# DISTRIBUTED PROXIMAL GRADIENT ALGORITHM FOR PARTIALLY ASYNCHRONOUS COMPUTER CLUSTERS

YI ZHOU*, YAOLIANG YU†, WEI DAI‡, YINGBIN LIANG*, AND ERIC P. XING‡§

**Abstract.** With ever growing data volume and model size, an error-tolerant, communication efficient, yet versatile distributed algorithm has become vital for the success of many large-scale machine learning applications. In this work we propose m-PAPG, an implementation of the flexible proximal gradient algorithm in model parallel systems equipped with the partially asynchronous communication protocol. The worker machines communicate asynchronously with a controlled staleness bound $s$ and operate at different frequencies. We characterize various convergence properties of m-PAPG: 1) Under a general non-smooth and non-convex setting, we prove that every limit point of the sequence generated by m-PAPG is a critical point of the objective function; 2) Under an error bound condition, we prove that the function value decays linearly for every $s$ steps; 3) Under the Kurdyka-Łojasiewicz inequality, we prove that the sequences generated by m-PAPG converge to the same critical point, provided that a proximal Lipschitz condition is satisfied.

**Key words.** Proximal gradient, distributed system, model parallel, partially asynchronous, machine learning

**AMS subject classifications.**

**1. Introduction.** The composite minimization problem

$$(1) \qquad \min_{\mathbf{x} \in \mathbb{R}^d} \; f(\mathbf{x}) + g(\mathbf{x})$$

has drawn a lot of recent attention due to its ubiquity in machine learning and statistical applications. Typically, the first term

$$(2) \qquad f(\mathbf{x}) := \frac{1}{n} \sum_{i=1}^{n} f_i(\mathbf{x})$$

is a smooth loss function over $n$ training samples that describes the fitness to data, and the second term $g$ is a nonsmooth regularization function that encodes *a priori* information. We list below some popular examples under this framework.

- Lasso: least squares loss $f_i(\mathbf{x}) = (y_i - \mathbf{a}_i^\top \mathbf{x})^2$ and $\ell_1$ norm regularizer $g(\mathbf{x}) = \|\mathbf{x}\|_1$;
- Logistic regression: logistic loss $f_i = \log(1 + \exp(-y_i \mathbf{a}_i^\top \mathbf{x}_i))$;
- Boosting: exponential loss $f_i(\mathbf{x}) = \exp(-y_i \mathbf{a}_i^\top \mathbf{x})$;
- Support vector machines: hinge loss $f_i(\mathbf{x}) = \max\{0, 1 - y_i \mathbf{a}_i^\top \mathbf{x}\}$ and (squared) $\ell_2$ norm regularizer $g(\mathbf{x}) = \|\mathbf{x}\|_2^2$.

Over the years there is also a rising interest in using nonconvex losses $f$ (mainly for robustness against outlying observations) [12, 34, 35, 36] and nonconvex regularizers $g$ (mainly for smaller bias in statistical estimation) [15, 39].

Due to the apparent importance of the composite minimization framework and the rapidly growing size in both dimension ($d$) and volume ($n$) of data, there is a strong need to develop a practical *parallel* system that can solve the problem in (1) efficiently

*Syracuse University. (yzhou35@syr.edu, yliang06@syr.edu).
†University of Waterloo. (yaoliang.yu@uwaterloo.ca).
‡Carnegie Mellon University. (wdai@cs.cmu.edu, epxing@cs.cmu.edu).

and in a scale that is impossible for a single machine [2, 8, 13, 17, 19, 21, 23, 37]. Existing systems can be categorized by how communication among worker machines is managed: bulk synchronous (also called fully synchronous) [13, 33, 37], totally asynchronous [5, 8, 23], and partially asynchronous (a.k.a. stale synchronous or chaotic) [2, 8, 11, 17, 19, 21, 32]. Bulk synchronous parallel (BSP) systems explicitly force synchronization barriers so that the worker machines can stay on the same page to ensure correctness. However, in a real deployed parallel system, BSP usually suffers from the straggler problem, that is, the performance of the whole system is bottlenecked at the bandwidth of communication and the *slowest* worker machine. On the other hand, totally asynchronous systems do not put any constraint on synchronization, hence achieve much greater throughputs by potentially sacrificing the correctness of the algorithm. Partially asynchronous parallel (PAP) systems [8, 11] are a compromise between the previous two: it allows the worker machines to communicate asynchronously up to a controlled staleness and to perform updates at different paces. PAP is particularly suitable for machine learning applications, where iterative algorithms that are robust to small computational errors are usually favored for finding an appropriate solution. Due to its flexibility, the PAP mechanism has been the method of choice in many recent practical implementations [2, 17, 19, 21, 22, 28].

Existing parallel systems can also be categorized by how computation is divided among worker machines: data parallel and model parallel. Data parallel systems usually distribute the computation involving each component function $f_i$ in (2) into different worker machines, which is suitable when $n \gg d$, i.e., large data volume but moderate model size. In this setting the stochastic proximal gradient algorithm, along with the PAP protocol, has been shown to be quite effective in solving the composite problem (1) [2, 17, 19, 21]. In this work, we focus on the "dual" model parallel regime where $d \gg n$, i.e., large model size but moderate data volume. In modern machine learning and statistics applications, it is not uncommon that the dimensionality of data largely exceeds its volume, for example, in computational biology, conducting an experimental study that involves many patients can be very expensive but for each patient, technology (e.g. next-generation genome sequencing) has advanced to a stage where taking a large number of measurements (model parameters) is relatively cheap. Deep neural networks are another example that calls for model parallelism. Not surprisingly, the design of a model parallel system is fundamentally different from that of a data parallel system, and so is the subsequent analysis.

To achieve model parallelism, the model $\mathbf{x}$ is partitioned into different (disjoint) blocks and is distributed among many worker machines. In this setting, the block proximal gradient algorithm has been proposed to solve the composite problem (1) [16, 24, 29], although under the more restrictive BSP protocol. Under the PAP protocol, the only work that we are aware of is [8] which focused on a special case of (1) where $g$ is an indicator function of a convex set, and [32] which established a periodic linear rate of convergence under an error bound condition. Our main goal in this work is to provide a formal convergence analysis of the model parallel proximal gradient algorithm under the more flexible PAP communication protocol, and our results naturally extend those in [8, 32] to allow nonsmooth and nonconvex functions.

Our main contributions in this work are: 1). We propose m-PAPG, an extension of the proximal gradient algorithm to the model parallel and partially asynchronous setting. 2). We provide a rigorous analysis of the convergence properties of m-PAPG, allowing both *nonsmooth* and *nonconvex* functions. In particular, we prove in Theorem 6 that any limit point of the sequences generated by m-PAPG is a critical point. 3) Under an additional error bound condition, we prove in Theorem 8 that the

function values generated by m-PAPG decays periodically linearly. 4) Lastly, using the Kurdyka-Łojasiewicz (KŁ) inequality [10], we prove in Theorem 10 that for functions that satisfy a proximal Lipschitz condition the whole sequences of m-PAPG converge to a single critical point.

This paper proceeds as follows: We first set up the notations and definitions in Section 2. The proposed algorithm m-PAPG is presented in Section 3, and convergence analysis are detailed in Sections 4 to 6. Section 7 concludes our work.

**2. Preliminaries.** We first recall some fundamental definitions that will be needed in our analysis. Throughout, $h : \mathbb{R}^d \to (-\infty, +\infty]$ denotes an extended real-valued function that is proper and closed, i.e., its domain $\operatorname{dom} h := \{\mathbf{x} : h(\mathbf{x}) < +\infty\}$ is nonempty and its sublevel set $\{\mathbf{x} : h(\mathbf{x}) \leq \alpha\}$ is closed for all $\alpha \in \mathbb{R}$. Since the function $h$ may not be smooth or convex, we need the following generalized notion of "derivative."

DEFINITION 1 (Subdifferential and critical point, e.g. [30]). *The Frechét subdifferential $\hat{\partial}h$ of $h$ at $\mathbf{x} \in \operatorname{dom} h$ is the set of $\mathbf{u}$ such that*

$$(3) \qquad \liminf_{\mathbf{z} \neq \mathbf{x}, \mathbf{z} \to \mathbf{x}} \frac{h(\mathbf{z}) - h(\mathbf{x}) - \mathbf{u}^\top (\mathbf{z} - \mathbf{x})}{\|\mathbf{z} - \mathbf{x}\|} \geq 0,$$

*while the (limiting) subdifferential $\partial h$ at $\mathbf{x} \in \operatorname{dom} h$ is the "closure" of $\hat{\partial}h$:*

$$(4) \qquad \{\mathbf{u} : \exists \mathbf{x}^k \to \mathbf{x}, h(\mathbf{x}^k) \to h(\mathbf{x}), \mathbf{u}^k \in \hat{\partial}h(\mathbf{x}^k), \mathbf{u}^k \to \mathbf{u}\}.$$

*The critical points of $h$ are $\operatorname{crit} h := \{\mathbf{x} : \mathbf{0} \in \partial h(\mathbf{x})\}$.*

When $h$ is continuously differentiable or convex, the subdifferential $\partial h$ and the set of critical points $\operatorname{crit} h$ coincide with the usual notions. For a closed function $h$, its subdifferential is either nonempty at any point in its domain or the subgradient diverges to some "direction" [30, Corollary 8.10].

DEFINITION 2 (Distance and projection). *The distance function w.r.t. a closed set $\Omega \subseteq \mathbb{R}^d$ is defined as:*

$$(5) \qquad \operatorname{dist}_\Omega(\mathbf{x}) := \min_{\mathbf{y} \in \Omega} \|\mathbf{y} - \mathbf{x}\|,$$

*while the metric projection onto $\Omega$ is defined as:*

$$(6) \qquad \operatorname{proj}_\Omega(\mathbf{x}) := \operatorname*{argmin}_{\mathbf{y} \in \Omega} \|\mathbf{y} - \mathbf{x}\|,$$

*where $\| \cdot \|$ is the usual Euclidean norm.*

Note that $\operatorname{proj}_\Omega(\mathbf{x})$ is single-valued for all $\mathbf{x} \in \mathbb{R}^d$ if and only if $\Omega$ is convex.

DEFINITION 3 (Proximal map, e.g. [30]). *The proximal map of a closed and proper function $h$ is (with parameter $\eta > 0$):*

$$(7) \qquad \operatorname{prox}_h^\eta(\mathbf{x}) := \operatorname*{argmin}_{\mathbf{z} \in \mathbb{R}^d} h(\mathbf{z}) + \tfrac{1}{2\eta}\|\mathbf{z} - \mathbf{x}\|^2.$$

*Occasionally, we will write $\operatorname{prox}_h$ instead of $\operatorname{prox}_h^1$.*

Clearly, for the indicator function $h(\mathbf{x}) = \iota_\Omega(\mathbf{x})$, which takes the value 0 for $\mathbf{x} \in \Omega$ and $\infty$ otherwise, its proximal map (with any $\eta > 0$) reduces to the metric

projection $\mathrm{proj}_\Omega$. If $h$ decreases slower than a quadratic function (in particular, when $h$ is bounded below), then its proximal map is well-defined for all (small) $\eta$ [30]. If $h$ is convex, then its proximal map is always a singleton while for nonconvex $h$, the proximal map can be set-valued. In the latter case we will also abuse the notation $\mathrm{prox}_h^\eta(\mathbf{x})$ for an arbitrary element from that set. For convex functions, the proximal map is nonexpansive:

$$(8) \qquad\qquad \forall \mathbf{x}, \forall \mathbf{y}, \ \|\mathrm{prox}_h^\eta(\mathbf{x}) - \mathrm{prox}_h^\eta(\mathbf{y})\| \le \|\mathbf{x} - \mathbf{y}\|,$$

while for nonconvex functions this may not hold everywhere.

