



Roubust Learning

Amir H. Payberah
payberah@kth.se
2020-11-09





The Course Web Page

`https://fid3024.github.io`



Adversarial Goals



Adversarial Goals

- ▶ Confidentiality and privacy
 - Confidentiality of the **model** itself (e.g., intellectual property)
 - Privacy of the **training** or **test data** (e.g., medical records)



Adversarial Goals

- ▶ Confidentiality and privacy
 - Confidentiality of the **model** itself (e.g., intellectual property)
 - Privacy of the **training** or **test data** (e.g., medical records)
- ▶ Integrity
 - Integrity of the **predictions**



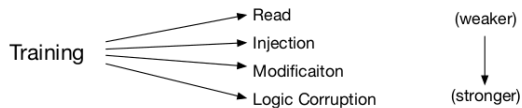
Adversarial Goals

- ▶ **Confidentiality** and **privacy**
 - Confidentiality of the **model** itself (e.g., intellectual property)
 - Privacy of the **training** or **test data** (e.g., medical records)
- ▶ **Integrity**
 - Integrity of the **predictions**
- ▶ **Availability**
 - Availability of the **system** deploying machine learning



Adversarial Capabilities for Integrity Attacks

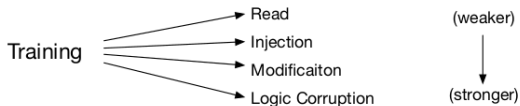
► Training phase



[Papernot et al., SoK: Security and Privacy in Machine Learning, 2018]

Adversarial Capabilities for Integrity Attacks

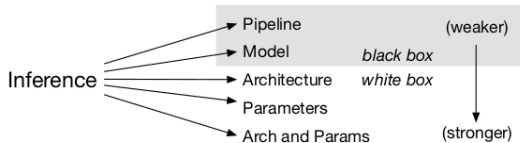
▶ Training phase



[Papernot et al., SoK: Security and Privacy in Machine Learning, 2018]

▶ Inference phase

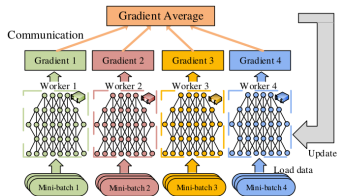
- White box
- Black box



[Papernot et al., SoK: Security and Privacy in Machine Learning, 2018]

Our Focus and Goal

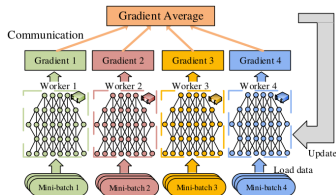
- ▶ Data parallelization



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

Our Focus and Goal

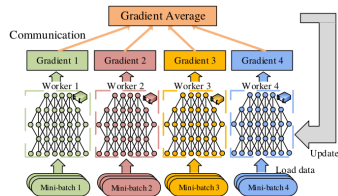
- ▶ Data parallelization
- ▶ Each worker is prone to adversarial attack.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

Our Focus and Goal

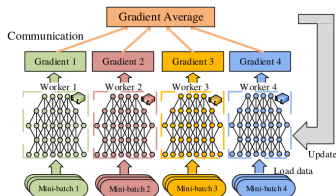
- ▶ Data parallelization
- ▶ Each worker is prone to **adversarial attack**.
- ▶ **Adversarial attacks**: some unknown subset of computing devices are **compromised and behave adversarially** (e.g., sending out malicious messages)



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

Our Focus and Goal

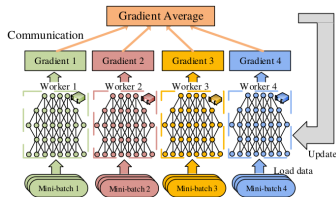
- ▶ Data parallelization
- ▶ Each worker is prone to **adversarial attack**.
- ▶ **Adversarial attacks**: some unknown subset of computing devices are **compromised and behave adversarially** (e.g., sending out malicious messages)
- ▶ Our goal: **integrity** of the model in the **training** phase



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

Distributed Stochastic Gradient Descent (1/3)

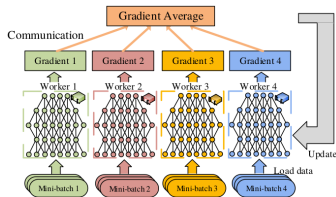
- ▶ One parameter server, and n workers.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

Distributed Stochastic Gradient Descent (1/3)

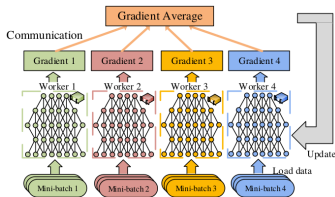
- ▶ One **parameter server**, and **n** workers.
- ▶ Computation is divided into **synchronous rounds**.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

Distributed Stochastic Gradient Descent (1/3)

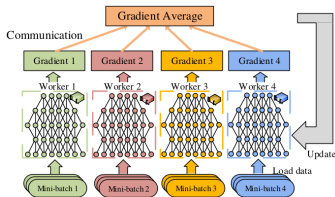
- ▶ One **parameter server**, and **n** workers.
- ▶ Computation is divided into **synchronous rounds**.
- ▶ During round **t** , the **parameter server** broadcasts its parameter vector $\mathbf{w} \in \mathbb{R}^d$ to all the **workers**.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

Distributed Stochastic Gradient Descent (2/3)

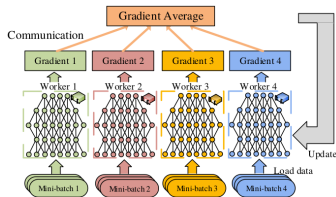
- ▶ At each round t , each correct worker i computes $G_i(\mathbf{w}_t, \beta)$.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

Distributed Stochastic Gradient Descent (2/3)

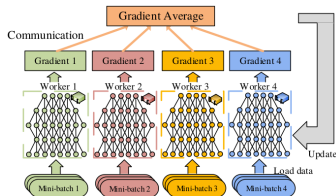
- ▶ At each round t , each **correct worker** i computes $G_i(\mathbf{w}_t, \beta)$.
- ▶ $G_i(\mathbf{w}_t, \beta)$: the **local estimate** of the gradient of the loss function $\nabla J(\mathbf{w}_t)$.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

Distributed Stochastic Gradient Descent (2/3)

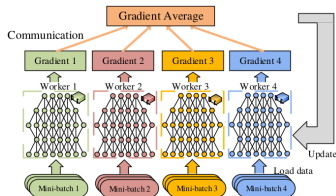
- ▶ At each round t , each **correct worker** i computes $G_i(\mathbf{w}_t, \beta)$.
- ▶ $G_i(\mathbf{w}_t, \beta)$: the **local estimate** of the gradient of the loss function $\nabla J(\mathbf{w}_t)$.
- ▶ β : a mini-batch of **i.i.d. samples** drawn from the dataset.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

Distributed Stochastic Gradient Descent (2/3)

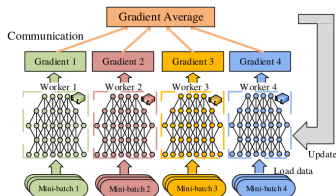
- ▶ At each round t , each **correct worker** i computes $G_i(\mathbf{w}_t, \beta)$.
- ▶ $G_i(\mathbf{w}_t, \beta)$: the **local estimate** of the gradient of the loss function $\nabla J(\mathbf{w}_t)$.
- ▶ β : a mini-batch of **i.i.d. samples** drawn from the dataset.
- ▶ $G_i(\mathbf{w}_t, \beta) = \frac{1}{|\beta|} \sum_{\mathbf{x} \in \beta} \nabla l_i(\mathbf{w}_t, \mathbf{x})$



