

# Quantization in Federated Learning

Tao Shen

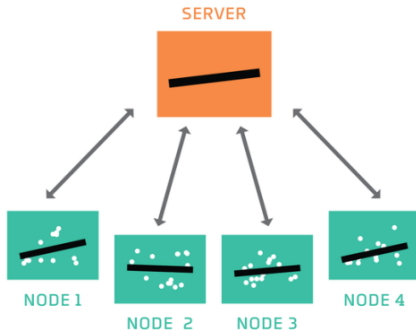
Zhejiang University

November 23, 2019

# FedPAQ: A Communication-Efficient Federated Learning Method with Periodic Averaging and Quantization

# Federated Learning

## Federated Learning



## Communication bottleneck

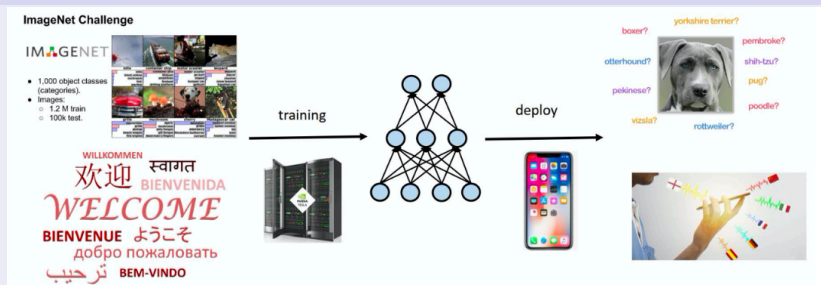
Communication-efficient? Large number of parameters

# Model Compress

## Small model

better generalization, faster prediction, smaller memory footprint

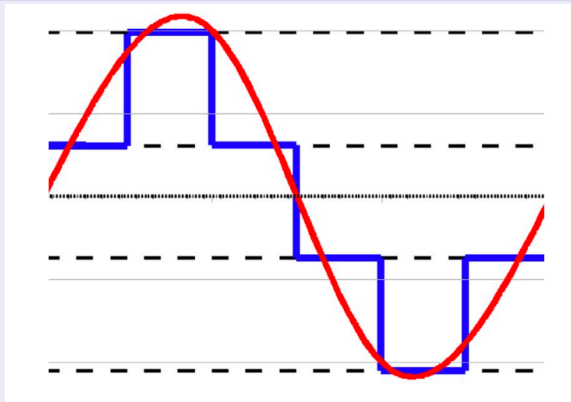
## Computation and memory intensive



capacity of deep network is usually larger than needed

# Quantization

## Quantization



## Weight Quantization

32-bit → fewer bits

# FedPAQ

Periodic Update

...

Partial participation

...

Weight quantization

!!! how?

# Quantizer

## Quantizer in FedPAQ

**Example 1** (Low-precision quantizer [Alistarh et al., 2017]). For any variable  $\mathbf{x} = [x_1, \dots, x_p]^\top$ , the low precision quantizer  $Q^{LP} : \mathbb{R}^p \rightarrow \mathbb{R}^p$  is defined as below

$$Q^{LP}(\mathbf{x}) = [Q_1^{LP}(\mathbf{x}), \dots, Q_p^{LP}(\mathbf{x})]^\top, \quad \text{and} \quad Q_i^{LP}(\mathbf{x}) = \|\mathbf{x}\| \cdot \text{sign}(x_i) \cdot \xi_i(\mathbf{x}, s), \quad (5)$$

where  $\xi_i(\mathbf{x}, s)$  is a random variable defined as

$$\xi_i(\mathbf{x}, s) = \begin{cases} \frac{l}{s} & \text{w.p. } 1 - \left( \frac{|x_i|}{\|\mathbf{x}\|} s - l \right), \\ \frac{l+1}{s} & \text{w.p. } \frac{|x_i|}{\|\mathbf{x}\|} s - l. \end{cases} \quad (6)$$

In above, the tuning parameter  $s$  corresponds to the number of quantization levels and  $l \in [0, s)$  is an integer such that  $|x_i|/\|\mathbf{x}\| \in [l/s, (l+1)/s)$ .

