|x^{t+1} - x^{t}\right\|^{2}$$

$$+\left(\frac{2\left(1 - \frac{p_{a}B}{m}\right)b^{2}}{\frac{p_{a}B}{m}} + (1 - b)^{2}\right)\left\|h_{ij}^{t} - \nabla f_{ij}(x^{t})\right\|^{2}, \quad \forall i \in [n], \forall j \in [m].$$

4.

$$E_{B} \left[\left\| k_{i}^{t+1} \right\|^{2} \right] \\
\leq \left(\frac{2L_{\max}^{2}}{B} + 2L_{i}^{2} \right) \left\| x^{t+1} - x^{t} \right\|^{2} \\
+ \frac{2b^{2}}{Bm} \sum_{j=1}^{m} \left\| h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right\|^{2} + 2b^{2} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}, \quad \forall i \in [n].$$

613 Proof. We start by proving the first inequality. Note that

$$\begin{split} & \mathbf{E}_{B} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[h_{i}^{t+1} \right] \right] \\ & = p_{\mathbf{a}} \left(h_{i}^{t} + \frac{1}{p_{\mathbf{a}}} \mathbf{E}_{B} \left[k_{i}^{t+1} \right] \right) + (1 - p_{\mathbf{a}}) h_{i}^{t} \\ & = h_{i}^{t} + \frac{1}{m} \sum_{j=1}^{m} \frac{B}{m} \cdot \frac{m}{B} \left(\nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) - b \left(h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right) \right) + \left(1 - \frac{B}{m} \right) \cdot 0 \\ & = \nabla f_{i}(x^{t+1}) + (1 - b) \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right), \end{split}$$

614 thus

$$\begin{split} & \mathbf{E}_{B} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \\ & \stackrel{\text{(13)}}{=} \mathbf{E}_{B} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h^{t+1} - \mathbf{E}_{B} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[h^{t+1} \right] \right] \right\|^{2} \right] \right] + (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}. \end{split}$$

We can use Lemma 1 with $r_i = h_i^t$ and $s_i = k_i^{t+1}$ to obtain

$$\begin{split} & \mathbf{E}_{B} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \\ & \leq \frac{1}{n^{2}p_{\mathbf{a}}} \sum_{i=1}^{n} \mathbf{E}_{B} \left[\left\| k_{i}^{t+1} - \mathbf{E}_{B} \left[k_{i}^{t+1} \right] \right\|^{2} \right] + \frac{p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}}{n^{2}p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| \mathbf{E}_{B} \left[k_{i}^{t+1} \right] \right\|^{2} \\ & + (1 - b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \\ & = \frac{1}{n^{2}p_{\mathbf{a}}} \sum_{i=1}^{n} \mathbf{E}_{B} \left[\left\| \frac{1}{m} \sum_{j=1}^{m} k_{ij}^{t+1} - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right) \right\|^{2} \right] \\ & + \frac{p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}}{n^{2}p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right\|^{2} \\ & + (1 - b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}. \end{split}$$

Next, we again use Lemma 1 with $r_i = 0, \, s_i = \nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^t) - b \left(h_{ij}^t - \nabla f_{ij}(x^t)\right)$

617
$$p_{\rm a} = \frac{B}{m}, \, {\rm and} \, \, p_{\rm aa} = \frac{B(B-1)}{m(m-1)}$$
:

$$\mathbf{E}_{B}\left[\mathbf{E}_{p_{\mathbf{a}}}\left[\left\|h^{t+1}-\nabla f(x^{t+1})\right\|^{2}\right]\right]$$

$$\leq \frac{1}{n^{2}p_{a}} \sum_{i=1}^{n} \left(\frac{m-B}{Bm(m-1)} \sum_{j=1}^{m} \left\| \nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) - b \left(h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right) \right\|^{2} \right)$$

$$+ \frac{p_{a} - p_{aa}}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right\|^{2}$$

$$+ (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}$$

$$\leq \frac{1}{n^{2}p_{a}Bm} \sum_{i=1}^{n} \sum_{j=1}^{m} \left\| \nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) - b \left(h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right) \right\|^{2}$$

$$+ \frac{p_{a} - p_{aa}}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right\|^{2}$$

$$+ (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}$$

$$\leq \frac{2}{n^{2}p_{a}Bm} \sum_{i=1}^{n} \sum_{j=1}^{m} \left\| \nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) \right\|^{2} + \frac{2b^{2}}{n^{2}p_{a}Bm} \sum_{i=1}^{n} \sum_{j=1}^{m} \left\| h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right\|^{2}$$

$$+ \frac{2(p_{a} - p_{aa})}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} + \frac{2(p_{a} - p_{aa})}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}$$

$$+ (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2} .$$

Due to Assumptions 3 and 4, we have

$$\begin{split} & \mathbf{E}_{B} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \\ & \leq \left(\frac{2L_{\max}^{2}}{np_{\mathbf{a}}B} + \frac{2\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)\widehat{L}^{2}}{np_{\mathbf{a}}^{2}} \right) \left\| x^{t+1} - x^{t} \right\|^{2} \\ & + \frac{2\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)b^{2}}{n^{2}p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} + \frac{2b^{2}}{n^{2}p_{\mathbf{a}}Bm} \sum_{i=1}^{n} \sum_{j=1}^{m} \left\| h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right\|^{2} \\ & + (1 - b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}. \end{split}$$

619 Let us get the bound for the second inequality:

$$\begin{split} & \mathbf{E}_{B} \left[\mathbf{E}_{p_{a}} \left[\left\| h_{i}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \right] \\ & \stackrel{\textbf{(13)}}{=} \mathbf{E}_{B} \left[\mathbf{E}_{p_{a}} \left[\left\| h_{i}^{t+1} - \left(\nabla f_{i}(x^{t+1}) + (1-b)(h_{i}^{t} - \nabla f_{i}(x^{t})) \right) \right\|^{2} \right] \right] \\ & + (1-b)^{2} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \\ & = p_{a} \mathbf{E}_{B} \left[\left\| h_{i}^{t} + \frac{1}{p_{a}} k_{i}^{t+1} - \left(\nabla f_{i}(x^{t+1}) + (1-b)(h_{i}^{t} - \nabla f_{i}(x^{t})) \right) \right\|^{2} \right] \\ & + (1-p_{a}) \left\| h_{i}^{t} - \left(\nabla f_{i}(x^{t+1}) + (1-b)(h_{i}^{t} - \nabla f_{i}(x^{t})) \right) \right\|^{2} \\ & + (1-b)^{2} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \\ & \stackrel{\textbf{(13)}}{=} \frac{1}{p_{a}} \mathbf{E}_{B} \left[\left\| k_{i}^{t+1} - \mathbf{E}_{B} \left[k_{i}^{t+1} \right] \right\|^{2} \right] \\ & + \frac{1-p_{a}}{p_{a}} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b(h_{i}^{t} - \nabla f_{i}(x^{t})) \right\|^{2} \\ & + (1-b)^{2} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}. \end{split}$$

Let us use Lemma 1 with $r_i=0,$ $s_i=\nabla f_{ij}(x^{t+1})-\nabla f_{ij}(x^t)-b\left(h_{ij}^t-\nabla f_{ij}(x^t)\right),$ $p_{\mathbf{a}}=\frac{B}{m},$ and

621 $p_{aa} = \frac{B(B-1)}{m(m-1)}$:

$$\mathbf{E}_{B}\left[\mathbf{E}_{p_{\mathbf{a}}}\left[\left\|h_{i}^{t+1}-\nabla f_{i}(x^{t+1})\right\|^{2}\right]\right]$$

$$\leq \frac{1}{p_{a}} \left(\frac{m - B}{Bm(m - 1)} \sum_{j=1}^{m} \left\| \nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) - b \left(h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right) \right\|^{2} \right) \\
+ \frac{1 - p_{a}}{p_{a}} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right\|^{2} \\
+ (1 - b)^{2} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \\
\leq \frac{1}{p_{a}Bm} \sum_{j=1}^{m} \left\| \nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) - b \left(h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right) \right\|^{2} \\
+ \frac{1 - p_{a}}{p_{a}} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right\|^{2} \\
+ (1 - b)^{2} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \\
\leq \frac{2}{p_{a}Bm} \sum_{j=1}^{m} \left\| \nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) \right\|^{2} + \frac{2(1 - p_{a})}{p_{a}} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} \\
+ \frac{2b^{2}}{p_{a}Bm} \sum_{j=1}^{m} \left\| h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right\|^{2} + \left(\frac{2(1 - p_{a})b^{2}}{p_{a}} + (1 - b)^{2} \right) \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \\
\leq \left(\frac{2L_{\max}^{2}}{p_{a}B} + \frac{2(1 - p_{a})L_{i}^{2}}{p_{a}} \right) \left\| x^{t+1} - x^{t} \right\|^{2} \\
+ \frac{2b^{2}}{p_{a}Bm} \sum_{j=1}^{m} \left\| h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right\|^{2} + \left(\frac{2(1 - p_{a})b^{2}}{p_{a}} + (1 - b)^{2} \right) \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2},$$

where we used Assumptions 3 and 4. We continue the proof by considering E_B $\left[\mathbb{E}_{p_{a}}\left[\left\|h_{ij}^{t+1}-\nabla f_{ij}(x^{t+1})\right\|^{2}\right]\right]$:

$$\begin{split} & \mathbf{E}_{B} \left[\mathbf{E}_{p_{a}} \left[\left\| h_{ij}^{t+1} - \nabla f_{ij}(x^{t+1}) \right\|^{2} \right] \right] \\ & \stackrel{\textbf{(13)}}{=} \mathbf{E}_{B} \left[\mathbf{E}_{p_{a}} \left[\left\| h_{ij}^{t+1} - \left(\nabla f_{ij}(x^{t+1}) + (1-b)(h_{ij}^{t} - \nabla f_{ij}(x^{t})) \right) \right\|^{2} \right] \right] \\ & + (1-b)^{2} \left\| h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right\|^{2} \\ & = \frac{p_{a}B}{m} \mathbf{E}_{B} \left[\left\| h_{ij}^{t} + \frac{m}{Bp_{a}} \left(\nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) - b\left(h_{ij}^{t} - \nabla f_{ij}(x^{t})\right) \right) - \left(\nabla f_{ij}(x^{t+1}) + (1-b)(h_{ij}^{t} - \nabla f_{ij}(x^{t})) \right) \right\|^{2} \right] \\ & + \left(1 - \frac{p_{a}B}{m} \right) \left\| h_{ij}^{t} - \left(\nabla f_{ij}(x^{t+1}) + (1-b)(h_{ij}^{t} - \nabla f_{ij}(x^{t})) \right) \right\|^{2} \\ & + (1-b)^{2} \left\| h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right\|^{2} \end{split}$$

$$= \frac{\left(1 - \frac{p_{a}B}{m}\right)^{2}}{\frac{p_{a}B}{m}} \left\| \nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) - b(h_{ij}^{t} - \nabla f_{ij}(x^{t})) \right\|^{2} + \left(1 - \frac{p_{a}B}{m}\right) \left\| \nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) - b(h_{ij}^{t} - \nabla f_{ij}(x^{t})) \right\|^{2} + (1 - b)^{2} \left\| h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right\|^{2}$$

$$= \frac{\left(1 - \frac{p_{a}B}{m}\right)}{\frac{p_{a}B}{m}} \left\| \nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) - b(h_{ij}^{t} - \nabla f_{ij}(x^{t})) \right\|^{2}$$

$$+ (1-b)^2 \|h_{ij}^t - \nabla f_{ij}(x^t)\|^2$$

$$\stackrel{\text{(12)}}{\leq} \frac{2\left(1-\frac{p_{\mathbf{a}}B}{m}\right)}{\frac{p_{\mathbf{a}}B}{m}} \left\|\nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^t)\right\|^2 + \left(\frac{2\left(1-\frac{p_{\mathbf{a}}B}{m}\right)b^2}{\frac{p_{\mathbf{a}}B}{m}} + (1-b)^2\right) \left\|h_{ij}^t - \nabla f_{ij}(x^t)\right\|^2.$$

624 It is left to consider Assumption 4:

$$\mathbf{E}_{B} \left[\mathbf{E}_{p_{a}} \left[\left\| h_{ij}^{t+1} - \nabla f_{ij}(x^{t+1}) \right\|^{2} \right] \right] \\
\leq \frac{2 \left(1 - \frac{p_{a}B}{m} \right) L_{\max}^{2}}{\frac{p_{a}B}{m}} \left\| x^{t+1} - x^{t} \right\|^{2} + \left(\frac{2 \left(1 - \frac{p_{a}B}{m} \right) b^{2}}{\frac{p_{a}B}{m}} + (1 - b)^{2} \right) \left\| h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right\|^{2}.$$

Finally, we obtain the bound for the last inequality of the lemma:

$$\begin{aligned} & \mathbf{E}_{B} \left[\left\| k_{i}^{t+1} \right\|^{2} \right] \\ & \stackrel{\text{(13)}}{=} \mathbf{E}_{B} \left[\left\| k_{i}^{t+1} - \mathbf{E}_{B} \left[k_{i}^{t+1} \right] \right\|^{2} \right] \\ & + \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b(h_{i}^{t} - \nabla f_{i}(x^{t})) \right\|^{2}. \end{aligned}$$

626 Using Lemma 1, we get

$$\begin{split} & \mathbf{E}_{B} \left[\left\| k_{i}^{t+1} \right\|^{2} \right] \\ & \leq \frac{m-B}{Bm(m-1)} \sum_{j=1}^{m} \left\| \nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) - b \left(h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right) \right\|^{2} \\ & + \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b (h_{i}^{t} - \nabla f_{i}(x^{t})) \right\|^{2} \\ & \leq \frac{1}{Bm} \sum_{j=1}^{m} \left\| \nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) - b \left(h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right) \right\|^{2} \\ & + \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b (h_{i}^{t} - \nabla f_{i}(x^{t})) \right\|^{2} \\ & \leq \frac{2}{Bm} \sum_{j=1}^{m} \left\| \nabla f_{ij}(x^{t+1}) - \nabla f_{ij}(x^{t}) \right\|^{2} + 2 \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} \\ & + \frac{2b^{2}}{Bm} \sum_{j=1}^{m} \left\| h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right\|^{2} + 2b^{2} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \\ & \leq \left(\frac{2L_{\max}^{2}}{B} + 2L_{i}^{2} \right) \left\| x^{t+1} - x^{t} \right\|^{2} \\ & + \frac{2b^{2}}{Bm} \sum_{j=1}^{m} \left\| h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right\|^{2} + 2b^{2} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}, \end{split}$$

where we used Assumptions 3 and 4.

Theorem 7. Suppose that Assumptions 1, 2, 3, 4, 7, and 8 hold. Let us take $a = \frac{p_a}{2\omega+1}$, $b = \frac{\frac{p_aB}{m}}{2-\frac{p_aB}{m}}$.

$$\gamma \leq \left(L + \sqrt{\frac{148\omega(2\omega+1)}{np_{\mathrm{a}}^2}\left(\widehat{L}^2 + \frac{L_{\mathrm{max}}^2}{B}\right) + \frac{72m}{np_{\mathrm{a}}^2B}\left(\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)\widehat{L}^2 + \frac{L_{\mathrm{max}}^2}{B}\right)}\right)^{-1},$$

 $g_{i}^{0} = h_{i}^{0} = \nabla f_{i}(x^{0}) \text{ for all } i \in [n] \text{ and } h_{ij}^{0} = \nabla f_{ij}(x^{0}) \text{ for all } i \in [n], j \in [m] \text{ in Algorithm 1}$ $\text{629} \quad (\mathsf{DASHA-PP-FINITE-MVR}) \text{ then } \mathrm{E}\left[\left\|\nabla f(\widehat{x}^{T})\right\|^{2}\right] \leq \frac{2\Delta_{0}}{\gamma T}.$

Proof. Let us fix constants $\nu, \rho, \delta \in [0, \infty)$ that we will define later. Considering Lemma 6, Lemma 9, and the law of total expectation, we obtain

$$\mathrm{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathrm{a}}}\mathrm{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathrm{a}} - p_{\mathrm{aa}})}{np_{\mathrm{a}}^{2}}\mathrm{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right]$$

$$\begin{split} &+\nu \mathbb{E}\left[\left\|h_{i}^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ &+\delta \mathbb{E}\left[\frac{1}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}\left\|h_{ij}^{t+1} - \nabla f_{ij}(x^{t+1})\right\|^{2}\right] \\ &\leq \mathbb{E}\left[f(x^{t}) - \frac{2}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ &+\frac{\gamma(2\omega+1)}{p_{a}}\mathbb{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)p_{a} - p_{ab})}{np_{a}^{2}}\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f(x^{t})\right\|^{2}\right] \\ &+\frac{4\gamma\omega(2\omega+1)}{p_{a}}\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1}\right\|^{2}\right] \\ &+\nu\mathbb{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ &+\delta\mathbb{E}\left[\frac{1}{mm}\sum_{i=1}^{n}\sum_{j=1}^{m}\left\|h_{i}^{t+1} - \nabla f_{ij}(x^{t+1})\right\|^{2}\right] \\ &=\mathbb{E}\left[f(x^{t}) - \frac{2}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ &+\frac{\gamma(2\omega+1)}{p_{a}}\mathbb{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)p_{a} - p_{ab})}{np_{a}^{2}}\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g^{t} - h^{t}_{i}\right\|^{2}\right] \\ &+\nu\mathbb{E}\left[\mathbb{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right]\right] \\ &+\nu\mathbb{E}\left[\mathbb{E}\left[\mathbb{E}\left[\frac{1}{p_{i}}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right]\right] \right] \\ &+\rho\mathbb{E}\left[\mathbb{E}\left[\mathbb{E}\left[\frac{1}{p_{i}}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right]\right] \right] \\ &+\rho\mathbb{E}\left[\mathbb{E}\left[\mathbb{E}\left[\frac{1}{p_{i}}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right]\right] \right] \\ &+\rho\mathbb{E}\left[\mathbb{E}\left[\frac{1}{p_{i}}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \right] \\ &+\frac{\gamma(2\omega+1)}{p_{a}}\mathbb{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)p_{a} - p_{ab})}{np_{a}^{2}}\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ &+\frac{\gamma(2\omega+1)}{p_{a}}\mathbb{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)p_{a} - p_{ab})}{np_{a}^{2}}\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ &+\frac{\gamma(2\omega+1)}{p_{a}}\mathbb{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma(2\omega+1)p_{a} - p_{ab}}{np_{a}^{2}}\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ &+\frac{\gamma(2\omega+1)}{p_{a}}\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ &+\frac{\gamma(2\omega+1)}{p_{a}}\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ &$$

$$+ \rho \mathbf{E} \left(\left(\frac{2L_{\max}^{2}}{p_{\mathbf{a}}B} + \frac{2(1-p_{\mathbf{a}})\widehat{L}^{2}}{p_{\mathbf{a}}} \right) \|x^{t+1} - x^{t}\|^{2} \right)$$

$$+ \frac{2b^{2}}{p_{\mathbf{a}}Bnm} \sum_{i=1}^{n} \sum_{j=1}^{m} \|h_{ij}^{t} - \nabla f_{ij}(x^{t})\|^{2} + \left(\frac{2(1-p_{\mathbf{a}})b^{2}}{p_{\mathbf{a}}} + (1-b)^{2} \right) \frac{1}{n} \sum_{i=1}^{n} \|h_{i}^{t} - \nabla f_{i}(x^{t})\|^{2}$$

$$+ \delta \mathbf{E} \left(\frac{2\left(1 - \frac{p_{\mathbf{a}}B}{m}\right) L_{\max}^{2}}{\frac{p_{\mathbf{a}}B}{m}} \|x^{t+1} - x^{t}\|^{2} \right)$$

$$+ \left(\frac{2\left(1 - \frac{p_{\mathbf{a}}B}{m}\right)b^{2}}{\frac{p_{\mathbf{a}}B}{m}} + (1-b)^{2} \right) \frac{1}{nm} \sum_{i=1}^{n} \sum_{j=1}^{m} \|h_{ij}^{t} - \nabla f_{ij}(x^{t})\|^{2} \right).$$

Due to $b=\frac{\frac{p_{\mathrm{a}}B}{m}}{2-\frac{p_{\mathrm{a}}B}{m}}\leq \frac{p_{\mathrm{a}}}{2-p_{\mathrm{a}}},$ we have

$$\left(\frac{2\left(1 - \frac{p_a B}{m}\right)b^2}{\frac{p_a B}{m}} + (1 - b)^2\right) \le 1 - b$$

and

$$\left(\frac{2(1-p_{\rm a})b^2}{p_{\rm a}} + (1-b)^2\right) \le 1 - b.$$

Moreover, we consider that $1 - \frac{p_a B}{m} \le 1$, therefore

$$\begin{split} & \mathbb{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbb{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1) \, p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|g^{t+1}_{i} - h^{t+1}_{i}\right\|^{2}\right] \\ & + \nu \mathbb{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|h^{t+1}_{i} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & + \delta \mathbb{E}\left[\frac{1}{nm} \sum_{i=1}^{n} \sum_{j=1}^{m} \left\|h^{t+1}_{ij} - \nabla f_{ij}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbb{E}\left[f(x^{t}) - \frac{\gamma}{2} \left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right) \left\|x^{t+1} - x^{t}\right\|^{2} + \gamma \left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbb{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1) \, p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|g^{t}_{i} - h^{t}_{i}\right\|^{2}\right] \\ & + \frac{4\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} \mathbb{E}\left[\left(\frac{2L_{\max}^{2}}{B} + 2\hat{L}^{2}\right) \left\|x^{t+1} - x^{t}\right\|^{2} + \frac{2b^{2}}{Bmn} \sum_{i=1}^{n} \sum_{j=1}^{m} \left\|h^{t}_{ij} - \nabla f_{ij}(x^{t})\right\|^{2} + \frac{2b^{2}}{n} \sum_{i=1}^{n} \left\|h^{t}_{i} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \nu \mathbb{E}\left(\left(\frac{2L_{\max}^{2}}{np_{\mathbf{a}}B} + \frac{2(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})\hat{L}^{2}}{np_{\mathbf{a}}^{2}}\right) \left\|x^{t+1} - x^{t}\right\|^{2} \\ & + \frac{2(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})b^{2}}{n^{2}p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\|h^{t}_{i} - \nabla f_{i}(x^{t})\right\|^{2} + \frac{2b^{2}}{n^{2}p_{\mathbf{a}}Bm} \sum_{i=1}^{n} \sum_{j=1}^{m} \left\|h^{t}_{ij} - \nabla f_{ij}(x^{t})\right\|^{2} \\ & + (1 - b)^{2} \left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right) \\ & + \rho \mathbb{E}\left(\left(\frac{2L_{\max}^{2}}{p_{\mathbf{a}}B} + \frac{2(1 - p_{\mathbf{a}})\hat{L}^{2}}{p_{\mathbf{a}}}\right) \left\|x^{t+1} - x^{t}\right\|^{2}\right) \end{aligned}$$

$$+ \frac{2b^{2}}{p_{a}Bnm} \sum_{i=1}^{n} \sum_{j=1}^{m} \left\| h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right\|^{2} + (1-b) \frac{1}{n} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}$$

$$+ \delta E \left(\frac{2mL_{\max}^{2}}{p_{a}B} \left\| x^{t+1} - x^{t} \right\|^{2} + (1-b) \frac{1}{nm} \sum_{i=1}^{n} \sum_{j=1}^{m} \left\| h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right\|^{2} \right).$$

633 After rearranging the terms, we get

$$\begin{split} & \mathbb{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbb{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \nu \mathbb{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & + \delta \mathbb{E}\left[\frac{1}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}\left\|h_{ij}^{t+1} - \nabla f_{ij}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbb{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbb{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbb{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{4\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}}\left(\frac{2L_{\max}^{2}}{B} + 2\widehat{L}^{2}\right) - \rho\left(\frac{2L_{\max}^{2}}{p_{\mathbf{a}}B} + \frac{2(1-p_{\mathbf{a}})\widehat{L}^{2}}{p_{\mathbf{a}}}\right) - \delta\frac{2mL_{\max}^{2}}{p_{\mathbf{a}}B}\right)\mathbb{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \left(\gamma + \nu\left(1 - b\right)^{2}\right)\mathbb{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\nu\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)b^{2}}{np_{\mathbf{a}}^{2}} + \rho\left(1 - b\right)\right)\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}B} + \frac{2\nu b^{2}}{np_{\mathbf{a}}B} + \delta\left(1 - b\right)\right)\mathbb{E}\left[\frac{1}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}\left\|h_{ij}^{t} - \nabla f_{ij}(x^{t})\right\|^{2}\right]. \end{split}$$

Thus, if we take $\nu=\frac{\gamma}{b}$, then $\gamma+\nu\left(1-b\right)^2\leq \nu$ and

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & + \delta \mathbf{E}\left[\frac{1}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}\left\|h_{ij}^{t+1} - \nabla f_{ij}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{4\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}}\left(\frac{2L_{\max}^{2}}{B} + 2\widehat{L}^{2}\right) \end{split}$$

$$-\left(\frac{2\gamma L_{\max}^{2}}{bnp_{a}B} + \frac{2\gamma\left(p_{a} - p_{aa}\right)\widehat{L}^{2}}{bnp_{a}^{2}}\right) - \rho\left(\frac{2L_{\max}^{2}}{p_{a}B} + \frac{2(1 - p_{a})\widehat{L}^{2}}{p_{a}}\right) - \delta\frac{2mL_{\max}^{2}}{p_{a}B}\right) \mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] + \frac{\gamma}{b}\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(\frac{8b^{2}\gamma\omega(2\omega + 1)}{np_{a}^{2}} + \frac{2\gamma\left(p_{a} - p_{aa}\right)b}{np_{a}^{2}} + \rho\left(1 - b\right)\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] + \left(\frac{8b^{2}\gamma\omega(2\omega + 1)}{np_{a}^{2}B} + \frac{2\gamma b}{np_{a}B} + \frac{2\rho b^{2}}{p_{a}B} + \delta\left(1 - b\right)\right)\mathbf{E}\left[\frac{1}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}\left\|h_{ij}^{t} - \nabla f_{ij}(x^{t})\right\|^{2}\right].$$

Next, if we take $ho=rac{8b\gamma\omega(2\omega+1)}{np_a^2}+rac{2\gamma(p_a-p_{aa})}{np_a^2},$ then

$$\left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{a}^{2}}+\frac{2\gamma\left(p_{a}-p_{aa}\right)b}{np_{a}^{2}}+\rho\left(1-b\right)\right)=\rho,$$

635 therefore

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & + \delta \mathbf{E}\left[\frac{1}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}\left\|h_{ij}^{t+1} - \nabla f_{ij}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{4\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} \left(\frac{2L_{\max}^{2}}{B} + 2\hat{L}^{2}\right) - \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\right) \left(\frac{2L_{\max}^{2}}{p_{\mathbf{a}}B} + \frac{2(1-p_{\mathbf{a}})\hat{L}^{2}}{p_{\mathbf{a}}}\right) \\ & - \delta\frac{2mL_{\max}^{2}}{p_{\mathbf{a}}B} \mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}B} + \frac{2\gamma b}{np_{\mathbf{a}}B} + \frac{16b^{3}\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}B} + \frac{4b^{2}\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{nBp_{\mathbf{a}}^{3}} + \delta\left(1-b\right)\right) \mathbf{E}\left[\frac{1}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}\left\|h_{ij}^{t} - \nabla f_{ij}(x^{t})\right\|^{2}\right]. \end{split}$$

Due to $b \le p_{\rm a}$ and $\frac{p_{\rm a}-p_{\rm aa}}{p_{\rm a}} \le 1$, we have

$$\begin{split} & \frac{8b^{2}\gamma\omega(2\omega+1)}{np_{\mathrm{a}}^{2}B} + \frac{2\gamma b}{np_{\mathrm{a}}B} + \frac{16b^{3}\gamma\omega(2\omega+1)}{np_{\mathrm{a}}^{3}B} + \frac{4b^{2}\gamma\left(p_{\mathrm{a}} - p_{\mathrm{aa}}\right)}{nBp_{\mathrm{a}}^{3}} \\ & \leq \frac{8b^{2}\gamma\omega(2\omega+1)}{np_{\mathrm{a}}^{2}B} + \frac{2\gamma b}{np_{\mathrm{a}}B} + \frac{16b^{2}\gamma\omega(2\omega+1)}{np_{\mathrm{a}}^{2}B} + \frac{4\gamma b}{np_{\mathrm{a}}B} \end{split}$$

$$=\frac{24b^2\gamma\omega(2\omega+1)}{np_a^2B}+\frac{6\gamma b}{np_aB}$$

Easy Let us take $\delta=\frac{24b\gamma\omega(2\omega+1)}{np_a^2B}+\frac{6\gamma}{np_aB}.$ Thus

$$\left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{a}^{2}B}+\frac{2\gamma b}{np_{a}B}+\frac{16b^{3}\gamma\omega(2\omega+1)}{np_{a}^{3}B}+\frac{4b^{2}\gamma\left(p_{a}-p_{aa}\right)}{nBp_{a}^{3}}+\delta\left(1-b\right)\right)\leq\delta$$

638 and

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{\gamma}{b}\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & + \left(\frac{24b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}B} + \frac{6\gamma}{np_{\mathbf{a}}B}\right)\mathbf{E}\left[\frac{1}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}\left\|h_{ij}^{t+1} - \nabla f_{ij}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{4\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}}\left(\frac{2L_{\max}^{2}}{B} + 2\hat{L}^{2}\right)\right) \\ & - \left(\frac{2h\gamma\omega(2\omega+1)}{np_{\mathbf{a}}B} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)\hat{L}^{2}}{np_{\mathbf{a}}B}\right) - \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right)\left(\frac{2L_{\max}^{2}}{p_{\mathbf{a}}B} + \frac{2(1-p_{\mathbf{a}})\hat{L}^{2}}{p_{\mathbf{a}}}\right) \\ & - \left(\frac{24b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}B} + \frac{6\gamma}{np_{\mathbf{a}}B}\right)\frac{2mL_{\max}^{2}}{p_{\mathbf{a}}B}\right)\mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}B} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{24b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}B} + \frac{6\gamma}{np_{\mathbf{a}}B}\right)\mathbf{E}\left[\frac{1}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}\left\|h_{ij}^{t} - \nabla f_{ij}(x^{t})\right\|^{2}\right]. \end{split}$$

Let us simplify the term near $\mathrm{E}\left[\left\|x^{t+1}-x^{t}\right\|^{2}\right]$. Due to $b\leq p_{\mathrm{a}},\,\frac{p_{\mathrm{a}}-p_{\mathrm{aa}}}{p_{\mathrm{a}}}\leq 1,\,\mathrm{and}\,\,1-p_{\mathrm{a}}\leq 1,\,\mathrm{we}$

$$\begin{split} &\frac{4\gamma\omega(2\omega+1)}{np_{\mathrm{a}}^2}\left(\frac{2L_{\mathrm{max}}^2}{B}+2\widehat{L}^2\right)\\ &+\left(\frac{2\gamma L_{\mathrm{max}}^2}{bnp_{\mathrm{a}}B}+\frac{2\gamma\left(p_{\mathrm{a}}-p_{\mathrm{aa}}\right)\widehat{L}^2}{bnp_{\mathrm{a}}^2}\right)\\ &+\left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathrm{a}}^2}+\frac{2\gamma\left(p_{\mathrm{a}}-p_{\mathrm{aa}}\right)}{np_{\mathrm{a}}^2}\right)\left(\frac{2L_{\mathrm{max}}^2}{p_{\mathrm{a}}B}+\frac{2(1-p_{\mathrm{a}})\widehat{L}^2}{p_{\mathrm{a}}}\right)\\ &+\left(\frac{24b\gamma\omega(2\omega+1)}{np_{\mathrm{a}}^2B}+\frac{6\gamma}{np_{\mathrm{a}}B}\right)\frac{2mL_{\mathrm{max}}^2}{p_{\mathrm{a}}B}\\ &\leq \frac{12\gamma\omega(2\omega+1)}{np_{\mathrm{a}}^2}\left(\frac{2L_{\mathrm{max}}^2}{B}+2\widehat{L}^2\right) \end{split}$$

$$\begin{split} & + \left(\frac{6\gamma L_{\max}^2}{bnp_{a}B} + \frac{6\gamma\left(p_{a} - p_{aa}\right)\widehat{L}^2}{bnp_{a}^2}\right) \\ & + \left(\frac{24b\gamma\omega(2\omega + 1)}{np_{a}^2B} + \frac{6\gamma}{np_{a}B}\right)\frac{2mL_{\max}^2}{p_{a}B} \end{split}$$

Considering that $b \leq \frac{p_{\mathrm{a}}B}{m}$ and $b \geq \frac{p_{\mathrm{a}}B}{2m}$, we obtain

$$\begin{split} & \frac{4\gamma\omega(2\omega+1)}{np_{\rm a}^2} \left(\frac{2L_{\rm max}^2}{B} + 2\hat{L}^2\right) \\ & + \left(\frac{2\gamma L_{\rm max}^2}{bnp_{\rm a}B} + \frac{2\gamma\left(p_{\rm a} - p_{\rm aa}\right)\hat{L}^2}{bnp_{\rm a}^2}\right) \\ & + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\rm a}^2} + \frac{2\gamma\left(p_{\rm a} - p_{\rm aa}\right)}{np_{\rm a}^2}\right) \left(\frac{2L_{\rm max}^2}{p_{\rm a}B} + \frac{2(1-p_{\rm a})\hat{L}^2}{p_{\rm a}}\right) \\ & + \left(\frac{24b\gamma\omega(2\omega+1)}{np_{\rm a}^2B} + \frac{6\gamma}{np_{\rm a}B}\right) \frac{2mL_{\rm max}^2}{p_{\rm a}B} \\ & \leq \frac{36\gamma\omega(2\omega+1)}{np_{\rm a}^2} \left(\frac{2L_{\rm max}^2}{B} + 2\hat{L}^2\right) + \left(\frac{18\gamma L_{\rm max}^2}{bnp_{\rm a}B} + \frac{6\gamma\left(p_{\rm a} - p_{\rm aa}\right)\hat{L}^2}{bnp_{\rm a}^2}\right) \\ & \leq \frac{36\gamma\omega(2\omega+1)}{np_{\rm a}^2} \left(\frac{2L_{\rm max}^2}{B} + 2\hat{L}^2\right) + \left(\frac{36m\gamma L_{\rm max}^2}{np_{\rm a}^2B^2} + \frac{12m\gamma\left(p_{\rm a} - p_{\rm aa}\right)\hat{L}^2}{Bnp_{\rm a}^3}\right). \end{split}$$

642 All in all, we have

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & + \left(\frac{24b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}B} + \frac{6\gamma}{np_{\mathbf{a}}B}\right) \mathbf{E}\left[\frac{1}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}\left\|h_{ij}^{t+1} - \nabla f_{ij}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{36\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} \left(\frac{2L_{\max}^{2}}{B} + 2\widehat{L}^{2}\right) - \left(\frac{36m\gamma L_{\max}^{2}}{np_{\mathbf{a}}^{2}B^{2}} + \frac{12m\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)\widehat{L}^{2}}{Bnp_{\mathbf{a}}^{3}}\right)\right) \mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{24b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}B} + \frac{6\gamma}{np_{\mathbf{a}}B}\right) \mathbf{E}\left[\frac{1}{nm}\sum_{i=1}^{n}\sum_{j=1}^{m}\left\|h_{ij}^{t} - \nabla f_{ij}(x^{t})\right\|^{2}\right]. \end{split}$$

Using Lemma 4 and the assumption about γ , we get

$$\mathrm{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathrm{a}}}\mathrm{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathrm{a}} - p_{\mathrm{aa}})}{np_{\mathrm{a}}^{2}}\mathrm{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right]$$

$$\begin{split} & + \frac{\gamma}{b} \mathbf{E} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}} \right) \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| h_{i}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \\ & + \left(\frac{24b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}B} + \frac{6\gamma}{np_{\mathbf{a}}B} \right) \mathbf{E} \left[\frac{1}{nm} \sum_{i=1}^{n} \sum_{j=1}^{m} \left\| h_{ij}^{t+1} - \nabla f_{ij}(x^{t+1}) \right\|^{2} \right] \\ & \leq \mathbf{E} \left[f(x^{t}) \right] - \frac{\gamma}{2} \mathbf{E} \left[\left\| \nabla f(x^{t}) \right\|^{2} \right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E} \left[\left\| g^{t} - h^{t} \right\|^{2} \right] + \frac{\gamma((2\omega+1)p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \right] \\ & + \frac{\gamma}{b} \mathbf{E} \left[\left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \right] \\ & + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}} \right) \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \right] \\ & + \left(\frac{24b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}B} + \frac{6\gamma}{np_{\mathbf{a}}B} \right) \mathbf{E} \left[\frac{1}{nm} \sum_{i=1}^{n} \sum_{j=1}^{m} \left\| h_{ij}^{t} - \nabla f_{ij}(x^{t}) \right\|^{2} \right]. \end{split}$$

644 It is left to apply Lemma 3 with

$$\Psi^{t} = \frac{(2\omega + 1)}{p_{a}} E\left[\|g^{t} - h^{t}\|^{2} \right] + \frac{((2\omega + 1) p_{a} - p_{aa})}{np_{a}^{2}} E\left[\frac{1}{n} \sum_{i=1}^{n} \|g_{i}^{t} - h_{i}^{t}\|^{2} \right]
+ \frac{1}{b} E\left[\|h^{t} - \nabla f(x^{t})\|^{2} \right]
+ \left(\frac{8b\omega(2\omega + 1)}{np_{a}^{2}} + \frac{2(p_{a} - p_{aa})}{np_{a}^{2}} \right) E\left[\frac{1}{n} \sum_{i=1}^{n} \|h_{i}^{t} - \nabla f_{i}(x^{t})\|^{2} \right]
+ \left(\frac{24b\omega(2\omega + 1)}{np_{a}^{2}B} + \frac{6}{np_{a}B} \right) E\left[\frac{1}{nm} \sum_{i=1}^{n} \sum_{j=1}^{m} \|h_{ij}^{t} - \nabla f_{ij}(x^{t})\|^{2} \right]$$

645 to conclude the proof.

