# Distributed Roubust Learning

Amir H. Payberah
payberah@kth.se
2021-12-15
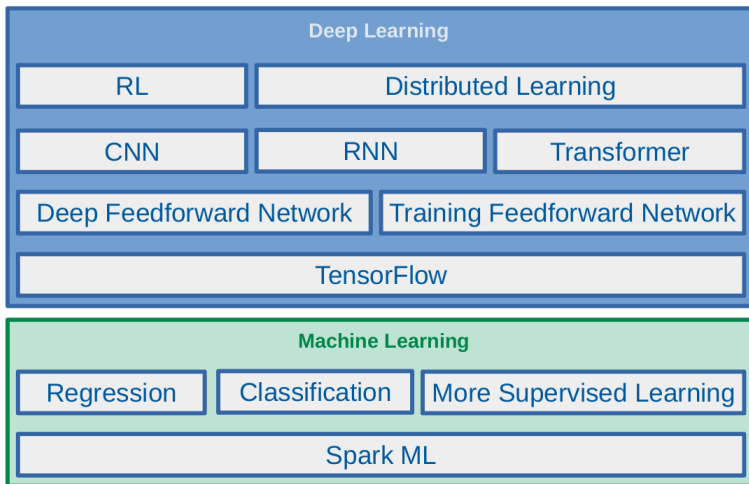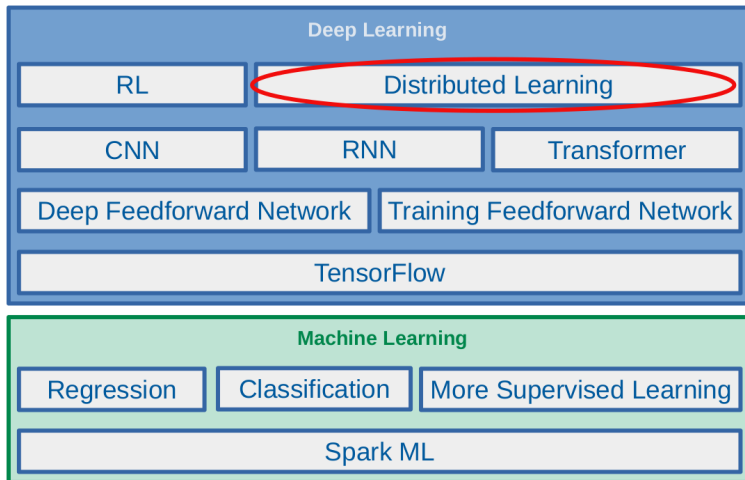
https://id2223kth.github.io

https://tinyurl.com/6s5jy46a

# Where Are We?

- Confidentiality and privacy
  - Confidentiality of the model or the data.

# Adversarial Goals

- ▶ Confidentiality and privacy
  - • Confidentiality of the model or the data.

- ▶ Integrity
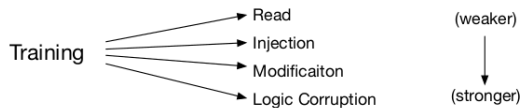  - • Integrity of the predictions

# Adversarial Goals

- ▶ Confidentiality and privacy
  - Confidentiality of the model or the data.

- ▶ Integrity
  - Integrity of the predictions

- ▶ Availability
  - Availability of the system deploying machine learning

▶ Training phase



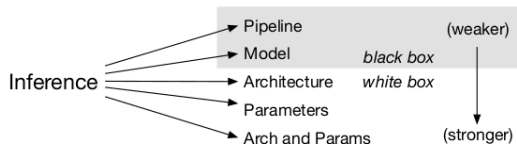[Papernot et al., SoK: Security and Privacy in Machine Learning, 2018]

▶ Training phase



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

▶ Inference phase
  • White box
  • Black box



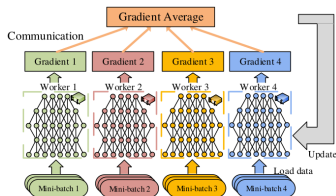[Papernot et al., SoK: Security and Privacy in Machine Learning, 2018]