The proximal map is the key component of the proximal gradient algorithm [18] (a.k.a. forward-backward splitting):

$$(9) \qquad\qquad \forall\, t = 0, 1, \ldots, \quad \mathbf{x}(t+1) = \mathrm{prox}_g^\eta\big(\mathbf{x}(t) - \eta \nabla f(\mathbf{x}(t))\big),$$

where $\nabla f$ is the (sub)gradient of $f$, and $\eta$ is a suitable step size (that may change with $t$). It is known that when $f$ is convex with $L$-Lipschitz continuous gradient and $0 < \eta < 2/L$, then $F_t := f(\mathbf{x}(t)) + g(\mathbf{x}(t))$ converges to the minimum at the rate $O(1/t)$ and $\mathbf{x}(t)$ converges to some minimizer $\mathbf{x}^*$. Accelerated versions [6, 27] where $F_t$ converges at the faster rate $O(1/t^2)$ are also well-known. Recently, [10] proved that $\mathbf{x}(t)$ converges to a critical point even for nonconvex $f$ and nonconvex and nonsmooth $g$ as long as together they satisfy a certain KŁ inequality.

**3. Formulation of m-PAPG.** Recall the composite minimization problem:

$$(\mathrm{P}) \qquad\qquad \min_{\mathbf{x} \in \mathbb{R}^d} F(\mathbf{x}), \quad \text{where} \quad F(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x}).$$

We are interested in the case where $d$ is so large that implementing the proximal gradient algorithm (9) on a single machine is no longer feasible, hence distributed computation is necessary.

We consider a **model** parallel system with $p$ machines in total, and decompose the $d$ model parameters into $p$ disjoint groups. Formally, consider the decomposition $\mathbb{R}^d = \mathbb{R}^{d_1} \times \mathbb{R}^{d_2} \times \cdots \times \mathbb{R}^{d_p}$, and denote $x_i$ and $\nabla_i f(\mathbf{x}) : \mathbb{R}^d \to \mathbb{R}^{d_i}$ as the $i$-th component of $\mathbf{x}$ and $\nabla f(\mathbf{x})$, respectively. Clearly, $\mathbf{x} = (x_1, x_2, \ldots, x_p)$ and $\nabla f = (\nabla_1 f, \nabla_2 f, \cdots, \nabla_p f)$. The $i$-th machine is responsible for updating the component $x_i \in \mathbb{R}^{d_i}$, and for the purpose of evaluating the partial gradient $\nabla_i f(\mathbf{x})$ we assume the $i$-th machine also has access to a local, full model parameter $\mathbf{x}^i \in \mathbb{R}^d$. The last assumption is made only to simplify our presentation; it can be removed for many machine learning problems, see for instance [29, 41].

We make the following standard assumptions regarding problem (P):

ASSUMPTION 1 (Bounded Below). *The function $F = f + g$ is bounded below.*

ASSUMPTION 2 (Smooth). *The gradient $\nabla f$ of $f$ is $L$-Lipschitz continuous:*

$$(10) \qquad\qquad \forall \mathbf{x}, \forall \mathbf{y}, \ \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \le L\|\mathbf{x} - \mathbf{y}\|.$$

ASSUMPTION 3 (Separable). *The function $g$ is closed and separable, i.e., $g(\mathbf{x}) = \sum_{i=1}^p g_i(x_i)$.*

Assumption 1 simply allows us to have a finite minimum value and is usually satisfied in practice. The smoothness assumption is critical in two aspects: (1) It allows us to upper bound $f$ by its quadratic expansion at the current iterate—a standard step in the convergence proof of gradient type algorithms:

$$(11) \qquad\qquad \forall \mathbf{x}, \forall \mathbf{y}, \ f(\mathbf{x}) \le f(\mathbf{y}) + \langle \mathbf{x} - \mathbf{y}, \nabla f(\mathbf{y}) \rangle + \tfrac{L}{2}\|\mathbf{x} - \mathbf{y}\|^2.$$

(2) It allows us to bound the inconsistencies in different machines due to asynchronous updates, see Lemma 4 below. The separable assumption is what makes model parallelism interesting and feasible. We remark that both Assumption 2 and Assumption 3 can be relaxed using techniques in [7] and [36], respectively. For brevity we do not pursue these extensions here. Note that we do *not* assume convexity on either $f$ or $g$, and $g$ need not even be continuous.

We now specify the m-PAPG algorithm for solving (P) under model parallelism and the PAP protocol. The separable assumption on $g$ implies that

$$(12) \qquad \mathrm{prox}_g^\eta(\mathbf{x}) = \big(\mathrm{prox}_{g_1}^\eta(x_1), \ \ldots \ , \ \mathrm{prox}_{g_p}^\eta(x_p)\big).$$

Then, the update on machine $i$ is defined as:

$$(13) \qquad x_i \leftarrow \mathrm{prox}_{g_i}^\eta(x_i - \eta \nabla_i f(\mathbf{x}^i)).$$

That is, machine $i$ computes a partial gradient mapping [27] w.r.t. the $i$-th component using the local component $x_i$ and the local full model $\mathbf{x}^i$. To define the latter, consider a global clock shared by all machines and denote $T_i$ as the set of active clocks when machine $i$ performs an update. Note that the global clock is introduced solely for the purpose of our analysis, and the machines need not maintain it in a practical implementation. Formally, the $t$-th iteration on machine $i$ can be written as:

$$(\text{m-PAPG}) \qquad \begin{cases} \forall i, \ x_i(t+1) = \begin{cases} x_i(t), & t \notin T_i \\ \mathrm{prox}_{g_i}^\eta(x_i(t) - \eta \nabla_i f(\mathbf{x}^i(t))), & t \in T_i \end{cases}, \\ (\text{local}) \quad \mathbf{x}^i(t) = \big(x_1(\tau_1^i(t)), \ \ldots \ , \ x_p(\tau_p^i(t))\big), \\ (\text{global}) \quad \mathbf{x}(t) = \big(x_1(t), \ \ldots \ , \ x_p(t)\big). \end{cases}$$

That is, machine $i$ only performs its update operator at its active clocks. The local full model $\mathbf{x}^i(t)$ assembles all components from other machines, and is possibly a delayed version of the global model $\mathbf{x}(t)$, which assembles the most up-to-date component in each machine. Note that the global model is introduced for our analysis, and is not accessible in a real implementation. More specifically, $\tau_j^i(t) \le t$ models the communication delay among machines: when machine $i$ conducts its $t$-th update it only has access to $x_j(\tau_j^i(t))$, a delayed version of the component $x_j(t)$ on the $j$-th machine. We refer to the above algorithm as m-PAPG (for **m**odel parallel, **P**artially **A**synchronous, **P**roximal **G**radient).

In a practical distributed system, communication among machines is much slower than local computations, and the performance of a *synchronous* system is often bottlenecked at the *slowest* machine, due to the need of synchronization in every step. The delays $\tau_j^i(t)$ and active clocks $T_i$ that we introduced in m-PAPG aim to address such issues. For our convergence proofs, we need the following assumptions:

ASSUMPTION 4 (Bounded Delay). $\exists s \in \mathbb{N}, \ \forall i, \forall j, \forall t, \ 0 \le t - \tau_j^i(t) \le s, \ \tau_i^i(t) \equiv t.$

ASSUMPTION 5 (Frequent Update). $\exists s \in \mathbb{N}, \ \forall i, \forall t, T_i \cap \{t, t+1, \cdots, t+s\} \ne \emptyset.$

Intuitively, Assumption 4 guarantees the information that machine $i$ gathered from other machines at the $t$-th iteration are not too obsolete (bounded by at most $s$ clocks apart). The assumption $\tau_i^i(t) \equiv t$ is natural since the $i$-th worker machine is maintaining $x_i$ hence would always have the latest copy. Assumption 5 requires each machine to update at least once in every $s+1$ iterations, for otherwise some component $x_i$ may not be updated at all. We remark that Assumption 4 and Assumption 5 are

very natural and have been widely adopted in previous works [5, 8, 11, 17, 32]. Clearly, when $s = 0$ (i.e., no delay), m-PAPG reduces to the fully synchronous, model parallel proximal gradient algorithm.

Before closing this section, we provide a technical tool to control the inconsistency between the local models $\mathbf{x}^i(t)$ and the global model $\mathbf{x}(t)$. Recall that $(t)_+ = \max\{t, 0\}$ is the positive part of $t$.

LEMMA 4. *Let Assumption 4 hold, then the global model* $\mathbf{x}(t)$ *and the local models* $\{\mathbf{x}^i(t)\}_{i=1}^p$ *satisfy:*

$$(14) \qquad \forall i = 1, \cdots, p, \quad \|\mathbf{x}(t) - \mathbf{x}^i(t)\| \leq \sum_{k=(t-s)_+}^{t-1} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|$$

$$(15) \qquad \|\mathbf{x}^i(t+1) - \mathbf{x}^i(t)\| \leq \sum_{k=(t-s)_+}^{t} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|.$$

*Proof.* Indeed, by the definitions in (m-PAPG):

$$\|\mathbf{x}(t) - \mathbf{x}^i(t)\|^2 = \sum_{j=1}^p \|x_j(t) - x_j(\tau_j^i(t))\|^2$$

$$\leq \sum_{j=1}^p \left( \sum_{k=\tau_j^i(t)}^{t-1} \|x_j(k+1) - x_j(k)\| \right)^2$$

$$\leq \sum_{j=1}^p \left( \sum_{k=(t-s)_+}^{t-1} \|x_j(k+1) - x_j(k)\| \right)^2$$

$$= \sum_{j=1}^p \sum_{k=(t-s)_+}^{t-1} \sum_{k'=(t-s)_+}^{t-1} \|x_j(k+1) - x_j(k)\| \|x_j(k'+1) - x_j(k')\|$$

$$= \sum_{k=(t-s)_+}^{t-1} \sum_{k'=(t-s)_+}^{t-1} \sum_{j=1}^p \|x_j(k+1) - x_j(k)\| \|x_j(k'+1) - x_j(k')\|$$

$$\leq \sum_{k=(t-s)_+}^{t-1} \sum_{k'=(t-s)_+}^{t-1} \|\mathbf{x}(k+1) - \mathbf{x}(k)\| \|\mathbf{x}(k'+1) - \mathbf{x}(k')\|$$

$$= \left( \sum_{k=(t-s)_+}^{t-1} \|\mathbf{x}(k+1) - \mathbf{x}(k)\| \right)^2,$$

where the first inequality is due to the triangle inequality; the second inequality is due to Assumption 4; and the last inequality follows from the Cauchy-Schwarz inequality.

