# Approximate Byzantine Fault-Tolerance in Distributed Optimization [*]

Shuo Liu [†]      Nirupam Gupta [‡]      Nitin H. Vaidya [§]

## Abstract

This paper considers the problem of Byzantine fault-tolerance in distributed multi-agent optimization. In this problem, each agent has a local cost function, and in the fault-free case, the goal is to design a distributed algorithm that allows all the agents to find a minimum point of all the agents' aggregate cost function. We consider a scenario where some agents might be Byzantine faulty that renders the original goal of computing a minimum point of all the agents' aggregate cost vacuous. A more reasonable objective for an algorithm in this scenario is to allow all the non-faulty agents to compute the minimum point of only the non-faulty agents' aggregate cost. Prior work [26] shows that if there are up to $f$ (out of $n$) Byzantine agents then a minimum point of the non-faulty agents' aggregate cost can be computed *exactly* if and only if the non-faulty agents' costs satisfy a certain redundancy property called $2f$-*redundancy*. However, $2f$-redundancy is an ideal property that can be satisfied only in systems free from noise or uncertainties, which can make the goal of exact fault-tolerance *unachievable* in some applications. Thus, we introduce the notion of $(f, \epsilon)$-resilience, a generalization of exact fault-tolerance wherein the objective is to find an approximate minimum point of the non-faulty aggregate cost, with $\epsilon$ accuracy. This approximate fault-tolerance can be achieved under a weaker condition that is easier to satisfy in practice, compared to $2f$-redundancy. We obtain necessary and sufficient conditions for achieving $(f, \epsilon)$-resilience characterizing the correlation between relaxation in redundancy and approximation in resilience. In case when the agents' cost functions are differentiable, we obtain conditions for $(f, \epsilon)$-resilience of the distributed gradient-descent method when equipped with *robust gradient aggregation*; such as *comparative gradient elimination* or *coordinate-wise trimmed mean*.

**Keywords**— Distributed optimization; Approximate fault-tolerance; Distributed gradient-descent

---

# Contents

# 1 Introduction

The problem of distributed optimization in multi-agent systems has gained significant attention in recent years [9, 19, 36]. In this problem, each agent has a *local cost function* and, when the agents are fault-free, the goal is to design algorithms that allow the agents to collectively minimize the aggregate of their cost functions. To be precise, suppose that there are $n$ agents in the system and let $Q_i(x)$ denote the local cost function of agent $i$, where $x$ is a $d$-dimensional vector of real values, i.e., $x \in \mathbb{R}^d$. A traditional distributed optimization algorithm outputs a *global minimum* $x^*$ such that

$$x^* \in \arg \min_{x \in \mathbb{R}^d} \sum_{i=1}^n Q_i(x). \tag{1}$$

As a simple example, $Q_i(x)$ may denote the cost for an agent $i$ (which may be a robot or a person) to travel to location $x$ from their current location, and $x^*$ is a location that minimizes the total cost of meeting for all the agents. Such multi-agent optimization is of interest in many practical applications, including distributed machine learning [9], swarm robotics [42], and distributed sensing [41].

We consider the distributed optimization problem in the presence of up to $f$ Byzantine faulty agents, originally introduced by Su and Vaidya [48]. The Byzantine faulty agents may behave arbitrarily [30]. In particular, the non-faulty agents may share arbitrary incorrect and inconsistent information in order to bias the output of a distributed optimization algorithm. For example, consider an application of multi-agent optimization in the case of distributed sensing where the agents (or *sensors*) observe a common *object* in order to collectively identify the object. However, the faulty agents may send arbitrary observations concocted to prevent the non-faulty agents from making the correct identification [13, 15, 38, 49]. Similarly, in the case of distributed learning, which is another application of distributed optimization, the faulty agents may send incorrect information based on *mislabelled* or arbitrary concocted data points to prevent the non-faulty agents from learning a *good* classifier [1, 3, 6, 11, 12, 14, 25, 52].

## 1.1 Background: Exact Fault-Tolerance

In the *exact fault-tolerance* problem, the goal is to design a distributed algorithm that allows all the non-faulty agents to compute a minimum point of the aggregate cost of only the non-faulty agents [26]. Specifically, suppose that in a given execution, set $\mathcal{B}$ with $|\mathcal{B}| \leq f$ is the set of Byzantine agents, where notation $|\cdot|$ denotes the set cardinality, and $\mathcal{H} = \{1, \ldots, n\} \setminus \mathcal{B}$ denotes the set of non-faulty (i.e., honest) agents. Then, a distributed optimization algorithm has exact fault-tolerance if it outputs a point $x_{\mathcal{H}}^*$ such that

$$x_{\mathcal{H}}^* \in \arg \min_{x \in \mathbb{R}^d} \sum_{i \in \mathcal{H}} Q_i(x). \tag{2}$$

However, since the identity of the Byzantine agents is a priori unknown, in general, exact fault-tolerance is unachievable [48]. Specifically, as shown in [26, 27], exact fault-tolerance can be achieved *if and only if* the agents' cost functions satisfy the $2f$-*redundancy* property defined below.

**Definition 1** ($2f$-**redundancy**). *The agents' cost functions are said to have $2f$-redundancy property if and only if for every pair of subsets $S, \widehat{S} \subseteq \{1, \ldots, n\}$ with $\widehat{S} \subseteq S$, $|S| = n - f$, and $\left|\widehat{S}\right| \geq n - 2f$,*

$$\arg \min_{x \in \mathbb{R}^d} \sum_{i \in \widehat{S}} Q_i(x) = \arg \min_{x \in \mathbb{R}^d} \sum_{i \in S} Q_i(x).$$

3

In principle, the $2f$-redundancy property can be realized by design for many applications of multi-agent distributed optimization including distributed sensing and distributed learning (see [24, 26]). However, practical realization of $2f$-redundancy can be difficult in the presence of *noise* in the real-world systems. Therefore, we propose a pragmatic generalization of exact fault-tolerance, namely $(f, \epsilon)$-*resilience*.

## 1.2 $(f, \epsilon)$-Resilience: A Relaxation of Exact Fault-Tolerance

Intuitively, the proposed notion of $(f, \epsilon)$-*resilience* requires an algorithm to output *approximation* of a minimum point of the aggregate of the cost functions of sufficiently large subsets of non-faulty agents. We define $(f, \epsilon)$-*resilience* below, where $\epsilon \in \mathbb{R}_{\geq 0}$ is the measure of approximation and $\|\cdot\|$ denotes the Euclidean norm. The Euclidean distance between a point $x$ and a non-empty set $X$ in space $\mathbb{R}^d$ is denoted by $\mathrm{dist}\,(x,\, X)$, and is defined as

$$\mathrm{dist}\,(x,\, X) = \inf_{y \in X} \|x - y\|. \tag{3}$$

**Definition 2** (($f, \epsilon$)**-resilience**)**.** *A distributed optimization algorithm is said to be* $(f, \epsilon)$-*resilient if it outputs a point $\widehat{x} \in \mathbb{R}^d$ such that for every subset $S$ of non-faulty agents with $|S| = n - f$,*

$$\mathrm{dist}\left(\widehat{x},\, \arg\min_{x \in \mathbb{R}^d} \sum_{i \in S} Q_i(x)\right) \leq \epsilon,$$

*despite the presence of up to $f$ Byzantine agents.*

Thus, with $(f, \epsilon)$-*resilience*, the output is within distance $\epsilon$ of a minimum point of the aggregate cost function of any $n - f$ non-faulty agents. As there can be at most $f$ Byzantine faulty agents whose identity remains unknown, the following two scenarios are indistinguishable in general: (1) there are exactly $f$ Byzantine agents, and (2) there are less than $f$ Byzantine agents. Thus, estimation for the minimum point of the aggregate cost functions of $n - f$ non-faulty agents is indeed a reasonable goal [48]. Analogous resilience requirements have been previously studied in other contexts as well, such as robust statistics (e.g., robust mean estimation [12, 46]) and fault-tolerant linear state estimation [5, 21, 22, 34, 37, 45]. In this work, we address resilience in the context of distributed optimization.

In this paper, we only consider *deterministic* algorithms which, given a fixed set of inputs from the agents, always output the same point in $\mathbb{R}^d$. Thus, a deterministic $(f, \epsilon)$-resilient algorithm produces a unique output point in all of its executions with identical inputs from all the agents (including the faulty ones). **Note that** in the deterministic framework, exact fault-tolerance is equivalent to $(f, 0)$-resilience, i.e., a deterministic $(f, 0)$-resilient algorithm achieves exact fault-tolerance, and vice-versa.[1] Therefore, results on $(f, \epsilon)$-resilience for arbitrary $\epsilon \geq 0$ have a wider application compared to results applicable only to exact fault-tolerance, e.g., [6, 20, 26, 47].

We show that $(f, \epsilon)$-*resilience* requires a *weaker redundancy* condition, in comparison to $2f$-redundancy, named $(2f, \epsilon)$-*redundancy* defined in Definition 3 below. Recall that the *Euclidean Hausdorff distance* between two sets $X$ and $Y$ in $\mathbb{R}^d$, which we denote by $\mathrm{dist}\,(X,\, Y)$, is defined as follows [35]:

$$\mathrm{dist}\,(X,\, Y) \triangleq \max\left\{\sup_{x \in X} \mathrm{dist}\,(x,\, Y),\ \sup_{y \in Y} \mathrm{dist}\,(y,\, X)\right\}. \tag{4}$$

---

[1]Refer to Appendix B for proof.

**Definition 3** (($2f$, $\epsilon$)-**redundancy**). *The agents' cost functions are said to have* ($2f$, $\epsilon$)-*redundancy property if and only if for every pair of subsets* $S, \widehat{S} \subseteq \{1, \ldots, n\}$ *with* $|S| = n - f$, $\left|\widehat{S}\right| \geq n - 2f$ *and* $\widehat{S} \subseteq S$,

$$\mathrm{dist}\left(\arg\min_{x \in \mathbb{R}^d} \sum_{i \in S} Q_i(x), \arg\min_{x \in \mathbb{R}^d} \sum_{i \in \widehat{S}} Q_i(x)\right) \leq \epsilon. \tag{5}$$

It is easy to show that $2f$-redundancy (Definition 1) is equivalent to ($2f$, $0$)-redundancy (note that $\epsilon = 0$ here). It is also obvious that $2f$-redundancy implies ($2f$, $\epsilon$)-redundancy for all $\epsilon \geq 0$. However, the converse need not be true. Thus, the ($2f$, $\epsilon$)-redundancy property with $\epsilon > 0$ is *weaker* than $2f$-redundancy.

## 1.3 Applications

Our results are applicable to a large class of distributed optimization problems; including distributed sensing [15, 37, 38, 47], distributed machine learning [8, 9, 14, 54], and distributed linear regression (Section 5). We discuss below the specific case of distributed learning.

**Distributed Learning:** In this particular optimization problem, each agent has some local *data points* and the goal for the agents is to compute a learning parameter that best models the collective data points observed by all the agents [8]. Specifically, given a learning parameter $x$, for each data point $z$, we define a loss function $\ell(x; z)$. Suppose that the data generating distribution of agent $i$ is $\mathcal{D}_i$, and let $\mathbb{E}_{z \sim \mathcal{D}_i}$ denote the expectation with respect to the random data point $z$ over distribution $\mathcal{D}_i$. Then,

$$Q_i(x) \triangleq \mathbb{E}_{z \sim \mathcal{D}_i} \ell(x; z)$$

When the distribution of data points is identical for all the agents then the $2f$-redundancy property holds true. However, in practice this is rarely the case [14, 54]. Indeed, different agents may have different data distributions in practice. Therefore, exact fault-tolerance in a pragmatic distributed learning framework is an extremely difficult (if not impossible) goal. In the context of distributed learning, our results on approximate fault-tolerance characterize the relationship between the correlation amongst different agents' data (i.e., degree of redundancy), and the fault-tolerance achieved.

## 1.4 System architecture

We consider *synchronous* systems. Our results apply to the two architectures shown in Figure 1. In the server-based architecture, the server is assumed to be trustworthy, but up to $f$ agents may be Byzantine faulty. In the peer-to-peer architecture, the agents are connected by a complete network, and up to $f$ of these agents may be Byzantine faulty. Provided that $f < \frac{n}{3}$, an algorithm for the server-based architecture can be simulated in the peer-to-peer system using the well-known *Byzantine broadcast* primitive [33]. For simplicity of presentation, the rest of this paper considers the server-based architecture.

## 1.5 Summary of Our Contributions

In the **first part** of the paper, i.e., Section 3, we obtain conditions on feasibility and achievability of approximate fault-tolerance of desirable accuracy. Specifically, we show that

- ($f, \epsilon$)-resilience is feasible only if ($2f, \epsilon$)-redundancy property holds true.
- If ($2f, \epsilon$)-redundancy property holds true then ($f, 2\epsilon$)-resilience is achievable.
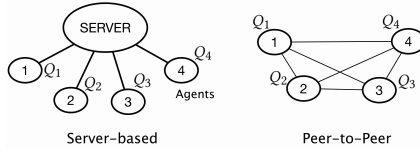
5

Figure 1: System architecture.

In the **second part**, i.e., Sections 4 and 5, we consider the case when agents' costs are differentiable, such as in machine learning [8, 53], or regression [24, 47, 49]. We consider the distributed gradient-descent (DGD) method - an iterative distributed optimization algorithm commonly used in this particular case.

- We propose a generic sufficient condition for convergence of the DGD method equipped with a *gradient-filter* (also referred as *robust gradient aggregation*), which is a common fault-tolerance mechanism, e.g., see [6, 14, 26, 54].

- Later, in Section 4.2, we utilize the above result to obtain approximate fault-tolerance properties of the following two specific gradient-filters, under $(2f, \epsilon)$-redundancy: (i) Comparative gradient elimination (CGE) [24], and (ii) Coordinate-wise trimmed mean (CWTM) [47]. These two gradient-filters are both easy to implement and versatile [23, 26, 47, 54].

- Finally, in Section 5, we present empirical comparisons between approximate fault-tolerance of the two gradient-filters by simulating a problem of distributed linear regression.

**Note:** As $(f, 0)$-resilience is equivalent to exact fault-tolerance (see Section 1.2), our results on $(f, \epsilon)$-resilience encapsulate all the existing results applicable only to exact fault-tolerance, such as the ones in [6, 20, 26, 47].

Compared to related works [28, 29], we present precise redundancy conditions needed for obtaining Byzantine fault-tolerance within a specified approximation error. Unlike them, our results on the impossibility and feasibility of approximate fault-tolerance are applicable to non-differentiable cost functions. Moreover, in the case when the cost functions are differentiable, we present a generic condition for convergence of the DGD method that can precisely model the approximate fault-tolerance property of a generic robust gradient-aggregation rule (a.k.a., gradient-filter).

This is the full version of the paper including proofs of theorems and additional experimental results and discussion.

## 2 Other Related Work

In the past, different notions of approximate fault-tolerance, besides $(f, \epsilon)$-resilience, have been used to analyze Byzantine fault-tolerance of different distributed optimization algorithms [18, 48]. As we discuss below in Section 2.1, the difference between these other definitions and our definition of $(f, \epsilon)$-resilience arises mainly due to the applicability of the distributed optimization problems. Later, in Section 2.2, we discuss some prior work on gradient-filters used for achieving Byzantine fault-tolerance in the distributed gradient-descent method.

## 2.1 Alternate Notions of Approximation in Fault-Tolerance

As proposed by Su and Vaidya, 2016 [48], instead of a minimum point of the *uniformly weighted* aggregate of non-faulty agents' cost functions, a distributed optimization algorithm may output a minimum point of a *non-uniformly weighted* aggregate of non-faulty costs, i.e., $\sum_{i \in \mathcal{H}} \alpha_i Q_i(x)$, where $\mathcal{H}$ denotes the set of at least $n - f$ non-faulty agents, and $\alpha_i \geq 0$ for all $i \in \mathcal{H}$. As is suggested in [48], upon re-scaling the coefficients such that $\sum_{i \in \mathcal{H}} \alpha_i = 1$, we can measure approximation in fault-tolerance using two metrics: (1) the number of coefficients in $\{\alpha_i, i \in \mathcal{H}\}$ that are positive, and (2) the minimum positive value amongst the coefficients: $\min \{\alpha_i; \alpha_i > 0, i \in \mathcal{H}\}$. Results on the achievability of this particular form of approximation for the scalar case (i.e., $d = 1$) can be found in [48, 50]. However, we are unaware of similar results for the case of higher-dimensional optimization problem, i.e., when $d > 1$. There is some work on this particular notion of approximate fault-tolerance in high-dimensions, such as [47, 53], however their results only apply to special cost functions, specifically, quadratic or strictly convex functions, as opposed to the generic cost functions (that need not even be differentiable) considered in this paper.

Another way of measuring approximation is by the value of the aggregate cost function, or its gradient. For instance, as discussed in [18], for the case of differentiable cost functions a resilient distributed optimization algorithm $\Pi$ may output a point $x_\Pi \in \mathbb{R}^d$ such that each element of the aggregate non-faulty gradient $\sum_{i \in \mathcal{H}} \nabla Q_i(x_\Pi)$ is bounded by $\epsilon$. As yet another alternative, a resilient algorithm $\Pi$ may aim to output a point $x_\Pi$ such that the non-faulty aggregate cost $\sum_{i \in \mathcal{H}} Q_i(x_\Pi)$ is within $\epsilon$ of the true minimum cost $\min_x \sum_{i \in \mathcal{H}} Q_i(x)$. However, these definitions of approximate resilience are sensitive to scaling of the cost functions. In particular, if the elements of $\sum_{i \in \mathcal{H}} \nabla Q_i(x_\Pi)$ are bounded by $\epsilon$ then the elements of $\sum_{i \in \mathcal{H}} \alpha \nabla Q_i(x_\Pi)$ are bounded by $\alpha\epsilon$, where $\alpha$ is a positive scalar value. On the other hand, both $\sum_{i \in \mathcal{H}} Q_i(x)$ and $\sum_{i \in \mathcal{H}} \alpha Q_i(x)$ have identical minimum point regardless of the value of $\alpha$. Therefore, when the objective is to approximate a minimum point of the non-faulty aggregate cost $\arg\min_x \sum_{i \in \mathcal{H}} Q_i(x)$, which is indeed the case in this paper, use of function (or gradient) values to measure approximation is not a suitable choice.