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

Distributed Stochastic Gradient Descent (3/3)

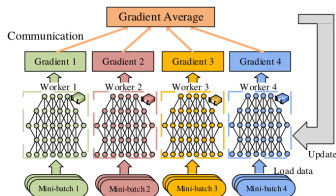
- ▶ The parameter server computes $F(G_1, G_2, \dots, G_n)$



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

Distributed Stochastic Gradient Descent (3/3)

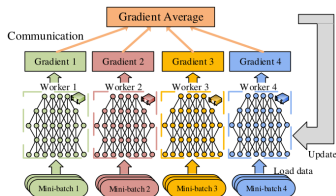
- ▶ The parameter server computes $F(G_1, G_2, \dots, G_n)$
- ▶ Gradient Aggregation Rule (GAR): $F(G_1, G_2, \dots, G_n) = \frac{1}{n} \sum_{i=1}^n G_i$



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

Distributed Stochastic Gradient Descent (3/3)

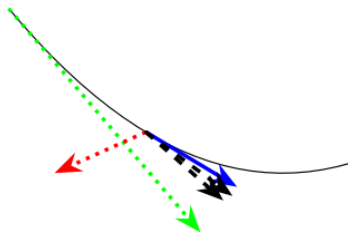
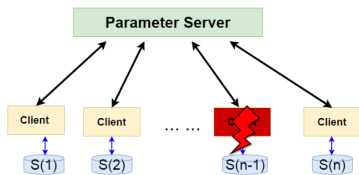
- ▶ The parameter server computes $F(G_1, G_2, \dots, G_n)$
- ▶ Gradient Aggregation Rule (GAR): $F(G_1, G_2, \dots, G_n) = \frac{1}{n} \sum_{i=1}^n G_i$
- ▶ The parameter server updates the parameter vector $\mathbf{w} \leftarrow \mathbf{w} - \gamma F(G_1, G_2, \dots, G_n)$



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]

Distributed SGD with Byzantine Workers

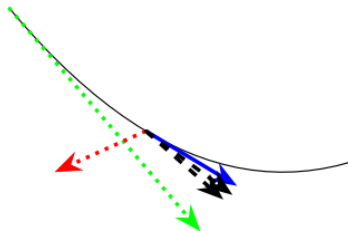
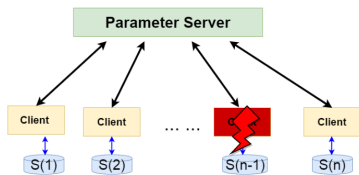
- ▶ Among the n workers, f of them are possibly Byzantine (behaving arbitrarily).



[El-Mhamdi et al., Fast and Secure Distributed Learning in High Dimension, 2019]

Distributed SGD with Byzantine Workers

- ▶ Among the n workers, f of them are possibly Byzantine (behaving arbitrarily).
- ▶ A Byzantine worker b proposes a vector G_b that can deviate arbitrarily from the vector it is supposed.



[El-Mhamdi et al., Fast and Secure Distributed Learning in High Dimension, 2019]



Averaging GAR and Byzantine Workers

- ▶ Averaging GAR: $F(G_1, G_2, \dots, G_n) = \frac{1}{n} \sum_{i=1}^n G_i$
- ▶ $\mathbf{w} \leftarrow \mathbf{w} - \gamma F(G_1, G_2, \dots, G_n)$



Averaging GAR and Byzantine Workers

- ▶ Averaging GAR: $F(G_1, G_2, \dots, G_n) = \frac{1}{n} \sum_{i=1}^n G_i$
- ▶ $\mathbf{w} \leftarrow \mathbf{w} - \gamma F(G_1, G_2, \dots, G_n)$
- ▶ Even a single Byzantine worker can prevent convergence.



Averaging GAR and Byzantine Workers

- ▶ Averaging GAR: $F(G_1, G_2, \dots, G_n) = \frac{1}{n} \sum_{i=1}^n G_i$
- ▶ $\mathbf{w} \leftarrow \mathbf{w} - \gamma F(G_1, G_2, \dots, G_n)$
- ▶ Even a **single Byzantine** worker can **prevent convergence**.
- ▶ **Proof:** if the Byzantine worker proposes $G_n = nU - \sum_{i=1}^{n-1} G_i$, then $F = U$.



(α, f) -Byzantine-Resilience (1/2)

- ▶ Assume n workers, where f of them are Byzantine workers.



(α, f) -Byzantine-Resilience (1/2)

- ▶ Assume n workers, where f of them are Byzantine workers.
- ▶ $\alpha \in [0, \pi/2]$ and $f \in \{0, \dots, n\}$.



(α, f) -Byzantine-Resilience (1/2)

- ▶ Assume n workers, where f of them are Byzantine workers.
- ▶ $\alpha \in [0, \pi/2]$ and $f \in \{0, \dots, n\}$.
- ▶ $(G_1, \dots, G_{n-f}) \in (\mathbb{R}^d)^{n-f}$ are i.i.d. random vectors



(α, f) -Byzantine-Resilience (1/2)

- ▶ Assume n workers, where f of them are Byzantine workers.
- ▶ $\alpha \in [0, \pi/2]$ and $f \in \{0, \dots, n\}$.
- ▶ $(\mathbf{G}_1, \dots, \mathbf{G}_{n-f}) \in (\mathbb{R}^d)^{n-f}$ are i.i.d. random vectors
 - $\mathbf{G}_i \sim g$
 - $\mathbb{E}[g] = \mathcal{J}$, where $\mathcal{J} = \nabla J(\mathbf{w})$



(α, f) -Byzantine-Resilience (1/2)

- ▶ Assume n workers, where f of them are Byzantine workers.
- ▶ $\alpha \in [0, \pi/2]$ and $f \in \{0, \dots, n\}$.
- ▶ $(\mathbf{G}_1, \dots, \mathbf{G}_{n-f}) \in (\mathbb{R}^d)^{n-f}$ are i.i.d. random vectors
 - $\mathbf{G}_i \sim g$
 - $\mathbb{E}[g] = \mathcal{J}$, where $\mathcal{J} = \nabla J(\mathbf{w})$
- ▶ $(\mathbf{B}_1, \dots, \mathbf{B}_f) \in (\mathbb{R}^d)^f$ are random vectors, possibly dependent between them and the vectors $(\mathbf{G}_1, \dots, \mathbf{G}_{n-f})$

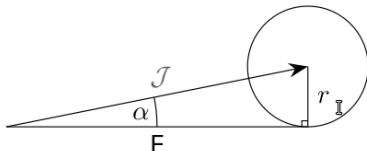


(α, f) -Byzantine-Resilience (2/2)

- ▶ A GAR F is said to be (α, f) -Byzantine-resilient if, for any $1 \leq j_1 < \dots < j_f \leq n$, the vector $F(G_1, \dots, B_1, \dots, B_f, \dots, G_n)$ satisfies:

(α, f) -Byzantine-Resilience (2/2)

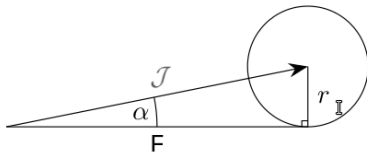
- A GAR \mathbf{F} is said to be (α, f) -Byzantine-resilient if, for any $1 \leq j_1 < \dots < j_f \leq n$, the vector $\mathbf{F}(\mathbf{G}_1, \dots, \mathbf{B}_1, \dots, \mathbf{B}_f, \dots, \mathbf{G}_n)$ satisfies:
1. Vector \mathbf{F} that is not too far from the real gradient \mathcal{J} , i.e., $\|\mathbb{E}[\mathbf{F}] - \mathcal{J}\| \leq r$.