Similarly,

$$\|\mathbf{x}^i(t) - \mathbf{x}^i(t+1)\|^2 = \sum_{j=1}^p \|x_j(\tau_j^i(t)) - x_j(\tau_j^i(t+1))\|^2$$

$$\leq \sum_{j=1}^p \left( \sum_{k=\tau_j^i(t)}^{\tau_j^i(t+1)-1} \|x_j(k+1) - x_j(k)\| \right)^2$$

$$\leq \sum_{j=1}^{p} \left( \sum_{k=(t-s)_+}^{t} \|x_j(k+1) - x_j(k)\| \right)^2,$$

and the rest of the proof is completely similar to the previous case.    □

**4. Characterizing the limit points.** In this section, we characterize the convergence property of the sequences generated by m-PAPG under very general conditions. Recall from Assumption 2 that $\nabla f$ is $L$-Lipschitz continuous. Our first result is as follows:

THEOREM 5. *Let Assumptions 1 to 5 hold. If the step size $\eta \in \left( 0, \frac{1}{L(1+2\sqrt{p}s)} \right)$, then the sequence generated by m-PAPG is square summable, i.e.*

$$(16) \qquad \sum_{t=0}^{\infty} \|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2 < \infty.$$

*In particular, $\lim_{t \to \infty} \|\mathbf{x}(t+1) - \mathbf{x}(t)\| = 0$ and $\lim_{t \to \infty} \|\mathbf{x}(t) - \mathbf{x}^i(t)\| = 0$.*

REMARK 1. *Our bound on the step size $\eta$ is natural: If $s = 0$, i.e., there is no asynchronism then we recover the standard step size rule $\eta < 1/L$ (we can increase $\eta$ by another factor of 2, had convexity on $g$ been assumed). As staleness $s$ increases, we need a smaller step size to "damp" the system to still ensure convergence. The factor $\sqrt{p}$ is another measurement of the degree of "dependency" among worker machines: Indeed, we can reduce $\sqrt{p}$ to $\sqrt{\sum_i L_i^2}/L$, where $L_i$ is the Lipschitz constant of $\nabla_i f$ (cf. (21)).*

*Proof.* The last claim follows immediately from (16) and (14), so we only need to prove (16).

Consider machine $i$ and any $t \in T_i$. Combining (13) with (m-PAPG) gives

$$(17) \qquad x_i(t+1) = \mathrm{prox}_{g_i}^{\eta} \big( x_i(t) - \eta \nabla_i f(\mathbf{x}^i(t)) \big).$$

Then, from Definition 3 of the proximal map we have for all $z \in \mathbb{R}^{d_i}$:

$$(18) \qquad g_i\big(x_i(t+1)\big) + \frac{1}{2\eta} \|x_i(t+1) - x_i(t) + \eta \nabla_i f\big(\mathbf{x}^i(t)\big)\|^2$$

$$\leq g_i(z) + \frac{1}{2\eta} \left\| z - x_i(t) + \eta \nabla_i f\big(\mathbf{x}^i(t)\big) \right\|^2.$$

Set $z = x_i(t)$ and simplify, we obtain:

$$(19) \quad g_i\big(x_i(t+1)\big) - g_i\big(x_i(t)\big)$$

$$\leq -\frac{1}{2\eta} \|x_i(t+1) - x_i(t)\|^2 - \big\langle \nabla_i f\big(\mathbf{x}^i(t)\big), x_i(t+1) - x_i(t) \big\rangle.$$

Note that if $t \notin T_i$, then $x_i(t+1) = x_i(t)$ and (19) still holds. On the other hand, Assumption 2 implies that for all $t$ (cf. (11)):

$$(20) \quad f\big(\mathbf{x}(t+1)\big) - f\big(\mathbf{x}(t)\big) \leq \big\langle \mathbf{x}(t+1) - \mathbf{x}(t), \nabla f\big(\mathbf{x}(t)\big) \big\rangle + \frac{L}{2} \|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2.$$

Adding up (20) and (19) (for all $i$) and recall $F = f + \sum_i g_i$, we have

$$F\big(\mathbf{x}(t+1)\big) - F\big(\mathbf{x}(t)\big) - \tfrac{1}{2}(L - 1/\eta)\|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2$$

$$\leq \sum_{i=1}^{p} \left\langle x_i(t+1) - x_i(t), \nabla_i f(\mathbf{x}(t)) - \nabla_i f(\mathbf{x}^i(t)) \right\rangle$$

$$\leq \sum_{i=1}^{p} \|x_i(t+1) - x_i(t)\| \cdot \|\nabla_i f(\mathbf{x}(t)) - \nabla_i f(\mathbf{x}^i(t))\|$$

(21)
$$\overset{(i)}{\leq} \sum_{i=1}^{p} \|x_i(t+1) - x_i(t)\| \cdot L \|\mathbf{x}(t) - \mathbf{x}^i(t)\|$$

$$\overset{(ii)}{\leq} L \cdot \sum_{i=1}^{p} \|x_i(t+1) - x_i(t)\| \cdot \sum_{k=(t-s)_+}^{t-1} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|$$

(22)
$$\overset{(iii)}{\leq} \sqrt{p}L \|\mathbf{x}(t+1) - \mathbf{x}(t)\| \cdot \sum_{k=(t-s)_+}^{t-1} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|$$

$$\overset{(iv)}{\leq} \frac{\sqrt{p}L}{2} \sum_{k=(t-s)_+}^{t-1} \left[ \|\mathbf{x}(k+1) - \mathbf{x}(k)\|^2 + \|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2 \right]$$

(23)
$$\leq \frac{\sqrt{p}Ls}{2} \|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2 + \frac{\sqrt{p}L}{2} \sum_{k=(t-s)_+}^{t-1} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|^2,$$

where (i) is due to the $L$-Lipschitz continuity of $\nabla f$, (ii) follows from (14), (iii) is the Cauchy-Schwarz inequality, and (iv) follows from the elementary inequality $ab \leq \frac{a^2+b^2}{2}$. Summing the above inequality over $t$ from 0 to $n-1$ and rearranging we obtain

$$F(\mathbf{x}(n)) - F(\mathbf{x}(0)) \leq \frac{1}{2}(L + \sqrt{p}Ls - 1/\eta) \sum_{t=0}^{n-1} \|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2$$

$$+ \frac{L}{2} \sum_{t=0}^{n-1} \sum_{k=(t-s)_+}^{t-1} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|^2$$

$$\leq \frac{1}{2}(L + 2\sqrt{p}Ls - 1/\eta) \sum_{t=0}^{n-1} \|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2.$$

Therefore, if we choose $0 < \eta < \frac{1}{L(1+2\sqrt{p}s)}$, then let $n \to \infty$ we deduce

(24)
$$\sum_{t=0}^{\infty} \|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2 \leq \frac{2}{1/\eta - L - 2\sqrt{p}Ls} [F(\mathbf{x}(0)) - \inf_{\mathbf{z}} F(\mathbf{z})].$$

By Assumption 1, $F$ is bounded from below, hence the right-hand side is finite.  □

The first assertion of the above theorem states that the global sequence $\mathbf{x}(t)$ has square summable successive differences, while the second assertion implies that both the successive difference of the global sequence and the inconsistency between the local sequences and the global sequence diminish as the number of iterations grows. These two conclusions provide a preliminary stability guarantee for m-PAPG.

Next, we prove that the limit points (if exist) of the sequences $\mathbf{x}(t)$ and $\mathbf{x}^i(t), i = 1, \ldots, p$ coincide, and they are critical points of $F$. Again, no convexity assumption is imposed on either $f$ or $g$.

THEOREM 6. *Consider the same setting as in Theorem 5. Then, the sequences* $\{\mathbf{x}(t)\}$ *and* $\{\mathbf{x}^i(t)\}, i = 1, \ldots, p$, *generated by* m-PAPG *share the same set of limit points, which is a subset of* $\text{crit} \, F$.

*Proof.* It is clear from Theorem 5 that $\{\mathbf{x}(t)\}$ and $\{\mathbf{x}^i(t)\}, i = 1, \ldots, p$, share the same set of limit points, and we need to show that any limit point of $\{\mathbf{x}(t)\}$ is also a critical point of $F$.

Let $\mathbf{x}^*$ be a limit point of $\{\mathbf{x}(t)\}$. By Definition 1 it suffices to exhibit a sequence $\mathbf{x}(k)$ satisfying[1]

$$(25) \qquad \mathbf{x}(k) \to \mathbf{x}^*, \ F(\mathbf{x}(k)) \to F(\mathbf{x}^*), \ \mathbf{0} \leftarrow \mathbf{u}(k) \in \partial F(\mathbf{x}(k)).$$

Let us first construct the subgradient sequence $\mathbf{u}(k)$. Consider machine $i$ and any $\hat{t} \in T_i$, the optimality condition of (17) gives

$$(26) \qquad u_i(\hat{t}+1) := -\tfrac{1}{\eta} \left[ x_i(\hat{t}+1) - x_i(\hat{t}) + \eta \nabla_i f(\mathbf{x}^i(\hat{t})) \right] \in \partial g_i(x_i(\hat{t}+1)).$$

It then follows that

$$\|u_i(\hat{t}+1) + \nabla_i f(\mathbf{x}(\hat{t}+1))\|$$
$$\leq \|u_i(\hat{t}+1) + \nabla_i f(\mathbf{x}(\hat{t}))\| + \|\nabla_i f(\mathbf{x}(\hat{t}+1)) - \nabla_i f(\mathbf{x}(\hat{t}))\|$$
$$\overset{(i)}{\leq} \left\| \tfrac{1}{\eta} \left[ x_i(\hat{t}+1) - x_i(\hat{t}) \right] + \nabla_i f(\mathbf{x}^i(\hat{t})) - \nabla_i f(\mathbf{x}(\hat{t})) \right\| + L\|\mathbf{x}(\hat{t}+1) - \mathbf{x}(\hat{t})\|$$
$$\overset{(ii)}{\leq} \tfrac{1}{\eta}\|x_i(\hat{t}+1) - x_i(\hat{t})\| + L\|\mathbf{x}^i(\hat{t}) - \mathbf{x}(\hat{t})\| + L\|\mathbf{x}(\hat{t}+1) - \mathbf{x}(\hat{t})\|$$
$$(27) \qquad \overset{(iii)}{\leq} \tfrac{1}{\eta}\|x_i(\hat{t}+1) - x_i(\hat{t})\| + L \sum_{k=(\hat{t}-s)_+}^{\hat{t}} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|,$$

where (i) and (ii) are due to the $L$-Lipschitz continuity of $\nabla f$, and (iii) follows from (14). Next, consider any other $t \notin T_i$ and $t \geq s$, we denote $\hat{t}$ as the *largest* element in the set $\{k \leq t : k \in T_i\}$. By Assumption 5 $\hat{t}$ always exists and $t - \hat{t} \leq s$. Since no update is performed on machine $i$ at any clock in $[\hat{t}+1, t]$, we have $x_i(t+1) = x_i(\hat{t}+1)$. Thus, we can choose $u_i(t+1) = u_i(\hat{t}+1) \in \partial g_i(x_i(\hat{t}+1)) = \partial g_i(x_i(t+1))$, and obtain

$$(28) \qquad \|u_i(t+1) + \nabla_i f(\mathbf{x}(t+1)) - u_i(\hat{t}+1) - \nabla_i f(\mathbf{x}(\hat{t}+1))\|$$
$$= \|\nabla_i f(\mathbf{x}(t+1)) - \nabla_i f(\mathbf{x}(\hat{t}+1))\|$$
$$\leq \sum_{k=\hat{t}+1}^{t} \|\nabla_i f(\mathbf{x}(k+1)) - \nabla_i f(\mathbf{x}(k))\|$$
$$\leq \sum_{k=(t-s+1)_+}^{t} \|\nabla_i f(\mathbf{x}(k+1)) - \nabla_i f(\mathbf{x}(k))\|$$
$$(29) \qquad \leq \sum_{k=(t-s+1)_+}^{t} L\|\mathbf{x}(k+1) - \mathbf{x}(k)\|.$$

---

[1] Technically, from Definition 1 we should have the Frechét subdifferential $\hat{\partial} F$ in (25), however, a standard argument allows us to use the more convenient subdifferential [30, Proposition 8.7].

