
Flag Aggregator: Distributed Training under Failures and Augmented Losses using Convex Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Modern ML applications increasingly rely on complex deep learning models and
2 large datasets. There has been an exponential growth in the amount of computa-
3 tion needed to train the largest models. Therefore, to scale computation and data,
4 these models are inevitably trained in a distributed manner in clusters of nodes,
5 and their updates are aggregated before being applied to the model. However, a
6 distributed setup is prone to Byzantine failures of individual nodes, components,
7 and software. With data augmentation added to these settings, there is a critical
8 need for robust and efficient aggregation systems. We define the quality of workers
9 as reconstruction ratios $\in (0, 1]$, and formulate aggregation as a Maximum Like-
10 likelihood Estimation procedure using Beta densities. We show that the Regularized
11 form of log-likelihood wrt subspace can be approximately solved using iterative
12 least squares solver, and provide convergence guarantees using recent Convex
13 Optimization landscape results. Our empirical findings demonstrate that our ap-
14 proach significantly enhances the robustness of state-of-the-art Byzantine resilient
15 aggregators. We evaluate our method in a distributed setup with a parameter server,
16 and show simultaneous improvements in communication efficiency and accuracy
17 across various tasks.

18 1 Introduction

19 **How to Design Aggregators?** We consider the problem of designing aggregation functions that can
20 be written as optimization problems of the form,

$$\mathcal{A}(g_1, \dots, g_p) \in \arg \min_{Y \in C} A_{g_1, \dots, g_p}(Y), \quad (1)$$

21 where $\{g_i\}_{i=1}^p \subseteq \mathbb{R}^n$ are given estimates of an unknown summary statistic used to compute the
22 Aggregator Y^* . If we choose A to be a quadratic function that decomposes over g_i 's, and $C = \mathbb{R}^n$,
23 then we can see \mathcal{A} is simply the standard mean operator. There is a mature literature of studying such
24 functions for various scientific computing applications [1]. More recently, from the machine learning
25 standpoint there has been a plethora of work [2, 3, 4, 5] on designing provably robust aggregators \mathcal{A}
26 for mean estimation tasks under various technical assumptions on the distribution or moments of g_i .

27 **Distributed ML Use Cases.** Consider training a model with a large dataset such as ImageNet-1K
28 [6] or its augmented version which would require data to be distributed over p workers and uses
29 back propagation. Indeed, in this case, g_i 's are typically the gradients computed by individual
30 workers at each iteration. In settings where the training objective is convex, the convergence and
31 generalization properties of distributed optimization can be achieved by defining \mathcal{A} as a weighted
32 combination of gradients facilitated by a simple consensus matrix, even if some g_i 's are noisy [7, 8].
33 In a distributed setup, as long as the model is convex we can simultaneously minimize the total
34 iteration or communication complexity to a significant extent i.e., it is possible to achieve convergence

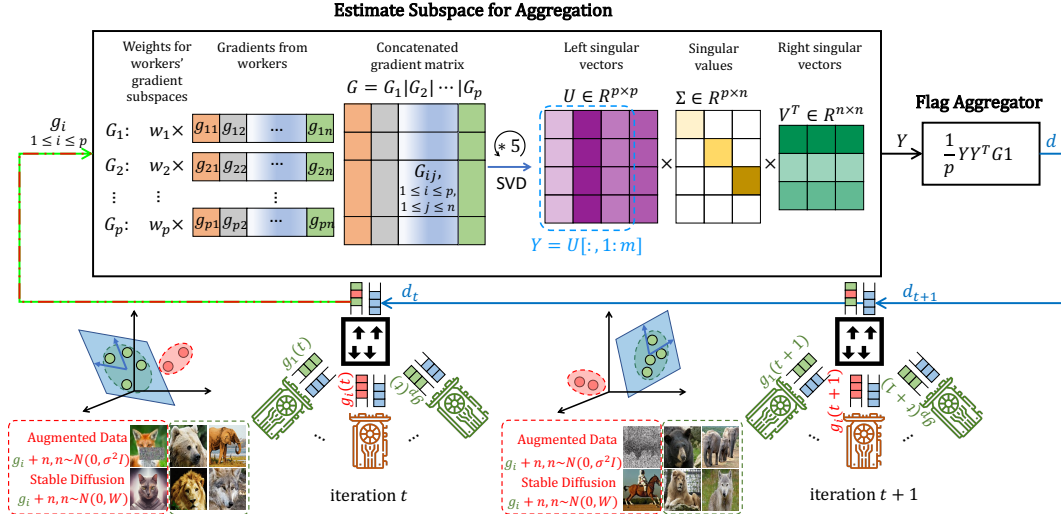


Figure 1: Robust gradient aggregation in our distributed training framework. In our applications, each of the p workers provides gradients computed using a random sample obtained from given training data, derived synthetic data from off-the-shelf Diffusion models, and random noise in each iteration. Our Flag Aggregator (FA) removes high frequency noise components by using few rounds of Singular Value Decomposition of the concatenated Gradient Matrix G , and provides new update Y^* .

35 *and* robustness under technical assumptions on the moments of (unknown) distribution from which
 36 g_i 's are drawn. However, it is still an open problem to determine the optimality of these procedures
 37 in terms of either convergence or robustness [9, 10].

38 **Potential Causes of Noise.** When data is distributed among workers, hardware and software failures
 39 in workers [11, 12, 13] can cause them to send incorrect gradients, which can significantly mislead
 40 the model [14]. To see this, let's consider a simple experiment with 15 workers, that f of them
 41 produce uniformly random gradients. Figure 2 shows that the model accuracy is heavily impacted
 42 when $f > 0$ when mean is used to aggregate the gradients.

43 The failures can occur due to component or software failures and
 44 their probability increases with the scale of the system [15, 16, 17].
 45 Reliability theory is used to analyze such failures, see Chapter 9
 46 in [18], but for large-scale training, the distribution of total system
 47 failures is not independent over workers, making the total noise in
 48 gradients dependent and a key challenge for large-scale training.
 49 Moreover, even if there are no issues with the infrastructure, our
 50 work is motivated by the prevalence of data augmentation, including
 51 hand-chosen augmentations. Since number of parameters n is often
 52 greater than number of samples, data augmentation improves the
 53 generalization capabilities of large-scale models under technical con-
 54 ditions [19, 20, 21]. In particular, Adversarial training is a common
 55 technique that finds samples that are close to training samples but
 56 classified as a different class at the current set of parameters, and
 57 then use such samples for parameter update purposes [22]. Unfortunately, computing adversarial
 58 samples is often difficult [23], done using randomized algorithms [24] and so may introduce depen-
 59 dent (across samples) noise themselves. In other words, using adversarial training paradigm, or the
 60 so-called inner optimization can lead to noise in gradients, which can cause or simulate dependent
 61 "Byzantine" failures in the distributed context.

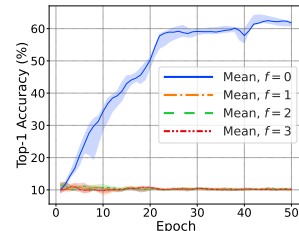


Figure 2: Tolerance to f Byzantine workers for a non-robust aggregator (mean).

62 **Available Computational Solutions.** Most existing open source implementations of \mathcal{A} rely just
 63 on (functions of) pairwise distances to filter gradients from workers using suitable neighborhood
 64 based thresholding schemes, based on moment conditions [25, 26, 27]. While these may be a good
 65 strategy when the noise in samples/gradients is somewhat independent, these methods are suboptimal
 66 when the noise is dependent or nonlinear, especially when n is large. Moreover, choosing discrete

67 hyperparameters such as number of neighbors is impractical in our use cases since they hamper
68 convergence of the overall training procedure. To mitigate the suboptimality of existing aggregation
69 schemes, we explicitly estimate a subspace Y spanned by “most” of the gradient workers, and then
70 use this subspace to estimate that a **sparse** linear combination of g_i gradients, achieving robustness.