## Theoretical guarantee

---

## Algorithm 1 FedPAQ

---

**Require:** number of periods  $K$ , period length  $\tau$ , stepsize  $\eta_{k,t}$

```

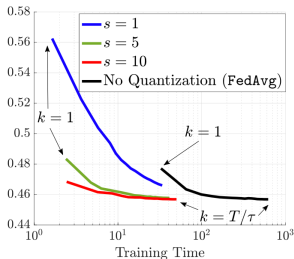
1: for  $k = 0, 1, \dots, K - 1$  do
2:   parameter server picks  $r$  nodes uniformly at random denoted by  $\mathcal{S}_k \subseteq [n]$ 
3:   parameter server sends  $\mathbf{x}_k$  to all the nodes in  $\mathcal{S}_k$ 
4:   for node  $i \in \mathcal{S}_k$  do
5:      $\mathbf{x}_{k,0}^{(i)} \leftarrow \mathbf{x}_k$ 
6:     for  $t = 0, 1, \dots, \tau - 1$  do
7:       randomly pick a data point  $\xi \in \mathcal{D}^i$  and
8:       compute the stochastic gradient  $\tilde{\nabla} f_i(\mathbf{x}) = \nabla \ell(\mathbf{x}, \xi)$ 
9:        $\mathbf{x}_{k,t+1}^{(i)} \leftarrow \mathbf{x}_{k,t}^{(i)} - \eta_{k,t} \tilde{\nabla} f_i(\mathbf{x}_{k,t}^{(i)})$ 
10:    end for
11:    send  $Q(\mathbf{x}_{k,\tau}^{(i)} - \mathbf{x}_k)$  to the parameter server
12:  end for
13:  parameter server updates  $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + \frac{1}{r} \sum_{i \in \mathcal{S}_k} Q(\mathbf{x}_{k,\tau}^{(i)} - \mathbf{x}_k)$ 
14: end for

```

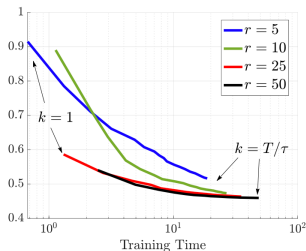
---



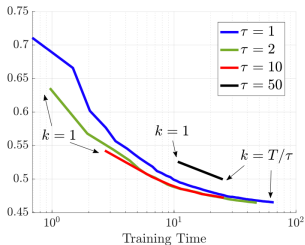
# FedPAQ



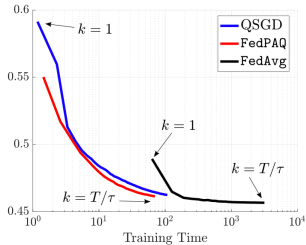
(a) varying the quantization level  $s$



(b) varying the number of active nodes  $r$



(c) varying the period length  $\tau$



(d) comparing FedPAQ, FedAvg, and QSGD

# Loss-aware Binarization of Deep Networks

# Loss-aware Binarization

## Loss-aware

explicitly consider the loss during quantization

## Binarization

$\{+1, -1\}$

# Loss-aware Binarization

## Loss-aware Binarization

$$\begin{aligned} \min_{\hat{\mathbf{w}}} \quad & \ell(\hat{\mathbf{w}}) \\ \text{s.t.} \quad & \hat{\mathbf{w}}_l = \alpha_l \mathbf{b}_l, \alpha_l > 0, \mathbf{b}_l \in \{\pm 1\}^{n_l}, l = 1, \dots, L \end{aligned}$$