Combining the two cases in (27) and (29) we have for all $t$:

$$(30) \qquad \|\mathbf{u}(t+1) + \nabla f(\mathbf{x}(t+1))\| \leq (\sqrt{p}/\eta + 2L) \sum_{k=(t-2s)_+}^{t} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|,$$

where $\mathbf{u}(t+1) = \big(u_1(t+1), \ \ldots \ , \ u_p(t+1)\big) \in \partial g(\mathbf{x}(t+1))$, and the factor $\sqrt{p} \geq 1$ is artificially introduced for the convenience of subsequent analysis. Therefore, by (30) and Theorem 5 we deduce

$$(31) \qquad \lim_{t\to\infty} \mathrm{dist}_{\partial F(\mathbf{x}(t+1))}(\mathbf{0}) \leq \lim_{t\to\infty} \|\mathbf{u}(t+1) + \nabla f(\mathbf{x}(t+1))\| = 0.$$

Recall that $\mathbf{x}^*$ is a limit point of $\{\mathbf{x}(t)\}$, thus there exists a subsequence $\mathbf{x}(t_m) \to \mathbf{x}^*$. Next we verify the function value convergence in (25). The challenge here is that the component function $g$ is only closed, hence may not be continuous. For any $t \in T_i$, applying (18) with $z = x_i^*$ and rearranging gives

$$g_i(x_i(t+1)) \leq g_i\big(x_i^*\big) + \frac{1}{2\eta}\|x_i^* - x_i(t)\|^2 - \frac{1}{2\eta}\|x_i(t+1) - x_i(t)\|^2$$
$$\qquad\qquad + \langle x_i^* - x_i(t+1), \nabla_i f(\mathbf{x}^i(t))\rangle$$
$$(32) \qquad = g_i\big(x_i^*\big) + \frac{1}{2\eta}\|x_i^* - x_i(t)\|^2 - \frac{1}{2\eta}\|x_i(t+1) - x_i(t)\|^2$$
$$\qquad\qquad + \langle x_i^* - x_i(t+1), \nabla_i f(\mathbf{x}^*)\rangle + \langle x_i^* - x_i(t+1), \nabla_i f(\mathbf{x}^i(t)) - \nabla_i f(\mathbf{x}^*)\rangle.$$

Observe from Theorem 5 that

$$\lim_{m\to\infty} \max_{t\in[t_m-s,t_m+s]} \|\mathbf{x}(t+1) - \mathbf{x}^*\| = 0, \qquad \lim_{m\to\infty} \max_{t\in[t_m-s,t_m+s]} \|\mathbf{x}^i(t+1) - \mathbf{x}^*\| = 0.$$

By Assumption 5, $[t_m-s, t_m+s]\cap T_i \neq \emptyset$ for all $i$. Then using the Lipschitz continuity of $\nabla f$ we deduce from (32) that

$$(33) \qquad \limsup_{m\to\infty} \max_{t\in[t_m-s,t_m+s]\cap T_i} g_i(x_i(t+1)) \leq g_i(x_i^*).$$

Since each machine updates at least once during $[t_m - s, t_m]$, let $\hat{t}_m$ be the largest element of $[t_m - s, t_m]\cap T_i$. Also notice that each machine $i$ only updates at its active clocks $T_i$, it then follows that

$$\max_{t\in[t_m,t_m+s]} g_i(x_i(t+1)) = \max_{t\in[\hat{t}_m,t_m+s]\cap T_i} g_i(x_i(t+1)) \leq \max_{t\in[t_m-s,t_m+s]\cap T_i} g_i(x_i(t+1)),$$

and hence by (33)

$$(34) \qquad \limsup_{m\to\infty} \max_{t\in[t_m,t_m+s]} g_i(x_i(t+1)) \leq g_i(x_i^*).$$

To complete the proof, choose any $k_m \in [t_m, t_m+s]$. Since $\mathbf{x}(t_m) \to \mathbf{x}^*$, Theorem 5 implies that

$$(35) \qquad\qquad\qquad\qquad \mathbf{x}(k_m) \to \mathbf{x}^*.$$

From (34) we know for all $i$, $\limsup_{m\to\infty} g_i(x_i(k_m)) \leq g_i(x_i^*)$. On the other hand, it follows from the closedness of the function $g_i$ (cf. Assumption 3) that $\liminf_{m\to\infty} g_i(x_i(k_m)) \geq g_i(x_i^*)$, thus in fact $\lim_{m\to\infty} g_i(x_i(k_m)) = g_i(x_i^*)$. Since $f$ is continuous, we know

$$(36) \qquad \lim_{m\to\infty} F(\mathbf{x}(k_m)) = \lim_{m\to\infty} f(\mathbf{x}(k_m)) + \sum_i g_i(x_i(k_m)) = F(\mathbf{x}^*).$$

Combining (31), (35) and (36) we know from Definition 1 that $\mathbf{x}^* \in \mathrm{crit}\, F$.            □

Theorem 6 further justifies m-PAPG by showing that any limit point it produces is necessarily a critical point. Of course, for convex functions any critical point is a global minimizer. The closest result to Theorem 5 and Theorem 6 we are aware of is [8, Proposition 7.5.3], where essentially the same conclusion was reached but under the much more restrictive assumption that $g$ is an indicator function of a product *convex* set. Thus, our result is new even when $g$ is a convex function such as the $\ell_1$ norm that is widely used to promote sparsity. Furthermore, we allow $g$ to be any closed separable function (convex or not), covering the many recent nonconvex regularization functions in machine learning and statistics (see e.g. [15, 26, 38, 39]). We also note that the proof of Theorem 6 (for nonconvex $g$) involves significantly new ideas beyond those of [8].

We note that the existence of limit points can be guaranteed, for instance, if $\{\mathbf{x}(t)\}$ is bounded or the sublevel set $\{\mathbf{x} \mid F(\mathbf{x}) \leq \alpha\}$ is bounded for all $\alpha \in \mathbb{R}$. However, we have yet to prove that the sequence $\{\mathbf{x}(t)\}$ generated by m-PAPG does converge to one of the critical points, and we fill this gap under two complementary sets of assumptions on the objective function in Sections 5 and 6, respectively.

**5. Convergence under Error Bound.** In this section we prove that the global sequence $\{\mathbf{x}(t)\}$ produced by m-PAPG converges periodically linearly to a global minimizer, by assuming an error bound condition on the objective function in (P) and a convexity assumption that serves to simplify the presentation:

ASSUMPTION 6 (Convex).   *The functions $f$ and $g$ in (P) are convex.*

Note that for convex functions $g$ the proximal mapping $\mathrm{prox}_g^\eta$ is single valued for any $\eta > 0$. The error bound condition we need is as follows:

ASSUMPTION 7 (Error Bound).   *For every $\alpha > 0$, there exist $\delta, \kappa > 0$ such that for all $\mathbf{x}$ with $f(\mathbf{x}) \leq \alpha$ and $\|\mathbf{x} - \mathrm{prox}_g(\mathbf{x} - \nabla f(\mathbf{x}))\| \leq \delta$,*

$$(37) \qquad \mathrm{dist}_{\mathrm{crit}\,F}(\mathbf{x}) \leq \kappa \|\mathbf{x} - \mathrm{prox}_g(\mathbf{x} - \nabla f(\mathbf{x}))\|,$$

*where recall that* $\mathrm{crit}\,F$ *is the set of critical points of* $F$.

Equation (37) is a proximal extension of the Luo-Tseng error bound [25] where $g$ is the indicator function of a closed convex set. A prototypic convex function $F$ satisfying (37) is the following:

$$(38) \qquad F(\mathbf{x}) = f(A\mathbf{x}) + g(\mathbf{x}),$$

where $f$ is strongly convex (i.e., $f - \frac{\mu}{2}\|\cdot\|^2$ is convex for some $\mu > 0$), $A$ is a linear map, and $g$ is either an indicator function of a convex set [25] or the $\ell_p$ norm for $p \in [1, 2] \cup \{\infty\}$ [42]. Many machine learning formulations such as Lasso and sparse logistic regression fit into this form. In fact, for convex functions $F$ taking such form, the error bound condition in (37) is recently shown to be equivalent to the following conditions [14, 40]:

$$\text{Restricted strong convexity}: \ \langle \mathbf{x} - \mathrm{prox}_g(\mathbf{x}), \mathbf{x} - \mathrm{proj}_{\mathrm{crit}\,F}(\mathbf{x}) \rangle \geq \mu \cdot \mathrm{dist}^2_{\mathrm{crit}\,F}(\mathbf{x}),$$
$$\text{Quadratic growth}: \ F(\mathbf{x}) - F^* \geq \mu \cdot \mathrm{dist}^2_{\mathrm{crit}\,F}(\mathbf{x}),$$

where $F^*$ is the minimum value of $F$ and $\mu > 0$ is a constant. In general, the error bound condition in (37) is not exclusive to convex functions. For instance, it holds for $f(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|^2$ and any function $g$ that has a unique global minimizer at 0 (such

as the cardinality function $g(\mathbf{x}) = \|\mathbf{x}\|_0$. However, it is often quite challenging to establish the error bound condition for a large family of nonconvex functions.

We define the following nonnegative quantities that measure the progress of m-PAPG:

$$(39) \qquad\qquad A(t) := F(\mathbf{x}(t)) - F^*, \ F^* := \inf_{\mathbf{x}} F(\mathbf{x}),$$

$$(40) \qquad\qquad B(t) := \sum_{k=(t-s-1)_+}^{t-1} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|^2,$$

In the following key lemma we relate the gap quantities defined above inductively.

LEMMA 7. *Let Assumptions 1 to 7 hold. Then, we have*

$$A(t+s+1) \le A(t) - \tfrac{1}{2}(\tfrac{1}{\eta} - L - 2sL\sqrt{p})B(t+s+1) + \tfrac{1}{2}sL\sqrt{p}B(t)$$
$$0 \le A(t+s+1) \le a_\eta B(t+s+1) + bB(t),$$

*where $a_\eta$ and $b$ are given in (50) below.*

*Proof.* The first inequality is obtained by summing the inequality (23) over $t, t + 1, \cdots, t + s$. So we need only prove the second inequality.

