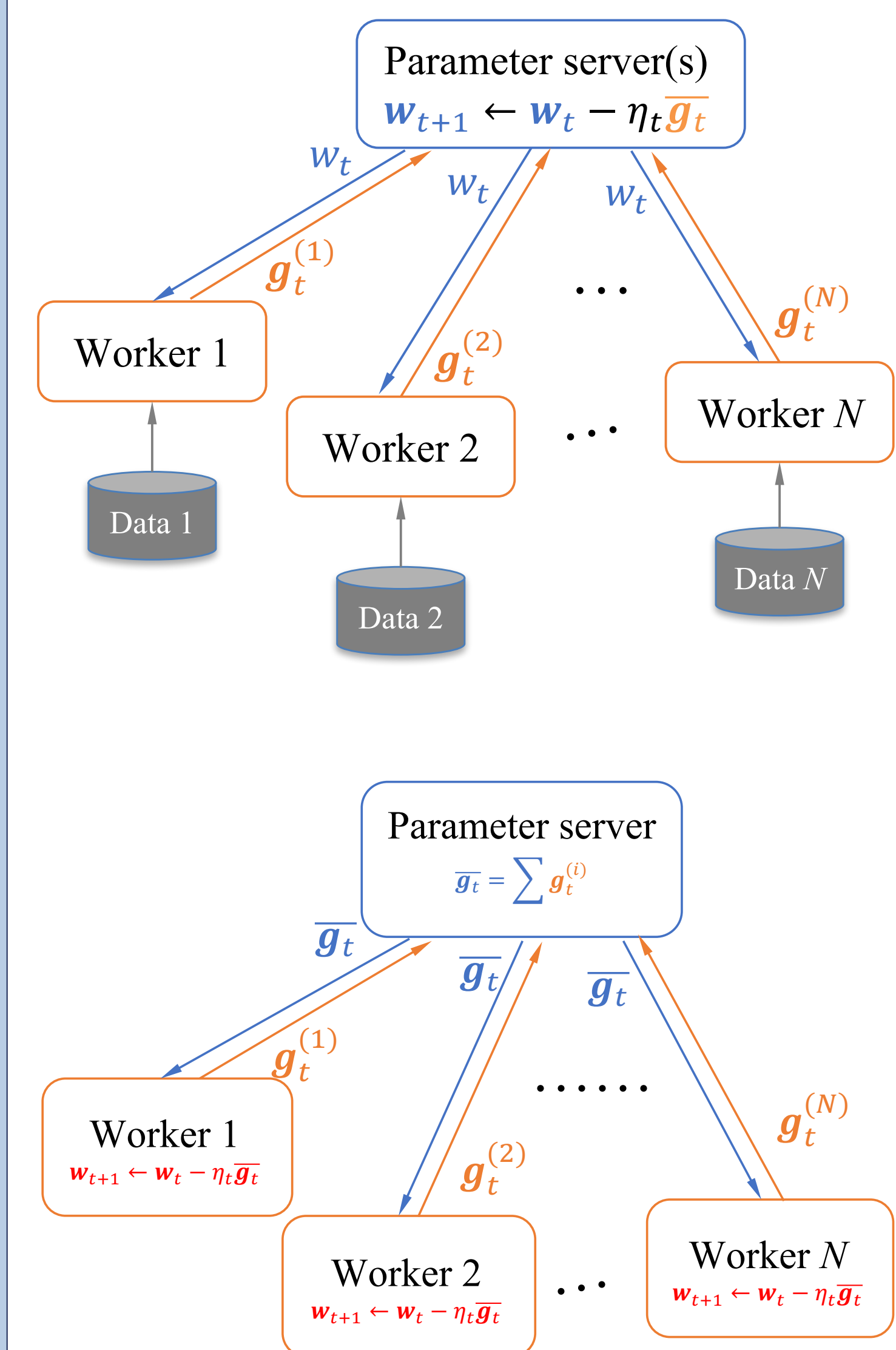


## Background & Motivation

### Goal

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