646 E.6 Proof for DASHA-PP-MVR

Let us denote $\nabla f_i(x^{t+1}; \xi_i^{t+1}) := \frac{1}{B} \sum_{j=1}^B \nabla f_i(x^{t+1}; \xi_{ij}^{t+1}).$

Lemma 10. Suppose that Assumptions 3, 5, 6 and 8 hold. For h_i^{t+1} and k_i^{t+1} from Algorithm 1

649 (DASHA-PP-MVR) we have

$$\begin{split} I. \\ & \qquad \qquad \mathbf{E}_{k} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \\ & \leq \frac{2b^{2}\sigma^{2}}{np_{\mathbf{a}}B} + \left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{np_{\mathbf{a}}B} + \frac{2\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)\widehat{L}^{2}}{np_{\mathbf{a}}^{2}} \right) \left\| x^{t+1} - x^{t} \right\|^{2} \\ & \qquad \qquad + \frac{2\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)b^{2}}{n^{2}p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} + (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}. \end{split}$$

2.
$$E_{k} \left[E_{p_{a}} \left[\left\| h_{i}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \right]$$

$$\leq \frac{2b^{2}\sigma^{2}}{p_{a}B} + \left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{p_{a}B} + \frac{2(1-p_{a})L_{i}^{2}}{p_{a}} \right) \left\| x^{t+1} - x^{t} \right\|^{2}$$

$$+ \left(\frac{2(1-p_{a})b^{2}}{p_{a}} + (1-b)^{2} \right) \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}, \quad \forall i \in [n].$$

$$\mathbb{E}_{k}\left[\left\|k_{i}^{t+1}\right\|^{2}\right] \leq \frac{2b^{2}\sigma^{2}}{B} + \left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{B} + 2L_{i}^{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + 2b^{2}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}, \quad \forall i \in [n].$$

Proof. First, let us proof the bound for $\mathbf{E}_k \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^2 \right] \right]$:

651 Using

$$\mathbf{E}_k\left[\mathbf{E}_{p_{\mathbf{a}}}\left[h_i^{t+1}\right]\right] = h_i^t + \mathbf{E}_k\left[k_i^{t+1}\right] = h_i^t + \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t))$$

and (13), we have

$$\begin{split} & \mathbf{E}_{k} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \\ & = \mathbf{E}_{k} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h^{t+1} - \mathbf{E}_{k} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[h^{t+1} \right] \right] \right\|^{2} \right] \right] + (1 - b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}. \end{split}$$

We can use Lemma 1 with $r_i = h_i^t$ and $s_i = k_i^{t+1}$ to obtain

$$\begin{split} & \operatorname{E}_{k} \left[\operatorname{E}_{p_{a}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \\ & \leq \frac{1}{n^{2}p_{a}} \sum_{i=1}^{n} \operatorname{E}_{k} \left[\left\| k_{i}^{t+1} - \operatorname{E}_{k} \left[k_{i}^{t+1} \right] \right\|^{2} \right] + \frac{p_{a} - p_{aa}}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \left\| \operatorname{E}_{k} \left[k_{i}^{t+1} \right] \right\|^{2} + (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \\ & = \frac{1}{n^{2}p_{a}} \sum_{i=1}^{n} \operatorname{E}_{k} \left[\left\| \nabla f_{i}(x^{t+1}; \xi_{i}^{t+1}) - \nabla f_{i}(x^{t}; \xi_{i}^{t+1}) - b \left(h_{i}^{t} - \nabla f_{i}(x^{t}; \xi_{i}^{t+1}) \right) \right. \\ & - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right) \right\|^{2} \right] \\ & + \frac{p_{a} - p_{aa}}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right\|^{2} \\ & + (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \\ & \leq \frac{2}{n^{2}p_{a}} \sum_{i=1}^{n} \operatorname{E}_{k} \left[\left\| b \left(\nabla f_{i}(x^{t+1}; \xi_{i}^{t+1}) - \nabla f_{i}(x^{t+1}) \right) \right\|^{2} \right] \\ & + \frac{2}{n^{2}p_{a}} \sum_{i=1}^{n} \operatorname{E}_{k} \left[\left\| (1-b) \left(\nabla f_{i}(x^{t+1}; \xi_{i}^{t+1}) - \nabla f_{i}(x^{t}; \xi_{i}^{t+1}) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right) \right\|^{2} \right] \\ & + \frac{p_{a} - p_{aa}}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right\|^{2} \\ & = \frac{2b^{2}}{n^{2}p_{a}} \sum_{i=1}^{n} \operatorname{E}_{k} \left[\left\| \nabla f_{i}(x^{t+1}; \xi_{i}^{t+1}) - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \\ & + \frac{2(1-b)^{2}}{n^{2}p_{a}} \sum_{i=1}^{n} \operatorname{E}_{k} \left[\left\| \nabla f_{i}(x^{t+1}; \xi_{i}^{t+1}) - \nabla f_{i}(x^{t}; \xi_{i}^{t+1}) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right] \\ & + \frac{p_{a} - p_{aa}}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}; \xi_{i}^{t+1}) - \nabla f_{i}(x^{t}; \xi_{i}^{t+1}) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \\ & + (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \\ & + (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \right. \end{aligned}$$

$$\begin{split} &= \frac{2b^2}{n^2 p_{\mathsf{a}} B^2} \sum_{i=1}^n \sum_{j=1}^B \mathbf{E}_k \left[\left\| \nabla f_i(x^{t+1}; \xi_{ij}^{t+1}) - \nabla f_i(x^{t+1}) \right\|^2 \right] \\ &\quad + \frac{2(1-b)^2}{n^2 p_{\mathsf{a}}} \sum_{i=1}^n \mathbf{E}_k \left[\left\| \nabla f_i(x^{t+1}; \xi_i^{t+1}) - \nabla f_i(x^t; \xi_i^{t+1}) - \left(\nabla f_i(x^{t+1}) - \nabla f_i(x^t) \right) \right\|^2 \right] \\ &\quad + \frac{p_{\mathsf{a}} - p_{\mathsf{a}\mathsf{a}}}{n^2 p_{\mathsf{a}}^2} \sum_{i=1}^n \left\| \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t)) \right\|^2 \\ &\quad + (1-b)^2 \left\| h^t - \nabla f(x^t) \right\|^2. \end{split}$$

In the last equality, we use the independence of elements in the mini-batches. Due to Assumption 5, we get

$$\begin{split} & \operatorname{E}_{k} \left[\operatorname{E}_{p_{a}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \\ & \leq \frac{2b^{2}\sigma^{2}}{np_{a}B} \\ & + \frac{2(1-b)^{2}}{n^{2}p_{a}} \sum_{i=1}^{n} \operatorname{E}_{k} \left[\left\| \nabla f_{i}(x^{t+1}; \xi_{i}^{t+1}) - \nabla f_{i}(x^{t}; \xi_{i}^{t+1}) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right] \\ & + \frac{p_{a} - p_{aa}}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b(h_{i}^{t} - \nabla f_{i}(x^{t})) \right\|^{2} \\ & + (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \\ & \leq \frac{2b^{2}\sigma^{2}}{np_{a}B} \\ & + \frac{2(1-b)^{2}}{n^{2}p_{a}} \sum_{i=1}^{n} \operatorname{E}_{k} \left[\left\| \nabla f_{i}(x^{t+1}; \xi_{i}^{t+1}) - \nabla f_{i}(x^{t}; \xi_{i}^{t+1}) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right] \\ & + \frac{2(p_{a} - p_{aa})}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} + \frac{2(p_{a} - p_{aa})b^{2}}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \\ & + (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}. \\ & = \frac{2b^{2}\sigma^{2}}{np_{a}B} \\ & + \frac{2(1-b)^{2}}{n^{2}p_{a}B^{2}} \sum_{i=1}^{n} \sum_{j=1}^{B} \operatorname{E}_{k} \left[\left\| \nabla f_{i}(x^{t+1}; \xi_{ij}^{t+1}) - \nabla f_{i}(x^{t}; \xi_{ij}^{t+1}) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right] \\ & + \frac{2(p_{a} - p_{aa})}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} + \frac{2(p_{a} - p_{aa})b^{2}}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \right) \\ & + (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}, \end{split}$$

where we use the independence of elements in the mini-batches. Using Assumptions 3 and 6, we obtain

$$\begin{split} & \mathbf{E}_{k} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \\ & \leq \frac{2b^{2}\sigma^{2}}{np_{\mathbf{a}}B} + \left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{np_{\mathbf{a}}B} + \frac{2(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})\widehat{L}^{2}}{np_{\mathbf{a}}^{2}} \right) \left\| x^{t+1} - x^{t} \right\|^{2} \\ & + \frac{2(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})b^{2}}{n^{2}p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} + (1-b)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}. \end{split}$$

Now, we prove the second inequality:

$$\begin{split} & E_k \left[\mathbb{E}_{p_k} \left[\left\| h_i^{t+1} - \nabla f_i(x^{t+1}) \right\|^2 \right] \right] \\ & = \mathbb{E}_k \left[\mathbb{E}_{p_k} \left[\left| h_i^{t+1} - \mathbb{E}_k \left[\mathbb{E}_{p_k} \left[h_i^{t+1} \right] \right] \right] \right] \\ & + \left\| \mathbb{E}_k \left[\mathbb{E}_{p_k} \left[\left| h_i^{t+1} - \mathbb{E}_k \left[\mathbb{E}_{p_k} \left[h_i^{t+1} \right] \right] \right] \right] \right] \\ & + \left\| \mathbb{E}_k \left[\mathbb{E}_{p_k} \left[\left| h_i^{t+1} - \left(h_i^t + \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t)) \right) \right] \right] \right] \\ & + \left\| h_i^t + \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t)) - \nabla f_i(x^{t+1}) \right\|^2 \\ & = \mathbb{E}_k \left[\mathbb{E}_{p_k} \left[\left\| h_i^{t+1} - \left(h_i^t + \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t)) \right) \right]^2 \right] \right] \\ & + (1 - b)^2 \left\| h_i^t - \nabla f_i(x^t) \right\|^2 \\ & = p_k \mathbb{E}_k \left[\left\| h_i^t + \frac{1}{p_k} h_i^{t+1} - \left(h_i^t + \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t)) \right) \right\|^2 \right] \\ & + (1 - p_k) \left\| h_i^t - \left(h_i^t + \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t)) \right) \right\|^2 \\ & = p_k \mathbb{E}_k \left[\left\| \frac{1}{p_k} h_i^{t+1} - \left(\nabla f_i(x^{t+1}) - \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t)) \right) \right\|^2 \right] \\ & + (1 - p_k) \left\| \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t)) \right\|^2 \\ & = p_k \mathbb{E}_k \left[\left\| h_i^{t+1} - \left(\nabla f_i(x^{t+1}) - \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t)) \right) \right\|^2 \right] \\ & + (1 - p_k) \left\| \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t)) \right\|^2 \\ & + (1 - b)^2 \left\| h_i^t - \nabla f_i(x^t) \right\|^2 - b(h_i^t - \nabla f_i(x^t)) \right\|^2 \\ & + (1 - b)^2 \left\| h_i^t - \nabla f_i(x^t) \right\|^2 \\ & + (1 - b)^2 \left\| h_i^t - \nabla f_i(x^t) \right\|^2 \\ & = \frac{1}{p_k} \mathbb{E}_k \left[\left\| \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t)) \right\|^2 \\ & + (1 - b)^2 \left\| h_i^t - \nabla f_i(x^t) \right\|^2 \\ & = \frac{1}{p_k} \mathbb{E}_k \left[\left\| \nabla f_i(x^{t+1}) + \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t)) \right\|^2 \\ & + (1 - b)^2 \left\| h_i^t - \nabla f_i(x^t) \right\|^2 \\ & + (1 - b)^2 \left\| h_i^t - \nabla f_i(x^t) \right\|^2 \\ & = \frac{1}{p_k} \mathbb{E}_k \left[\left\| \nabla f_i(x^{t+1}) + \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t)) \right\|^2 \\ & + \frac{1 - p_k}{p_k} \left\| \nabla f_i(x^{t+1}) + \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t)) \right\|^2 \\ & + \frac{1 - p_k}{p_k} \mathbb{E}_k \left[\left\| \nabla f_i(x^{t+1}) + \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t)) \right\|^2 \\ & + \frac{1 - p_k}{p_k} \mathbb{E}_k \left[\left\| \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - b(h_i^t - \nabla f_i(x^t) \right\|^2 \right] \\ & + \frac{1$$

659 Considering the independence of elements in the mini-batch, we obtain

$$\begin{split} & \mathbf{E}_{k} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h_{i}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \right] \\ & = \frac{2b^{2}}{p_{\mathbf{a}}B^{2}} \sum_{j=1}^{B} \mathbf{E}_{k} \left[\left\| \nabla f_{i}(x^{t+1}; \xi_{ij}^{t+1}) - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \\ & + \frac{2(1-b)^{2}}{p_{\mathbf{a}}B^{2}} \sum_{j=1}^{B} \mathbf{E}_{k} \left[\left\| \nabla f_{i}(x^{t+1}; \xi_{ij}^{t+1}) - \nabla f_{i}(x^{t}; \xi_{ij}^{t+1}) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right] \\ & + \frac{1-p_{\mathbf{a}}}{p_{\mathbf{a}}} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b(h_{i}^{t} - \nabla f_{i}(x^{t})) \right\|^{2} \\ & + (1-b)^{2} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \,. \end{split}$$

$$& \stackrel{(12)}{\leq} \frac{2b^{2}}{p_{\mathbf{a}}B^{2}} \sum_{j=1}^{B} \mathbf{E}_{k} \left[\left\| \nabla f_{i}(x^{t+1}; \xi_{ij}^{t+1}) - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \\ & + \frac{2(1-b)^{2}}{p_{\mathbf{a}}B^{2}} \sum_{j=1}^{B} \mathbf{E}_{k} \left[\left\| \nabla f_{i}(x^{t+1}; \xi_{ij}^{t+1}) - \nabla f_{i}(x^{t}; \xi_{ij}^{t+1}) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right] \\ & + \frac{2(1-p_{\mathbf{a}})}{p_{\mathbf{a}}} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} + \left(\frac{2(1-p_{\mathbf{a}})b^{2}}{p_{\mathbf{a}}} + (1-b)^{2} \right) \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \end{split}$$

Next, we use Assumptions 3, 6, 5, to get

$$\begin{split} & \mathbf{E}_{k} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h_{i}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \right] \\ & \leq \frac{2b^{2}\sigma^{2}}{p_{\mathbf{a}}B} + \left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{p_{\mathbf{a}}B} + \frac{2(1-p_{\mathbf{a}})L_{i}^{2}}{p_{\mathbf{a}}} \right) \left\| x^{t+1} - x^{t} \right\|^{2} \\ & + \left(\frac{2(1-p_{\mathbf{a}})b^{2}}{p_{\mathbf{a}}} + (1-b)^{2} \right) \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}. \end{split}$$

 $\text{ }^{\text{661}}\quad\text{It is left to prove the bound for }\mathrm{E}_{k}\left[\left\|k_{i}^{t+1}\right\|^{2}\right]\!:$

$$\begin{split} & \operatorname{E}_{k} \left[\left\| k_{i}^{t+1} \right\|^{2} \right] \\ & = \operatorname{E}_{k} \left[\left\| \nabla f_{i}(x^{t+1}; \xi_{i}^{t+1}) - \nabla f_{i}(x^{t}; \xi_{i}^{t+1}) - b \left(h_{i}^{t} - \nabla f_{i}(x^{t}; \xi_{i}^{t+1}) \right) \right\|^{2} \right] \\ & \stackrel{\text{(13)}}{=} \operatorname{E}_{k} \left[\left\| \nabla f_{i}(x^{t+1}; \xi_{i}^{t+1}) - \nabla f_{i}(x^{t}; \xi_{i}^{t+1}) - b \left(h_{i}^{t} - \nabla f_{i}(x^{t}; \xi_{i}^{t+1}) \right) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b (h_{i}^{t} - \nabla f_{i}(x^{t})) \right) \right\|^{2} \right] \\ & + \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b (h_{i}^{t} - \nabla f_{i}(x^{t})) \right\|^{2} \\ & = \operatorname{E}_{k} \left[\left\| b \left(\nabla f_{i}(x^{t+1}; \xi_{i}^{t+1}) - \nabla f_{i}(x^{t+1}) \right) + (1 - b) \left(\nabla f_{i}(x^{t+1}; \xi_{i}^{t+1}) - \nabla f_{i}(x^{t}; \xi_{i}^{t+1}) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right) \right\|^{2} \right] \\ & + \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - b (h_{i}^{t} - \nabla f_{i}(x^{t})) \right\|^{2} \\ & \leq 2b^{2} \operatorname{E}_{k} \left[\left\| \nabla f_{i}(x^{t+1}; \xi_{i}^{t+1}) - \nabla f_{i}(x^{t}; \xi_{i}^{t+1}) - \left(\nabla f_{i}(x^{t}; \xi_{i}^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right] \\ & + 2 \left\| \nabla f_{i}(x^{t+1}; \xi_{i}^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} + 2b^{2} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}. \end{split}$$

Using Assumptions 3, 6, 5 and the independence of elements in the mini-batch, we get

$$\mathbf{E}_k \left[\left\| k_i^{t+1} \right\|^2 \right]$$

$$\leq \frac{2b^{2}\sigma^{2}}{B} + \left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{B} + 2L_{i}^{2}\right) \left\|x^{t+1} - x^{t}\right\|^{2} + 2b^{2} \left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}.$$

663

Theorem 4. Suppose that Assumptions 1, 2, 3, 5, 6, 7 and 8 hold. Let us take
$$a = \frac{p_a}{2\omega + 1}$$
, $b \in \left(0, \frac{p_a}{2-p_a}\right]$, $\gamma \le \left(L + \left[\frac{48\omega(2\omega + 1)}{np_a^2}\left(\hat{L}^2 + \frac{(1-b)^2L_\sigma^2}{B}\right) + \frac{12}{np_ab}\left(\left(1 - \frac{p_{aa}}{p_a}\right)\hat{L}^2 + \frac{(1-b)^2L_\sigma^2}{B}\right)\right]^{1/2}\right)^{-1}$, and $g_i^0 = h_i^0$ for all $i \in [n]$ in Algorithm 1 (DASHA-PP-MVR). Then

$$\begin{split} & \mathbf{E}\left[\left\|\nabla f(\widehat{x}^T)\right\|^2\right] \leq \frac{1}{T}\left[\frac{2\Delta_0}{\gamma} + \frac{2}{b}\left\|h^0 - \nabla f(x^0)\right\|^2 + \left(\frac{32b\omega(2\omega+1)}{np_{\mathrm{a}}^2} + \frac{4\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)}{np_{\mathrm{a}}}\right)\left(\frac{1}{n}\sum_{i=1}^n\left\|h_i^0 - \nabla f_i(x^0)\right\|^2\right)\right] \\ & + \left(\frac{48b^2\omega(2\omega+1)}{p_{\mathrm{a}}^2} + \frac{12b}{p_{\mathrm{a}}}\right)\frac{\sigma^2}{nB}. \end{split}$$

Proof. Let us fix constants $\nu, \rho \in [0, \infty)$ that we will define later. Considering Lemma 6, Lemma 10, and the law of total expectation, we obtain

$$\begin{split} & \operatorname{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \operatorname{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1) p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \operatorname{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \nu \operatorname{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho \operatorname{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \operatorname{E}\left[f(x^{t}) - \frac{\gamma}{2} \left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right) \left\|x^{t+1} - x^{t}\right\|^{2} + \gamma \left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \operatorname{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1) p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \operatorname{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \frac{4\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} \operatorname{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|k_{i}^{t+1}\right\|^{2}\right] \\ & + \nu \operatorname{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho \operatorname{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & = \operatorname{E}\left[f(x^{t}) - \frac{\gamma}{2} \left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right) \left\|x^{t+1} - x^{t}\right\|^{2} + \gamma \left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \operatorname{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1) p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \operatorname{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \nu \operatorname{E}\left[\operatorname{E}_{B}\left[\operatorname{E}_{p_{\mathbf{a}}}\left\|h_{i}^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right]\right] \\ & + \nu \operatorname{E}\left[\operatorname{E}_{B}\left[\operatorname{E}_{p_{\mathbf{a}}}\left\|h_{i}^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right]\right] \\ & + \nu \operatorname{E}\left[\operatorname{E}_{B}\left[\operatorname{E}_{p_{\mathbf{a}}}\left\|h_{i}^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right]\right] \\ & + \left(\operatorname{E}_{B}\left[\operatorname{E}_{p_{\mathbf{a}}}\left\|h_{i}^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right]\right] \\ & + \left(\operatorname{E}_{B}\left[\operatorname{E}_{p_{\mathbf{a}}}\left\|h_{i}^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}}\operatorname{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1) p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\operatorname{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \frac{4\gamma\omega(2\omega+1)}{np_{\mathbf{a}}}\operatorname{E}\left[\frac{2b^{2}\sigma^{2}}{B} + \left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{B} + 2\widehat{L}^{2}\right)\right] \left\|x^{t+1} - x^{t}\right\|^{2} + 2b^{2} \frac{1}{n} \sum_{i=1}^{n} \left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \end{aligned}$$

$$\begin{split} &+\nu \mathbf{E} \Bigg(\frac{2b^{2}\sigma^{2}}{np_{\mathbf{a}}B} + \Bigg(\frac{2(1-b)^{2}L_{\sigma}^{2}}{np_{\mathbf{a}}B} + \frac{2\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)\widehat{L}^{2}}{np_{\mathbf{a}}^{2}} \Bigg) \left\| x^{t+1} - x^{t} \right\|^{2} \\ &+ \frac{2\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)b^{2}}{n^{2}p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} + \left(1-b\right)^{2} \left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \Bigg) \\ &+ \rho \mathbf{E} \Bigg(\frac{2b^{2}\sigma^{2}}{p_{\mathbf{a}}B} + \Bigg(\frac{2(1-b)^{2}L_{\sigma}^{2}}{p_{\mathbf{a}}B} + \frac{2(1-p_{\mathbf{a}})\widehat{L}^{2}}{p_{\mathbf{a}}} \Bigg) \left\| x^{t+1} - x^{t} \right\|^{2} \\ &+ \Bigg(\frac{2(1-p_{\mathbf{a}})b^{2}}{p_{\mathbf{a}}} + (1-b)^{2} \Bigg) \frac{1}{n} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \Bigg). \end{split}$$

669 After rearranging the terms, we get

$$\begin{split} & \mathbb{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbb{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \nu\mathbb{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbb{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbb{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbb{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}}\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{4\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}}\left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{B} + 2\hat{L}^{2}\right) - \rho\left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{p_{\mathbf{a}}B} + \frac{2(1-p_{\mathbf{a}})\hat{L}^{2}}{p_{\mathbf{a}}}\right)\right)\mathbb{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \left(\gamma + \nu\left(1 - b\right)^{2}\right)\mathbb{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\nu\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)b^{2}}{np_{\mathbf{a}}^{2}} + \rho\left(\frac{2(1-p_{\mathbf{a}})b^{2}}{p_{\mathbf{a}}} + (1-b)^{2}\right)\right)\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \nu\frac{2b^{2}}{np_{\mathbf{a}}} + \rho\frac{2b^{2}}{p_{\mathbf{a}}}\right)\frac{\sigma^{2}}{B}. \end{split}$$

670 By taking $\nu=\frac{\gamma}{b},$ one can show that $\left(\gamma+\nu(1-b)^2\right)\leq \nu,$ and

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{4\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}}\left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{B} + 2\widehat{L}^{2}\right) - \rho\left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{p_{\mathbf{a}}B} + \frac{2(1-p_{\mathbf{a}})\widehat{L}^{2}}{p_{\mathbf{a}}}\right)\right) \mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \end{split}$$

$$\begin{split} & + \left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{a}^{2}} + \frac{2\gamma\left(p_{a} - p_{aa}\right)b}{np_{a}^{2}} + \rho\left(\frac{2(1-p_{a})b^{2}}{p_{a}} + (1-b)^{2}\right)\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{a}^{2}} + \frac{2\gamma b}{np_{a}} + \rho\frac{2b^{2}}{p_{a}}\right)\frac{\sigma^{2}}{B}. \end{split}$$

Note that $b \leq \frac{p_a}{2-p_a}$, thus

$$\left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{a}^{2}} + \frac{2\gamma(p_{a} - p_{aa})b}{np_{a}^{2}} + \rho\left(\frac{2(1-p_{a})b^{2}}{p_{a}} + (1-b)^{2}\right)\right) \\
\leq \left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{a}^{2}} + \frac{2\gamma(p_{a} - p_{aa})b}{np_{a}^{2}} + \rho(1-b)\right).$$

And if we take $ho=rac{8b\gamma\omega(2\omega+1)}{np_{
m a}^2}+rac{2\gamma(p_{
m a}-p_{
m aa})}{np_{
m a}^2},$ then

$$\left(\frac{8b^2\gamma\omega(2\omega+1)}{np_{\rm a}^2} + \frac{2\gamma\left(p_{\rm a}-p_{\rm aa}\right)b}{np_{\rm a}^2} + \rho\left(1-b\right)\right) \leq \rho,$$

673 and

$$\begin{split} & \mathbb{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbb{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g^{t+1}_{i} - h^{t+1}_{i}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbb{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h^{t+1}_{i} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbb{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbb{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbb{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g^{t}_{i} - h^{t}_{i}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{4\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}}\left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{B} + 2\hat{L}^{2}\right) - \frac{\gamma}{np_{\mathbf{a}}b}\left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{B} + 2\left(1 - \frac{p_{\mathbf{aa}}}{p_{\mathbf{a}}}\right)\hat{L}^{2}\right) \\ & - \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{3}} + \frac{2\gamma\left(1 - \frac{p_{\mathbf{aa}}}{p_{\mathbf{a}}}\right)\hat{L}^{2}\right) - \frac{2(1-b)^{2}L_{\sigma}^{2}}{B} + 2(1-p_{\mathbf{a}})\hat{L}^{2}\right) \mathbb{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \frac{\gamma}{b}\mathbb{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\right)\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h^{t}_{i} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma b}{np_{\mathbf{a}}} + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\right)\frac{2b^{2}}{p_{\mathbf{a}}}\right)\frac{\sigma^{2}}{B}. \end{split}$$

Let us simplify the inequality. First, due to $b \le p_{\rm a}$ and $(1-p_{\rm a}) \le \left(1-\frac{p_{\rm aa}}{p_{\rm a}}\right)$, we have

$$\begin{split} &\left(\frac{8b\gamma\omega(2\omega+1)}{np_{\rm a}^3} + \frac{2\gamma\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)}{np_{\rm a}^2}\right) \left(\frac{2(1-b)^2L_\sigma^2}{B} + 2(1-p_{\rm a})\widehat{L}^2\right) \\ &= \frac{8b\gamma\omega(2\omega+1)}{np_{\rm a}^3} \left(\frac{2(1-b)^2L_\sigma^2}{B} + 2(1-p_{\rm a})\widehat{L}^2\right) \\ &\quad + \frac{2\gamma\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)}{np_{\rm a}^2} \left(\frac{2(1-b)^2L_\sigma^2}{B} + 2(1-p_{\rm a})\widehat{L}^2\right) \end{split}$$

$$\leq \frac{8\gamma\omega(2\omega+1)}{np_a^2} \left(\frac{2(1-b)^2L_\sigma^2}{B} + 2\widehat{L}^2\right) + \frac{2\gamma}{np_ab} \left(\frac{2(1-b)^2L_\sigma^2}{B} + 2\left(1 - \frac{p_{aa}}{p_a}\right)\widehat{L}^2\right),$$

675 therefore

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g^{t+1}_{i} - h^{t+1}_{i}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h^{t+1}_{i} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g^{t}_{i} - h^{t}_{i}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{12\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} \left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{B} + 2\hat{L}^{2}\right) - \frac{3\gamma}{np_{\mathbf{a}}b} \left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{B} + 2\left(1 - \frac{p_{\mathbf{a}\mathbf{a}}}{p_{\mathbf{a}}}\right)\hat{L}^{2}\right)\right) \mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h^{t}_{i} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma b}{np_{\mathbf{a}}} + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right) \frac{2b^{2}}{B}\right) \frac{\sigma^{2}}{B} \\ & = \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{np_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g^{t}_{i} - h^{t}_{i}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{24\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} \left(\frac{(1-b)^{2}L_{\sigma}^{2}}{B} + \hat{L}^{2}\right) - \frac{6\gamma}{np_{\mathbf{a}}b} \left(\frac{(1-b)^{2}L_{\sigma}^{2}}{B} + \left(1 - \frac{p_{\mathbf{a}\mathbf{a}}}{np_{\mathbf{a}}}\right)\hat{L}^{2}\right) \mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \frac{\beta}{b} \mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h^{t}_{i} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h^{t}_{i} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{8b^{2}\gamma\omega(2\omega+1)}{np_{\mathbf{a}}} + \frac{2\gamma(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a})}{np_{\mathbf{a}}^{2}}$$

676 Also, we can simplify the last term

$$\begin{split} &\left(\frac{8b\gamma\omega(2\omega+1)}{np_{\rm a}^2} + \frac{2\gamma\left(p_{\rm a} - p_{\rm aa}\right)}{np_{\rm a}^2}\right)\frac{2b^2}{p_{\rm a}} \\ &= \frac{16b^3\gamma\omega(2\omega+1)}{np_{\rm a}^3} + \frac{4b^2\gamma\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)}{np_{\rm a}^2} \\ &\leq \frac{16b^2\gamma\omega(2\omega+1)}{np_{\rm a}^2} + \frac{4b\gamma}{np_{\rm a}}, \end{split}$$

677 thus

$$\mathrm{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathrm{a}}}\mathrm{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathrm{a}} - p_{\mathrm{aa}})}{np_{\mathrm{a}}^{2}}\mathrm{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right]$$

$$\begin{split} & + \frac{\gamma}{b} \mathbf{E} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{a}^{2}} + \frac{2\gamma\left(p_{a} - p_{aa}\right)}{np_{a}^{2}} \right) \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| h_{i}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \\ & \leq \mathbf{E} \left[f(x^{t}) \right] - \frac{\gamma}{2} \mathbf{E} \left[\left\| \nabla f(x^{t}) \right\|^{2} \right] \\ & + \frac{\gamma(2\omega+1)}{p_{a}} \mathbf{E} \left[\left\| g^{t} - h^{t} \right\|^{2} \right] + \frac{\gamma((2\omega+1)p_{a} - p_{aa})}{np_{a}^{2}} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{24\gamma\omega(2\omega+1)}{np_{a}^{2}} \left(\frac{(1-b)^{2}L_{\sigma}^{2}}{B} + \hat{L}^{2} \right) \right) \\ & - \frac{6\gamma}{np_{a}b} \left(\frac{(1-b)^{2}L_{\sigma}^{2}}{B} + \left(1 - \frac{p_{aa}}{p_{a}} \right) \hat{L}^{2} \right) \right) \mathbf{E} \left[\left\| x^{t+1} - x^{t} \right\|^{2} \right] \\ & + \frac{\gamma}{b} \mathbf{E} \left[\left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{a}^{2}} + \frac{2\gamma\left(p_{a} - p_{aa}\right)}{np_{a}^{2}} \right) \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \right] \\ & + \left(\frac{24b^{2}\gamma\omega(2\omega+1)}{np_{a}^{2}} + \frac{6\gamma b}{np_{a}} \right) \frac{\sigma^{2}}{B}. \end{split}$$

Using Lemma 4 and the assumption about γ , we get

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(\frac{8b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{24b^{2}\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{6\gamma b}{np_{\mathbf{a}}}\right)\frac{\sigma^{2}}{B}. \end{split}$$

It is left to apply Lemma 3 with

$$\Psi^{t} = \frac{(2\omega + 1)}{p_{a}} E\left[\|g^{t} - h^{t}\|^{2} \right] + \frac{((2\omega + 1)p_{a} - p_{aa})}{np_{a}^{2}} E\left[\frac{1}{n} \sum_{i=1}^{n} \|g_{i}^{t} - h_{i}^{t}\|^{2} \right]$$

$$+ \frac{1}{b} E\left[\|h^{t} - \nabla f(x^{t})\|^{2} \right] + \left(\frac{8b\omega(2\omega + 1)}{np_{a}^{2}} + \frac{2(p_{a} - p_{aa})}{np_{a}^{2}} \right) E\left[\frac{1}{n} \sum_{i=1}^{n} \|h_{i}^{t} - \nabla f_{i}(x^{t})\|^{2} \right]$$

and $C=\left(\frac{24b^2\omega(2\omega+1)}{p_{\rm a}^2}+\frac{6b}{p_{\rm a}}\right)\frac{\sigma^2}{nB}$ to conclude the proof.