# Our Focus and Goal
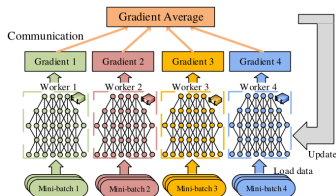
▶ Data parallelization



[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]
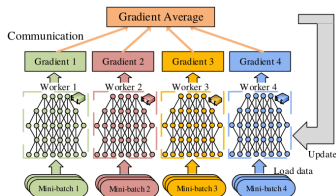
# Our Focus and Goal
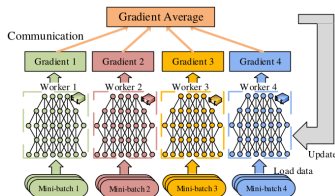
- Data parallelization

- 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]

- One parameter server, and n workers.

- Computation is divided into synchronous rounds.

- During round t, the parameter server broadcasts its parameter vector $w \in \mathbb{R}^d$ to all the workers.



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

▶ At each round `t`, each correct worker `i` computes $G_i(w_t, \beta)$.



[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 $G_i(w_t, \beta)$.

▶ $G_i(w_t, \beta)$: the local estimate of the gradient of the loss function $\nabla J(w_t)$.



[Tang et al., Communication-Efficient Distributed Deep Learning: A Comprehensive Survey, 2020]
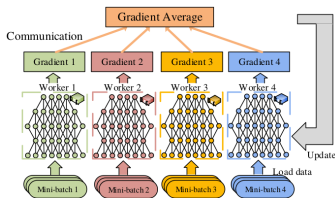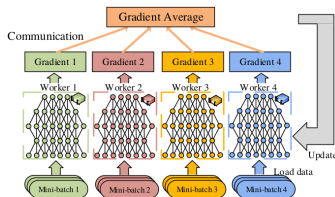
# Distributed Stochastic Gradient Descent (2/3)

- At each round `t`, each correct worker `i` computes $G_i(w_t, \beta)$.

- $G_i(w_t, \beta)$: the local estimate of the gradient of the loss function $\nabla J(w_t)$.

- $\beta$: a mini-batch of i.i.d. samples drawn from the dataset.



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

- At each round `t`, each correct worker `i` computes $G_i(w_t, \beta)$.

- $G_i(w_t, \beta)$: the local estimate of the gradient of the loss function $\nabla J(w_t)$.

- $\beta$: a mini-batch of i.i.d. samples drawn from the dataset.

- $G_i(w_t, \beta) = \frac{1}{|\beta|} \sum_{x \in \beta} \nabla l_i(w_t, x)$



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

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



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

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

▶ Gradient Aggregation Rule (GAR): $F(G_1, G_2, \cdots, G_n) = \frac{1}{n} \sum_{i=1}^{n} G_i$



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

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



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

► 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]

- 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: $F(G_1, G_2, \cdots, G_n) = \frac{1}{n} \sum_{i=1}^{n} G_i$

- $w \leftarrow w - \gamma F(G_1, G_2, \cdots, G_n)$

- Averaging GAR: $\mathtt{F}(\mathtt{G_1}, \mathtt{G_2}, \cdots, \mathtt{G_n}) = \frac{1}{\mathtt{n}} \sum_{\mathtt{i}=1}^{\mathtt{n}} \mathtt{G_i}$

- $\mathtt{w} \leftarrow \mathtt{w} - \gamma \mathtt{F}(\mathtt{G_1}, \mathtt{G_2}, \cdots, \mathtt{G_n})$

- Even a single Byzantine worker can prevent convergence.

- Averaging GAR: $F(G_1, G_2, \cdots, G_n) = \frac{1}{n} \sum_{i=1}^{n} G_i$

- $w \leftarrow w - \gamma F(G_1, G_2, \cdots, G_n)$

- 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 $F = U$.