## 2.2 Gradient-Filters

In the past, several gradient-filters have been proposed to *robustify* the distributed gradient-descent (DGD) method against Byzantine faulty agents in a server-based architecture, e.g., see [1, 6, 17, 18, 23, 40, 48, 54]. A gradient-filter refers to *Byzantine robust aggregation* of agents' gradients that mitigates the detrimental impact of incorrect gradients sent by the Byzantine agents to the server. To name a few gradient-filters, that are provably effective against Byzantine agents, we have the comparative gradient elimination (CGE) [23, 24], coordinate-wise trimmed mean (CWTM) [48, 54], geometric median-of-means (GMoM) [14], KRUM [6], Bulyan [20], and other spectral gradient-filters [18]. Different gradient-filters guarantee some fault-tolerance under different assumptions on non-faulty agents' cost functions.

In this paper, we propose a generic result, in Theorem 3 in Section 4, on the convergence of the DGD method equipped with a gradient-filter. The result holds true regardless of the gradient-filter used, and thus, can be utilized to obtain formal fault-tolerance property of a gradient-filter in context of the considered distributed optimization problem. We demonstrate this, in Section 4.2, by obtaining $(f, \epsilon)$-resilience properties of two specific gradient-filters; CGE and CWTM. As exact fault-tolerance is equivalent to $(f, 0)$-resilience (see Section 1.2), our results generalize the prior work on exact fault-tolerance of these two filters, see [23, 24, 47]. Moreover, until now, exact fault-tolerance of the CWTM gradient-filter was only studied for special optimization problems of state estimation [47], and machine learning [54]. Our result presents the fault-tolerance property of CWTM for a much larger

class of optimization problems.

## 2.3 Robust Statistics with Arbitrary Outliers

As noted earlier, there has been work on the problem of robust statistics with arbitrary outliers [12, 22, 46]. In this problem, we are given a finite set of data points; $\alpha$ fraction of which are sampled independently and identically from a common distribution $\mathcal{D}$ in $\mathbb{R}^d$, and the remaining $1 - \alpha$ fraction of data points may be arbitrary. The identity of arbitrary data points is a priori unknown, otherwise the problem is trivialized. The objective in this problem is to estimate statistical measures of distribution $\mathcal{D}$, such as mean, or variance, despite the presence of arbitrary outliers. The problem robust mean estimation can potentially be modelled as a fault-tolerant distributed optimization problem where for each non-faulty agent $i$, $Q_i(x) : (x, x_i) \mapsto y \in \mathbb{R}$ for all $x \in \mathbb{R}^d$ where $x_i \sim \mathcal{D}$. A faulty agent may choose an arbitrary cost function. The cost functions can be designed in a manner such that the minimum point of the aggregate of non-faulty cost functions is equal to the mean for the non-faulty data points. In particular, suppose that for each non-faulty agent $i$, $Q_i(x) \triangleq \|x - x_i\|^2$ where $x_i \sim \mathcal{D}$. In this case, the minimum point of the non-faulty aggregate cost function is equal to the average of the non-faulty data points sampled from distribution $\mathcal{D}$.

Prior work on robust statistics considers a centralized setting wherein, unlike a distributed setting, all the data points are accessible to a single machine. In this report, we also present distributed algorithms that do not require the agents to share their local data points. Moreover, in the centralized setting, our results are applicable to a larger class of cost functions, including non-convex functions.

## 2.4 Fault-tolerance in State Estimation

The problem of distributed optimization finds direct application in distributed state estimation [41]. In this problem, the system comprises multiple sensors, and each sensor makes partial observations about the system's state. The goal is to compute the entire state of the system using collective observations from all the sensors. However, if a sensor is faulty then it may share incorrect observations, preventing correct state estimation. The special case of distributed state estimation when the observations are *linear* in the system's state has gained significant attention in the past, e.g. see [5, 15, 34, 37, 38, 44, 45, 47]. These works have shown that the state can be determined despite up to $f$ (out of $n$) faulty observations *if and only if* the system is $2f$-*sparse observable*, i.e., the complete state can be determined using observations of only $n - 2f$ non-faulty sensors. We note that, in this particular case, $2f$-*sparse observability* is equivalent to $2f$-redundancy. Additionally, some of these works, such as [34, 47], also consider the case of *approximate* linear state estimation when the observations are noisy. Our work is more general in that we consider the problem setting of distributed optimization, and our results apply to a larger class of cost functions.

# 3 Necessary and Sufficient Conditions for $(f, \epsilon)$-Resilience

Throughout this paper we assume, as stated below, that the non-faulty agents' cost functions and their aggregates have well-defined minimum points. Otherwise, the problem of optimization is rendered vacuous.

**Assumption 1.** *For every non-empty set of non-faulty agents $S$, we assume that the set* $\arg\min_{x \in \mathbb{R}^d} \sum_{i \in S} Q_i(x)$ *is non-empty and closed.*

We also assume that $f < n/2$. Lemma 1 below shows that $(f, \epsilon)$-resilience is impossible in general when $f \geq n/2$. Proof of Lemma 1 is easy, and can be found in Appendix A.

**Lemma 1.** *If $f \geq n/2$ then there cannot exist a deterministic $(f, \epsilon)$-resilient algorithm for any $\epsilon \geq 0$.*

## 3.1 Necessary Condition

**Theorem 1.** *Suppose that Assumption 1 holds true. There exists a deterministic $(f, \epsilon)$-resilient distributed optimization algorithm where $\epsilon \geq 0$ only if the agents' cost functions satisfy the $(2f, \epsilon)$-redundancy property.*

*Proof.* To prove the theorem we present a scenario when the agents' cost functions (if non-faulty) are scalar functions, i.e., $d = 1$ and for all $i$, $Q_i : \mathbb{R} \to \mathbb{R}$, and the minimum point of an aggregate of one or more agents' cost functions is uniquely defined. Obviously, if a condition is necessary in this particular scenario then it is so in the general case involving vector functions with non-unique minimum points.

To prove the necessary condition, we also assume that the server has full knowledge of all the agents' cost functions. This may not hold true in practice, where instead the server may only have partial information about the agents' cost functions. Indeed, this assumption forces the Byzantine faulty agents to a priori fix their cost functions. However, in reality the Byzantine agents may send arbitrary information over time to the server that need not be consistent with a fixed cost function. Thus, necessity of $(2f, \epsilon)$-redundancy under this strong assumption implies its necessity in general.

The proof is by contradiction. Specifically, we show that *If the cost functions of non-faulty agents do not satisfy the $(2f, \epsilon)$-redundancy property then there cannot exist a deterministic $(f, \epsilon)$-resilient distributed optimization algorithm.*

Recall that we have assumed that for a non-empty set of agents $T$ the aggregate cost function $\sum_{i \in T} Q_i(x)$ has a unique minimum point. To be precise, for each non-empty subset of agents $T$, we define

$$x_T = \arg\min_x \sum_{i \in T} Q_i(x).$$

Suppose that the agents' cost functions **do not** satisfy the $(2f, \epsilon)$-redundancy property stated in Definition 3. Then, there exists a real number $\delta > 0$ and a pair of subsets $S, \widehat{S}$ with $\widehat{S} \subset S$, $|S| = n - f$, and $n - 2f \leq \left|\widehat{S}\right| < n - f$ such that

$$\left\|x_{\widehat{S}} - x_S\right\| \geq \epsilon + \delta. \tag{6}$$

Now, suppose that $n - f - \left|\widehat{S}\right|$ agents in the remainder set $\{1, \ldots, n\} \setminus S$ are Byzantine faulty. Let us denote the set of faulty agents by $\mathcal{B}$. Note that $\mathcal{B}$ is non-empty with $|\mathcal{B}| = n - f - \left|\widehat{S}\right| \leq f$. Similar to the non-faulty agents, the faulty agents send to the server cost functions that are scalar, and the aggregate of one or more agents' cost functions in the set $S \cup \mathcal{B}$ is unique. However, the aggregate cost function of the agents in the set $\mathcal{B} \cup \widehat{S}$ minimizes at a unique point $x_{\mathcal{B} \cup \widehat{S}}$ which is $\left\|x_{\widehat{S}} - x_S\right\|$ distance away from $x_{\widehat{S}}$, similar to $x_S$, but lies on the other side of $x_{\widehat{S}}$ as shown in the figure below. Note that it is always possible to pick such functions for the faulty agents.



Note that the distance between the two points $x_S$ and $x_{\mathcal{B} \cup \widehat{S}}$ is $2\epsilon + 2\delta$, i.e.,

$$\left\|x_S - x_{\mathcal{B} \cup \widehat{S}}\right\| = 2\epsilon + 2\delta. \tag{7}$$

Now, suppose, toward a contradiction, that there exists an $(f, \epsilon)$-resilient deterministic optimization algorithm named $\Pi$. As the identity of Byzantine faulty agents is a priori unknown to the server, and the cost functions sent by the Byzantine faulty agents have similar properties as the non-faulty agents, the server cannot distinguish between the following two possible scenarios; i) $S$ is the set of non-faulty agents, and ii) $\mathcal{B} \cup \widehat{S}$ is the set of non-faulty agents. Note that both the sets $S$ and $\mathcal{B} \cup \widehat{S}$ contain $n - f$ agents.

As the cost functions received by the server are identical in both of the above scenarios, being a deterministic algorithm, $\Pi$ should have identical output in both the cases. We let $\widehat{x}$ denote the output of $\Pi$. In scenario (i) when the set of honest agents is given by $S$ with $|S| = n - f$, as $\Pi$ is assumed $(f, \epsilon)$-resilient, by Definition 2 the output

$$\widehat{x} \in [x_S - \epsilon, \, x_S + \epsilon] \tag{8}$$

as shown in the figure above. Similarly, in scenario (ii) when the set of honest agents is $\mathcal{B} \cup \widehat{S}$ with $\left| \mathcal{B} \cup \widehat{S} \right| = n - f$,

$$\widehat{x} \in [x_{\mathcal{B} \cup \widehat{S}} - \epsilon, \, x_{\mathcal{B} \cup \widehat{S}} + \epsilon]. \tag{9}$$

However, (7) implies that (8) and (9) cannot be satisfied simultaneously. That is, if $\Pi$ is $(f, \epsilon)$-resilient in scenario (i) then it cannot be so in scenario (ii), and vice-versa. This contradicts the assumption that $\Pi$ is $(f, \epsilon)$-resilient. $\qquad \square$

## 3.2 Sufficient Condition

**Theorem 2.** *Suppose that Assumption 1 holds true. For a real value $\epsilon \geq 0$, if the agents' cost functions satisfy the $(2f, \epsilon)$-redundancy property then $(f, 2\epsilon)$-resilience is achievable.*

*Proof.* The proof is constructive where we assume that all the agents send their individual cost functions to the server. We assume that $f > 0$ to avoid the trivial case of $f = 0$. Throughout the proof we write the notation $\arg\min_{x \in \mathbb{R}^d}$ simply as $\arg\min$, unless otherwise stated. We begin by presenting an algorithm below, comprising three steps.

**Step 1:** Each agent sends their cost function to the server. An honest agent sends its actual cost function, while a faulty agent may send an arbitrary function.

**Step 2:** For each set $T$ of received functions, $|T| = n - f$, the server computes a point

$$x_T \in \arg\min \sum_{i \in T} Q_i(x).$$

For each subset $\widehat{T} \subset T$, $\left| \widehat{T} \right| = n - 2f$, the server computes

$$r_{T\widehat{T}} \triangleq \mathrm{dist}\left( x_T, \, \arg\min \sum_{i \in \widehat{T}} Q_i(x) \right), \tag{10}$$

and

$$r_T = \max_{\substack{\widehat{T} \subset T, \\ |\widehat{T}| = n - 2f}} r_{T\widehat{T}}. \tag{11}$$

**Step 3:** The server outputs $x_S$ such that

$$S = \arg\min_{\substack{T \subset \{1, \dots, n\}, \\ |T| = n - f}} r_T. \tag{12}$$

10

We show that above algorithm is $(f, 2\epsilon)$-resilient under $(2f, \epsilon)$-redundancy. For a non-empty set of agents $T$, we denote
$$X_T = \arg\min \sum_{i \in T} Q_i(x).$$

Consider an arbitrary set of non-faulty agents $G$ with $|G| = n - f$. Such a set is guaranteed to exist as there are at most $f$ faulty agents, and therefore, at least $n - f$ non-faulty agents exist in the system. Consider an arbitrary set $\widehat{T}$ such that $\widehat{T} \subset G$ and $\left|\widehat{T}\right| = n - 2f$. By Definition 3 of $(2f, \epsilon)$-redundancy,

$$\text{dist}\left(X_G, X_{\widehat{T}}\right) \le \epsilon. \tag{13}$$

Recall from (10) that $r_{G\widehat{T}} = \text{dist}\left(x_G, X_{\widehat{T}}\right)$. As $x_G \in X_G$, by Definition (4) of Hausdorff set distance, $\text{dist}\left(x_G, X_{\widehat{T}}\right) \le \text{dist}\left(X_G, X_{\widehat{T}}\right)$. Therefore, $r_{G\widehat{T}} \le \text{dist}\left(X_G, X_{\widehat{T}}\right)$, and substituting from (13) implies that

$$r_{G\widehat{T}} \le \epsilon. \tag{14}$$

Now, recall from (11) that $r_G = \max\left\{r_{G\widehat{T}} \,\Big|\, \widehat{T} \subset G, \left|\widehat{T}\right| = n - 2f\right\}$. As $\widehat{T}$ in (14) is an arbitrary subset of $G$ with $\left|\widehat{T}\right| = n - 2f$,

$$r_G = \max_{\substack{\widehat{T} \subset G, \\ |\widehat{T}| = n-2f}} r_{G\widehat{T}} \le \epsilon. \tag{15}$$

From (12) and (15) we obtain that

$$r_S \le r_G \le \epsilon. \tag{16}$$

As $|G| = n - f$, for every set of agents $T$ with $|T| = n - f$, $|T \cap G| \ge n - 2f$. Therefore, for the set $S$ defined in (12), there exists a subset $\widehat{G}$ of $G$ such that $\widehat{G} \subset S$ and $\left|\widehat{G}\right| = n - 2f$. For such a set $\widehat{G}$, by definition of $r_S$ in (11), we obtain that

$$r_{S\widehat{G}} \triangleq \text{dist}\left(x_S, X_{\widehat{G}}\right) \le r_S.$$

Substituting from (16) above, we obtain that

$$\text{dist}\left(x_S, X_{\widehat{G}}\right) \le \epsilon. \tag{17}$$

As $\widehat{G}$ is a subset of $G$, all the agents in $\widehat{G}$ are non-faulty. Therefore, by Assumption 1, $X_{\widehat{G}}$ is a closed set. Recall that $\text{dist}\left(x_S, X_{\widehat{G}}\right) = \inf_{x \in X_{\widehat{G}}} \|x_S - x\|$. The closedness of $X_{\widehat{G}}$ implies that there exists a point $z \in X_{\widehat{G}}$ such that

$$\|x_S - z\| = \inf_{x \in X_{\widehat{G}}} \|x_S - x\| = \text{dist}\left(x_S, X_{\widehat{G}}\right).$$

The above, in conjunction with (17), implies that

$$\|x_S - z\| \le \epsilon. \tag{18}$$

Moreover, as $z \in X_{\widehat{G}}$ where $\widehat{G} \subset G$ with $\left|\widehat{G}\right| = n - 2f$ and $|G| = n - f$, the $(2f, \epsilon)$-redundancy condition stated in Definition 3 implies that $\text{dist}\left(z, X_G\right) \le \epsilon$. Similar to an argument made above, under Assumption 1, $X_G$ is a closed set, and therefore, there exists $x^* \in X_G$ such that

$$\|z - x^*\| = \text{dist}\left(z, X_G\right) \le \epsilon. \tag{19}$$

By triangle inequality, (18) and (19) implies that $\|x_S - x^*\| \le \|x_S - z\| + \|z - x^*\| \le 2\epsilon$. Recall that set $G$ here is an arbitrary set of $n - f$ non-faulty agents. $\qquad \square$

It is worth noting that the algorithm constructed in the proof of Theorem 2 only shows sufficiency; it is not a very practical algorithm due to being computationally expensive.

In the next part of the paper, i.e., Sections 4 and 5, we consider the case when the (non-faulty) agents' cost functions are differentiable. Specifically, we study approximate fault-tolerance in the distributed gradient-descent (DGD) method.

# 4 Distributed Gradient-Descent (DGD) Method

In this section, we consider a setting wherein the non-faulty agents' cost functions are differentiable. In this particular case, we study the approximate fault-tolerance of the distributed gradient-descent method coupled with a *gradient-filter*, described below. We consider the server-based system architecture, shown in Fig. 1, assuming a synchronous system.

The DGD method is an iterative algorithm wherein the server maintains an *estimate* of a minimum point, and updates it iteratively using gradients sent by the agents. Specifically, in each iteration $t \in \{0, 1, \ldots\}$, the server starts with an estimate $x^t$ and broadcasts to all the agents. Each non-faulty agent $i$ sends back to the sever the gradient of its cost function at $x^t$, i.e., $\nabla Q_i(x^t)$. However, Byzantine faulty agents may send arbitrary incorrect vectors as their gradients to the server. The initial estimate, named $x^0$, is chosen arbitrarily by the server.

A *gradient-filter* is a vector function, denoted by GradFilter, that maps the $n$ gradients received by the server from all the $n$ agents to a $d$-dimensional vector, i.e., GradFilter : $\mathbb{R}^{d \times n} \to \mathbb{R}^d$. For example, an average of all the gradients as in the case of the traditional distributed gradient-descent method is technically a gradient-filter. However, averaging is not quite robust against Byzantine faulty agents [6, 48]. The real purpose of a gradient-filter is to mitigate the detrimental impact of incorrect gradients sent by the Byzantine faulty agents. In other words, a gradient-filter *robustifies* the traditional gradient-descent method against Byzantine faults. We show that if a gradient-filter satisfies a certain property then it can confer fault-tolerance to the distributed gradient-descent method.