Corollary 3. Suppose that assumptions from Theorem 4 hold, momentum
$$b = \Theta\left(\min\left\{\frac{p_a}{\omega}\sqrt{\frac{n\varepsilon B}{\sigma^2}},\frac{p_a n\varepsilon B}{\sigma^2}\right\}\right), \frac{\sigma^2}{n\varepsilon B} \geq 1, \text{ and } h_i^0 = \frac{1}{B_{\text{init}}}\sum_{k=1}^{B_{\text{init}}}\nabla f_i(x^0;\xi_{ik}^0) \text{ for all } i \in [n],$$

and batch size $B_{ ext{init}} = \Theta\left(\frac{\sqrt{p_a}B}{b}\right)$, then Algorithm 1 (DASHA-PP-MVR) needs

$$T := \mathcal{O}\left(\frac{\Delta_0}{\varepsilon} \left[L + \frac{\omega}{p_{\rm a}\sqrt{n}} \left(\widehat{L} + \frac{L_\sigma}{\sqrt{B}}\right) + \frac{\sigma}{p_{\rm a}\sqrt{\varepsilon}n} \left(\frac{\mathbbm{1}_{p_{\rm a}}\widehat{L}}{\sqrt{B}} + \frac{L_\sigma}{B}\right)\right] + \frac{\sigma^2}{\sqrt{p_{\rm a}}n\varepsilon B}\right)$$

communication rounds to get an ε -solution and the number of stochastic gradient calculations per node equals $\mathcal{O}(B_{\text{init}} + BT)$

686 *Proof.* Using the result from Theorem 4, we have

$$\begin{split} & \mathbf{E}\left[\left\|\nabla f(\widehat{x}^{T})\right\|^{2}\right] \\ & \leq \frac{1}{T}\left[2\Delta_{0}\left(L+\sqrt{\frac{48\omega(2\omega+1)}{np_{a}^{2}}\left(\widehat{L}^{2}+\frac{(1-b)^{2}L_{\sigma}^{2}}{B}\right)+\frac{12}{np_{a}b}\left(\left(1-\frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)\widehat{L}^{2}+\frac{(1-b)^{2}L_{\sigma}^{2}}{B}\right)\right) \\ & + \frac{2}{b}\left\|h^{0}-\nabla f(x^{0})\right\|^{2}+\left(\frac{32b\omega(2\omega+1)}{np_{a}^{2}}+\frac{4\left(1-\frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)}{np_{\mathrm{a}}}\right)\left(\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{0}-\nabla f_{i}(x^{0})\right\|^{2}\right)\right] \\ & + \left(\frac{48b^{2}\omega(2\omega+1)}{p_{\mathrm{a}}^{2}}+\frac{12b}{p_{\mathrm{a}}}\right)\frac{\sigma^{2}}{nB} \end{split}$$

We choose b to ensure $\left(\frac{48b^2\omega(2\omega+1)}{p_a^2}+\frac{12b}{p_a}\right)\frac{\sigma^2}{nB}=\Theta\left(\varepsilon\right)$. Note that $\frac{1}{b}=0$ 0 for $\left\{\max\left\{\frac{\omega}{p_a}\sqrt{\frac{\sigma^2}{n\varepsilon B}},\frac{\sigma^2}{p_an\varepsilon B}\right\}\right\} \leq \Theta\left(\max\left\{\frac{\omega^2}{p_a},\frac{\sigma^2}{p_an\varepsilon B}\right\}\right)$, thus

$$\begin{split} & \mathbf{E}\left[\left\|\nabla f(\widehat{x}^T)\right\|^2\right] \\ & = \mathcal{O}\left(\frac{1}{T}\left[\Delta_0\left(L + \frac{\omega}{p_{\mathrm{a}}\sqrt{n}}\left(\widehat{L} + \frac{L_\sigma}{\sqrt{B}}\right) + \sqrt{\frac{\sigma^2}{p_{\mathrm{a}}^2\varepsilon n^2B}}\left(\mathbbm{1}_{p_{\mathrm{a}}}\widehat{L} + \frac{L_\sigma}{\sqrt{B}}\right)\right) \right. \\ & \left. + \frac{1}{b}\left\|h^0 - \nabla f(x^0)\right\|^2 + \left(\frac{b\omega^2}{np_{\mathrm{a}}^2} + \frac{1}{np_{\mathrm{a}}}\right)\left(\frac{1}{n}\sum_{i=1}^n\left\|h_i^0 - \nabla f_i(x^0)\right\|^2\right)\right] + \varepsilon\right), \end{split}$$

where $\mathbb{1}_{p_{\rm a}}=\sqrt{1-rac{p_{\rm aa}}{p_{\rm a}}}.$ It enough to take the following T to get arepsilon-solution.

$$T = \mathcal{O}\left(\frac{1}{\varepsilon} \left[\Delta_0 \left(L + \frac{\omega}{p_a \sqrt{n}} \left(\widehat{L} + \frac{L_\sigma}{\sqrt{B}} \right) + \sqrt{\frac{\sigma^2}{p_a^2 \varepsilon n^2 B}} \left(\mathbb{1}_{p_a} \widehat{L} + \frac{L_\sigma}{\sqrt{B}} \right) \right) + \frac{1}{b} \left\| h^0 - \nabla f(x^0) \right\|^2 + \left(\frac{b\omega^2}{np_a^2} + \frac{1}{np_a} \right) \left(\frac{1}{n} \sum_{i=1}^n \left\| h_i^0 - \nabla f_i(x^0) \right\|^2 \right) \right] \right).$$

690 Let us bound the norms:

Using the same reasoning, one cat get $\frac{1}{n}\sum_{i=1}^n \mathrm{E}\left[\left\|h_i^0 - \nabla f_i(x^0)\right\|^2\right] \leq \frac{\sigma^2}{B_{\mathrm{init}}}$. Combining all inequal-

692 ities, we have

$$\begin{split} T &= \mathcal{O} \Bigg(\frac{1}{\varepsilon} \Bigg[\Delta_0 \left(L + \frac{\omega}{p_{\rm a} \sqrt{n}} \left(\widehat{L} + \frac{L_\sigma}{\sqrt{B}} \right) + \sqrt{\frac{\sigma^2}{p_{\rm a}^2 \varepsilon n^2 B}} \left(\mathbbm{1}_{p_{\rm a}} \widehat{L} + \frac{L_\sigma}{\sqrt{B}} \right) \right) \\ &+ \frac{\sigma^2}{b n B_{\rm init}} + \frac{b \omega^2 \sigma^2}{n p_{\rm a}^2 B_{\rm init}} + \frac{\sigma^2}{n p_{\rm a} B_{\rm init}} \Bigg] \Bigg). \end{split}$$

Using the choice of B_{init} and b, we obtain

$$\begin{split} T &= \mathcal{O} \left(\frac{1}{\varepsilon} \left[\Delta_0 \left(L + \frac{\omega}{p_{\rm a} \sqrt{n}} \left(\hat{L} + \frac{L_\sigma}{\sqrt{B}} \right) + \sqrt{\frac{\sigma^2}{p_{\rm a}^2 \varepsilon n^2 B}} \left(\mathbbm{1}_{p_{\rm a}} \hat{L} + \frac{L_\sigma}{\sqrt{B}} \right) \right) \right. \\ &+ \frac{\sigma^2}{\sqrt{p_{\rm a}} n B} + \frac{b^2 \omega^2 \sigma^2}{n p_{\rm a}^{5/2} B} + \frac{b \sigma^2}{p_{\rm a}^{3/2} n B} \right] \right) \\ &= \mathcal{O} \left(\frac{1}{\varepsilon} \left[\Delta_0 \left(L + \frac{\omega}{p_{\rm a} \sqrt{n}} \left(\hat{L} + \frac{L_\sigma}{\sqrt{B}} \right) + \sqrt{\frac{\sigma^2}{p_{\rm a}^2 \varepsilon n^2 B}} \left(\mathbbm{1}_{p_{\rm a}} \hat{L} + \frac{L_\sigma}{\sqrt{B}} \right) \right) \right. \\ &+ \frac{\sigma^2}{\sqrt{p_{\rm a}} n B} + \frac{\varepsilon}{\sqrt{p_{\rm a}}} \right] \right) \\ &= \mathcal{O} \left(\frac{\Delta_0}{\varepsilon} \left[L + \frac{\omega}{p_{\rm a} \sqrt{n}} \left(\hat{L} + \frac{L_\sigma}{\sqrt{B}} \right) + \sqrt{\frac{\sigma^2}{p_{\rm a}^2 \varepsilon n^2 B}} \left(\mathbbm{1}_{p_{\rm a}} \hat{L} + \frac{L_\sigma}{\sqrt{B}} \right) \right] + \frac{\sigma^2}{\sqrt{p_{\rm a}} n \varepsilon B} + \frac{1}{\sqrt{p_{\rm a}}} \right). \end{split}$$

Using $\frac{\sigma^2}{n \in B} \ge 1$, we can conclude the proof of the inequality. The number of stochastic gradients that each node calculates equals $B_{\text{init}} + 2BT = \mathcal{O}(B_{\text{init}} + BT)$.

Corollary 4. Suppose that assumptions of Corollary 3 hold, batch size $B \leq \min \left\{ \frac{\sigma}{p_n \sqrt{\varepsilon} n}, \frac{L_{\sigma}^2}{1_{p_n}^2 \widehat{L}^2} \right\}$,

we take RandK compressors with $K=\Theta\left(\frac{Bd\sqrt{arepsilon n}}{\sigma}\right)$. Then the communication complexity equals

 $\mathcal{O}\left(\frac{d\sigma}{\sqrt{p_a}\sqrt{n\varepsilon}} + \frac{L_{\sigma}\Delta_0 d}{p_a\sqrt{n\varepsilon}}\right)$, and the expected number of stochastic gradient calculations per node equals $\mathcal{O}\left(\frac{\sigma^2}{\sqrt{p_a}n\varepsilon} + \frac{L_{\sigma}\Delta_0 \sigma}{p_a\varepsilon^{3/2}n}\right)$.

Proof. The communication complexity equals

$$\mathcal{O}\left(d+KT\right) \ = \ \mathcal{O}\left(d+\frac{\Delta_0}{\varepsilon}\left[KL+K\frac{\omega}{p_{\rm a}\sqrt{n}}\left(\widehat{L}+\frac{L_\sigma}{\sqrt{B}}\right)+K\sqrt{\frac{\sigma^2}{p_{\rm a}^2\varepsilon n^2B}}\left(\mathbbm{1}_{p_{\rm a}}\widehat{L}+\frac{L_\sigma}{\sqrt{B}}\right)\right]+K\frac{\sigma^2}{\sqrt{p_{\rm a}}n\varepsilon B}\right).$$

701 Due to $B \leq \frac{L_\sigma^2}{1_{p_a}^2 \widehat{L}^2}$, we have $\mathbb{1}_{p_a} \widehat{L} + \frac{L_\sigma}{\sqrt{B}} \leq \frac{2L_\sigma}{\sqrt{B}}$ and

$$\mathcal{O}\left(d+KT\right) = \mathcal{O}\left(d+\frac{\Delta_0}{\varepsilon}\left[KL+K\frac{\omega}{p_{\rm a}\sqrt{n}}\left(\widehat{L}+\frac{L_\sigma}{\sqrt{B}}\right)+K\sqrt{\frac{\sigma^2}{p_{\rm a}^2\varepsilon n^2B}}\frac{L_\sigma}{\sqrt{B}}\right]+K\frac{\sigma^2}{\sqrt{p_{\rm a}}n\varepsilon B}\right).$$

From Theorem 6, we have $\omega+1=\frac{d}{K}$. Since $K=\Theta\left(\frac{Bd\sqrt{\varepsilon n}}{\sigma}\right)=\mathcal{O}\left(\frac{d}{p_a\sqrt{n}}\right)$, the communication complexity equals

$$\mathcal{O}(d+KT) = \mathcal{O}\left(d + \frac{\Delta_0}{\varepsilon} \left[\frac{d}{p_{\rm a}\sqrt{n}} L + \frac{d}{p_{\rm a}\sqrt{n}} \left(\widehat{L} + \frac{L_\sigma}{\sqrt{B}}\right) + \frac{d}{p_{\rm a}\sqrt{n}} L_\sigma \right] + \frac{d\sigma}{\sqrt{p_{\rm a}}\sqrt{n\varepsilon}} \right)$$

$$= \mathcal{O}\left(\frac{d\sigma}{\sqrt{p_{\rm a}}\sqrt{n\varepsilon}} + \frac{L_\sigma\Delta_0 d}{p_{\rm a}\sqrt{n\varepsilon}}\right)$$

And the expected number of stochastic gradient calculations per node equals

$$\mathcal{O}\left(B_{\text{init}} + BT\right)$$

$$\begin{split} &=\mathcal{O}\left(\frac{\sigma^2}{\sqrt{p_{\rm a}}n\varepsilon}+\frac{B\omega}{\sqrt{p_{\rm a}}}\sqrt{\frac{\sigma^2}{n\varepsilon B}}+\frac{\Delta_0}{\varepsilon}\left[BL+B\frac{\omega}{p_{\rm a}\sqrt{n}}\left(\widehat{L}+\frac{L_\sigma}{\sqrt{B}}\right)+B\sqrt{\frac{\sigma^2}{p_{\rm a}^2\varepsilon n^2 B}}\left(\mathbbm{1}_{p_{\rm a}}\widehat{L}+\frac{L_\sigma}{\sqrt{B}}\right)\right]+B\frac{\sigma^2}{\sqrt{p_{\rm a}}n\varepsilon B}\right)\\ &=\mathcal{O}\left(\frac{\sigma^2}{\sqrt{p_{\rm a}}n\varepsilon}+\frac{Bd}{K\sqrt{p_{\rm a}}}\sqrt{\frac{\sigma^2}{n\varepsilon B}}+\frac{\Delta_0}{\varepsilon}\left[BL+B\frac{d}{Kp_{\rm a}\sqrt{n}}\left(\widehat{L}+\frac{L_\sigma}{\sqrt{B}}\right)+B\sqrt{\frac{\sigma^2}{p_{\rm a}^2\varepsilon n^2 B}}\frac{L_\sigma}{\sqrt{B}}\right]+\frac{\sigma^2}{\sqrt{p_{\rm a}}n\varepsilon}\right)\\ &=\mathcal{O}\left(\frac{\sigma^2}{\sqrt{p_{\rm a}}n\varepsilon}+\frac{\sigma^2}{\sqrt{p_{\rm a}}n\varepsilon\sqrt{B}}+\frac{\Delta_0}{\varepsilon}\left[\frac{\sigma}{p_{\rm a}\sqrt{\varepsilon n}}L+\frac{\sigma}{p_{\rm a}\sqrt{\varepsilon n}}\left(\widehat{L}+\frac{L_\sigma}{\sqrt{B}}\right)+\frac{\sigma}{p_{\rm a}\sqrt{\varepsilon n}}L_\sigma\right]\right)\\ &=\mathcal{O}\left(\frac{\sigma^2}{\sqrt{p_{\rm a}}n\varepsilon}+\frac{L_\sigma\Delta_0\sigma}{p_{\rm a}\varepsilon^{3/2}n}\right). \end{split}$$

705

Analysis of DASHA-PP under Polyak-Łojasiewicz Condition

In this section, we provide the theoretical convergence rates of DASHA-PP under Polyak-Łojasiewiczc 707 Condition. 708

Assumption 9. The function f satisfy (Polyak-Łojasiewicz) PŁ-condition: 709

$$\|\nabla f(x)\|^2 \ge 2\mu(f(x) - f^*), \quad \forall x \in \mathbb{R},\tag{27}$$

where $f^* = \inf_{x \in \mathbb{R}^d} f(x) > -\infty$. 710

Under Polyak-Łojasiewicz condition, a (random) point \hat{x} is ε -solution, if $E[f(\hat{x})] - f^* \leq \varepsilon$.

We now provide the convergence rates of DASHA-PP under PŁ-condition.

F.1 Gradient Setting 713

Theorem 8. Suppose that Assumption 1, 2, 3, 7, 8 and 9 hold. Let us take $a = \frac{p_a}{2\omega+1}$, $b = \frac{p_a}{2-p_a}$.

$$\gamma \leq \min \left\{ \left(L + \sqrt{\frac{200\omega\left(2\omega + 1\right)}{np_{\mathrm{a}}^2} + \frac{48}{np_{\mathrm{a}}^2}\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)}\widehat{L}\right)^{-1}, \frac{a}{4\mu} \right\},$$

714 and $h_i^0=g_i^0=\nabla f_i(x^0)$ for all $i\in[n]$ in Algorithm I (DASHA-PP), then $\mathrm{E}\left[f(x^T)\right]-f^*\leq 1$

Let us provide bounds up to logarithmic factors and use $\mathcal{O}(\cdot)$ notation. The provided theorem states that to get ε -solution DASHA-PP have to run

$$\widetilde{\mathcal{O}}\left(rac{\omega+1}{p_{\mathrm{a}}}+rac{L}{\mu}+rac{\omega\widehat{L}}{p_{\mathrm{a}}\mu\sqrt{n}}+rac{\widehat{L}}{p_{\mathrm{a}}\mu\sqrt{n}}
ight),$$

communication rounds. The method DASHA from (Tyurin and Richtárik, 2023), have to run

$$\widetilde{\mathcal{O}}\left(\omega + \frac{L}{\mu} + \frac{\omega \widehat{L}}{\mu \sqrt{n}}\right),\,$$

communication rounds to get ε -solution. The difference is the same as in the general nonconvex case

(see Section 6.1). Up to Lipschitz constants factors, we get the degeneration up to $1/p_a$ factor due to

the partial participation.

F.2 Finite-Sum Setting 722

Theorem 9. Suppose that Assumption 1, 2, 3, 7, 4, 8, and 9 hold. Let us take $a = \frac{p_a}{2\omega+1}$, probability $p_{page} = \frac{B}{m+B}, b = \frac{p_{page}p_a}{2-p_a},$

$$\gamma \leq \min \left\{ \left(L + \sqrt{\frac{200\omega(2\omega+1)}{np_{\mathrm{a}}^2} \left(\widehat{L}^2 + \frac{(1-p_{\mathrm{page}})L_{\mathrm{max}}^2}{B}\right) + \frac{48}{np_{\mathrm{a}}^2p_{\mathrm{page}}} \left(\left(1-\frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)\widehat{L}^2 + \frac{(1-p_{\mathrm{page}})L_{\mathrm{max}}^2}{B}\right)\right)^{-1}, \frac{a}{2\mu}, \frac{b}{2\mu} \right\},$$

and $h_i^0 = g_i^0 = \nabla f_i(x^0)$ for all $i \in [n]$ in Algorithm 1 (DASHA-PP-PAGE), then $\mathrm{E}\left[f(x^T)\right] - f^* \leq 1$

 $(1-\gamma\mu)^T\Delta_0$.

The provided theorem states that to get ε -solution DASHA-PP have to run

$$\widetilde{\mathcal{O}}\left(\frac{\omega+1}{p_{\rm a}} + \frac{m}{p_{\rm a}B} + \frac{L}{\mu} + \frac{\omega}{p_{\rm a}\mu\sqrt{n}}\left(\widehat{L} + \frac{L_{\rm max}}{\sqrt{B}}\right) + \frac{\sqrt{m}}{p_{\rm a}\mu\sqrt{nB}}\left(\widehat{L} + \frac{L_{\rm max}}{\sqrt{B}}\right)\right),$$

communication rounds. The method DASHA-PAGE from (Tyurin and Richtárik, 2023), have to run

$$\widetilde{\mathcal{O}}\left(\omega + \frac{m}{B} + \frac{L}{\mu} + \frac{\omega}{\mu\sqrt{n}}\left(\widehat{L} + \frac{L_{\max}}{\sqrt{B}}\right) + \frac{\sqrt{m}}{\mu\sqrt{nB}}\left(\frac{L_{\max}}{\sqrt{B}}\right)\right),\,$$

communication rounds to get ε -solution. We can guarantee the degeneration up to $^1/p_a$ factor due to the partial participation only if $B=\mathcal{O}\left(\frac{L_{\max}^2}{\hat{L}^2}\right)$. The same conclusion we have in Section 6.2.

729 F.3 Stochastic Setting

Theorem 10. Suppose that Assumption 1, 2, 3, 7, 5, 6, 8 and 9 hold. Let us take $a = \frac{p_a}{2\omega+1}$, $b \in \left(0, \frac{p_a}{2-p_a}\right]$,

$$\gamma \leq \min \left\{ \left(L + \sqrt{\frac{200\omega(2\omega+1)}{np_{\mathrm{a}}^2} \left(\frac{(1-b)^2 L_\sigma^2}{B} + \widehat{L}^2 \right) + \frac{40}{np_{\mathrm{a}}b} \left(\frac{(1-b)^2 L_\sigma^2}{B} + \left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}} \right) \widehat{L}^2 \right)} \right)^{-1}, \frac{a}{2\mu}, \frac{b}{2\mu} \right\},$$

and $h_i^0 = g_i^0$ for all $i \in [n]$ in Algorithm I (DASHA-PP-MVR), then $F_i \lceil f(x^T) - f^* \rceil$

$$\leq (1 - \gamma \mu)^T \left(\Delta_0 + \frac{2\gamma}{b} \left\| h^0 - \nabla f(x^0) \right\|^2 + \left(\frac{40\gamma b\omega(2\omega + 1)}{np_{\rm a}^2} + \frac{8\gamma \left(p_{\rm a} - p_{\rm aa} \right)}{np_{\rm a}^2} \right) \frac{1}{n} \sum_{i=1}^n \left\| h_i^0 - \nabla f_i(x^0) \right\|^2 \right)$$

$$+ \frac{1}{\mu} \left(\frac{100b^2\omega(2\omega + 1)}{p_{\rm a}^2} + \frac{20b}{p_{\rm a}} \right) \frac{\sigma^2}{nB}.$$

The provided theorems states that to get ε -solution DASHA-PP have to run

$$\widetilde{\mathcal{O}}\left(\frac{\omega+1}{p_{a}} + \underbrace{\frac{\omega}{p_{a}}\sqrt{\frac{\sigma^{2}}{\mu n \varepsilon B}}}_{\widetilde{\mathcal{P}}_{2}} + \underbrace{\frac{\sigma^{2}}{p_{a}\mu n \varepsilon B}}_{+} + \underbrace{\frac{L}{\mu} + \frac{\omega}{p_{a}\mu \sqrt{n}}\left(\widehat{L} + \frac{L_{\sigma}}{\sqrt{B}}\right)}_{p_{a}\mu \sqrt{n}} + \underbrace{\frac{\sigma}{p_{a}\mu \sqrt{n}}\left(\widehat{L} + \frac{L_{\sigma}}{\sqrt{B}}\right)}_{\widetilde{\mathcal{P}}_{1}}\right)$$
(28)

732 communication rounds. We take $b = \Theta\left(\min\left\{\frac{p_a}{\omega}\sqrt{\frac{\mu n \varepsilon B}{\sigma^2}}, \frac{p_a \mu n \varepsilon B}{\sigma^2}\right\}\right) \ge$

733 $\Theta\left(\min\left\{\frac{p_a}{\omega^2}, \frac{p_a\mu n\varepsilon B}{\sigma^2}\right\}\right)$.

The method DASHA-SYNC-MVR from (Tyurin and Richtárik, 2023), have to run

$$\widetilde{\mathcal{O}}\left(\omega + \frac{\sigma^2}{\mu n \varepsilon B} + \frac{L}{\mu} + \frac{\omega}{\mu \sqrt{n}} \left(\widehat{L} + \frac{L_{\sigma}}{\sqrt{B}}\right) + \frac{\sigma}{n \mu^{3/2} \sqrt{\varepsilon B}} \left(\frac{L_{\sigma}}{\sqrt{B}}\right)\right) \tag{29}$$

735 communication rounds to get ε -solution⁸.

In the stochastic setting, the comparison is a little bit more complicated. As in the finite-sum setting,

we have to take $B=\mathcal{O}\left(\frac{L_{\sigma}^2}{\widehat{L}^2}\right)$ to guarantee the degeneration up to $1/p_a$ of the term \mathcal{P}_1 from (28).

However, DASHA-PP-MVR has also suboptimal term \mathcal{P}_2 . This suboptimality is tightly connected with

the suboptimality of $B_{\rm init}$ in the general nonconvex case, which we discuss in Section 6.3, and it also

appears in the analysis of DASHA-MVR (Tyurin and Richtárik, 2023). Let us provide the counterpart

of Corollary 4. The corollary reveals that we can escape regimes when \mathcal{P}_2 is the bottleneck by

choosing the parameters of the compressors.

Corollary 5. Suppose that assumptions of Theorem 10 hold, batch size $B \leq \min\left\{\frac{\sigma}{p_a\sqrt{\mu\varepsilon}n}, \frac{L_{\sigma}^2}{\widehat{L}^2}\right\}$,

we take RandK compressors with $K=\Theta\left(rac{Bd\sqrt{\muarepsilon n}}{\sigma}
ight)$. Then the communication complexity equals

$$\widetilde{\mathcal{O}}\left(rac{d\sigma}{p_{\mathrm{a}}\sqrt{\muarepsilon n}}+rac{dL_{\sigma}}{p_{\mathrm{a}}\mu\sqrt{n}}
ight),$$

and the expected number of stochastic gradient calculations per node equals

$$\widetilde{\mathcal{O}}\left(rac{\sigma^2}{p_{
m a}\mu narepsilon}+rac{\sigma L_{\sigma}}{p_{
m a}n\mu^{3/2}\sqrt{arepsilon}}
ight).$$

condition (see (Li et al., 2021a)). Moreover, DASHA-PP-MVR has the state-of-the-art communication

complexity of DASHA for a small enough μ .

⁸For simplicity, we omitted $\frac{d}{\zeta_C}$ term from the complexity in the stochastic setting, where ζ_C is defined in Definition 12. For instance, for the RandK compressor (see Definition 5 and Theorem 6), $\zeta_C = K$ and $\frac{d}{\zeta_C} = \Theta(\omega)$.

749 F.4 Proofs of Theorems

The following proofs almost repeat the proofs from Section E. And one of the main changes is that instead of Lemma 3, we use the following lemma.

752 F.4.1 Standard Lemma under Polyak-Łojasiewicz Condition

753 **Lemma 11.** Suppose that Assumptions 1 and 9 hold and

$$\mathrm{E}\left[f(x^{t+1})\right] + \gamma \Psi^{t+1} \le \mathrm{E}\left[f(x^t)\right] - \frac{\gamma}{2} \mathrm{E}\left[\left\|\nabla f(x^t)\right\|^2\right] + (1 - \gamma \mu) \gamma \Psi^t + \gamma C_{t+1}$$

where Ψ^t is a sequence of numbers, $\Psi^t \geq 0$ for all $t \in [T]$, constant $C \geq 0$, constant $\mu > 0$, and constant $\gamma \in (0, 1/\mu)$. Then

$$E[f(x^{T}) - f^{*}] \le (1 - \gamma \mu)^{T} ((f(x^{0}) - f^{*}) + \gamma \Psi^{0}) + \frac{C}{\mu}.$$
 (30)

756 *Proof.* We subtract f^* and use PŁ-condition (27) to get

$$\begin{split} \mathbf{E}\left[f(x^{t+1}) - f^*\right] + \gamma \Psi^{t+1} &\leq \mathbf{E}\left[f(x^t) - f^*\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^t)\right\|^2\right] + \gamma \Psi^t + \gamma C \\ &\leq (1 - \gamma \mu) \mathbf{E}\left[f(x^t) - f^*\right] + (1 - \gamma \mu) \gamma \Psi^t + \gamma C \\ &= (1 - \gamma \mu) \left(\mathbf{E}\left[f(x^t) - f^*\right] + \gamma \Psi^t\right) + \gamma C. \end{split}$$

757 Unrolling the inequality, we have

$$\mathbb{E}\left[f(x^{t+1}) - f^*\right] + \gamma \Psi^{t+1} \leq (1 - \gamma \mu)^{t+1} \left(\left(f(x^0) - f^* \right) + \gamma \Psi^0 \right) + \gamma C \sum_{i=0}^t (1 - \gamma \mu)^i$$

$$\leq (1 - \gamma \mu)^{t+1} \left(\left(f(x^0) - f^* \right) + \gamma \Psi^0 \right) + \frac{C}{\mu}.$$

758 It is left to note that $\Psi^t \geq 0$ for all $t \in [T]$.

759 F.4.2 Generic Lemma

We now provide the counterpart of Lemma 6.