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

# $(\alpha, f)$-Byzantine-Resilience (1/2)

- Assume `n` workers, where `f` of them are <span style="color:green">Byzantine</span> workers.

- $\alpha \in [0, \pi/2]$ and $f \in \{0, \cdots, n\}$.

# $(\alpha, f)$-Byzantine-Resilience (1/2)

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

- $\alpha \in [0, \pi/2]$ and $\mathtt{f} \in \{0, \cdots, \mathtt{n}\}$.

- $(\mathtt{G_1}, \cdots, \mathtt{G_{n-f}}) \in (\mathbb{R}^\mathtt{d})^\mathtt{n-f}$ are i.i.d. random vectors

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

- $\alpha \in [0, \pi/2]$ and $\texttt{f} \in \{0, \cdots, \texttt{n}\}$.

- $(\texttt{G}_1, \cdots, \texttt{G}_{\texttt{n}-\texttt{f}}) \in (\mathbb{R}^\texttt{d})^{\texttt{n}-\texttt{f}}$ are i.i.d. random vectors
  - $\texttt{G}_\texttt{i} \sim \texttt{g}$
  - $\mathbb{E}[\texttt{g}] = \mathcal{J}$, where $\mathcal{J} = \nabla\texttt{J}(\texttt{w})$

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

- $\alpha \in [0, \pi/2]$ and $\texttt{f} \in \{0, \cdots, \texttt{n}\}$.

- $(\texttt{G}_1, \cdots, \texttt{G}_{\texttt{n-f}}) \in (\mathbb{R}^{\texttt{d}})^{\texttt{n-f}}$ are i.i.d. random vectors
  - $\texttt{G}_{\texttt{i}} \sim \texttt{g}$
  - $\mathbb{E}[\texttt{g}] = \mathcal{J}$, where $\mathcal{J} = \nabla \texttt{J}(\texttt{w})$

- $(\texttt{B}_1, \cdots, \texttt{B}_{\texttt{f}}) \in (\mathbb{R}^{\texttt{d}})^{\texttt{f}}$ are random vectors, possibly dependent between them and the vectors $(\texttt{G}_1, \cdots, \texttt{G}_{\texttt{n-f}})$

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

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

  1. Vector F that is not too far from the real gradient $\mathcal{J}$, i.e., $||\mathbb{E}[F] - \mathcal{J}|| \leq r$.



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

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

1. Vector $F$ that is not too far from the real gradient $\mathcal{J}$, i.e., $||\mathbb{E}[F] - \mathcal{J}|| \leq r$.

2. Moments of $F$ should be controlled by the moments of the (correct) gradient estimator $g$, where $\mathbb{E}[g] = \mathcal{J}$.



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

- Median
- Krum
- Multi-Krum
- Brute

- $n \geq 2f + 1$

# Median

- $n \geq 2f + 1$

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

- $n \geq 2f + 1$

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

- $d$: the gradient vectors dimension.

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

- $n \geq 2f + 1$

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

- $d$: the gradient vectors dimension.

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

- Marginal median
$$F = \texttt{MarMed}(G_1, \cdots, G_n) = \begin{pmatrix} \texttt{median}(G_1[1], \cdots, G_n[1]) \\ \vdots \\ \texttt{median}(G_1[d], \cdots, G_n[d]) \end{pmatrix} \tag{1}$$

- $n \geq 2f + 3$

- $n \geq 2f + 3$

- Idea: to preclude the vectors that are too far away.

- $n \geq 2f + 3$

- 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 $G_i$.

- $F(G_1, \cdots, G_n) = G_{i_*}$

# Krum

- $n \geq 2f + 3$

- Idea: to preclude the vectors that are too far away.

- $s(i) = \sum_{i \to j} ||G_i - G_j||^2$, the score of the worker $i$.