We first formally describe below the steps in each iteration of the distributed gradient-descent method implemented on a synchronous server-based system. Note that we constrain the estimates computed by the server to a compact convex set $\mathcal{W} \subset \mathbb{R}^d$. The set $\mathcal{W}$ can be arbitrarily large. For a vector $x \in \mathbb{R}^d$, its projection onto $\mathcal{W}$, denoted by $[x]_{\mathcal{W}}$, is defined to be

$$[x]_{\mathcal{W}} = \arg \min_{y \in \mathcal{W}} \|x - y\|. \tag{20}$$

As $\mathcal{W}$ is a convex and compact set, $[x]_{\mathcal{W}}$ is unique for each $x$ (see [10]).

## 4.1 Steps in $t$-th iteration

In each iteration $t \in \{0, 1, \ldots\}$ the server updates its current estimate $x^t$ to $x^{t+1}$ using Steps S1 and S2 described as follows.

**S1:** The server requests from each agent the gradient of its local cost function at the current estimate $x^t$. Each non-faulty agent $i$ will then send to the server the gradient $\nabla Q_i(x^t)$, whereas a faulty agent may send an incorrect arbitrary value for the gradient.

The gradient received by the server from agent $i$ is denoted as $g_i^t$. If no gradient is received from some agent $i$, agent $i$ must be faulty (because the system is assumed to be synchronous) – in this case, the server eliminates the agent $i$ from the system, updates the values of $n$, $f$, and re-assigns the agents indices from 1 to $n$.

**S2:** **[Gradient-filtering]** The server applies a gradient-filter GradFilter to the $n$ received gradients and computes

GradFilter $(g_1^t, \ldots, g_n^t) \in \mathbb{R}^d$. Then, the server updates its estimate to

$$x^{t+1} = \left[ x^t - \eta_t \, \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right]_{\mathcal{W}} \tag{21}$$

where $\eta_t$ is the step-size of positive value for iteration $t$.

We propose, in Theorem 3, a generic convergence result for the above algorithm.

**Theorem 3.** *Consider the update rule* (21) *in the above iterative algorithm, with diminishing step-sizes* $\{\eta_t, \, t = 0, 1, \ldots\}$ *satisfying* $\sum_{t=0}^{\infty} \eta_t = \infty$ *and* $\sum_{t=0}^{\infty} \eta_t^2 < \infty$. *Suppose that*

$$\left\| \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\| < \infty$$

*for all* $t$. *For some point* $x^* \in \mathcal{W}$, *if there exists real-valued constants* $\mathsf{D}^* \in [0, \max_{x \in \mathcal{W}} \|x - x^*\|)$ *and* $\xi > 0$ *such that for each iteration* $t$,

$$\phi_t = \left\langle x^t - x^*, \, \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\rangle \geq \xi \text{ when } \left\| x^t - x^*\right\| \geq \mathsf{D}^*, \tag{22}$$

*then* $\lim_{t \to \infty} \|x^t - x^*\| \leq \mathsf{D}^*$.

The values $\mathsf{D}^*$ and $\xi$ in Theorem 3 may be interdependent. Proof of Theorem 3 is deferred to Appendix D.

Using Theorem 3 we can obtain conditions under which a gradient-filter guarantees the approximate fault-tolerance property of $(f, \epsilon)$-resilience with $\epsilon \geq 0$, of which exact fault-tolerance is a special case. On the other hand, the prior results on the convergence of DGD method with a gradient-filter, e.g., see [6, 20], apply only to exact fault-tolerance.

We demonstrate below the utility of Theorem 3 to obtain the fault-tolerance properties of two commonly used gradient-filters in the literature; namely *Comparative Gradient Elimination* [23] and *Coordinate-Wise Trimmed Mean* [47].

## 4.2 Gradient-Filters and their Fault-Tolerance Properties

In this subsection, we present precise approximate fault-tolerance properties of two specific gradient-filters; the Comparative Gradient Elimination (CGE) [23, 24], and the Coordinate-Wise Trimmed Mean (CWTM) [47, 54]. Note that differentiability of non-faulty agents' cost functions, which is already assumed for the DGD method, implies Assumption 1 (see [10]). We additionally make Assumptions 2, 3 and 4 about the non-faulty agents' cost functions. Similar assumptions are made in prior work on fault-free distributed optimization [4, 9, 36].

**Assumption 2** (Lipschitz smoothness)**.** *For each non-faulty agent* $i$, *we assume that the gradient of its cost function* $\nabla Q_i(x)$ *is Lipschitz continuous, i.e., there exists a finite real value* $\mu > 0$ *such that*

$$\|\nabla Q_i(x) - \nabla Q_i(x')\| \leq \mu \|x - x'\|, \quad \forall x, x' \in \mathcal{W}.$$

**Assumption 3** (Strong convexity)**.** *For a non-empty set of non-faulty agents* $\mathcal{H}$, *let* $Q_{\mathcal{H}}(x)$ *denote the average cost function of the agents in* $\mathcal{H}$, *i.e.,*

$$Q_{\mathcal{H}}(x) = \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} Q_i(x).$$

*For each such set* $\mathcal{H}$ *with* $|\mathcal{H}| = n - f$, *we assume that* $Q_{\mathcal{H}}(x)$ *is strongly convex, i.e., there exists a finite real value* $\gamma > 0$ *such that*

$$\langle \nabla Q(x) - \nabla Q(x'), \, x - x' \rangle \geq \gamma \|x - x'\|^2, \quad \forall x, x' \in \mathcal{W}.$$

13

Note that, under Assumptions 2 and 3, $\gamma \leq \mu$. This inequality is proved in Appendix C. Now, recall that the iterative estimates of the algorithm in Section 4.1 are constrained to a compact convex set $\mathcal{W} \subset \mathbb{R}^d$.

**Assumption 4** (Existence). *For each set of non-faulty agents $\mathcal{H}$ with $|\mathcal{H}| = n - f$, we assume that there exists a point $x_{\mathcal{H}} \in \arg\min_{x \in \mathbb{R}^d} \sum_{i \in \mathcal{H}} Q_i(x)$ such that $x_{\mathcal{H}} \in \mathcal{W}$.*

We describe below the CGE and CWTM gradient-filters. Later, we obtain the fault-tolerance properties of these filters using the result stated in Theorem 3, under $(2f, \epsilon)$-redundancy.

**CGE Gradient-Filter:** To apply the CGE gradient-filter in Step S2, the server sorts the $n$ gradients received from the $n$ agents at the completion of Step S1 as per their Euclidean norms (ties broken arbitrarily):

$$\left\| g_{i_1}^t \right\| \leq \ldots \leq \left\| g_{i_{n-f}}^t \right\| \leq \left\| g_{i_{n-f+1}}^t \right\| \leq \ldots \leq \left\| g_{i_n}^t \right\|.$$

That is, the gradient with the smallest norm, $g_{i_1}^t$, is received from agent $i_1$, and the gradient with the largest norm, $g_{i_n}^t$, is received from agent $i_n$. Then, the output of the CGE gradient-filter is the vector sum of the $n - f$ gradients with smallest $n - f$ Euclidean norms. Specifically,

$$\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right) = \sum_{j=1}^{n-f} g_{i_j}^t. \tag{23}$$

**CWTM Gradient-Filter:** To implement this particular gradient-filter in Step S2, the server sorts the $n$ gradients received from the $n$ agents at the completion of Step S1 as per their individual elements. For a vector $v \in \mathbb{R}^d$, we let $v[k]$ denote its $k$-th element. Specifically, for each $k \in \{1, \ldots, d\}$, the server sorts the $k$-th elements of the gradients by breaking ties arbitrarily:

$$g_{i_1[k]}^t[k] \leq \ldots \leq g_{i_{f+1}[k]}^t[k] \leq \ldots \leq g_{i_{n-f}[k]}^t[k] \leq \ldots \leq g_{i_n[k]}^t[k].$$

The gradient with the smallest of the $k$-th element, $g_{i_1[k]}^t$, is received from agent $i_1[k]$, and the gradient with the largest of the $k$-th element, $g_{i_n[k]}^t$, is received from agent $i_n[k]$. For each $k$, the server eliminates the largest $f$ and the smallest $f$ of the $k$-th elements of the gradients received. Then, the output of the CWTM gradient-filter is a vector whose $k$-th element is equal to the average of the remaining $n - 2f$ gradients' $k$-th elements. That is, for each $k \in \{1, \ldots, d\}$,

$$\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)[k] = \frac{1}{n - 2f} \sum_{j=f+1}^{n-f} g_{i_j[k]}^t[k]. \tag{24}$$

We present the precise fault-tolerance properties of the two gradient-filters in Theorems 4 and 5 below. However, the **reader may skip to Section 5 without loss of continuity**. Proofs of the theorems are deferred to Appendices E and F, respectively.

Note that, under Assumptions 3 and 4, for each non-empty set of non-faulty agents $\mathcal{H}$ with $|\mathcal{H}| = n - f$, the aggregate cost function $\sum_{i \in \mathcal{H}} Q_i(x)$ has a unique minimum point, denoted by $x_{\mathcal{H}}$, in the set $\mathcal{W}$. Specifically,

$$\{x_{\mathcal{H}}\} = \mathcal{W} \cap \arg\min_{x \in \mathbb{R}^d} \sum_{i \in \mathcal{H}} Q_i(x). \tag{25}$$

We first show below, in Theorem 4, that when the fraction of Byzantine faulty agents $f/n$ is bounded then the DGD method with the CGE gradient-filter is $(f, \mathcal{O}(\epsilon))$-resilient, under $(2f, \epsilon)$-redundancy and the above assumptions.

**Theorem 4.** *Suppose that the non-faulty agents' cost functions satisfy the $(2f, \epsilon)$-redundancy property, and the Assumptions 2, 3 and 4 hold true. Consider the algorithm in Section 4.1 with the CGE gradient-filter defined in* (23). *The following holds true:*

1. $\|\mathsf{GradFilter}\,(g_1^t, \ldots, g_n^t)\| < \infty$ *for all $t$.*

2. *If*

$$\alpha = 1 - \frac{f}{n}\left(1 + \frac{2\mu}{\gamma}\right) > 0 \tag{26}$$

*then for each set of $n - f$ non-faulty agents $\mathcal{H}$, for each $\delta > 0$,*

$$\phi_t = \langle x^t - x_{\mathcal{H}},\, \mathsf{GradFilter}\,(g_1^t, \ldots, g_n^t)\rangle \geq \alpha n \gamma \delta \left(\left(\frac{4\mu f}{\alpha\gamma}\right)\epsilon + \delta\right)$$

*when* $\|x^t - x^*\| \geq \left(\frac{4\mu f}{\alpha\,\gamma}\right)\epsilon + \delta.$

Let $\mathcal{H}$ denote an arbitrary set of $n - f$ non-faulty agents. If the step-size $\eta_t$ in (21) is diminishing, i.e., $\sum_{t=0}^{\infty} \eta_t = \infty$ and $\sum_{t=0}^{\infty} \eta_t^2 < \infty$, then Theorem 4, in conjunction with Theorem 3 implies that, under the said conditions,

$$\lim_{t\to\infty} \|x^t - x_{\mathcal{H}}\| \leq \left(\frac{4\mu f}{\alpha\,\gamma}\right)\epsilon + \delta, \quad \forall \delta > 0.$$

The above implies that $\lim_{t\to\infty} \|x^t - x_{\mathcal{H}}\| \leq (4\mu f/\alpha\gamma)\,\epsilon$. Thus, Theorem 4 shows that under $(2f, \epsilon)$-redundancy, and Assumptions 2, 3 and 4, if $\alpha > 0$, or the fraction of Byzantine faulty agents $f/n$ is less than $1/(1 + 2(\mu/\gamma))$, then the DGD method with the CGE gradient-filter is asymptotically $(f, \mathsf{D}\epsilon)$-*resilient* (by Definition 2) where

$$\mathsf{D} = \frac{4\mu f}{\alpha\gamma} = \frac{4\mu\,n}{(n/f)\,\gamma - (\gamma + 2\mu)}. \tag{27}$$

A smaller number $f$ of Byzantine faulty agents implies a smaller value of $\mathsf{D}$, and therefore, better fault-tolerance of the algorithm. Moreover, $\mathsf{D} = 0$ when $f = 0$, i.e., the algorithm indeed converges to the actual minimum point of all the agents' aggregate cost function in the fault-free case. Note that under Assumptions 2 and 3, $\gamma \leq \mu$. So, the fault-tolerance guarantee of the CGE gradient-filter, presented in Theorem 4, requires $f/n < 1/3$, or $f < n/3$.

Next, we show that when the separation between the gradients of the non-faulty agents' cost functions is sufficiently small then the CWTM gradient-filter can guarantee some approximate fault-tolerance under $(2f, \epsilon)$-redundancy. To present the fault-tolerance of the CWTM gradient-filter, we make the following additional assumption.

**Assumption 5.** *For two non-faulty agents $i$ and $j$, we assume that there exists $\lambda > 0$ such that for all $x \in \mathcal{W}$,*

$$\|\nabla Q_i(x) - \nabla Q_j(x)\| \leq \lambda \max\{\|\nabla Q_i(x)\|, \|\nabla Q_j(x)\|\}.$$

Due to the triangle triangle inequality, Assumption 5 trivially holds true when $\lambda = 2$. However, we can presently guarantee fault-tolerance of CWTM gradient-filter when $\lambda < \gamma/(\mu\sqrt{d})$ where $\mu$ and $\gamma$ are the Lipschitz smoothness and strong convexity coefficients, respectively defined in Assumption 2 and 3. Recall the definition of point $x_{\mathcal{H}} \in \mathbb{R}^d$ from (25) where $\mathcal{H}$ denotes an arbitrary set of $n - f$ non-faulty agents.

**Theorem 5.** *Suppose that the non-faulty agents' cost functions satisfy the $(2f, \epsilon)$-redundancy property, and the Assumptions 2, 3, 4 and 5 hold true. Consider the algorithm in Section 4.1 with the CWTM gradient-filter defined in* (24). *The following holds true:*

1. $\|\mathsf{GradFilter}\,(g_1^t, \ldots, g_n^t)\| < \infty$ *for all t.*

2. *If $\lambda < \gamma/(\mu\sqrt{d})$ then for each set of $n - f$ non-faulty agents $\mathcal{H}$, for each $\delta > 0$,*

$$\phi_t = \left\langle x^t - x_{\mathcal{H}}, \, \mathsf{GradFilter}\,(g_1^t, \ldots, g_n^t) \right\rangle \geq \left(2\sqrt{d}n\mu\lambda\epsilon + \left(\gamma - \sqrt{d}\lambda\mu\right)\delta\right)\delta$$

*when $\|x^t - x_{\mathcal{H}}\| \geq \dfrac{2\sqrt{d}n\mu\lambda}{(\gamma - \sqrt{d}\mu\lambda)}\epsilon + \delta.$*

By similar arguments as in the case of CGE, Theorem 5, in conjunction with Theorem 3, implies that the DGD method with CWTM gradient-filter and diminishing step-sizes is asymptotically $(f,\ \mathsf{D}'\,\epsilon)$-*resilient* where

$$\mathsf{D}' = \frac{2\sqrt{d}n\mu\lambda}{(\gamma - \sqrt{d}\mu\lambda)} = \left(\frac{2n}{(\gamma/\mu\lambda\sqrt{d}) - 1}\right),$$

under the conditions stated in Theorem 5. The smaller the value of $\lambda$ is, i.e., the closer non-faulty gradients to each other are, the smaller is the value of $\mathsf{D}'$, and therefore, better is the approximate fault-tolerance guarantee of the CWTM gradient-filter. Unlike the CGE gradient-filter, resilience of CWTM presented in Theorem 5 is independent of $f$, as long as $\lambda < \gamma/(\mu\sqrt{d})$. However, the condition on $\lambda$ to guarantee the resilience of CWTM gradient-filter depends upon the dimension $d$ of the optimization problem. Larger dimension result in a tighter bound on $\lambda$.

## 5    Numerical Experiments

We present simulation results for an empirical evaluation of the CGE and CWTM gradient-filters applied to the problem of *distributed linear regression* [2]. More details about the experiments can be found in Appendix G.

We consider the synchronous server-based system in Figure 1. We assume that $n = 6$ and $f = 1$. Each agent $i$ knows a row vector $A_i$ of dimension $d = 2$. Each agent $i$ makes a real-valued (scalar) observation denoted by $B_i$ such that $B_i = A_i x^* + N_i$, where $x^* = (1,1)^T$ for all $i$, and $N_i$ is a randomly chosen noise. The value of $A_i$, $B_i$ and $N_i$ are omitted here for brevity. To solve the linear regression problem distributedly, each agent $i$'s cost function is defined as $Q_i(x) = (B_i - A_i x)^2$. For a non-empty set of agents $S$, we denote by $A_S$ a matrix of dimension $|S| \times 2$ obtained by stacking rows $\{A_i,\ i \in S\}$. Similarly, we obtain column vector $B_S$ by stacking the values $\{B_i,\ i \in S\}$. Thus for every such non-empty set $S$, $Q_S(x) \triangleq \sum_{i \in S}(B_i - A_i x)^2 = \|B_S - A_S x\|^2$. The rows $A_1, \ldots, A_n$ are chosen specifically to ensure that the system has $2f$-redundancy if $N_i = 0$ for all $i$. That is, each matrix $A_S$ with $|S| \geq n - 2f = 4$ is column full-rank or $\mathrm{rank}\,(A_S) = d = 2$. Consequentially, the cost function $Q_S(x)$ has a unique minimum point for each set $S$ with $|S| \geq 4$.