## proximal Newton algorithm

$$\ell(\hat{\mathbf{w}}^{t-1}) + \nabla \ell(\hat{\mathbf{w}}^{t-1})^\top (\hat{\mathbf{w}}^t - \hat{\mathbf{w}}^{t-1}) + \frac{1}{2} (\hat{\mathbf{w}}^t - \hat{\mathbf{w}}^{t-1})^\top \mathbf{H}^{t-1} (\hat{\mathbf{w}}^t - \hat{\mathbf{w}}^{t-1})$$

## Approximate diagonal Hessian

$$\begin{aligned} \min_{\hat{\mathbf{w}}^t} \quad & \nabla \ell(\hat{\mathbf{w}}^{t-1})^\top (\hat{\mathbf{w}}^t - \hat{\mathbf{w}}^{t-1}) + \frac{1}{2} (\hat{\mathbf{w}}^t - \hat{\mathbf{w}}^{t-1})^\top \mathbf{D}^{t-1} (\hat{\mathbf{w}}^t - \hat{\mathbf{w}}^{t-1}) \\ \text{s.t.} \quad & \hat{\mathbf{w}}_l^t = \alpha_l^t \mathbf{b}_l^t, \alpha_l^t > 0, \mathbf{b}_l^t \in \{\pm 1\}^{n_l}, \quad l = 1, \dots, L \end{aligned}$$

# Loss-aware Binarization

## Closed-form optimal solution

**Proposition 3.1** Let  $\mathbf{d}_l^{t-1} \equiv \text{diag}(\mathbf{D}_l^{t-1})$ , and

$$\mathbf{w}_l^t \equiv \hat{\mathbf{w}}_l^{t-1} - \nabla_l \ell(\hat{\mathbf{w}}^{t-1}) \oslash \mathbf{d}_l^{t-1}. \quad (7)$$

The optimal solution of (6) can be obtained in closed-form as

$$\alpha_l^t = \frac{\|\mathbf{d}_l^{t-1} \odot \mathbf{w}_l^t\|_1}{\|\mathbf{d}_l^{t-1}\|_1}, \quad \mathbf{b}_l^t = \text{sign}(\mathbf{w}_l^t). \quad (8)$$

loss-aware quantized weights can be easily computed

# Loss-aware Binarization

---

**Algorithm 1** Loss-Aware Binarization (LAB) for training a feedforward neural network.

---

**Input:** Minibatch  $\{(\mathbf{x}_0^t, \mathbf{y}^t)\}$ , current full-precision weights  $\{\mathbf{w}_l^t\}$ , first moment  $\{\mathbf{m}_l^{t-1}\}$ , second moment  $\{\mathbf{v}_l^{t-1}\}$ , and learning rate  $\eta^t$ .

1: **Forward Propagation**

2: **for**  $l = 1$  to  $L$  **do**

3:  $\alpha_l^t = \frac{\|\mathbf{d}_l^{t-1} \odot \mathbf{w}_l^t\|_1}{\|\mathbf{d}_l^{t-1}\|_1};$

4:  $\mathbf{b}_l^t = \text{sign}(\mathbf{w}_l^t);$

5: rescale the layer- $l$  input:  $\tilde{\mathbf{x}}_{l-1}^t = \alpha_l^t \mathbf{x}_{l-1}^t;$

6: compute  $\mathbf{z}_l^t$  with input  $\tilde{\mathbf{x}}_{l-1}^t$  and binary weight  $\mathbf{b}_l^t;$

7: apply batch-normalization and nonlinear activation to  $\mathbf{z}_l^t$  to obtain  $\mathbf{x}_l^t;$

8: **end for**

9: compute the loss  $\ell$  using  $\mathbf{x}_L^t$  and  $\mathbf{y}^t;$

10: **Backward Propagation**

11: initialize output layer's activation's gradient  $\frac{\partial \ell}{\partial \mathbf{x}_L^t};$

12: **for**  $l = L$  to 2 **do**

13: compute  $\frac{\partial \ell}{\partial \mathbf{x}_{l-1}^t}$  using  $\frac{\partial \ell}{\partial \mathbf{x}_l^t}, \alpha_l^t$  and  $\mathbf{b}_l^t;$