Lemma 12. Suppose that Assumptions 2, 7, 8 and 9 hold and let us take $a=\frac{p_a}{2\omega+1}$, then

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t}) - \frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + (1 - \gamma\mu)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu)\frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \frac{10\gamma(2\omega+1)\omega}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|k_{i}^{t+1}\right\|^{2}\right]. \end{split}$$

Proof. Let us fix some constants $\kappa, \eta \in [0, \infty)$ that we will define later. Using the same reasoning as in Lemma 6, we can get

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] \\ & + \kappa \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \eta \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t}) - \frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \end{split}$$

$$\begin{split} & + \left(\gamma + \kappa \left(1 - a\right)^{2}\right) \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] \\ & + \left(\frac{\kappa a^{2}(\left(2\omega + 1\right)p_{a} - p_{aa})}{np_{a}^{2}} + \eta \left(\frac{a^{2}(2\omega + 1 - p_{a})}{p_{a}} + \left(1 - a\right)^{2}\right)\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \left(\frac{2\kappa\omega}{np_{a}} + \frac{2\eta\omega}{p_{a}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|k_{i}^{t+1}\right\|^{2}\right]. \end{split}$$

Let us take $\kappa = \frac{2\gamma}{a}$. One can show that $\gamma + \kappa \left(1 - a\right)^2 \leq \left(1 - \frac{a}{2}\right)\kappa$, and thus

$$\begin{split} & \operatorname{E}\left[f(x^{t+1})\right] \\ & + \frac{2\gamma}{a}\operatorname{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \eta\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & \leq \operatorname{E}\left[f(x^{t}) - \frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \left(1 - \frac{a}{2}\right)\frac{2\gamma}{a}\operatorname{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] \\ & + \left(\frac{2\gamma a((2\omega + 1)\,p_{a} - p_{aa})}{np_{a}^{2}} + \eta\left(\frac{a^{2}(2\omega + 1 - p_{a})}{p_{a}} + (1 - a)^{2}\right)\right)\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \left(\frac{4\gamma\omega}{anp_{a}} + \frac{2\eta\omega}{p_{a}}\right)\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|k_{i}^{t+1}\right\|^{2}\right]. \end{split}$$

Considering the choice of a, one can show that $\left(\frac{a^2(2\omega+1-p_a)}{p_a}+(1-a)^2\right)\leq 1-a$. If we take $\eta=\frac{4\gamma((2\omega+1)p_a-p_{aa})}{np_a^2}$, then $\left(\frac{2\gamma a((2\omega+1)p_a-p_{aa})}{np_a^2}+\eta\left(\frac{a^2(2\omega+1-p_a)}{p_a}+(1-a)^2\right)\right)\leq \left(1-\frac{a}{2}\right)\eta$ and

$$\begin{split} & \quad E\left[f(x^{t+1})\right] \\ & \quad + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t}) - \frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & \quad + \left(1 - \frac{a}{2}\right)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \left(1 - \frac{a}{2}\right)\frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & \quad + \left(\frac{2\gamma(2\omega+1)\omega}{np_{\mathbf{a}}^{2}} + \frac{8\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})\omega}{np_{\mathbf{a}}^{3}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|k_{i}^{t+1}\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t}) - \frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & \quad + \left(1 - \frac{a}{2}\right)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \left(1 - \frac{a}{2}\right)\frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & \quad + \frac{10\gamma(2\omega+1)\omega}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|k_{i}^{t+1}\right\|^{2}\right]. \end{split}$$

It it left to consider that $\gamma \leq \frac{a}{2\mu}$, and therefore $1 - \frac{a}{2} \leq 1 - \gamma \mu$.

768 F.4.3 Proof for DASHA-PP under PŁ-condition

Theorem 8. Suppose that Assumption 1, 2, 3, 7, 8 and 9 hold. Let us take $a = \frac{p_a}{2\omega+1}$, $b = \frac{p_a}{2-p_a}$,

$$\gamma \le \min \left\{ \left(L + \sqrt{\frac{200\omega \left(2\omega + 1\right)}{np_{\mathrm{a}}^2} + \frac{48}{np_{\mathrm{a}}^2} \left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)} \widehat{L} \right)^{-1}, \frac{a}{4\mu} \right\},$$

769 and $h_i^0=g_i^0=\nabla f_i(x^0)$ for all $i\in[n]$ in Algorithm I (DASHA-PP), then $\mathrm{E}\left[f(x^T)\right]-f^*\leq$ 770 $(1-\gamma\mu)^T\Delta_0$.

Proof. Let us fix constants $\nu, \rho \in [0, \infty)$ that we will define later. Considering Lemma 12, Lemma 7,

and the law of total expectation, we obtain

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \nu\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t}) - \frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + (1 - \gamma\mu)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu)\frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \frac{10\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[2\widehat{L}^{2}\left\|x^{t+1} - x^{t}\right\|^{2} + 2b^{2}\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \nu\mathbf{E}\left[\frac{2\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)\widehat{L}^{2}}{np_{\mathbf{a}}^{2}}\left\|x^{t+1} - x^{t}\right\|^{2} + \frac{2b^{2}\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{n^{2}p_{\mathbf{a}}^{2}}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2} + (1 - b)^{2}\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \rho\mathbf{E}\left[\frac{2(1 - p_{\mathbf{a}})}{p_{\mathbf{a}}}\widehat{L}^{2}\left\|x^{t+1} - x^{t}\right\|^{2} + \left(\frac{2b^{2}(1 - p_{\mathbf{a}})}{p_{\mathbf{a}}} + (1 - b)^{2}\right)\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right]. \end{split}$$

773 After rearranging the terms, we get

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \nu\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + (1 - \gamma\mu)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu)\frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{20\gamma\omega\left(2\omega+1\right)\,\hat{L}^{2}}{np_{\mathbf{a}}^{2}} - \nu\frac{2\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)\,\hat{L}^{2}}{np_{\mathbf{a}}^{2}} - \rho\frac{2(1 - p_{\mathbf{a}})\hat{L}^{2}}{p_{\mathbf{a}}}\right)\mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \left(\gamma + \nu(1 - b)^{2}\right)\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \left(\frac{20b^{2}\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \nu\frac{2b^{2}\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}} + \rho\left(\frac{2b^{2}(1 - p_{\mathbf{a}})}{p_{\mathbf{a}}} + (1 - b)^{2}\right)\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right]. \end{split}$$

774 By taking $\nu=\frac{2\gamma}{b},$ one can show that $\left(\gamma+\nu(1-b)^2\right)\leq \left(1-\frac{b}{2}\right)\nu,$ and

$$\mathrm{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathrm{a}}}\mathrm{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)\,p_{\mathrm{a}} - p_{\mathrm{aa}})}{np_{\mathrm{a}}^{2}}\mathrm{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right]$$

$$\begin{split} & + \frac{2\gamma}{b} \mathbf{E} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^2 \right] + \rho \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^n \left\| h_i^{t+1} - \nabla f_i(x^{t+1}) \right\|^2 \right] \\ & \leq \mathbf{E} \left[f(x^t) \right] - \frac{\gamma}{2} \mathbf{E} \left[\left\| \nabla f(x^t) \right\|^2 \right] \\ & + (1 - \gamma\mu) \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbf{E} \left[\left\| g^t - h^t \right\|^2 \right] + (1 - \gamma\mu) \frac{4\gamma((2\omega + 1) \, p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^2} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^n \left\| g_i^t - h_i^t \right\|^2 \right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{20\gamma\omega\left(2\omega + 1\right) \, \hat{L}^2}{np_{\mathbf{a}}^2} - \frac{4\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right) \, \hat{L}^2}{bnp_{\mathbf{a}}^2} - \rho \frac{2(1 - p_{\mathbf{a}}) \hat{L}^2}{p_{\mathbf{a}}} \right) \mathbf{E} \left[\left\| x^{t+1} - x^t \right\|^2 \right] \\ & + \left(1 - \frac{b}{2} \right) \frac{2\gamma}{b} \mathbf{E} \left[\left\| h^t - \nabla f(x^t) \right\|^2 \right] \\ & + \left(\frac{20b^2\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^2} + \frac{4\gamma b \left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^2} + \rho \left(\frac{2b^2(1 - p_{\mathbf{a}})}{p_{\mathbf{a}}} + (1 - b)^2 \right) \right) \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^n \left\| h_i^t - \nabla f_i(x^t) \right\|^2 \right]. \end{split}$$

Note that $b = \frac{p_a}{2-p_a}$, thus

$$\begin{split} &\left(\frac{20b^2\gamma\omega(2\omega+1)}{np_{\mathrm{a}}^2} + \frac{4\gamma b\left(p_{\mathrm{a}} - p_{\mathrm{aa}}\right)}{np_{\mathrm{a}}^2} + \rho\left(\frac{2b^2(1-p_{\mathrm{a}})}{p_{\mathrm{a}}} + (1-b)^2\right)\right) \\ &\leq \left(\frac{20b^2\gamma\omega(2\omega+1)}{np_{\mathrm{a}}^2} + \frac{4\gamma b\left(p_{\mathrm{a}} - p_{\mathrm{aa}}\right)}{np_{\mathrm{a}}^2} + \rho\left(1-b\right)\right). \end{split}$$

776 And if we take $ho=rac{40b\gamma\omega(2\omega+1)}{np_{\rm a}^2}+rac{8\gamma(p_{\rm a}-p_{\rm aa})}{np_{\rm a}^2},$ then

$$\left(\frac{20b^2\gamma\omega(2\omega+1)}{np_{\rm a}^2} + \frac{4\gamma b\left(p_{\rm a}-p_{\rm aa}\right)}{np_{\rm a}^2} + \rho\left(1-b\right)\right) \leq \left(1-\frac{b}{2}\right)\rho,$$

777 and

$$\begin{split} & \operatorname{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\operatorname{a}}}\operatorname{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)\,p_{\operatorname{a}} - p_{\operatorname{aa}})}{np_{\operatorname{a}}^{2}}\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{2\gamma}{b}\operatorname{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{40b\gamma\omega(2\omega+1)}{np_{\operatorname{a}}^{2}} + \frac{8\gamma\left(p_{\operatorname{a}} - p_{\operatorname{aa}}\right)}{np_{\operatorname{a}}^{2}}\right)\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \operatorname{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\operatorname{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \left(1 - \gamma\mu\right)\frac{2\gamma(2\omega+1)}{p_{\operatorname{a}}}\operatorname{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \left(1 - \gamma\mu\right)\frac{4\gamma((2\omega+1)\,p_{\operatorname{a}} - p_{\operatorname{aa}})}{np_{\operatorname{a}}^{2}}\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{20\gamma\omega\left(2\omega+1\right)\widehat{L}^{2}}{np_{\operatorname{a}}^{2}} - \frac{4\gamma\left(p_{\operatorname{a}} - p_{\operatorname{aa}}\right)\widehat{L}^{2}}{bnp_{\operatorname{a}}^{2}} - \frac{16\gamma\left(p_{\operatorname{a}} - p_{\operatorname{aa}}\right)\left(1 - p_{\operatorname{a}}\right)\widehat{L}^{2}}{np_{\operatorname{a}}^{3}}\right)\operatorname{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & - \left(1 - \frac{b}{2}\right)\frac{2\gamma}{b}\operatorname{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(1 - \frac{b}{2}\right)\left(\frac{40b\gamma\omega(2\omega+1)}{np_{\operatorname{a}}^{2}} + \frac{8\gamma\left(p_{\operatorname{a}} - p_{\operatorname{aa}}\right)}{np_{\operatorname{a}}^{2}}\right)\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right]. \end{split}$$

778 Due to $\frac{p_a}{2} \le b \le p_a$, we have

$$\mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] + \frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{40b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{8\gamma(p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}}\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right]$$

$$\leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ + \left(1 - \gamma\mu\right)\frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \left(1 - \gamma\mu\right)\frac{4\gamma((2\omega + 1)p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{100\gamma\omega\left(2\omega + 1\right)\widehat{L}^{2}}{np_{\mathbf{a}}^{2}} - \frac{24\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)\widehat{L}^{2}}{np_{\mathbf{a}}^{3}}\right)\mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ + \left(1 - \frac{b}{2}\right)\frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(1 - \frac{b}{2}\right)\left(\frac{40b\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right].$$

Using Lemma 4 and the assumption about γ , we ge

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{40b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \left(1 - \gamma\mu\right)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \left(1 - \gamma\mu\right)\frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \left(1 - \frac{b}{2}\right)\frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(1 - \frac{b}{2}\right)\left(\frac{40b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right]. \end{split}$$

Note that $\gamma \leq \frac{a}{4\mu} \leq \frac{p_a}{4\mu} \leq \frac{b}{2\mu}$, thus $1 - \frac{b}{2} \leq 1 - \gamma\mu$ and

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{40b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + (1 - \gamma\mu)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu)\frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + (1 - \gamma\mu)\frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + (1 - \gamma\mu)\left(\frac{40b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right]. \end{split}$$

781 In the view of Lemma 11 with

$$\Psi^{t} = \frac{2(2\omega + 1)}{p_{a}} E\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{4((2\omega + 1)p_{a} - p_{aa})}{np_{a}^{2}} E\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] + \frac{2}{b} E\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(\frac{40b\omega(2\omega + 1)}{np_{a}^{2}} + \frac{8\gamma(p_{a} - p_{aa})}{np_{a}^{2}}\right) E\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right],$$

we can conclude the proof of the theorem.

783 F.4.4 Proof for DASHA-PP-PAGE under PŁ-condition

Theorem 9. Suppose that Assumption 1, 2, 3, 7, 4, 8, and 9 hold. Let us take $a = \frac{p_a}{2\omega + 1}$, probability $p_{page} = \frac{B}{m+B}$, $b = \frac{p_{page}p_a}{2-p_a}$,

$$\gamma \leq \min \left\{ \left(L + \sqrt{\frac{200\omega(2\omega+1)}{np_{\mathrm{a}}^2} \left(\widehat{L}^2 + \frac{(1-p_{\mathrm{page}})L_{\mathrm{max}}^2}{B} \right) + \frac{48}{np_{\mathrm{a}}^2p_{\mathrm{page}}} \left(\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}} \right) \widehat{L}^2 + \frac{(1-p_{\mathrm{page}})L_{\mathrm{max}}^2}{B} \right) \right)^{-1}, \frac{a}{2\mu}, \frac{b}{2\mu} \right\},$$

784 and $h_i^0=g_i^0=\nabla f_i(x^0)$ for all $i\in[n]$ in Algorithm I (DASHA-PP-PAGE), then $\mathrm{E}\left[f(x^T)\right]-f^*\leq (1-\gamma\mu)^T\Delta_0$.

786 *Proof.* Let us fix constants $\nu, \rho \in [0, \infty)$ that we will define later. Considering Lemma 12, Lemma 8, and the law of total expectation, we obtain

$$\begin{split} & \operatorname{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\operatorname{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g^{t+1}_{i} - h^{t+1}_{i}\right\|^{2}\right] \\ & + \nu\operatorname{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h^{t+1}_{i} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \operatorname{E}\left[f(x^{t}) - \frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + (1 - \gamma\mu)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\operatorname{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu)\frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g^{t}_{i} - h^{t}_{i}\right\|^{2}\right] \\ & + \frac{10\gamma(2\omega+1)\omega}{np_{\mathbf{a}}}\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|k^{t+1}_{i}\right\|^{2}\right] \\ & + \nu\operatorname{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h^{t+1}_{i} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \operatorname{E}\left[f(x^{t}) - \frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + (1 - \gamma\mu)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\operatorname{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu)\frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h^{t}_{i} - h^{t}_{i}\right\|^{2}\right] \\ & + \frac{10\gamma(2\omega+1)\omega}{np_{\mathbf{a}}^{2}}\operatorname{E}\left[\left(2\hat{L}^{2} + \frac{(1 - p_{\mathbf{p}\mathbf{a}\mathbf{g}\mathbf{e}})L^{2}_{\mathbf{max}}}{B}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \frac{2b^{2}}{p_{\mathbf{p}\mathbf{a}\mathbf{g}\mathbf{e}}\frac{1}{n}\sum_{i=1}^{n}\left\|h^{t}_{i} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \nu\operatorname{E}\left(\left(\frac{2(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})\hat{L}^{2}}{np_{\mathbf{a}}^{2}} + \frac{(1 - p_{\mathbf{p}\mathbf{a}\mathbf{g}\mathbf{e}})L^{2}_{\mathbf{max}}}{np_{\mathbf{a}}B}\right)\left\|x^{t+1} - x^{t}\right\|^{2} \\ & + \frac{2(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})\hat{L}^{2}}{np_{\mathbf{a}}^{2}} + \frac{(1 - p_{\mathbf{p}\mathbf{a}\mathbf{g}\mathbf{e}})L^{2}_{\mathbf{max}}}{p_{\mathbf{a}}B}\right)\left\|x^{t+1} - x^{t}\right\|^{2} \\ & + \rho\operatorname{E}\left(\left(\frac{2(1 - p_{\mathbf{a}})\hat{L}^{2}}{p_{\mathbf{a}}} + \frac{(1 - p_{\mathbf{p}\mathbf{a}\mathbf{g}\mathbf{e}})L^{2}_{\mathbf{max}}}{p_{\mathbf{a}}B}\right)\left\|x^{t+1} - x^{t}\right\|^{2} \\ & + \left(\frac{2(1 - p_{\mathbf{a}})\hat{L}^{2}}{p_{\mathbf{a}}} + \frac{(1 - p_{\mathbf{p}\mathbf{a}\mathbf{g}\mathbf{e}})L^{2}_{\mathbf{max}}}{p_{\mathbf{a}}B}\right)\left\|x^{t+1} - x^{t}\right\|^{2} \\ & + \left(\frac{2(1 - p_{\mathbf{a}})\hat{L}^{2}}{p_{\mathbf{a}}} + \frac{(1 - p_{\mathbf{p}\mathbf{a}\mathbf{g}\mathbf{e}})L^{2}_{\mathbf{max}}}{p_{\mathbf{a}}}\right)\left\|x^{t+1} - x^{t}\right\|^{2} \\ & + \left(\frac{2(1 - p_{\mathbf{a}})\hat{L}^$$

After rearranging the terms, we get

$$E\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{a}}E\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)p_{a} - p_{aa})}{np_{a}^{2}}E\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] + \nu E\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho E\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ \leq E\left[f(x^{t})\right] - \frac{\gamma}{2}E\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ + (1 - \gamma\mu)\frac{2\gamma(2\omega+1)}{p_{a}}E\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu)\frac{4\gamma((2\omega+1)p_{a} - p_{aa})}{np_{a}^{2}}E\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right]$$

$$\begin{split} &-\left(\frac{1}{2\gamma}-\frac{L}{2}-\frac{10\gamma\omega(2\omega+1)}{np_{\rm a}^2}\left(2\hat{L}^2+\frac{(1-p_{\rm page})L_{\rm max}^2}{B}\right)\right.\\ &-\nu\left(\frac{2\left(p_{\rm a}-p_{\rm aa}\right)\hat{L}^2}{np_{\rm a}^2}+\frac{(1-p_{\rm page})L_{\rm max}^2}{np_{\rm a}B}\right)-\rho\left(\frac{2\left(1-p_{\rm a}\right)\hat{L}^2}{p_{\rm a}}+\frac{(1-p_{\rm page})L_{\rm max}^2}{p_{\rm a}B}\right)\right)\mathbf{E}\left[\left\|x^{t+1}-x^t\right\|^2\right]\\ &+\left(\gamma+\nu\left(p_{\rm page}\left(1-\frac{b}{p_{\rm page}}\right)^2+(1-p_{\rm page})\right)\right)\mathbf{E}\left[\left\|h^t-\nabla f(x^t)\right\|^2\right]\\ &+\left(\frac{20b^2\gamma\omega(2\omega+1)}{np_{\rm a}^2p_{\rm page}}+\frac{2\nu\left(p_{\rm a}-p_{\rm aa}\right)b^2}{np_{\rm a}^2p_{\rm page}}\right.\\ &+\rho\left(\frac{2\left(1-p_{\rm a}\right)b^2}{p_{\rm a}p_{\rm page}}+p_{\rm page}\left(1-\frac{b}{p_{\rm page}}\right)^2+(1-p_{\rm page})\right)\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^n\left\|h_i^t-\nabla f_i(x^t)\right\|^2\right]. \end{split}$$

Due to $b=\frac{p_{\text{page}}p_{\text{a}}}{2-p_{\text{a}}}\leq p_{\text{page}},$ one can show that $\left(p_{\text{page}}\left(1-\frac{b}{p_{\text{page}}}\right)^2+(1-p_{\text{page}})\right)\leq 1-b.$ Thus, if we take $\nu=\frac{2\gamma}{b}$, then

$$\left(\gamma + \nu \left(p_{\text{page}}\left(1 - \frac{b}{p_{\text{page}}}\right)^2 + (1 - p_{\text{page}})\right)\right) \le \gamma + \nu(1 - b) = \left(1 - \frac{b}{2}\right)\nu,$$

789 therefore

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + (1-\gamma\mu)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1-\gamma\mu)\frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{10\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}}\left(2\hat{L}^{2} + \frac{(1-p_{\mathsf{page}})L_{\mathsf{max}}^{2}}{B}\right) - \rho\left(\frac{2\left(1-p_{\mathbf{a}}\right)\hat{L}^{2}}{p_{\mathbf{a}}} + \frac{(1-p_{\mathsf{page}})L_{\mathsf{max}}^{2}}{p_{\mathbf{a}}B}\right)\right)\mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \left(1 - \frac{b}{2}\right)\frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \left(\frac{20b^{2}\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}p_{\mathsf{page}}} + \frac{4\gamma\left(p_{\mathbf{a}} - p_{\mathsf{a}\mathbf{a}}\right)b}{np_{\mathbf{a}}^{2}p_{\mathsf{page}}} \\ & + \rho\left(\frac{2\left(1-p_{\mathbf{a}}\right)b^{2}}{p_{\mathbf{a}}p_{\mathsf{page}}} + p_{\mathsf{page}}\left(1 - \frac{b}{p_{\mathsf{page}}}\right)^{2} + (1-p_{\mathsf{page}})\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right]. \end{split}$$

Next, with the choice of $b = \frac{p_{\text{page}}p_{\text{a}}}{2-p_{\text{a}}}$, we ensure that

$$\left(\frac{2(1-p_{\rm a})\,b^2}{p_{\rm a}p_{\rm page}} + p_{\rm page}\left(1 - \frac{b}{p_{\rm page}}\right)^2 + (1-p_{\rm page})\right) \le 1 - b.$$

If we take $\rho = \frac{40b\gamma\omega(2\omega+1)}{np_s^2p_{\text{page}}} + \frac{8\gamma(p_a-p_{\text{aa}})}{np_s^2p_{\text{page}}}$, then

$$\left(\frac{20b^2\gamma\omega(2\omega+1)}{np_{\mathrm{a}}^2p_{\mathrm{page}}} + \frac{4\gamma\left(p_{\mathrm{a}} - p_{\mathrm{aa}}\right)b}{np_{\mathrm{a}}^2p_{\mathrm{page}}} + \rho\left(\frac{2\left(1 - p_{\mathrm{a}}\right)b^2}{p_{\mathrm{a}}p_{\mathrm{page}}} + p_{\mathrm{page}}\left(1 - \frac{b}{p_{\mathrm{page}}}\right)^2 + \left(1 - p_{\mathrm{page}}\right)\right)\right) \leq \left(1 - \frac{b}{2}\right)\rho,$$

790 therefore

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{2\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{40b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}p_{\mathsf{page}}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathsf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}p_{\mathsf{page}}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \left(1 - \gamma\mu\right)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \left(1 - \gamma\mu\right)\frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathsf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{10\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}}\left(2\hat{L}^{2} + \frac{(1 - p_{\mathsf{page}})L_{\mathsf{max}}^{2}}{B}\right) \\ & - \frac{2\gamma}{bmp_{\mathbf{a}}}\left(2\left(1 - \frac{p_{\mathsf{a}\mathbf{a}}}{p_{\mathbf{a}}}\right)\hat{L}^{2} + \frac{(1 - p_{\mathsf{page}})L_{\mathsf{max}}^{2}}{B}\right) \\ & - \left(\frac{40b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}p_{\mathsf{page}}} + \frac{8\gamma\left(1 - \frac{p_{\mathsf{a}\mathbf{a}}}{p_{\mathbf{a}}}\right)}{np_{\mathbf{a}}^{2}p_{\mathsf{page}}}\right)\left(2\left(1 - p_{\mathbf{a}}\right)\hat{L}^{2} + \frac{(1 - p_{\mathsf{page}})L_{\mathsf{max}}^{2}}{B}\right)\right) \mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \left(1 - \frac{b}{2}\right)\frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(1 - \frac{b}{2}\right)\left(\frac{40b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}p_{\mathsf{page}}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathsf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}p_{\mathsf{page}}}\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right]. \end{split}$$

Let us simplify the inequality. First, due to $b \geq \frac{p_{\text{page}}p_{\text{a}}}{2}$, we have

$$\frac{2\gamma}{bnp_{\mathrm{a}}}\left(2\left(1-\frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)\widehat{L}^2+\frac{(1-p_{\mathrm{page}})L_{\mathrm{max}}^2}{B}\right)\leq \frac{8\gamma}{np_{\mathrm{a}}^2p_{\mathrm{page}}}\left(\left(1-\frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)\widehat{L}^2+\frac{(1-p_{\mathrm{page}})L_{\mathrm{max}}^2}{B}\right).$$

Second, due to $b \leq p_{\rm a} p_{\rm page}$ and $p_{\rm aa} \leq p_{\rm a}^2$, we get

$$\begin{split} &\left(\frac{40b\gamma\omega(2\omega+1)}{np_{\rm a}^{3}p_{\rm page}} + \frac{8\gamma\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)}{np_{\rm a}^{2}p_{\rm page}}\right) \left(2\left(1 - p_{\rm a}\right)\widehat{L}^{2} + \frac{(1 - p_{\rm page})L_{\rm max}^{2}}{B}\right) \\ &\leq \left(\frac{40\gamma\omega(2\omega+1)}{np_{\rm a}^{2}} + \frac{8\gamma\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)}{np_{\rm a}^{2}p_{\rm page}}\right) \left(2\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)\widehat{L}^{2} + \frac{(1 - p_{\rm page})L_{\rm max}^{2}}{B}\right) \\ &\leq \frac{80\gamma\omega(2\omega+1)}{np_{\rm a}^{2}} \left(\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)\widehat{L}^{2} + \frac{(1 - p_{\rm page})L_{\rm max}^{2}}{B}\right) \\ &+ \frac{16\gamma\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)}{np_{\rm a}^{2}p_{\rm page}} \left(\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)\widehat{L}^{2} + \frac{(1 - p_{\rm page})L_{\rm max}^{2}}{B}\right) \\ &\leq \frac{80\gamma\omega(2\omega+1)}{np_{\rm a}^{2}} \left(\widehat{L}^{2} + \frac{(1 - p_{\rm page})L_{\rm max}^{2}}{B}\right) \\ &+ \frac{16\gamma}{np_{\rm a}^{2}p_{\rm page}} \left(\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)\widehat{L}^{2} + \frac{(1 - p_{\rm page})L_{\rm max}^{2}}{B}\right). \end{split}$$

Combining all bounds together, we obtain the following inequality:

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{40b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}p_{\mathsf{page}}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathsf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}p_{\mathsf{page}}}\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \end{split}$$

$$\begin{split} &+ (1 - \gamma \mu) \, \frac{2 \gamma (2 \omega + 1)}{p_{\mathrm{a}}} \mathbf{E} \left[\left\| g^{t} - h^{t} \right\|^{2} \right] + (1 - \gamma \mu) \, \frac{4 \gamma ((2 \omega + 1) \, p_{\mathrm{a}} - p_{\mathrm{aa}})}{n p_{\mathrm{a}}^{2}} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \right] \\ &- \left(\frac{1}{2 \gamma} - \frac{L}{2} - \frac{100 \gamma \omega (2 \omega + 1)}{n p_{\mathrm{a}}^{2}} \left(\widehat{L}^{2} + \frac{(1 - p_{\mathrm{page}}) L_{\mathrm{max}}^{2}}{B} \right) \right. \\ &- \left. \frac{24 \gamma}{n p_{\mathrm{a}}^{2} p_{\mathrm{page}}} \left(\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}} \right) \widehat{L}^{2} + \frac{(1 - p_{\mathrm{page}}) L_{\mathrm{max}}^{2}}{B} \right) \right) \mathbf{E} \left[\left\| x^{t+1} - x^{t} \right\|^{2} \right] \\ &+ \left(1 - \frac{b}{2} \right) \frac{2 \gamma}{b} \mathbf{E} \left[\left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \right] + \left(1 - \frac{b}{2} \right) \left(\frac{40 b \gamma \omega (2 \omega + 1)}{n p_{\mathrm{a}}^{2} p_{\mathrm{page}}} + \frac{8 \gamma \left(p_{\mathrm{a}} - p_{\mathrm{aa}} \right)}{n p_{\mathrm{a}}^{2} p_{\mathrm{page}}} \right) \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \right]. \end{split}$$

Using Lemma 4 and the assumption about γ , we get

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{40b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}p_{\mathrm{page}}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{page}}}\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \left(1 - \gamma\mu\right)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \left(1 - \gamma\mu\right)\frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \left(1 - \frac{b}{2}\right)\frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(1 - \frac{b}{2}\right)\left(\frac{40b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}p_{\mathbf{page}}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}p_{\mathbf{page}}}\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right]. \end{split}$$

Note that $\gamma \leq \frac{b}{2\mu}$, thus $1 - \frac{b}{2} \leq 1 - \gamma \mu$ and

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{40b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}p_{\mathsf{page}}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathsf{aa}}\right)}{np_{\mathbf{a}}^{2}p_{\mathsf{page}}}\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + (1-\gamma\mu)\,\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1-\gamma\mu)\,\frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathsf{aa}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + (1-\gamma\mu)\,\frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + (1-\gamma\mu)\left(\frac{40b\gamma\omega(2\omega+1)}{np_{\mathbf{a}}^{2}p_{\mathsf{page}}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathsf{aa}}\right)}{np_{\mathbf{a}}^{2}p_{\mathsf{page}}}\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right]. \end{split}$$

795 It is left to apply Lemma 11 with

$$\begin{split} \Psi^{t} &= \frac{2(2\omega + 1)}{p_{a}} E\left[\left\| g^{t} - h^{t} \right\|^{2} \right] + \frac{4((2\omega + 1) p_{a} - p_{aa})}{np_{a}^{2}} E\left[\frac{1}{n} \sum_{i=1}^{n} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \right] \\ &+ \frac{2}{b} E\left[\left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \right] + \left(\frac{40b\omega(2\omega + 1)}{np_{a}^{2}p_{page}} + \frac{8\left(p_{a} - p_{aa}\right)}{np_{a}^{2}p_{page}} \right) E\left[\frac{1}{n} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \right] \end{split}$$

796 to conclude the proof.