We simulate the distributed gradient-descent algorithm described in Section 4 by assuming agent 1 to be Byzantine faulty, i.e., the set of honest agents is $\mathcal{H} = \{2, 3, 4, 5, 6\}$. The minimum point of $\sum_{i \in \mathcal{H}} Q_i(x)$, denoted by $x_{\mathcal{H}}$, can be obtained by solving $B_{\mathcal{H}} = A_{\mathcal{H}} x$. Specifically, $x_{\mathcal{H}} = (1.0780, 0.9825)^T$. The goal of fault-tolerant distributed linear regression is to estimate $x_{\mathcal{H}}$. In our simulations, it can be verified that the agents' cost functions satisfy the $(2f,\ \epsilon)$-*redundancy* property, stated in Definition 3, with $\epsilon = 0.0890$. It can also be verified that the non-faulty agents' cost functions satisfy Assumptions 2 and 3 with $\mu = 2$ and $\gamma = 0.712$, respectively. We simulate the following fault behaviors for the Byzantine agents:

16

Table 1: *For the distributed linear regression problem, our algorithm's outputs with gradient-filters CGE and CWTM, and the approximation errors, corresponding to executions when the faulty agent 1 exhibits two different types of Byzantine faults; gradient-reverse and random.*

| | gradient-reverse | | random | |
|---|---|---|---|---|
| | $x_{\text{out}}$ | dist $(x_{\mathcal{H}}, x_{\text{out}})$ | $x_{\text{out}}$ | dist $(x_{\mathcal{H}}, x_{\text{out}})$ |
| **CGE** | $\begin{pmatrix} 1.0541 \\ 0.9826 \end{pmatrix}$ | 0.0239 | $\begin{pmatrix} 1.0779 \\ 0.9826 \end{pmatrix}$ | $4.72 \times 10^{-5}$ |
| **CWTM** | $\begin{pmatrix} 1.0645 \\ 0.9924 \end{pmatrix}$ | 0.0167 | $\begin{pmatrix} 1.0775 \\ 0.9840 \end{pmatrix}$ | $1.51 \times 10^{-3}$ |

- *gradient-reverse*: the faulty agent *reverses* its true gradient. Suppose the correct gradient of a faulty agent $i$ at step $t$ is $s_i^t$, the agent $i$ will send the incorrect gradient $g_i^t = -s_i^t$ to the server.

- *random*: the faulty agent sends a randomly chosen vector in $\mathbb{R}^d$. In our experiments, the faulty agent in each step chooses i.i.d. Gaussian random vector with mean 0 and an isotropic covariance matrix of standard deviation 200.

In the simulations, we apply a diminishing step size $\eta_t$, and a convex compact $\mathcal{W}$ as described in previous sections. For comparison purpose, all experiments have the same initial estimate $x^0 = (-0.0085, -0.5643)^T$. In every execution, the estimates practically converge after 400 iterations. Thus, we document the output of the algorithm to be $x_{\text{out}} = x^{500}$. The outputs for the two gradient-filters, CGE and CWTM, under different faulty behaviors, are shown in Table 1. Note that dist $(x_{\mathcal{H}}, x_{\text{out}}) = \|x_{\mathcal{H}} - x_{\text{out}}\|$. In all the executions, the distance between $\|x_{\mathcal{H}} - x_{\text{out}}\| < \epsilon$.

For the said executions, we plot in Figure 2 the values of the aggregate cost function $\sum_{i \in \mathcal{H}} Q_i(x^t)$ (referred as *loss*) and the approximation error $\|x^t - x_{\mathcal{H}}\|$ (referred as *distance*) for iteration $t$ ranging from 0 to 500. We also show the plots of the fault-free DGD method where the faulty agent is omitted, and the DGD method without any gradient-filter when agent 1 is Byzantine faulty. The details for iteration $t$ ranging from 0 to 80 are also highlighted in Figure 3.

We also conducted experiments for distributed learning with support vector machine with faulty agents in the distributed learning system (see Section 1.3). We observed that the DGD method with the said gradient-filters reaches comparable performance to the fault-free case, and that the accuracy of the learning process depends upon the correlation between the data points of non-faulty agents. Details of those results can be found in Appendix H.

# 6 Summary

We have considered the problem of *approximate* Byzantine fault-tolerance – a generalization of the *exact* fault-tolerance problem studied in prior work [26]. Unlike the exact fault-tolerance, the goal in approximate fault-tolerance is to design a distributed optimization algorithm that approximates a minimum point of the aggregate cost function of (at least $n - f$) non-faulty agents, in the presence of up to $f$ (out of $n$) Byzantine faulty agents. We have defined approximate fault-tolerance formally as $(f, \epsilon)$-resilience where $\epsilon \in \mathbb{R}_{\geq 0}$ represents the approximation error. In the first part of the paper, i.e, Section 3, we have obtained necessary and sufficient conditions for achieving $(f, \epsilon)$-resilience. In the second part of the paper, i.e., Sections 4 and 5, we have considered the case when agents' cost functions are differentiable. In this particular case, we have obtained a generic approximate fault-tolerance
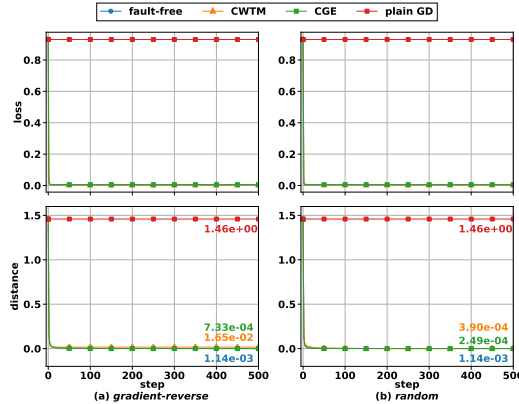
Figure 2: The losses, i.e., $\sum_{i \in \mathcal{H}} Q_i(x^t)$, and distances, i.e., $\|x^t - x_{\mathcal{H}}\|$, versus the number of iterations in the algorithm. The final approximation errors, i.e., $\|x^{500} - x_{\mathcal{H}}\|$, are annotated in the same colors of their corresponding plots. For the executions shown, agent 1 is assumed to be Byzantine faulty. The different columns show the results when the faulty agent exhibits the different types of faults: *(a)* gradient-reverse, and *(b)* random. Apart from the plots with CGE (in *green*) and CWTM (in *yellow*) gradient-filters, we also plot the fault-free DGD method where the faulty agent is omitted (in *blue*), and the DGD method without any gradient-filters when agent 1 is Byzantine faulty (in *red*).

property of the distributed gradient-descent method equipped with Byzantine robust gradient aggregation or *gradient-filter*, and have demonstrated the utility of this property by considering two specific well-known gradient-filters; comparative gradient elimination and coordinate-wise trimmed mean. In Section 5, we have demonstrated the applicability of our results through experiments.

# Acknowlegements

# References

[1] ALISTARH, D., ALLEN-ZHU, Z., AND LI, J. Byzantine stochastic gradient descent. In *Advances in Neural Information Processing Systems* (2018), S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds., vol. 31, Curran Associates, Inc.

[2] AMEMIYA, T. *Advanced econometrics*. Harvard university press, 1985.

[3] BERNSTEIN, J., ZHAO, J., AZIZZADENESHELI, K., AND ANANDKUMAR, A. signsgd with majority vote is communication efficient and fault tolerant, 2019.

[4] BERTSEKAS, D. P., AND TSITSIKLIS, J. N. *Parallel and distributed computation: numerical methods*, vol. 23. Prentice hall Englewood Cliffs, NJ, 1989.

Figure 3: The losses, i.e., $\sum_{i \in \mathcal{H}} Q_i(x^t)$, and distances, i.e., $\|x^t - x_{\mathcal{H}}\|$, versus the number of iterations in the algorithm, magnified for the initial 80 iterations in the training process. The interpretation of the plots is same as that in Figure 2.

[5] BHATIA, K., JAIN, P., AND KAR, P. Robust regression via hard thresholding. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1* (Cambridge, MA, USA, 2015), NIPS'15, MIT Press, p. 721–729.

[6] BLANCHARD, P., EL MHAMDI, E. M., GUERRAOUI, R., AND STAINER, J. Machine learning with adversaries: Byzantine tolerant gradient descent. In *Proceedings of the 31st International Conference on Neural Information Processing Systems* (Red Hook, NY, USA, 2017), NIPS'17, Curran Associates Inc., p. 118–128.

[7] BOTTOU, L. Online learning and stochastic approximations. *On-line learning in neural networks 17*, 9 (1998), 142.

[8] BOTTOU, L., CURTIS, F. E., AND NOCEDAL, J. Optimization methods for large-scale machine learning. *Siam Review 60*, 2 (2018), 223–311.

[9] BOYD, S., PARIKH, N., CHU, E., PELEATO, B., AND ECKSTEIN, J. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning 3*, 1 (Jan. 2011), 1–122.

[10] BOYD, S., AND VANDENBERGHE, L. *Convex optimization.* Cambridge university press, 2004.

[11] CAO, X., AND LAI, L. Distributed gradient descent algorithm robust to an arbitrary number of byzantine attackers. *IEEE Transactions on Signal Processing 67*, 22 (2019), 5850–5864.

[12] CHARIKAR, M., STEINHARDT, J., AND VALIANT, G. Learning from untrusted data. In *Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing* (New York, NY, USA, 2017), STOC 2017, Association for Computing Machinery, p. 47–60.

[13] CHEN, Y., KAR, S., AND MOURA, J. M. F. Resilient distributed estimation through adversary detection. *IEEE Transactions on Signal Processing 66*, 9 (2018), 2455–2469.

[14] CHEN, Y., SU, L., AND XU, J. Distributed statistical machine learning in adversarial settings: Byzantine gradient descent. *Proceedings of the ACM on Measurement and Analysis of Computing Systems 1*, 2 (2017), 44.

[15] CHONG, M. S., WAKAIKI, M., AND HESPANHA, J. P. Observability of linear systems under adversarial attacks. In *American Control Conference* (2015), IEEE, pp. 2439–2444.

[16] DALCIN, L. D., PAZ, R. R., KLER, P. A., AND COSIMO, A. Parallel distributed computing using python. *Advances in Water Resources 34*, 9 (2011), 1124–1139.

[17] DAMASKINOS, G., EL MHAMDI, E. M., GUERRAOUI, R., PATRA, R., AND TAZIKI, M. Asynchronous Byzantine machine learning (the case of SGD). In *Proceedings of the 35th International Conference on Machine Learning* (10–15 Jul 2018), J. Dy and A. Krause, Eds., vol. 80 of *Proceedings of Machine Learning Research*, PMLR, pp. 1145–1154.

[18] DIAKONIKOLAS, I., KAMATH, G., KANE, D. M., LI, J., STEINHARDT, J., AND STEWART, A. Sever: A robust meta-algorithm for stochastic optimization, 2019.

[19] DUCHI, J. C., AGARWAL, A., AND WAINWRIGHT, M. J. Dual averaging for distributed optimization: Convergence analysis and network scaling. *IEEE Transactions on Automatic control 57*, 3 (2011), 592–606.

[20] EL MHAMDI, E. M., GUERRAOUI, R., AND ROUAULT, S. The hidden vulnerability of distributed learning in Byzantium. In *Proceedings of the 35th International Conference on Machine Learning* (10–15 Jul 2018), J. Dy and A. Krause, Eds., vol. 80 of *Proceedings of Machine Learning Research*, PMLR, pp. 3521–3530.

[21] FAWZI, H., TABUADA, P., AND DIGGAVI, S. Secure estimation and control for cyber-physical systems under adversarial attacks. *IEEE Transactions on Automatic control 59*, 6 (2014), 1454–1467.

[22] FENG, J., XU, H., AND MANNOR, S. Distributed robust learning, 2015.

[23] GUPTA, N., LIU, S., AND VAIDYA, N. H. Byzantine fault-tolerant distributed machine learning using stochastic gradient descent (sgd) and norm-based comparative gradient elimination (cge), 2021.

[24] GUPTA, N., AND VAIDYA, N. H. Byzantine fault tolerant distributed linear regression, 2019.

[25] GUPTA, N., AND VAIDYA, N. H. Byzantine fault-tolerant parallelized stochastic gradient descent for linear regression. In *2019 57th Annual Allerton Conference on Communication, Control, and Computing (Allerton)* (2019), IEEE, pp. 415–420.

[26] GUPTA, N., AND VAIDYA, N. H. Fault-tolerance in distributed optimization: The case of redundancy. In *Proceedings of the 39th Symposium on Principles of Distributed Computing* (New York, NY, USA, 2020), PODC '20, Association for Computing Machinery, p. 365–374.

[27] GUPTA, N., AND VAIDYA, N. H. Resilience in collaborative optimization: Redundant and independent cost functions, 2020.

[28] KARIMIREDDY, S. P., HE, L., AND JAGGI, M. Learning from history for byzantine robust optimization, 2020.

[29] KUWARANANCHAROEN, K., XIN, L., AND SUNDARAM, S. Byzantine-resilient distributed optimization of multi-dimensional functions. In *2020 American Control Conference (ACC)* (2020), IEEE, pp. 4399–4404.

[30] LAMPORT, L., SHOSTAK, R., AND PEASE, M. *The Byzantine Generals Problem.* Association for Computing Machinery, New York, NY, USA, 2019, p. 203–226.

[31] LeCun, Y., Bottou, L., Bengio, Y., and Haffner, P. Gradient-based learning applied to document recognition. *Proceedings of the IEEE 86*, 11 (1998), 2278–2324.

[32] Liu, S., Gupta, N., and Vaidya, N. H. Approximate byzantine fault-tolerance in distributed optimization. In *Proceedings of the 2021 ACM Symposium on Principles of Distributed Computing* (New York, NY, USA, 2021), PODC '21, Association for Computing Machinery.

[33] Lynch, N. A. *Distributed algorithms.* Elsevier, 1996.

[34] Mishra, S., Shoukry, Y., Karamchandani, N., Diggavi, S. N., and Tabuada, P. Secure state estimation against sensor attacks in the presence of noise. *IEEE Transactions on Control of Network Systems 4*, 1 (2016), 49–59.

[35] Munkres, J. R. *Topology.* Prentice Hall Upper Saddle River, NJ, 2000.

[36] Nedic, A., and Ozdaglar, A. Distributed subgradient methods for multi-agent optimization. *IEEE Transactions on Automatic Control 54*, 1 (2009), 48–61.

[37] Pajic, M., Lee, I., and Pappas, G. J. Attack-resilient state estimation for noisy dynamical systems. *IEEE Transactions on Control of Network Systems 4*, 1 (2017), 82–92.

[38] Pajic, M., Weimer, J., Bezzo, N., Tabuada, P., Sokolsky, O., Lee, I., and Pappas, G. J. Robustness of attack-resilient state estimators. In *ICCPS'14: ACM/IEEE 5th International Conference on Cyber-Physical Systems (with CPS Week 2014)* (2014), IEEE, pp. 163–174.

[39] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al. Pytorch: An imperative style, high-performance deep learning library. *arXiv preprint arXiv:1912.01703* (2019).

[40] Prasad, A., Suggala, A. S., Balakrishnan, S., and Ravikumar, P. Robust estimation via robust gradient estimation, 2018.

[41] Rabbat, M., and Nowak, R. Distributed optimization in sensor networks. In *Proceedings of the 3rd international symposium on Information processing in sensor networks* (2004), IEEE, pp. 20–27.

[42] Raffard, R. L., Tomlin, C. J., and Boyd, S. P. Distributed optimization for cooperative agents: Application to formation flight. In *2004 43rd IEEE Conference on Decision and Control (CDC)* (2004), vol. 3, IEEE, pp. 2453–2459.

[43] Rudin, W. *Principles of mathematical analysis*, vol. 3. McGraw-hill New York, 1964.

[44] Shoukry, Y., Nuzzo, P., Puggelli, A., Sangiovanni-Vincentelli, A. L., Seshia, S. A., Srivastava, M., and Tabuada, P. Imhotep-smt: A satisfiability modulo theory solver for secure state estimation. In *Proc. Int. Workshop on Satisfiability Modulo Theories* (2015).

[45] Shoukry, Y., Nuzzo, P., Puggelli, A., Sangiovanni-Vincentelli, A. L., Seshia, S. A., and Tabuada, P. Secure state estimation for cyber-physical systems under sensor attacks: A satisfiability modulo theory approach. *IEEE Transactions on Automatic Control 62*, 10 (2017), 4917–4932.

[46] Steinhardt, J., Charikar, M., and Valiant, G. Resilience: A criterion for learning in the presence of arbitrary outliers, 2017.

[47] Su, L., and Shahrampour, S. Finite-time guarantees for byzantine-resilient distributed state estimation with noisy measurements, 2018.

[48] Su, L., and Vaidya, N. H. Fault-tolerant multi-agent optimization: Optimal iterative distributed algorithms. PODC '16, Association for Computing Machinery, p. 425–434.

[49] Su, L., and Vaidya, N. H. Non-bayesian learning in the presence of byzantine agents. In *Distributed Computing* (Berlin, Heidelberg, 2016), Springer Berlin Heidelberg, pp. 414–427.

[50] Su, L., and Vaidya, N. H. Byzantine-resilient multiagent optimization. *IEEE Transactions on Automatic Control 66*, 5 (2021), 2227–2233.

[51] Xiao, H., Rasul, K., and Vollgraf, R. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747* (2017).

[52] Xie, C., Koyejo, O., and Gupta, I. Generalized byzantine-tolerant sgd, 2018.

[53] Yang, Z., and Bajwa, W. U. Byrdie: Byzantine-resilient distributed coordinate descent for decentralized learning, 2017.

[54] Yin, D., Chen, Y., Kannan, R., and Bartlett, P. Byzantine-robust distributed learning: Towards optimal statistical rates. In *Proceedings of the 35th International Conference on Machine Learning* (10–15 Jul 2018), J. Dy and A. Krause, Eds., vol. 80 of *Proceedings of Machine Learning Research*, PMLR, pp. 5650–5659.