[Blanchard et al., Machine Learning with Adversaries: Byzantine Tolerant Gradient Descent, 2017]

(α, f) -Byzantine-Resilience (2/2)

- A GAR \mathbf{F} is said to be (α, f) -Byzantine-resilient if, for any $1 \leq j_1 < \dots < j_f \leq n$, the vector $\mathbf{F}(\mathbf{G}_1, \dots, \mathbf{B}_1, \dots, \mathbf{B}_f, \dots, \mathbf{G}_n)$ satisfies:
1. Vector \mathbf{F} that is **not too far** from the **real gradient** \mathcal{J} , i.e., $\|\mathbb{E}[\mathbf{F}] - \mathcal{J}\| \leq r$.
 2. Moments of \mathbf{F} should be **controlled** by the moments of the (correct) gradient estimator \mathbf{g} , where $\mathbb{E}[\mathbf{g}] = \mathcal{J}$.



[Blanchard et al., Machine Learning with Adversaries: Byzantine Tolerant Gradient Descent, 2017]



Byzantine-Resilience GAR

- ▶ Median
- ▶ Krum
- ▶ Multi-Krum
- ▶ Brute



Median

▶ $n \geq 2f + 1$



Median

▶ $n \geq 2f + 1$

▶ $\text{median}(x_1, \dots, x_n) = \arg \min_{x \in \mathbb{R}} \sum_{i=1}^n |x_i - x|$



Median

▶ $n \geq 2f + 1$

▶ $\text{median}(x_1, \dots, x_n) = \arg \min_{x \in \mathbb{R}} \sum_{i=1}^n |x_i - x|$

▶ d : the gradient vectors **dimension**.

▶ **Geometric** median

$$F = \text{GeoMed}(G_1, \dots, G_n) = \arg \min_{G \in \mathbb{R}^d} \sum_{i=1}^n \|G_i - G\|$$

- ▶ $n \geq 2f + 1$
- ▶ $\text{median}(x_1, \dots, x_n) = \arg \min_{x \in \mathbb{R}} \sum_{i=1}^n |x_i - x|$
- ▶ d : the gradient vectors **dimension**.

- ▶ **Geometric** median

$$F = \text{GeoMed}(G_1, \dots, G_n) = \arg \min_{G \in \mathbb{R}^d} \sum_{i=1}^n \|G_i - G\|$$

- ▶ **Marginal** median

$$F = \text{MarMed}(G_1, \dots, G_n) = \begin{pmatrix} \text{median}(G_1[1], \dots, G_n[1]) \\ \vdots \\ \text{median}(G_1[d], \dots, G_n[d]) \end{pmatrix} \quad (1)$$



Krum

► $n \geq 2f + 3$



Krum

- ▶ $n \geq 2f + 3$
- ▶ Idea: to preclude the vectors that are too far away.



Krum

- ▶ $n \geq 2f + 3$
- ▶ Idea: to **preclude** the vectors that are **too far away**.
- ▶ $s(i) = \sum_{i \rightarrow j} \|G_i - G_j\|^2$, the score of the worker i .



- ▶ $n \geq 2f + 3$
- ▶ Idea: to **preclude** the vectors that are **too far away**.
- ▶ $s(i) = \sum_{i \rightarrow j} \|G_i - G_j\|^2$, the score of the worker i .
- ▶ $i \rightarrow j$ denotes that G_j belongs to the $n - f - 2$ closest vectors to G_i .



- ▶ $n \geq 2f + 3$
- ▶ Idea: to **preclude** the vectors that are **too far away**.
- ▶ $s(i) = \sum_{i \rightarrow j} \|G_i - G_j\|^2$, the score of the worker i .
- ▶ $i \rightarrow j$ denotes that G_j belongs to the $n - f - 2$ closest vectors to G_i .
- ▶ $F(G_1, \dots, G_n) = G_{i^*}$



- ▶ $n \geq 2f + 3$
- ▶ Idea: to **preclude** the vectors that are **too far away**.
- ▶ $s(i) = \sum_{i \rightarrow j} \|G_i - G_j\|^2$, the score of the worker i .
- ▶ $i \rightarrow j$ denotes that G_j belongs to the $n - f - 2$ closest vectors to G_i .
- ▶ $F(G_1, \dots, G_n) = G_{i_*}$
- ▶ G_{i_*} refers to the worker minimizing the score, $s(i_*) \leq s(i)$ for all i .



Multi-Krum

- ▶ **Multi-Krum** computes the **score** for each vector proposed (as in Krum).



Multi-Krum

- ▶ **Multi-Krum** computes the **score** for each vector proposed (as in Krum).
- ▶ It selects **m** vectors G_{1*}, \dots, G_{m*} , which score the **best** ($1 \leq m \leq n - f - 2$).

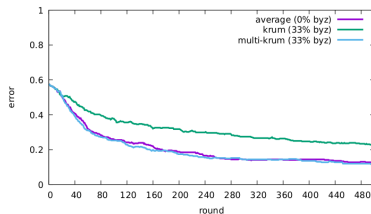


Multi-Krum

- ▶ **Multi-Krum** computes the **score** for each vector proposed (as in Krum).
- ▶ It selects **m** vectors G_{1*}, \dots, G_{m*} , which score the **best** ($1 \leq m \leq n - f - 2$).
- ▶ It outputs their average $\frac{1}{m} \sum_i G_{i*}$.

Multi-Krum

- ▶ **Multi-Krum** computes the **score** for each vector proposed (as in Krum).
- ▶ It selects m vectors G_{1*}, \dots, G_{m*} , which score the **best** ($1 \leq m \leq n - f - 2$).
- ▶ It outputs their average $\frac{1}{m} \sum_i G_{i*}$.
- ▶ The cases $m = 1$ and $m = n$ correspond to **Krum** and **averaging**, respectively.



[Blanchard et al., Machine Learning with Adversaries: Byzantine Tolerant Gradient Descent, 2017]



Brute

▶ $n \geq 2f + 1$



Brute

▶ $n \geq 2f + 1$

▶ $Q = \{G_1, G_2, \dots, G_n\}$



Brute

- ▶ $n \geq 2f + 1$
- ▶ $\mathcal{Q} = \{G_1, G_2, \dots, G_n\}$
- ▶ $\mathcal{R} = \{\mathcal{X} \mid \mathcal{X} \subset \mathcal{Q}, |\mathcal{X}| = n - f\}$
 - The set of all the **subsets** of $n - f$



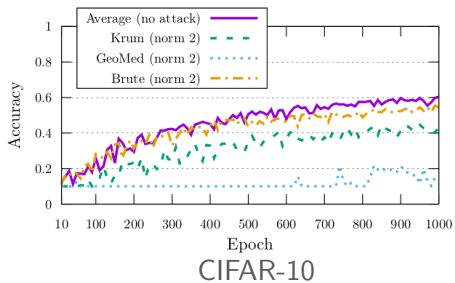
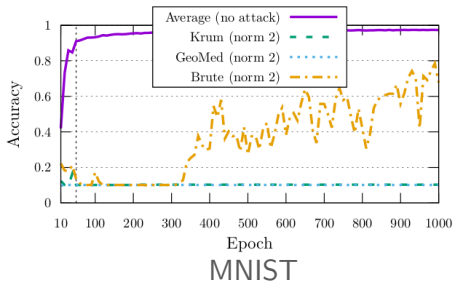
Brute

- ▶ $n \geq 2f + 1$
- ▶ $\mathcal{Q} = \{G_1, G_2, \dots, G_n\}$
- ▶ $\mathcal{R} = \{\mathcal{X} \mid \mathcal{X} \subset \mathcal{Q}, |\mathcal{X}| = n - f\}$
 - The set of all the **subsets** of $n - f$
- ▶ $\mathcal{S} = \arg \min_{\mathcal{X} \in \mathcal{R}} (\max_{(G_i, G_j) \in \mathcal{X}^2} (\|G_i - G_j\|))$
 - Selects the $n - f$ **most clumped gradients** among the submitted ones.