71 We present a new optimization based formulation for generalized gradient aggregation purposes in
72 the context of distributed training of deep learning architectures, as shown in Figure 1.

73 **Summary of our Contributions.** From the theoretical perspective, we present a simple Maximum
74 Likelihood Based estimation procedure for aggregation purposes, with novel regularization functions.
75 Algorithmically, we argue that any procedure used to solve Flag Optimization can be directly used to
76 obtain the optimal summary statistic Y^* for our aggregation purposes. **Experimentally**, our results
77 show resilience against Byzantine attacks, encompassing physical failures, while effectively managing
78 the stochasticity arising from data augmentation schemes. In practice, we achieve a *significantly*
79 ($\approx 20\%$) better accuracy on standard datasets. Our **implementation** offers substantial advantages in
80 reducing communication complexity across diverse noise settings through the utilization of our novel
81 aggregation function, making it applicable in numerous scenarios.

82 2 Robust Aggregators as Orthogonality Constrained Optimization

83 In this section, we first provide the basic intuition of our proposed approach to using subspaces for
84 aggregation purposes using linear algebra, along with connections of our approach standard eigende-
85 composition based denoising approaches. We then present our overall optimization formulation in
86 two steps, and argue that it can be optimized using existing methods.

87 2.1 Optimal Subspace Hypothesis for Distributed Descent

88 We will use lowercase letters y, g to denote vectors, and uppercase letters Y, G to denote ma-
89 trices. We will use **boldfont 1** to denote the vector of all ones in appropriate dimensions.
90 Let $g_i \in \mathbb{R}^n$ is the gradient vector from worker i , and $Y \in \mathbb{R}^{n \times m}$
91 is an orthogonal matrix representation of a subspace that gradients
92 could live in such that $m \leq p$. Now, we may interpret each column
93 of Y as a basis function that act on $g_i \in \mathbb{R}^n$, i.e., j -th coordinate of
94 $(Y^T g)_j$ for $1 \leq j \leq m$ is the application of j -th basis or column
95 of Y on g . Recall that by definition of dot product, we have that
96 if $Y_{:,j} \perp x$, then $(Y^T g)_j$ will be close to zero. Equivalently, if
97 $g \in \text{span}(Y)$, then $(Y^T g)^T Y^T g$ will be bounded away from zero,
98 see Chapter 2 in [28]. Assuming that $G \in \mathbb{R}^{n \times p}$ is the gradient
99 matrix of p workers, $Y Y^T G \in \mathbb{R}^{n \times p}$ is the reconstruction of G
100 using Y as basis. That is, i^{th} column of $Y^T G$ specifies the amount
101 of gradient from worker i as a function of Y , and high l_2 norm of
102 $Y^T g_i$ implies that there is a basis in Y such that $Y \not\perp g_i$. So it is
103 easy to see that the average over columns of $Y Y^T G$ would give the final gradient for update.

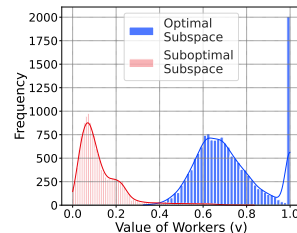


Figure 3: Distributions of Explained Variances on Minibatches

104 **Explained Variance of worker i .** If we denote $z_i = Y^T g_i \in \mathbb{R}^m$ representing the transformation
105 of gradient g_i to z_i using Y , then, $0 \leq \|z_i\|_2^2 = z_i^T z_i = (Y^T g)^T Y^T g = g_i^T Y Y^T g_i$ is a scalar,
106 and so is equal to its trace $\text{tr}(g_i^T Y Y^T g_i)$. Moreover, when Y is orthogonal, we have $0 \leq \|z_i\|_2 =$
107 $\|Y^T g_i\|_2 \leq \|Y\|_2 \|g_i\|_2 \leq \|g_i\|_2$ since the operator norm (or largest singular value) $\|Y\|_2$ of Y is at
108 most 1. Our main idea is to use $\|z_i\|_2^2, \|g_i\|_2^2$ to define the quality of the subspace Y for aggregation,
109 as is done in some previous works for Robust Principal Component Estimation [29] – the quantity
110 $\|z_i\|_2^2 / \|g_i\|_2^2$ is called as *Explained/Expressed* variance of subspace Y wrt i -th worker [30, 31] – we
111 refer to $\|z_i\|_2^2 / \|g_i\|_2^2$ as the “value” of i -th worker. In Figure 3, we can see from the spike near 1.0
112 that if we choose the subspace carefully (blue) as opposed to merely choosing the mean gradient
113 (with unit norm) of all workers, then we can increase the value of workers.

114 **Advantages of Subspace based Aggregation.** We can see that using subspace Y , we can easily: 1.
115 handle different number of gradients from each worker, 2. compute gradient reconstruction $Y Y^T G$
116 efficiently whenever Y is constrained to be orthogonal $Y = \sum_i y_i y_i^T$ where y_i is the i -th column
117 of Y , otherwise have to use eigendecomposition of Y to measure explained variance which can
118 be time consuming. In (practical) distributed settings, the quality (or noise level) of gradients in

119 each worker may be different, **and/or** each worker may use a different batch size. In such cases,
 120 handcrafted aggregation schemes may be difficult to maintain, and fine-tune. For these purposes with
 121 an Orthogonal Subspace Y , we can simply reweigh gradients of worker i according to its noise level,
 122 **and/or** use $g_i \in \mathbb{R}^{n \times b_i}$ where b_i is the batch size of i -th worker with $\text{tr}(z_i^T z_i)$ instead.

123 **Why is optimizing over subspaces called “Flag” Optimization?** Recent optimization results
 124 suggest that we can exploit the finer structure available in Flag Manifold to specify Y more precisely
 125 [32]. For example, $Y \in \mathbb{R}^{m \times n}$ can be parametrized directly as a subspace of dimension m or
 126 as a nested sequence of $Y_k \in \mathbb{R}^{m_k \times n}$, $k = 1, \dots, K$ where $m_k < m_{k+1} \leq p \leq n$ such that
 127 $\text{span}(Y_k) \subseteq \text{span}(Y_{k+1})$ with $Y_K \in \mathbb{R}^{m \times n}$. When $m_{k+1} = m_k = 1$, we have the usual (real)
 128 Grassmanian Manifold (quotient of orthogonal group) whose coordinates can be used for optimization,
 129 please see Section 5 in [33] for details. In fact, [34] used this idea to extend median in one-dimensional
 130 vector spaces to different finite dimensional *subspaces* using the so-called chordal distance between
 131 them. In our distributed training context, we use the explained variance of each worker instead. Here,
 132 workers may specify dimensions along which gradient information is relevant for faster convergence
 133 – an advantage currently not available in existing aggregation implementations – which may be used
 134 for smart initialization also. We use “Flag” to emphasize this additional nested structure available in
 135 our formulation for distributed training purposes.