- $i \to j$ denotes that $G_j$ belongs to the $n - f - 2$ closest vectors to $G_i$.

- $F(G_1, \cdots, G_n) = G_{i_*}$

- $G_{i_*}$ 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 $\mathtt{m}$ vectores $G_{1_*}, \cdots, G_{m_*}$, which score the best ($1 \leq \mathtt{m} \leq \mathtt{n} - \mathtt{f} - 2$).

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

- It selects $\mathtt{m}$ vectores $G_{1_*}, \cdots, G_{m_*}$, which score the best ($1 \le \mathtt{m} \le \mathtt{n} - \mathtt{f} - 2$).

- It outputs their average $\frac{1}{\mathtt{m}} \sum_i G_{i_*}$.

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

- ▶ It selects `m` vectors $G_{1_*}, \cdots, G_{m_*}$, which score the best ($1 \le \mathtt{m} \le \mathtt{n} - \mathtt{f} - 2$).

- ▶ It outputs their average $\frac{1}{\mathtt{m}} \sum_{\mathtt{i}} G_{\mathtt{i}_*}$.

- ▶ The cases $\mathtt{m} = 1$ and $\mathtt{m} = \mathtt{n}$ correspond to Krum and averaging, respectively.



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

- $n \geq 2f + 1$

# Brute

- $n \geq 2f + 1$

- $\mathcal{Q} = \{G_1, G_2, \cdots, G_n\}$

- $n \geq 2f + 1$

- $\mathcal{Q} = \{G_1, G_2, \cdots, G_n\}$

- $\mathcal{R} = \{\mathcal{X} | \mathcal{X} \subset \mathcal{Q}, |\mathcal{X}| = n - f\}$
  - The set of all the subsets of $n - f$

- $n \geq 2f + 1$

- $\mathcal{Q} = \{G_1, G_2, \cdots, G_n\}$

- $\mathcal{R} = \{\mathcal{X} | \mathcal{X} \subset \mathcal{Q}, |\mathcal{X}| = n - f\}$
  - 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

- $n \geq 2f + 1$

- $\mathcal{Q} = \{G_1, G_2, \cdots, G_n\}$

- $\mathcal{R} = \{\mathcal{X} | \mathcal{X} \subset \mathcal{Q}, |\mathcal{X}| = n - f\}$
  - 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.

- $F(G_1, \cdots, G_n) = \frac{1}{n-f} \sum_{G \in \mathcal{S}} G$

MNIST          CIFAR-10

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

- 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:

- 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?

- 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?

▶ 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

- $n \geq 4f + 3$

- A two step process.

- $n \geq 4f + 3$

- A two step process.

- The first one is to recursively use A to select $\theta = n - 2f$ gradients:

- $n \geq 4f + 3$

- 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).

- $n \geq 4f + 3$

- 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 $|S| < \theta$.

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

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

- For each $i \in [1..d]$, the median of the $\theta$ coordinates $i$ of the selected gradients is always bounded by coordinates from non-Byzantine submissions.

▶ The second step is to generate the resulting gradient $F = (F[1], \cdots, F[d])$.

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

▶ $\forall \mathtt{i} \in [1..\mathtt{d}], \mathtt{F[i]} = \frac{1}{\beta} \sum_{\mathtt{X} \in \mathtt{M[i]}} \mathtt{X[i]}$

# Bulyan - Step Two

- The second step is to generate the resulting gradient $F = (F[1], \cdots, F[d])$.