Brute

- ▶ $n \geq 2f + 1$
- ▶ $\mathcal{Q} = \{G_1, G_2, \dots, G_n\}$
- ▶ $\mathcal{R} = \{\mathcal{X} \mid \mathcal{X} \subset \mathcal{Q}, |\mathcal{X}| = n - f\}$
 - The set of all the **subsets** of $n - f$
- ▶ $\mathcal{S} = \arg \min_{\mathcal{X} \in \mathcal{R}} (\max_{(G_i, G_j) \in \mathcal{X}^2} (\|G_i - G_j\|))$
 - Selects the $n - f$ **most clumped gradients** among the submitted ones.
- ▶ $F(G_1, \dots, G_n) = \frac{1}{n-f} \sum_{G \in \mathcal{S}} G$



[El Mhamdi et al., The Hidden Vulnerability of Distributed Learning in Byzantium, 2018]



Weak Byzantine Resilience

- ▶ **Limitation** of previous aggregation methods.
- ▶ If gradient dimension $d \gg 1$, then the distance function between two vectors $\|X - Y\|_p$, cannot distinguish these two cases:



Weak Byzantine Resilience

- ▶ **Limitation** of previous aggregation methods.
- ▶ If gradient dimension $d \gg 1$, then the distance function between two vectors $\|X - Y\|_p$, cannot distinguish these two cases:
 - ▶ 1. Does X and Y **disagree** a **bit** on each coordinate?



Weak Byzantine Resilience

- ▶ **Limitation** of previous aggregation methods.
- ▶ If gradient dimension $d \gg 1$, then the distance function between two vectors $\|X - Y\|_p$, cannot distinguish these two cases:
 - ▶ 1. Does X and Y **disagree** a **bit** on each coordinate?
 - ▶ 2. Does X and Y **disagree** a **lot** on **only one**?



Strong Byzantine Resilience

- ▶ Ensuring convergence (as in weak Byzantine resilience functions).



Strong Byzantine Resilience

- ▶ Ensuring convergence (as in weak Byzantine resilience functions).
- ▶ Ensures that each coordinate is agreed on by a majority of vectors that were selected by a Byzantine resilient aggregation rule A .



Strong Byzantine Resilience

- ▶ Ensuring convergence (as in weak Byzantine resilience functions).
- ▶ Ensures that each coordinate is agreed on by a majority of vectors that were selected by a Byzantine resilient aggregation rule A .
- ▶ A can be Brute, Krum, Median, etc.



Strong Byzantine Resilience

- ▶ Ensuring **convergence** (as in **weak Byzantine resilience** functions).
- ▶ Ensures that **each coordinate** is agreed on by a **majority of vectors** that were selected by a **Byzantine resilient** aggregation rule **A**.
- ▶ **A** can be Brute, Krum, Median, etc.
- ▶ **Bulyan** is a strong Byzantine-resilience algorithm.



The Hidden Vulnerability of Distributed Learning in Byzantium



Bulyan - Step One (1/2)

- ▶ $n \geq 4f + 3$
- ▶ A **two step** process.



Bulyan - Step One (1/2)

- ▶ $n \geq 4f + 3$
- ▶ A **two step** process.
- ▶ The first one is to **recursively** use **A** to select $\theta = n - 2f$ gradients:



Bulyan - Step One (1/2)

- ▶ $n \geq 4f + 3$
- ▶ A **two step** process.
- ▶ The first one is to **recursively** use **A** to select $\theta = n - 2f$ gradients:
 1. With **A**, choose, among the proposed vectors, the closest one to **A**'s output (for Krum this would be the exact output of **A**).



Bulyan - Step One (1/2)

- ▶ $n \geq 4f + 3$
- ▶ A **two step** process.
- ▶ The first one is to **recursively** use **A** to select $\theta = n - 2f$ gradients:
 1. With **A**, choose, among the proposed vectors, the closest one to **A**'s output (for Krum this would be the exact output of **A**).
 2. Remove the chosen gradient from the **received set** and add it to the **selection set S**.



Bulyan - Step One (1/2)

- ▶ $n \geq 4f + 3$
- ▶ A **two step** process.
- ▶ The first one is to **recursively** use **A** to select $\theta = n - 2f$ gradients:
 1. With **A**, choose, among the proposed vectors, the closest one to **A**'s output (for Krum this would be the exact output of **A**).
 2. Remove the chosen gradient from the **received set** and add it to the **selection set S**.
 3. Loop back to step 1 if $|S| < \theta$.



Bulyan - Step One (2/2)

- ▶ $\theta = n - 2f \geq 2f + 3$, thus $S = (S_1, \dots, S_\theta)$ contains a majority of non-Byzantine gradients.



Bulyan - Step One (2/2)

- ▶ $\theta = n - 2f \geq 2f + 3$, thus $S = (S_1, \dots, S_\theta)$ contains a majority of non-Byzantine gradients.
- ▶ For each $i \in [1..d]$, the median of the θ coordinates i of the selected gradients is always bounded by coordinates from non-Byzantine submissions.



Bulyan - Step Two

- ▶ The second step is to generate the resulting gradient $\mathbf{F} = (F[1], \dots, F[d])$.



Bulyan - Step Two

- ▶ The second step is to generate the resulting gradient $\mathbf{F} = (F[1], \dots, F[d])$.
- ▶ $\forall i \in [1..d], F[i] = \frac{1}{\beta} \sum_{\mathbf{x} \in \mathcal{M}[i]} \mathbf{x}[i]$



Bulyan - Step Two

- ▶ The second step is to generate the resulting gradient $\mathbf{F} = (F[1], \dots, F[d])$.
- ▶ $\forall i \in [1..d], F[i] = \frac{1}{\beta} \sum_{\mathbf{x} \in \mathcal{M}[i]} \mathbf{x}[i]$
- ▶ $\beta = \theta - 2f \geq 3$



Bulyan - Step Two

- ▶ The second step is to generate the resulting gradient $\mathbf{F} = (F[1], \dots, F[d])$.
- ▶ $\forall i \in [1..d], F[i] = \frac{1}{\beta} \sum_{\mathbf{x} \in M[i]} \mathbf{x}[i]$
- ▶ $\beta = \theta - 2f \geq 3$
- ▶ $M[i] = \arg \min_{RCS, |R|=\beta} (\sum_{\mathbf{x} \in R} |\mathbf{x}[i] - \text{median}[i]|)$