136 2.2 Approximate Maximum Likelihood Estimation of Optimal Subspace

137 Now that we can evaluate a subspace Y on individual gradients g_i , we now show that finding subspace
 138 Y can be formulated using standard maximum likelihood estimation principles [35]. Our formulation
 139 reveals that regularization is critical for aggregation especially in distributed training. In order to
 140 write down the objective function for finding optimal Y , we proceed in the following two steps:

141 **Step 1.** Assume that each worker provides a single gradient for simplicity. Now, denoting the value of
 142 information v of worker i by $v_i = \frac{z_i^T z_i}{g_i^T g_i}$, we have $v_i \in [0, 1]$. Now by assuming that v_i ’s are observed
 143 from Beta distribution with $\alpha = 1$ and $\beta = \frac{1}{2}$ (for simplicity), we can see that the likelihood $\mathbb{P}(v_i)$ is,

$$\mathbb{P}(v_i) := \frac{(1 - v_i)^{-\frac{1}{2}}}{B(1, \frac{1}{2})} = \frac{\left(1 - \frac{z_i^T z_i}{g_i^T g_i}\right)^{-\frac{1}{2}}}{B(1, \frac{1}{2})}, \quad (2)$$

144 where $B(a, b)$ is the normalization constant. Then, the total log-likelihood of observing gradients g_i
 145 as a function of Y (or v_i ’s) is given by taking the log of product of $\mathbb{P}(v_i)$ ’s as (ignoring constants),

$$\log \left(\prod_{i=1}^p \mathbb{P}(v_i) \right) = \sum_{i=1}^p \log(\mathbb{P}(v_i)) = -\frac{1}{2} \sum_{i=1}^p \log(1 - v_i). \quad (3)$$

146 **Step 2.** Now we use Taylor’s series with constant $a > 0$ to approximate individual worker log-
 147 likelihoods $\log(1 - v_i) \approx a(1 - v_i)^{\frac{1}{a}} - a$ as follows: first, we know that $\exp\left(\frac{\log(v_i)}{a}\right) = v_i^{\frac{1}{a}}$. On
 148 the other hand, using Taylor expansion of \exp about the origin (so large $a > 1$ is better), we have that
 149 $\exp\left(\frac{\log(v_i)}{a}\right) \approx 1 + \frac{\log(v_i)}{a}$. Whence, we have that $1 + \frac{\log(v_i)}{a} \approx v_i^{\frac{1}{a}}$ which immediately implies
 150 that $\log(v_i) \approx a v_i^{\frac{1}{a}} - a$. So, by substituting the Taylor series approximation of \log in Equation 3, we
 151 obtain the *negative* log-likelihood approximation to be *minimized* for robust aggregation purposes as,

$$-\log \left(\prod_{i=1}^p \mathbb{P}(v_i) \right) \approx \frac{1}{2} \sum_{i=1}^p \left(a(1 - v_i)^{\frac{1}{a}} - a \right), \quad (4)$$

152 where $a > 1$ is a sufficiently large constant. In the above mentioned steps, the first step is standard.
 153 Our key insight is using Taylor expansion in (4) with a sufficiently large a to eliminate log optimization
 154 which are known to be computationally expensive to solve, and instead solve *smooth* ℓ_a , $a > 1$ norm
 155 based optimization problems which can be done efficiently by modifying existing procedures [36].

156 **Extension to general beta distributions, and gradients** $\alpha > 0, \beta > 0, g_i \in \mathbb{R}^{n \times k}$. Note that our
 157 derivation in the above two steps can be extended to any beta shape parameters $\alpha > 0, \beta > 0$ – there
 158 will be two terms in the final negative log-likelihood expression in our formulation (4), one for each
 159 α, β . Similarly, by simply using $v_i = \text{tr}(g_i^T Y Y^T g_i)$ to define value of worker i in equation (2), and
 160 then in our estimator in (4), we can easily handle multiple k gradients from a single worker i for Y .

Algorithm 1 Distributed SGD with proposed Flag Aggregator (FA) at the Parameter Server

Input: Number of workers p , loss functions l_1, l_2, \dots, l_p , per-worker minibatch size B , learning rate schedule α_t , initial parameters w_0 , number of iterations T

Output: Updated parameters w_T from any worker

```
1 for  $t = 1$  to  $T$  do
2   for  $p = 1$  to  $p$  in parallel on machine  $p$  do
3     Select a minibatch:  $i_{p,1,t}, i_{p,2,t}, \dots, i_{p,B,t}$   $g_{p,t} \leftarrow \frac{1}{B} \sum_{b=1}^B \nabla l_{i_{p,b,t}}(w_{t-1})$ 
4      $G_t \leftarrow \{g_{1,t}, \dots, g_{p,t}\}$  // Parameter Server receives gradients from  $p$  workers
5      $\hat{Y}_t \leftarrow$  IRLS( $\hat{G}_t$ ) with  $\hat{G}_t = G_t + \lambda \nabla \mathcal{R}(Y) \mathbf{1}^T$  // Do IRLS at the Parameter Server for  $\hat{Y}$ 
6     Obtain gradient direction  $d_t$ :  $d_t = \frac{1}{p} \hat{Y}_t \hat{Y}_t^T G_t \mathbf{1}$  // Compute, Send  $d_t$  to all  $p$  machines
7     for  $p = 1$  to  $p$  in parallel on machine  $p$  do
8       update model:  $w_t \leftarrow w_{t-1} - \alpha_t \cdot d_t$ 
9 Return  $w_T$ 
```

161 2.3 Flag Aggregator for Distributed Optimization

162 It is now easy to see that by choosing $a = 2$, in equation (4), we obtain the negative loglikelihood
163 (ignoring constants) as $(\sum_{i=1}^p \sqrt{1 - g_i^T Y Y^T g_i})$ showing that Flag Median can indeed be seen as
164 an Maximum Likelihood Estimator (MLE). In particular, Flag Median can be seen as an MLE of
165 Beta Distribution with parameters $\alpha = 1$ and $\beta = \frac{1}{2}$. Recent results suggest that in many cases, MLE
166 is ill-posed, and regularization is necessary, even when the likelihood distribution is Gaussian [37].
167 So, based on the Flag Median estimator for subspaces, we propose an optimization based subspace
168 estimator Y^* for aggregation purposes. We formulate our Flag Aggregator (FA) objective function
169 with respect to Y as a *regularized* sum of likelihood based (or data) terms in (4) using trace operators
170 $\text{tr}(\cdot)$ as the solution to the following constrained optimization problem:

$$\min_{Y: Y^T Y = I} A(Y) := \sum_{i=1}^p \sqrt{\left(1 - \frac{\text{tr}(Y^T g_i g_i^T Y)}{\|g_i\|_2^2}\right)} + \lambda \mathcal{R}(Y) \quad (5)$$

171 where $\lambda > 0$ is a regularization hyperparameter. In our analysis, and implementation, we provide
172 support for two possible choices for $\mathcal{R}(Y)$:

173 (1) **Mathematical norms:** $\mathcal{R}(Y)$ can be a form of norm-based regularization other than $\|Y\|_{\text{Fro}}^2$ since
174 it is constant over the feasible set in (5). For example, it could be convex norm with efficient
175 subgradient oracle such as, i.e. element-wise: $\sum_{i=1}^n \sum_{j=1}^m \|Y_{ij}\|_1$ or $\sum_{i=1}^m \|Y_{:,i}\|_1$,

176 (2) **Data-dependent norms:** Following our subspace construction in Section 2.1, we may choose

177 $\mathcal{R}(Y) = \frac{1}{p-1} \sum_{i,j=1, i \neq j}^p \sqrt{\left(1 - \frac{\text{tr}(Y^T (g_i - g_j)(g_i - g_j)^T Y)}{D_{ij}^2}\right)}$ where $D_{ij}^2 = \|g_i - g_j\|_2^2$ denotes the

178 distance between gradient vectors g_i, g_j from workers i, j . Intuitively, the pairwise terms in our
179 loss function (5) favors subspace Y that also reconstructs the pairwise vectors $g_i - g_j$ that are close
180 to each other. So, by setting $\lambda = \Theta(p)$, that is, the pairwise terms dominate the objective function
181 in (5). Hence, λ regularizes optimal solutions Y^* of (5) to contain g_i 's with low pairwise distance
182 in its span – similar in spirit to AggreaThor in [38].

