

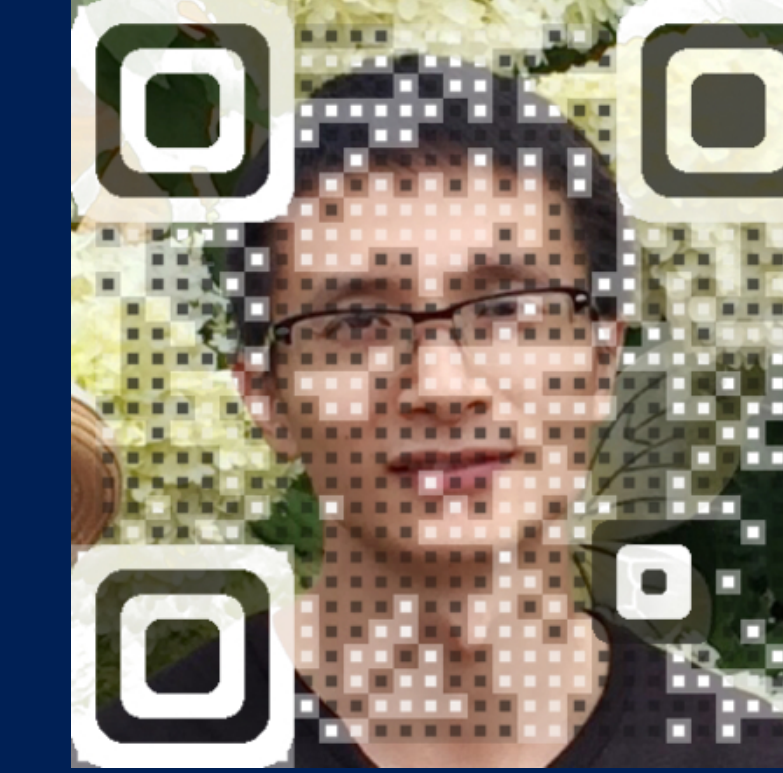
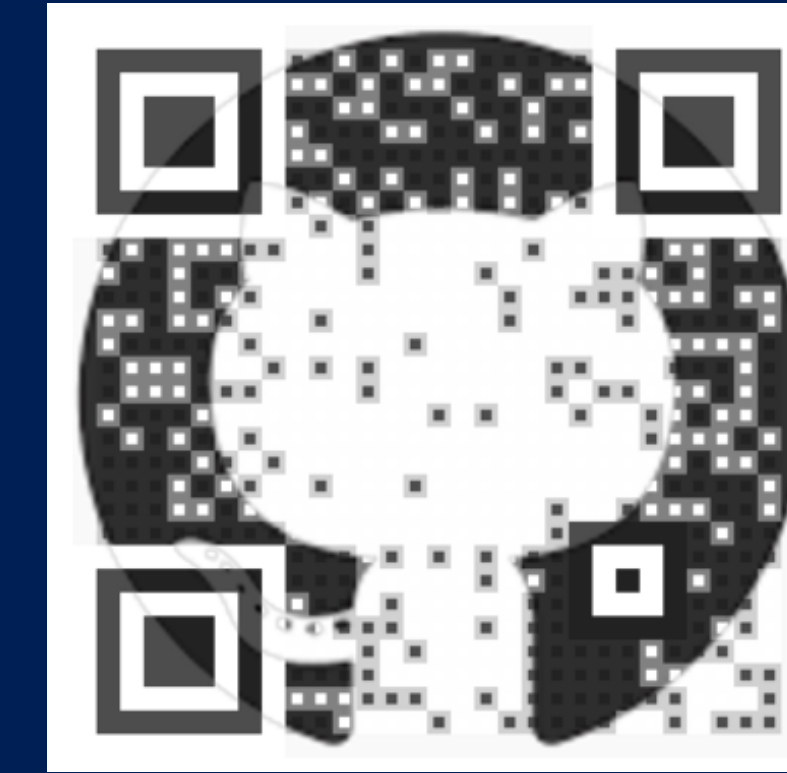


# TernGrad: Ternary Gradients to Reduce Communication in Distributed Deep Learning

Wei Wen<sup>1</sup>, Cong Xu<sup>2</sup>, Feng Yan<sup>3</sup>, Chunpeng Wu<sup>1</sup>, Yandan Wang<sup>4</sup>, Yiran Chen<sup>1</sup>, Hai (Helen) Li<sup>1</sup>