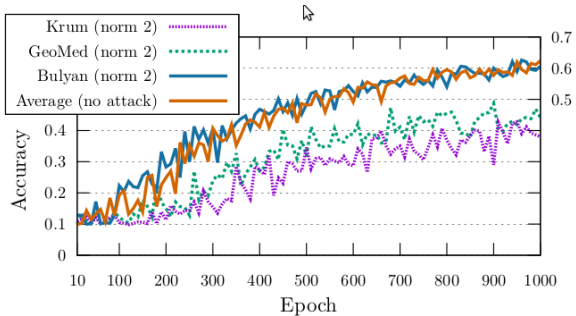
Bulyan - Step Two

- ▶ The second step is to generate the resulting gradient $\mathbf{F} = (F[1], \dots, F[d])$.
- ▶ $\forall i \in [1..d], F[i] = \frac{1}{\beta} \sum_{\mathbf{x} \in M[i]} \mathbf{x}[i]$
- ▶ $\beta = \theta - 2f \geq 3$
- ▶ $M[i] = \arg \min_{R \subseteq S, |R|=\beta} (\sum_{\mathbf{x} \in R} |\mathbf{x}[i] - \text{median}[i]|)$
- ▶ $\text{median}[i] = \arg \min_{m=Y[i], Y \subseteq S} (\sum_{Z \in S} |Z[i] - m|)$



Bulyan - Step Two

- ▶ The second step is to generate the resulting gradient $\mathbf{F} = (F[1], \dots, F[d])$.
- ▶ $\forall i \in [1..d], F[i] = \frac{1}{\beta} \sum_{\mathbf{x} \in M[i]} \mathbf{x}[i]$
- ▶ $\beta = \theta - 2f \geq 3$
- ▶ $M[i] = \arg \min_{R \subseteq S, |R|=\beta} (\sum_{\mathbf{x} \in R} |\mathbf{x}[i] - \text{median}[i]|)$
- ▶ $\text{median}[i] = \arg \min_{m=Y[i], Y \subseteq S} (\sum_{Z \in S} |Z[i] - m|)$
- ▶ Each i th coordinate of \mathbf{F} is equal to the average of the β closest i th coordinates to the median i th coordinate of the θ selected gradients.



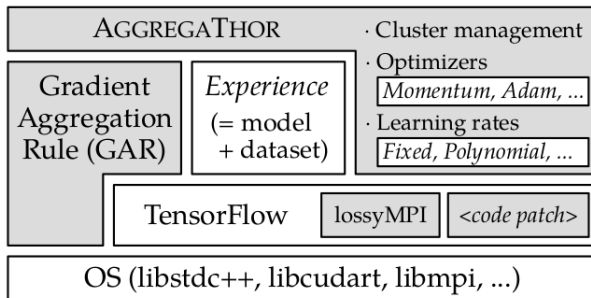
[El Mhamdi et al., The Hidden Vulnerability of Distributed Learning in Byzantium, 2018]



AggregaThor: Byzantine Machine Learning via Robust Gradient Aggregation

AggregaThor (1/2)

- ▶ A framework that handles the **distribution** of the training of a **TensorFlow neural network** graph over a cluster of machines.
- ▶ This distribution is **robust** to **Byzantine** cluster nodes.



[Damaskinos et al., AggregaThor: Byzantine Machine Learning via Robust Gradient Aggregation, 2019]



AggregaThor (2/2)

- ▶ Relies on [Multi-Krum](#) and [Bulyan](#).



AggregaThor (2/2)

- ▶ Relies on Multi-Krum and Bulyan.
- ▶ Multi-Krum selects m gradients that deviate less from the majority
 - Based on their relative distances.



AggregaThor (2/2)

- ▶ Relies on **Multi-Krum** and **Bulyan**.
- ▶ **Multi-Krum** selects m gradients that **deviate less** from the **majority**
 - Based on their **relative distances**.
- ▶ **Bulyan** takes the aforementioned m vectors.



AggregaThor (2/2)

- ▶ Relies on **Multi-Krum** and **Bulyan**.
- ▶ **Multi-Krum** selects m gradients that **deviate less** from the **majority**
 - Based on their **relative distances**.
- ▶ **Bulyan** takes the aforementioned m vectors.
 - Computes their **coordinate-wise median**.



AggregaThor (2/2)

- ▶ Relies on **Multi-Krum** and **Bulyan**.
- ▶ **Multi-Krum** selects m gradients that **deviate less** from the **majority**
 - Based on their **relative distances**.
- ▶ **Bulyan** takes the aforementioned m vectors.
 - Computes their **coordinate-wise median**.
 - Produces a gradient that coordinates are the **average** of the $m - 2f$ closest values to the **median**.



TensorFlow Limitation

- ▶ In TensorFlow, Byzantine resilience **cannot** be achieved **solely** through the use of a **Byzantine-resilient GAR**.



TensorFlow Limitation

- ▶ In TensorFlow, Byzantine resilience **cannot** be achieved **solely** through the use of a **Byzantine-resilient GAR**.
- ▶ TensorFlow allows **any node** in the cluster to execute **arbitrary operations anywhere** in the cluster.



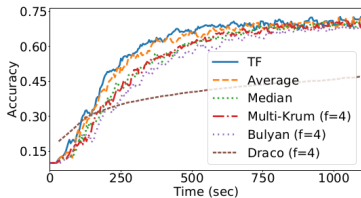
TensorFlow Limitation

- ▶ In TensorFlow, Byzantine resilience **cannot** be achieved **solely** through the use of a **Byzantine-resilient GAR**.
- ▶ TensorFlow allows **any node** in the cluster to execute **arbitrary operations anywhere** in the cluster.
- ▶ A **single** Byzantine worker could **continually overwrite** the **shared parameters** with arbitrary values.

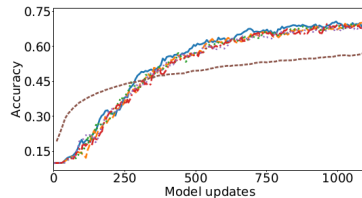


TensorFlow Limitation

- ▶ In TensorFlow, Byzantine resilience **cannot** be achieved **solely** through the use of a **Byzantine-resilient GAR**.
- ▶ TensorFlow allows **any node** in the cluster to execute **arbitrary operations anywhere** in the cluster.
- ▶ A **single** Byzantine worker could **continually overwrite** the **shared parameters** with arbitrary values.
- ▶ AggregaThor **patches** TensorFlow to overcome the above issues.

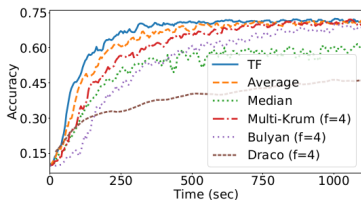


(a)

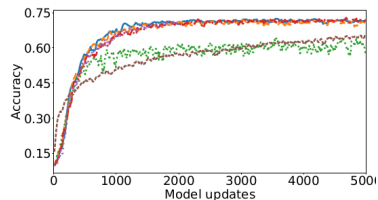


(b)

mini-batch size = 250



(c)



(d)

mini-batch size = 20

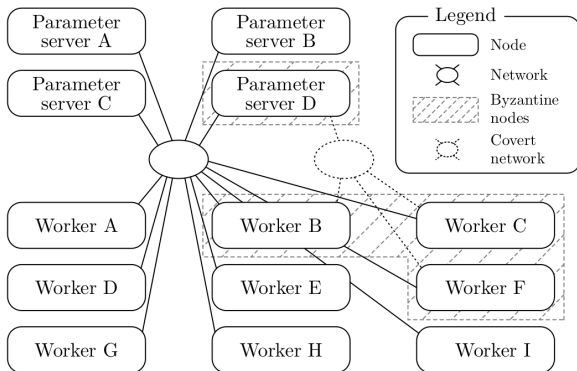
$$n = 19, f = 4 \quad (n \geq 4f + 3)$$

[Damaskinos et al., AggregaThor: Byzantine Machine Learning via Robust Gradient Aggregation, 2019]

What if parameter servers are Byzantine?