# A   Appendix: Proof of Lemma 1

**Lemma** (Restated)**.** *If $f \geq n/2$ then there cannot exist a deterministic $(f, \epsilon)$-resilient algorithm for any $\epsilon \geq 0$.*

*Proof.* We prove the lemma by contradiction. We consider a case when $n = 2$, $d = 1$, i.e., $Q_i : \mathbb{R} \to \mathbb{R}$ for all $i$, and all the cost functions have unique minimum points. Suppose that $f = 1$, and that there exists a deterministic $(f, \epsilon)$-resilient algorithm $\Pi$ for some $\epsilon \geq 0$. Without loss of generality, we suppose that agent 2 is Byzantine faulty. We denote $x_1 = \arg\min_{x \in \mathbb{R}} Q_1(x)$.

The Byzantine agent 2 can choose to behave as a non-faulty agent with cost function $\widetilde{Q}_2(x) = Q_1(x - x_1 - 2\epsilon - \delta)$ where $\delta$ is some positive real value. Now, note that the minimum point of $\widetilde{Q}_2(x)$, which we denoted by $x_2$, is unique and equal to $x_1 + 2\epsilon + \delta$. Therefore, $|x_1 - x_2| = 2\epsilon + \delta > 2\epsilon$. As the identity of Byzantine agent is a priori unknown, the server cannot distinguish between scenarios; (i) agent 1 is non-faulty, and (ii) agent 2 is non-faulty. Now, being deterministic algorithm, $\Pi$ should produce the same output in both the scenarios. In scenario (i), as $\Pi$ is assumed $(f, \epsilon)$-resilient, its output must lie in the interval $[x_1 - \epsilon, x_1 + \epsilon]$. Similarly, in scenario (ii), the output of $\Pi$ must lie in the interval $[x_2 - \epsilon, x_2 + \epsilon]$. However, as $|x_1 - x_2| > 2\epsilon$, the two intervals $[x_1 - \epsilon, x_1 + \epsilon]$ and $[x_2 - \epsilon, x_2 + \epsilon]$ do not overlap. Therefore, $\Pi$ cannot be $(f, \epsilon)$-resilient in both the scenarios simultaneously, which is a contradiction to the assumption that $\Pi$ is $(f, \epsilon)$-resilient.

The above argument extends easily for the case when $n > 2$, and $f > n/2$.   $\square$

# B    Appendix: The Special Case of $(f, 0)$-Resilience

We show that $(f, 0)$-*resilience*, stated in Definition 2, and *exact fault-tolerance*, defined in Section 1.1, are equivalent in the deterministic framework. Specifically, we show that a deterministic $(f, 0)$-resilient algorithm achieves exact fault-tolerance, and a deterministic exact fault-tolerant algorithm is $(f, 0)$-resilient. We consider the server-based system architecture. Notation $\arg\min_{x\in\mathbb{R}^d}$ is simply written as $\arg\min$, unless otherwise stated.

First, we show that $(f, 0)$-resilience implies exact fault-tolerance. Suppose that there exists a deterministic $(f, 0)$-resilient algorithm $\Pi$. Consider an arbitrary execution $E_{\mathcal{H}}$ of $\Pi$ wherein $\mathcal{H} \subseteq \{1, \ldots, n\}$ denotes the set of all the non-faulty agents, and let $\widehat{x}$ denote the output. Recall that, as there are at most $f$ faulty agents, $|\mathcal{H}| \geq n - f$. To prove that $\Pi$ has exact fault-tolerance, it suffices to show that, in execution $E_{\mathcal{H}}$, $\widehat{x}$ is a minimum point of the aggregate cost function of all non-faulty agents $\sum_{i \in \mathcal{H}} Q_i(x)$.

By Definition 2 of $(f, 0)$-resilience, for every set $S \subseteq \mathcal{H}$ with $|S| = n - f$,

$$\widehat{x} \in \arg\min \sum_{i \in S} Q_i(x).$$

Therefore, for every set $S$ with $S \subseteq \mathcal{H}$ and $|S| = n - f$,

$$\sum_{i \in S} Q_i(\widehat{x}) \leq \sum_{i \in S} Q_i(x), \quad \forall x \in \mathbb{R}^d. \tag{28}$$

Now, note that there are $\binom{|\mathcal{H}|}{n-f}$ subsets in $\mathcal{H}$ of size $n - f$, and each agent $i \in \mathcal{H}$ is contained in $\binom{|\mathcal{H}|-1}{n-f-1}$ of those subsets. Therefore,

$$\sum_{\substack{S \subseteq \mathcal{H}, \\ |S|=n-f}} \sum_{i \in S} Q_i(x) = \binom{|\mathcal{H}| - 1}{n - f - 1} \sum_{i \in \mathcal{H}} Q_i(x). \tag{29}$$

Substituting from (28) in (29) we obtain that

$$\sum_{i \in \mathcal{H}} Q_i(\widehat{x}) \leq \sum_{i \in \mathcal{H}} Q_i(x), \quad \forall x \in \mathbb{R}^d$$

The above implies that

$$\widehat{x} \in \arg\min \sum_{i \in \mathcal{H}} Q_i(x).$$

The above proves that $\Pi$ has exact fault-tolerance in execution $E_{\mathcal{H}}$.

Now, we show that exact fault-tolerance implies $(f, 0)$-resilience. Suppose that $\Pi$ is a deterministic algorithm with exact fault-tolerance. Similar to above, consider an arbitrary execution $E_{\mathcal{H}}$ of $\Pi$ wherein set $\mathcal{H}$ comprises all the non-faulty agents, and $\widehat{x}$ is its output. Therefore,

$$\widehat{x} \in \arg\min \sum_{i \in \mathcal{H}} Q_i(x).$$

To prove that $\Pi$ is $(f, 0)$-resilient, it suffices to show that in execution $E_{\mathcal{H}}$ for every set $S \subseteq \mathcal{H}$ with $|S| = n - f$, $\widehat{x}$ is a minimum point of the aggregate cost function $\sum_{i \in S} Q_i(x)$. This is trivially true when $|\mathcal{H}| = n - f$. We assume below that $|\mathcal{H}| > n - f$.

Consider an arbitrary subset $S$ of $\mathcal{H}$ with $|S| = n - f$. Consider an execution $E_S$ wherein $S$ is the set of all non-faulty agents, with the remaining agents in $\{1, \ldots, n\} \setminus S$ being Byzantine faulty. Suppose that the inputs from all the agents to the server in $E_S$ are identical to their inputs in $E_{\mathcal{H}}$. Therefore, as $\Pi$ is a deterministic algorithm, its output in execution $E_S$ is same as that in execution $E_{\mathcal{H}}$, i.e., $\widehat{x}$. Moreover, as $\Pi$ is assumed to have exact fault-tolerance,

$$\widehat{x} \in \arg\min \sum_{i \in S} Q_i(x).$$

As $S$ is an arbitrary subset $S$ of $\mathcal{H}$ with $|S| = n - f$, the above proves that $\Pi$ is $(f, 0)$-resilient in execution $E_{\mathcal{H}}$.

# C  Appendix: Proof of $\gamma \leq \mu$

We show below that if Assumptions 2 and 3 hold true simultaneously then $\gamma \leq \mu$.

Consider an arbitrary set of $n - f$ non-faulty agents $\mathcal{H}$, and two arbitrary non-identical points $x, y \in \mathbb{R}^d$, i.e., $x \neq y$. If Assumption 2 holds true then

$$\|\nabla Q_i(x) - \nabla Q_i(y)\| \leq \mu \|x - y\|, \quad \forall i \in \mathcal{H}.$$

Therefore, owing to the Cauchy-Schwartz inequality, for all $i \in \mathcal{H}$,

$$\langle x - y, \nabla Q_i(x) - \nabla Q_i(y) \rangle \leq \|x - y\| \|\nabla Q_i(x) - \nabla Q_i(y)\| \leq \mu \|x - y\|^2. \tag{30}$$

From (30) we obtain that

$$\sum_{i \in \mathcal{H}} \langle x - y, \nabla Q_i(x) - \nabla Q_i(y) \rangle \leq \mu |\mathcal{H}| \|x - y\|^2. \tag{31}$$

If Assumption 3 holds true then

$$\sum_{i \in \mathcal{H}} \langle x - y, \nabla Q_i(x) - \nabla Q_i(y) \rangle \geq \gamma |\mathcal{H}| \|x - y\|^2. \tag{32}$$

As $x, y$ are arbitrary non-identical points, (31) and (32) together imply that $\gamma \leq \mu$.

# D  Appendix: Proof of Theorem 3

The proof of Theorem 3 relies on the following sufficient criterion for the convergence of non-negative sequences.

---

**Lemma 2** (Bottou, 1998 [7]). *Consider a sequence of real values $\{u_t, t = 0, 1, \ldots\}$. If $u_t \geq 0, \forall t$ then*

$$\sum_{t=0}^{\infty} (u_{t+1} - u_t)_+ < \infty \implies \begin{cases} u_t \underset{t \to \infty}{\longrightarrow} u_\infty < \infty \\[2mm] \sum_{t=0}^{\infty} (u_{t+1} - u_t)_- > -\infty \end{cases} \tag{33}$$

*where the operators $(\cdot)_+$ and $(\cdot)_-$ are defined as follows for a real scalar $x$,*

$$(x)_+ = \begin{cases} x &, \quad x > 0 \\ 0 &, \quad otherwise \end{cases}, \text{ and } (x)_- = \begin{cases} 0 &, \quad x > 0 \\ x &, \quad otherwise \end{cases}$$

---

Recall from the statement of Theorem 3 that $x^* \in \mathcal{W}$ where $\mathcal{W}$ is a compact convex set. We define, for all $t \in \{0, 1, \ldots\}$,

$$e_t = \left\| x^t - x^* \right\|. \tag{34}$$

Next, we define a univariate real-valued function $\psi : \mathbb{R} \to \mathbb{R}$:

$$\psi(y) = \begin{cases} 0, & y < (\mathsf{D}^*)^2, \\ \left(y - (\mathsf{D}^*)^2\right)^2, & y \geq (\mathsf{D}^*)^2. \end{cases} \tag{35}$$

Let $\psi'(y)$ denote the derivative of $\psi$ at $y \in \mathbb{R}$. Specifically,

$$\psi'(y) = \max\left\{0, 2\left(y - (\mathsf{D}^*)^2\right)\right\}. \tag{36}$$

We show below that $\psi'(\cdot)$ is a Lipschitz continuous function with Lipschitz coefficient of 2. From (36), we obtain that

$$|\psi'(y) - \psi'(z)| = \begin{cases} 2\,|y - z| & , \quad \text{both } y, z \geq (\mathsf{D}^*)^2 \\ 2\,|y - (D^*)^2| & , \quad y \geq (\mathsf{D}^*)^2, z < (\mathsf{D}^*)^2 \\ 0 & , \quad \text{both } y, z < (\mathsf{D}^*)^2 \end{cases} \tag{37}$$

Note from (37) that for the case when $y \geq (\mathsf{D}^*)^2$, $z < (\mathsf{D}^*)^2$,

$$|\psi'(y) - \psi'(z)| = 2\left|y - (D^*)^2\right| < 2\,|y - z|.$$

Similarly, due to symmetry, when $y < (\mathsf{D}^*)^2$, $z \geq (\mathsf{D}^*)^2$ then $|\psi'(y) - \psi'(z)| = 22\left|z - (D^*)^2\right| < 2\,|y - z|$. Therefore, from (36) we obtain that

$$|\psi'(y) - \psi'(z)| \leq 2\,|y - z|, \quad \forall y, z \in \mathbb{R}. \tag{38}$$

The Lipschitz continuity of $\psi'(\cdot)$, shown in (38), implies that [8, Section 4.1]

$$\psi(z) - \psi(y) \leq (z - y)\psi'(y) + (z - y)^2, \quad \forall y, z \in \mathbb{R}. \tag{39}$$

Now, for each $t \in \{0, 1, \ldots\}$, we define

$$h_t = \psi(e_t^2). \tag{40}$$

From (39) and (40), for all $t$, we obtain that

$$h_{t+1} - h_t = \psi\left(e_{t+1}^2\right) - \psi\left(e_t^2\right) \leq \left(e_{t+1}^2 - e_t^2\right) \cdot \psi'\left(e_t^2\right) + \left(e_{t+1}^2 - e_t^2\right)^2.$$

From now on we use $\psi_t'$ as a shorthand for $\psi'(e_t^2)$. From above we have

$$h_{t+1} - h_t \leq \left(e_{t+1}^2 - e_t^2\right)\psi_t' + \left(e_{t+1}^2 - e_t^2\right)^2. \tag{41}$$

Now, recall from (21) that for all $t \in \{0, 1, \ldots\}$,

$$x^{t+1} = \left[x^t - \eta_t\, \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right]_{\mathcal{W}} \tag{42}$$

By the non-expansion property of Euclidean projection onto a closed convex set,

$$\left\| x^{t+1} - x^* \right\| \leq \left\| x^t - x^* - \eta_t\, \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\|.$$

Recall from (34) that $e_t$ denotes $\|x^t - x^*\|$ for all $t$. Upon squaring the both sides in the above inequality, we obtain that

$$e_{t+1}^2 \leq e_t^2 - 2\eta_t \left\langle x_t - x^*, \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\rangle$$
$$+ \eta_t^2 \left\|\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\|^2 . \tag{43}$$

Recall, from (22) in the statement of Theorem 3, that

$$\phi_t = \left\langle x_t - x^*, \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\rangle, \quad \forall t.$$

Substituting from the above in (43), we obtain that

$$e_{t+1}^2 \leq e_t^2 - 2\eta_t \phi_t + \eta_t^2 \left\|\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\|^2 . \tag{44}$$

As $\psi_t' \geq 0$, $\forall t$, substituting from (44) in (41) we get

$$h_{t+1} - h_t \leq \left(-2\eta_t \phi_t + \eta_t^2 \left\|\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\|^2\right)\psi_t' + \left(e_{t+1}^2 - e_t^2\right)^2 . \tag{45}$$

Note that for an arbitrary $t$,

$$\left|e_{t+1}^2 - e_t^2\right| = (e_{t+1} + e_t)\left|e_{t+1} - e_t\right| . \tag{46}$$

As $\mathcal{W}$ is assumed compact, there exists $\Gamma = \max_{x \in \mathcal{W}} \|x - x^*\| < \infty$. We assume $\Gamma > 0$, otherwise $\mathcal{W} = \{x^*\}$ and the theorem is trivial. Recall from the update rule (21), which is re-stated above in (42), that $x^t \in \mathcal{W}$ for all $t$, and that $x^* \in \mathcal{W}$. Therefore,

$$e_t = \left\|x^t - x^*\right\| \leq \max_{x \in \mathcal{W}} \|x - x^*\| = \Gamma, \quad \forall t. \tag{47}$$

From (47), for all $t$, we obtain that

$$e_{t+1} + e_t \leq 2\Gamma.$$

Substituting from above in (46) implies that

$$\left|e_{t+1}^2 - e_t^2\right| \leq 2\Gamma \left|e_{t+1} - e_t\right|, \quad \forall t. \tag{48}$$

From triangle inequality, we get

$$\left|e_{t+1} - e_t\right| = \left|\left\|x^{t+1} - x^*\right\| - \left\|x^t - x^*\right\|\right| \leq \left\|x^{t+1} - x^t\right\| .$$

Substituting from above in (48) we obtain that

$$\left|e_{t+1}^2 - e_t^2\right| \leq 2\Gamma \left\|x^{t+1} - x^t\right\| . \tag{49}$$

Due to the non-expansion property of Euclidean projection onto a closed convex set, from (42) we obtain that

$$\left\|x^{t+1} - x^t\right\| = \left\|\left[x^t - \eta_t \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right]_{\mathcal{W}} - x^t\right\| \leq \eta_t \left\|\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\| .$$

Substituting from above in (49) we obtain that

$$\left|e_{t+1}^2 - e_t^2\right| \leq 2\eta_t \Gamma \left\|\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\| .$$

Thus,

$$\left(e_{t+1}^2 - e_t^2\right)^2 \leq 4\eta_t^2 \Gamma^2 \left\|\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\|^2 . \tag{50}$$

Substituting from (50) in (45) we obtain that, for all $t$,

$$h_{t+1} - h_t \leq \left(-2\eta_t \phi_t + \eta_t^2 \left\|\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\|^2\right) \psi_t' + 4\eta_t^2 \Gamma^2 \left\|\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\|^2$$
$$= -2\eta_t \phi_t \psi_t' + \eta_t^2 \left(\psi_t' + 4\Gamma^2\right) \left\|\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\|^2. \tag{51}$$

Recall from (47) that $e_t \leq \Gamma$. Also, by assumption, $\mathsf{D}^* < \max_{x \in \mathcal{W}} \|x - x^*\| = \Gamma$. Recall that $\psi_t'$ is short for $\psi'(e_t^2)$. Therefore, from (36) we obtain that

$$0 \leq \psi_t' \leq 2 \left(\Gamma^2 - (\mathsf{D}^*)^2\right) \leq 2\Gamma^2, \quad \forall t. \tag{52}$$

As the statement of Theorem 3 assumes that $\left\|\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\| < \infty$ for all $t$, there exists a real value $\mathsf{M} < \infty$ such that

$$\left\|\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\| \leq \mathsf{M}, \quad \forall t. \tag{53}$$

Substituting from (52) and (53) in (51) we obtain that

$$h_{t+1} - h_t \leq -2\eta_t \phi_t \psi_t' + 6\eta_t^2 \Gamma^2 \mathsf{M}^2. \tag{54}$$

We now use Lemma 2 to prove that $h_\infty = 0$ as follows.

For an iteration $t$, we consider below two possible cases: (i) $e_t < \mathsf{D}^*$, and (ii) $e_t = \mathsf{D}^* + \delta$ for some $\delta \geq 0$.