14: **end for**

15: **Update parameters using Adam**

16: **for**  $l = 1$  to  $L$  **do**

17: compute gradients  $\nabla_l \ell(\hat{\mathbf{w}}^t)$  using  $\frac{\partial \ell}{\partial \mathbf{x}_l^t}$  and  $\mathbf{x}_{l-1}^t;$

18: update first moment  $\mathbf{m}_l^t = \beta_1 \mathbf{m}_l^{t-1} + (1 - \beta_1) \nabla_l \ell(\hat{\mathbf{w}}^t);$

19: update second moment  $\mathbf{v}_l^t = \beta_2 \mathbf{v}_l^{t-1} + (1 - \beta_2) (\nabla_l \ell(\hat{\mathbf{w}}^t) \odot \nabla_l \ell(\hat{\mathbf{w}}^t));$

20: compute unbiased first moment  $\hat{\mathbf{m}}_l^t = \mathbf{m}_l^t / (1 - \beta_1^t);$

21: compute unbiased second moment  $\hat{\mathbf{v}}_l^t = \mathbf{v}_l^t / (1 - \beta_2^t);$

22: compute current curvature matrix  $\mathbf{d}_l^t = \frac{1}{\eta^t} (\epsilon \mathbf{1} + \sqrt{\hat{\mathbf{v}}_l^t});$

23: update full-precision weights  $\mathbf{w}_l^{t+1} = \mathbf{w}_l^t - \hat{\mathbf{m}}_l^t \odot \mathbf{d}_l^t;$

24: update learning rate  $\eta^{t+1} = \text{UpdateRule}(\eta^t, t + 1);$

25: **end for**

---

# Loss-aware Binarization

Table 1: Test error rates (%) for feedforward neural network models.

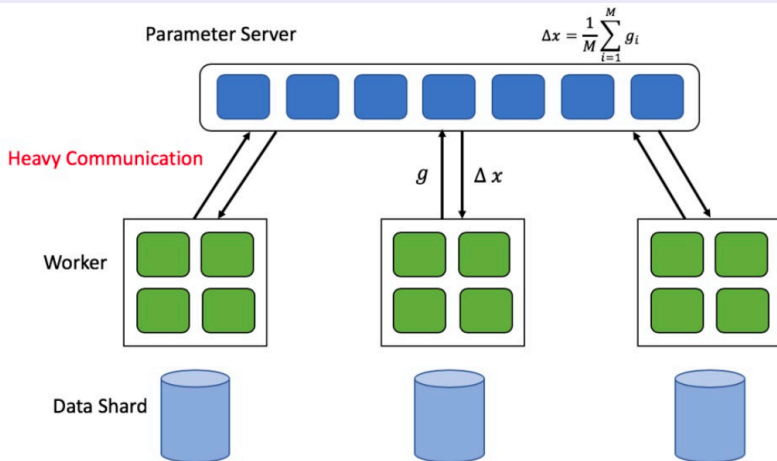
		<i>MNIST</i>	<i>CIFAR-10</i>	<i>SVHN</i>
(no binarization)	full-precision	1.190	11.900	2.277
(binarize weights)	BinaryConnect	1.280	<b>9.860</b>	2.450
	BWN	1.310	10.510	2.535
	LAB	<b>1.180</b>	10.500	<b>2.354</b>
(binarize weights and activations)	BNN	1.470	12.870	3.500
	XNOR	1.530	12.620	3.435
	LAB2	<b>1.380</b>	<b>12.280</b>	<b>3.362</b>

# Analysis of Quantized Models



# Gradient Quantization?

## Distributed Machine Learning



quantize gradients to m bits

# Weight or Gradient?

## Weight or Gradient?

1. full-precision gradients and quantized weights (li et al., 2017, de et al., 2018);
2. full-precision weights and quantized gradients (alistarh et al., 2017, wen et al., 2017, bernstein et al., 2018);

## Both Weight & Gradient

Awesome!

