# PowerSGD:
# Practical Low-Rank Gradient Compression for Distributed Optimization

Thijs Vogels, Sai Praneeth Karimireddy, Martin Jaggi

# Why distributed SGD?

## Why Distributed?

- Data Privacy

- Resource Distribution

- Data Distribution

## Why dSGD?

- Gradients are Linear

- Gradients are Optimal

- Are gradients the smallest statistic?

$$\sum_{s=1}^{S} \nabla_{\mathbf{W}} \mathcal{L}^{(s)} = \nabla_{\mathbf{W}} \mathcal{L}$$
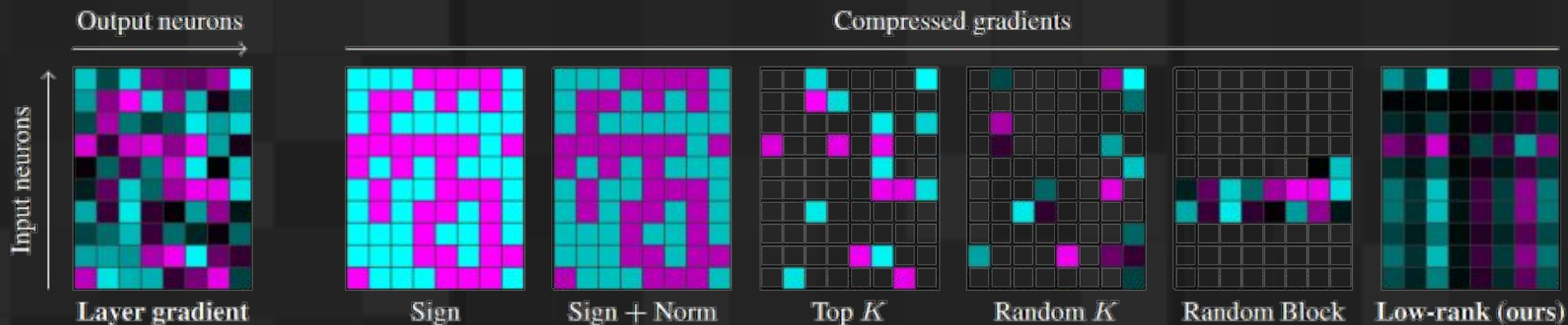
# Reducing Bandwidth



Figure 1: Compression schemes compared in this paper. Left: Interpretation of a layer's gradient as a matrix. Coordinate values are color coded (positive, negative). Right: The output of various compression schemes on the same input. Implementation details are in Appendix G.

# Low-Rank Decompositions



$$C_k = U \quad \Sigma_k \quad V^T$$

▶ **Figure 18.2** Illustration of low rank approximation using the singular-value decomposition. The dashed boxes indicate the matrix entries affected by "zeroing out" the smallest singular values.

# Power Iterations

Algorithm

- initial approximation - random unit vector $x_0$
- $x_1 = A x_0$
- $x_2 = A A x_0 = A^2 x_0$
- $x_3 = A A A x_0 = A^3 x_0$
- ...
- until converges

For large powers of $k$, we will obrain a good approximation of the dominant eigenvector

# Power Iterations (contd.)

We can just remove the dominant direction from the matrix and repeat

So:

- $A = Q\Lambda Q^T$, so $A = \sum_{i=1}^{n} \lambda_i \mathbf{q}_i \mathbf{q}_i^T$
- use power iteration to find $\mathbf{q}_1$ and $\lambda_1$
- then let $A_2 \leftarrow A - \lambda_1 \mathbf{q}_1 \mathbf{q}_1^T$
- repeat power iteration on $A_2$ to find $\mathbf{q}_2$ and $\lambda_2$
- continue like this for $\lambda_3, \ldots, \lambda_n$

## QR Decomposition

Algorithm:

- choose $Q_0$ such that $Q_0^T Q_0 = I$
- for $k = 1, 2, \ldots$:
  - $Z_k = A Q_{k-1}$
  - $Q_k R_k = Z_k$ (QR decomposition)