Case i) In this particular case, $\psi_t' = 0$. Therefore, due to Cauchy-Schwartz inequality,

$$|\phi_t| = \left|\left\langle x^t - x^*, \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\rangle\right| \leq e_t \left\|\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)\right\|.$$

Substituting from (53) above, we obtain that $|\phi_t| \leq \Gamma\mathsf{M} < \infty$. Therefore,

$$\phi_t \psi_t' = 0. \tag{55}$$

Case ii) In this particular case, from (36), we obtain that

$$\psi_t' = 2 \left((\mathsf{D}^* + \delta)^2 - (\mathsf{D}^*)^2\right) = 2\delta(2\mathsf{D}^* + \delta).$$

Now, by assumption, $\phi_t \geq \xi$ when $e_t \geq \mathsf{D}^*$ where $\xi > 0$. Therefore,

$$\phi_t \psi_t' \geq 2\xi\delta(2\mathsf{D}^* + \delta). \tag{56}$$

From (55) and (56) above, we obtain that

$$\phi_t \psi_t' \geq 0, \quad \forall t. \tag{57}$$

Substituting the above in (54) implies that

$$h_{t+1} - h_t \leq 6\eta_t^2 \Gamma^2 \mathsf{M}^2, \quad \forall t.$$

Recall that notation $(\cdot)_+$ from Lemma 2. The above inequality implies that

$$(h_{t+1} - h_t)_+ \leq 6\eta_t^2 \Gamma^2 \mathsf{M}^2.$$

As $\eta_t^2 < \infty$, $\forall t$, and constants $\mathsf{L}, \mathsf{M} < \infty$, the above implies that

$$\sum_{t=0}^{\infty} (h_{t+1} - h_t)_+ \leq 6\Gamma^2 \mathsf{M}^2 \sum_{t=0}^{\infty} \eta_t^2 < \infty.$$

27

As $h_t \geq 0$ for all $t$, the above in conjunction with Lemma 2 implies that

$$h_t \xrightarrow{t \to \infty} h_\infty < \infty, \text{ and}$$

$$\sum_{t=0}^{\infty} (h_{t+1} - h_t)_- > -\infty. \tag{58}$$

Note that $h_\infty - h_0 = \sum_{t=0}^{\infty} (h_{t+1} - h_t)$. Thus, from (54) we obtain that

$$h_\infty - h_0 \leq -2 \sum_{t=0}^{\infty} \eta_t \phi_t \psi_t' + 6\Gamma^2 \mathsf{M}^2 \sum_{t=0}^{\infty} \eta_t^2. \tag{59}$$

By Definition (40), $h_t \geq 0$ for all $t$. Therefore, from (59) above we obtain that

$$2 \left| \sum_{t=0}^{\infty} \eta_t \phi_t \psi_t' \right| \leq h_0 + h_\infty + 6\Gamma^2 \mathsf{M}^2 \sum_{t=0}^{\infty} \eta_t^2. \tag{60}$$

By assumption, $\sum_{t=0}^{\infty} \eta_t^2 < \infty$. From (58), $h_\infty < \infty$. Substituting from (47) that $e_t < \infty$, $\forall t$ in Definition of $h_t$ (40), we obtain that $h_0 = \psi(e_0^2) < \infty$. Therefore, (60) implies that

$$2 \left| \sum_{t=0}^{\infty} \eta_t \phi_t \psi_t' \right| < \infty.$$

Recall from (57) that $\phi_t \psi_t' \geq 0$ for all $t$. Thus, from above we obtain that

$$\sum_{t=0}^{\infty} \eta_t \phi_t \psi_t' < \infty. \tag{61}$$

Finally, we reason below by contradiction that $h_\infty = 0$.

Note that for any $\zeta > 0$, there exists a unique positive value $\beta$ such that $\zeta = 2\beta \left(2\mathsf{D}^* + \sqrt{\beta}\right)^2$. Suppose that $h_\infty = 2\beta(2\mathsf{D}^* + \sqrt{\beta})^2$ for some positive value $\beta$. As the sequence $\{h_t\}_{t=0}^{\infty}$ converges to $h_\infty$ (see (58)), there exists some finite $\tau \in \mathbb{Z}_{\geq 0}$ such that for all $t \geq \tau$,

$$|h_t - h_\infty| \leq \beta \left(2\mathsf{D}^* + \sqrt{\beta}\right)^2 \implies h_t \geq h_\infty - \beta \left(2\mathsf{D}^* + \sqrt{\beta}\right)^2.$$

As $h_\infty = 2\beta(2\mathsf{D}^* + \sqrt{\beta})^2$, the above implies that

$$h_t \geq \beta \left(2\mathsf{D}^* + \sqrt{\beta}\right)^2, \quad \forall t \geq \tau. \tag{62}$$

Therefore (cf. (35) and (40)), for all $t \geq \tau$,

$$\left(e_t^2 - (\mathsf{D}^*)^2\right)^2 \geq \beta \left(2\mathsf{D}^* + \sqrt{\beta}\right)^2, \text{ or}$$

$$\left|e_t^2 - (\mathsf{D}^*)^2\right| \geq \sqrt{\beta} \left(2\mathsf{D}^* + \sqrt{\beta}\right).$$

Thus, for each $t \geq \tau$, either

$$e_t^2 \geq (\mathsf{D}^*)^2 + \sqrt{\beta} \left(2\mathsf{D}^* + \sqrt{\beta}\right) = \left(\mathsf{D}^* + \sqrt{\beta}\right)^2, \tag{63}$$

or

$$e_t^2 \leq (\mathsf{D}^*)^2 - \sqrt{\beta} \left(2\mathsf{D}^* + \sqrt{\beta}\right) < (\mathsf{D}^*)^2. \tag{64}$$

If the latter, i.e. (64), holds true for some $t' \geq \tau$ then

$$h_{t'} = \psi \left(e_{t'}^2\right) = 0,$$

which contradicts (62). Therefore, (62) implies (63) for all $t \geq \tau$.

From above we obtain that if $h_\infty = 2\beta(2\mathsf{D}^* + \sqrt{\beta})^2$ then there exists $\tau < \infty$ such that for all $t \geq \tau$,

$$e_t \geq \mathsf{D}^* + \sqrt{\beta}.$$

Thus, from (56) we obtain that

$$\phi_t \psi'_t \geq 2\xi\sqrt{\beta}(2\mathsf{D}^* + \sqrt{\beta}), \quad \forall t \geq \tau.$$

Therefore,

$$\sum_{t=\tau}^{\infty} \eta_t \phi_t \psi'_t \geq 2\xi\sqrt{\beta}(2\mathsf{D}^* + \sqrt{\beta}) \sum_{t=\tau}^{\infty} \eta_t = \infty.$$

This is a contradiction to (61). Therefore, $h_\infty = 0$, and by definition of $h_t$ in (40),

$$h_\infty = \lim_{t \to \infty} \psi(e_t^2) = 0.$$

Hence, by definition of $\psi(\cdot)$ in (35),

$$\lim_{t \to \infty} \left\| x^t - x^* \right\| \leq \mathsf{D}^*.$$

# E   Appendix: Proof of Theorem 4

Throughout, we assume $f > 0$ to ignore the trivial case of $f = 0$.

Consider an arbitrary set $\mathcal{H}$ of non-faulty agents with $|\mathcal{H}| = n - f$. Recall that under Assumptions 3 and 4, the aggregate cost function $\sum_{i \in \mathcal{H}} Q_i(x)$ has a unique minimum point in set $\mathcal{W}$, which we denote by $x_{\mathcal{H}}$. Specifically,

$$\{x_{\mathcal{H}}\} = \mathcal{W} \cap \arg\min_{x \in \mathbb{R}^d} \sum_{i \in \mathcal{H}} Q_i(x). \tag{65}$$

Recall from (23) that for CGE gradient-filter, in update rule (21),

$$\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right) = \sum_{j=1}^{n-f} g_{i_j}^t, \quad \forall t. \tag{66}$$

First, we show that $\left\| \sum_{j=1}^{n-f} g_{i_j}^t \right\| < \infty$, $\forall t$. Consider a subset $S_1 \subset \mathcal{H}$ with $|S_1| = n - 2f$. From triangle inequality, $\forall x \in \mathbb{R}^d$,

$$\left\| \sum_{j \in S_1} \nabla Q_j(x) - \sum_{j \in S_1} \nabla Q_j(x_{\mathcal{H}}) \right\| \leq \sum_{j \in S_1} \left\| \nabla Q_j(x) - \nabla Q_j(x_{\mathcal{H}}) \right\|.$$

Under Assumption 2, i.e., Lipschitz continuity of non-faulty gradients, for each non-faulty agent $j$, $\left\| \nabla Q_j(x) - \nabla Q_j(x_{\mathcal{H}}) \right\| \leq \mu \left\| x - x_{\mathcal{H}} \right\|$. Substituting this above implies that

$$\left\| \sum_{j \in S_1} \nabla Q_j(x) - \sum_{j \in S_1} \nabla Q_j(x_{\mathcal{H}}) \right\| \leq |S_1| \mu \left\| x - x_{\mathcal{H}} \right\|. \tag{67}$$

As $|S_1| = n - 2f$, the $(2f, \epsilon)$-redundancy property defined in Definition 3 implies that for all $x_1 \in \arg\min_x \sum_{j \in S_1} Q_j(x)$,

$$\|x_1 - x_{\mathcal{H}}\| \leq \epsilon.$$

Substituting from above in (67) implies that, for all $x_1 \in \arg\min_x \sum_{j \in S_1} Q_j(x)$,

$$\left\| \sum_{j \in S_1} \nabla Q_j(x_1) - \sum_{j \in S_1} \nabla Q_j(x_{\mathcal{H}}) \right\| \leq |S_1| \, \mu \, \|x_1 - x_{\mathcal{H}}\| \leq |S_1| \, \mu\epsilon. \tag{68}$$

For all $x_1 \in \arg\min_x \sum_{j \in S_1} Q_j(x)$, $\nabla Q_j(x_1) = 0$. Thus, (68) implies that

$$\left\| \sum_{j \in S_1} \nabla Q_j(x_{\mathcal{H}}) \right\| \leq |S_1| \, \mu\epsilon. \tag{69}$$

Now, consider an arbitrary non-faulty agent $i \in \mathcal{H} \setminus S_1$. Let $S_2 = S_1 \cup \{i\}$. Using similar arguments as above we obtain that under the $(2f, \epsilon)$-redundancy property and Assumption 2, for all $x_2 \in \arg\min_x \sum_{j \in S_2} Q_j(x)$,

$$\left\| \sum_{j \in S_2} \nabla Q_j(x_{\mathcal{H}}) \right\| = \left\| \sum_{j \in S_2} \nabla Q_j(x_2) - \sum_{j \in S_2} \nabla Q_j(x_{\mathcal{H}}) \right\| \leq |S_2| \, \mu\epsilon. \tag{70}$$

Note that $\sum_{j \in S_2} \nabla Q_j(x) = \sum_{j \in S_1} \nabla Q_j(x) + \nabla Q_i(x)$. From triangle inequality,

$$\|\nabla Q_i(x_{\mathcal{H}})\| - \left\| \sum_{j \in S_1} \nabla Q_j(x_{\mathcal{H}}) \right\| \leq \left\| \sum_{j \in S_1} \nabla Q_j(x_{\mathcal{H}}) + \nabla Q_i(x_{\mathcal{H}}) \right\|. \tag{71}$$

Therefore, for each non-faulty agent $i \in \mathcal{H}$,

$$\|\nabla Q_i(x_{\mathcal{H}})\| \leq \left\| \sum_{j \in S_1} \nabla Q_j(x_{\mathcal{H}}) + \nabla Q_i(x_{\mathcal{H}}) \right\| + \left\| \sum_{j \in S_1} \nabla Q_j(x_{\mathcal{H}}) \right\| \leq |S_2| \, \mu\epsilon + |S_1| \, \mu\epsilon$$
$$= (n - 2f + 1)\mu\epsilon + (n - 2f)\mu\epsilon = (2n - 4f + 1)\mu\epsilon. \tag{72}$$

Now, for all $x$ and $i \in \mathcal{H}$, by Assumption 2,

$$\|\nabla Q_i(x) - \nabla Q_i(x_{\mathcal{H}})\| \leq \mu \, \|x - x_{\mathcal{H}}\|.$$

By triangle inequality,

$$\|\nabla Q_i(x)\| \leq \|\nabla Q_i(x_{\mathcal{H}})\| + \mu \, \|x - x_{\mathcal{H}}\|.$$

Substituting from (72) above we obtain that

$$\|\nabla Q_i(x)\| \leq (2n - 4f + 1)\mu\epsilon + \mu \, \|x - x_{\mathcal{H}}\| \leq 2n\mu\epsilon + \mu \, \|x - x_{\mathcal{H}}\|. \tag{73}$$

We use the above inequality (73) to show below that $\left\| \sum_{j=1}^{n-f} g_{i_j}^t \right\|$ is bounded for all $t$. Recall that for each iteration $t$,

$$\left\| g_{i_1}^t \right\| \leq \dots \leq \left\| g_{i_{n-f}}^t \right\| \leq \left\| g_{i_{n-f+1}}^t \right\| \leq \dots \leq \left\| g_{i_n}^t \right\|.$$

As there are at most $f$ Byzantine agents, for each $t$ there exists $\sigma_t \in \mathcal{H}$ such that

$$\left\| g_{i_{n-f}}^t \right\| \leq \left\| g_{i_{\sigma_t}}^t \right\|. \tag{74}$$

As $g_j^t = \nabla Q_j(x^t)$ for all $j \in \mathcal{H}$, from (74) we obtain that

$$\left\| g_{i_j}^t \right\| \leq \left\| \nabla Q_{\sigma_t}(x^t) \right\|, \quad \forall j \in \{1, \dots, n - f\}, \; t.$$

Substituting from (73) above we obtain that for every $j \in \{1, \ldots, n - f\}$,

$$\left\| g_{i_j}^t \right\| \leq \left\| g_{i_{n-f}}^t \right\| \leq 2n\mu\epsilon + \mu \left\| x^t - x_{\mathcal{H}} \right\|.$$

Therefore, from triangle inequality,

$$\left\| \sum_{j=1}^{n-f} g_{i_j}^t \right\| \leq \sum_{j=1}^{n-f} \left\| g_{i_j}^t \right\| \leq (n - f) \left( 2n\mu\epsilon + \mu \left\| x^t - x_{\mathcal{H}} \right\| \right). \tag{75}$$

Recall from (65) that $x_{\mathcal{H}} \in \mathcal{W}$. Let $\Gamma = \max_{x \in \mathcal{W}} \|x - x_{\mathcal{H}}\|$. As $\mathcal{W}$ is a compact set, $\Gamma < \infty$. Recall from the update rule (21) that $x^t \in \mathcal{W}$ for all $t$. Thus, $\|x^t - x_{\mathcal{H}}\| \leq \max_{x \in \mathcal{W}} \|x - x_{\mathcal{H}}\| = \Gamma < \infty$. Substituting this in (75) implies that

$$\left\| \sum_{j=1}^{n-f} g_{i_j}^t \right\| \leq (n - f) \left( 2n\mu\epsilon + \mu\Gamma \right) < \infty. \tag{76}$$

Recall that in this particular case, $\sum_{j=1}^{n-f} g_{i_j}^t = \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)$ (see (66)). Therefore, from above we obtain that

$$\left\| \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right) \right\| < \infty, \quad \forall t. \tag{77}$$

Next, we show that for an arbitrary $\delta > 0$ there exists $\xi > 0$ such that

$$\phi_t \triangleq \left\langle x^t - x_{\mathcal{H}}, \sum_{j=1}^{n-f} g_{i_j}^t \right\rangle \geq \xi \quad \text{when} \quad \left\| x^t - x_{\mathcal{H}} \right\| \geq \mathsf{D}\,\epsilon + \delta.$$

Consider an arbitrary iteration $t$. Note that, as $|\mathcal{H}| = n - f$, there are at least $n - 2f$ agents that are common to both sets $\mathcal{H}$ and $\{i_1, \ldots, i_{n-f}\}$. We let $\mathcal{H}^t = \{i_1, \ldots, i_{n-f}\} \cap \mathcal{H}$. The remaining set of agents $\mathcal{B}^t = \{i_1, \ldots, i_{n-f}\} \setminus \mathcal{H}^t$ comprises of only faulty agents. Note that $|\mathcal{H}^t| \geq n - 2f$ and $|\mathcal{B}^t| \leq f$. Therefore,

$$\phi_t = \left\langle x^t - x_{\mathcal{H}}, \sum_{j \in \mathcal{H}^t} g_j^t \right\rangle + \left\langle x^t - x_{\mathcal{H}}, \sum_{k \in \mathcal{B}^t} g_k^t \right\rangle. \tag{78}$$

Consider the first term in the right-hand side of (78). Note that

$$\left\langle x^t - x_{\mathcal{H}}, \sum_{j \in \mathcal{H}^t} g_j^t \right\rangle = \left\langle x^t - x_{\mathcal{H}}, \sum_{j \in \mathcal{H}^t} g_j^t + \sum_{j \in \mathcal{H} \setminus \mathcal{H}^t} g_j^t - \sum_{j \in \mathcal{H} \setminus \mathcal{H}^t} g_j^t \right\rangle$$

$$= \left\langle x^t - x_{\mathcal{H}}, \sum_{j \in \mathcal{H}} g_j^t \right\rangle - \left\langle x^t - x_{\mathcal{H}}, \sum_{j \in \mathcal{H} \setminus \mathcal{H}^t} g_j^t \right\rangle.$$

Recall that $g_j^t = \nabla Q_j(x^t)$, $\forall j \in \mathcal{H}$. Therefore,

$$\left\langle x^t - x_{\mathcal{H}}, \sum_{j \in \mathcal{H}^t} g_j^t \right\rangle = \left\langle x^t - x_{\mathcal{H}}, \sum_{j \in \mathcal{H}} \nabla Q_j(x^t) \right\rangle - \left\langle x^t - x_{\mathcal{H}}, \sum_{j \in \mathcal{H} \setminus \mathcal{H}^t} \nabla Q_j(x^t) \right\rangle. \tag{79}$$