# PRELIMINARIES

## Loss-aware Weight Quantization

$$\begin{aligned} \min_{\hat{\mathbf{w}}} \quad & \|\mathbf{w}_{t+1} - \hat{\mathbf{w}}\|_{\text{Diag}(\sqrt{\gamma_t})}^2 \\ \text{s.t.} \quad & \hat{\mathbf{w}} = \alpha \mathbf{b}, \alpha > 0, \mathbf{b} \in (\mathcal{S}_w)^d \end{aligned}$$

## Stochastic Gradient Quantization

**Example 1** (Low-precision quantizer [Alistarh et al., 2017]). For any variable  $\mathbf{x} = [x_1, \dots, x_p]^\top$ , the low precision quantizer  $Q^{LP} : \mathbb{R}^p \rightarrow \mathbb{R}^p$  is defined as below

$$Q^{LP}(\mathbf{x}) = [Q_1^{LP}(\mathbf{x}), \dots, Q_p^{LP}(\mathbf{x})]^\top, \quad \text{and} \quad Q_i^{LP}(\mathbf{x}) = \|\mathbf{x}\| \cdot \text{sign}(x_i) \cdot \xi_i(\mathbf{x}, s), \quad (5)$$

where  $\xi_i(\mathbf{x}, s)$  is a random variable defined as

$$\xi_i(\mathbf{x}, s) = \begin{cases} \frac{l}{s} & \text{w.p. } 1 - \left(\frac{|x_i|}{\|\mathbf{x}\|} s - l\right), \\ \frac{l+1}{s} & \text{w.p. } \frac{|x_i|}{\|\mathbf{x}\|} s - l. \end{cases} \quad (6)$$

In above, the tuning parameter  $s$  corresponds to the number of quantization levels and  $l \in [0, s)$  is an integer such that  $|x_i|/\|\mathbf{x}\| \in [l/s, (l+1)/s)$ .

# Weight Quantization with **full-precision** gradients

Update

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \text{Diag}(\sqrt{\hat{\mathbf{v}}_t})^{-1} \hat{\mathbf{g}}_t$$

Speed

$$O(1/\sqrt{T})$$

Error

$$LD\sqrt{D^2 + \frac{d\alpha^2\Delta_w^2}{4}}$$

# Weight Quantization with quantized gradients

## Update

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \text{Diag}(\sqrt{\mathbf{v}_t})^{-1} \tilde{\mathbf{g}}_t$$

## Speed

slow down by a factor

$$\sqrt{\frac{1 + \sqrt{2d - 1}}{2}} \Delta_g + 1$$

## Error

No change

## Problems

1. deep networks typically have a large  $d$ ;
2. distributed learning prefers a large  $\Delta_g$

# Weight Quantization with **clipped quantized** gradients

## Update

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \text{Diag}(\sqrt{\mathbf{v}_t})^{-1} \check{\mathbf{g}}_t$$

## Speed

slow down by a factor

$$\sqrt{(2/\pi)^{\frac{1}{2}} c \Delta_g + 1}$$

## Error

introduce extra error

$$\sqrt{d} D\sigma (2/\pi)^{\frac{1}{4}} \sqrt{F(c)}$$

# Weight Quantization with **clipped quantized** gradients

## Speed & Error

slow down by a factor

$$\sqrt{(2/\pi)^{\frac{1}{2}} c \Delta_g + 1}$$

introduce extra error

$$\sqrt{d} D\sigma (2/\pi)^{\frac{1}{4}} \sqrt{F(c)}$$

A larger  $c$  leads to

smaller  $F(c)$ , and thus smaller error

slower convergence

# Test

weight gradient	FP	LAB	LAQ2	LAQ3	LAQ4
FP	83.74	80.37	82.11	83.14	83.35
SQ2 (no clipping)	81.40	78.67	80.27	81.27	81.38
SQ2 (clip, $c = 3$ )	82.99	80.25	81.59	83.14	83.40
SQ3 (no clipping)	83.24	80.18	81.63	82.75	83.17
SQ3 (clip, $c = 3$ )	83.89	80.13	81.77	82.97	83.43
SQ4 (no clipping)	83.64	80.44	81.88	83.13	83.47
SQ4 (clip, $c = 3$ )	83.80	79.27	81.42	82.77	83.43