183 **Convergence of Flag Aggregator (FA) Algorithm 1.** With these, we can state our main algorithmic
184 result showing that our FA (5) can be solved efficiently using standard convex optimization proof
185 techniques. In particular, in supplement, we present a smooth Semi-Definite Programming (SDP)
186 relaxation of FA in equation (5) using the Flag structure. This allows us to view the IRLS procedure
187 in 1 as solving the low rank parametrization of the smooth SDP relaxation, thus guaranteeing fast
188 convergence to second order optimal (local) solutions. Importantly, our SDP based proof works for
189 any degree of approximation of the constant a in equation (4) and only relies on smoothness of the
190 loss function wrt Y , although speed of convergence is reduced for higher values of $a \neq 2$, see [39].
191 We leave determining the exact dependence of a on rate of convergence for future work.

192 **How is FA aggregator different from (Bulyan and Multi-Krum)?** Bulyan is a strong Byzantine
193 resilient gradient aggregation rule for $p \geq 4f + 3$ where p is the total number of workers and f is

194 the number of Byzantine workers. Bulyan is a two-stage algorithm. In the first stage, a gradient
 195 aggregation rule R like coordinate-wise median [40] or Krum [9] is recursively used to select
 196 $\theta = p - 2f$ gradients. The process uses R to select gradient vector g_i which is closest to R 's output
 197 (e.g. for Krum, this would be the gradient with the top score, and hence the exact output of R). The
 198 chosen gradient is removed from the received set and added to the selection set S repeatedly until
 199 $|S| = \theta$. The second stage produces the resulting gradient. If $\beta = \theta - 2f$, each coordinate would
 200 be the average of β -nearest to the median coordinate of the θ gradients in S . In matrix terms, if we
 201 consider $S \in \mathbb{R}^{p \times m}$ as a matrix with each column having one non-zero entry summing to 1, Bulyan
 202 would return $\frac{1}{m} \text{ReLU}(GS) \mathbf{1}_m$, where $\mathbf{1}_m \in \mathbb{R}^m$ is the vector of all ones, while FA would return
 203 $\frac{1}{p} YY^T G \mathbf{1}_p$. Importantly, the gradient matrix is being right-multiplied in Bulyan, but left-multiplied
 204 in FA, before getting averaged. While this may seem like a discrepancy, in supplement we show that
 205 by observing the optimality conditions of (5) wrt Y , we show that $\frac{1}{m} YY^T G$ can be seen as a right
 206 multiplication by a matrix parametrized by lagrangian multipliers associated with the orthogonality
 207 constraints in (5). This means it should be possible to combine both approaches for faster aggregation.

208 3 Experiments

209 In this section, we conduct experiments to test our proposed FA in the context of distributed training
 210 in two testbeds. First, to test the performance of our FA scheme solved using IRLS (Flag Mean) on
 211 standard Byzantine benchmarks. Then, to evaluate the ability of existing state-of-the-art gradient
 212 aggregators we augment data via two techniques that can be implemented with Sci-kit package.

213 **Implementation Details.** We implement FA in Pytorch [41], which is popular but does not support
 214 Byzantine resilience natively. We adopt the parameter server architecture and employ Pytorch's
 215 distributed RPC framework with TensorPipe backend for machine-to-machine communication. We
 216 extend Garfield's Pytorch library [42] with FA and limit our IRLS convergence criteria to a small
 217 error, 10^{-10} , or 5 iterations of flag mean for SVD calculation. We set $m = \lceil \frac{p+1}{2} \rceil$.

218 3.1 Setup

219 **Baselines:** We compare FA to several existing aggregation rules: (1) coordinate-wise **Trimmed**
 220 **Mean** [40] (2) coordinate-wise **Median** [40] (3) mean-around-median (**MeaMed**) [43] (4) **Phocas**
 221 [44] (5) **Multi-Krum** [9] (6) **Bulyan** [45].

222 **Accuracy:** The fraction of correct predictions among all predictions, using the test dataset (top-1
 223 cross-accuracy).

224 **Testbed:** We used 4 servers as our experimental platform. Each server has 2 Intel(R) Xeon(R) Gold
 225 6240 18-core CPU @ 2.60GHz with Hyper-Threading and 384GB of RAM. Servers have a Tesla
 226 V100 PCIe 32GB GPU and employ a Mellanox ConnectX-5 100Gbps NIC to connect to a switch.
 227 We use one of the servers as the parameter server and instantiate 15 workers on other servers, each
 228 hosting 5 worker nodes, unless specified differently in specific experiments. For the experiments
 229 designed to show scalability, we instantiate 60 workers.

230 **Dataset and model:** We focus on the image classification task since it is a widely used task for
 231 benchmarking in distributed training [46]. We train ResNet-18 [47] on CIFAR-10 [48] which has
 232 60,000 32×32 color images in 10 classes. For the scalability experiment, we train a CNN with two
 233 convolutional layers followed by two fully connected layers on MNIST [49] which has 70,000 $28 \times$
 234 28 grayscale images in 10 classes. We also run another set of experiments on Tiny ImageNet [50] in
 235 the supplement. We use SGD as the optimizer, and cross-entropy to measure loss. The batch size
 236 for each worker is 128 unless otherwise stated. Also, we use a learning decay strategy where we
 237 decrease the learning rate by a factor of 0.2 every 10 epochs.

238 **Threat models:** We evaluate FA under two classes of Byzantine workers. They can send uniformly
 239 random gradients that are representative of errors in the physical setting, or use non-linear augmented
 240 data described as below.

241 **Evaluating resilience against nonlinear data augmentation:** In order to induce Byzantine behavior
 242 in our workers we utilize ODE solvers to approximately solve 2 non-linear processes, Lotka Volterra

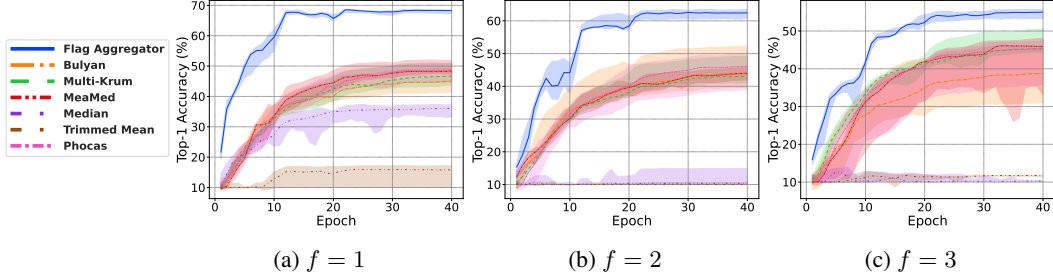


Figure 4: Tolerance to the number of Byzantine workers for robust aggregators for batch size 128.

243 [51] and Arnold’s Cat Map [52], as augmentation methods. Since the augmented samples are
 244 deterministic, albeit nonlinear functions of training samples, the “noise” is dependent across samples.
 245 In **Lotka Volterra**, we use the following linear gradient transformation of 2D pixels:

$$(x, y) \rightarrow (\alpha x - \beta xy, \delta xy - \gamma y),$$

246 where α, β, γ and δ are hyperparameters. We choose them to be $\frac{2}{3}, \frac{4}{3}, -1$ and -1 respectively.