797 F.4.5 Proof for DASHA-PP-MVR under PŁ-condition

Theorem 10. Suppose that Assumption 1, 2, 3, 7, 5, 6, 8 and 9 hold. Let us take $a = \frac{p_a}{2\omega+1}$, $b \in \left(0, \frac{p_a}{2-p_a}\right]$,

$$\gamma \leq \min \left\{ \left(L + \sqrt{\frac{200\omega(2\omega+1)}{np_{\mathrm{a}}^2} \left(\frac{(1-b)^2L_\sigma^2}{B} + \widehat{L}^2 \right) + \frac{40}{np_{\mathrm{a}}b} \left(\frac{(1-b)^2L_\sigma^2}{B} + \left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}} \right) \widehat{L}^2 \right)} \right)^{-1}, \frac{a}{2\mu}, \frac{b}{2\mu} \right\},$$

798 and $h_i^0=g_i^0$ for all $i\in[n]$ in Algorithm ${\it 1\hspace{-0.07cm} l}$ (DASHA-PP-MVR), then

$$\begin{split} & E\left[f(x^{T}) - f^{*}\right] \\ & \leq (1 - \gamma\mu)^{T} \left(\Delta_{0} + \frac{2\gamma}{b} \left\|h^{0} - \nabla f(x^{0})\right\|^{2} + \left(\frac{40\gamma b\omega(2\omega + 1)}{np_{a}^{2}} + \frac{8\gamma\left(p_{a} - p_{aa}\right)}{np_{a}^{2}}\right) \frac{1}{n} \sum_{i=1}^{n} \left\|h_{i}^{0} - \nabla f_{i}(x^{0})\right\|^{2} \right) \\ & + \frac{1}{\mu} \left(\frac{100b^{2}\omega(2\omega + 1)}{p_{a}^{2}} + \frac{20b}{p_{a}}\right) \frac{\sigma^{2}}{nB}. \end{split}$$

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t}) - \frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + (1 - \gamma\mu)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu)\frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \frac{10\gamma(2\omega+1)\omega}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|k_{i}^{t+1}\right\|^{2}\right]. \end{split}$$

Proof. Let us fix constants $\nu, \rho \in [0, \infty)$ that we will define later. Considering Lemma 12, Lemma 10, and the law of total expectation, we obtain

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \nu\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t}) - \frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + (1 - \gamma\mu)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu)\frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \frac{10\gamma(2\omega+1)\omega}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|k_{i}^{t+1}\right\|^{2}\right] \\ & + \nu\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t}) - \frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + (1 - \gamma\mu)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu)\frac{4\gamma((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \end{split}$$

$$\begin{split} &+\frac{10\gamma\omega(2\omega+1)}{np_{\mathrm{a}}^{2}}\mathrm{E}\left[\frac{2b^{2}\sigma^{2}}{B}+\left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{B}+2\widehat{L}^{2}\right)\left\|x^{t+1}-x^{t}\right\|^{2}+2b^{2}\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t}-\nabla f_{i}(x^{t})\right\|^{2}\right]\\ &+\nu\mathrm{E}\left(\frac{2b^{2}\sigma^{2}}{np_{\mathrm{a}}B}+\left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{np_{\mathrm{a}}B}+\frac{2\left(p_{\mathrm{a}}-p_{\mathrm{aa}}\right)\widehat{L}^{2}}{np_{\mathrm{a}}^{2}}\right)\left\|x^{t+1}-x^{t}\right\|^{2}\\ &+\frac{2\left(p_{\mathrm{a}}-p_{\mathrm{aa}}\right)b^{2}}{n^{2}p_{\mathrm{a}}^{2}}\sum_{i=1}^{n}\left\|h_{i}^{t}-\nabla f_{i}(x^{t})\right\|^{2}+\left(1-b\right)^{2}\left\|h^{t}-\nabla f(x^{t})\right\|^{2}\right)\\ &+\rho\mathrm{E}\left(\frac{2b^{2}\sigma^{2}}{p_{\mathrm{a}}B}+\left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{p_{\mathrm{a}}B}+\frac{2(1-p_{\mathrm{a}})\widehat{L}^{2}}{p_{\mathrm{a}}}\right)\left\|x^{t+1}-x^{t}\right\|^{2}\\ &+\left(\frac{2(1-p_{\mathrm{a}})b^{2}}{p_{\mathrm{a}}}+\left(1-b\right)^{2}\right)\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t}-\nabla f_{i}(x^{t})\right\|^{2}\right). \end{split}$$

801 After rearranging the terms, we get

$$\begin{split} & \mathbb{E}\left[f(x^{t+1})\right] + (1 - \gamma\mu) \, \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbb{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + (1 - \gamma\mu) \, \frac{4\gamma((2\omega + 1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \nu\mathbb{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbb{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbb{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + (1 - \gamma\mu) \, \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbb{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu) \, \frac{4\gamma((2\omega + 1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{10\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} \left(\frac{2(1 - b)^{2}L_{\sigma}^{2}}{B} + 2\hat{L}^{2}\right) - \rho\left(\frac{2(1 - b)^{2}L_{\sigma}^{2}}{p_{\mathbf{a}}B} + \frac{2(1 - p_{\mathbf{a}})\hat{L}^{2}}{p_{\mathbf{a}}}\right)\right) \mathbb{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \left(\gamma + \nu\left(1 - b\right)^{2}\right) \mathbb{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \left(\frac{20b^{2}\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \frac{2\nu\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)\hat{L}^{2}}{np_{\mathbf{a}}^{2}} + \rho\left(\frac{2(1 - p_{\mathbf{a}})b^{2}}{p_{\mathbf{a}}} + (1 - b)^{2}\right)\right) \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{20b^{2}\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \nu\frac{2b^{2}}{np_{\mathbf{a}}} + \rho\frac{2b^{2}}{p_{\mathbf{a}}}\right)\frac{\sigma^{2}}{B}. \end{split}$$

By taking $\nu=\frac{2\gamma}{b}$, one can show that $\left(\gamma+\nu(1-b)^2\right)\leq \left(1-\frac{b}{2}\right)\nu$, and

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + (1 - \gamma\mu) \, \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + (1 - \gamma\mu) \, \frac{4\gamma((2\omega + 1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & \quad + \frac{2\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho \mathbf{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \quad \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & \quad + (1 - \gamma\mu) \, \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu) \, \frac{4\gamma((2\omega + 1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & \quad - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{10\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} \left(\frac{2(1 - b)^{2}L_{\sigma}^{2}}{B} + 2\widehat{L}^{2}\right) \end{split}$$

$$\begin{split} & -\frac{2\gamma}{b}\left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{np_{a}B} + \frac{2\left(p_{a} - p_{aa}\right)\widehat{L}^{2}}{np_{a}^{2}}\right) - \rho\left(\frac{2(1-b)^{2}L_{\sigma}^{2}}{p_{a}B} + \frac{2(1-p_{a})\widehat{L}^{2}}{p_{a}}\right)\right) \mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \left(1 - \frac{b}{2}\right)\frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \left(\frac{20b^{2}\gamma\omega(2\omega+1)}{np_{a}^{2}} + \frac{4\gamma\left(p_{a} - p_{aa}\right)b}{np_{a}^{2}} + \rho\left(\frac{2(1-p_{a})b^{2}}{p_{a}} + (1-b)^{2}\right)\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{20b^{2}\gamma\omega(2\omega+1)}{np_{a}^{2}} + \frac{4\gamma b}{np_{a}} + \rho\frac{2b^{2}}{p_{a}}\right)\frac{\sigma^{2}}{B}. \end{split}$$

Note that $b \leq \frac{p_a}{2-p_a}$, thus

$$\begin{split} &\left(\frac{20b^{2}\gamma\omega(2\omega+1)}{np_{a}^{2}} + \frac{4\gamma\left(p_{a} - p_{aa}\right)b}{np_{a}^{2}} + \rho\left(\frac{2(1-p_{a})b^{2}}{p_{a}} + (1-b)^{2}\right)\right) \\ &\leq \left(\frac{20b^{2}\gamma\omega(2\omega+1)}{np_{a}^{2}} + \frac{4\gamma\left(p_{a} - p_{aa}\right)b}{np_{a}^{2}} + \rho\left(1-b\right)\right). \end{split}$$

And if we take $ho=rac{40b\gamma\omega(2\omega+1)}{np_{\rm a}^2}+rac{8\gamma(p_{\rm a}-p_{\rm aa})}{np_{\rm a}^2},$ then

$$\left(\frac{20b^2\gamma\omega(2\omega+1)}{np_a^2} + \frac{4\gamma\left(p_a - p_{aa}\right)b}{np_a^2} + \rho\left(1 - b\right)\right) \le \rho,$$

805 and

$$\begin{split} & \mathbb{E}\left[f(x^{t+1})\right] + (1 - \gamma\mu) \, \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbb{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + (1 - \gamma\mu) \, \frac{4\gamma((2\omega + 1) \, p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & \quad + \frac{2\gamma}{b} \mathbb{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{40b\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \quad \leq \mathbb{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbb{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & \quad + (1 - \gamma\mu) \, \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbb{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu) \, \frac{4\gamma((2\omega + 1) \, p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & \quad - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{10\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} \left(\frac{2(1 - b)^{2}L_{\sigma}^{2}}{B} + 2\hat{L}^{2}\right) \right. \\ & \quad - \left(\frac{40b\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{3}} + \frac{8\gamma\left(1 - \frac{p_{\mathbf{aa}}}{p_{\mathbf{a}}}\right)\hat{L}^{2}\right) \\ & \quad - \left(\frac{40b\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{3}} + \frac{8\gamma\left(1 - \frac{p_{\mathbf{aa}}}{p_{\mathbf{a}}}\right)\hat{L}^{2}\right) \\ & \quad + \left(1 - \frac{b}{2}\right) \frac{2\gamma}{b} \mathbb{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(1 - \frac{b}{2}\right) \left(\frac{40b\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & \quad + \left(\frac{20b^{2}\gamma\omega(2\omega + 1)}{np^{2}} + \frac{4\gamma b}{np_{\mathbf{a}}} + \left(\frac{40b\gamma\omega(2\omega + 1)}{np^{2}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np^{2}}\right) \frac{2b^{2}}{p_{\mathbf{a}}} \frac{\sigma^{2}}{B}. \end{split}$$

Let us simplify the inequality. First, due to $b \le p_{\rm a}$ and $(1-p_{\rm a}) \le \left(1-\frac{p_{\rm aa}}{p_{\rm a}}\right)$, we have

$$\left(\frac{40b\gamma\omega(2\omega+1)}{np_{\rm a}^3} + \frac{2\gamma\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)}{np_{\rm a}^2}\right)\left(\frac{2(1-b)^2L_\sigma^2}{B} + 8(1-p_{\rm a})\widehat{L}^2\right)$$

$$\begin{split} &= \frac{40b\gamma\omega(2\omega+1)}{np_{\rm a}^3} \left(\frac{2(1-b)^2L_\sigma^2}{B} + 2(1-p_{\rm a})\widehat{L}^2\right) \\ &+ \frac{8\gamma\left(1-\frac{p_{\rm aa}}{p_{\rm a}}\right)}{np_{\rm a}^2} \left(\frac{2(1-b)^2L_\sigma^2}{B} + 2(1-p_{\rm a})\widehat{L}^2\right) \\ &\leq \frac{40\gamma\omega(2\omega+1)}{np_{\rm a}^2} \left(\frac{2(1-b)^2L_\sigma^2}{B} + 2\widehat{L}^2\right) \\ &+ \frac{8\gamma}{np_{\rm a}b} \left(\frac{2(1-b)^2L_\sigma^2}{B} + 2\left(1-\frac{p_{\rm aa}}{p_{\rm a}}\right)\widehat{L}^2\right), \end{split}$$

807 therefore

$$\begin{split} & \mathbb{E}\left[f(x^{t+1})\right] + (1 - \gamma\mu) \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbb{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + (1 - \gamma\mu) \frac{4\gamma((2\omega + 1)p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g^{t+1}_{i} - h^{t+1}_{i}\right\|^{2}\right] \\ & \quad + \frac{2\gamma}{b} \mathbb{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{40b\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \frac{8\gamma(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\right) \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h^{t+1}_{i} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \quad \leq \mathbb{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbb{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & \quad + (1 - \gamma\mu) \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbb{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu) \frac{4\gamma((2\omega + 1)p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g^{t}_{i} - h^{t}_{i}\right\|^{2}\right] \\ & \quad - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{50\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} \left(\frac{2(1 - b)^{2}L_{\sigma}^{2}}{B} + 2\hat{L}^{2}\right) \\ & \quad - \frac{10\gamma}{np_{\mathbf{a}}b} \left(\frac{2(1 - b)^{2}L_{\sigma}^{2}}{B} + 2\left(1 - \frac{p_{\mathbf{a}\mathbf{a}}}{p_{\mathbf{a}}}\right)\hat{L}^{2}\right)\right) \mathbb{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & \quad + \left(1 - \frac{b}{2}\right) \frac{2\gamma}{b} \mathbb{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(1 - \frac{b}{2}\right) \left(\frac{40b\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \frac{8\gamma(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\right) \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h^{t}_{i} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & \quad + \left(\frac{20b^{2}\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}} + \frac{4\gamma b}{np_{\mathbf{a}}} + \left(\frac{40b\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \frac{8\gamma(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}}\right) \frac{2b^{2}}{b} \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g^{t}_{i} - h^{t}_{i}\right\|^{2}\right] \\ & \quad + (1 - \gamma\mu) \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbb{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu) \frac{4\gamma((2\omega + 1)p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g^{t}_{i} - h^{t}_{i}\right\|^{2}\right] \\ & \quad - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{100\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \left(\frac{1 - p_{\mathbf{a}\mathbf{a}}}{b}\right)\hat{L}^{2}\right) \mathbb{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & \quad - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{100\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \left(\frac{1 - b}{p_{\mathbf{a}}}\right)\hat{L}^{2}\right) \mathbb{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & \quad + \left(\frac{1}{2} - \frac{b}{2}\right) \mathbb{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(\frac{1}{2} - \frac{b}{2}\right) \mathbb{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & \quad + \left(\frac{1}{2} - \frac{b}{2}\right) \mathbb{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(\frac{1}{2} - \frac{b}{2}\right) \mathbb{E}$$

808 Also, we can simplify the last term:

$$\left(\frac{40b\gamma\omega(2\omega+1)}{np_{\rm a}^2} + \frac{8\gamma\left(p_{\rm a} - p_{\rm aa}\right)}{np_{\rm a}^2}\right)\frac{2b^2}{p_{\rm a}}$$

$$= \frac{80b^3\gamma\omega(2\omega+1)}{np_{\rm a}^3} + \frac{16b^2\gamma\left(1 - \frac{p_{\rm aa}}{p_{\rm a}}\right)}{np_{\rm a}^2}$$

$$\leq \frac{80b^2\gamma\omega(2\omega+1)}{np_{\rm a}^2} + \frac{16b\gamma}{np_{\rm a}},$$

809 thus

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + (1 - \gamma\mu) \, \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + (1 - \gamma\mu) \, \frac{4\gamma((2\omega + 1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g^{t+1}_{i} - h^{t+1}_{i}\right\|^{2}\right] \\ & + \frac{2\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{40b\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h^{t+1}_{i} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + (1 - \gamma\mu) \, \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu) \, \frac{4\gamma((2\omega + 1)\,p_{\mathbf{a}} - p_{\mathbf{aa}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g^{t}_{i} - h^{t}_{i}\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{100\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} \left(\frac{(1 - b)^{2}L_{\sigma}^{2}}{B} + \hat{L}^{2}\right) - \frac{20\gamma}{np_{\mathbf{a}}b} \left(\frac{(1 - b)^{2}L_{\sigma}^{2}}{B} + \left(1 - \frac{p_{\mathbf{aa}}}{p_{\mathbf{a}}}\right)\hat{L}^{2}\right)\right) \mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \left(1 - \frac{b}{2}\right) \, \frac{2\gamma}{b} \mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(1 - \frac{b}{2}\right) \left(\frac{40b\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h^{t}_{i} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{100b^{2}\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \frac{20\gamma b}{np_{\mathbf{a}}}\right) \frac{\sigma^{2}}{B}. \end{split}$$

Using Lemma 4 and the assumption about γ , we get

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + (1 - \gamma\mu) \, \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + (1 - \gamma\mu) \, \frac{4\gamma((2\omega + 1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & \quad + \frac{2\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{40b\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \quad \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & \quad + (1 - \gamma\mu) \, \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu) \, \frac{4\gamma((2\omega + 1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & \quad + \left(1 - \frac{b}{2}\right) \frac{2\gamma}{b} \mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(1 - \frac{b}{2}\right) \left(\frac{40b\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & \quad + \left(\frac{100b^{2}\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \frac{20\gamma b}{np_{\mathbf{a}}}\right) \frac{\sigma^{2}}{B}. \end{split}$$

Note that $\gamma \leq \frac{b}{2\mu}$, thus $1 - \frac{b}{2} \leq 1 - \gamma \mu$ and

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + (1 - \gamma\mu) \, \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + (1 - \gamma\mu) \, \frac{4\gamma((2\omega + 1) \, p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & \quad + \frac{2\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \left(\frac{40b\gamma\omega(2\omega + 1)}{np_{\mathbf{a}}^{2}} + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\right) \mathbf{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \quad \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & \quad + (1 - \gamma\mu) \, \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu) \, \frac{4\gamma((2\omega + 1) \, p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}})}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n} \sum_{i=1}^{n} \left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \end{split}$$

$$\begin{split} & + \left(1 - \gamma \mu\right) \frac{2\gamma}{b} \mathbf{E} \left[\left\| h^t - \nabla f(x^t) \right\|^2 \right] + \left(1 - \gamma \mu\right) \left(\frac{40b\gamma \omega (2\omega + 1)}{np_{\rm a}^2} + \frac{8\gamma \left(p_{\rm a} - p_{\rm aa} \right)}{np_{\rm a}^2} \right) \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^n \left\| h_i^t - \nabla f_i(x^t) \right\|^2 \right] \\ & + \left(\frac{100b^2 \gamma \omega (2\omega + 1)}{np_{\rm a}^2} + \frac{20\gamma b}{np_{\rm a}} \right) \frac{\sigma^2}{B}. \end{split}$$

812 It is left to apply Lemma 11 with

$$\Psi^{t} = \frac{2(2\omega + 1)}{p_{a}} E\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{4((2\omega + 1)p_{a} - p_{aa})}{np_{a}^{2}} E\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] + \frac{2}{b} E\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(\frac{40b\omega(2\omega + 1)}{np_{a}^{2}} + \frac{8(p_{a} - p_{aa})}{np_{a}^{2}}\right) E\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right]$$

and $C=\left(\frac{100b^2\omega(2\omega+1)}{p_a^2}+\frac{20b}{p_a}\right)\frac{\sigma^2}{nB}$ to conclude the proof.

Corollary 5. Suppose that assumptions of Theorem 10 hold, batch size $B \leq \min\left\{\frac{\sigma}{p_a\sqrt{\mu\epsilon}n}, \frac{L_\sigma^2}{\hat{f}_2^2}\right\}$,

we take RandK compressors with $K=\Theta\left(rac{Bd\sqrt{\mu arepsilon n}}{\sigma}
ight)$. Then the communication complexity equals

$$\widetilde{\mathcal{O}}\left(\frac{d\sigma}{p_{\mathbf{a}}\sqrt{\mu\varepsilon n}} + \frac{dL_{\sigma}}{p_{\mathbf{a}}\mu\sqrt{n}}\right),\,$$

and the expected number of stochastic gradient calculations per node equals

$$\widetilde{\mathcal{O}}\left(rac{\sigma^2}{p_{
m a}\mu narepsilon}+rac{\sigma L_{\sigma}}{p_{
m a}n\mu^{3/2}\sqrt{arepsilon}}
ight).$$

817 *Proof.* In the view of Theorem 10, DASHA-PP have to run

$$\widetilde{\mathcal{O}}\left(\frac{\omega+1}{p_{\rm a}} + \frac{\omega}{p_{\rm a}}\sqrt{\frac{\sigma^2}{\mu n \varepsilon B}} + \frac{\sigma^2}{p_{\rm a}\mu n \varepsilon B} + \frac{L}{\mu} + \frac{\omega}{p_{\rm a}\mu \sqrt{n}}\left(\widehat{L} + \frac{L_{\sigma}}{\sqrt{B}}\right) + \frac{\sigma}{p_{\rm a}n\mu^{3/2}\sqrt{\varepsilon B}}\left(\widehat{L} + \frac{L_{\sigma}}{\sqrt{B}}\right)\right)$$

- sommunication rounds in the stochastic settings to get ε -solution. Note that $K = \mathcal{O}\left(\frac{d}{p_0\sqrt{n}}\right)$.
- Moreover, we can skip the initialization procedure and initialize h_i^0 and g_i^0 , for instance, with zeros
- because the initialization error is under a logarithm. Considering Theorem 6, the communication
- 821 complexity equals

$$\begin{split} &\widetilde{\mathcal{O}}\left(K\frac{\omega+1}{p_{\mathrm{a}}}+K\frac{\omega}{p_{\mathrm{a}}}\sqrt{\frac{\sigma^{2}}{\mu n \varepsilon B}}+K\frac{\sigma^{2}}{p_{\mathrm{a}}\mu n \varepsilon B}+K\frac{L}{\mu}+K\frac{\omega}{p_{\mathrm{a}}\mu \sqrt{n}}\left(\widehat{L}+\frac{L_{\sigma}}{\sqrt{B}}\right)+K\frac{\sigma}{p_{\mathrm{a}}n\mu^{3/2}\sqrt{\varepsilon B}}\left(\widehat{L}+\frac{L_{\sigma}}{\sqrt{B}}\right)\right)\\ &=\widetilde{\mathcal{O}}\left(K\frac{\omega+1}{p_{\mathrm{a}}}+K\frac{\omega}{p_{\mathrm{a}}}\sqrt{\frac{\sigma^{2}}{\mu n \varepsilon B}}+K\frac{\sigma^{2}}{p_{\mathrm{a}}\mu n \varepsilon B}+K\frac{L}{\mu}+K\frac{\omega}{p_{\mathrm{a}}\mu \sqrt{n}}\left(\widehat{L}+\frac{L_{\sigma}}{\sqrt{B}}\right)+K\frac{\sigma L_{\sigma}}{p_{\mathrm{a}}n\mu^{3/2}\sqrt{\varepsilon} B}\right)\\ &=\widetilde{\mathcal{O}}\left(\frac{d}{p_{\mathrm{a}}}+\frac{d}{p_{\mathrm{a}}}\sqrt{\frac{\sigma^{2}}{\mu n \varepsilon B}}+\frac{K\sigma^{2}}{p_{\mathrm{a}}\mu n \varepsilon B}+\frac{dL}{p_{\mathrm{a}}\mu \sqrt{n}}+\frac{d}{p_{\mathrm{a}}\mu \sqrt{n}}\left(\widehat{L}+\frac{L_{\sigma}}{\sqrt{B}}\right)+\frac{K\sigma L_{\sigma}}{p_{\mathrm{a}}n\mu^{3/2}\sqrt{\varepsilon} B}\right)\\ &=\widetilde{\mathcal{O}}\left(\frac{d}{p_{\mathrm{a}}}+\frac{d\sigma}{p_{\mathrm{a}}\sqrt{\mu n \varepsilon B}}+\frac{d\sigma}{p_{\mathrm{a}}\sqrt{\mu \varepsilon n}}+\frac{dL}{p_{\mathrm{a}}\mu \sqrt{n}}+\frac{d}{p_{\mathrm{a}}\mu \sqrt{n}}\left(\widehat{L}+\frac{L_{\sigma}}{\sqrt{B}}\right)+\frac{dL_{\sigma}}{p_{\mathrm{a}}\mu \sqrt{n}}\right)\\ &=\widetilde{\mathcal{O}}\left(\frac{d\sigma}{p_{\mathrm{a}}\sqrt{\mu \varepsilon n}}+\frac{dL_{\sigma}}{p_{\mathrm{a}}\mu \sqrt{n}}\right). \end{split}$$

The expected number of stochastic gradient calculations per node equals

$$\widetilde{\mathcal{O}}\left(B\frac{\omega+1}{p_{\rm a}}+B\frac{\omega}{p_{\rm a}}\sqrt{\frac{\sigma^2}{\mu n\varepsilon B}}+B\frac{\sigma^2}{p_{\rm a}\mu n\varepsilon B}+B\frac{L}{\mu}+B\frac{\omega}{p_{\rm a}\mu\sqrt{n}}\left(\widehat{L}+\frac{L_\sigma}{\sqrt{B}}\right)+B\frac{\sigma}{p_{\rm a}n\mu^{3/2}\sqrt{\varepsilon B}}\left(\widehat{L}+\frac{L_\sigma}{\sqrt{B}}\right)\right)$$

$$\begin{split} &=\widetilde{\mathcal{O}}\left(B\frac{\omega+1}{p_{\mathrm{a}}}+B\frac{\omega}{p_{\mathrm{a}}}\sqrt{\frac{\sigma^{2}}{\mu n \varepsilon B}}+B\frac{\sigma^{2}}{p_{\mathrm{a}}\mu n \varepsilon B}+B\frac{L}{\mu}+B\frac{\omega}{p_{\mathrm{a}}\mu \sqrt{n}}\left(\widehat{L}+\frac{L_{\sigma}}{\sqrt{B}}\right)+B\frac{\sigma}{p_{\mathrm{a}}n\mu^{3/2}\sqrt{\varepsilon B}}\left(\frac{L_{\sigma}}{\sqrt{B}}\right)\right)\\ &=\widetilde{\mathcal{O}}\left(\frac{Bd}{Kp_{\mathrm{a}}}+\frac{Bd}{Kp_{\mathrm{a}}}\sqrt{\frac{\sigma^{2}}{\mu n \varepsilon B}}+\frac{\sigma^{2}}{p_{\mathrm{a}}\mu n \varepsilon}+B\frac{L}{\mu}+\frac{Bd}{Kp_{\mathrm{a}}\mu \sqrt{n}}\left(\widehat{L}+\frac{L_{\sigma}}{\sqrt{B}}\right)+\frac{\sigma L_{\sigma}}{p_{\mathrm{a}}n\mu^{3/2}\sqrt{\varepsilon}}\right)\\ &=\widetilde{\mathcal{O}}\left(\frac{\sigma}{p_{\mathrm{a}}\sqrt{\mu \varepsilon n}}+\frac{\sigma^{2}}{p_{\mathrm{a}}\mu \varepsilon n\sqrt{B}}+\frac{\sigma^{2}}{p_{\mathrm{a}}\mu n \varepsilon}+\frac{\sigma L}{p_{\mathrm{a}}\mu^{3/2}\sqrt{\varepsilon n}}+\frac{\sigma}{p_{\mathrm{a}}\mu^{3/2}\sqrt{\varepsilon n}}\left(\widehat{L}+\frac{L_{\sigma}}{\sqrt{B}}\right)+\frac{\sigma L_{\sigma}}{p_{\mathrm{a}}n\mu^{3/2}\sqrt{\varepsilon}}\right)\\ &=\widetilde{\mathcal{O}}\left(\frac{\sigma^{2}}{p_{\mathrm{a}}\mu n \varepsilon}+\frac{\sigma L_{\sigma}}{p_{\mathrm{a}}n\mu^{3/2}\sqrt{\varepsilon}}\right). \end{split}$$

G Description of DASHA-PP-SYNC-MVR

By analogy to (Tyurin and Richtárik, 2023), we provide a "synchronized" version of the algorithm.
With a small probability, participating nodes calculate and send a mega batch without compression.
This helps us to resolve the suboptimality of DASHA-PP-MVR w.r.t. ω. Note that this suboptimality is not a problem. We show in Corollary 4 that DASHA-PP-MVR can have the optimal oracle complexity and SOTA communication complexity with the particular choices of parameters of the compressors.

Algorithm 8 DASHA-PP-SYNC-MVR

```
1: Input: starting point x^0 \in \mathbb{R}^d, stepsize \gamma > 0, momentum a \in (0,1], momentum b \in
         (0,1], probability p_{\text{mega}} \in (0,1], batch size B' and B, probability p_{\text{a}} \in (0,1] that a node is
 \begin{array}{l} \textit{participating}^{(\mathbf{a})}, \text{ number of iterations } T \geq 1. \\ 2 : \text{ Initialize } g_i^0, h_i^0 \text{ on the nodes and } g^0 = \frac{1}{n} \sum_{i=1}^n g_i^0 \text{ on the server} \\ 3 : \textbf{ for } t = 0, 1, \dots, T-1 \textbf{ do} \\ 4 : \quad x^{t+1} = x^t - \gamma g^t \end{array}
             c^{t+1} = \begin{cases} 1, \text{with probability } p_{\text{mega}}, \\ 0, \text{with probability } 1 - p_{\text{mega}} \end{cases}
               Broadcast x^{t+1}, x^t to all participating^{(a)} nodes
  6:
  7:
               for i = 1, \ldots, n in parallel do
  8:
                    if i<sup>th</sup> node is participating<sup>(a)</sup> then
                          if c^{t+1} = \hat{1} then
  9:
                                Generate i.i.d. samples \{\xi_{ik}^{t+1}\}_{k=1}^{B'} of size B' from \mathcal{D}_i. k_i^{t+1} = \frac{1}{B'} \sum_{k=1}^{B'} \nabla f_i(x^{t+1}; \xi_{ik}^{t+1}) - \frac{1}{B'} \sum_{k=1}^{B'} \nabla f_i(x^t; \xi_{ik}^{t+1}) - \frac{b}{p_{\text{mega}}} \left( h_i^t - \frac{1}{B'} \sum_{k=1}^{B'} \nabla f_i(x^t; \xi_{ik}^{t+1}) \right)
10:
11:
                                m_i^{t+1} = \frac{1}{p_a} k_i^{t+1} - \frac{a}{p_a} \left( g_i^t - h_i^t \right)
12:
13:
                               Generate i.i.d. samples \{\xi_{ij}^{t+1}\}_{j=1}^{B} of size B from \mathcal{D}_{i}. k_{i}^{t+1} = \frac{1}{B} \sum_{j=1}^{B} \nabla f_{i}(x^{t+1}; \xi_{ij}^{t+1}) - \frac{1}{B} \sum_{j=1}^{B} \nabla f_{i}(x^{t}; \xi_{ij}^{t+1}) \\ m_{i}^{t+1} = \mathcal{C}_{i} \left( \frac{1}{p_{a}} k_{i}^{t+1} - \frac{a}{p_{a}} (g_{i}^{t} - h_{i}^{t}) \right)
14:
15:
16:
                         end if h_i^{t+1} = h_i^t + \frac{1}{p_a} k_i^{t+1} g_i^{t+1} = g_i^t + m_i^{t+1} Send m_i^{t+1} to the server
17:
18:
20:
                    else h_i^{t+1} = h_i^t
m_i^{t+1} = 0
g_i^{t+1} = g_i^t
end if
21:
22:
23:
24:
25:
26:
               g^{t+1} = g^t + \frac{1}{n} \sum_{i=1}^n m_i^{t+1}
27:
29: Output: \hat{x}^T chosen uniformly at random from \{x^t\}_{k=0}^{T-1}
         (a): For the formal description see Section 2.2.
```

In the following theorem, we provide the convergence rate of DASHA-PP-SYNC-MVR.

Theorem 11. Suppose that Assumptions 1, 2, 3, 5, 6, 7 and 8 hold. Let us take $a = \frac{p_a}{2\omega+1}$, $b = \frac{p_{mega}p_a}{2-p_s}$, probability $p_{mega} \in (0,1]$, batch size $B' \geq B \geq 1$

$$\gamma \leq \left(L + \sqrt{\frac{8\left(2\omega + 1\right)\omega}{np_{\mathrm{a}}^2}\left(\widehat{L}^2 + \frac{L_\sigma^2}{B}\right) + \frac{16}{np_{\mathrm{mega}}p_{\mathrm{a}}^2}\left(\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)\widehat{L}^2 + \frac{L_\sigma^2}{B}\right)}\right)^{-1},$$

and $h_i^0 = g_i^0$ for all $i \in [n]$ in Algorithm 8. Then

$$\mathbb{E}\left[\left\|\nabla f(\widehat{x}^{T})\right\|^{2}\right] \leq \frac{1}{T} \left[\frac{2\Delta_{0}}{\gamma} + \frac{4}{p_{mega}p_{a}} \left\|h^{0} - \nabla f(x^{0})\right\|^{2} + \frac{4\left(1 - \frac{p_{aa}}{p_{a}}\right)}{np_{mega}p_{a}} \frac{1}{n} \sum_{i=1}^{n} \left\|h_{i}^{0} - \nabla f_{i}(x^{0})\right\|^{2} \right] + \frac{12\sigma^{2}}{nB'}.$$

First, we introduce the expected density of compressors (Gorbunov et al., 2021; Tyurin and Richtárik,

833 2023).

- Definition 12. The expected density of the compressor C_i is $\zeta_{C_i} := \sup_{x \in \mathbb{R}^d} \mathrm{E}[\|C_i(x)\|_0]$, where
- 835 $\|x\|_0$ is the number of nonzero components of $x \in \mathbb{R}^d$. Let $\zeta_{\mathcal{C}} = \max_{i \in [n]} \zeta_{\mathcal{C}_i}$
- Note that $\zeta_{\mathcal{C}}$ is finite and $\zeta_{\mathcal{C}} \leq d$.
- In the next corollary, we choose particular algorithm parameters to reveal the communication and
- 838 oracle complexity.

 $\begin{array}{l} \textbf{Corollary 6. Suppose that assumptions from Theorem 11 hold, probability $p_{mega} = \min\left\{\frac{\zeta_{\mathcal{C}}}{d}, \frac{n\varepsilon B}{\sigma^2}\right\}$, }\\ batch \textit{ size } B' = \Theta\left(\frac{\sigma^2}{n\varepsilon}\right), \textit{ and } h_i^0 = g_i^0 = \frac{1}{B_{\text{init}}} \sum_{k=1}^{B_{\text{init}}} \nabla f_i(x^0; \xi_{ik}^0) \textit{ for all } i \in [n], \textit{ initial batch size } \\ B_{\text{init}} = \Theta\left(\frac{B}{p_{\text{mega}}\sqrt{p_a}}\right) = \Theta\left(\max\left\{\frac{Bd}{\sqrt{p_a}\zeta_{\mathcal{C}}}, \frac{\sigma^2}{\sqrt{p_a}n\varepsilon}\right\}\right), \textit{ then } \text{DASHA-PP-SYNC-MVR } \textit{ needs} \\ \end{array}$

$$T := \mathcal{O}\left(\frac{\Delta_0}{\varepsilon}\left[L + \left(\frac{\omega}{p_{\rm a}\sqrt{n}} + \sqrt{\frac{d}{p_{\rm a}^2\zeta_{\mathcal{C}}n}}\right)\left(\widehat{L} + \frac{L_\sigma}{\sqrt{B}}\right) + \frac{\sigma}{p_{\rm a}\sqrt{\varepsilon}n}\left(\frac{\widehat{L}}{\sqrt{B}} + \frac{L_\sigma}{B}\right)\right] + \frac{\sigma^2}{\sqrt{p_{\rm a}}n\varepsilon B}\right).$$

- communication rounds to get an ε -solution, the expected communication complexity is equal to
- 840 $\mathcal{O}(d+\zeta_{\mathcal{C}}T)$, and the expected number of stochastic gradient calculations per node equals $\mathcal{O}(B_{\mathrm{init}}+$
- 841 BT), where ζ_C is the expected density from Definition 12.
- The main improvement of Corollary 6 over Corollary 3 is the size of the initial batch size B_{init} .
- However, Corollary 4 reveals that we can avoid regimes when DASHA-PP-MVR is suboptimal.
- We also provide a theorem under PŁ-condition (see Assumption 9).

Theorem 13. Suppose that Assumptions 1, 2, 3, 5, 6, 7, 8 and 9 hold. Let us take $a = \frac{p_a}{2\omega+1}$, $b = \frac{p_{mega}p_a}{2-p_a}$, probability $p_{mega} \in (0,1]$, batch size $B' \geq B \geq 1$,

$$\gamma \leq \min \left\{ \left(L + \sqrt{\frac{16\left(2\omega + 1\right)\omega}{np_{\mathrm{a}}^2} \left(\frac{L_{\sigma}^2}{B} + \widehat{L}^2\right) + \left(\frac{48L_{\sigma}^2}{np_{\mathrm{mega}}p_{\mathrm{a}}^2B} + \frac{24\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)\widehat{L}^2}{np_{\mathrm{mega}}p_{\mathrm{a}}^2} \right)} \right)^{-1}, \frac{a}{2\mu}, \frac{b}{2\mu} \right\},$$

and $h_i^0 = g_i^0$ for all $i \in [n]$ in Algorithm 8. Then

$$\begin{split} & \mathbf{E}\left[f(x^{T}) - f^{*}\right] \\ & \leq (1 - \gamma\mu)^{T} \left(\Delta_{0} + \frac{2\gamma}{b} \left\|h^{0} - \nabla f(x^{0})\right\|^{2} + \frac{8\gamma\left(p_{a} - p_{aa}\right)}{np_{a}^{2}p_{mega}} \frac{1}{n} \sum_{i=1}^{n} \left\|h_{i}^{0} - \nabla f_{i}(x^{0})\right\|^{2}\right) + \frac{20\sigma^{2}}{\mu nB'}. \end{split}$$

Let us provide bounds up to logarithmic factors and use $\widetilde{\mathcal{O}}\left(\cdot\right)$ notation.