SGD: Decentralized Byzantine Resilience



[El Mhamdi et al., SGD: Decentralized Byzantine Resilience, 2019]



- ▶ Byzantine tolerant learning algorithm that is



- ▶ Byzantine tolerant learning algorithm that is
 1. Resilience to Byzantine workers.



- ▶ Byzantine tolerant learning algorithm that is
 1. Resilience to Byzantine workers.
 2. Resilience to Byzantine parameter servers.



- ▶ Byzantine tolerant learning algorithm that is
 1. Resilience to Byzantine workers.
 2. Resilience to Byzantine parameter servers.

- ▶ GuanYu tolerates up to $\frac{1}{3}$ Byzantine servers and $\frac{1}{3}$ Byzantine workers.



GuanYu

- ▶ Byzantine tolerant learning algorithm that is
 1. Resilience to Byzantine workers.
 2. Resilience to Byzantine parameter servers.
- ▶ GuanYu tolerates up to $\frac{1}{3}$ Byzantine servers and $\frac{1}{3}$ Byzantine workers.
- ▶ GuanYu uses a GAR for aggregating workers' gradients and Median for aggregating models received from servers.



Assumptions and Notations (1/2)

- ▶ **Asynchronous network**: the **lack of any bound** on communication delays.



Assumptions and Notations (1/2)

- ▶ **Asynchronous network**: the **lack of any bound** on communication delays.
- ▶ **Synchronous training**: **bulk-synchronous** training.



Assumptions and Notations (1/2)

- ▶ **Asynchronous network**: the **lack of any bound** on communication delays.
- ▶ **Synchronous training**: **bulk-synchronous** training.
 - The parameter server does **not need to wait** for all the workers' gradients to make progress, and vice versa.



Assumptions and Notations (1/2)

- ▶ **Asynchronous network**: the **lack of any bound** on communication delays.
- ▶ **Synchronous training**: **bulk-synchronous** training.
 - The parameter server does **not need to wait** for all the workers' gradients to make progress, and vice versa.
 - The **quorums** indicate the **number of messages to wait** before aggregating them.



Assumptions and Notations (2/2)

- ▶ $n_{ps} \geq 3f_{ps} + 3$ the total number of parameter servers, among which f_{ps} are Byzantine.



Assumptions and Notations (2/2)

- ▶ $n_{ps} \geq 3f_{ps} + 3$ the total number of parameter servers, among which f_{ps} are Byzantine.
- ▶ $n_{wr} \geq 3f_{wr} + 3$ the total number of workers, among which f_{wr} are Byzantine.



Assumptions and Notations (2/2)

- ▶ $n_{ps} \geq 3f_{ps} + 3$ the total number of parameter servers, among which f_{ps} are Byzantine.
- ▶ $n_{wr} \geq 3f_{wr} + 3$ the total number of workers, among which f_{wr} are Byzantine.
- ▶ M the coordinate-wise median (used in both workers and servers).



Assumptions and Notations (2/2)

- ▶ $n_{ps} \geq 3f_{ps} + 3$ the **total number of parameter servers**, among which f_{ps} are Byzantine.
- ▶ $n_{wr} \geq 3f_{wr} + 3$ the **total number of workers**, among which f_{wr} are Byzantine.
- ▶ M the **coordinate-wise median** (used in both workers and servers).
- ▶ F the **GAR** function (used in the servers)



Assumptions and Notations (2/2)

- ▶ $n_{ps} \geq 3f_{ps} + 3$ the **total number of parameter servers**, among which f_{ps} are Byzantine.
- ▶ $n_{wr} \geq 3f_{wr} + 3$ the **total number of workers**, among which f_{wr} are Byzantine.
- ▶ M the **coordinate-wise median** (used in both workers and servers).
- ▶ F the **GAR** function (used in the servers)
- ▶ $2f_{ps} + 3 \leq q_{ps} \leq n_{ps} - f_{ps}$ the **quorum** used for M .



Assumptions and Notations (2/2)

- ▶ $n_{ps} \geq 3f_{ps} + 3$ the **total number of parameter servers**, among which f_{ps} are Byzantine.
- ▶ $n_{wr} \geq 3f_{wr} + 3$ the **total number of workers**, among which f_{wr} are Byzantine.
- ▶ M the **coordinate-wise median** (used in both workers and servers).
- ▶ F the **GAR** function (used in the servers)
- ▶ $2f_{ps} + 3 \leq q_{ps} \leq n_{ps} - f_{ps}$ the **quorum** used for M .
- ▶ $2f_{wr} + 3 \leq q_{wr} \leq n_{wr} - f_{wr}$ the **quorum** used for F .



Assumptions and Notations (2/2)

- ▶ $n_{ps} \geq 3f_{ps} + 3$ the total number of parameter servers, among which f_{ps} are Byzantine.
- ▶ $n_{wr} \geq 3f_{wr} + 3$ the total number of workers, among which f_{wr} are Byzantine.
- ▶ M the coordinate-wise median (used in both workers and servers).
- ▶ F the GAR function (used in the servers)
- ▶ $2f_{ps} + 3 \leq q_{ps} \leq n_{ps} - f_{ps}$ the quorum used for M .
- ▶ $2f_{wr} + 3 \leq q_{wr} \leq n_{wr} - f_{wr}$ the quorum used for F .
- ▶ d the dimension of the parameter space \mathbb{R}^d .



GuanYu Algorithm - Step 1

- ▶ At each step t , each non-Byzantine server i broadcasts its current parameter vector w_i^t to every worker.



GuanYu Algorithm - Step 1

- ▶ At each step t , each non-Byzantine server i broadcasts its current parameter vector \mathbf{w}_i^t to every worker.
- ▶ Each non-Byzantine worker j aggregates with M the q_{ps} first received \mathbf{w}^t .



GuanYu Algorithm - Step 1

- ▶ At each step t , each non-Byzantine server i broadcasts its current parameter vector \mathbf{w}_i^t to every worker.
- ▶ Each non-Byzantine worker j aggregates with M the q_{ps} first received \mathbf{w}^t .
- ▶ And computes an estimate G_j^t of the gradient at the aggregated parameters.



GuanYu Algorithm - Step 2

- ▶ Each non-Byzantine worker j broadcasts its computed gradient estimation G_j^t to every parameter server.



GuanYu Algorithm - Step 2

- ▶ Each non-Byzantine worker j broadcasts its computed gradient estimation G_j^t to every parameter server.
- ▶ Each non-Byzantine parameter server i aggregates with F the q_{wr} first received G^t .



GuanYu Algorithm - Step 2

- ▶ Each **non-Byzantine worker j** broadcasts its **computed gradient estimation G_j^t** to every **parameter server**.
- ▶ Each **non-Byzantine parameter server i** aggregates with **F** the **q_{wr}** first received **G^t** .
- ▶ And performs a **local parameter update** with the aggregated gradient, resulting in **\bar{w}_i^t** .



GuanYu Algorithm - Step 3

- ▶ Each non-Byzantine parameter server i broadcasts $\bar{\mathbf{w}}_i^{t+1}$ to every other parameter servers.