Due to the strong convexity assumption (i.e., Assumption 3), for all $x, y \in \mathbb{R}^d$,

$$\left\langle x - y, \nabla \sum_{j \in \mathcal{H}} Q_j(x) - \nabla \sum_{j \in \mathcal{H}} Q_j(y) \right\rangle \geq |\mathcal{H}| \, \gamma \, \|x - y\|^2.$$

31

As $x_{\mathcal{H}}$ is minimum point of $\sum_{j \in \mathcal{H}} Q_j(x)$, $\nabla \sum_{j \in \mathcal{H}} Q_j(x_{\mathcal{H}}) = 0$. Thus,

$$\left\langle x^t - x_{\mathcal{H}}, \sum_{j \in \mathcal{H}} \nabla Q_j(x^t) \right\rangle = \left\langle x^t - x_{\mathcal{H}}, \nabla \sum_{j \in \mathcal{H}} Q_j(x^t) - \nabla \sum_{j \in \mathcal{H}} Q_j(x_{\mathcal{H}}) \right\rangle \geq |\mathcal{H}| \, \gamma \left\| x^t - x_{\mathcal{H}} \right\|^2 . \tag{80}$$

Now, due to the Cauchy-Schwartz inequality,

$$\left\langle x^t - x_{\mathcal{H}}, \sum_{j \in \mathcal{H} \backslash \mathcal{H}^t} \nabla Q_j(x^t) \right\rangle = \sum_{j \in \mathcal{H} \backslash \mathcal{H}^t} \left\langle x^t - x_{\mathcal{H}}, \nabla Q_j(x^t) \right\rangle \leq \sum_{j \in \mathcal{H} \backslash \mathcal{H}^t} \left\| x^t - x_{\mathcal{H}} \right\| \, \left\| \nabla Q_j(x^t) \right\| . \tag{81}$$

Substituting from (80) and (81) in (79) we obtain that

$$\left\langle x^t - x_{\mathcal{H}}, \sum_{j \in \mathcal{H}^t} g_j^t \right\rangle \geq \gamma \, |\mathcal{H}| \, \left\| x^t - x_{\mathcal{H}} \right\|^2 - \sum_{j \in \mathcal{H} \backslash \mathcal{H}^t} \left\| x^t - x_{\mathcal{H}} \right\| \, \left\| \nabla Q_j(x^t) \right\| . \tag{82}$$

Next, we consider the second term in the right-hand side of (78). From the Cauchy-Schwartz inequality,

$$\left\langle x^t - x_{\mathcal{H}}, g_k^t \right\rangle \geq - \left\| x^t - x_{\mathcal{H}} \right\| \, \left\| g_k^t \right\| .$$

Substituting from (82) and above in (78) we obtain that

$$\phi_t \geq \gamma \, |\mathcal{H}| \, \left\| x^t - x_{\mathcal{H}} \right\|^2 - \sum_{j \in \mathcal{H} \backslash \mathcal{H}^t} \left\| x^t - x_{\mathcal{H}} \right\| \, \left\| \nabla Q_j(x^t) \right\| - \sum_{k \in \mathcal{B}^t} \left\| x^t - x_{\mathcal{H}} \right\| \, \left\| g_k^t \right\| . \tag{83}$$

Recall that, due to the sorting of the gradients, for an arbitrary $k \in \mathcal{B}^t$ and an arbitrary $j \in \mathcal{H} \backslash \mathcal{H}^t$,

$$\left\| g_k^t \right\| \leq \left\| g_j^t \right\| = \left\| \nabla Q_j(x^t) \right\| . \tag{84}$$

Recall that $\mathcal{B}^t = \{i_1, \ldots, i_{n-f}\} \backslash \mathcal{H}^t$. Thus, $|\mathcal{B}^t| = n - f - |\mathcal{H}^t|$. Also, as $|\mathcal{H}| = n - f$, $|\mathcal{H} \backslash \mathcal{H}^t| = n - f - |\mathcal{H}^t|$. That is, $|\mathcal{B}^t| = |\mathcal{H} \backslash \mathcal{H}^t|$. Therefore, (84) implies that

$$\sum_{k \in \mathcal{B}^t} \left\| g_k^t \right\| \leq \sum_{j \in \mathcal{H} \backslash \mathcal{H}^t} \left\| \nabla Q_j(x^t) \right\| .$$

Substituting from above in (83), we obtain that

$$\phi_t \geq \gamma \, |\mathcal{H}| \, \left\| x^t - x_{\mathcal{H}} \right\|^2 - 2 \sum_{j \in \mathcal{H} \backslash \mathcal{H}^t} \left\| x^t - x_{\mathcal{H}} \right\| \, \left\| \nabla Q_j(x^t) \right\| .$$

Substituting from (73), i.e., $\left\| \nabla Q_i(x) \right\| \leq 2n\mu\epsilon + \mu \left\| x - x_{\mathcal{H}} \right\|$, above we obtain that

$$\phi_t \geq \gamma \, |\mathcal{H}| \, \left\| x^t - x_{\mathcal{H}} \right\|^2 - 2 \left| \mathcal{H} \backslash \mathcal{H}^t \right| \, \left\| x^t - x_{\mathcal{H}} \right\| \, \left( 2n\mu\epsilon + \mu \left\| x^t - x_{\mathcal{H}} \right\| \right)$$
$$\geq \left( \gamma \, |\mathcal{H}| - 2\mu \left| \mathcal{H} \backslash \mathcal{H}^t \right| \right) \left\| x^t - x_{\mathcal{H}} \right\|^2 - 4n\mu\epsilon \left| \mathcal{H} \backslash \mathcal{H}^t \right| \, \left\| x^t - x_{\mathcal{H}} \right\| .$$

As $|\mathcal{H}| = n - f$ and $\left| \mathcal{H} \backslash \mathcal{H}^t \right| \leq f$, the above implies that

$$\phi_t \geq \left( \gamma(n-f) - 2\mu f \right) \left\| x^t - x_{\mathcal{H}} \right\|^2 - 4n\mu\epsilon f \left\| x^t - x_{\mathcal{H}} \right\|$$
$$= \left( \gamma(n-f) - 2\mu f \right) \left\| x^t - x_{\mathcal{H}} \right\| \left( \left\| x^t - x_{\mathcal{H}} \right\| - \frac{4n\mu\epsilon f}{\gamma(n-f) - 2\mu f} \right)$$
$$= n\gamma \left( 1 - \frac{f}{n} \left( 1 + \frac{2\mu}{\gamma} \right) \right) \left\| x^t - x_{\mathcal{H}} \right\| \left( \left\| x^t - x_{\mathcal{H}} \right\| - \frac{4\mu f \, \epsilon}{\gamma \left( 1 - \frac{f}{n} \left( 1 + \frac{2\mu}{\gamma} \right) \right)} \right) . \tag{85}$$

Recall from (26) that

$$\alpha = 1 - \frac{f}{n}\left(1 + \frac{2\mu}{\gamma}\right).$$

Substituting from above in (85) we obtain that

$$\phi_t \geq \alpha\, n\gamma\, \|x^t - x_{\mathcal{H}}\| \left(\|x^t - x_{\mathcal{H}}\| - \left(\frac{4\mu f}{\alpha\gamma}\right)\epsilon\right). \tag{86}$$

As it is assumed that $\alpha > 0$, (86) implies that for an arbitrary $\delta > 0$,

$$\phi_t \geq \alpha\, n\gamma\, \delta\left(\left(\frac{4\mu f}{\alpha\gamma}\right)\epsilon + \delta\right) > 0 \quad\text{when}\quad \|x^t - x_{\mathcal{H}}\| \geq \mathsf{D}\epsilon + \delta.$$

Hence, the proof.

# F    Appendix: Proof of Theorem 5

In this section we present the proof of Theorem 5. Throughout this section we assume $f > 0$ to ignore the trivial case of $f = 0$. The proof closely follows that of Theorem 4, and we may borrow some notation and results directly from Appendix E.

Recall from (24) that for CWTM gradient-filter, for all $l \in \{1, \ldots, d\}$,

$$\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)[l] = \frac{1}{n - 2f} \sum_{j=f+1}^{n-f} g_{i_j[l]}^t[l], \quad \forall t. \tag{87}$$

First, we show that $\sum_{j=f+1}^{n-f} g_{i_j[l]}^t[l]$ is finite for all $l$ and $t$. From (73) in Appendix E, we know that if the $(2f,\ \epsilon)$-redundancy property and Assumption 2 hold true then for each non-faulty agent $i \in \mathcal{H}$,

$$\|\nabla Q_i(x)\| \leq 2n\mu\epsilon + \mu\|x - x_{\mathcal{H}}\|. \tag{88}$$

The above implies that for all $i \in \mathcal{H}$, $l \in \{1, \ldots, d\}$ and $x$,

$$|\nabla Q_i(x)[l]| \leq 2n\mu\epsilon + \mu\|x - x_{\mathcal{H}}\|. \tag{89}$$

Recall that for all $l$ and $t$,

$$g_{i_1[l]}^t[l] \leq \ldots \leq g_{i_{f+1}[l]}^t[l] \leq \ldots \leq g_{i_{n-f}[l]}^t[l] \leq \ldots \leq g_{i_n[l]}^t[l].$$

As there are at most $f$ Byzantine agents and $|\mathcal{H}| = n - f$, for all $l$ and $t$ there exists a pair of non-faulty agents $\sigma_t^1[l]$, $\sigma_t^2[l] \in \mathcal{H}$ such that

$$g_{i_{n-f}[l]}^t[l] \leq g_{i_{\sigma_t^1[l]}}^t[l], \text{ and } g_{i_{f+1}[l]}^t[l] \geq g_{i_{\sigma_t^2[l]}}^t[l]. \tag{90}$$

As $g_j^t = \nabla Q_j(x^t)$ for all $j \in \mathcal{H}$, from (90) we obtain that for all $j \in \{f+1, \ldots, n-f\}$, $l$ and $t$,

$$\left|g_{i_j[l]}^t[l]\right| \leq \max\left\{\left|\nabla Q_{\sigma_t^1[l]}(x^t)[l]\right|, \left|\nabla Q_{\sigma_t^1[l]}(x^t)[l]\right|\right\}.$$

Substituting from (89) above we obtain that for all $j \in \{f+1, \ldots, n-f\}$, $l$ and $t$,

$$\left|g_{i_j[l]}^t[l]\right| \leq 2n\mu\epsilon + \mu\|x^t - x_{\mathcal{H}}\|.$$

Therefore, owing to the triangle inequality,

$$\left|\sum_{j=f+1}^{n-f} g_{i_j[l]}^t[l]\right| \leq \sum_{j=f+1}^{n-2f} \left|g_{i_j[l]}^t[l]\right| \leq (n - f)\left(2n\mu\epsilon + \mu\|x^t - x_{\mathcal{H}}\|\right). \tag{91}$$

33

Let $\Gamma = \max_{x \in \mathcal{W}} \|x - x_{\mathcal{H}}\|$. As $\mathcal{W}$ is a compact set, $\Gamma < \infty$. Recall from the update rule (21) that $x^t \in \mathcal{W}$ for all $t$. Thus, $\|x^t - x_{\mathcal{H}}\| \leq \max_{x \in \mathcal{W}} \|x - x_{\mathcal{H}}\| = \Gamma < \infty$. Substituting this in (91) implies that for all $l \in \{1, \ldots, d\}$,

$$\left| \sum_{j=f+1}^{n-f} g_{i_j[l]}^t[l] \right| \leq (n - 2f)(2n\mu\epsilon + \mu\Gamma) < \infty.$$

Substituting from above in (87) we obtain that

$$\left\| \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right) \right\| < \infty, \quad \forall t. \tag{92}$$

Now, consider an arbitrary iteration $t$ and $l \in \{1, \ldots, d\}$. From prior works on CWTM gradient-filter for the scalar case [48], i.e., when $d = 1$, we know that trimmed mean of the $l$-th elements of the gradients lies in the convex hull of $l$-th elements of the non-faulty agents' gradients in set $\mathcal{H}$. Specifically,

$$\min_{i \in \mathcal{H}} g_i^t[l] \leq \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)[l] \leq \max_{i \in \mathcal{H}} g_i^t[l]. \tag{93}$$

Obviously,

$$\min_{i \in \mathcal{H}} g_i^t[l] \leq \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} g_i^t[l] \leq \max_{i \in \mathcal{H}} g_i^t[l]. \tag{94}$$

Therefore, from (93) and (94) we obtain that

$$\left| \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)[l] - \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} g_i^t[l] \right| \leq \max_{i \in \mathcal{H}} g_i^t[l] - \min_{i \in \mathcal{H}} g_i^t[l].$$

As $\max_{i \in \mathcal{H}} g_i^t[l] - \min_{i \in \mathcal{H}} g_i^t[l] = \max_{i, j \in \mathcal{H}} \left| g_i^t[l] - g_j^t[l] \right|$, and $g_i^t = \nabla Q_i(x^t)$ for all $i \in \mathcal{H}$, the above can be re-written as follows.

$$\left| \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)[l] - \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \nabla Q_i(x^t)[l] \right| \leq \max_{i, j \in \mathcal{H}} \left| \nabla Q_i(x^t)[l] - \nabla Q_j(x^t)[l] \right|. \tag{95}$$

Note that for any two $i, j \in \mathcal{H}$,

$$\left| \nabla Q_i(x^t)[l] - \nabla Q_j(x^t)[l] \right| \leq \left\| \nabla Q_i(x^t) - \nabla Q_j(x^t) \right\|. \tag{96}$$

Substituting from Assumption 5,

$$\left\| \nabla Q_i(x) - \nabla Q_j(x) \right\| \leq \lambda \max \left\{ \left\| \nabla Q_i(x) \right\|, \left\| \nabla Q_j(x) \right\| \right\},$$

in (96) we obtain that

$$\left| \nabla Q_i(x^t)[l] - \nabla Q_j(x^t)[l] \right| \leq \lambda \max \left\{ \left\| \nabla Q_i(x^t) \right\|, \left\| \nabla Q_j(x^t) \right\| \right\}. \tag{97}$$

Substituting from (88) above we obtain that

$$\left| \nabla Q_i(x^t)[l] - \nabla Q_j(x^t)[l] \right| \leq \lambda \left( 2n\mu\epsilon + \mu \left\| x^t - x_{\mathcal{H}} \right\| \right). \tag{98}$$

Finally, substituting from (98) in (95) we obtain that, for all $l$,

$$\left| \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right)[l] - \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \nabla Q_i(x^t)[l] \right| \leq \lambda \left( 2n\mu\epsilon + \mu \left\| x^t - x_{\mathcal{H}} \right\| \right).$$

34

As $\|x\| = \sqrt{\sum_{l=1}^{d} |x[l]|^2}$ for $x \in \mathbb{R}^d$, the above implies that

$$\left\| \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right) - \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \nabla Q_i(x^t) \right\| \leq \sqrt{d}\lambda \left(2n\mu\,\epsilon + \mu\, \|x^t - x_{\mathcal{H}}\|\right). \tag{99}$$

Now, note that

$$\mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right) = \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \nabla Q_i(x^t) + \left( \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right) - \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \nabla Q_i(x^t) \right). \tag{100}$$

Recall from Theorem 3 that $\phi_t$, for each $t$, is defined to be

$$\phi_t = \left\langle x^t - x_{\mathcal{H}}, \, \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right) \right\rangle.$$

Substituting from (100) above we obtain that

$$\phi_t = \left\langle x^t - x_{\mathcal{H}}, \, \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \nabla Q_i(x^t) \right\rangle + \left\langle x^t - x_{\mathcal{H}}, \, \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right) - \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \nabla Q_i(x^t) \right\rangle. \tag{101}$$

Recall from Assumption 3 that $Q_{\mathcal{H}}(x) = (1/|\mathcal{H}|) \sum_{i \in \mathcal{H}} Q_i(x)$. Thus, the first term on the right-hand side of (101),

$$\left\langle x^t - x_{\mathcal{H}}, \, \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \nabla Q_i(x^t) \right\rangle = \left\langle x^t - x_{\mathcal{H}}, \, \nabla Q_{\mathcal{H}}(x^t) \right\rangle.$$

Substituting from the Assumption 3 above, and recalling that $\nabla Q_{\mathcal{H}}(x_{\mathcal{H}}) = 0$, we obtain that

$$\left\langle x^t - x_{\mathcal{H}}, \, \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \nabla Q_i(x^t) \right\rangle \geq \gamma \, \|x^t - x_{\mathcal{H}}\|^2. \tag{102}$$

Next, we consider the second term on the right-hand side of (101). From Cauchy-Schwartz inequality,

$$\left\langle x^t - x_{\mathcal{H}}, \, \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right) - \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \nabla Q_i(x^t) \right\rangle$$

$$\geq - \|x^t - x_{\mathcal{H}}\| \left\| \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right) - \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \nabla Q_i(x^t) \right\|.$$

Substituting from (99) above we obtain that

$$\left\langle x^t - x_{\mathcal{H}}, \, \mathsf{GradFilter}\left(g_1^t, \ldots, g_n^t\right) - \frac{1}{|\mathcal{H}|} \sum_{i \in \mathcal{H}} \nabla Q_i(x^t) \right\rangle$$

$$\geq -\sqrt{d}\lambda \, \|x^t - x_{\mathcal{H}}\| \left(2n\mu\,\epsilon + \mu\, \|x^t - x_{\mathcal{H}}\|\right). \tag{103}$$

Substituting from (102) and (103) in (101) we obtain that

$$\phi_t \geq \gamma \, \|x^t - x_{\mathcal{H}}\|^2 - \sqrt{d}\lambda \, \|x^t - x_{\mathcal{H}}\| \left(2n\mu\,\epsilon + \mu\, \|x^t - x_{\mathcal{H}}\|\right)$$

$$= \left(\gamma - \sqrt{d}\lambda\mu\right) \|x^t - x_{\mathcal{H}}\| \left( \|x^t - x_{\mathcal{H}}\| - \frac{2\sqrt{d}n\mu\lambda}{(\gamma - \sqrt{d}\mu\lambda)} \epsilon \right). \tag{104}$$

Therefore, if, for an arbitrary $\delta > 0$,

$$\left\| x^t - x_{\mathcal{H}} \right\| \geq \frac{2\sqrt{d}n\mu\lambda}{(\gamma - \sqrt{d}\mu\lambda)} \epsilon + \delta$$

then

$$\phi_t \geq \left(\gamma - \sqrt{d}\lambda\mu\right) \left(\frac{2\sqrt{d}n\mu\lambda}{(\gamma - \sqrt{d}\mu\lambda)} \epsilon + \delta\right) \delta = \left(2\sqrt{d}n\mu\lambda\,\epsilon + \left(\gamma - \sqrt{d}\lambda\mu\right)\delta\right)\delta.$$

Hence, the proof.