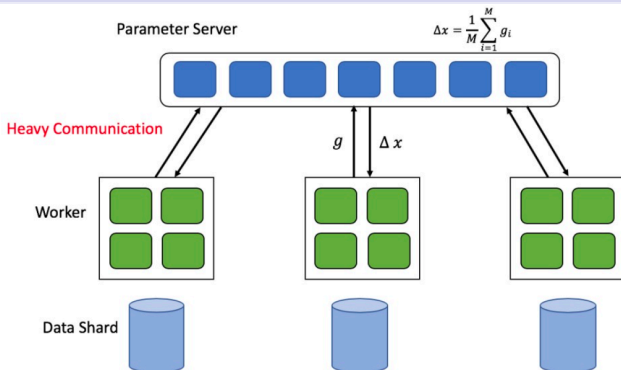
# Main Findings

- ▶ An error related to the weight quantization resolution  $\Delta_w$  and dimension  $d$ .
- ▶ Slow convergence by a factor related to gradient quantization resolution  $\Delta_g$  and  $d$ .
- ▶ Gradient clipping renders the speed degradation mentioned above dimension-free.
- ▶ Distributed training of weight-quantized networks is much faster, while comparable accuracy with the use of full-precision gradients is maintained.

# Communication-Efficient Distributed Blockwise Momentum SGD with Error-Feedback

# Gradient Quantization

## Distributed Machine Learning



two-way compression

same convergence rates as full-precision distributed SGD

# Error-feedback

## Error-feedback

---

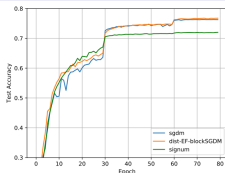
**Algorithm 1** SGD with Error-Feedback (EF-SGD) [11]

---

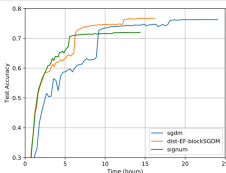
- 1: **Input:** stepsize  $\eta$ ; compressor  $\mathcal{C}(\cdot)$ .
  - 2: **Initialize:**  $x_0 \in \mathbb{R}^d$ ;  $e_0 = 0 \in \mathbb{R}^d$
  - 3: **for**  $t = 0, \dots, T - 1$  **do**
  - 4:    $p_t = \eta g_t + e_t$  {stochastic gradient  $g_t = \nabla f(x_t, \xi_t)$ }
  - 5:    $\Delta_t = \mathcal{C}(p_t)$  {compressed value output}
  - 6:    $x_{t+1} = x_t - \Delta_t$
  - 7:    $e_{t+1} = p_t - \Delta_t$
  - 8: **end for**
-

# Error-feedback

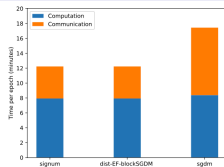
## Error-feedback



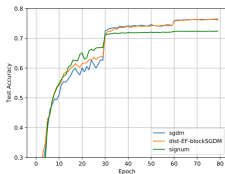
(a) Test accuracy w.r.t. epoch.



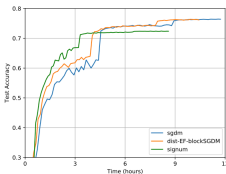
(b) Test accuracy w.r.t. time.



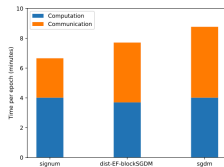
(c) Workload breakdown.



(d) Test accuracy w.r.t. epoch.



(e) Test accuracy w.r.t. time.



(f) Workload breakdown.

Figure 3: Distributed training results on the ImageNet dataset. Top: 7 workers; Bottom: 15 workers.

# The End