Let us introduce some notations to simplify the proof. For each machine $i$ let $t_i$ be the largest clock in $[t, t + s] \cap T_i$, and denote

$$(41) \qquad \mathbf{z} = \big(x_1(t_1), \ \ldots \ , x_p(t_p)\big)$$
$$(42) \qquad \mathbf{z}^+ = \big(x_1(t_1 + 1), \ \ldots \ , x_p(t_p + 1)\big) = \big(x_1(t+s+1), \ \ldots \ , x_p(t+s+1)\big),$$

where the last equality is due to the maximality of each $t_i$. From the optimality condition of the proximal map $z_i^+ = \mathrm{prox}_{g_i}^\eta(z_i - \eta\nabla_i f(\mathbf{x}^i(t_i)))$ we deduce

$$(43) \qquad\qquad \eta^{-1}(z_i - z_i^+) - \nabla_i f(\mathbf{x}^i(t_i)) \in \partial g_i(z_i^+).$$

Since the gradient of $f$ is $L$-Lipschitz continuous and the function $g$ is convex, we obtain

$$f(\mathbf{z}^+) - f(\bar{\mathbf{z}}) \le \sum_{i=1}^p \langle z_i^+ - \bar{z}_i, \nabla_i f(\bar{\mathbf{z}})\rangle + \frac{L}{2}\|\mathbf{z}^+ - \bar{\mathbf{z}}\|^2,$$

$$g(\mathbf{z}^+) - g(\bar{\mathbf{z}}) \le \sum_{i=1}^p \langle z_i^+ - \bar{z}_i, \eta^{-1}(z_i - z_i^+) - \nabla_i f(\mathbf{x}^i(t_i))\rangle,$$

where we define $\bar{\mathbf{z}} := \mathrm{proj}_{\mathrm{crit}\,F}(\mathbf{z})$, i.e., the projection of $\mathbf{z}$ onto the set of critical points of $F$, and the last inequality follows from (43). Adding up the above two inequalities we obtain

$$F(\mathbf{z}^+) - F^* - \tfrac{L}{2}\|\mathbf{z}^+ - \bar{\mathbf{z}}\|^2 \le \sum_{i=1}^p \langle z_i^+ - \bar{z}_i, \nabla_i f(\bar{\mathbf{z}}) + \eta^{-1}(z_i - z_i^+) - \nabla_i f(\mathbf{x}^i(t_i))\rangle$$

$$\overset{(i)}{\le} \sum_{i=1}^p [\|z_i^+ - z_i\| + \|z_i - \bar{z}_i\|][\|\nabla_i f(\mathbf{x}^i(t_i)) - \nabla_i f(\bar{\mathbf{z}})\| + \eta^{-1}\|z_i - z_i^+\|]$$

$$\overset{(ii)}{\le} \sum_{i=1}^p 4\left[\|z_i^+ - z_i\|^2 + \|z_i - \bar{z}_i\|^2 + \eta^{-2}\|z_i^+ - z_i\|^2 + \|\nabla_i f(\mathbf{x}^i(t_i)) - \nabla_i f(\bar{\mathbf{z}})\|^2\right]$$

$$\leq 4\left[\|\bar{\mathbf{z}}-\mathbf{z}\|^2 + (1+\eta^{-2})\|\mathbf{z}^+ - \mathbf{z}\|^2 + \sum_{i=1}^{p}L^2\|\mathbf{x}^i(t_i)-\bar{\mathbf{z}}\|^2\right],$$

where (i) is due to the Cauchy-Schwarz inequality and the triangle inequality, (ii) is due to the elementary inequality $(a+b)(c+d) \leq 4(a^2+b^2+c^2+d^2)$, and the last inequality is due to the $L$-Lipschitz continuity of $\nabla f$. Using again the triangle inequality we obtain from the above inequality that

$$F(\mathbf{z}^+) - F^* \leq (L+4)\|\bar{\mathbf{z}}-\mathbf{z}\|^2 + (L+4+\tfrac{4}{\eta^2})\|\mathbf{z}^+ - \mathbf{z}\|^2 + 4L^2\sum_{i=1}^{p}\|\mathbf{x}^i(t_i)-\mathbf{z}\|^2$$

$$\overset{(i)}{=} (L+4)\|\bar{\mathbf{z}}-\mathbf{z}\|^2 + \sum_{i=1}^{p}[(L+4+\tfrac{4}{\eta^2})\|x_i(t_i+1)-x_i(t_i)\|^2 + 4L^2\|\mathbf{x}^i(t_i)-\mathbf{z}\|^2],$$

$$(44) \qquad \overset{(ii)}{\leq} (L+4)\|\bar{\mathbf{z}}-\mathbf{z}\|^2 + (L+4+\tfrac{4}{\eta^2})B(t+s+1) + 4L^2\sum_{i=1}^{p}\|\mathbf{x}^i(t_i)-\mathbf{z}\|^2,$$

where (i) is due to our definition of $\mathbf{z}$ and $\mathbf{z}^+$ in (41) and (42), and (ii) is due to the fact that $t_i \in [t, t+s]$ for all $i$.

We next bound the terms $\|\bar{\mathbf{z}}-\mathbf{z}\|^2$ and $\|\mathbf{x}^i(t_i)-\mathbf{z}\|^2$. First, note that

$$\|x_i(t_i+1)-x_i(t_i)\| = \|\mathrm{prox}^{\eta}_{g_i}(x_i(t_i)-\eta\nabla_i f(\mathbf{x}^i(t_i))) - x_i(t_i)\|$$

$$\geq \|\mathrm{prox}^{\eta}_{g_i}(x_i(t_i)-\eta\nabla_i f(\mathbf{z})) - x_i(t_i)\|$$

$$\quad - \|\mathrm{prox}^{\eta}_{g_i}(x_i(t_i)-\eta\nabla_i f(\mathbf{x}^i(t_i))) - \mathrm{prox}^{\eta}_{g_i}(x_i(t_i)-\eta\nabla_i f(\mathbf{z}))\|$$

$$\overset{(i)}{\geq} \|\mathrm{prox}^{\eta}_{g_i}(x_i(t_i)-\eta\nabla_i f(\mathbf{z})) - x_i(t_i)\| - \eta L\|\mathbf{z}-\mathbf{x}^i(t_i)\|,$$

where (i) follows from the non-expansiveness of $\mathrm{prox}^{\eta}_g$ (recall that $g$ is convex) and the $L$-Lipschitz continuity of $\nabla f$. Rearranging the above inequality and summing over all $i$, we obtain

$$\|\mathrm{prox}^{\eta}_g(\mathbf{z}-\eta\nabla f(\mathbf{z})) - \mathbf{z}\|^2 \leq \sum_{i=1}^{p}\left[\|x_i(t_i+1)-x_i(t_i)\| + \eta L\|\mathbf{z}-\mathbf{x}^i(t_i)\|\right]^2$$

$$(45) \qquad \leq 2\sum_{i=1}^{p}\left[\|x_i(t_i+1)-x_i(t_i)\|^2 + \eta^2 L^2\|\mathbf{z}-\mathbf{x}^i(t_i)\|^2\right].$$

The last term $\|\mathbf{z}-\mathbf{x}^i(t_i)\|^2$ can be further bounded as follows:

$$\|\mathbf{z}-\mathbf{x}^i(t_i)\|^2 = \sum_{j=1}^{p}\|x_j(t_j)-x_j(\tau^i_j(t_i))\|^2$$

$$= \sum_{j=1}^{p}\left\|\sum_{k=\min\{t_j,\tau^i_j(t_i)\}}^{\max\{t_j,\tau^i_j(t_i)\}-1}x_j(k+1)-x_j(k)\right\|^2$$

$$\leq \sum_{j=1}^{p}\left[\sum_{k=\min\{t_j,\tau^i_j(t_i)\}}^{\max\{t_j,\tau^i_j(t_i)\}-1}\|x_j(k+1)-x_j(k)\|\right]^2$$

$$\overset{(i)}{\leq} \sum_{j=1}^{p}2s\sum_{k=t-s}^{t+s-1}\|x_j(k+1)-x_j(k)\|^2$$

$$= 2s \sum_{k=t-s}^{t+s-1} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|^2$$

$$(46) \qquad\qquad \leq 2s[B(t) + B(t+s+1)],$$

where (i) is due to the fact that $t_j \in [t, t+s]$ and $\tau_j^i(t_i) \in [t-s, t+s]$. Combining (45) and (46) we obtain

$$(47) \quad \|\mathrm{prox}_g^\eta(\mathbf{z} - \eta\nabla f(\mathbf{z})) - \mathbf{z}\|^2 \leq 2B(t+s+1) + 4ps\eta^2 L^2[B(t) + B(t+s+1)].$$

Thanks to Theorem 5, we know for $t$ sufficiently large, $\|\mathrm{prox}_g^\eta(\mathbf{z} - \eta\nabla f(\mathbf{z})) - \mathbf{z}\| \leq \eta\delta$. Since the function $\eta \mapsto \frac{1}{\eta}\|\mathrm{prox}_g^\eta(\mathbf{z} - \eta\nabla f(\mathbf{z})) - \mathbf{z}\|$ is monotonically decreasing [31], we can apply the error bound condition in Assumption 7 for $\eta < 1$ and $t$ sufficiently large, and obtain

$$(48) \qquad \|\bar{\mathbf{z}} - \mathbf{z}\|^2 \leq \kappa\|\mathbf{z} - \mathrm{prox}_g(\mathbf{z} - \nabla f(\mathbf{z}))\|^2 \leq \kappa\eta^{-2}\|\mathbf{z} - \mathrm{prox}_g^\eta(\mathbf{z} - \eta\nabla f(\mathbf{z}))\|^2.$$

Finally, combining (44), (46), (47) and (48) we arrive at:

$$F(\mathbf{x}(t+s+1)) - F^* = F(\mathbf{z}^+) - F^*$$

$$\leq (L+4)\|\bar{\mathbf{z}} - \mathbf{z}\|^2 + (L+4+\tfrac{4}{\eta^2})B(t+s+1) + 4L^2\sum_{i=1}^{p}\|\mathbf{x}^i(t_i) - \mathbf{z}\|^2,$$

$$(49) \qquad\qquad \leq a_\eta B(t+s+1) + bB(t),$$

where the coefficients are

$$(50) \qquad a_\eta = L + 4 + 8psL^2 + 4ps\kappa L^2(L+4) + \tfrac{2}{\eta^2}(2 + 4\kappa + \kappa L),$$

$$(51) \qquad b = 8psL^2 + 4ps\kappa L^2(L+4). \qquad\qquad\qquad\qquad\qquad \square$$

Lemma 7 improves the analysis of [32] in three aspects: (1) it is shorter and simpler; (2) it allows any convex function $g$; and (3) the leading coefficient for $B(t)$ is reduced from $O(1/\eta)$ to $O(1)$. The two recursive relations in Lemma 7, as shown in [32, Lemma 4.5], easily imply the following convergence guarantee:

THEOREM 8. *Let Assumptions 1 to 7 hold. Then, there exists some $\eta_0 > 0$ such that if $0 < \eta < \eta_0$, then the sequences $\{A(t), B(t)\}$ generated by m-PAPG satisfy for all $r = 0, 1, 2, \cdots$*

$$(52) \qquad A(r(s+1)) \leq C_1(1 - \gamma\eta)^r, \ \ B(r(s+1)) \leq C_2(1 - \gamma\eta)^r,$$

*where $C_1, C_2, \gamma < 1/\eta$ are positive constants.*

Hence, the gaps $A(t)$ and $B(t)$ that measure the progress of m-PAPG decrease by a constant factor $(1 - \gamma\eta)$ for every $s + 1$ steps, which makes intuitive sense since in the worst case each worker machine only performs one update in every $s + 1$ steps. In other words, $(s + 1)$ is the natural time scale for measuring progress here. Note that since $\|\mathbf{x}(t+s+1) - \mathbf{x}(t)\|^2 \leq (s+1)B(t+s+1)$, it follows easily that the global sequence $\mathbf{x}(t)$ and consequently also the local sequences $\{\mathbf{x}^i(t)\}$ all converge to the same limit point in crit $F$ at a $(s+1)$-periodically linear rate.

**6. Convergence with KŁ inequality.** The error bound condition considered in the previous section is not easy to verify in general. It has been discovered recently that the error bound condition is equivalent to other notions in optimization that can be verified in alternative ways [14, 40], see e.g. (38). However, for nonconvex functions, sometimes even the simple ones, it remains a challenging task to verify if the error bound condition holds. This failure motivates us to investigate another property, the Kurdyka-Łojasiewicz (KŁ) inequality, that has been shown to be quite effective in dealing with nonconvex functions.