Duke University<sup>1</sup>, Hewlett Packard Labs<sup>2</sup>, University of Nevada - Reno<sup>3</sup>, University of Pittsburgh<sup>4</sup>

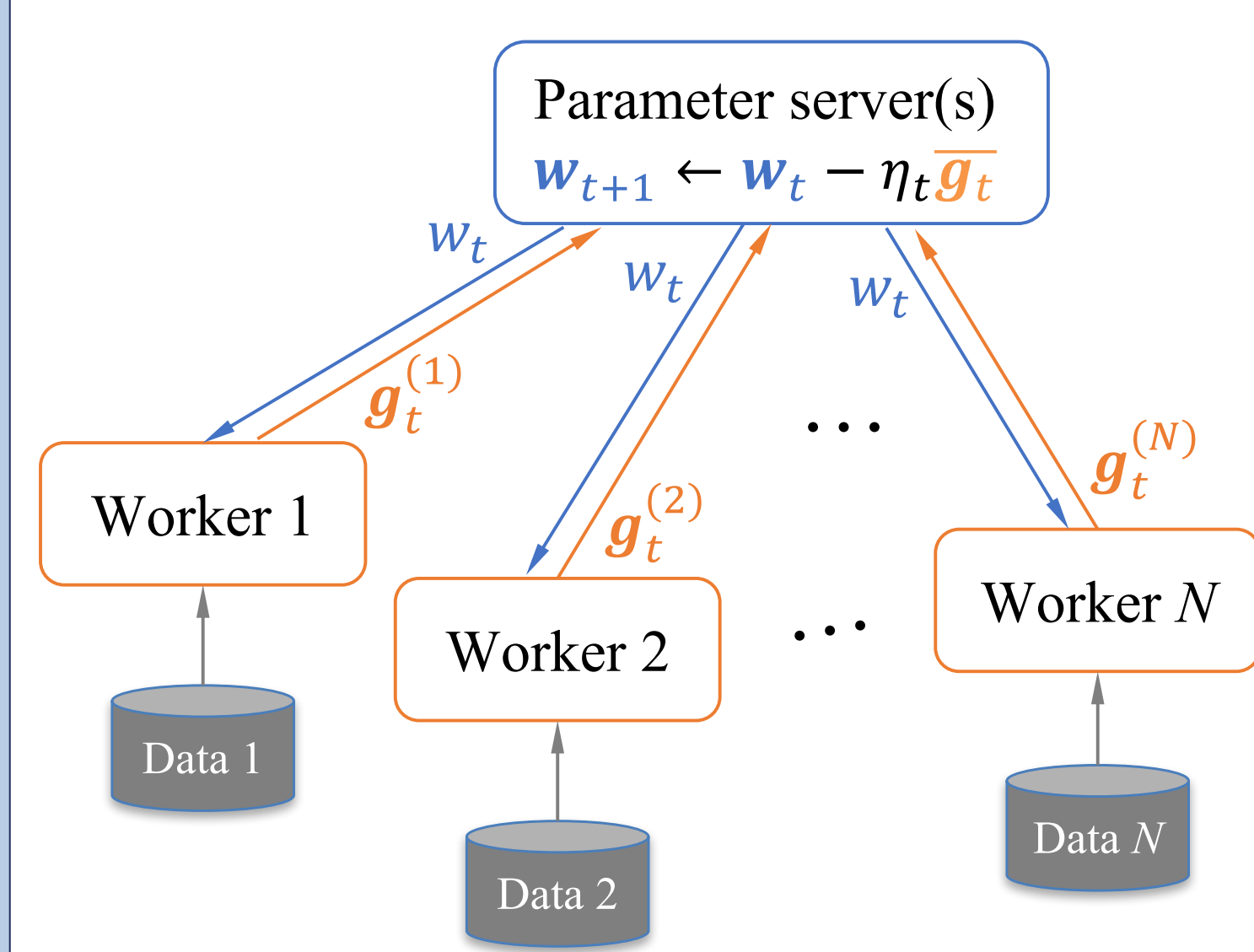
<sup>1</sup>{wei.wen, chunpeng.wu, yiran.chen, hai.li}@duke.edu, <sup>2</sup>cong.xu@hpe.com, <sup>3</sup>fyan@unr.edu, <sup>4</sup>yaw46@pitt.edu



## Background & Motivation

### Goal

Accelerate distributed training (of deep neural networks) by overcoming communication bottleneck.