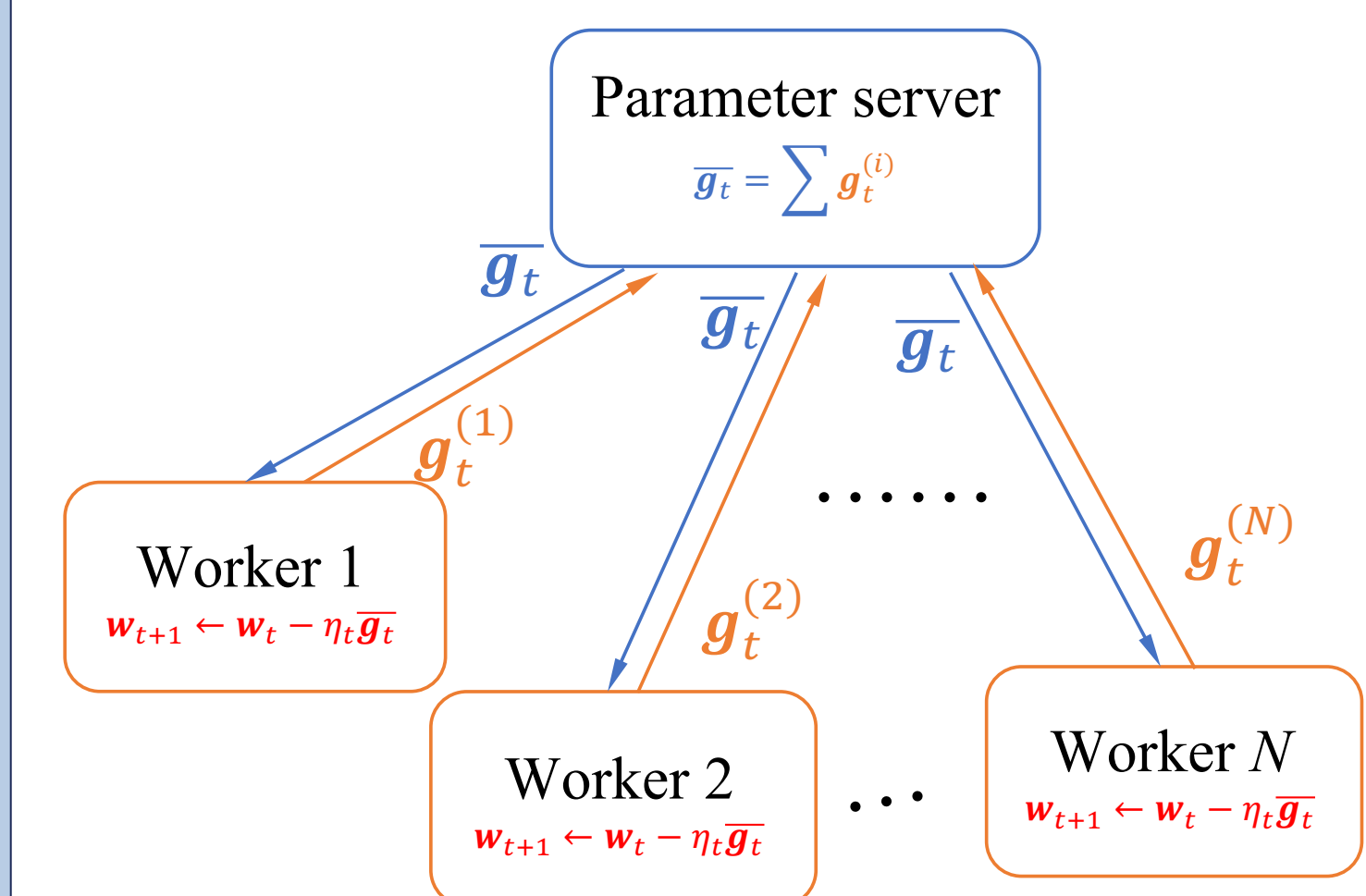


Synchronized Data Parallelism for Stochastic Gradient Descent (SGD):