247 Second, we use a *nonsmooth* transformation called **Arnold’s Cat Map** as a data augmentation scheme.
 248 Once again, the map can be specified using a two-dimensional matrix as,

$$(x, y) \rightarrow \left(\frac{2x + y}{N}, \frac{x + y}{N} \right) \bmod 1,$$

249 where mod represents the modulus operation, x and y are the coordinates or pixels of images and N
 250 is the height/width of images (assumed to be square). We also used a smooth approximation of the
 251 Cat Map obtained by approximating the mod function as,

$$(x, y) \rightarrow \frac{1}{n} \left(\frac{2x + y}{(1 + \exp(-m \log(\alpha_1)))}, \frac{x + y}{(1 + \exp(-m \log(\alpha_2)))} \right),$$

252 where $\alpha_1 = \frac{2x+y}{n}$, $\alpha_2 = \frac{x+y}{n}$, and m is the degree of approximation, which we choose to be 0.95 in
 253 our data augmentation experiments.

254 **How to perform nonlinear data augmentation?** In all three cases, we used SciPy’s [53] `solve_ivp`
 255 method to solve the differential equations, by using the LSODA solver. In addition to the setup
 256 described above, we also added a varying level of Gaussian noise to each of the training images. All
 257 the images in the training set are randomly chosen to be augmented with varying noise levels of the
 258 above mentioned augmentation schemes. We have provided the code that implements all our data
 259 augmentation schemes in the supplement zipped folder.

260 3.2 Results

261 **Tolerance to the number of Byzantine workers:** In this experiment, we show the effect of Byzantine
 262 behavior on the convergence of different gradient aggregation rules in comparison to FA. Byzantine
 263 workers send random gradients and we vary the number of them from 1 to 3. Figure 4 shows that for
 264 some rules, i.e. Trimmed Mean, the presence of even a single Byzantine worker has a catastrophic
 265 impact. For other rules, as the number of Byzantine workers increases, filtering out the outliers
 266 becomes more challenging because the amount of noise increases. Regardless, FA remains more
 267 robust compared to other approaches.

268 Marginal utility of larger batch sizes under a fixed noise level:

269 We empirically verified the batch size required to identify our optimal Y^* - the FA matrix at each
 270 iteration. In particular, we fixed the noise level to $f = 3$ Byzantine workers and varied batch sizes.
 271 We show the results in Figure 5. **Our results indicate that, in cases where a larger batch size is
 272 a training requirement, FA achieves a significantly better accuracy compared to the existing
 273 state of the art aggregators.** This may be useful in some large scale vision applications, see [54, 55]
 274 for more details. Empirically, we can already see that our spectral relaxation to identify gradient
 275 subspace is effective in practice in all our experiments.

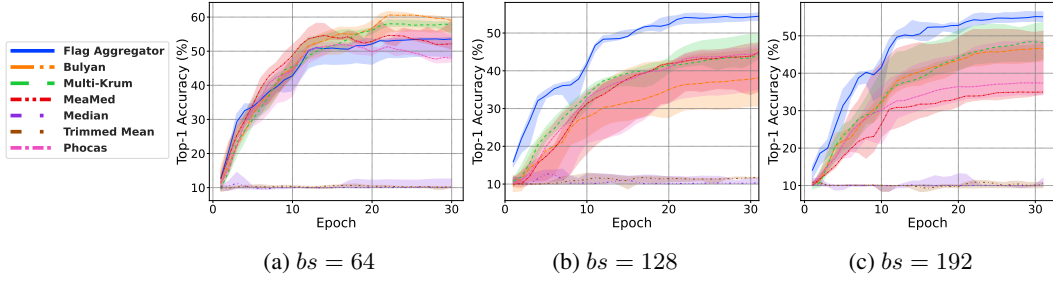


Figure 5: Marginal utility of larger batch sizes under a fixed noise level $f = 3$.

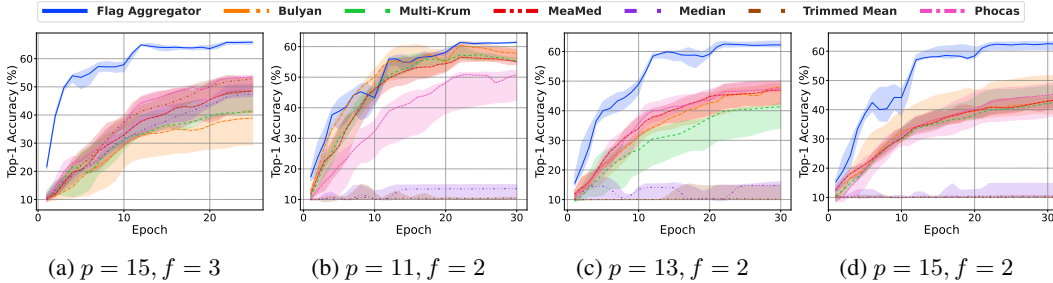


Figure 6: We present results under two different gradient attacks. The attack in (a) corresponds to simply dropping 10% of gradients from f workers. The attacks in (b)-(d) correspond to generic f workers sending random gradient vectors, i.e. we simply fix noise level while adding more workers.

276 **Tolerance to communication loss:** To analyze the effect of unreliable communication channels
 277 between the workers and the parameter server on convergence, we design an experiment where the
 278 physical link between some of the workers and the parameter server randomly drops a percentage of
 279 packets. Here, we set the loss rate of three links to 10% i.e., there are 3 Byzantine workers in our
 280 setting. The loss is introduced using the *netem* queuing discipline in Linux designed to emulate the
 281 properties of wide area networks [56]. The two main takeaways in Figure 6a are:

1. FA converges to a significantly higher accuracy than other aggregators, and thus is more robust to unreliable underlying network transports.
2. Considering time-to-accuracy for comparison, FA reaches a similar accuracy in less total number of training iterations, and thus is more robust to slow underlying network transports.

282 **Analyzing the marginal utility of additional workers.** To see the effect of adding more workers
 283 to a fixed number of Byzantine workers, we ran experiments where we fixed f , and increased p .
 284 Our experimental results shown in Figures 6b-6d indicate that our FA algorithm possesses strong
 285 resilience property for reasonable choices of p .

286 **The effect of having augmented data during training in Byzantine workers:** Figure 7 shows FA
 287 can handle nonlinear data augmentation in a much more stable fashion. Please see supplement for
 288 details on the level of noise, and exact solver settings that were used to obtain augmented images.

289 **The effect of the regularization parameter in FA:** The data-dependent regularization parameter λ
 290 in FA provides flexibility in the loss function to cover aggregators that benefit from pairwise distances
 291 such as Bulyan and Multi-Krum. To verify whether varying λ can interpolate Bulyan and Multi-Krum,
 292 we change λ in Figure 8. We can see when FA improves or performs similarly for a range of λ . Here,
 293 we set p and f to satisfy the strong Byzantine resilience condition of Bulyan, i.e. $p \geq 4f + 3$.

294 **Scaling out to real-world situations with more workers:** In distributed ML, p and f are usually
 295 large. To test high-dimensional settings commonly dealt in Semantic Vision with our FA, we used
 296 ResNet-18. Now, to specifically test the scalability of FA, we fully utilized our available GPU servers
 297 and set up to $p = 60$ workers (up to $f = 14$ Byzantine) with the MNIST dataset and a simple CNN
 298 with two convolutional layers followed by two fully connected layers (useful for simple detection).
 299 Figure 9 shows evidence that FA is feasible for larger setups.