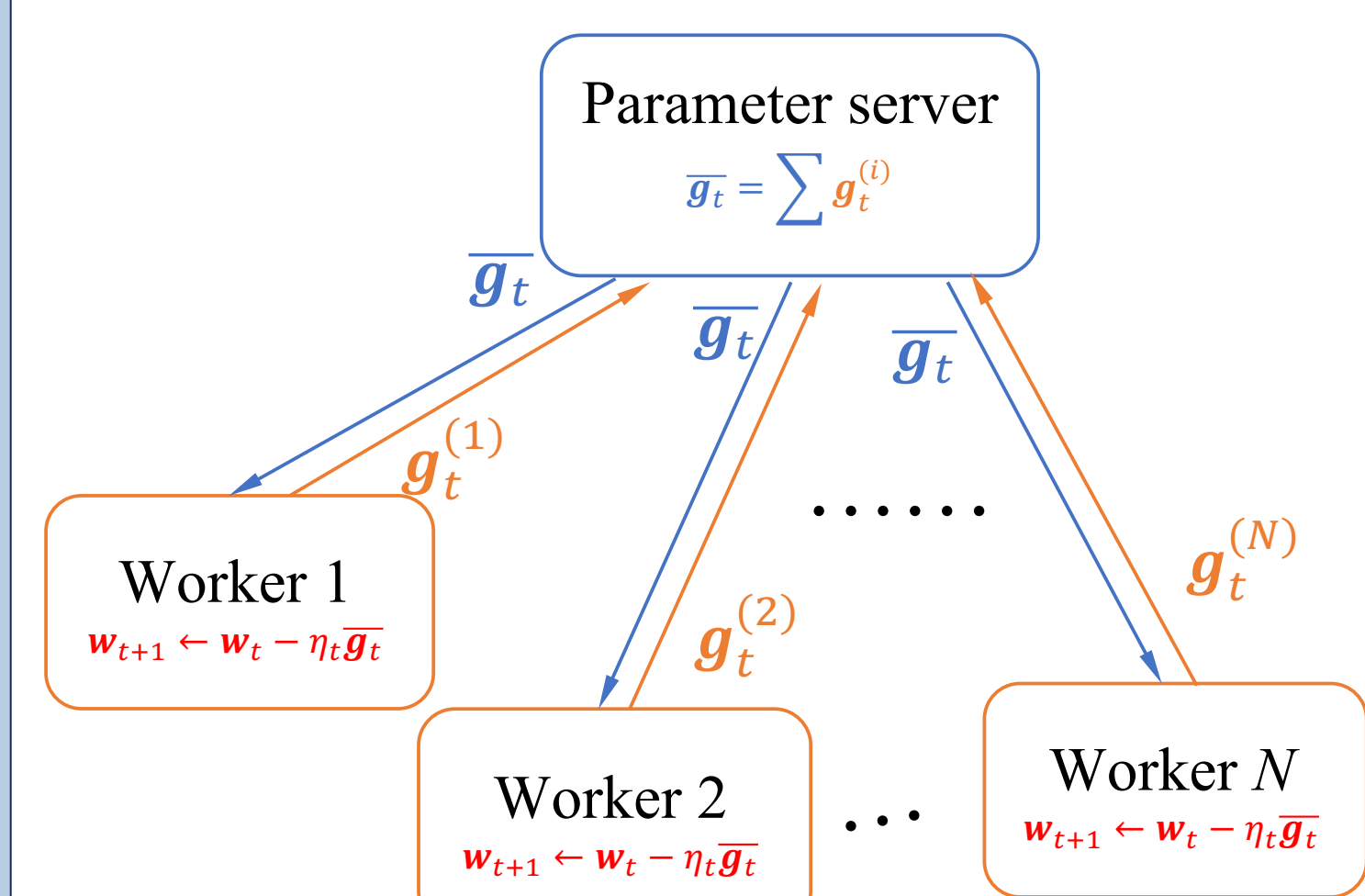


Synchronized Data Parallelism for Stochastic Gradient Descent (SGD):

1. Training data is split to  $N$  subsets
2. Each worker has a model replica (copy)
3. Each replica is trained on its subset of training data
4. Synchronization in parameter server(s)

Scalability (larger  $N$ ):