DEFINITION 9 (KŁ property, [10, Lemma 6]).    *Let $\Omega \subset \mathrm{dom}h$ be a compact set on which the function $h$ is a constant. We say that $h$ satisfies the KŁ property if there exist $\varepsilon, \lambda > 0$ such that for all $\bar{\mathbf{x}} \in \Omega$ and all $\mathbf{x} \in \{\mathbf{z} \in \mathbb{R}^d : \mathrm{dist}_\Omega(\mathbf{z}) < \varepsilon\} \cap [\mathbf{z} : h(\bar{\mathbf{x}}) < h(\mathbf{z}) < h(\bar{\mathbf{x}}) + \lambda]$, it holds that*

$$(53) \qquad \varphi'(h(\mathbf{x}) - h(\bar{\mathbf{x}})) \cdot \mathrm{dist}_{\partial h(\mathbf{x})}(\mathbf{0}) \geq 1,$$

*where the function $\varphi : [0, \lambda) \to \mathbb{R}_+, 0 \mapsto 0$, is continuous, concave, and has continuous and positive derivative $\varphi'$ on $(0, \lambda)$.*

The KŁ inequality in (53) is an important tool to bound the trajectory length of a dynamical system (see [9, 20] and the references therein for some historic developments). It has recently been used to analyze discrete-time algorithms in [1] and proximal algorithms in [3, 4, 10]. As we shall see, the function $\varphi$ will serve as a Lyapunov potential function. Quite conveniently, most practical functions, in particular, the quasi-norm $\|\cdot\|_p$ for positive rational $p$, as well as convex functions with certain growth conditions, are KŁ. For a more detailed discussion of KŁ functions, including many familiar examples, see [10, Section 5] and [4, Section 4].

Following the recipe in [10], we need the following assumption to guarantee the algorithm is making *sufficient* progress:

ASSUMPTION 8 (Sufficient decrease).    *There exists $\alpha > 0$ such that for all large $t$,*

$$(54) \qquad F(\mathbf{x}(t+1)) \leq F(\mathbf{x}(t)) - \alpha\|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2.$$

The sufficient decrease assumption is automatically satisfied in many descent algorithms, e.g., the proximal gradient algorithm. However, in the partially asynchronous parallel (PAP) setting, it is highly nontrivial to satisfy the sufficient decrease assumption because of the complication due to communication delays and update skips. Note also that none of the worker machines actually has access to the global sequence $\mathbf{x}(t)$, so even verifying the sufficient decrease property is not trivial. To simplify the presentation, we first analyze the performance of m-PAPG using the KŁ inequality and taking the sufficient decrease property for granted, and later we we will give some verifiable conditions to justify this simplification.

Our first result in this section strengthens the convergence properties in Theorems 5 and 6 for m-PAPG:

THEOREM 10 (Finite Length).    *Let Assumptions 1 to 5 and 8 hold for m-PAPG, and let $F$ satisfy the KŁ property in Definition 9.   Then, with step size $\eta \in \left(0, \frac{1}{L(1+2\sqrt{\bar{p}}s)}\right)$, every bounded sequence $\{\mathbf{x}(t)\}$ generated by m-PAPG satisfies*

$$(55) \qquad \sum_{t=0}^{\infty} \|\mathbf{x}(t+1) - \mathbf{x}(t)\| < \infty,$$

$$(56) \qquad \forall i = 1, \ldots, p, \ \sum_{t=0}^{\infty} \|\mathbf{x}^i(t+1) - \mathbf{x}^i(t)\| < \infty.$$

*Furthermore, $\{\mathbf{x}(t)\}$ and $\{\mathbf{x}^i(t)\}_{i=1}^{p}$ converge to the same critical point of $F$.*

*Proof.* We first show that (55) implies (56). Indeed, recall from (15):

$$\|\mathbf{x}^i(t+1) - \mathbf{x}^i(t)\| \leq \sum_{k=(t-s)_+}^{t} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|.$$

Therefore, summing for $t = 0, 1, \cdots, n$ gives

$$\sum_{t=0}^{n} \|\mathbf{x}^i(t+1) - \mathbf{x}^i(t)\| \leq \sum_{t=0}^{n} \sum_{k=(t-s)_+}^{t} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|$$

$$\leq (2s+1) \sum_{t=0}^{n} \|\mathbf{x}(t+1) - \mathbf{x}(t)\|.$$

The claim then follows by letting $n$ tend to infinity.

By Theorem 5, the limit points of $\{\mathbf{x}(t)\}$ and $\{\mathbf{x}^i(t)\}_{i=1}^{p}$ coincide and are critical points of $F$. Thus, the only thing left to prove is the finite length property in (55). By Assumption 8 and Assumption 1, the objective value $F(\mathbf{x}(t))$ decreases to a finite limit $F^*$. Since $\{\mathbf{x}(t)\}$ is assumed to be bounded, the set of its limit points $\Omega$ is nonempty and compact. Summing (18) over all $i$ and set $\mathbf{z} \in \Omega$, we obtain

$$g(\mathbf{x}(t+1)) \leq g(\mathbf{z}) - \frac{1}{2\eta} \|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2 - \sum_{i=1}^{p} \langle \nabla_i f(\mathbf{x}^i(t)), \mathbf{x}(t+1) - \mathbf{x}(t) \rangle.$$

Note that $\mathbf{x}(t+1) - \mathbf{x}(t) \to 0$. Also, since $\{\mathbf{x}(t)\}$ is bounded and $\mathbf{x}(t) - \mathbf{x}^i(t) \to 0$ for all $i$, $\{\mathbf{x}^i(t)\}_{i=1}^{p}$ are all bounded. we then take limsup on both sides and obtain that $\limsup_{t\to\infty} g(\mathbf{x}(t+1)) \leq g(\mathbf{z})$. Together with the closedness of $g$ we further obtain that $\lim_{t\to\infty} g(\mathbf{x}(t+1)) = g(\mathbf{z})$. Note that $f$ is continuous, we thus conclude that $\lim_{t\to\infty} F(\mathbf{x}(t+1)) = F(\mathbf{z})$ for all $\mathbf{z} \in \Omega$. Note that $F(\mathbf{x}(t)) \downarrow F^*$. Thus for all $\mathbf{x}^* \in \Omega$, we have $F(\mathbf{x}^*) \equiv F^*$. Now fix $\varepsilon > 0$. Since $\Omega$ is compact, for $t$ sufficiently large we have $\mathrm{dist}_\Omega(\mathbf{x}(t)) \leq \varepsilon$. We now have all ingredients to apply the KŁ inequality in Definition 9: for all sufficiently large $t$,

$$(57) \qquad \varphi'\big(F(\mathbf{x}(t)) - F^*\big) \cdot \mathrm{dist}_{\partial F(\mathbf{x}(t))}(\mathbf{0}) \geq 1.$$

Since $\varphi$ is concave, we obtain

$$\Delta_{t,t+1} := \varphi\big(F(\mathbf{x}(t)) - F^*\big) - \varphi\big(F(\mathbf{x}(t+1)) - F^*\big)$$

$$\geq \varphi'\big(F(\mathbf{x}(t)) - F^*\big)\big(F(\mathbf{x}(t)) - F(\mathbf{x}(t+1))\big)$$

$$(58) \qquad \overset{(i)}{\geq} \frac{\alpha \|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2}{\mathrm{dist}_{\partial F(\mathbf{x}(t))}(\mathbf{0})},$$

where (i) follows from Assumption 8 and (57). It is clear that the function $\varphi$ (composed with $F$) serves as a Lyapunov function. Using the elementary inequality $2\sqrt{ab} \leq a + b$ we obtain from (58) that for $t$ sufficiently large,

$$2\|\mathbf{x}(t+1) - \mathbf{x}(t)\| \leq \tfrac{\delta}{\alpha} \Delta_{t,t+1} + \tfrac{1}{\delta} \mathrm{dist}_{\partial F(\mathbf{x}(t))}(\mathbf{0}),$$

where $\delta > 0$ will be specified later. Recalling the bound for $\partial F(\mathbf{x}(t))$ in (30), and summing over $t$ from $m$ (sufficiently large) to $n$ gives:

$$2 \sum_{t=m}^{n} \|\mathbf{x}(t+1) - \mathbf{x}(t)\| \leq \sum_{t=m}^{n} \frac{\delta}{\alpha} \Delta_{t,t+1} + \sum_{t=m}^{n} \frac{1}{\delta} \text{dist}_{\partial F(\mathbf{x}(t))}(\mathbf{0})$$

$$\overset{(i)}{\leq} \frac{\delta}{\alpha} \varphi\big(F(\mathbf{x}(m)) - F^*\big) + \sum_{t=m}^{n} \frac{\sqrt{p}/\eta + 2L}{\delta} \sum_{k=(t-2s)_+}^{t} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|$$

$$\leq \frac{\delta}{\alpha} \varphi\big(F(\mathbf{x}(m)) - F^*\big) + \frac{(2s+1)(\sqrt{p}/\eta + 2L)}{\delta} \sum_{k=(m-2s)_+}^{m-1} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|$$

$$+ \frac{(2s+1)(\sqrt{p}/\eta + 2L)}{\delta} \sum_{t=m}^{n} \|\mathbf{x}(t+1) - \mathbf{x}(t)\|,$$

where (i) is due to (30). Setting $\delta = (2s+1)(\sqrt{p}/\eta + 2L)$ and rearranging gives

$$\sum_{t=m}^{n} \|\mathbf{x}(t+1) - \mathbf{x}(t)\| \leq \frac{(2s+1)(\sqrt{p}/\eta + 2L)}{\alpha} \varphi\big(F(\mathbf{x}(m)) - F^*\big)$$

$$+ \sum_{k=(m-2s)_+}^{m-1} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|.$$

Since the right-hand side is finite, let $n$ tend to infinity completes the proof for (55). □

Compared with (16) in Theorem 5, we now have the successive differences to be absolutely summable (instead of square summable). This is a significantly stronger result as it immediately implies that the whole sequence is Cauchy and hence convergent, whereas we cannot get the same conclusion from the square summable property in Theorem 5. We note that local maxima are excluded from being the limit in Theorem 10, due to Assumption 8. Also, the boundedness assumption on the trajectory $\{\mathbf{x}(t)\}$ is easy to satisfy, for instance, when $F$ has bounded sublevel sets. We refer to [4, Remark 3.3] for more conditions that imply the boundedness condition. Moreover, following similar arguments in [4] we can also determine the local convergence rates of the sequences generated by m-PAPG.

In the remaining part of this section we provide some justifications for the sufficient decrease property in Assumption 8. For simplicity we assume all worker machines perform updates in each time step $t$:

ASSUMPTION 9. $\forall i = 1, \cdots, p, \forall t, t \in T_i$.

Note that Assumption 9 is commonly adopted in the analysis of many recent parallel systems [2, 17, 19, 21, 22, 28]. Put it differently, under Assumption 9 we measure the performance of the system w.r.t. the minimum number of updates among all worker machines whereas under the more relaxed Assumption 5 we measure the performance w.r.t. the total number of updates among all worker machines.

We will replace the sufficient decrease property in Assumption 8 with the following key property that turns out to be easier to verify:

ASSUMPTION 10 (Proximal Lipschitz). *We say a pair of functions $f$ and $g$ satisfy the proximal Lipschitz property on a sequence $\{\mathbf{x}(t)\}$ if for all $\eta$ sufficiently small, there exists $L_\eta \in o(1)$, i.e. $L_\eta \to 0$ as $\eta \to 0$, such that for all large $t$,*

$$(59) \qquad \|\Delta_\eta(\mathbf{x}(t)) - \Delta_\eta(\mathbf{x}(t+1))\| \leq L_\eta \|\mathbf{x}(t) - \mathbf{x}(t+1)\|,$$

$where^2 \; \Delta_\eta(\mathbf{x}) \in \mathrm{prox}_g^\eta(\mathbf{x} - \eta\nabla f(\mathbf{x})) - \mathbf{x}$.