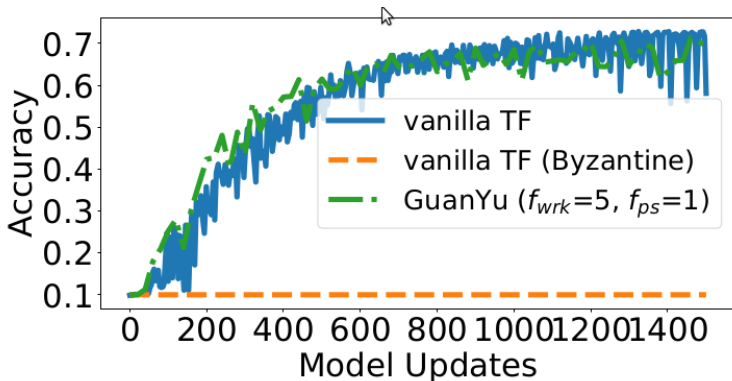
GuanYu Algorithm - Step 3

- ▶ Each non-Byzantine parameter server i broadcasts $\bar{\mathbf{w}}_i^{t+1}$ to every other parameter servers.
- ▶ They aggregate with M the q_{ps} first received $\bar{\mathbf{w}}_k^{t+1}$.



GuanYu Algorithm - Step 3

- ▶ Each non-Byzantine parameter server i broadcasts $\bar{\mathbf{w}}_i^{t+1}$ to every other parameter servers.
- ▶ They aggregate with M the q_{ps} first received $\bar{\mathbf{w}}_k^{t+1}$.
- ▶ This aggregated parameter vector is $\bar{\mathbf{w}}_i^{t+1}$.



[El Mhamdi et al., SGD: Decentralized Byzantine Resilience, 2019]



GuanYu Limitations

- ▶ Network asynchrony assumption is costly:
 - It requires three communication rounds, instead of two in the vanilla case.



GuanYu Limitations

- ▶ Network asynchrony assumption is costly:
 - It requires three communication rounds, instead of two in the vanilla case.
- ▶ It requires a large number of compute nodes and server replicas to work, as one cannot differentiate between a Byzantine machine and a slow one in such network.



Fast Machine Learning with Byzantine Workers and Servers



LiuBei

- ▶ Uses a GAR to aggregate workers' gradients.



LiuBei

- ▶ Uses a GAR to aggregate workers' gradients.
- ▶ Tolerating Byzantine servers using a filtering technique and the scatter/gather protocol (both assume network synchrony).



- ▶ Uses a GAR to aggregate workers' gradients.
- ▶ Tolerating Byzantine servers using a filtering technique and the scatter/gather protocol (both assume network synchrony).
 - Scatter phase: servers work independently and do not communicate with each other.



- ▶ Uses a GAR to aggregate workers' gradients.
- ▶ Tolerating Byzantine servers using a filtering technique and the scatter/gather protocol (both assume network synchrony).
 - Scatter phase: servers work independently and do not communicate with each other.
 - Gather phase: correct servers communicate to bring their view of models back close to each other.



- ▶ Uses a GAR to aggregate workers' gradients.
- ▶ Tolerating Byzantine servers using a filtering technique and the scatter/gather protocol (both assume network synchrony).
 - Scatter phase: servers work independently and do not communicate with each other.
 - Gather phase: correct servers communicate to bring their view of models back close to each other.
 - The number of gather steps is usually very small and hence, their overhead is insignificant.



Assumptions and Notations

- ▶ Network synchrony: an upper bound on communication machines.
- ▶ $n_{ps} \geq 3f_{ps} + 1$ the total number of parameter servers, among which f_{ps} are Byzantine.
- ▶ $n_{wr} \geq 2f_{wr} + 1$ the total number of workers, among which f_{wr} are Byzantine.



LiuBei - Byzantine Workers

- ▶ LiuBei can use **any existing synchronous GAR** that follows the robustness definition (α, f) -Byzantine-Resilience.



LiuBei - Byzantine Servers (1/2)

- ▶ Tolerating Byzantine servers using **robust aggregation** requires **communication with all servers** in **each round**: big communication **overhead**.



LiuBei - Byzantine Servers (1/2)

- ▶ Tolerating Byzantine servers using **robust aggregation** requires **communication with all servers** in **each round**: big communication **overhead**.
- ▶ LiuBei lets **each worker** pull **only one model** from **any of the server** replicas and then checks if the pulled model is **suspicious or not**.



LiuBei - Byzantine Servers (2/2)

- ▶ A *worker* does this check by applying *two filters* on the pulled model: the *Lipschitz filter* and the *models filter*.



LiuBei - Byzantine Servers (2/2)

- ▶ A **worker** does this check by applying **two filters** on the pulled model: the **Lipschitz filter** and the **models filter**.
- ▶ If the **model is suspicious**, the worker **discards it and pulls a new model** from another parameter server.



LiuBei - Byzantine Servers (2/2)

- ▶ A **worker** does this check by applying **two filters** on the pulled model: the **Lipschitz filter** and the **models filter**.
- ▶ If the **model is suspicious**, the worker **discards it and pulls a new model** from another parameter server.
- ▶ The **maximum number** of models that can be pulled by a worker in **one iteration** is $f_{ps} + 1$.



Lipschitz Filter

- ▶ Assume at time t worker j owns a model w_j^t and computes gradient G_j^t based on that model.



Lipschitz Filter

- ▶ Assume at time t worker j owns a model \mathbf{w}_j^t and computes gradient G_j^t based on that model.
- ▶ A correct server i should include G_j^t while updating its model \mathbf{w}_i^t , given network synchrony.



Lipschitz Filter

- ▶ Assume at time t worker j owns a model \mathbf{w}_j^t and computes gradient G_j^t based on that model.
- ▶ A correct server i should include G_j^t while updating its model \mathbf{w}_i^t , given network synchrony.
- ▶ The worker j then does two steps in parallel: \mathbf{w}_i^{t+1}



Lipschitz Filter

- ▶ Assume at time t worker j owns a model \mathbf{w}_j^t and computes gradient G_j^t based on that model.
- ▶ A correct server i should include G_j^t while updating its model \mathbf{w}_i^t , given network synchrony.
- ▶ The worker j then does two steps in parallel: \mathbf{w}_i^{t+1}
 1. Estimates the updated model locally based on its own gradient: $\mathbf{w}_{j(1)}^{t+1}$



Lipschitz Filter

- ▶ Assume at time t worker j owns a model \mathbf{w}_j^t and computes gradient G_j^t based on that model.
- ▶ A correct server i should include G_j^t while updating its model \mathbf{w}_i^t , given network synchrony.
- ▶ The worker j then does two steps in parallel: \mathbf{w}_i^{t+1}
 1. Estimates the updated model locally based on its own gradient: $\mathbf{w}_{j(1)}^{t+1}$
 2. Pulls a model \mathbf{w}_i^{t+1} from a parameter server i .



Lipschitz Filter

- ▶ Assume at time t worker j owns a model \mathbf{w}_j^t and computes gradient G_j^t based on that model.
- ▶ A correct server i should include G_j^t while updating its model \mathbf{w}_i^t , given network synchrony.
- ▶ The worker j then does two steps in parallel: \mathbf{w}_i^{t+1}
 1. Estimates the updated model locally based on its own gradient: $\mathbf{w}_{j(1)}^{t+1}$
 2. Pulls a model \mathbf{w}_i^{t+1} from a parameter server i .
- ▶ If server i is correct, then, the growth of the pulled model \mathbf{w}_i^{t+1} should be close to that of the estimated local model $\mathbf{w}_{j(1)}^{t+1}$.



Model Filter

- ▶ LiuBei uses the **model filter** to **bound the distance** between **models** in any two **successive iterations**.



Model Filter

- ▶ LiuBei uses the **model filter** to **bound the distance** between **models** in any two **successive iterations**.
- ▶ We assume all **correct machines initialize models** with the **same state**.



Model Filter

- ▶ LiuBei uses the **model filter** to **bound the distance** between **models** in any two **successive iterations**.
- ▶ We assume all **correct machines initialize models** with the **same state**.
- ▶ Building upon the **guarantees given by the used GAR**, at iteration t , a worker can **estimate an upper bound** on the **distance between two successive states** of a correct model.