1. Computing time decreases as  $N$  increases
2. Communication time becomes the bottleneck
3. This work: ternary gradients to reduce communication



An alternative setting of Synchronized Data Parallelism :

1. Only exchange gradients
2. Each worker pulls average gradients from server and update weights locally
3. Gradient quantization can reduce communication in both directions
4. An identical model across workers (ensured by using the same random seed to initialize parameters)

## TernGrad and Convergence

### Stochastic Gradients without Bias

#### Batch Gradient Descent

$$C(w) \triangleq \frac{1}{n} \sum_{i=1}^n Q(z_i, w)$$

$$w_{t+1} = w_t - \frac{\eta_t}{n} \sum_{i=1}^n g_t^{(i)}$$

#### SGD

$$w_{t+1} = w_t - \eta_t \cdot g_t^{(I)}$$

$I$  is randomly drawn from  $[1, n]$

$$\mathbb{E}\{g_t^{(I)}\} = \nabla C(w) \quad \text{No bias}$$

#### TernGrad

$$w_{t+1} = w_t - \eta_t \cdot \text{ternarize}(g_t^{(I)})$$

$$\mathbb{E}\{\text{ternarize}(g_t^{(I)})\} = \nabla C(w) \quad \text{No bias}$$

$$\tilde{g}_t = \text{ternarize}(g_t) = s_t \cdot \text{sign}(g_t) \circ b_t$$

$$s_t \triangleq \|g_t\|_\infty \triangleq \max(\text{abs}(g_t))$$

$$\begin{cases} P(b_{tk} = 1 | g_t) = |g_{tk}| / s_t \\ P(b_{tk} = 0 | g_t) = 1 - |g_{tk}| / s_t \end{cases}$$

$$\begin{aligned} \mathbf{E}_{z,b}\{\tilde{g}_t\} &= \mathbf{E}_{z,b}\{s_t \cdot \text{sign}(g_t) \circ b_t\} \\ &= \mathbf{E}_z\{s_t \cdot \text{sign}(g_t) \circ \mathbf{E}_b\{b_t | z_t\}\} = \mathbf{E}_z\{g_t\} = \nabla_w C(w_t) \quad \text{No bias} \end{aligned}$$

Example:

$$g_t^{(i)}: [0.30, -1.20, \dots, 0.9]$$

$$s_t: 1.20$$

$$\text{Signs: } [1, -1, \dots, 1]$$

$$P(b_{tk} = 1 | g_t): [\frac{0.3}{1.2}, \frac{1.2}{1.2}, \dots, \frac{0.9}{1.2}]$$

$$\tilde{g}_t^{(i)}: [0, -1, \dots, 1] * 1.20$$

Convergence of standard SGD and TernGrad (Fisk 1965, Metivier 1981&1983, Bottou 1998)

**Assumption 1:**  $C(w)$  has a single minimum  $w^*$  and  $\forall \epsilon > 0, \inf_{\|w - w^*\|^2 > \epsilon} (w - w^*)^T \nabla_w C(w) > 0$

**Assumption 2:** Learning rate  $\gamma_t$  decreases neither very fast nor very slow  $\begin{cases} \sum_{t=0}^{+\infty} \gamma_t^2 < +\infty \\ \sum_{t=0}^{+\infty} \gamma_t = +\infty \end{cases}$

**Assumption 3:**  $\mathbb{E}\{\|g\|^2\} \leq A + B\|w - w^*\|^2$   $\mathbb{E}\{\|g\|_\infty \cdot \|g\|_1\} \leq A + B\|w - w^*\|^2$   
(gradient bound) SGD almost truly converges TernGrad almost truly converges

$\mathbb{E}\{\|g\|^2\} \leq \mathbb{E}\{\|g\|_\infty \cdot \|g\|_1\} \leq A + B\|w - w^*\|^2$  Stronger gradient bound in TernGrad