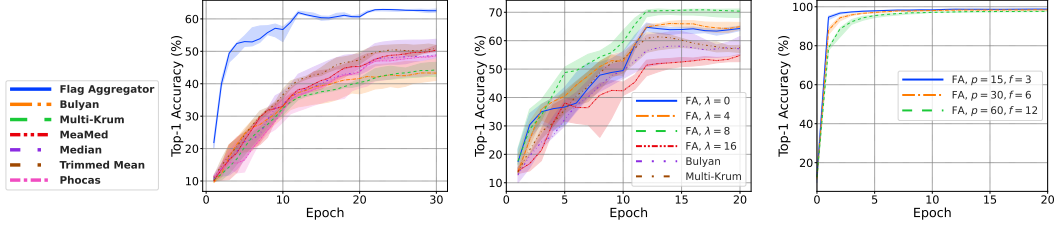


Figure 7: Accuracy of using augmented data in $f =$ ResNet-18, $p = 7$, and $f = 1$ 3 workers
 Figure 8: CIFAR10 with ResNet-18, $p = 7$, and larger setups
 Figure 9: Scaling FA to larger setups

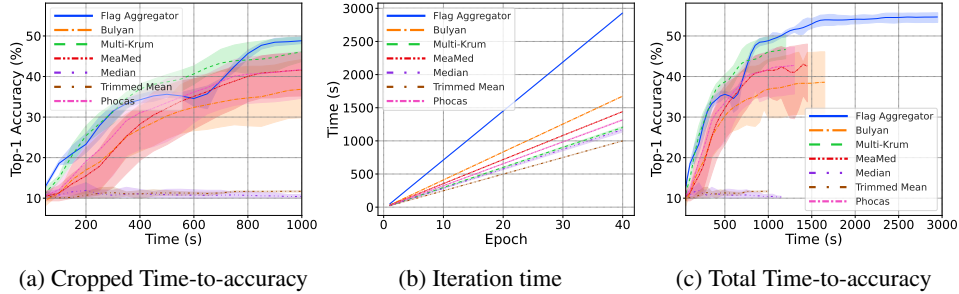


Figure 10: Wall clock time comparison

300 4 Discussion and Limitation

301 **Is it possible to fully “offload” FA computation to switches?** Recent work propose that aggregation
 302 be performed entirely on network infrastructure to alleviate any communication bottleneck that may
 303 arise [57, 58]. However, to the best of our knowledge, switches that are in use today only allow
 304 limited computation to be performed on gradient g_i as packets whenever they are transmitted [59, 60].
 305 That is, *programmability* is restrictive at the moment— switches used in practice have no floating
 306 point, or loop support, and are severely memory/state constrained. Fortunately, solutions seem near.
 307 For instance, [61] have already introduced support for floating point arithmetic in programmable
 308 switches. We may use quantization approaches for SVD calculation with some accuracy loss [62] to
 309 approximate floating point arithmetic. Offloading FA to switches has great potential in improving
 310 its computational complexity because the switch would perform as a high-throughput streaming
 311 parameter server to synchronize gradients over the network. Considering that FA’s accuracy currently
 312 outperforms its competition in several experiments, an offloaded FA can reach their accuracy even
 313 faster or it could reach a higher accuracy in the same amount of time.

314 **Potential Limitation.** Because in every iteration of FA, we perform SVD, the complexity of the
 315 algorithm would be $O(nN_\delta(\sum_{i=1}^p k_i)^2)$ with N_δ being the number of iterations for the algorithm.
 316 Figure 10 show the wall clock time it takes for FA to reach a certain accuracy (10a) or epoch(10b)
 317 compared to other methods under a fixed amount of random noise $f = 3$ with $p = 15$ workers.
 318 Although the iteration complexity of FA is higher, here each iteration has a higher utility as reflected in
 319 the time-to-accuracy measures. This makes FA comparable to others in a shorter time span, however,
 320 if there is more wall clock time to spare, FA converges to a better state as shown in Figure 10c where
 321 we let the same number of total iterations finish for all methods.

322 5 Conclusion

323 In this paper we proposed Flag Aggregator (FA) that can be used for robust aggregation of gradients
 324 in distributed training. FA is an optimization-based subspace estimator that formulates aggregation as
 325 a Maximum Likelihood Estimation procedure using Beta densities. We perform extensive evaluations
 326 of FA and show it can be effectively used in providing Byzantine resilience for gradient aggregation.
 327 Using techniques from convex optimization, we theoretically analyze FA and with tractable relaxations
 328 show its amenability to be solved by off-the-shelf solvers or first-order reweighing methods.

References

- 329
- 330 [1] Michel Grabisch, Jean-Luc Marichal, Radko Mesiar, and Endre Pap. *Aggregation functions*,
331 volume 127. Cambridge University Press, 2009.
- 332 [2] Sivaraman Balakrishnan, Simon S. Du, Jerry Li, and Aarti Singh. Computationally effi-
333 cient robust sparse estimation in high dimensions. In Satyen Kale and Ohad Shamir, edi-
334 tors, *Proceedings of the 2017 Conference on Learning Theory*, volume 65 of *Proceedings*
335 *of Machine Learning Research*, pages 169–212. PMLR, 07–10 Jul 2017. URL <https://proceedings.mlr.press/v65/balakrishnan17a.html>.
336
- 337 [3] Ilias Diakonikolas, Daniel Kane, Sushrut Karmalkar, Eric Price, and Alistair Stewart. Outlier-
338 robust high-dimensional sparse estimation via iterative filtering. *Advances in Neural Information*
339 *Processing Systems*, 32, 2019.
- 340 [4] Yu Cheng, Ilias Diakonikolas, Rong Ge, Shivam Gupta, Daniel M. Kane, and Mahdi
341 Soltanolkotabi. Outlier-robust sparse estimation via non-convex optimization. *Advances*
342 *in Neural Information Processing Systems*, 2022.
- 343 [5] Ilias Diakonikolas, Daniel M. Kane, Sushrut Karmalkar, Ankit Pensia, and Thanasis Pittas.
344 Robust sparse mean estimation via sum of squares. In Po-Ling Loh and Maxim Raginsky,
345 editors, *Proceedings of Thirty Fifth Conference on Learning Theory*, volume 178 of *Proceedings*
346 *of Machine Learning Research*, pages 4703–4763. PMLR, 02–05 Jul 2022. URL <https://proceedings.mlr.press/v178/diakonikolas22e.html>.
347
- 348 [6] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
349 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei.
350 ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*
351 (*IJCV*), 115(3):211–252, 2015. doi: 10.1007/s11263-015-0816-y.
- 352 [7] Konstantinos I Tsianos and Michael G Rabbat. Distributed strongly convex optimization. In
353 *2012 50th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*,
354 pages 593–600. IEEE, 2012.
- 355 [8] Tao Yang, Xinlei Yi, Junfeng Wu, Ye Yuan, Di Wu, Ziyang Meng, Yiguang Hong, Hong Wang,
356 Zongli Lin, and Karl H Johansson. A survey of distributed optimization. *Annual Reviews in*
357 *Control*, 47:278–305, 2019.
- 358 [9] Peva Blanchard, El Mahdi El Mhamdi, Rachid Guerraoui, and Julien Stainer. Machine learning
359 with adversaries: Byzantine tolerant gradient descent. In *Proceedings of the 31st Interna-*
360 *tional Conference on Neural Information Processing Systems, NIPS’17*, page 118–128. Curran
361 Associates Inc., 2017. ISBN 9781510860964.
- 362 [10] Sadegh Farhadkhani, Rachid Guerraoui, Nirupam Gupta, Rafael Pinot, and John Stephan.
363 Byzantine machine learning made easy by resilient averaging of momentums. In Kamalika
364 Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, edi-
365 tors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of
366 *Proceedings of Machine Learning Research*, pages 6246–6283. PMLR, 17–23 Jul 2022. URL
367 <https://proceedings.mlr.press/v162/farhadkhani22a.html>.
- 368 [11] Leonardo Bautista-Gomez, Ferad Zyulkyarov, Osman Unsal, and Simon McIntosh-Smith.
369 Unprotected computing: A large-scale study of dram raw error rate on a supercomputer. In *SC*
370 *’16: Proceedings of the International Conference for High Performance Computing, Networking,*
371 *Storage and Analysis*, pages 645–655, 2016. doi: 10.1109/SC.2016.54.
- 372 [12] Bianca Schroeder and Garth A. Gibson. Disk failures in the real world: What does an mttf
373 of 1,000,000 hours mean to you? In *Proceedings of the 5th USENIX Conference on File and*
374 *Storage Technologies, FAST ’07*, page 1–es, USA, 2007. USENIX Association.
- 375 [13] Phillipa Gill, Navendu Jain, and Nachiappan Nagappan. Understanding network failures
376 in data centers: Measurement, analysis, and implications. *SIGCOMM Comput. Commun.*
377 *Rev.*, 41(4):350–361, aug 2011. ISSN 0146-4833. doi: 10.1145/2043164.2018477. URL
378 <https://doi.org/10.1145/2043164.2018477>.