LiuBei Algorithm (1/2)

- ▶ LiuBei operates in two phases: **scatter** and **gather**.

Algorithm 1 Worker Algorithm

```
1: Calculate the value of  $T$  and a value for  $seed$ 
2:  $model \leftarrow init\_model(seed)$ 
3:  $r \leftarrow random\_int(1, n_{ps})$ 
4:  $t \leftarrow 0$ 
5:  $grad \leftarrow model.backprop()$ 
6: repeat
7:    $local\_model \leftarrow apply\_grad(model, grad)$ 
8:   if  $t \bmod T = 0$  then
9:      $models \leftarrow read\_models()$ 
10:     $model \leftarrow MeaMed(models)$ 
11:   else
12:      $i \leftarrow 0$ 
13:     repeat
14:        $new\_model \leftarrow read\_model($   

            $(r + t + i) \bmod n_{ps})$ 
15:        $new\_grad \leftarrow new\_model.backprop()$ 
16:        $i \leftarrow i + 1$ 
17:     until  $pass\_filters(new\_model)$ 
18:      $model \leftarrow new\_model$ 
19:      $grad \leftarrow new\_grad$ 
20:   end if
21:    $t \leftarrow t + 1$ 
22: until  $t > num\_iterations$ 
```



LiuBei Algorithm (1/2)

- ▶ LiuBei operates in two phases: **scatter** and **gather**.
- ▶ One **gather step** is executed every **T** iterations (line 8 to 11).

Algorithm 1 Worker Algorithm

```
1: Calculate the value of  $T$  and a value for seed
2:  $\text{model} \leftarrow \text{init\_model}(\text{seed})$ 
3:  $r \leftarrow \text{random\_int}(1, n_{ps})$ 
4:  $t \leftarrow 0$ 
5:  $\text{grad} \leftarrow \text{model.backprop}()$ 
6: repeat
7:    $\text{local\_model} \leftarrow \text{apply\_grad}(\text{model}, \text{grad})$ 
8:   if  $t \bmod T = 0$  then
9:      $\text{models} \leftarrow \text{read\_models}()$ 
10:     $\text{model} \leftarrow \text{MeaMed}(\text{models})$ 
11:   else
12:      $i \leftarrow 0$ 
13:     repeat
14:        $\text{new\_model} \leftarrow \text{read\_model}(\text{r} + t + i \bmod n_{ps})$ 
15:        $\text{new\_grad} \leftarrow \text{new\_model.backprop}()$ 
16:        $i \leftarrow i + 1$ 
17:     until  $\text{pass\_filters}(\text{new\_model})$ 
18:      $\text{model} \leftarrow \text{new\_model}$ 
19:      $\text{grad} \leftarrow \text{new\_grad}$ 
20:   end if
21:    $t \leftarrow t + 1$ 
22: until  $t > \text{num\_iterations}$ 
```



LiuBei Algorithm (1/2)

- ▶ LiuBei operates in two phases: **scatter** and **gather**.
- ▶ One **gather step** is executed every T iterations (line 8 to 11).
- ▶ We call the whole T iterations a **scatter step**.

Algorithm 1 Worker Algorithm

```
1: Calculate the value of  $T$  and a value for seed
2:  $\text{model} \leftarrow \text{init\_model}(\text{seed})$ 
3:  $r \leftarrow \text{random\_int}(1, n_{ps})$ 
4:  $t \leftarrow 0$ 
5:  $\text{grad} \leftarrow \text{model.backprop}()$ 
6: repeat
7:    $\text{local\_model} \leftarrow \text{apply\_grad}(\text{model}, \text{grad})$ 
8:   if  $t \bmod T = 0$  then
9:      $\text{models} \leftarrow \text{read\_models}()$ 
10:     $\text{model} \leftarrow \text{MeaMed}(\text{models})$ 
11:   else
12:      $i \leftarrow 0$ 
13:     repeat
14:        $\text{new\_model} \leftarrow \text{read\_model}(\text{r} + t + i \bmod n_{ps})$ 
15:        $\text{new\_grad} \leftarrow \text{new\_model.backprop}()$ 
16:        $i \leftarrow i + 1$ 
17:     until  $\text{pass\_filters}(\text{new\_model})$ 
18:      $\text{model} \leftarrow \text{new\_model}$ 
19:      $\text{grad} \leftarrow \text{new\_grad}$ 
20:   end if
21:    $t \leftarrow t + 1$ 
22: until  $t > \text{num\_iterations}$ 
```

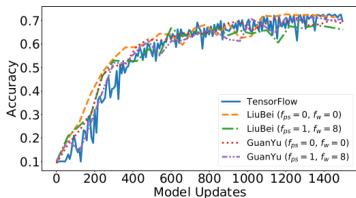


LiuBei Algorithm (2/2)

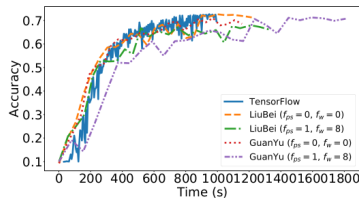
- ▶ LiuBei operates in two phases: **scatter** and **gather**.
- ▶ One **gather step** is executed every T iterations (line 8 to 11).
- ▶ We call the whole T iterations a **scatter step**.

Algorithm 2 Parameter Server Algorithm

```
1: Calculate the value of  $T$  and a value for seed
2:  $model \leftarrow \text{init\_model}(\text{seed})$ 
3:  $t \leftarrow 0$ 
4: repeat
5:    $grads \leftarrow \text{read\_gradients}()$ 
6:    $grad \leftarrow \text{MDA}(grads)$ 
7:    $model.\text{update}(grad)$ 
8:   if  $t \bmod T = 0$  then
9:      $model \leftarrow \text{read\_models}()$ 
10:     $model \leftarrow \text{MeaMed}(models)$ 
11:   end if
12:    $t \leftarrow t + 1$ 
13: until  $t > \text{num\_iterations}$ 
```

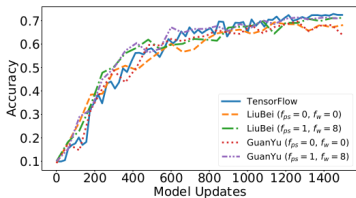


(a)

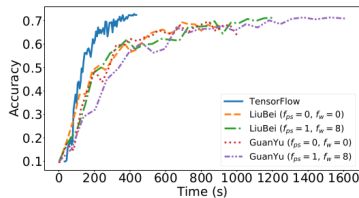


(b)

mini-batch size = 250



(c)



(d)

mini-batch size = 100

[El Mhamdi et al., Fast Machine Learning with Byzantine Workers and Servers, 2019]

Summary



Summary

- ▶ Integrity in data-parallel learning
- ▶ Weak Byzantine resilience
- ▶ Strong Byzantine resilience
- ▶ Byzantine parameter servers



Reference

- ▶ Xie et al., Generalized Byzantine-tolerant SGD, 2018
- ▶ Blanchard et al., Machine Learning with Adversaries: Byzantine Tolerant Gradient Descent, 2017
- ▶ El Mhamdi et al., The Hidden Vulnerability of Distributed Learning in Byzantium, 2018
- ▶ Damaskinos et al., AGGREGATHOR: Byzantine Machine Learning via Robust Gradient Aggregation, 2019
- ▶ El Mhamdi et al., SGD: Decentralized Byzantine Resilience, 2019
- ▶ El Mhamdi et al., Fast Machine Learning with Byzantine Workers and Servers, 2019

Questions?