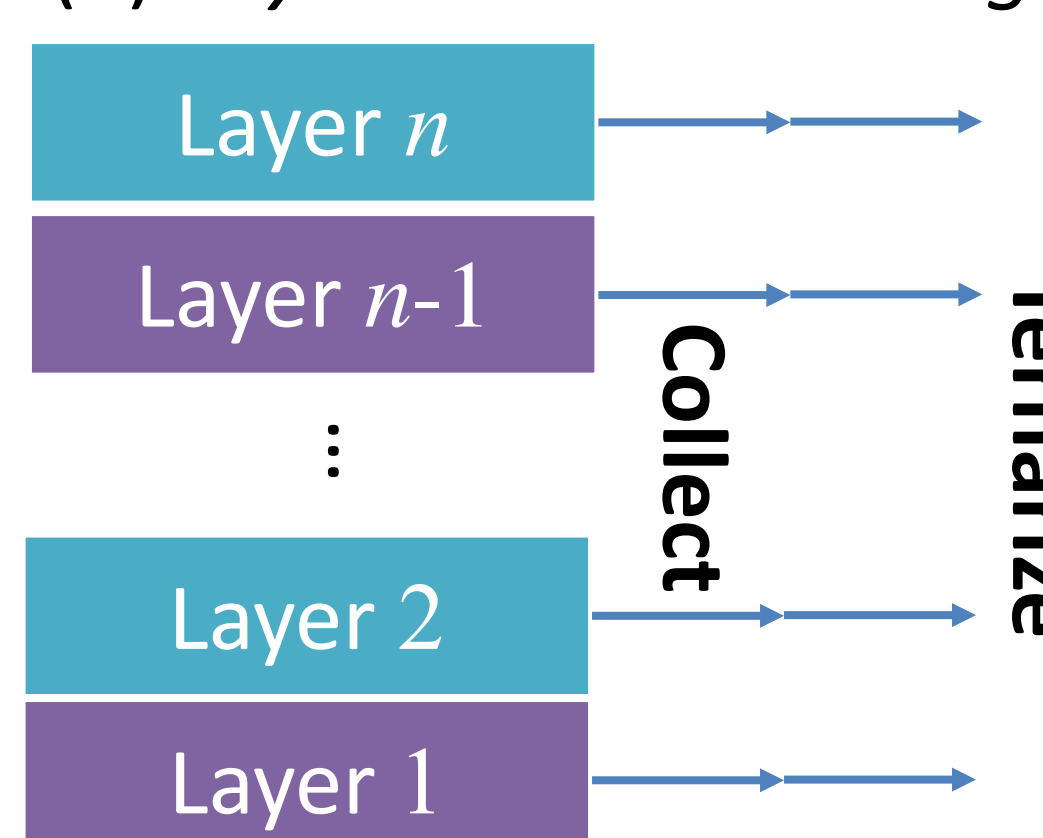
## Gradient Bound and Variance

Two views to improve the convergence of TernGrad:

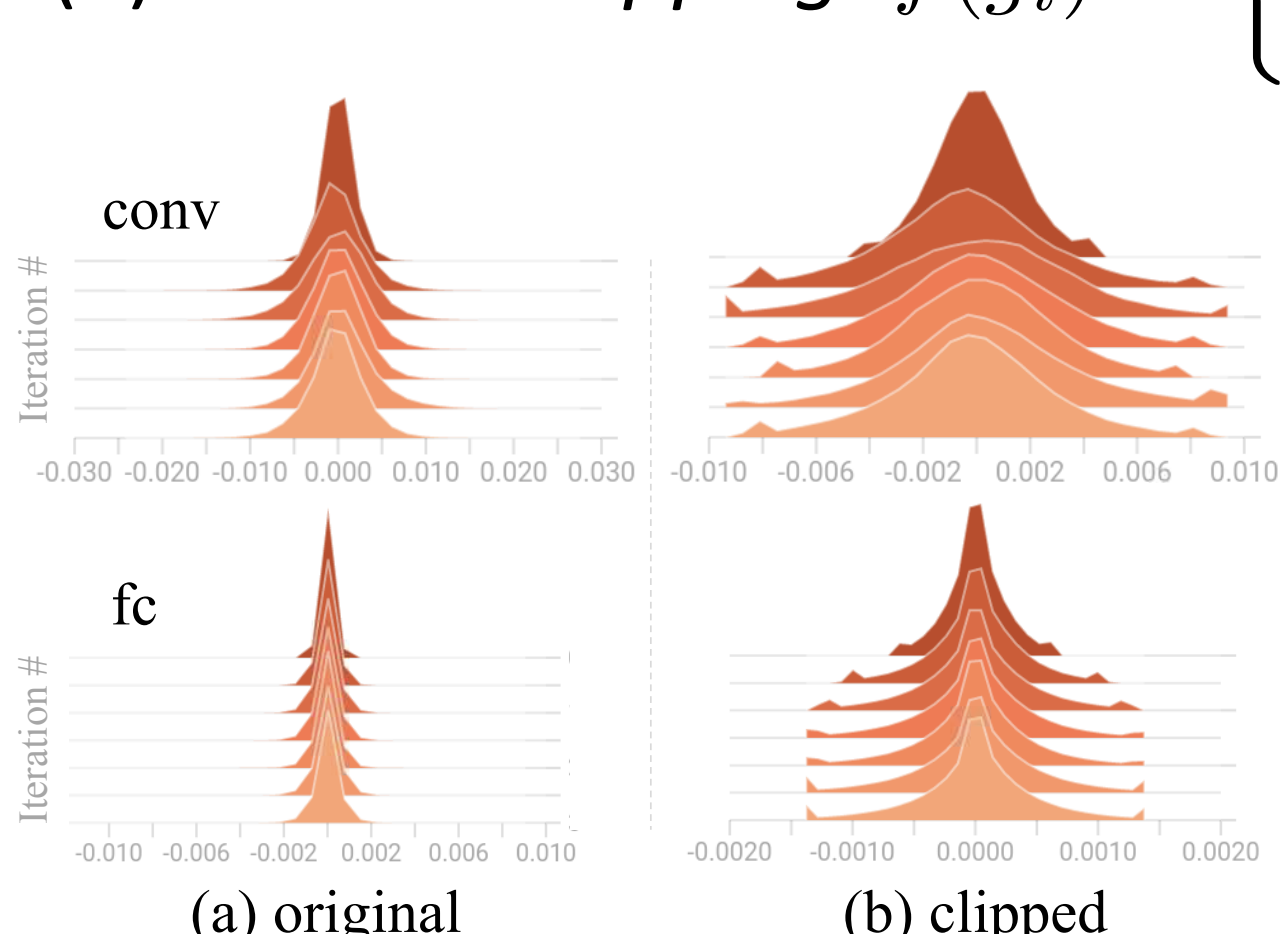
1. Variance
  - ✓ Relatively smaller  $\|g_t\|_\infty$  can achieve smaller gradient variance
2. Gradient Bound
  - ✓ Relatively smaller  $\|g_t\|_\infty$  pushes gradient bound (**Assumption 3**) of TernGrad closer to the one of standard SGD

Two methods to push gradient bound closer:

### (1) Layer-wise ternarizing



### (2) Gradient clipping



$$f(g_i) = \begin{cases} g_i & |g_i| \leq c\sigma \\ \text{sign}(g_i) \cdot c\sigma & |g_i| > c\sigma \end{cases}$$

$c=2.5$  works well for all tested datasets, DNNs and optimizers. Why? It

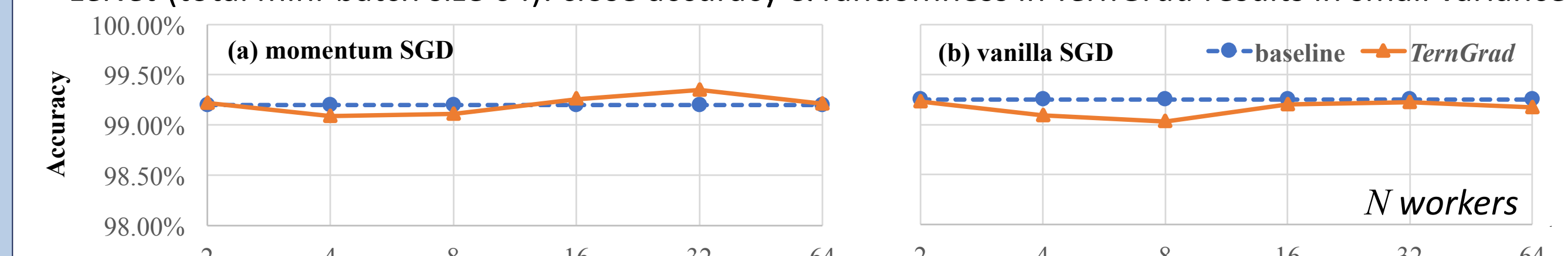
1. changes length 1%-1.5%
2. changes direction  $2^\circ - 3^\circ$
3. is variance reduction method with very small bias