Corollary 7. Suppose that assumptions from Theorem 13 hold, probability $p_{mega} = \min\left\{\frac{\zeta_{\mathcal{C}}}{d}, \frac{\mu n \varepsilon B}{\sigma^2}\right\}$, batch size $B' = \Theta\left(\frac{\sigma^2}{\mu n \varepsilon}\right)$ then DASHA-PP-SYNC-MVR needs

$$T := \widetilde{\mathcal{O}}\left(\frac{\omega+1}{p_{\mathrm{a}}} + \frac{d}{p_{\mathrm{a}}\zeta_{\mathcal{C}}} + \frac{\sigma^2}{p_{\mathrm{a}}\mu n\varepsilon B} + \frac{L}{\mu} + \frac{\omega}{p_{\mathrm{a}}\mu\sqrt{n}}\left(\frac{L_{\sigma}}{\sqrt{B}} + \widehat{L}\right) + \left(\frac{\sqrt{d}}{p_{\mathrm{a}}\mu\sqrt{\zeta_{\mathcal{C}}n}} + \frac{\sigma}{p_{\mathrm{a}}n\mu^{3/2}\sqrt{\varepsilon B}}\right)\left(\frac{L_{\sigma}}{\sqrt{B}} + \widehat{L}\right)\right).$$

communication rounds to get an ε -solution, the expected communication complexity is equal to $\widetilde{\mathcal{O}}(\zeta_{\mathcal{C}}T)$, and the expected number of stochastic gradient calculations per node equals $\widetilde{\mathcal{O}}(BT)$, where $\zeta_{\mathcal{C}}$ is the expected density from Definition 12.

The proof of this corollary almost repeats the proof of Corollary 6. Note that we can skip the initialization procedure and initialize h_i^0 and g_i^0 , for instance, with zeros because the initialization error is under a logarithm.

Let us assume that $\frac{d}{\zeta_C} = \Theta(\omega)$ (holds for the RandK compressor), then the convergence rate of DASHA-PP-SYNC-MVR is

$$\widetilde{\mathcal{O}}\left(\frac{\omega+1}{p_{\rm a}} + \frac{\sigma^2}{p_{\rm a}\mu n\varepsilon B} + \frac{L}{\mu} + \frac{\omega}{p_{\rm a}\mu\sqrt{n}}\left(\frac{L_{\sigma}}{\sqrt{B}} + \widehat{L}\right) + \frac{\sigma}{p_{\rm a}n\mu^{3/2}\sqrt{\varepsilon B}}\left(\frac{L_{\sigma}}{\sqrt{B}} + \widehat{L}\right)\right). \tag{31}$$

Comparing (31) with the rate of DASHA-PP-MVR (28), one can see that DASHA-PP-SYNC-MVR improves the suboptimal term \mathcal{P}_2 from (28). However, Corollary 5 reveals that we can escape these suboptimal regimes by choosing the parameter K of RandK compressors in a particular way.

G.1 Proof for DASHA-PP-SYNC-MVR

In this section, we provide the proof of the convergence rate for DASHA-PP-SYNC-MVR. There are four different sources of randomness in Algorithm 8: the first one from random samples ξ^{t+1} , the second one from compressors $\{\mathcal{C}_i\}_{i=1}^n$, the third one from availability of nodes, and the fourth one from c^{t+1} . We define $\mathbf{E}_k\left[\cdot\right]$, $\mathbf{E}_{\mathcal{C}}\left[\cdot\right]$, $\mathbf{E}_{p_a}\left[\cdot\right]$ and $\mathbf{E}_{p_{\text{mega}}}\left[\cdot\right]$ to be conditional expectations w.r.t. ξ^{t+1} , $\{\mathcal{C}_i\}_{i=1}^n$, availability, and c^{t+1} , accordingly, conditioned on all previous randomness. Moreover, we define $\mathbf{E}_{t+1}\left[\cdot\right]$ to be a conditional expectation w.r.t. all randomness in iteration t+1 conditioned on all previous randomness.

Let us denote

858

$$\begin{split} k_{i,1}^{t+1} &:= \frac{1}{B'} \sum_{k=1}^{B'} \nabla f_i(x^{t+1}; \xi_{ik}^{t+1}) - \frac{1}{B'} \sum_{k=1}^{B'} \nabla f_i(x^t; \xi_{ik}^{t+1}) - \frac{b}{p_{\text{mega}}} \left(h_i^t - \frac{1}{B'} \sum_{k=1}^{B'} \nabla f_i(x^t; \xi_{ik}^{t+1}) \right) \\ k_{i,2}^{t+1} &:= \frac{1}{B} \sum_{j=1}^{B} \nabla f_i(x^{t+1}; \xi_{ij}^{t+1}) - \frac{1}{B} \sum_{j=1}^{B} \nabla f_i(x^t; \xi_{ij}^{t+1}), \\ h_{i,1}^{t+1} &:= \begin{cases} h_i^t + \frac{1}{p_a} k_{i,1}^{t+1}, & i^{\text{th node is } participating,} \\ h_i^t, & \text{otherwise,} \end{cases} \\ h_{i,2}^{t+1} &:= \begin{cases} h_i^t + \frac{1}{p_a} k_{i,2}^{t+1}, & i^{\text{th node is } participating,} \\ h_i^t, & \text{otherwise,} \end{cases} \\ g_{i,1}^{t+1} &:= \begin{cases} g_i^t + \frac{1}{p_a} k_{i,1}^{t+1} - \frac{a}{p_a} \left(g_i^t - h_i^t \right), & i^{\text{th node is } participating,} \\ g_i^t, & \text{otherwise,} \end{cases} \\ g_{i,2}^{t+1} &:= \begin{cases} g_i^t + \mathcal{C}_i \left(\frac{1}{p_a} k_{i,2}^{t+1} - \frac{a}{p_a} \left(g_i^t - h_i^t \right) \right), & i^{\text{th node is } participating,} \\ g_i^t, & \text{otherwise,} \end{cases} \\ g_{i,2}^t &:= \begin{cases} g_i^t + \mathcal{C}_i \left(\frac{1}{p_a} k_{i,2}^{t+1} - \frac{a}{p_a} \left(g_i^t - h_i^t \right) \right), & i^{\text{th node is } participating,} \\ g_i^t, & \text{otherwise,} \end{cases} \\ \end{split}$$

867 $h_1^{t+1}:=\frac{1}{n}\sum_{i=1}^n h_{i,1}^{t+1}, h_2^{t+1}:=\frac{1}{n}\sum_{i=1}^n h_{i,2}^{t+1}, g_1^{t+1}:=\frac{1}{n}\sum_{i=1}^n g_{i,1}^{t+1}, \text{ and } g_2^{t+1}:=\frac{1}{n}\sum_{i=1}^n g_{i,2}^{t+1}.$ 868 Note, that

$$h^{t+1} = \begin{cases} h_1^{t+1}, & c^{t+1} = 1, \\ h_2^{t+1}, & c^{t+1} = 0, \end{cases}$$

869 and

$$g^{t+1} = \begin{cases} g_1^{t+1}, & c^{t+1} = 1, \\ g_2^{t+1}, & c^{t+1} = 0 \end{cases}$$

First, we will prove two lemmas.

Lemma 13. Suppose that Assumptions 3, 5, 7 and 8 hold and let us consider sequences $\{g_i^{t+1}\}_{i=1}^n$ and $\{h_i^{t+1}\}_{i=1}^n$ from Algorithm 8, then

$$\begin{split} & \mathbf{E}_{\mathcal{C}}\left[\mathbf{E}_{p_{a}}\left[\mathbf{E}_{p_{mega}}\left[\left\|g^{t+1}-h^{t+1}\right\|^{2}\right]\right]\right] \\ & \leq \frac{2\left(1-p_{mega}\right)\omega}{n^{2}p_{a}}\sum_{i=1}^{n}\left\|k_{i,2}^{t+1}\right\|^{2} + \left(\frac{\left(p_{a}-p_{aa}\right)a^{2}}{n^{2}p_{a}^{2}} + \frac{2\left(1-p_{mega}\right)a^{2}\omega}{n^{2}p_{a}}\right)\sum_{i=1}^{n}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} \\ & + (1-a)^{2}\left\|g^{t}-h^{t}\right\|^{2}, \end{split}$$

873 *and*

$$\begin{split} &\mathbf{E}_{\mathcal{C}}\left[\mathbf{E}_{p_{\mathbf{a}}}\left[\mathbf{E}_{p_{\textit{mega}}}\left[\left\|\boldsymbol{g}_{i}^{t+1}-\boldsymbol{h}_{i}^{t+1}\right\|^{2}\right]\right]\right] \\ &\leq \frac{2\left(1-p_{\textit{mega}}\right)\omega}{p_{\mathbf{a}}}\left\|\boldsymbol{k}_{i,2}^{t+1}\right\|^{2}+\left(\frac{(1-p_{\mathbf{a}})a^{2}}{p_{\mathbf{a}}}+\frac{2\left(1-p_{\textit{mega}}\right)a^{2}\omega}{p_{\mathbf{a}}}\right)\left\|\boldsymbol{g}_{i}^{t}-\boldsymbol{h}_{i}^{t}\right\|^{2} \\ &+(1-a)^{2}\left\|\boldsymbol{g}_{i}^{t}-\boldsymbol{h}_{i}^{t}\right\|^{2}, \quad \forall i \in [n]. \end{split}$$

874 *Proof.* First, we get the bound for $\mathbf{E}_{t+1}\left[\left\|g^{t+1}-h^{t+1}\right\|^2\right]$:

$$\begin{split} & \mathbf{E}_{\mathcal{C}}\left[\mathbf{E}_{p_{\text{a}}}\left[\left\|g^{t+1}-h^{t+1}\right\|^{2}\right]\right]\right] \\ & = p_{\text{mega}}\mathbf{E}_{p_{\text{a}}}\left[\left\|g^{t+1}_{1}-h^{t+1}_{1}\right\|^{2}\right] + (1-p_{\text{mega}})\,\mathbf{E}_{\mathcal{C}}\left[\mathbf{E}_{p_{\text{a}}}\left[\left\|g^{t+1}_{2}-h^{t+1}_{2}\right\|^{2}\right]\right]. \end{split}$$

875 Using

$$\mathbf{E}_{p_{\mathbf{a}}}\left[g_{i,1}^{t+1}-h_{i,1}^{t+1}\right] = g_{i}^{t}+k_{i,1}^{t+1}-a\left(g_{i}^{t}-h_{i}^{t}\right)-h_{i}^{t}-k_{i,1}^{t+1} = (1-a)\left(g_{i}^{t}-h_{i}^{t}\right)$$

876 and

$$\mathbf{E}_{\mathcal{C}}\left[\mathbf{E}_{p_{a}}\left[g_{i,2}^{t+1}-h_{i,2}^{t+1}\right]\right] = g_{i}^{t}+k_{i,2}^{t+1}-a\left(g_{i}^{t}-h_{i}^{t}\right)-h_{i}^{t}-k_{i,2}^{t+1} = (1-a)\left(g_{i}^{t}-h_{i}^{t}\right)$$

877 we have

$$\begin{split} & \mathbf{E}_{\mathcal{C}} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\| g^{t+1} - h^{t+1} \|^{2} \right] \right] \\ & \stackrel{\text{(13)}}{=} p_{\text{mega}} \mathbf{E}_{p_{\mathbf{a}}} \left[\| g_{1}^{t+1} - h_{1}^{t+1} - \mathbf{E}_{p_{\mathbf{a}}} \left[g_{1}^{t+1} - h_{1}^{t+1} \right] \|^{2} \right] \\ & + (1 - p_{\text{mega}}) \, \mathbf{E}_{\mathcal{C}} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\| g_{2}^{t+1} - h_{2}^{t+1} - \mathbf{E}_{p_{\mathbf{a}}} \left[g_{2}^{t+1} - h_{2}^{t+1} \right] \|^{2} \right] \right] \\ & + (1 - a)^{2} \, \| g^{t} - h^{t} \|^{2} \, . \end{split}$$

We can use Lemma 1 two times with i) $r_i = g_i^t - h_i^t$ and $s_i = -a (g_i^t - h_i^t)$ and ii) $r_i = g_i^t - h_i^t$ and

втэ
$$s_i=p_{\mathrm{a}}\mathcal{C}_i\left(rac{1}{p_{\mathrm{a}}}k_{i,2}^{t+1}-rac{a}{p_{\mathrm{a}}}\left(g_i^t-h_i^t
ight)
ight)-k_{i,2}^{t+1}$$
, to obtain

$$\begin{split} & \mathbf{E}_{\mathcal{C}} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\mathbf{E}_{p_{\text{mega}}} \left[\left\| g^{t+1} - h^{t+1} \right\|^{2} \right] \right] \right] \\ & \leq \frac{p_{\text{mega}} a^{2} \left(p_{\mathbf{a}} - p_{\mathbf{aa}} \right)}{n^{2} p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \\ & + \left(1 - p_{\text{mega}} \right) \left(\frac{1}{n^{2} p_{\mathbf{a}}} \sum_{i=1}^{n} \mathbf{E}_{\mathcal{C}} \left[\left\| p_{\mathbf{a}} \mathcal{C}_{i} \left(\frac{1}{p_{\mathbf{a}}} k_{i,2}^{t+1} - \frac{a}{p_{\mathbf{a}}} \left(g_{i}^{t} - h_{i}^{t} \right) \right) - \left(k_{i,2}^{t+1} - a \left(g_{i}^{t} - h_{i}^{t} \right) \right) \right\|^{2} \right] \right) \\ & + \left(1 - p_{\text{mega}} \right) \left(\frac{a^{2} \left(p_{\mathbf{a}} - p_{\mathbf{aa}} \right)}{n^{2} p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \right) \\ & + \left(1 - a \right)^{2} \left\| g^{t} - h^{t} \right\|^{2} \\ & = \frac{a^{2} \left(p_{\mathbf{a}} - p_{\mathbf{aa}} \right)}{n^{2} p_{\mathbf{a}}^{2}} \sum_{i=1}^{n} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \end{split}$$

$$\begin{split} &+ (1 - p_{\text{mega}}) \left(\frac{p_{\text{a}}}{n^2} \sum_{i=1}^n \mathbf{E}_{\mathcal{C}} \left[\left\| \mathcal{C}_i \left(\frac{1}{p_{\text{a}}} k_{i,2}^{t+1} - \frac{a}{p_{\text{a}}} \left(g_i^t - h_i^t \right) \right) - \left(\frac{1}{p_{\text{a}}} k_{i,2}^{t+1} - \frac{a}{p_{\text{a}}} \left(g_i^t - h_i^t \right) \right) \right\|^2 \right] \right) \\ &+ (1 - a)^2 \left\| g^t - h^t \right\|^2 \\ &\leq \frac{a^2 \left(p_{\text{a}} - p_{\text{aa}} \right)}{n^2 p_{\text{a}}^2} \sum_{i=1}^n \left\| g_i^t - h_i^t \right\|^2 \\ &+ \frac{\left(1 - p_{\text{mega}} \right) p_{\text{a}} \omega}{n^2} \sum_{i=1}^n \left\| \frac{1}{p_{\text{a}}} k_{i,2}^{t+1} - \frac{a}{p_{\text{a}}} \left(g_i^t - h_i^t \right) \right\|^2 \\ &+ (1 - a)^2 \left\| g^t - h^t \right\|^2 \\ &= \frac{a^2 \left(p_{\text{a}} - p_{\text{aa}} \right)}{n^2 p_{\text{a}}^2} \sum_{i=1}^n \left\| g_i^t - h_i^t \right\|^2 \\ &+ \frac{\left(1 - p_{\text{mega}} \right) \omega}{n^2 p_{\text{a}}} \sum_{i=1}^n \left\| k_{i,2}^{t+1} - a \left(g_i^t - h_i^t \right) \right\|^2 \\ &+ (1 - a)^2 \left\| g^t - h^t \right\|^2. \end{split}$$

880 In the last inequality, we use Assumption 7. Next, using (12), we have

$$\begin{split} & \mathbf{E}_{\mathcal{C}}\left[\mathbf{E}_{p_{\text{a}}}\left[\mathbf{E}_{p_{\text{mega}}}\left[\left\|g^{t+1}-h^{t+1}\right\|^{2}\right]\right]\right] \\ & \leq \frac{2\left(1-p_{\text{mega}}\right)\omega}{n^{2}p_{\text{a}}}\sum_{i=1}^{n}\left\|k_{i,2}^{t+1}\right\|^{2} + \left(\frac{\left(p_{\text{a}}-p_{\text{aa}}\right)a^{2}}{n^{2}p_{\text{a}}^{2}} + \frac{2\left(1-p_{\text{mega}}\right)\omega a^{2}}{n^{2}p_{\text{a}}}\right)\sum_{i=1}^{n}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} \\ & + (1-a)^{2}\left\|g^{t}-h^{t}\right\|^{2}. \end{split}$$

The second inequality can be proved almost in the same way:

$$\begin{split} & \mathbf{E}_{\mathcal{C}}\left[\mathbf{E}_{p_{a}}\left[\mathbf{E}_{p_{mega}}\left[\left\|g_{i,1}^{t+1}-h_{i,1}^{t+1}\right\|^{2}\right]\right]\right] \\ & = p_{mega}\mathbf{E}_{p_{a}}\left[\left\|g_{i,1}^{t+1}-h_{i,1}^{t+1}\right\|^{2}\right] + (1-p_{mega})\,\mathbf{E}_{\mathcal{C}}\left[\mathbf{E}_{p_{a}}\left[\left\|g_{i,2}^{t+1}-h_{i,2}^{t+1}\right\|^{2}\right]\right] \\ & = \frac{1}{2}p_{mega}\mathbf{E}_{p_{a}}\left[\left\|g_{i,1}^{t+1}-h_{i,1}^{t+1}-(1-a)\left(g_{i}^{t}-h_{i}^{t}\right)\right\|^{2}\right] + (1-p_{mega})\,\mathbf{E}_{\mathcal{C}}\left[\mathbf{E}_{p_{a}}\left[\left\|g_{i,2}^{t+1}-h_{i,2}^{t+1}\right\|^{2}\right]\right] \\ & + p_{mega}(1-a)^{2}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} \\ & = \frac{p_{mega}(1-p_{a})a^{2}}{p_{a}}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} + (1-p_{mega})\,\mathbf{E}_{\mathcal{C}}\left[\mathbf{E}_{p_{a}}\left[\left\|g_{i,2}^{t+1}-h_{i,2}^{t+1}\right\|^{2}\right]\right] \\ & + p_{mega}(1-a)^{2}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} \\ & = \frac{p_{mega}(1-p_{a})a^{2}}{p_{a}}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} + (1-p_{mega})\,\mathbf{E}_{\mathcal{C}}\left[\mathbf{E}_{p_{a}}\left[\left\|g_{i,2}^{t+1}-h_{i,2}^{t+1}-(1-a)\left(g_{i}^{t}-h_{i}^{t}\right)\right\|^{2}\right]\right] \\ & + (1-a)^{2}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} \\ & = \frac{p_{mega}(1-p_{a})a^{2}}{p_{a}}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} \\ & + (1-p_{mega})\left(1-p_{a}\right)\left\|g_{i}^{t}-h_{i}^{t}-(1-a)\left(g_{i}^{t}-h_{i}^{t}\right)\right) - \left(h_{i}^{t}+\frac{1}{p_{a}}k_{i,2}^{t+1}\right) - (1-a)\left(g_{i}^{t}-h_{i}^{t}\right)\right\|^{2} \\ & + (1-p_{mega})\left(1-p_{a}\right)\left\|g_{i}^{t}-h_{i}^{t}-(1-a)\left(g_{i}^{t}-h_{i}^{t}\right)\right\|^{2} \\ & = \frac{p_{mega}(1-p_{a})a^{2}}{p_{a}}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2} \end{aligned}$$

$$\begin{split} &+ (1-p_{\text{mega}}) \left(1-p_{\text{a}}\right) a^{2} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \\ &+ (1-a)^{2} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \\ &+ (1-a)^{2} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \\ &+ \left(1-p_{\text{mega}}\right) a^{2} + \frac{\left(1-p_{\text{mega}}\right) \left(1-p_{\text{a}}\right) a^{2}}{p_{\text{a}}} \right) \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \\ &+ \left(1-p_{\text{mega}}\right) p_{\text{a}} \mathbf{E}_{\mathcal{C}} \left[\left\| \mathcal{C}_{i} \left(\frac{1}{p_{\text{a}}} k_{i,2}^{t+1} - \frac{a}{p_{\text{a}}} \left(g_{i}^{t} - h_{i}^{t} \right) \right) - \left(\frac{1}{p_{\text{a}}} k_{i,2}^{t+1} - \frac{a}{p_{\text{a}}} \left(g_{i}^{t} - h_{i}^{t} \right) \right) \right\|^{2} \right] \\ &+ \left(1-a\right)^{2} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \\ &+ \left(1-p_{\text{mega}}\right) p_{\text{a}} \mathbf{E}_{\mathcal{C}} \left[\left\| \mathcal{C}_{i} \left(\frac{1}{p_{\text{a}}} k_{i,2}^{t+1} - \frac{a}{p_{\text{a}}} \left(g_{i}^{t} - h_{i}^{t} \right) \right) - \left(\frac{1}{p_{\text{a}}} k_{i,2}^{t+1} - \frac{a}{p_{\text{a}}} \left(g_{i}^{t} - h_{i}^{t} \right) \right) \right\|^{2} \right] \\ &+ \left(1-a\right)^{2} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \\ &\leq \frac{\left(1-p_{\text{a}}\right) a^{2}}{p_{\text{a}}} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \\ &+ \frac{\left(1-p_{\text{mega}}\right) \omega}{p_{\text{a}}} \left\| k_{i,2}^{t+1} - a \left(g_{i}^{t} - h_{i}^{t} \right) \right\|^{2} \\ &+ \left(1-a\right)^{2} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \\ &\leq \frac{2\left(1-p_{\text{mega}}\right) \omega}{p_{\text{a}}} \left\| k_{i,2}^{t+1} \right\|^{2} + \left(\frac{\left(1-p_{\text{a}}\right) a^{2}}{p_{\text{a}}} + \frac{2\left(1-p_{\text{mega}}\right) a^{2} \omega}{p_{\text{a}}} \right) \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \\ &+ \left(1-a\right)^{2} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2}. \end{split}$$

Lemma 14. Suppose that Assumptions 3, 5, 6 and 8 hold and let us consider sequence $\{h_i^{t+1}\}_{i=1}^n$ from Algorithm 8, then

$$\begin{split} & \mathbf{E}_{k} \left[\mathbf{E}_{p_{a}} \left[\mathbf{E}_{p_{mega}} \left[\left\| h^{t+1} - \nabla f(\boldsymbol{x}^{t+1}) \right\|^{2} \right] \right] \right] \\ & \leq \frac{2b^{2}\sigma^{2}}{np_{mega}p_{\mathbf{a}}B'} + \left(\frac{2p_{mega}L_{\sigma}^{2}}{np_{\mathbf{a}}B'} \left(1 - \frac{b}{p_{mega}} \right)^{2} + \frac{(1 - p_{mega})L_{\sigma}^{2}}{np_{\mathbf{a}}B} + \frac{2\left(p_{\mathbf{a}} - p_{\mathbf{aa}} \right) \hat{L}^{2}}{np_{\mathbf{a}}^{2}} \right) \left\| \boldsymbol{x}^{t+1} - \boldsymbol{x}^{t} \right\|^{2} \\ & + \frac{2\left(p_{\mathbf{a}} - p_{\mathbf{aa}} \right)b^{2}}{n^{2}p_{\mathbf{a}}^{2}p_{mega}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(\boldsymbol{x}^{t}) \right\|^{2} + \left(p_{mega} \left(1 - \frac{b}{p_{mega}} \right)^{2} + (1 - p_{mega}) \right) \left\| h^{t} - \nabla f(\boldsymbol{x}^{t}) \right\|^{2}, \end{split}$$

$$\begin{split} & \mathbf{E}_{k} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\mathbf{E}_{p_{\textit{mega}}} \left[\left\| h_{i}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \right] \right] \\ & \leq \frac{2b^{2}\sigma^{2}}{p_{\mathbf{a}}p_{\textit{mega}}B'} + \left(\frac{2p_{\textit{mega}}L_{\sigma}^{2}}{p_{\mathbf{a}}B'} \left(1 - \frac{b}{p_{\textit{mega}}} \right)^{2} + \frac{(1 - p_{\textit{mega}})L_{\sigma}^{2}}{p_{\mathbf{a}}B} + \frac{2(1 - p_{\mathbf{a}})L_{i}^{2}}{p_{\mathbf{a}}} \right) \left\| x^{t+1} - x^{t} \right\|^{2} \\ & + \frac{2\left(1 - p_{\mathbf{a}} \right)b^{2}}{p_{\textit{mega}}p_{\mathbf{a}}} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} + \left(p_{\textit{mega}} \left(1 - \frac{b}{p_{\textit{mega}}} \right)^{2} + (1 - p_{\textit{mega}}) \right) \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}, \quad \forall i \in [n], \end{split}$$

and

$$\mathbb{E}_{k}\left[\left\|k_{i,2}^{t+1}\right\|^{2}\right] \leq \left(\frac{L_{\sigma}^{2}}{B} + L_{i}^{2}\right) \left\|x^{t+1} - x^{t}\right\|^{2}, \quad \forall i \in [n],$$

Proof. First, we prove the bound for $\mathbf{E}_k\left[\mathbf{E}_{p_{\mathrm{a}}}\left[\mathbf{E}_{p_{\mathrm{mega}}}\left[\left\|h^{t+1}-\nabla f(x^{t+1})\right\|^2\right]\right]\right]$. Using $\mathbf{E}_k\left[\mathbf{E}_{p_{\mathrm{a}}}\left[h_{i,1}^{t+1}\right]\right]$

$$= h_i^t + \mathbf{E}_k \left[\frac{1}{B'} \sum_{k=1}^{B'} \nabla f_i(x^{t+1}; \xi_{ik}^{t+1}) - \frac{1}{B'} \sum_{k=1}^{B'} \nabla f_i(x^t; \xi_{ik}^{t+1}) - \frac{b}{p_{\text{mega}}} \left(h_i^t - \frac{1}{B'} \sum_{k=1}^{B'} \nabla f_i(x^t; \xi_{ik}^{t+1}) \right) \right]$$

$$= h_i^t + \nabla f_i(x^{t+1}) - \nabla f_i(x^t) - \frac{b}{p_{\text{mega}}} \left(h_i^t - \nabla f_i(x^t) \right)$$

888 and

$$\begin{split} & \mathbf{E}_{k} \left[\mathbf{E}_{p_{a}} \left[h_{i,2}^{t+1} \right] \right] \\ & = h_{i}^{t} + \mathbf{E}_{k} \left[\frac{1}{B} \sum_{j=1}^{B} \nabla f_{i}(x^{t+1}; \xi_{ij}^{t+1}) - \frac{1}{B} \sum_{j=1}^{B} \nabla f_{i}(x^{t}; \xi_{ij}^{t+1}) \right] \\ & = h_{i}^{t} + \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}), \end{split}$$

889 we have

$$\begin{split} & \mathbf{E}_{k} \left[\mathbf{E}_{p_{\text{a}}} \left[\mathbf{E}_{p_{\text{mega}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \right] \\ &= p_{\text{mega}} \mathbf{E}_{k} \left[\mathbf{E}_{p_{\text{a}}} \left[\left\| h^{t+1}_{1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] + (1 - p_{\text{mega}}) \mathbf{E}_{k} \left[\mathbf{E}_{p_{\text{a}}} \left[\left\| h^{t+1}_{2} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \\ &\stackrel{\text{(13)}}{=} p_{\text{mega}} \mathbf{E}_{k} \left[\mathbf{E}_{p_{\text{a}}} \left[\left\| h^{t+1}_{1} - \mathbf{E}_{k} \left[\mathbf{E}_{p_{\text{a}}} \left[h^{t+1}_{1} \right] \right] \right\|^{2} \right] \right] + (1 - p_{\text{mega}}) \mathbf{E}_{k} \left[\mathbf{E}_{p_{\text{a}}} \left[\left\| h^{t+1}_{2} - \mathbf{E}_{k} \left[\mathbf{E}_{p_{\text{a}}} \left[h^{t+1}_{2} \right] \right] \right\|^{2} \right] \right] \\ &+ \left(p_{\text{mega}} \left(1 - \frac{b}{p_{\text{mega}}} \right)^{2} + (1 - p_{\text{mega}}) \right) \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}. \end{split}$$

- We can use Lemma 1 two times with i) $r_i=h_i^t$ and $s_i=k_{i,1}^{t+1}$ and ii) $r_i=h_i^t$ and $s_i=k_{i,2}^{t+1}$, to
- 891 obtain

$$\begin{split} & \mathbb{E}_{k} \left[\mathbb{E}_{p_{a}} \left[\mathbb{E}_{p_{\text{mega}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \right] \\ & \leq p_{\text{mega}} \left(\frac{1}{n^{2}p_{a}} \sum_{i=1}^{n} \mathbb{E}_{k} \left[\left\| k_{i,1}^{t+1} - \mathbb{E}_{k} \left[k_{i,1}^{t+1} \right] \right\|^{2} \right] + \frac{p_{a} - p_{\text{aa}}}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - \frac{b}{p_{\text{mega}}} \left(h_{i}^{t} - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right) \\ & + \left(1 - p_{\text{mega}} \right) \left(\frac{1}{n^{2}p_{a}} \sum_{i=1}^{n} \mathbb{E}_{k} \left[\left\| k_{i,2}^{t+1} - \mathbb{E}_{k} \left[k_{i,2}^{t+1} \right] \right\|^{2} \right] + \frac{p_{a} - p_{\text{aa}}}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} \right) \\ & + \left(p_{\text{mega}} \left(1 - \frac{b}{p_{\text{mega}}} \right)^{2} + \left(1 - p_{\text{mega}} \right) \right) \left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \\ & \leq \frac{p_{\text{mega}}}{n^{2}p_{a}} \sum_{i=1}^{n} \mathbb{E}_{k} \left[\left\| k_{i,1}^{t+1} - \mathbb{E}_{k} \left[k_{i,2}^{t+1} \right] \right\|^{2} \right] \\ & + \frac{1 - p_{\text{mega}}}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \mathbb{E}_{k} \left[\left\| k_{i,2}^{t+1} - \mathbb{E}_{k} \left[k_{i,2}^{t+1} \right] \right\|^{2} \right] \\ & + \frac{2 \left(p_{a} - p_{\text{aa}} \right) b^{2}}{n^{2}p_{a}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} \\ & + \frac{2 \left(p_{a} - p_{\text{aa}} \right) b^{2}}{n^{2}p_{a}^{2}p_{\text{mega}}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} + \left(p_{\text{mega}} \left(1 - \frac{b}{p_{\text{mega}}} \right)^{2} + \left(1 - p_{\text{mega}} \right) \right) \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}. \end{split}$$

Let us consider $\mathbf{E}_k \left[\left\| k_{i,1}^{t+1} - \mathbf{E}_k \left[k_{i,1}^{t+1} \right] \right\|^2 \right]$.

$$\begin{split} & \mathbf{E}_{k} \left[\left\| k_{i,1}^{t+1} - \mathbf{E}_{k} \left[k_{i,1}^{t+1} \right] \right\|^{2} \right] \\ & = \mathbf{E}_{k} \left[\left\| \frac{1}{B'} \sum_{k=1}^{B'} \nabla f_{i}(x^{t+1}; \xi_{ik}^{t+1}) - \frac{1}{B'} \sum_{k=1}^{B'} \nabla f_{i}(x^{t}; \xi_{ik}^{t+1}) - \frac{b}{p_{\text{mega}}} \left(h_{i}^{t} - \frac{1}{B'} \sum_{k=1}^{B'} \nabla f_{i}(x^{t}; \xi_{ik}^{t+1}) \right) \right] \end{split}$$

$$-\left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) - \frac{b}{p_{\text{mega}}} \left(h_{i}^{t} - \nabla f_{i}(x^{t})\right)\right) \right\|^{2} \right]$$

$$= \mathbf{E}_{k} \left[\left\| \frac{1}{B'} \sum_{k=1}^{B'} \nabla f_{i}(x^{t+1}; \xi_{ik}^{t+1}) - \frac{1}{B'} \sum_{k=1}^{B'} \nabla f_{i}(x^{t}; \xi_{ik}^{t+1}) + \frac{b}{p_{\text{mega}}} \left(\frac{1}{B'} \sum_{k=1}^{B'} \nabla f_{i}(x^{t}; \xi_{ik}^{t+1}) \right) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) + \frac{b}{p_{\text{mega}}} \left(\nabla f_{i}(x^{t}) \right) \right) \right\|^{2} \right]$$

$$= \frac{1}{B'^{2}} \sum_{k=1}^{B'} \mathbf{E}_{k} \left[\left\| \frac{b}{p_{\text{mega}}} \left(\nabla f_{i}(x^{t+1}; \xi_{ik}^{t+1}) - \nabla f_{i}(x^{t}) \right) + \left(1 - \frac{b}{p_{\text{mega}}} \right) \left(\nabla f_{i}(x^{t+1}; \xi_{ik}^{t+1}) - \nabla f_{i}(x^{t}; \xi_{ik}^{t+1}) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right) \right\|^{2} \right],$$

where we used independence of the mini-batch samples. Using (12), we get

$$\begin{split} & \mathbf{E}_{k} \left[\left\| k_{i,1}^{t+1} - \mathbf{E}_{k} \left[k_{i,1}^{t+1} \right] \right\|^{2} \right] \\ & \leq \frac{2b^{2}}{B'^{2} p_{\text{mega}}^{2}} \sum_{k=1}^{B'} \mathbf{E}_{k} \left[\left\| \nabla f_{i}(x^{t+1}; \xi_{ik}^{t+1}) - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \\ & + \frac{2}{B'^{2}} \left(1 - \frac{b}{p_{\text{mega}}} \right)^{2} \sum_{k=1}^{B'} \mathbf{E}_{k} \left[\left\| \nabla f_{i}(x^{t+1}; \xi_{ik}^{t+1}) - \nabla f_{i}(x^{t}; \xi_{ik}^{t+1}) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right]. \end{split}$$