# G  Details on the Numerical Experiments

In this section, we present the details of the simulation results to empirically compare the approximate fault-tolerance achieved by the aforementioned gradient-filters; CGE and CWTM, extending Section 5. For the simulation, we consider the problem of distributed linear regression, which is a special distributed optimization problem with quadratic cost functions [24].

## G.1  Problem description

We consider a synchronous server-based system, as shown in Figure 1, wherein $n = 6$, $d = 2$, and $f = 1$. Each agent $i \in \{1, \ldots, n\}$ has a data point represented by a triplet $(A_i, B_i, N_i)$ where $A_i$ is a $d$-dimensional row vector, $B_i \in \mathbb{R}$ as the response, and a noise value $N_i \in \mathbb{R}$. Specifically, for all $i \in \{1, \ldots, n\}$,

$$B_i = A_i x^* + N_i \quad \text{where} \quad x^* = \begin{pmatrix} 1 \\ 1 \end{pmatrix}. \tag{105}$$

The collective data is represented by a triplet of matrices $(A, B, N)$ where the $i$-th row of $A$, $B$, and $N$ are equal to $A_i$, $B_i$ and $N_i$, respectively. The specific values are as follows.

$$A = \begin{pmatrix} 1 & 0 \\ 0.8 & 0.5 \\ 0.5 & 0.8 \\ 0 & 1 \\ -0.5 & 0.8 \\ -0.8 & 0.5 \end{pmatrix}, \quad B = \begin{pmatrix} 0.9108 \\ 1.3349 \\ 1.3376 \\ 1.0033 \\ 0.2142 \\ -0.3615 \end{pmatrix}, \quad \text{and } N = \begin{pmatrix} -0.0892 \\ 0.0349 \\ 0.0376 \\ 0.0033 \\ -0.0858 \\ -0.0615 \end{pmatrix}. \tag{106}$$

It should be noted that

$$B = Ax^* + N. \tag{107}$$

We let $A_S$, $B_S$ and $N_S$ represent matrices of dimensions $|S| \times 2$, $|S| \times 1$ and $|S| \times 1$ obtained by stacking the rows $\{A_i,\ i \in S\}$, $\{B_i,\ i \in S\}$ and $\{N_i,\ i \in S\}$, respectively, in the increasing order. From (107), observe that for every non-empty set $S$,

$$B_S = A_S x^* + N_S. \tag{108}$$

Recall from basic linear algebra that if $A_S$ is full-column rank, i.e., rank $(A_S) = d = 2$ then $x^*$ is the unique solution of the set of equations in (108). Note that for every set $S$ with $|S| \geq n - 2f = 6 - 2 = 4$, the matrix $A_S$ is full rank. Specifically,

$$\text{rank}\,(A_S) = d = 2, \quad \forall S \subseteq \{1, \ldots, 6\},\ |S| \geq 4. \tag{109}$$

In this particular distributed optimization problem, each agent $i$ has a quadratic cost function defined to be
$$Q_i(x) = (B_i - A_i x)^2, \quad \forall x \in \mathbb{R}^2.$$

For an arbitrary non-empty set of agents $S$, we define

$$Q_S(x) = \sum_{i \in S} Q_i(x) = \sum_{i \in S} (B_i - A_i x)^2 = \|B_S - A_S x\|^2, \quad \forall x \in \mathbb{R}^2. \tag{110}$$

As matrix $A_S$ is full rank for every $S$ with $|S| \geq 4$,

$$\arg \min_{x \in \mathbb{R}^2} Q_S(x) = \arg \min_{x \in \mathbb{R}^2} \|B_S - A_S x\|^2 = (x \mid A_S x = B_S). \tag{111}$$

Therefore, $Q_S(x)$ has a unique minimum point when $|S| \geq 4$. Henceforth, we write notation $\arg \min_{x \in \mathbb{R}^2}$ simply as $\arg \min$, unless otherwise stated.

## G.2  Simulations

Due to the rank condition (109), the agents' cost functions satisfy the $(2f, \epsilon)$-*redundancy* property, stated in Definition 3, with $\epsilon = 0.0890$. The steps for computing $\epsilon$ are described below.

1. For each set $S \subset \{1, \ldots, 6\}$ with $|S| = n - f = 5$, compute $x_S \in \mathbb{R}^2$ such that $B_S = A_S x_S$. Note that, due to (111), $x_S = \arg \min Q_S(x)$.

2. For each set $S \subset \{1, \ldots, 6\}$ with $|S| = n - f = 5$ do the following:

   (a) For each set $\widehat{S} \subseteq S$ with $\left|\widehat{S}\right| \geq n - 2f = 4$, compute $x_{\widehat{S}}$ such that $B_{\widehat{S}} = A_{\widehat{S}} x_{\widehat{S}}$. Note that, due to (111), $x_{\widehat{S}} = \arg \min Q_{\widehat{S}}(x)$.

   (b) Compute
   $$\epsilon_S = \max_{\widehat{S} \subseteq S, |\widehat{S}| \geq 4} \left\| x_S - x_{\widehat{S}} \right\|.$$

   In this particular case, both the sets of minimum points $\arg \min Q_S(x)$ and $\arg \min Q_{\widehat{S}}(x)$ are singleton with points $x_S$ and $x_{\widehat{S}}$, respectively. Therefore,
   $$\left\| x_S - x_{\widehat{S}} \right\| = \mathrm{dist}\left(\arg \min Q_S(x), \arg \min Q_{\widehat{S}}(x)\right).$$

3. In the final step, we compute
   $$\epsilon = \max_{|S| = n - f} \epsilon_S.$$

For each agent $i$, its cost function $Q_i(x)$ has Lipschitz continuous gradients, i.e., satisfy Assumption 2, with Lipschitz coefficient

$$\mu = \overline{v}_i \tag{112}$$

where $\overline{v}_i$ denotes the largest eigenvalue of $A_i^T A_i$. Also, for every set of agents $S$ with $|S| = n - f = 5$, their average cost function $(1/|S|)Q_S(x)$ is strongly convex, i.e., satisfy Assumption 3, with the strong convexity coefficient

$$\gamma = \frac{1}{|S|} \underline{v}_S \tag{113}$$

where $\underline{v}_S$ is the smallest eigenvalue of $A_S^T A_S$. Derivations of (112) and (113) can found in [24, Section 10].

## G.3    Simulation

In our experiments, we simulate the following fault behaviors for the faulty agent.

- *gradient-reverse*: the faulty agent *reverses* its true gradient. Suppose the correct gradient of a faulty agent $i$ at iteration $t$ is $s_i^t$, the agent $i$ will send the incorrect gradient $g_i^t = -s_i^t$ to the server.

- *random*: the faulty agent sends a randomly chosen vector in $\mathbb{R}^d$. In our experiments, the faulty agent in each iteration chooses i.i.d. Gaussian random vector with mean 0 and a isotropic covariance matrix with standard deviation of 200.

We simulate the distributed gradient-descent algorithm described in Section 4.1 by assuming agent 1 to be Byzantine faulty. It should be noted that the identity of the faulty agent is not used in any way during the simulations. Here, the set of non-faulty agents is $\mathcal{H} = \{2, \ldots, 6\}$ and $|\mathcal{H}| = n - f = 5$. Therefore, in this particular case $\mathcal{H}$ is the only set of $n - f$ non-faulty agents. From (111), we obtain that the minimum point of the aggregate cost function $\sum_{i \in \mathcal{H}} Q_i(x)$, denoted by $x_{\mathcal{H}}$, is equal to the solution of the following set of linear equations:

$$B_{\mathcal{H}} = A_{\mathcal{H}} x_{\mathcal{H}}.$$

Specifically, $x_{\mathcal{H}} = \begin{pmatrix} 1.0780 \\ 0.9825 \end{pmatrix}$. Also, note from our earlier deductions in (112) and (113) that in this particular case, the non-faulty agents' cost functions satisfy Assumptions 2 and 3 with $\mu = 1$ and $\gamma = 0.356$, respectively.

**Parameters:** We use the following parameters for implementing the algorithm. In the update rule (21), we use step-size $\eta_t = 1.5/(t+1)$ for iteration $t = 0, 1, \ldots$. Note that this particular step-size is diminishing and satisfies the conditions: $\sum_{t=0}^{\infty} \eta_t = \infty$ and $\sum_{t=0}^{\infty} \eta_t^2 = 3\pi^2/8 < \infty$ (see [43]). We assume the convex compact $\mathcal{W} \subset \mathbb{R}^d$ to be a 2-dimensional hypercube $[-1000, 1000]^2$. Note that $x_{\mathcal{H}} \in \mathcal{W}$, i.e., Assumption 4 holds true. In all the simulation results presented below, the initial estimate $x^0 = (0, 0)^T$.

In every execution, we observe that the iterative estimates produced by the algorithm practically converge after 400 iterations. Thus, to measure the approximate fault-tolerance achieved by the different gradient-filter, i.e., CGE and CWTM, we define the output of the algorithm to be $x_{\text{out}} = x^{500}$. The outputs for the two gradient-filters, under different faulty behaviors, are shown in Table 1. Note that $\text{dist}\,(x_{\mathcal{H}}, x_{\text{out}}) = \|x_{\mathcal{H}} - x_{\text{out}}\|$. The results for the case when the faulty agent sends *random* faulty gradients are only shown for a randomly chosen execution.

**Conclusion:** As shown in Table 1, in all executions, the distances between $x_{\mathcal{H}}$ the output of the algorithm $x_{\text{out}}$ in case of both CGE and CWTM gradient-filters are smaller than $\epsilon$. For the said executions, we plot in Figure 2 the values of the aggregate cost function $\sum_{i \in \mathcal{H}} Q_i(x^t)$ (referred as *loss*) and the approximation error $\|x^t - x_{\mathcal{H}}\|$ (referred as *distance*) for iteration $t$ ranging from 0 to 500. We also show the plots of the fault-free distributed gradient-descent (DGD) method where the faulty agent is omitted, and the DGD method without any gradient-filter when agent 1 is Byzantine faulty. The details for iteration $t$ ranging from 0 to 80 are also highlighted in Figure 3.

# H    Empirical results on machine learning tasks

In this section, we present some experiment results applying our algorithm with DGD proposed in Section 4.1 on distributed machine learning tasks where in the system there are Byzantine faulty agents. The applicability of our algorithm with to distributed learning problems is discussed in Section 1.3. Furthermore, our algorithm can be adapted to use

distributed stochastic gradient descent (D-SGD) instead of DGD. Suppose each agent $i$ samples $b$ data points in each iteration $t$ from its local data generating distribution $\mathcal{D}_i$, and let those data points be $\boldsymbol{z}_i^t = \{z_{i_1}^t, ..., z_{i_b}^t\}$. Each agent computes its stochastic gradient $g_i^t$ by

$$g_i^t = \frac{1}{b} \sum_{z \in \boldsymbol{z}_i^t} \nabla \ell(x^t; z),$$

and sends it to the server. Thus, $g_i^t$ is expected to be the gradient of cost function $Q_i(x)$ at $x^t$, since

$$\mathbb{E}_{\boldsymbol{z}_i^t} \frac{1}{b} \sum_{z \in \boldsymbol{z}_i^t} \ell(x; z) = \frac{1}{b} \sum_{z \in \boldsymbol{z}_i^t} \mathbb{E}_{z \sim \mathcal{D}_i} \ell(x; z) = \frac{1}{b} \sum_{z \in \boldsymbol{z}_i^t} Q_i(x) = Q_i(x).$$

Therefore, the correctness of our algorithm holds in expectation.

It is worth noting that to distributed machine learning problems, it is difficult to compute the exact approximation parameters $\epsilon$, or to ensure the cost functions hold the convexity and smoothness assumptions (i.e., to compute the value of $\mu$ and $\gamma$), as we did in Section 5 or Appendix G. However, the algorithm and our theoretical analysis are still worth considering, since prior research has pointed out that cost functions of many machine learning problems are strongly-convex in the neighborhood of local minimizers [8]. We conduct these experiments to show empirically the broader applicability of our algorithm, and that the redundancy in cost functions we discussed in this paper indeed exists in real-world scenarios.

In our experiments, we simulate a multi-agent distributed learning system in the server-based architecture using multiple threads. There is one thread representing the server, and others each representing one agent. The inter-thread communication is handled via a message passing interface (MPI). The simulator is built in Python using PyTorch [39] and MPI4py [16]. The simulator is deployed on a Google Cloud Platform cluster with 14 vCPUs and 100 GB memory.

We conduct our experiments on two benchmark datasets, MNIST [7] and Fashion-MNIST [51]. Both datasets are monochrome image-classification tasks, and comprise of 60.000 training and 10,000 testing data points, with each image of size $28 \times 28$. The training and testing data points in each dataset are evenly divided into 10 non-overlapping classes. For each dataset, we train a benchmark neural network Lenet [31] with 431,080 learnable parameters, i.e., $d = 431,080$.

In each of our experiments, we simulate a multi-agent system with $n = 10$ agents. Out of $n$ agents, we randomly select $f = 3$ to be Byzantine faulty. For faulty agents in our experiments, we consider two types of faults, listed as follows:
- **Label-flipping** (LF): we scramble the classification labels of data possessed by designated faulty agents. Specifically, suppose a data point of a faulty agent has a label $y \in \{0, ..., 9\}$, it is changed to $\widetilde{y} = 9 - y$.
- **Gradient-reverse** (GR): the faulty agents reverse their true gradients. Specifically, suppose the correct gradient of a faulty agent $i$ at iteration $t$ is $s_i^t$, the agent $i$ will send the incorrect gradient $g_i^t = -s_i^t$ to the server, the same as it is defined in Section 5 and Appendix G.

We also choose the following parameters: batch size $b = 128$, step-size $\eta = 0.01$. In each experiment, the training dataset of 60,000 data points are randomly and evenly divided into 10 subsets, one for each agent.

For each dataset, we compare the two gradient filters discussed in Section 4.2, CGE and coordinate-wise trimmed mean (CWTM), against two types of faulty agents, LF and GR. To make the results comparable, the random seed is fixed across executions, so that
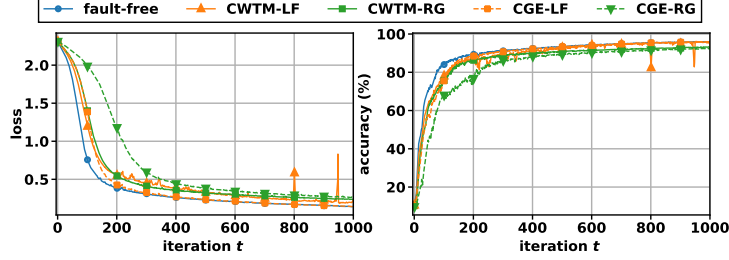
Figure 4: The cross-entropy loss and model accuracy, versus the number of iterations in the algorithm, using our algorithm with D-SGD on MNIST with $n = 10$ and $f = 3$. The two experiments using CWTM are in *solid plots, and the two using CGE are in* dashed plots. The two experiments against LF faults are plotted in *yellow*, while the two against RG are plotted in *green*. We also show the performance of fault-free D-SGD method where the faulty agent is omitted in *blue solid* plots.
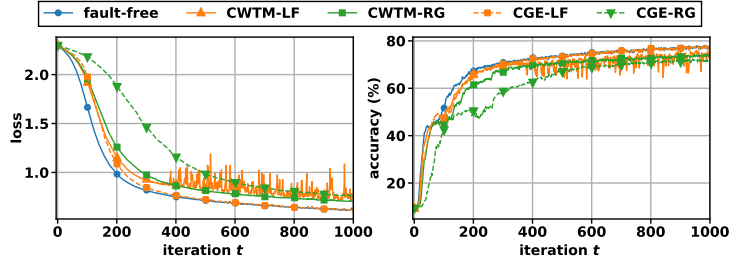


Figure 5: The cross-entropy loss and model accuracy, versus the number of iterations in the algorithm, using our algorithm with D-SGD on Fashion-MNIST with $n = 10$ and $f = 3$. The two experiments using CWTM are in *solid plots, and the two using CGE are in* dashed plots. The two experiments against LF faults are plotted in *yellow*, while the two against RG are plotted in *green*. We also show the performance of fault-free D-SGD method where the faulty agent is omitted in *blue solid* plots.

the data points distributed to agents and the selected faulty agents are fixed. The experiments are also compared by a fault-free setting, where the would-have-been faulty agents do not participate the computation. The performance of our algorithm is measured by cross-entropy loss and model accuracy at each step. The results are shown in Figure 4 for MNIST and Figure 5 for Fashion-MNIST. Since there is a clear trend of converging by the end of 1,000 iterations, we only show performance of the first 1,000 iterations in both figures.

As is shown in the two figures, the losses in all our experiments are converging to within a close range of the fault-free results. Still, the performance of the two gradient filters vary when facing different fault types. The differences between performance in presence of Byzantine faulty agents and fault-free case is more significant comparing to numerical experiments in Section 5, which could be the result of larger value of $\epsilon$ and higher value of $f/n$. Specifically, the performance of CWTM when facing LF faults tabulates during the training process – which is a possible situation according to our theoretical analysis, since the theoretical results only guarantees that the iterative estimates converges a bounded area around the true minimum. The differences in accuracies between gradient filters and fault-free case are less significant, possibly due to the nature of the loss functions in the two machine learning tasks. Overall, these empirical results indicate that our algorithm can be deployed to solve real-world problem, and the redundancy in cost functions indeed exists in some machine learning problems.