Suppose Gaussian distribution with standard deviation  $\sigma$

Extreme case: ternarizing gradient one by one is just the floating SGD

## Experiments

① LeNet (total mini-batch size 64): close accuracy & randomness in TernGrad results in small variance

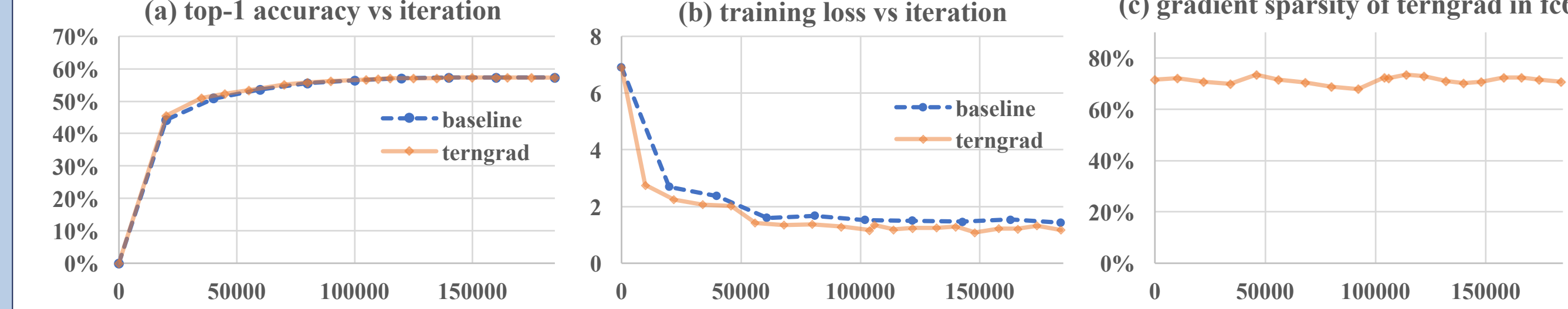


CIFAR-10, mini-batch size 64 per worker (all hyper-parameters are the same)

SGD	base LR	total mini-batch size	iterations	gradients	workers	accuracy
Adam	0.0002	128	300K	floating	2	86.56%
TernGrad	0.0002	128	300K	floating	2	85.64% (-0.92%)
Adam	0.0002	2048	18.75K	floating	16	83.19%
TernGrad	0.0002	2048	18.75K	floating	16	82.80% (-0.39%)

D. P. Kingma, 2014

③ AlexNet trained on 4 workers with mini-batch size 512



GoogLeNet (<2% accuracy loss on average)

base LR	mini-batch size	workers	iterations	gradients	weight decay	DR	top-5
0.04	128	2	600K	floating	4e-5	0.2	88.30%
TernGrad	0.04	2	600K	floating	1e-5	0.08	86.77%
0.08	256	4	300K	floating	4e-5	0.2	87.82%
TernGrad	0.08	4	300K	floating	1e-5	0.08	85.96%
0.10	512	8	300K	floating	4e-5	0.2	89.00%
TernGrad	0.10	8	300K	floating	2e-5	0.08	86.47%

All other hyper-parameters are tuned for standard SGD and fixed in TernGrad. Tune hyper-parameters specifically for TernGrad may reduce accuracy loss

② TernGrad: Randomness & regularization (1) decrease randomness in dropout or (2) use smaller weight decay (in large-scale dataset like ImageNet) No new hyper-parameters added

AlexNet

base LR	mini-batch size	workers	iterations	gradients	weight decay	DR†	top-1	top-5
0.01	256	2	370K	floating	0.0005	0.5	57.33%	80.56%
TernGrad	0.01	2	370K	floating	0.0005	0.2	57.61%	80.47%
TernGrad-noclip	0.01	2	370K	floating	0.0005	0.2	54.63%	78.16%
0.02	512	4	185K	floating	0.0005	0.5	57.32%	80.73%
TernGrad	0.02	4	185K	floating	0.0005	0.2	57.28%	80.23%
0.04	1024	8	92.5K	floating	0.0005	0.5	56.62%	80.28%
TernGrad	0.04	8	92.5K	floating	0.0005	0.2	57.54%	80.25%

† DR: dropout ratio, the ratio of dropped neurons. ‡ TernGrad without gradient clipping. Noise in TernGrad helps SGD escaping from sharp minima when batch size is large (N. S. Keskar, ICLR 2017)

④ A performance model to evaluate the speedup