- $\forall i \in [1..d], F[i] = \frac{1}{\beta} \sum_{X \in M[i]} X[i]$

- $\beta = \theta - 2f \geq 3$

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

- $\forall \mathtt{i} \in [1..\mathtt{d}], \mathtt{F[i]} = \frac{1}{\beta} \sum_{\mathtt{X} \in \mathtt{M[i]}} \mathtt{X[i]}$

- $\beta = \theta - 2\mathtt{f} \geq 3$

- $\mathtt{M[i]} = \arg\min_{\mathtt{R} \subset \mathtt{S}, |\mathtt{R}| = \beta} (\sum_{\mathtt{X} \in \mathtt{R}} |\mathtt{X[i]} - \mathtt{median[i]}|)$

# Bulyan - Step Two

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

- $\forall \mathtt{i} \in [1..\mathtt{d}], \mathtt{F[i]} = \frac{1}{\beta} \sum_{\mathtt{X} \in \mathtt{M[i]}} \mathtt{X[i]}$

- $\beta = \theta - 2\mathtt{f} \geq 3$

- $\mathtt{M[i]} = \arg\min_{\mathtt{R} \subset \mathtt{S}, |\mathtt{R}| = \beta} (\sum_{\mathtt{X} \in \mathtt{R}} |\mathtt{X[i]} - \mathtt{median[i]}|)$

- $\mathtt{median[i]} = \arg\min_{\mathtt{m} = \mathtt{Y[i]}, \mathtt{Y} \in \mathtt{S}} (\sum_{\mathtt{Z} \in \mathtt{S}} |\mathtt{Z[i]} - \mathtt{m}|)$
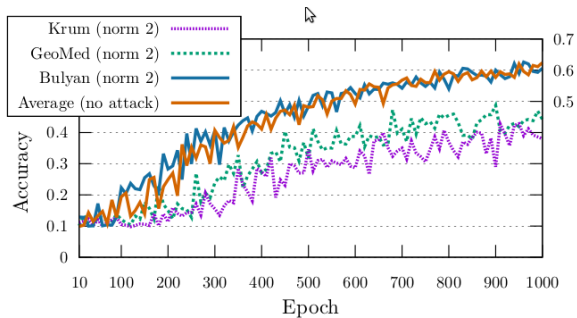
# Bulyan - Step Two

- The second step is to generate the resulting gradient $F = (F[1], \cdots, F[d])$.

- $\forall i \in [1..d], F[i] = \frac{1}{\beta} \sum_{X \in M[i]} X[i]$

- $\beta = \theta - 2f \geq 3$

- $M[i] = \arg \min_{R \subset S, |R| = \beta} (\sum_{X \in R} |X[i] - \text{median}[i]|)$

- $\text{median}[i] = \arg \min_{m = Y[i], Y \in S} (\sum_{Z \in S} |Z[i] - m|)$

- Each $i$th coordinate of $F$ is equal to the average of the $\beta$ closest $i$th coordinates to the median $i$th coordinate of the $\theta$ selected gradients.

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

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.

# 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.

- 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.

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

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

- Synchronous training: bulk-synchronous training.

▶ 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.

▶ 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.

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

- $n_{ps} \geq 3f_{ps} + 3$ 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.

- $n_{ps} \geq 3f_{ps} + 3$ 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).

- $n_{ps} \geq 3f_{ps}+3$ 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)

- $n_{ps} \geq 3f_{ps}+3$ 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.

- $n_{ps} \geq 3f_{ps} + 3$ 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.

- $n_{ps} \geq 3f_{ps} + 3$ 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 $\mathbb{R}^d$.

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

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

- Each non-Byzantine worker `j` aggregates with `M` the `q_ps` first received $w^t$.

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

- Each non-Byzantine worker `j` aggregates with `M` the $q_{ps}$ first received $w^t$.

- And computes an estimate $G_j^t$ of the gradient at the aggregated parameters.

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

- 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 $G^t$.

▶ 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 $\overline{w}_i^t$.

► Each non-Byzantine parameter server `i` broadcasts $\overline{w}_i^{t+1}$ to every other parameter servers.

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

▶ They aggregate with `M` the `q_ps` first received $\overline{w}_k^{t+1}$.

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

▶ They aggregate with `M` the $q_{ps}$ first received $\overline{w}_k^{t+1}$.

▶ This aggregated parameter vector is $\overline{\overline{w}}_i^{t+1}$.

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