- 379 [14] Gilad Baruch, Moran Baruch, and Yoav Goldberg. A little is enough: Circumventing defenses
380 for distributed learning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox,
381 and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32.
382 Curran Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper/2019/
383 file/ec1c59141046cd1866bbcbdfb6ae31d4-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/ec1c59141046cd1866bbcbdfb6ae31d4-Paper.pdf).
- 384 [15] Guosai Wang, Lifei Zhang, and Wei Xu. What can we learn from four years of data center
385 hardware failures? In *2017 47th Annual IEEE/IFIP International Conference on Dependable
386 Systems and Networks (DSN)*, pages 25–36, 2017. doi: 10.1109/DSN.2017.26.
- 387 [16] Devesh Tiwari, Saurabh Gupta, James Rogers, Don Maxwell, Paolo Rech, Sudharshan Vazhku-
388 dai, Daniel Oliveira, Dave Londo, Nathan DeBardeleben, Philippe Navaux, Luigi Carro, and
389 Arthur Bland. Understanding gpu errors on large-scale hpc systems and the implications for
390 system design and operation. In *2015 IEEE 21st International Symposium on High Performance
391 Computer Architecture (HPCA)*, pages 331–342, 2015. doi: 10.1109/HPCA.2015.7056044.
- 392 [17] Bin Nie, Devesh Tiwari, Saurabh Gupta, Evgenia Smiri, and James H. Rogers. A large-scale
393 study of soft-errors on gpus in the field. In *2016 IEEE International Symposium on High
394 Performance Computer Architecture (HPCA)*, pages 519–530, 2016. doi: 10.1109/HPCA.2016.
395 7446091.
- 396 [18] Sheldon M Ross. *Introduction to probability models*. Academic press, 2014.
- 397 [19] Fanny Yang, Zuowen Wang, and Christina Heinze-Deml. Invariance-inducing regulariza-
398 tion using worst-case transformations suffices to boost accuracy and spatial robustness.
399 In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Gar-
400 nett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran
401 Associates, Inc., 2019. URL [https://proceedings.neurips.cc/paper/2019/file/
402 1d01bd2e16f57892f0954902899f0692-Paper.pdf](https://proceedings.neurips.cc/paper/2019/file/1d01bd2e16f57892f0954902899f0692-Paper.pdf).
- 403 [20] Christina Heinze-Deml and Nicolai Meinshausen. Conditional variance penalties and domain
404 shift robustness, 2017. URL <https://arxiv.org/abs/1710.11469>.
- 405 [21] Saeid Motiian, Marco Piccirilli, Donald A. Adjeroh, and Gianfranco Doretto. Unified deep su-
406 pervised domain adaptation and generalization. In *IEEE International Conference on Computer
407 Vision (ICCV)*, 2017.
- 408 [22] Sravanti Addepalli, Samyak Jain, et al. Efficient and effective augmentation strategy for
409 adversarial training. *Advances in Neural Information Processing Systems*, 35:1488–1501, 2022.
- 410 [23] Eric Wong, Leslie Rice, and J Zico Kolter. Fast is better than free: Revisiting adversarial
411 training. *arXiv preprint arXiv:2001.03994*, 2020.
- 412 [24] Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized
413 smoothing. In *international conference on machine learning*, pages 1310–1320. PMLR, 2019.
- 414 [25] Youssef Allouah, Rachid Guerraoui, Nirupam Gupta, Rafael Pinot, and John Stephan. Dis-
415 tributed learning with curious and adversarial machines. *arXiv preprint arXiv:2302.04787*,
416 2023.
- 417 [26] Youssef Allouah, Sadegh Farhadkhani, Rachid Guerraoui, Nirupam Gupta, Rafaël Pinot, and
418 John Stephan. Fixing by mixing: A recipe for optimal byzantine ml under heterogeneity. In
419 *International Conference on Artificial Intelligence and Statistics*, pages 1232–1300. PMLR,
420 2023.
- 421 [27] Sadegh Farhadkhani, Rachid Guerraoui, Nirupam Gupta, Rafael Pinot, and John Stephan.
422 Byzantine machine learning made easy by resilient averaging of momentums. In *International
423 Conference on Machine Learning*, pages 6246–6283. PMLR, 2022.
- 424 [28] P-A Absil. *Optimization algorithms on matrix manifolds*. Princeton University Press, 2008.
- 425 [29] John Wright, Arvind Ganesh, Shankar Rao, Yigang Peng, and Yi Ma. Robust principal
426 component analysis: Exact recovery of corrupted low-rank matrices via convex optimization.
427 *Advances in neural information processing systems*, 22, 2009.