The proximal Lipschitz assumption is motivated by the special case where $g \equiv 0$ and hence $\Delta_\eta(\mathbf{x}) = -\eta\nabla f(\mathbf{x})$ is $\eta$-Lipschitz, thanks to Assumption 2. As we have seen in previous sections, Lipschitz continuity plays a crucial role in our proof where a major difficulty is to control the inconsistencies among different worker machines due to communication delays. Similarly here, the proximal Lipschitz property, as we show next, allows us to remove the sufficient decrease property in Assumption 8—the seemingly strong assumption that we needed in proving our main result Theorem 10.

Let us first present a quick justification for Assumption 10.

LEMMA 11. *Suppose the functions $f$ and $g$ both have Lipschitz continuous gradient, then Assumption 10 holds for any sequence $\{\mathbf{x}(t)\}$.*

*Proof.* Let us denote $L_f$ and $L_g$ as the Lipschitz constant of the gradient $\nabla f$ and $\nabla g$, respectively. Since $\Delta_\eta(\mathbf{x}) \in \mathrm{prox}_g^\eta(\mathbf{x} - \eta\nabla f(\mathbf{x})) - \mathbf{x}$, using the optimality condition for the proximal map, see for instance [36, Proposition 7(iii)], we have

$$\mathbf{x} + \Delta_\eta(\mathbf{x}) + \eta\nabla g\big(\mathbf{x} + \Delta_\eta(\mathbf{x})\big) = \mathbf{x} - \eta\nabla f(\mathbf{x}),$$

and similarly

$$\mathbf{z} + \Delta_\eta(\mathbf{z}) + \eta\nabla g\big(\mathbf{z} + \Delta_\eta(\mathbf{z})\big) = \mathbf{z} - \eta\nabla f(\mathbf{z}).$$

Subtracting one inequality from another, we obtain

$$\begin{aligned}
\|\Delta_\eta(\mathbf{x}) - \Delta_\eta(\mathbf{z})\| &= \|\eta\nabla g\big(\mathbf{z} + \Delta_\eta(\mathbf{z})\big) - \eta\nabla g\big(\mathbf{x} + \Delta_\eta(\mathbf{x})\big) + \eta\nabla f(\mathbf{z}) - \eta\nabla f(\mathbf{x})\| \\
&\le \eta L_g\|\mathbf{z} - \mathbf{x} + \Delta_\eta(\mathbf{z}) - \Delta_\eta(\mathbf{x})\| + \eta L_f\|\mathbf{z} - \mathbf{x}\| \\
&\le \eta L_g\|\Delta_\eta(\mathbf{z}) - \Delta_\eta(\mathbf{x})\| + \eta(L_f + L_g)\|\mathbf{z} - \mathbf{x}\|.
\end{aligned}$$

Rearranging we obtain

$$\|\Delta_\eta(\mathbf{x}) - \Delta_\eta(\mathbf{z})\| \le \frac{\eta(L_f + L_g)}{1 - \eta L_g}\|\mathbf{z} - \mathbf{x}\|,$$

when $0 < \eta < 1/L_g$. Clearly, when $\eta$ is mall, the leading coefficient $\frac{\eta(L_f+L_g)}{1-\eta L_g} \in \mathcal{O}(\eta) \subseteq o(1)$, and our proof is complete.                          □

It is clear that Lemma 11 captures the motivating case $g \equiv 0$, but also many other important functions, such as the widely-used regularization function $g = \|\cdot\|_p^p$ for any $p > 1$. We can now continue with our next result in this section.

THEOREM 12. *Let Assumptions 1 to 4 and 9 hold for m-PAPG, and let $F$ satisfy the KL property in Definition 9. Fix any $r > 1$ with $C = \frac{r^{s+1}-1}{r-1}$ and step size $\eta$ such that $\eta < \frac{1}{L(1+2\sqrt{p}C+2\sqrt{p}s)}$. If for each local sequence $\{\mathbf{x}^i(t)\}$ generated by m-PAPG, Assumption 10 holds with $L_\eta \le \frac{r^2-1}{2pr^2C^2}$, and the global sequence $\{\mathbf{x}(t)\}$ is bounded, then the finite length properties in (55) and (56) hold. In particular, $\{\mathbf{x}(t)\}$ and $\{\mathbf{x}^i(t)\}_{i=1}^p$ converge to the same critical point of $F$.*

*Proof.* Using the elementary inequality $\|a\|^2 - \|b\|^2 \le 2\|a\|\|a - b\|$, we have for all $t$:

$$\|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2 - \|\mathbf{x}(t+2) - \mathbf{x}(t+1)\|^2$$

---

$^2$Should the proximal map be multi-valued, we contend with any single-valued selection.

$$\leq 2 \left\| \mathbf{x}(t+1) - \mathbf{x}(t) \right\| \cdot \left\| (\mathbf{x}(t+1) - \mathbf{x}(t)) - (\mathbf{x}(t+2) - \mathbf{x}(t+1)) \right\|$$

$$\leq 2 \left\| \mathbf{x}(t+1) - \mathbf{x}(t) \right\| \cdot \sum_{i=1}^{p} \left\| (x_i(t+1) - x_i(t)) - (x_i(t+2) - x_i(t+1)) \right\|$$

$$\overset{(i)}{\leq} 2 \left\| \mathbf{x}(t+1) - \mathbf{x}(t) \right\| \cdot \sum_{i=1}^{p} \left\| \Delta_\eta(\mathbf{x}^i(t)) - \Delta_\eta(\mathbf{x}^i(t+1)) \right\|$$

$$\overset{(ii)}{\leq} 2 \left\| \mathbf{x}(t+1) - \mathbf{x}(t) \right\| \left( \sum_{i=1}^{p} L_\eta \|\mathbf{x}^i(t) - \mathbf{x}^i(t+1)\| \right)$$

$$(60) \qquad \overset{(iii)}{\leq} 2pL_\eta \left\| \mathbf{x}(t+1) - \mathbf{x}(t) \right\| \cdot \sum_{k=(t-s)_+}^{t} \left\| \mathbf{x}(k+1) - \mathbf{x}(k) \right\|,$$

where (i) is due to Assumption 9 hence $t \in T_i$ for all $t$, (ii) follows from Assumption 10, and (iii) is due to (15).

If for some $r > 1$ there exists some $T$ such that for all $t \geq T$,

$$(61) \qquad \sum_{k=(t-s)_+}^{t} \left\| \mathbf{x}(k+1) - \mathbf{x}(k) \right\| \geq C \|\mathbf{x}(t+1) - \mathbf{x}(t)\|,$$

where $C = \frac{r^{s+1}-1}{r-1} > s+1$ (since $r > 1$ and w.l.o.g. $s > 0$). Summing the index $t$ from $T$ to $n$ yields

$$C \sum_{t=T}^{n} \left\| \mathbf{x}(t+1) - \mathbf{x}(t) \right\| \leq \sum_{t=T}^{n} \sum_{k=(t-s)_+}^{t} \left\| \mathbf{x}(k+1) - \mathbf{x}(k) \right\|$$

$$\leq (s+1) \sum_{t=(T-s)_+}^{n} \left\| \mathbf{x}(t+1) - \mathbf{x}(t) \right\|,$$

which after rearranging terms becomes

$$(C - s - 1) \sum_{t=T}^{n} \left\| \mathbf{x}(t+1) - \mathbf{x}(t) \right\| \leq (s+1) \sum_{t=(T-s)_+}^{T-1} \left\| \mathbf{x}(t+1) - \mathbf{x}(t) \right\|.$$

Since the right hand side does not depend on $n$, letting $n$ tend to infinity we conclude

$$(62) \qquad \sum_{t=0}^{\infty} \left\| \mathbf{x}(t+1) - \mathbf{x}(t) \right\| < \infty,$$

and the proof of the finite length property would be complete.

Therefore, in the remaining part of the proof, we can assume (61) fails for infinitely many $t$. Take any such $t = \hat{t}$, we have

$$(63) \qquad \sum_{k=(t-s)_+}^{t} \left\| \mathbf{x}(k+1) - \mathbf{x}(k) \right\| \leq C \|\mathbf{x}(t+1) - \mathbf{x}(t)\| \leq C^2 \|\mathbf{x}(t+1) - \mathbf{x}(t)\|,$$

since $C > 1$. Combining (60) and (63) we have for $t = \hat{t}$:

$$\left\| \mathbf{x}(t+1) - \mathbf{x}(t) \right\|^2 - \left\| \mathbf{x}(t+2) - \mathbf{x}(t+1) \right\|^2 \leq 2pL_\eta C^2 \left\| \mathbf{x}(t+1) - \mathbf{x}(t) \right\|^2$$

$$\leq \left(1 - \frac{1}{r^2}\right)\|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2,$$

if $\eta$ is small enough (recall that $L_\eta = o(1)$). After rearranging terms we conclude that for $t = \hat{t}$:

(64) $$\|\mathbf{x}(t+1) - \mathbf{x}(t)\| \leq r\|\mathbf{x}(t+2) - \mathbf{x}(t+1)\|.$$

Using induction we can continue the same process for any $t \geq \hat{t}$. Indeed, suppose (64) is true for any $t \leq m-1$, then (60) holds (for any $t$), and (63) also holds: If $m \leq \hat{t} + s$, then

$$\sum_{k=(m-s)_+}^{m} \|\mathbf{x}(k+1) - \mathbf{x}(k)\| = \sum_{k=(m-s)_+}^{\hat{t}} \|\mathbf{x}(k+1) - \mathbf{x}(k)\| + \sum_{k=\hat{t}+1}^{m} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|$$

$$\overset{(i)}{\leq} \sum_{k=(\hat{t}-s)_+}^{\hat{t}} \|\mathbf{x}(k+1) - \mathbf{x}(k)\| + \sum_{k=\hat{t}+1}^{m} r^{m-k}\|\mathbf{x}(m+1) - \mathbf{x}(m)\|$$

$$\overset{(ii)}{\leq} C \left[ \|\mathbf{x}(\hat{t}+1) - \mathbf{x}(\hat{t})\| + \sum_{k=\hat{t}+1}^{m} r^{m-k}\|\mathbf{x}(m+1) - \mathbf{x}(m)\| \right]$$

$$\overset{(iii)}{\leq} C \sum_{k=\hat{t}}^{m} r^{m-k}\|\mathbf{x}(m+1) - \mathbf{x}(m)\|$$

$$\overset{(iv)}{\leq} C^2 \|\mathbf{x}(m+1) - \mathbf{x}(m)\|,$$

where (i) is due to the induction hypothesis, (ii) is due to the definition of $\hat{t}$ and the fact that $C > 1$, (iii) is due to again the induction hypothesis, and finally (iv) is due to the definition of $C$ (recall $m \leq \hat{t} + s$). If $m > \hat{t} + s$, the same inequality, with $C^2$ replaced by $C$, would still hold (essentially dropping all the first terms on the right hand side of the above inequalities). Thus, (60) and (63) would imply again (64) for $t = m$.