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

Scalability (large  $N$ ):

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



An alternative setting Synchronized Data Parallelism :

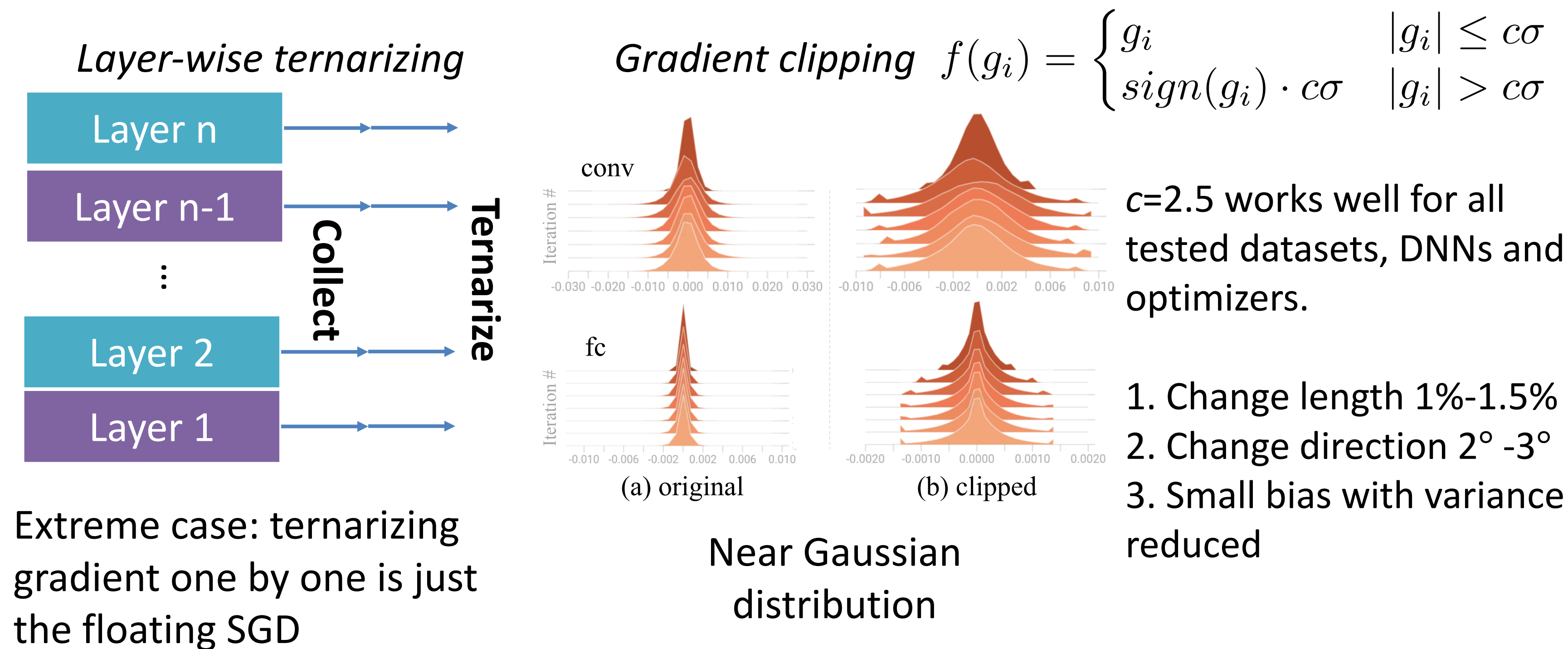
1. Only exchange gradients
2. Quantization can reduce communication in both directions
3. An identical model across workers (use the same initialization seed)

## Gradient Bound and Variance

Two views to improve the convergence of *TernGrad*:

1. Variance
  - ✓ Relatively smaller  $\|g_t\|_\infty$  can result in smaller variance
2. Gradient Bound
  - ✓ Relatively smaller  $\|g_t\|_\infty$  pushes gradient bound assumption of *TernGrad* closer to the gradient bound in standard SGD

Two methods:



## 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

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

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

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

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

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

$$\mathbf{E}_{z,b}\{\tilde{g}_t\} = \mathbf{E}_{z,b}\{s_t \cdot \text{sign}(g_t) \circ b_t\}$$

$$= \mathbf{E}_z\{s_t \cdot \text{sign}(g_t) \circ \mathbf{E}_b\{b_t | z_t\}\} = \mathbf{E}_z\{g_t\} = \nabla_w C(w_t) \quad \text{No bias}$$

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  $\begin{cases} \sum_{t=0}^{+\infty} \gamma_t^2 < +\infty \\ \sum_{t=0}^{+\infty} \gamma_t = +\infty \end{cases}$

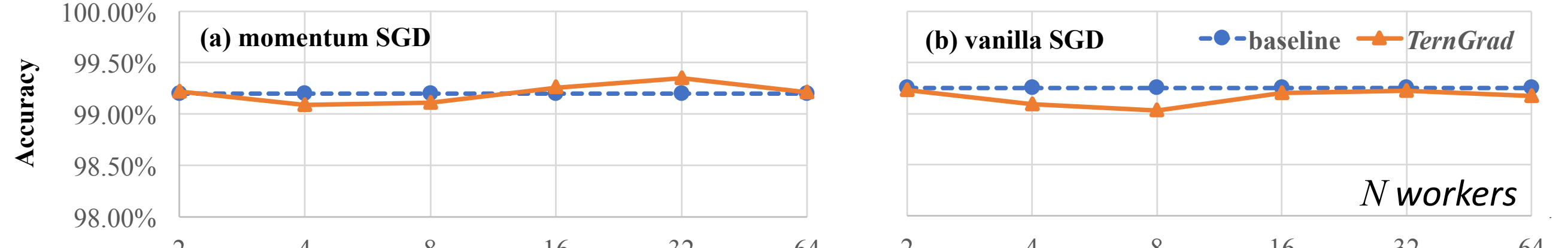
**Assumption 3:**  $\mathbb{E}\{\|g\|^2\} \leq A + B\|w - w^*\|^2$   $\mathbb{E}\{\|g\|_\infty \cdot \|g\|_1\} \leq A + B\|w - w^*\|^2$

(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*

## Experiments

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
D. P. Kingma, 2014	Adam	0.0002	128	300K	floating	2	86.56%
					TernGrad	2	85.64% (-0.92%)
	Adam	0.0002	2048	18.75K	floating	16	83.19%
					TernGrad	16	82.80% (-0.39%)

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	0.08	88.30%
					TernGrad	1e-5		86.77%
0.08	256	4	300K		floating	4e-5	0.2	87.82%
					TernGrad	1e-5		85.96%
0.10	512	8	300K		floating	4e-5	0.2	89.00%
					TernGrad	2e-5		86.47%

Tune hyper-parameters specifically for *TernGrad* may reduce accuracy loss

*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 <sup>†</sup>	top-1	top-5
0.01	256	2	370K	floating	0.0005	0.5	57.33%	80.56%
				TernGrad	0.0005	0.2	57.61%	80.47%
				TernGrad-noclip	0.0005	0.2	54.63%	78.16%
0.02	512	4	185K	floating	0.0005	0.5	57.32%	80.73%
				TernGrad	0.0005	0.2	57.28%	80.23%
0.04	1024	8	92.5K	floating	0.0005	0.5	56.62%	80.28%
				TernGrad	0.0005	0.2	57.54%	80.25%

<sup>†</sup> DR: dropout ratio, the ratio of dropped neurons. <sup>‡</sup> *TernGrad* without gradient clipping.

A performance model to evaluate the speedup