Due to Assumptions 5 and 6, we have

$$E_{k}\left[\left\|k_{i,1}^{t+1} - E_{k}\left[k_{i,1}^{t+1}\right]\right\|^{2}\right] \leq \frac{2b^{2}\sigma^{2}}{B'p_{\text{mega}}^{2}} + \frac{2L_{\sigma}^{2}}{B'}\left(1 - \frac{b}{p_{\text{mega}}}\right)^{2}\left\|x^{t+1} - x^{t}\right\|^{2}.$$
 (33)

Next, we estimate the bound for $\mathrm{E}_k\left[\left\|k_{i,2}^{t+1}-\mathrm{E}_k\left[k_{i,2}^{t+1}\right]\right\|^2
ight]$.

$$\begin{split} & \mathbf{E}_{k} \left[\left\| k_{i,2}^{t+1} - \mathbf{E}_{k} \left[k_{i,2}^{t+1} \right] \right\|^{2} \right] \\ & = \mathbf{E}_{k} \left[\left\| \frac{1}{B} \sum_{j=1}^{B} \nabla f_{i}(x^{t+1}; \xi_{ij}^{t+1}) - \frac{1}{B} \sum_{j=1}^{B} \nabla f_{i}(x^{t}; \xi_{ij}^{t+1}) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right] \\ & = \frac{1}{B^{2}} \sum_{j=1}^{B} \mathbf{E}_{k} \left[\left\| \nabla f_{i}(x^{t+1}; \xi_{ij}^{t+1}) - \nabla f_{i}(x^{t}; \xi_{ij}^{t+1}) - \left(\nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right) \right\|^{2} \right]. \end{split}$$

Bue to Assumptions 6, we have

$$E_{k}\left[\left\|k_{i,2}^{t+1} - E_{k}\left[k_{i,2}^{t+1}\right]\right\|^{2}\right] \le \frac{L_{\sigma}^{2}}{B} \left\|x^{t+1} - x^{t}\right\|^{2}.$$
(34)

Plugging (33) and (34) into (32), we obtain

$$\begin{split} & \mathbf{E}_{k} \left[\mathbf{E}_{p_{\text{a}}} \left[\mathbf{E}_{p_{\text{mega}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \right] \\ & \leq \frac{p_{\text{mega}}}{n p_{\text{a}}} \left(\frac{2b^{2} \sigma^{2}}{B' p_{\text{mega}}^{2}} + \frac{2L_{\sigma}^{2}}{B'} \left(1 - \frac{b}{p_{\text{mega}}} \right)^{2} \left\| x^{t+1} - x^{t} \right\|^{2} \right) \\ & + \frac{(1 - p_{\text{mega}}) L_{\sigma}^{2}}{n p_{\text{a}} B} \left\| x^{t+1} - x^{t} \right\|^{2} \\ & + \frac{2 \left(p_{\text{a}} - p_{\text{aa}} \right)}{n^{2} p_{\text{a}}^{2}} \sum_{i=1}^{n} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} \end{split}$$

$$+\left.\frac{2\left(p_{\mathrm{a}}-p_{\mathrm{aa}}\right)b^{2}}{n^{2}p_{\mathrm{a}}^{2}p_{\mathrm{mega}}}\sum_{i=1}^{n}\left\|h_{i}^{t}-\nabla f_{i}(x^{t})\right\|^{2}+\left(p_{\mathrm{mega}}\left(1-\frac{b}{p_{\mathrm{mega}}}\right)^{2}+\left(1-p_{\mathrm{mega}}\right)\right)\left\|h^{t}-\nabla f(x^{t})\right\|^{2}.$$

898 Using Assumption 3, we get

$$\begin{split} & \mathbf{E}_{k} \left[\mathbf{E}_{p_{\text{a}}} \left[\mathbf{E}_{p_{\text{mega}}} \left[\left\| h^{t+1} - \nabla f(x^{t+1}) \right\|^{2} \right] \right] \right] \\ & \leq \frac{2b^{2}\sigma^{2}}{np_{\text{mega}}p_{\text{a}}B'} + \left(\frac{2p_{\text{mega}}L_{\sigma}^{2}}{np_{\text{a}}B'} \left(1 - \frac{b}{p_{\text{mega}}} \right)^{2} + \frac{(1 - p_{\text{mega}}) L_{\sigma}^{2}}{np_{\text{a}}B} + \frac{2 \left(p_{\text{a}} - p_{\text{aa}} \right) \widehat{L}^{2}}{np_{\text{a}}^{2}} \right) \left\| x^{t+1} - x^{t} \right\|^{2} \\ & + \frac{2 \left(p_{\text{a}} - p_{\text{aa}} \right) b^{2}}{n^{2}p_{\text{a}}^{2}p_{\text{mega}}} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} + \left(p_{\text{mega}} \left(1 - \frac{b}{p_{\text{mega}}} \right)^{2} + (1 - p_{\text{mega}}) \right) \left\| h^{t} - \nabla f(x^{t}) \right\|^{2}. \end{split}$$

899 Using almost the same derivations, we can prove the second inequality:

$$\begin{split} & E_k \left[E_{p_k} \left[E_{p_{\text{asg}}} \left[\left\| h_{t,1}^{t+1} - \nabla f_i(x^{t+1}) \right\|^2 \right] \right] + (1 - p_{\text{mega}}) E_k \left[E_{p_k} \left[\left\| h_{i,1}^{t+1} - \nabla f_i(x^{t+1}) \right\|^2 \right] \right] + (1 - p_{\text{mega}}) E_k \left[E_{p_k} \left[\left\| h_{i,2}^{t+1} - \nabla f_i(x^{t+1}) \right\|^2 \right] \right] \\ & = p_{\text{mega}} E_k \left[E_{p_k} \left[\left\| h_{i,1}^{t+1} - E_k \left[E_{p_k} \left[h_{i,1}^{t+1} \right] \right] \right]^2 \right] + (1 - p_{\text{mega}}) E_k \left[E_{p_k} \left[\left\| h_{i,2}^{t+1} - E_k \left[E_{p_k} \left[h_{i,2}^{t+1} \right] \right] \right]^2 \right] \right] \\ & + \left(p_{\text{mega}} \left(1 - \frac{b}{p_{\text{mega}}} \right)^2 + (1 - p_{\text{mega}}) \right) \left\| h_i^t - \nabla f_i(x^t) \right\|^2 \\ & = p_{\text{mega}} p_a E_k \left[\left\| h_i^t + \frac{1}{p_a} k_{i,1}^{t+1} - \left(h_i^t + E_k \left[k_{i,1}^{t+1} \right] \right) \right\|^2 \right] \\ & + p_{\text{mega}} (1 - p_a) \left\| h_i^t - \left(h_i^t + E_k \left[k_{i,1}^{t+1} \right] \right) \right\|^2 \\ & + (1 - p_{\text{mega}}) p_a E_k \left[\left\| h_i^t + \frac{1}{p_a} k_{i,2}^{t+1} - \left(h_i^t + E_k \left[k_{i,2}^{t+1} \right] \right) \right\|^2 \right] \\ & + \left(p_{\text{mega}} \left(1 - p_a \right) \right\| h_i^t - \left(h_i^t + E_k \left[k_{i,2}^{t+1} \right] \right) \right\|^2 \\ & + \left(p_{\text{mega}} \left(1 - p_a \right) \right\| \nabla f_i(x^t) - \nabla f_i(x^t) - \frac{b}{p_{\text{mega}}} \left(h_i^t - \nabla f_i(x^t) \right) \right\|^2 \\ & + \left(1 - p_{\text{mega}} \right) p_a E_k \left[\left\| \frac{1}{p_a} k_{i,1}^{t+1} - E_k \left[k_{i,2}^{t+1} \right] \right\|^2 \right] \\ & + \left(1 - p_{\text{mega}} \right) \left(1 - p_a \right) \left\| \nabla f_i(x^{t+1}) - \nabla f_i(x^t) \right\|^2 \\ & + \left(p_{\text{mega}} \left(1 - \frac{b}{p_{\text{mega}}} \right)^2 + \left(1 - p_{\text{mega}} \right) \right) \left\| h_i^t - \nabla f_i(x^t) \right\|^2 \\ & + \frac{1 - p_{\text{mega}}}{p_a} E_k \left[\left\| k_{i,1}^{t+1} - E_k \left[k_{i,1}^{t+1} \right] \right\|^2 \right] \\ & + \frac{1 - p_{\text{mega}}}{p_a} E_k \left[\left\| k_{i,1}^{t+1} - E_k \left[k_{i,1}^{t+1} \right] \right\|^2 \right] \\ & + \frac{1 - p_{\text{mega}}}{p_a} \left[E_k \left[\left\| k_{i,1}^{t+1} - E_k \left[k_{i,1}^{t+1} \right] \right\|^2 \right] \\ & + \frac{1 - p_{\text{mega}}}{p_a} \left[E_k \left[\left\| k_{i,1}^{t+1} - E_k \left[k_{i,1}^{t+1} \right] \right\|^2 \right] \\ & + \frac{1 - p_{\text{mega}}}{p_a} \left[E_k \left[\left\| k_{i,1}^{t+1} - E_k \left[k_{i,1}^{t+1} \right] \right\|^2 \right] \\ & + \frac{1 - p_{\text{mega}}}{p_a} \left[E_k \left[\left\| k_{i,1}^{t+1} - E_k \left[k_{i,1}^{t+1} \right] \right\|^2 \right] \\ & + \frac{1 - p_{\text{mega}}}{p_a} \left[E_k \left[\left| k_{i,1}^{t+1} - E_k \left[k_{i,1}^{t+1} \right] \right] \right] \\ & + \frac{1 - p_{\text{mega}}}{p_a} \left[E_k \left[E_k \left[E_k \left[E_k \left[$$

$$\begin{split} & + \left(p_{\text{mega}} \left(1 - \frac{b}{p_{\text{mega}}} \right)^2 + (1 - p_{\text{mega}}) \right) \left\| h_i^t - \nabla f_i(x^t) \right\|^2 \\ & \leq \frac{p_{\text{mega}}}{p_{\text{a}}} \mathbf{E}_k \left[\left\| k_{i,1}^{t+1} - \mathbf{E}_k \left[k_{i,1}^{t+1} \right] \right\|^2 \right] \\ & + \frac{(1 - p_{\text{mega}})}{p_{\text{a}}} \mathbf{E}_k \left[\left\| k_{i,2}^{t+1} - \mathbf{E}_k \left[k_{i,2}^{t+1} \right] \right\|^2 \right] \\ & + \frac{2(1 - p_{\text{a}})}{p_{\text{a}}} \left\| \nabla f_i(x^{t+1}) - \nabla f_i(x^t) \right\|^2 \\ & + \frac{2(1 - p_{\text{a}})}{p_{\text{mega}}} \left\| h_i^t - \nabla f_i(x^t) \right\|^2 + \left(p_{\text{mega}} \left(1 - \frac{b}{p_{\text{mega}}} \right)^2 + (1 - p_{\text{mega}}) \right) \left\| h_i^t - \nabla f_i(x^t) \right\|^2. \end{split}$$

900 Using (33) and (34), we get

$$\begin{split} & \mathbf{E}_{k} \left[\mathbf{E}_{p_{\text{a}}} \left[\mathbf{E}_{p_{\text{mega}}} \left[\left\| h_{i}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \right] \right] \\ & \leq \frac{2b^{2}\sigma^{2}}{p_{\text{a}}p_{\text{mega}}B'} + \frac{2p_{\text{mega}}L_{\sigma}^{2}}{p_{\text{a}}B'} \left(1 - \frac{b}{p_{\text{mega}}} \right)^{2} \left\| x^{t+1} - x^{t} \right\|^{2} \\ & + \frac{(1 - p_{\text{mega}})L_{\sigma}^{2}}{p_{\text{a}}B} \left\| x^{t+1} - x^{t} \right\|^{2} \\ & + \frac{2(1 - p_{\text{a}})}{p_{\text{a}}} \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} \\ & + \frac{2\left(1 - p_{\text{a}} \right)b^{2}}{p_{\text{mega}}p_{\text{a}}} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} + \left(p_{\text{mega}} \left(1 - \frac{b}{p_{\text{mega}}} \right)^{2} + (1 - p_{\text{mega}}) \right) \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}. \end{split}$$

Next, due to Assumption 3, we obtain

$$\begin{split} & \mathbf{E}_{k} \left[\mathbf{E}_{p_{\mathbf{a}}} \left[\left\| h_{i}^{t+1} - \nabla f_{i}(x^{t+1}) \right\|^{2} \right] \right] \right] \\ & \leq \frac{2b^{2}\sigma^{2}}{p_{\mathbf{a}}p_{\mathrm{mega}}B'} + \left(\frac{2p_{\mathrm{mega}}L_{\sigma}^{2}}{p_{\mathbf{a}}B'} \left(1 - \frac{b}{p_{\mathrm{mega}}} \right)^{2} + \frac{(1 - p_{\mathrm{mega}})L_{\sigma}^{2}}{p_{\mathbf{a}}B} + \frac{2(1 - p_{\mathbf{a}})L_{i}^{2}}{p_{\mathbf{a}}} \right) \left\| x^{t+1} - x^{t} \right\|^{2} \\ & + \frac{2\left(1 - p_{\mathbf{a}} \right)b^{2}}{p_{\mathrm{mega}}p_{\mathbf{a}}} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} + \left(p_{\mathrm{mega}} \left(1 - \frac{b}{p_{\mathrm{mega}}} \right)^{2} + (1 - p_{\mathrm{mega}}) \right) \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2}. \end{split}$$

The third inequality can be proved with the help of (34) and Assumption 3.

$$\begin{aligned}
& \mathbf{E}_{k} \left[\left\| k_{i,2}^{t+1} \right\|^{2} \right] \\
& \stackrel{(13)}{=} \mathbf{E}_{k} \left[\left\| k_{i,2}^{t+1} - \mathbf{E}_{k} \left[k_{i,2}^{t+1} \right] \right\|^{2} \right] + \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} \\
& \leq \frac{L_{\sigma}^{2}}{B} \left\| x^{t+1} - x^{t} \right\|^{2} + \left\| \nabla f_{i}(x^{t+1}) - \nabla f_{i}(x^{t}) \right\|^{2} \\
& \leq \left(\frac{L_{\sigma}^{2}}{B} + L_{i}^{2} \right) \left\| x^{t+1} - x^{t} \right\|^{2}.
\end{aligned}$$

903

Theorem 11. Suppose that Assumptions 1, 2, 3, 5, 6, 7 and 8 hold. Let us take $a = \frac{p_a}{2\omega+1}$, $b = \frac{p_{mega}p_a}{2-p_a}$, probability $p_{mega} \in (0,1]$, batch size $B' \geq B \geq 1$

$$\gamma \leq \left(L + \sqrt{\frac{8\left(2\omega + 1\right)\omega}{np_{\mathrm{a}}^{2}}\left(\widehat{L}^{2} + \frac{L_{\sigma}^{2}}{B}\right) + \frac{16}{np_{\mathrm{mega}}p_{\mathrm{a}}^{2}}\left(\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)\widehat{L}^{2} + \frac{L_{\sigma}^{2}}{B}\right)}\right)^{-1},$$

904 and $h_i^0=g_i^0$ for all $i\in[n]$ in Algorithm 8. Then

$$\mathbb{E}\left[\left\|\nabla f(\widehat{x}^{T})\right\|^{2}\right] \leq \frac{1}{T} \left[\frac{2\Delta_{0}}{\gamma} + \frac{4}{p_{mega}p_{a}} \left\|h^{0} - \nabla f(x^{0})\right\|^{2} + \frac{4\left(1 - \frac{p_{aa}}{p_{a}}\right)}{np_{mega}p_{a}} \frac{1}{n} \sum_{i=1}^{n} \left\|h_{i}^{0} - \nabla f_{i}(x^{0})\right\|^{2} \right]$$

$$+\frac{12\sigma^2}{nB'}$$

Proof. Due to Lemma 2 and the update step from Line 5 in Algorithm 8, we have $\mathrm{E}_{t+1}\left[f(x^{t+1})\right]$

$$\leq \mathbf{E}_{t+1} \left[f(x^{t}) - \frac{\gamma}{2} \|\nabla f(x^{t})\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right) \|x^{t+1} - x^{t}\|^{2} + \frac{\gamma}{2} \|g^{t} - \nabla f(x^{t})\|^{2} \right] \\
= \mathbf{E}_{t+1} \left[f(x^{t}) - \frac{\gamma}{2} \|\nabla f(x^{t})\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right) \|x^{t+1} - x^{t}\|^{2} + \frac{\gamma}{2} \|g^{t} - h^{t} + h^{t} - \nabla f(x^{t})\|^{2} \right] \\
\leq \mathbf{E}_{t+1} \left[f(x^{t}) - \frac{\gamma}{2} \|\nabla f(x^{t})\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right) \|x^{t+1} - x^{t}\|^{2} + \gamma \left(\|g^{t} - h^{t}\|^{2} + \|h^{t} - \nabla f(x^{t})\|^{2}\right) \right].$$

Let us fix constants $\kappa, \eta, \nu, \rho \in [0, \infty)$ that we will define later. Considering Lemma 13, Lemma 14, and the law of total expectation, we obtain

$$\begin{split} & \mathbb{E}\left[f(x^{t+1})\right] + \kappa \mathbb{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^2\right] + \eta \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^n \left\|g_i^{t+1} - h_i^{t+1}\right\|^2\right] \\ & + \nu \mathbb{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^2\right] + \rho \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^n \left\|h_i^{t+1} - \nabla f_i(x^{t+1})\right\|^2\right] \\ & \leq \mathbb{E}\left[f(x^t) - \frac{\gamma}{2} \left\|\nabla f(x^t)\right\|^2 - \left(\frac{1}{2\gamma} - \frac{L}{2}\right) \left\|x^{t+1} - x^t\right\|^2 + \gamma \left(\left\|g^t - h^t\right\|^2 + \left\|h^t - \nabla f(x^t)\right\|^2\right)\right] \\ & + \kappa \mathbb{E}\left[\mathbb{E}_k \left[\mathbb{E}_C \left[\mathbb{E}_{p_k} \left[\mathbb{E}_{p_{\text{mega}}} \left[\left\|g^{t+1} - h^{t+1}\right\|^2\right]\right]\right]\right]\right] \\ & + \eta \mathbb{E}\left[\mathbb{E}_k \left[\mathbb{E}_C \left[\mathbb{E}_{p_k} \left[\mathbb{E}_{p_{\text{mega}}} \left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^2\right]\right]\right]\right]\right] \\ & + \nu \mathbb{E}\left[\mathbb{E}_k \left[\mathbb{E}_{p_k} \left[\mathbb{E}_{p_{\text{mega}}} \left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^2\right]\right]\right]\right] \\ & + \rho \mathbb{E}\left[\mathbb{E}_k \left[\mathbb{E}_{p_k} \left[\mathbb{E}_{p_{\text{mega}}} \left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^2\right]\right]\right]\right] \\ & \leq \mathbb{E}\left[f(x^t) - \frac{\gamma}{2} \left\|\nabla f(x^t)\right\|^2 - \left(\frac{1}{2\gamma} - \frac{L}{2}\right) \left\|x^{t+1} - x^t\right\|^2 + \gamma \left(\left\|g^t - h^t\right\|^2 + \left\|h^t - \nabla f(x^t)\right\|^2\right)\right] \\ & + \kappa \mathbb{E}\left(\frac{2\left(1 - p_{\text{mega}}\right)\omega}{np_a} \left(\frac{L_\sigma^2}{B} + \widehat{L}^2\right) \left\|x^{t+1} - x^t\right\|^2 \\ & + \left(\frac{(p_a - p_{\text{au}})a^2}{n^2p_a^2} + \frac{2\left(1 - p_{\text{mega}}\right)a^2\omega}{n^2p_a}\right) \sum_{i=1}^n \left\|g_i^t - h_i^t\right\|^2 + (1 - a)^2 \left\|g_i^t - h_i^t\right\|^2\right) \\ & + \eta \mathbb{E}\left(\frac{2\left(1 - p_{\text{mega}}\right)\omega}{p_a} \left(\frac{L_\sigma^2}{B} + \widehat{L}^2\right) \left\|x^{t+1} - x^t\right\|^2 \\ & + \left(\frac{(1 - p_a)a^2}{p_a} + \frac{2\left(1 - p_{\text{mega}}\right)a^2\omega}{np_a}\right) \frac{1}{n} \sum_{i=1}^n \left\|g_i^t - h_i^t\right\|^2 + (1 - a)^2 \left\|g_i^t - h_i^t\right\|^2\right) \\ & + \nu \mathbb{E}\left(\frac{2b^2\sigma^2}{np_{\text{mega}}p_a B^T} + \left(\frac{2p_{\text{mega}}L_\sigma^2}{np_a B^T} \left(1 - \frac{b}{p_{\text{mega}}}\right)^2 + \frac{(1 - p_{\text{mega}})L_\sigma^2}{np_a B}} + \frac{2\left((1 - p_{\text{mega}})\widehat{L}^2}{p_a}\right) \left\|k^t - \nabla f(x^t)\right\|^2\right) \right) \\ & + \rho \mathbb{E}\left(\frac{2b^2\sigma^2}{p_a p_{\text{mega}}p_{\text{expa}}} + \left(\frac{2p_{\text{mega}}L_\sigma^2}{np_a B^T} \left(1 - \frac{b}{p_{\text{mega}}}\right)^2 + \frac{(1 - p_{\text{mega}})L_\sigma^2}{np_a B}} + \frac{2\left((1 - p_{\text{mega}})\widehat{L}^2}{p_a}\right) \left\|x^{t+1} - x^t\right\|^2\right) \\ & + \rho \mathbb{E}\left(\frac{2b^2\sigma^2}{p_a p_{\text{mega}}} + \frac{2p_{\text{mega}}L_\sigma^2}{n^2} \left(1 - \frac{b}{p_{\text{mega}}}\right)^2 + \frac{(1 - p_{\text{mega}})L_\sigma^2}{n^2} + \frac{2\left((1 - p_{\text{mega}})\widehat{L}^2}{n^2}\right)} \left\|x^{t+1} - x^t\right\|^2\right) \\ & + \rho \mathbb{E}\left(\frac{2b^2\sigma^2}{p_a p_{\text{mega}}} + \frac{2p_{\text{mega}}L_\sigma^2}{n^2}$$

$$+ \left. \frac{2 \left(1 - p_{\rm a}\right) b^2}{n p_{\rm mega} p_{\rm a}} \sum_{i=1}^n \left\| h_i^t - \nabla f_i(x^t) \right\|^2 + \left(p_{\rm mega} \left(1 - \frac{b}{p_{\rm mega}} \right)^2 + \left(1 - p_{\rm mega}\right) \right) \frac{1}{n} \sum_{i=1}^n \left\| h_i^t - \nabla f_i(x^t) \right\|^2 \right).$$

Let us simplify the last inequality. Since $B' \geq B$ and $b = \frac{p_{\text{mega}}p_{\text{a}}}{2-p_{\text{a}}} \leq p_{\text{mega}}$, we have $1-p_{\text{mega}} \leq 1$,

$$\begin{split} \frac{2p_{\text{mega}}L_{\sigma}^2}{p_{\text{a}}B'}\left(1-\frac{b}{p_{\text{mega}}}\right)^2 &\leq \frac{2p_{\text{mega}}L_{\sigma}^2}{p_{\text{a}}B},\\ \left(p_{\text{mega}}\left(1-\frac{b}{p_{\text{mega}}}\right)^2 + (1-p_{\text{mega}})\right) &\leq 1-b, \end{split}$$

and

$$\left(\frac{2\left(1-p_{\mathrm{a}}\right)b^{2}}{p_{\mathrm{mega}}p_{\mathrm{a}}}+p_{\mathrm{mega}}\left(1-\frac{b}{p_{\mathrm{mega}}}\right)^{2}+\left(1-p_{\mathrm{mega}}\right)\right)\leq1-b.$$

908 Thus

$$\begin{split} & \operatorname{E}\left[f(x^{t+1})\right] + \kappa \operatorname{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \eta \operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g^{t+1}_{i} - h^{t+1}_{i}\right\|^{2}\right] \\ & + \nu \operatorname{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho \operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h^{t+1}_{i} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \operatorname{E}\left[f(x^{t}) - \frac{\gamma}{2}\left\|\nabla f(x^{t})\right\|^{2} - \left(\frac{1}{2\gamma} - \frac{L}{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} + \gamma\left(\left\|g^{t} - h^{t}\right\|^{2} + \left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right)\right] \\ & + \kappa \operatorname{E}\left(\frac{2\omega}{np_{a}}\left(\frac{L_{\sigma}^{2}}{B} + \hat{L}^{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} \\ & + \frac{\left(\left(2\omega + 1\right)p_{a} - p_{aa}\right)a^{2}}{n^{2}p_{a}^{2}}\sum_{i=1}^{n}\left\|g^{t}_{i} - h^{t}_{i}\right\|^{2} + \left(1 - a\right)^{2}\left\|g^{t} - h^{t}\right\|^{2}\right) \\ & + \eta \operatorname{E}\left(\frac{2\omega}{p_{a}}\left(\frac{L_{\sigma}^{2}}{B} + \hat{L}^{2}\right)\left\|x^{t+1} - x^{t}\right\|^{2} \\ & + \frac{\left(2\omega + 1 - p_{a}\right)a^{2}}{p_{a}}\frac{1}{n}\sum_{i=1}^{n}\left\|g^{t}_{i} - h^{t}_{i}\right\|^{2} + \left(1 - a\right)^{2}\left\|g^{t}_{i} - h^{t}_{i}\right\|^{2}\right) \\ & + \nu \operatorname{E}\left(\frac{2b^{2}\sigma^{2}}{np_{\text{mega}}p_{a}B^{t}} + \left(\frac{2L_{\sigma}^{2}}{np_{a}B} + \frac{2\left(p_{a} - p_{\text{aa}}\right)\hat{L}^{2}}{np_{a}^{2}}\right)\left\|x^{t+1} - x^{t}\right\|^{2} \\ & + \frac{2\left(p_{a} - p_{\text{aa}}\right)b^{2}}{n^{2}p_{\text{a}}^{2}p_{\text{mega}}}\sum_{i=1}^{n}\left\|h^{t}_{i} - \nabla f_{i}(x^{t})\right\|^{2} + \left(1 - b\right)\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right) \\ & + \rho \operatorname{E}\left(\frac{2b^{2}\sigma^{2}}{p_{a}p_{\text{mega}}B^{t}} + \left(\frac{2L_{\sigma}^{2}}{p_{a}B} + \frac{2\left(1 - p_{a}\right)\hat{L}^{2}}{p_{a}}\right)\left\|x^{t+1} - x^{t}\right\|^{2} \\ & + \left(1 - b\right)\frac{1}{n}\sum_{i=1}^{n}\left\|h^{t}_{i} - \nabla f_{i}(x^{t})\right\|^{2}\right). \end{split}$$

909 After rearranging the terms, we get

$$E\left[f(x^{t+1})\right] + \kappa E\left[\|g^{t+1} - h^{t+1}\|^{2}\right] + \eta E\left[\frac{1}{n}\sum_{i=1}^{n}\|g_{i}^{t+1} - h_{i}^{t+1}\|^{2}\right] + \nu E\left[\|h^{t+1} - \nabla f(x^{t+1})\|^{2}\right] + \rho E\left[\frac{1}{n}\sum_{i=1}^{n}\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\|^{2}\right]$$

$$\leq \mathbf{E} \left[f(x^{t}) \right] - \frac{\gamma}{2} \mathbf{E} \left[\left\| \nabla f(x^{t}) \right\|^{2} \right]$$

$$- \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{2\kappa\omega}{np_{\mathbf{a}}} \left(\frac{L_{\sigma}^{2}}{B} + \widehat{L}^{2} \right) - \frac{2\eta\omega}{p_{\mathbf{a}}} \left(\frac{L_{\sigma}^{2}}{B} + \widehat{L}^{2} \right) \right)$$

$$- \nu \left(\frac{2L_{\sigma}^{2}}{np_{\mathbf{a}}B} + \frac{2\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)\widehat{L}^{2}}{np_{\mathbf{a}}^{2}} \right) - \rho \left(\frac{2L_{\sigma}^{2}}{p_{\mathbf{a}}B} + \frac{2(1 - p_{\mathbf{a}})\widehat{L}^{2}}{p_{\mathbf{a}}} \right) \right) \mathbf{E} \left[\left\| x^{t+1} - x^{t} \right\|^{2} \right]$$

$$+ \left(\gamma + \kappa \left(1 - a \right)^{2} \right) \mathbf{E} \left[\left\| g^{t} - h^{t} \right\|^{2} \right]$$

$$+ \left(\kappa \frac{\left((2\omega + 1) \ p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}} \right) a^{2}}{np_{\mathbf{a}}^{2}} + \eta \left(\frac{\left(2\omega + 1 - p_{\mathbf{a}} \right) a^{2}}{p_{\mathbf{a}}} + (1 - a)^{2} \right) \right) \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \right]$$

$$+ \left(\gamma + \nu \left(1 - b \right) \right) \mathbf{E} \left[\left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \right]$$

$$+ \left(\nu \frac{2\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right) b^{2}}{np_{\mathbf{a}}^{2} p_{\mathbf{mega}}} + \rho(1 - b) \right) \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \right]$$

$$+ \left(\frac{2\nu b^{2}}{np_{\mathbf{mega}}p_{\mathbf{a}}} + \frac{2\rho b^{2}}{p_{\mathbf{a}} p_{\mathbf{mega}}} \right) \frac{\sigma^{2}}{B'}.$$

910 Let us take $\kappa = \frac{\gamma}{a}$, thus $\gamma + \kappa (1-a)^2 \le \kappa$ and

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma}{a}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \eta\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \nu\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{2\gamma\omega}{anp_{a}}\left(\frac{L_{\sigma}^{2}}{B} + \hat{L}^{2}\right) - \frac{2\eta\omega}{p_{a}}\left(\frac{L_{\sigma}^{2}}{B} + \hat{L}^{2}\right) \\ & - \nu\left(\frac{2L_{\sigma}^{2}}{np_{a}B} + \frac{2\left(p_{a} - p_{aa}\right)\hat{L}^{2}}{np_{a}^{2}}\right) - \rho\left(\frac{2L_{\sigma}^{2}}{p_{a}B} + \frac{2(1 - p_{a})\hat{L}^{2}}{p_{a}}\right)\right)\mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \frac{\gamma}{a}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] \\ & + \left(\frac{\gamma\left((2\omega + 1)p_{a} - p_{aa}\right)a}{np_{a}^{2}} + \eta\left(\frac{(2\omega + 1 - p_{a})a^{2}}{p_{a}} + (1 - a)^{2}\right)\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + (\gamma + \nu\left(1 - b\right))\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \left(\nu\frac{2\left(p_{a} - p_{aa}\right)b^{2}}{np_{a}^{2}p_{mega}} + \rho(1 - b)\right)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{2\nu b^{2}}{np_{mega}p_{a}} + \frac{2\rho b^{2}}{p_{a}p_{mega}}\right)\frac{\sigma^{2}}{B^{i}}. \end{split}$$

911 Next, since
$$a=\frac{p_{\rm a}}{2\omega+1}$$
, we have $\left(\frac{(2\omega+1-p_{\rm a})a^2}{p_{\rm a}}+(1-a)^2\right)\leq 1-a$. We the choice $\eta=\frac{\gamma((2\omega+1)p_{\rm a}-p_{\rm aa})}{np_{\rm a}^2}$, we guarantee $\frac{\gamma((2\omega+1)p_{\rm a}-p_{\rm aa})a}{np_{\rm a}^2}+\eta\left(\frac{(2\omega+1-p_{\rm a})a^2}{p_{\rm a}}+(1-a)^2\right)\leq \eta$ and

$$\mathbb{E}\left[f(x^{t+1})\right] + \frac{\gamma(2\omega + 1)}{p_{a}} \mathbb{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma((2\omega + 1)p_{a} - p_{aa})}{np_{a}^{2}} \mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] + \nu\mathbb{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right]$$

$$\begin{split} &\leq \operatorname{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \operatorname{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ &- \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{2\gamma\left(2\omega + 1\right)\omega}{np_{a}^{2}}\left(\frac{L_{\sigma}^{2}}{B} + \hat{L}^{2}\right) - \frac{2\gamma\left((2\omega + 1)p_{a} - p_{aa}\right)\omega}{np_{a}^{3}}\left(\frac{L_{\sigma}^{2}}{B} + \hat{L}^{2}\right) \\ &- \nu\left(\frac{2L_{\sigma}^{2}}{np_{a}B} + \frac{2\left(p_{a} - p_{aa}\right)\hat{L}^{2}}{np_{a}^{2}}\right) - \rho\left(\frac{2L_{\sigma}^{2}}{p_{a}B} + \frac{2(1 - p_{a})\hat{L}^{2}}{p_{a}}\right)\right)\operatorname{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ &+ \frac{\gamma\left(2\omega + 1\right)}{p_{a}}\operatorname{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma\left((2\omega + 1)p_{a} - p_{aa}\right)}{np_{a}^{2}}\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ &+ \left(\gamma + \nu\left(1 - b\right)\right)\operatorname{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ &+ \left(\nu^{2}\frac{(p_{a} - p_{aa})b^{2}}{np_{a}^{2}p_{mega}} + \rho(1 - b)\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ &+ \left(\frac{2\nu b^{2}}{np_{mega}p_{a}} + \frac{2\rho b^{2}}{p_{a}p_{mega}}\right)\frac{\sigma^{2}}{B'} \\ &\leq \operatorname{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\operatorname{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ &- \nu\left(\frac{2L_{\sigma}^{2}}{np_{a}B} + \frac{2\left(p_{a} - p_{aa}\right)\hat{L}^{2}}{np_{a}^{2}}\right) - \rho\left(\frac{2L_{\sigma}^{2}}{p_{a}B} + \frac{2\left(1 - p_{a}\right)\hat{L}^{2}}{p_{a}}\right)\right)\operatorname{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ &+ \frac{\gamma\left(2\omega + 1\right)}{p_{a}}\operatorname{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma\left(\left(2\omega + 1\right)p_{a} - p_{aa}\right)}{np_{a}^{2}}\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ &+ \left(\gamma + \nu\left(1 - b\right)\right)\operatorname{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ &+ \left(\nu^{2}\frac{(p_{a} - p_{aa})b^{2}}{np_{a}^{2}p_{mega}} + \rho\left(1 - b\right)\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ &+ \left(\frac{2\nu b^{2}}{np_{mega}p_{a}} + \frac{2\rho b^{2}}{p_{a}p_{mega}}\right)\frac{\sigma^{2}}{B'}, \end{aligned}$$

where simplified the term using $p_{\rm aa} \geq 0$. Let us take $\nu = \frac{\gamma}{b}$ to obtain