Lastly, we recall from (22) that for large $t$,

$$F\big(\mathbf{x}(t+1)\big) - F\big(\mathbf{x}(t)\big) \leq \tfrac{1}{2}(L - 1/\eta)\|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2$$

$$+ \sqrt{p}L\|\mathbf{x}(t+1) - \mathbf{x}(t)\| \sum_{k=(t-s)_+}^{t-1} \|\mathbf{x}(k+1) - \mathbf{x}(k)\|.$$

$$\leq \tfrac{1}{2}(L - 1/\eta)\|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2 + \sqrt{p}CL\|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2.$$

$$\leq -\alpha\|\mathbf{x}(t+1) - \mathbf{x}(t)\|^2,$$

where $\alpha = \tfrac{1}{2}(1/\eta - L - 2\sqrt{p}CL) > 0$ if $\eta$ is small. Hence, the sufficient decrease property in Assumption 8 is verified and the finite length properties follow from Theorem 10. $\square$

Lastly, we show that Assumption 10 also holds for the important cardinality function $\|\mathbf{x}\|_0$ (number of nonzero entries).

LEMMA 13. *Consider the same setting as in* Theorem 5, *then* Assumption 10 *holds for any function $f$ and $g = \|\cdot\|_0$ on all local sequences $\{\mathbf{x}^i(t)\}$ of m-PAPG.*

*Proof.* The crucial observation here is that for the cardinality function $g = \|\cdot\|_0$, its proximal map on the $j$-th entry can be chosen as:

$$(65) \qquad \text{prox}_{g_j}^\eta(z_j) = \begin{cases} z_j, & \text{if } |z_j| > \sqrt{2\eta} \\ 0, & \text{otherwise} \end{cases} .$$

However, Theorem 5 implies that $\lim_{t\to\infty} \|\mathbf{x}^i(t+1) - \mathbf{x}^i(t)\| = 0$. Thus, for $t$ sufficiently large, the sequence $\{\mathbf{x}^i(t)\}$ will have the same support $\Omega$ (indices that have nonzero entries), for otherwise $\|\mathbf{x}^i(t+1) - \mathbf{x}^i(t)\| \geq \sqrt{2\eta}$ even if one index in the support changes. Therefore,

$$\|\Delta_\eta(\mathbf{x}^i(t+1)) - \Delta_\eta(\mathbf{x}^i(t))\| \overset{(i)}{\leq} \sum_{j\in\Omega} \|\text{prox}_{g_j}^\eta(x_j^i(t+1) - \eta\nabla_j f(\mathbf{x}^i(t+1))) - x_j^i(t+1)$$

$$- \text{prox}_{g_j}^\eta(x_j^i(t) - \eta\nabla_j f(\mathbf{x}^i(t))) - x_j^i(t)\|$$

$$\overset{(ii)}{\leq} \sum_{j\in\Omega} \|\eta\nabla_j f(\mathbf{x}^i(t+1)) - \eta\nabla_j f(\mathbf{x}^i(t))\|$$

$$\overset{(iii)}{\leq} \eta p L \|\mathbf{x}^i(t+1) - \mathbf{x}^i(t)\|,$$

where (i) is the triangle inequality, (ii) uses the property of the proximal map (65), and (iii) is due to Assumption 2. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ □

Note that similar results as Lemma 13 can be derived for the rank function, and more generally for functions whose proximal map is discontinuous with pieces satisfying Lemma 11.

**7. Conclusion.** We have proposed m-PAPG as an extension of the proximal gradient algorithm to the model parallel and partially asynchronous setting. m-PAPG allows worker machines to operate asynchronously as long as they are not too far apart, hence greatly improves the system throughput. The convergence properties of m-PAPG are thoroughly analyzed. In particular, we proved that: 1) every limit point of the sequences generated by m-PAPG is a critical point of the objective function; 2) under an additional error bound condition, the function values decay periodically linearly; 3) under the additional Kurdyka-Łojasiewicz inequality, the sequences generated by m-PAPG converge to the same critical point, provided that a proximal Lipschitz condition is satisfied. In the future we plan to further weaken the proximal Lipschitz condition so that our analysis can handle many more nonsmooth functions.

REFERENCES

[1] P.-A. ABSIL, R. MAHONY, AND B. ANDREWS, *Convergence of the iterates of descent methods for analytic cost functions*, SIAM Journal on Optimization, 16 (2005), pp. 531–547.
[2] A. AGARWAL AND J. C. DUCHI, *Distributed delayed stochastic optimization*, in Advances in Neural Information Processing Systems 24, 2011, pp. 873–881.
[3] H. ATTOUCH AND J. BOLTE, *On the convergence of the proximal algorithm for nonsmooth functions involving analytic features*, Mathematical Programming, 116 (2009), pp. 5–16.
[4] H. ATTOUCH, J. BOLTE, P. REDONT, AND A. SOUBEYRAN, *Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the Kurdyka-Łojasiewicz inequality*, Mathematics of Operations Research, 35 (2010), pp. 438–457.

[5]   G. M. Baudet, *Asynchronous iterative methods for multiprocessors*, Journal of the Association for Computing Machinery, 25 (1978), pp. 226–244.

[6]   A. Beck and M. Teboulle, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*, SIAM J. Img. Sci., 2 (2009), pp. 183–202.

[7]   A. Beck and M. Teboulle, *Smoothing and first order methods: A unified framework*, SIAM Journal on Optimization, 22 (2012), pp. 557–580.

[8]   D. P. Bertsekas and J. N. Tsitsiklis, *Parallel and Distributed Computation: Numerical Methods*, Prentice-Hall, Inc., Upper Saddle River, NJ, USA, 1989.

[9]   J. Bolte, A. Danilidis, O. Ley, and L. Mazet, *Characterizations of Łojasiewicz inequalities and applications: Subgradient flows, talweg, convexity*, Transactions of the American Mathematical Society, 362 (2010), pp. 3319–3363.

[10]  J. Bolte, S. Sabach, and M. Teboulle, *Proximal alternating linearized minimization for nonconvex and nonsmooth problems*, Mathematical Programming, 146 (2014), pp. 459–494.

[11]  D. Chazan and W. Miranker, *Chaotic relaxation*, Linear Algebra and Its Applications, 2 (1969), pp. 199–222.

[12]  R. Collobert, F. Sinz, J. Weston, and L. Bottou, *Trading convexity for scalability*, 2006, pp. 201–208.

[13]  J. Dean and S. Ghemawat, *Mapreduce: Simplified data processing on large clusters*, Communications of ACM, 51 (2008), pp. 107–113.

[14]  D. Drusvyatskiy and A. S. Lewis, *Error bounds, quadratic growth, and linear convergence of proximal methods*, 2016.

[15]  J. Fan and R. Li, *Variable selection via nonconcave penalized likelihood and its oracle properties*, Journal of the American Statistical Association, 96 (2001), pp. 1348–1360.

[16]  O. Fercoq and P. Richtárik, *Accelerated, parallel, and proximal coordinate descent*, SIAM Journal on Optimization, 25 (2015), pp. 1997–2023.

[17]  H. Feyzmahdavian, A. Aytekin, and M. Johansson, *A delayed proximal gradient method with linear convergence rate*, in 2014 IEEE International Workshop on Machine Learning for Signal Processing.

[18]  M. Fukushima and H. Mine, *A generalized proximal point algorithm for certain non-convex minimization problems*, International Journal of Systems Science, 12 (1981), pp. 989–1000.

[19]  Q. Ho, J. Cipar, H. Cui, S. Lee, J. K. Kim, P. B. Gibbons, G. A. Gibson, G. Ganger, and E. P. Xing, *More effective distributed ml via a stale synchronous parallel parameter server*, in Advances in Neural Information Processing Systems 26, 2013, pp. 1223–1231.

[20]  K. Kurdyka, *On gradients of functions definable in o-minimal structures*, Annales de l'institut Fourier, 48 (1998), pp. 769–783.

[21]  M. Li, D. G. Andersen, J. W. Park, A. J. Smola, A. Ahmed, V. Josifovski, J. Long, E. J. Shekita, and B.-Y. Su, *Scaling distributed machine learning with the parameter server*, in 11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14), 2014, pp. 583–598.

[22]  J. Liu and S. J. Wright, *Asynchronous stochastic coordinate descent: Parallelism and convergence properties*, SIAM Journal on Optimization, 25 (2015), pp. 351–376.

[23]  Y. Low, D. Bickson, J. Gonzalez, C. Guestrin, A. Kyrola, and J. M. Hellerstein, *Distributed graphlab: A framework for machine learning and data mining in the cloud*, Proc. VLDB Endow., 5 (2012), pp. 716–727.

[24]  Z. Lu and L. Xiao, *On the complexity analysis of randomized block-coordinate descent methods*, Mathematical Programming, 152 (2015), pp. 615–642.

[25]  Z.-Q. Luo and P. Tseng, *Error bounds and convergence analysis of feasible descent methods: A general approach*, Annals of Operations Research, 46 (1993), pp. 157–178.

[26]  R. Mazumder, J. H. Friedman, and T. Hastie, *Sparsenet: Coordinate descent with nonconvex penalties*, Journal of the American Statistical Association, 106 (2011), pp. 1125–1138.

[27]  Y. Nesterov, *Gradient methods for minimizing composite functions*, Mathematical Programming, Series B, 140 (2013), pp. 125–161.

[28]  B. Recht, C. Re, S. Wright, and F. Niu, *Hogwild: A lock-free approach to parallelizing stochastic gradient descent*, in Advances in Neural Information Processing Systems 24, 2011, pp. 693–701.

[29]  P. Richtárik and M. Takáč, *Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function*, Mathematical Programming, 144 (2014), pp. 1–38.

[30]  R. Rockafellar and R. Wets, *Variational Analysis*, Springer, 1997.

[31]  S. Sra, *Scalable nonconvex inexact proximal splitting*, in Advances of Neural Information Processing Systems, 2012.

[32] P. Tseng, *On the rate of convergence of a partially asynchronous gradient projection algorithm*, SIAM Journal on Optimization, 1 (1991), pp. 603–619.

[33] L. G. Valiant, *A bridging model for parallel computation*, Communications of ACM, 33 (1990), pp. 103–111.

[34] Y. Wu and Y. Liu, *Robust truncated hinge loss support vector machines*, Journal of the American Statistical Association, 102 (2007), pp. 974–983.

[35] L. Xu, K. Crammer, and D. Schuurmans, *Robust support vector machine training via convex outlier ablation*, 2006.

[36] Y. Yu, X. Zheng, M. Marchetti-Bowick, and E. P. Xing, *Minimizing nonconvex non-separable functions*, in The $17^{th}$ International Conference on Artificial Intelligence and Statistics (AISTATS), 2015.

[37] M. Zaharia, M. Chowdhury, M. J. Franklin, S. Shenker, and I. Stoica, *Spark: Cluster computing with working sets*, 2010, pp. 10–10.

[38] C.-H. Zhang, *Nearly unbaised variable selection under minimax concave penalty*, Annals of Statistics, 38 (2010), pp. 894–942.

[39] C.-H. Zhang and T. Zhang, *A general theory of concave regularization for high-dimensional sparse estimation problems*, Statistical Science, 27 (2012), pp. 576–593.

[40] H. Zhang, *The restricted strong convexity revisited: analysis of equivalence to error bound and quadratic growth*, Optimization Letters, (2016), pp. 1–17.

[41] Y. Zhou, Y. Yu, W. Dai, Y. Liang, and E. Xing, *On convergence of model parallel proximal gradient algorithm for stale synchronous parallel system*, in The $19^{th}$ International Conference on Artificial Intelligence and Statistics (AISTATS), 2016.

[42] Z. Zhou, Q. Zhang, and A. M.-C. So, $\ell_{1,p}$-*norm regularization: Error bounds and convergence rate analysis of first-order methods*, in Proceedings of the 32nd International Conference on Machine Learning, JMLR Workshop and Conference Proceedings, 2015, pp. 1501–1510.