# PowerSGD

**function** COMPRESS+AGGREGATE(update matrix $M \in \mathbb{R}^{n \times m}$, previous $Q \in \mathbb{R}^{m \times r}$)

$\quad P \leftarrow MQ$

$\quad P \leftarrow$ ALL REDUCE MEAN$(P)$ $\qquad\qquad\qquad\qquad\qquad$ ▷ Now, $P = \frac{1}{W}(M_1 + \ldots + M_W)Q$

$\quad \hat{P} \leftarrow$ ORTHOGONALIZE$(P)$ $\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ ▷ Orthonormal columns

$\quad Q \leftarrow M^\top \hat{P}$

$\quad Q \leftarrow$ ALL REDUCE MEAN$(Q)$ $\qquad\qquad\qquad\qquad\qquad$ ▷ Now, $Q = \frac{1}{W}(M_1 + \ldots + M_W)^\top \hat{P}$

**return** the compressed representation $(\hat{P}, Q)$.



(a) Gather $\qquad\qquad\qquad\qquad$ (b) Reduce

# PowerSGD + Error Feedback

**Algorithm 2** Distributed Error-feedback SGD with Momentum

1: **hyperparameters:** learning rate $\gamma$, momentum parameter $\lambda$
2: **initialize** model parameters $\mathbf{x} \in \mathbb{R}^d$, momentum $\mathbf{m} \leftarrow \mathbf{0} \in \mathbb{R}^d$, replicated across workers
3: **at** each worker $w = 1, \ldots, W$ **do**
4:      **initialize** memory $\mathbf{e}_w \leftarrow \mathbf{0} \in \mathbb{R}^d$
5:      **for** each iterate $t = 0, \ldots$ **do**
6:          Compute a stochastic gradient $\mathbf{g}_w \in \mathbb{R}^d$.
7:          $\Delta_w \quad \leftarrow \mathbf{g}_w + \mathbf{e}_w$            $\triangleright$ Incorporate error-feedback into update
8:          $\mathcal{C}(\Delta_w) \leftarrow \text{COMPRESS}(\Delta_w)$
9:          $\mathbf{e}_w \quad \leftarrow \Delta_w - \text{DECOMPRESS}(\mathcal{C}(\Delta_w))$       $\triangleright$ Memorize local errors
10:        $\mathcal{C}(\Delta) \quad \leftarrow \text{AGGREGATE}(\mathcal{C}(\Delta_1), \ldots, \mathcal{C}(\Delta_W))$     $\triangleright$ Exchange gradients
11:        $\Delta' \quad \leftarrow \text{DECOMPRESS}(\mathcal{C}(\Delta))$         $\triangleright$ Reconstruct an update $\in \mathbb{R}^d$
12:        $\mathbf{m} \quad \leftarrow \lambda \mathbf{m} + \Delta'$
13:        $\mathbf{x} \quad \leftarrow \mathbf{x} - \gamma (\Delta' + \mathbf{m})$
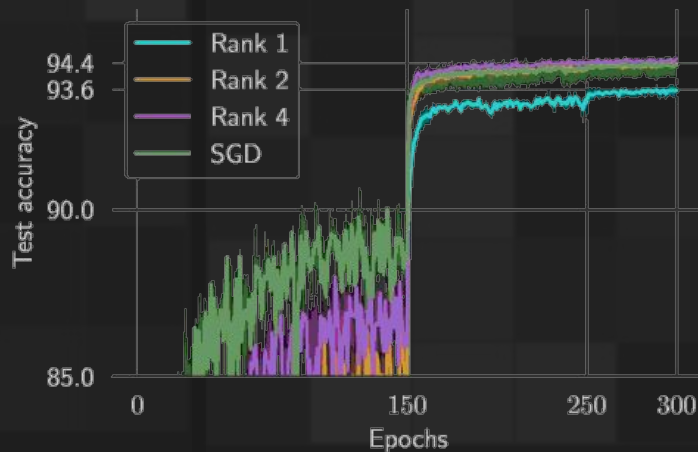14:      **end for**
15: **end at**

Results

# Our Approach

- Use AD statistics instead of gradients for automatic bandwidth reduction
- Further use structured power iterations