- 428 [30] Matthias Hein and Thomas Bühler. An inverse power method for nonlinear eigenproblems with
429 applications in 1-spectral clustering and sparse pca. *Advances in neural information processing*
430 *systems*, 23, 2010.
- 431 [31] Rudrasis Chakraborty, Soren Hauberg, and Baba C Vemuri. Intrinsic grassmann averages for
432 online linear and robust subspace learning. In *Proceedings of the IEEE Conference on Computer*
433 *Vision and Pattern Recognition*, pages 6196–6204, 2017.
- 434 [32] D. Monk. The geometry of flag manifolds. *Proceedings of the London Mathematical So-*
435 *ciety*, s3-9(2):253–286, 1959. doi: <https://doi.org/10.1112/plms/s3-9.2.253>. URL <https://londmathsoc.onlinelibrary.wiley.com/doi/abs/10.1112/plms/s3-9.2.253>.
436
- 437 [33] Ke Ye, Ken Sze-Wai Wong, and Lek-Heng Lim. Optimization on flag manifolds. *Mathematical*
438 *Programming*, 194(1-2):621–660, 2022.
- 439 [34] Nathan Mankovich, Emily J King, Chris Peterson, and Michael Kirby. The flag median
440 and flagirls. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
441 *Recognition*, pages 10339–10347, 2022.
- 442 [35] Kevin P Murphy. *Probabilistic machine learning: an introduction*. MIT press, 2022.
- 443 [36] Massimo Fornasier, Holger Rauhut, and Rachel Ward. Low-rank matrix recovery via iteratively
444 reweighted least squares minimization. *SIAM Journal on Optimization*, 21(4):1614–1640, 2011.
- 445 [37] Toni Karvonen and Chris J Oates. Maximum likelihood estimation in gaussian process regression
446 is ill-posed. *Journal of Machine Learning Research*, 24(120):1–47, 2023.
- 447 [38] Georgios Damaskinos, El-Mahdi El-Mhamdi, Rachid Guerraoui, Arsany Guirguis, and Sébastien
448 Rouault. Aggregathor: Byzantine machine learning via robust gradient aggregation. *Proceedings*
449 *of Machine Learning and Systems*, 1:81–106, 2019.
- 450 [39] Yuxin Chen, Yuejie Chi, Jianqing Fan, Cong Ma, and Yuling Yan. Noisy matrix completion:
451 Understanding statistical guarantees for convex relaxation via nonconvex optimization. *SIAM*
452 *journal on optimization*, 30(4):3098–3121, 2020.
- 453 [40] Dong Yin, Yudong Chen, Ramchandran Kannan, and Peter Bartlett. Byzantine-robust distributed
454 learning: Towards optimal statistical rates. In *Proceedings of the 35th International Conference*
455 *on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 5650–
456 5659. PMLR, 10–15 Jul 2018.
- 457 [41] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan,
458 Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas
459 Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy,
460 Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style,
461 high-performance deep learning library. In *Advances in Neural Information Processing Systems*,
462 volume 32, 2019.
- 463 [42] Rachid Guerraoui, Arsany Guirguis, Jérémy Plassmann, Anton Ragot, and Sébastien Rouault.
464 Garfield: System support for byzantine machine learning (regular paper). In *2021 51st Annual*
465 *IEEE/IFIP International Conference on Dependable Systems and Networks (DSN)*, pages 39–51,
466 2021. doi: 10.1109/DSN48987.2021.00021.
- 467 [43] Cong Xie, Oluwasanmi Koyejo, and Indranil Gupta. Generalized byzantine-tolerant sgd, 2018.
- 468 [44] Cong Xie, Oluwasanmi Koyejo, and Indranil Gupta. Phocas: dimensional byzantine-resilient
469 stochastic gradient descent, 2018.
- 470 [45] El Mahdi El Mhamdi, Rachid Guerraoui, and Sébastien Rouault. The hidden vulnerability of
471 distributed learning in Byzantium. In *Proceedings of the 35th International Conference on*
472 *Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 3521–3530.
473 PMLR, 10–15 Jul 2018.

- 474 [46] Trishul Chilimbi, Yutaka Suzue, Johnson Apacible, and Karthik Kalyanaraman. Project
475 adam: Building an efficient and scalable deep learning training system. In *Proceedings of the*
476 *11th USENIX Conference on Operating Systems Design and Implementation*, OSDI'14, page
477 571–582, USA, 2014.
- 478 [47] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
479 recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*,
480 pages 770–778, 2016. doi: 10.1109/CVPR.2016.90.
- 481 [48] Alex Krizhevsky. Learning multiple layers of features from tiny images, 2009. URL <https://www.cs.toronto.edu/~kriz/learning-features-2009-TR.pdf>.
482
- 483 [49] Yann LeCun and Corinna Cortes. MNIST handwritten digit database, 2010. URL <http://yann.lecun.com/exdb/mnist/>.
484
- 485 [50] Ya Le and Xuan S. Yang. Tiny imagenet visual recognition challenge, 2015.
- 486 [51] David Kelly. Rough path recursions and diffusion approximations. *The Annals of Applied*
487 *Probability*, 26(1):425–461, 2016.
- 488 [52] Jianghong Bao and Qigui Yang. Period of the discrete arnold cat map and general cat map.
489 *Nonlinear Dynamics*, 70(2):1365–1375, 2012.
- 490 [53] Fundamental algorithms for scientific computing in python. <https://scipy.org/>, 2023.
- 491 [54] Nitish Shirish Keskar, Dheevatsa Mudigere, Jorge Nocedal, Mikhail Smelyanskiy, and Ping
492 Tak Peter Tang. On large-batch training for deep learning: Generalization gap and sharp minima.
493 In *International Conference on Learning Representations*, 2017. URL <https://openreview.net/forum?id=H1oyRlygg>.
494
- 495 [55] Yang You, Jonathan Hseu, Chris Ying, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh.
496 Large-batch training for lstm and beyond. In *Proceedings of the International Conference for*
497 *High Performance Computing, Networking, Storage and Analysis*, SC '19, 2019.
- 498 [56] Kevin Hsieh, Aaron Harlap, Nandita Vijaykumar, Dimitris Konomis, Gregory R. Ganger,
499 Phillip B. Gibbons, and Onur Mutlu. Gaia: Geo-distributed machine learning approaching lan
500 speeds. In *Proceedings of the 14th USENIX Conference on Networked Systems Design and*
501 *Implementation*, page 629–647, 2017.
- 502 [57] Amedeo Sapio, Marco Canini, Chen-Yu Ho, Jacob Nelson, Panos Kalnis, Changhoon Kim,
503 Arvind Krishnamurthy, Masoud Moshref, Dan Ports, and Peter Richtárik. Scaling distributed
504 machine learning with {In-Network} aggregation. In *18th USENIX Symposium on Networked*
505 *Systems Design and Implementation (NSDI 21)*, pages 785–808, 2021.
- 506 [58] ChonLam Lao, Yanfang Le, Kshiteej Mahajan, Yixi Chen, Wenfei Wu, Aditya Akella, and
507 Michael Swift. {ATP}: In-network aggregation for multi-tenant learning. In *18th USENIX*
508 *Symposium on Networked Systems Design and Implementation (NSDI 21)*, pages 741–761,
509 2021.
- 510 [59] Pat Bosshart, Glen Gibb, Hun-Seok Kim, George Varghese, Nick McKeown, Martin Izzard,
511 Fernando Mujica, and Mark Horowitz. Forwarding metamorphosis: Fast programmable match-
512 action processing in hardware for sdn. In *Proceedings of the ACM SIGCOMM 2013 Conference*
513 *on SIGCOMM*, SIGCOMM '13, page 99–110, New York, NY, USA, 2013. Association for
514 Computing Machinery. ISBN 9781450320566. doi: 10.1145/2486001.2486011. URL <https://doi.org/10.1145/2486001.2486011>.
515
- 516 [60] N McKeown. Pisa: Protocol independent switch architecture. In *P4 Workshop*, 2015.
- 517 [61] Yifan Yuan, Omar Alama, Jiawei Fei, Jacob Nelson, Dan RK Ports, Amedeo Sapio, Marco
518 Canini, and Nam Sung Kim. Unlocking the power of inline {Floating-Point} operations on
519 programmable switches. In *19th USENIX Symposium on Networked Systems Design and*
520 *Implementation (NSDI 22)*, pages 683–700, 2022.

- 521 [62] Zhourui Song, Zhenyu Liu, and Dongsheng Wang. Computation error analysis of block floating
522 point arithmetic oriented convolution neural network accelerator design. In *Proceedings of the*
523 *Thirty-Second AAAI Conference on Artificial Intelligence and Thirtieth Innovative Applications*
524 *of Artificial Intelligence Conference and Eighth AAAI Symposium on Educational Advances in*
525 *Artificial Intelligence*, AAAI'18/IAAI'18/EAAI'18. AAAI Press, 2018. ISBN 978-1-57735-
526 800-8.