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma\left(2\omega+1\right)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1}-h^{t+1}\right\|^{2}\right] + \frac{\gamma\left((2\omega+1)\,p_{\mathbf{a}}-p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1}-h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{\gamma}{b}\mathbf{E}\left[\left\|h^{t+1}-\nabla f(x^{t+1})\right\|^{2}\right] + \rho\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1}-\nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma}-\frac{L}{2}-\frac{4\gamma\left(2\omega+1\right)\omega}{np_{\mathbf{a}}^{2}}\left(\frac{L_{\sigma}^{2}}{B}+\hat{L}^{2}\right)\right. \\ & - \left(\frac{2\gamma L_{\sigma}^{2}}{bnp_{\mathbf{a}}B}+\frac{2\gamma\left(p_{\mathbf{a}}-p_{\mathbf{a}\mathbf{a}}\right)\hat{L}^{2}}{bnp_{\mathbf{a}}^{2}}\right) - \rho\left(\frac{2L_{\sigma}^{2}}{p_{\mathbf{a}}B}+\frac{2(1-p_{\mathbf{a}})\hat{L}^{2}}{p_{\mathbf{a}}}\right)\right)\mathbf{E}\left[\left\|x^{t+1}-x^{t}\right\|^{2}\right] \\ & + \frac{\gamma\left(2\omega+1\right)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t}-h^{t}\right\|^{2}\right] + \frac{\gamma\left((2\omega+1)\,p_{\mathbf{a}}-p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2}\right] \\ & + \frac{\gamma}{b}\mathbf{E}\left[\left\|h^{t}-\nabla f(x^{t})\right\|^{2}\right] \end{split}$$

$$\begin{split} & + \left(\frac{2\gamma\left(p_{\mathrm{a}} - p_{\mathrm{aa}}\right)b}{np_{\mathrm{a}}^{2}p_{\mathrm{mega}}} + \rho(1 - b)\right) \mathrm{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{2\gamma b}{np_{\mathrm{mega}}p_{\mathrm{a}}} + \frac{2\rho b^{2}}{p_{\mathrm{a}}p_{\mathrm{mega}}}\right)\frac{\sigma^{2}}{B'}. \end{split}$$

914 Next, we take $ho=rac{2\gamma(p_{\rm a}-p_{\rm aa})}{np_{\rm a}^2p_{
m mega}},$ thus

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma\left(2\omega+1\right)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma\left((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{\gamma}{b}\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{4\gamma\left(2\omega+1\right)\omega}{np_{\mathbf{a}}^{2}}\left(\frac{L_{\sigma}^{2}}{B} + \hat{L}^{2}\right)\right. \\ & - \left(\frac{2\gamma L_{\sigma}^{2}}{bnp_{\mathbf{a}}B} + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)\hat{L}^{2}}{bnp_{\mathbf{a}}^{2}}\right) - \left(\frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}}\right)\left(\frac{2L_{\sigma}^{2}}{p_{\mathbf{a}}B} + \frac{2(1-p_{\mathbf{a}})\hat{L}^{2}}{p_{\mathbf{a}}}\right)\right)\mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \frac{\gamma\left(2\omega+1\right)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma\left((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \frac{\gamma}{b}\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{2\gamma b}{np_{\mathrm{mega}}p_{\mathbf{a}}} + \frac{4\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)b^{2}}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}}\right)\frac{\sigma^{2}}{B'}. \end{split}$$

915 Since $rac{p_{
m mega}p_{
m a}}{2} \le b \le p_{
m mega}p_{
m a}$ and $1-p_{
m a} \le 1-rac{p_{
m aa}}{p_{
m a}} \le 1,$ we get

$$\begin{split} & \operatorname{E}\left[f(x^{t+1})\right] + \frac{\gamma\left(2\omega+1\right)}{p_{\operatorname{a}}}\operatorname{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma\left((2\omega+1)p_{\operatorname{a}} - p_{\operatorname{aa}}\right)}{np_{\operatorname{a}}^{2}}\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{\gamma}{b}\operatorname{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \frac{2\gamma\left(p_{\operatorname{a}} - p_{\operatorname{aa}}\right)}{np_{\operatorname{a}}^{2}p_{\operatorname{mega}}}\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \operatorname{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\operatorname{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{4\gamma\left(2\omega+1\right)\omega}{np_{\operatorname{a}}^{2}}\left(\frac{L_{\sigma}^{2}}{B} + \hat{L}^{2}\right) - \left(\frac{4\gamma L_{\sigma}^{2}}{np_{\operatorname{mega}}p_{\operatorname{a}}^{2}B} + \frac{4\gamma(1-p_{\operatorname{a}})\hat{L}^{2}}{np_{\operatorname{mega}}p_{\operatorname{a}}^{2}}\right)\right)\operatorname{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \frac{\gamma\left(2\omega+1\right)}{p_{\operatorname{a}}}\operatorname{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma\left((2\omega+1)p_{\operatorname{a}} - p_{\operatorname{aa}}\right)}{np_{\operatorname{a}}^{2}}\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \frac{\gamma}{b}\operatorname{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \frac{2\gamma\left(p_{\operatorname{a}} - p_{\operatorname{aa}}\right)}{np_{\operatorname{a}}^{2}p_{\operatorname{mega}}}\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & \leq \operatorname{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\operatorname{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \end{split}$$

$$\begin{split} & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{4\gamma \left(2\omega + 1 \right)\omega}{np_{\rm a}^2} \left(\frac{L_{\sigma}^2}{B} + \widehat{L}^2 \right) - \left(\frac{8\gamma L_{\sigma}^2}{np_{\rm mega}p_{\rm a}^2 B} + \frac{8\gamma \left(1 - \frac{p_{\rm aa}}{p_{\rm a}} \right)\widehat{L}^2}{np_{\rm mega}p_{\rm a}^2} \right) \right) \mathbf{E} \left[\left\| x^{t+1} - x^t \right\|^2 \right] \\ & + \frac{\gamma \left(2\omega + 1 \right)}{p_{\rm a}} \mathbf{E} \left[\left\| g^t - h^t \right\|^2 \right] + \frac{\gamma \left(\left(2\omega + 1 \right)p_{\rm a} - p_{\rm aa} \right)}{np_{\rm a}^2} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^n \left\| g_i^t - h_i^t \right\|^2 \right] \\ & + \frac{\gamma}{b} \mathbf{E} \left[\left\| h^t - \nabla f(x^t) \right\|^2 \right] + \frac{2\gamma \left(p_{\rm a} - p_{\rm aa} \right)}{np_{\rm a}^2 p_{\rm mega}} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^n \left\| h_i^t - \nabla f_i(x^t) \right\|^2 \right] \\ & + \frac{6\gamma \sigma^2}{nB'}. \end{split}$$

Using Lemma 4 and the assumption about γ , we get

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{\gamma\left(2\omega+1\right)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{\gamma\left((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \frac{\gamma\left(2\omega+1\right)}{p_{\mathbf{a}}} \mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \frac{\gamma\left((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \frac{\gamma}{b} \mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \frac{2\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}} \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \frac{6\gamma\sigma^{2}}{nB'}. \end{split}$$

917 It is left to apply Lemma 3 with

$$\Psi^{t} = \frac{(2\omega + 1)}{p_{a}} E\left[\|g^{t} - h^{t}\|^{2} \right] + \frac{((2\omega + 1) p_{a} - p_{aa})}{np_{a}^{2}} E\left[\frac{1}{n} \sum_{i=1}^{n} \|g_{i}^{t} - h_{i}^{t}\|^{2} \right]$$

$$+ \frac{1}{b} E\left[\|h^{t} - \nabla f(x^{t})\|^{2} \right] + \frac{2\left(1 - \frac{p_{aa}}{p_{a}}\right)}{np_{a}p_{mega}} E\left[\frac{1}{n} \sum_{i=1}^{n} \|h_{i}^{t} - \nabla f_{i}(x^{t})\|^{2} \right]$$

and $C = \frac{6\sigma^2}{nB'}$ to conclude the proof.

Corollary 6. Suppose that assumptions from Theorem 11 hold, probability $p_{mega} = \min\left\{\frac{\zeta_{\mathcal{C}}}{d}, \frac{n\varepsilon B}{\sigma^2}\right\}$, batch size $B' = \Theta\left(\frac{\sigma^2}{n\varepsilon}\right)$, and $h_i^0 = g_i^0 = \frac{1}{B_{\text{init}}}\sum_{k=1}^{B_{\text{init}}} \nabla f_i(x^0; \xi_{ik}^0)$ for all $i \in [n]$, initial batch size $B_{\text{init}} = \Theta\left(\frac{B}{p_{\text{mega}}\sqrt{p_a}}\right) = \Theta\left(\max\left\{\frac{Bd}{\sqrt{p_a}\zeta_{\mathcal{C}}}, \frac{\sigma^2}{\sqrt{p_a}n\varepsilon}\right\}\right)$, then DASHA-PP-SYNC-MVR needs

$$T := \mathcal{O}\left(\frac{\Delta_0}{\varepsilon}\left[L + \left(\frac{\omega}{p_{\rm a}\sqrt{n}} + \sqrt{\frac{d}{p_{\rm a}^2\zeta_{\mathcal{C}}n}}\right)\left(\widehat{L} + \frac{L_\sigma}{\sqrt{B}}\right) + \frac{\sigma}{p_{\rm a}\sqrt{\varepsilon}n}\left(\frac{\widehat{L}}{\sqrt{B}} + \frac{L_\sigma}{B}\right)\right] + \frac{\sigma^2}{\sqrt{p_{\rm a}}n\varepsilon B}\right).$$

communication rounds to get an ε -solution, the expected communication complexity is equal to

920 $\mathcal{O}(d+\zeta_C T)$, and the expected number of stochastic gradient calculations per node equals $\mathcal{O}(B_{\text{init}}+$

BT), where $\zeta_{\mathcal{C}}$ is the expected density from Definition 12.

922 *Proof.* Due to the choice of B', we have

$$\begin{split} \mathbf{E} \left[\left\| \nabla f(\widehat{x}^T) \right\|^2 \right] &\leq \frac{1}{T} \left[2\Delta_0 \left(L + \sqrt{\frac{8 \left(2\omega + 1 \right) \omega}{n p_{\mathbf{a}}^2}} \left(\widehat{L}^2 + \frac{L_\sigma^2}{B} \right) + \frac{16}{n p_{\mathrm{mega}} p_{\mathbf{a}}^2} \left(\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathbf{a}}} \right) \widehat{L}^2 + \frac{L_\sigma^2}{B} \right) \right) \\ &+ \frac{4}{p_{\mathrm{mega}} p_{\mathbf{a}}} \left\| h^0 - \nabla f(x^0) \right\|^2 + \frac{4 \left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathbf{a}}} \right)}{n p_{\mathrm{mega}} p_{\mathbf{a}}} \frac{1}{n} \sum_{i=1}^n \left\| h_i^0 - \nabla f_i(x^0) \right\|^2 \\ &+ \frac{2\varepsilon}{3}. \end{split}$$

923 Using

$$E\left[\left\|h^{0} - \nabla f(x^{0})\right\|^{2}\right] = E\left[\left\|\frac{1}{n}\sum_{i=1}^{n}\frac{1}{B_{\text{init}}}\sum_{k=1}^{B_{\text{init}}}\nabla f_{i}(x^{0};\xi_{ik}^{0}) - \nabla f(x^{0})\right\|^{2}\right] \leq \frac{\sigma^{2}}{nB_{\text{init}}}$$

924 and

$$\frac{1}{n^2} \sum_{i=1}^n \mathbf{E} \left[\left\| h_i^0 - \nabla f_i(x^0) \right\|^2 \right] = \frac{1}{n^2} \sum_{i=1}^n \mathbf{E} \left[\left\| \frac{1}{B_{\text{init}}} \sum_{k=1}^{B_{\text{init}}} \nabla f_i(x^0; \xi_{ik}^0) - \nabla f_i(x^0) \right\|^2 \right] \le \frac{\sigma^2}{n B_{\text{init}}},$$

925 we have

$$\begin{split} & \mathbf{E}\left[\left\|\nabla f(\widehat{x}^T)\right\|^2\right] \leq \frac{1}{T}\left[2\Delta_0\left(L + \sqrt{\frac{8\left(2\omega + 1\right)\omega}{np_{\mathrm{a}}^2}\left(\widehat{L}^2 + \frac{L_\sigma^2}{B}\right) + \frac{16}{np_{\mathrm{mega}}p_{\mathrm{a}}^2}\left(\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)\widehat{L}^2 + \frac{L_\sigma^2}{B}\right)\right) \\ & + \frac{8\sigma^2}{np_{\mathrm{mega}}p_{\mathrm{a}}B_{\mathrm{init}}}\right] \\ & + \frac{2\varepsilon}{3}. \end{split}$$

Therefore, we can take the following T to get ε -solution.

$$T = \mathcal{O}\left(\frac{1}{\varepsilon}\left[\Delta_0\left(L + \sqrt{\frac{\omega^2}{np_{\rm a}^2}\left(\widehat{L}^2 + \frac{L_\sigma^2}{B}\right) + \frac{1}{np_{\rm mega}p_{\rm a}^2}\left(\widehat{L}^2 + \frac{L_\sigma^2}{B}\right)}\right) + \frac{\sigma^2}{np_{\rm mega}p_{\rm a}B_{\rm init}}\right]\right)$$

Considering the choice of p_{mega} and B_{init} , we obtain

$$T = \mathcal{O}\left(\frac{1}{\varepsilon}\left[\Delta_0\left(L + \left(\frac{\omega}{p_{\rm a}\sqrt{n}} + \sqrt{\frac{d}{p_{\rm a}^2\zeta_{\mathcal{C}}n}}\right)\left(\widehat{L} + \frac{L_\sigma}{\sqrt{B}}\right) + \frac{\sigma}{p_{\rm a}\sqrt{\varepsilon}n}\left(\frac{\widehat{L}}{\sqrt{B}} + \frac{L_\sigma}{B}\right)\right) + \frac{\sigma^2}{np_{\rm mega}p_{\rm a}B_{\rm init}}\right]\right)$$

$$= \mathcal{O}\left(\frac{\Delta_0}{\varepsilon}\left[L + \left(\frac{\omega}{p_{\rm a}\sqrt{n}} + \sqrt{\frac{d}{p_{\rm a}^2\zeta_{\mathcal{C}}n}}\right)\left(\widehat{L} + \frac{L_\sigma}{\sqrt{B}}\right) + \frac{\sigma}{p_{\rm a}\sqrt{\varepsilon}n}\left(\frac{\widehat{L}}{\sqrt{B}} + \frac{L_\sigma}{B}\right)\right] + \frac{\sigma^2}{\sqrt{p_{\rm a}}n\varepsilon B}\right).$$

The expected communication complexity equals $\mathcal{O}\left(d+p_{\mathrm{mega}}d+(1-p_{\mathrm{mega}})\zeta_{\mathcal{C}}\right)=0$

929 $\mathcal{O}(d+\zeta_{\mathcal{C}})$ and the expected number of stochastic gradient calculations per node equals

930
$$\mathcal{O}\left(B_{\mathsf{init}} + p_{\mathsf{mega}}B' + (1 - p_{\mathsf{mega}})B\right) = \mathcal{O}\left(B_{\mathsf{init}} + B\right)$$
.

Theorem 13. Suppose that Assumptions 1, 2, 3, 5, 6, 7, 8 and 9 hold. Let us take $a = \frac{p_a}{2\omega+1}$, $b = \frac{p_{mega}p_a}{2-p_a}$, probability $p_{mega} \in (0,1]$, batch size $B' \geq B \geq 1$,

$$\gamma \leq \min \left\{ \left(L + \sqrt{\frac{16\left(2\omega + 1\right)\omega}{np_{\mathrm{a}}^2} \left(\frac{L_{\sigma}^2}{B} + \widehat{L}^2\right) + \left(\frac{48L_{\sigma}^2}{np_{\mathrm{mega}}p_{\mathrm{a}}^2B} + \frac{24\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathrm{a}}}\right)\widehat{L}^2}{np_{\mathrm{mega}}p_{\mathrm{a}}^2} \right)} \right)^{-1}, \frac{a}{2\mu}, \frac{b}{2\mu} \right\},$$

and $h_i^0 = g_i^0$ for all $i \in [n]$ in Algorithm 8. Then

Proof. Let us fix constants $\kappa, \eta, \nu, \rho \in [0, \infty)$ that we will define later. As in the proof of Theorem 11, we can get

$$\begin{split} & \operatorname{E}\left[f(x^{t+1})\right] + \kappa \operatorname{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \eta \operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \nu \operatorname{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho \operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \operatorname{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\operatorname{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{2\kappa\omega}{np_{a}}\left(\frac{L_{\sigma}^{2}}{B} + \widehat{L}^{2}\right) - \frac{2\eta\omega}{p_{a}}\left(\frac{L_{\sigma}^{2}}{B} + \widehat{L}^{2}\right) \\ & - \nu\left(\frac{2L_{\sigma}^{2}}{np_{a}B} + \frac{2\left(p_{a} - p_{aa}\right)\widehat{L}^{2}}{np_{a}^{2}}\right) - \rho\left(\frac{2L_{\sigma}^{2}}{p_{a}B} + \frac{2(1 - p_{a})\widehat{L}^{2}}{p_{a}}\right)\right)\operatorname{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \left(\gamma + \kappa\left(1 - a\right)^{2}\right)\operatorname{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] \\ & + \left(\kappa\frac{\left(\left(2\omega + 1\right)p_{a} - p_{aa}\right)a^{2}}{np_{a}^{2}} + \eta\left(\frac{\left(2\omega + 1 - p_{a}\right)a^{2}}{p_{a}} + \left(1 - a\right)^{2}\right)\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \left(\gamma + \nu\left(1 - b\right)\operatorname{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \left(\nu\frac{2\left(p_{a} - p_{aa}\right)b^{2}}{np_{a}^{2}p_{mega}} + \rho(1 - b)\operatorname{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{2\nu b^{2}}{np_{mega}p_{a}} + \frac{2\rho b^{2}}{p_{a}p_{mega}}\right)\frac{\sigma^{2}}{B'}. \end{split}$$

Let us take $\kappa = \frac{2\gamma}{a}$, thus $\gamma + \kappa \left(1 - a\right)^2 \le \left(1 - \frac{a}{2}\right) \kappa$ and

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma}{a} \mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \eta \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \nu \mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho \mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2} \mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{4\gamma\omega}{anp_{a}}\left(\frac{L_{\sigma}^{2}}{B} + \widehat{L}^{2}\right) - \frac{2\eta\omega}{p_{a}}\left(\frac{L_{\sigma}^{2}}{B} + \widehat{L}^{2}\right) \end{split}$$

$$\begin{split} &-\nu\left(\frac{2L_{\sigma}^{2}}{np_{a}B}+\frac{2\left(p_{a}-p_{aa}\right)\widehat{L}^{2}}{np_{a}^{2}}\right)-\rho\left(\frac{2L_{\sigma}^{2}}{p_{a}B}+\frac{2(1-p_{a})\widehat{L}^{2}}{p_{a}}\right)\right)\mathbb{E}\left[\left\|x^{t+1}-x^{t}\right\|^{2}\right]\\ &+\left(1-\frac{a}{2}\right)\frac{2\gamma}{a}\mathbb{E}\left[\left\|g^{t}-h^{t}\right\|^{2}\right]\\ &+\left(\frac{2\gamma\left((2\omega+1)p_{a}-p_{aa}\right)a}{np_{a}^{2}}+\eta\left(\frac{(2\omega+1-p_{a})a^{2}}{p_{a}}+(1-a)^{2}\right)\right)\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2}\right]\\ &+\left(\gamma+\nu\left(1-b\right)\right)\mathbb{E}\left[\left\|h^{t}-\nabla f(x^{t})\right\|^{2}\right]\\ &+\left(\nu\frac{2\left(p_{a}-p_{aa}\right)b^{2}}{np_{a}^{2}p_{mega}}+\rho(1-b)\right)\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t}-\nabla f_{i}(x^{t})\right\|^{2}\right]\\ &+\left(\frac{2\nu b^{2}}{np_{mega}p_{a}}+\frac{2\rho b^{2}}{p_{a}p_{mega}}\right)\frac{\sigma^{2}}{B'}. \end{split}$$

935 Next, since $a=\frac{p_{\rm a}}{2\omega+1}$, we have $\left(\frac{(2\omega+1-p_{\rm a})a^2}{p_{\rm a}}+(1-a)^2\right)\leq 1-a$. We the choice $\eta=\frac{2\gamma((2\omega+1)p_{\rm a}-p_{\rm aa})}{np_{\rm a}^2}$, we guarantee $\frac{\gamma((2\omega+1)p_{\rm a}-p_{\rm aa})a}{np_{\rm a}^2}+\eta\left(\frac{(2\omega+1-p_{\rm a})a^2}{p_{\rm a}}+(1-a)^2\right)\leq \left(1-\frac{a}{2}\right)\eta$ and

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{2\gamma\left((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \nu\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \rho\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & - \nu\left(\frac{2L_{\sigma}^{2}}{2\gamma} - \frac{8\gamma\left(2\omega+1\right)\omega}{np_{\mathbf{a}}^{2}}\left(\frac{L_{\sigma}^{2}}{B} + \widehat{L}^{2}\right)\right. \\ & - \nu\left(\frac{2L_{\sigma}^{2}}{np_{\mathbf{a}}B} + \frac{2\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)\widehat{L}^{2}}{np_{\mathbf{a}}^{2}}\right) - \rho\left(\frac{2L_{\sigma}^{2}}{p_{\mathbf{a}}B} + \frac{2\left(1 - p_{\mathbf{a}}\right)\widehat{L}^{2}}{p_{\mathbf{a}}}\right)\right)\mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \left(1 - \frac{a}{2}\right)\frac{2\gamma\left((2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] \\ & + \left(\gamma + \nu\left(1 - b\right)\right)\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] \\ & + \left(\nu\frac{2\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)b^{2}}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}} + \rho(1 - b)\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{2\nu b^{2}}{np_{\mathrm{mega}}p_{\mathbf{a}}} + \frac{2\rho b^{2}}{p_{\mathbf{a}}p_{\mathrm{mega}}}\right)\frac{\sigma^{2}}{B^{t}}, \end{split}$$

where simplified the term using $p_{\rm aa} \geq 0$. Let us take $\nu = \frac{2\gamma}{b}$ to obtain

$$\begin{split} &-\left(\frac{1}{2\gamma}-\frac{L}{2}-\frac{8\gamma\left(2\omega+1\right)\omega}{np_{\mathrm{a}}^{2}}\left(\frac{L_{\sigma}^{2}}{B}+\widehat{L}^{2}\right)\right.\\ &-\left(\frac{4\gamma L_{\sigma}^{2}}{bnp_{\mathrm{a}}B}+\frac{4\gamma\left(p_{\mathrm{a}}-p_{\mathrm{aa}}\right)\widehat{L}^{2}}{bnp_{\mathrm{a}}^{2}}\right)-\rho\left(\frac{2L_{\sigma}^{2}}{p_{\mathrm{a}}B}+\frac{2(1-p_{\mathrm{a}})\widehat{L}^{2}}{p_{\mathrm{a}}}\right)\right)\mathrm{E}\left[\left\|x^{t+1}-x^{t}\right\|^{2}\right]\\ &+\left(1-\frac{a}{2}\right)\frac{2\gamma(2\omega+1)}{p_{\mathrm{a}}}\mathrm{E}\left[\left\|g^{t}-h^{t}\right\|^{2}\right]+\left(1-\frac{a}{2}\right)\frac{2\gamma\left((2\omega+1)p_{\mathrm{a}}-p_{\mathrm{aa}}\right)}{np_{\mathrm{a}}^{2}}\mathrm{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t}-h_{i}^{t}\right\|^{2}\right]\\ &+\left(1-\frac{b}{2}\right)\frac{2\gamma}{b}\mathrm{E}\left[\left\|h^{t}-\nabla f(x^{t})\right\|^{2}\right]\\ &+\left(\frac{4\gamma\left(p_{\mathrm{a}}-p_{\mathrm{aa}}\right)b}{np_{\mathrm{a}}^{2}p_{\mathrm{mega}}}+\rho(1-b)\right)\mathrm{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t}-\nabla f_{i}(x^{t})\right\|^{2}\right]\\ &+\left(\frac{4\gamma b}{np_{\mathrm{mega}}p_{\mathrm{a}}}+\frac{2\rho b^{2}}{p_{\mathrm{a}}p_{\mathrm{mega}}}\right)\frac{\sigma^{2}}{B'}, \end{split}$$

938 Next, we take $\rho = \frac{8\gamma(p_a - p_{aa})}{np_a^2 p_{\text{mega}}}$, thus

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{2\gamma\left((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{8\gamma\left(2\omega+1\right)\omega}{np_{\mathbf{a}}^{2}}\left(\frac{L_{\sigma}^{2}}{B} + \hat{L}^{2}\right)\right. \\ & - \left(\frac{4\gamma L_{\sigma}^{2}}{bnp_{\mathbf{a}}B} + \frac{4\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)\hat{L}^{2}}{bnp_{\mathbf{a}}^{2}}\right) - \left(\frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}}\right)\left(\frac{2L_{\sigma}^{2}}{p_{\mathbf{a}}B} + \frac{2(1-p_{\mathbf{a}})\hat{L}^{2}}{p_{\mathbf{a}}}\right)\right)\mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \\ & + \left(1 - \frac{a}{2}\right)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \left(1 - \frac{a}{2}\right)\frac{2\gamma\left((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \left(1 - \frac{b}{2}\right)\frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(1 - \frac{b}{2}\right)\frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \left(\frac{4\gamma b}{np_{\mathrm{mega}}p_{\mathbf{a}}} + \frac{16\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)b^{2}}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}}\right)\frac{\sigma^{2}}{B^{\prime}}, \end{split}$$

939 Since $rac{p_{
m mega}p_{
m a}}{2} \le b \le p_{
m mega}p_{
m a}$ and $1-p_{
m a} \le 1-rac{p_{
m aa}}{p_{
m a}} \le 1,$ we get

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{2\gamma\left((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{8\gamma\left(2\omega+1\right)\omega}{np_{\mathbf{a}}^{2}}\left(\frac{L_{\sigma}^{2}}{B} + \hat{L}^{2}\right) - \left(\frac{16\gamma L_{\sigma}^{2}}{np_{\mathrm{mega}}p_{\mathbf{a}}^{2}B} + \frac{16\gamma(1-p_{\mathbf{a}})\hat{L}^{2}}{np_{\mathrm{mega}}p_{\mathbf{a}}^{2}}\right)\right)\mathbf{E}\left[\left\|x^{t+1} - x^{t}\right\|^{2}\right] \end{split}$$

$$\begin{split} & + \left(1 - \frac{a}{2}\right) \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbf{E} \left[\left\| g^{t} - h^{t} \right\|^{2} \right] + \left(1 - \frac{a}{2}\right) \frac{2\gamma\left((2\omega + 1)\,p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \right] \\ & + \left(1 - \frac{b}{2}\right) \frac{2\gamma}{b} \mathbf{E} \left[\left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \right] + \left(1 - \frac{b}{2}\right) \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{aa}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \right] \\ & + \frac{20\gamma\sigma^{2}}{nB'} \\ & \leq \mathbf{E} \left[f(x^{t}) \right] - \frac{\gamma}{2} \mathbf{E} \left[\left\| \nabla f(x^{t}) \right\|^{2} \right] \\ & - \left(\frac{1}{2\gamma} - \frac{L}{2} - \frac{8\gamma\left(2\omega + 1\right)\omega}{np_{\mathbf{a}}^{2}} \left(\frac{L_{\sigma}^{2}}{B} + \hat{L}^{2} \right) - \left(\frac{24\gamma L_{\sigma}^{2}}{np_{\mathrm{mega}}p_{\mathbf{a}}^{2}B} + \frac{24\gamma\left(1 - \frac{p_{\mathrm{aa}}}{p_{\mathbf{a}}}\right)\hat{L}^{2}}{np_{\mathrm{mega}}p_{\mathbf{a}}^{2}} \right) \right) \mathbf{E} \left[\left\| x^{t+1} - x^{t} \right\|^{2} \right] \\ & + \left(1 - \frac{a}{2}\right) \frac{2\gamma(2\omega + 1)}{p_{\mathbf{a}}} \mathbf{E} \left[\left\| g^{t} - h^{t} \right\|^{2} \right] + \left(1 - \frac{a}{2}\right) \frac{2\gamma\left((2\omega + 1)\,p_{\mathbf{a}} - p_{\mathrm{aa}}\right)}{np_{\mathbf{a}}^{2}} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \right] \\ & + \left(1 - \frac{b}{2}\right) \frac{2\gamma}{b} \mathbf{E} \left[\left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \right] + \left(1 - \frac{b}{2}\right) \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathrm{aa}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \right] \\ & + \frac{20\gamma\sigma^{2}}{nB'}. \end{split}$$

Using Lemma 4 and the assumption about γ , we get

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{2\gamma\left((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + \left(1 - \frac{a}{2}\right)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + \left(1 - \frac{a}{2}\right)\frac{2\gamma\left((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + \left(1 - \frac{b}{2}\right)\frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + \left(1 - \frac{b}{2}\right)\frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \frac{20\gamma\sigma^{2}}{nB'}. \end{split}$$

941 Due to $\gamma \leq \frac{a}{2\mu}$ and $\gamma \leq \frac{b}{2\mu}$, we have

$$\begin{split} & \mathbf{E}\left[f(x^{t+1})\right] + \frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t+1} - h^{t+1}\right\|^{2}\right] + \frac{2\gamma\left((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t+1} - h_{i}^{t+1}\right\|^{2}\right] \\ & + \frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t+1} - \nabla f(x^{t+1})\right\|^{2}\right] + \frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t+1} - \nabla f_{i}(x^{t+1})\right\|^{2}\right] \\ & \leq \mathbf{E}\left[f(x^{t})\right] - \frac{\gamma}{2}\mathbf{E}\left[\left\|\nabla f(x^{t})\right\|^{2}\right] \\ & + (1 - \gamma\mu)\frac{2\gamma(2\omega+1)}{p_{\mathbf{a}}}\mathbf{E}\left[\left\|g^{t} - h^{t}\right\|^{2}\right] + (1 - \gamma\mu)\frac{2\gamma\left((2\omega+1)\,p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|g_{i}^{t} - h_{i}^{t}\right\|^{2}\right] \\ & + (1 - \gamma\mu)\frac{2\gamma}{b}\mathbf{E}\left[\left\|h^{t} - \nabla f(x^{t})\right\|^{2}\right] + (1 - \gamma\mu)\frac{8\gamma\left(p_{\mathbf{a}} - p_{\mathbf{a}\mathbf{a}}\right)}{np_{\mathbf{a}}^{2}p_{\mathrm{mega}}}\mathbf{E}\left[\frac{1}{n}\sum_{i=1}^{n}\left\|h_{i}^{t} - \nabla f_{i}(x^{t})\right\|^{2}\right] \\ & + \frac{20\gamma\sigma^{2}}{nB'}. \end{split}$$

942 It is left to apply Lemma 11 with

$$\begin{split} \Psi^{t} & = & \frac{2(2\omega+1)}{p_{\mathrm{a}}} \mathbf{E} \left[\left\| g^{t} - h^{t} \right\|^{2} \right] + \frac{2\left(\left(2\omega+1 \right) p_{\mathrm{a}} - p_{\mathrm{aa}} \right)}{np_{\mathrm{a}}^{2}} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| g_{i}^{t} - h_{i}^{t} \right\|^{2} \right] \\ & + & \frac{2}{b} \mathbf{E} \left[\left\| h^{t} - \nabla f(x^{t}) \right\|^{2} \right] + \frac{8\left(p_{\mathrm{a}} - p_{\mathrm{aa}} \right)}{np_{\mathrm{a}}^{2} p_{\mathrm{mega}}} \mathbf{E} \left[\frac{1}{n} \sum_{i=1}^{n} \left\| h_{i}^{t} - \nabla f_{i}(x^{t}) \right\|^{2} \right] \end{split}$$

and $C=\frac{20\sigma^2}{nB'}$ to conclude the proof.