

1-bit Adam: Communication Efficient Large-Scale Training with Adam's Convergence Speed

Authors: Hanlin Tang (presenter), Shaoduo Gan, Samyam Rajbhandari, Ammar Ahmad Awan, Conglong Li, Xiangru Lian, Ji Liu, Ce Zhang, Yuxiong He



Adam

Momentum SGD

Fast Convergence

Communication Compression

1-bit Adam

BERT-Large Pretrain (64 GPUs, 40 Gbps Ethernet):

174 hours

1-bit Adam

Efficient kernels

51 hours

3.4x Faster

Algorithm Design

Gradient Compression

Compression:

$$\begin{array}{c} \text{32 bits} \\ \underbrace{\left(\begin{array}{c} -2.2 \\ 3.3 \end{array} \right)}_{\text{Sign()}} \longrightarrow \underbrace{\left(\begin{array}{c} -1 \\ 1 \end{array} \right)}_{\text{1 bit}} \end{array}$$

32x of communication reduction

Slower convergence



Slower training speed

Error Compensation:

$$g_1 = 3.3, \quad g_2 = -1.1, \quad \gamma = 1$$

Step 1: $x = x + \text{Sign}(3.3)$; $error = 2.3$

Step 2: $x = x + \text{Sign}(-1.1 + \text{2.3})$; $error = 0.2$

32x of communication reduction

Same convergence



Faster training speed

Adam?

Adam:

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$$

$$= \sum_{s=0}^t \beta_1^{t-s} \mathbf{g}_s$$

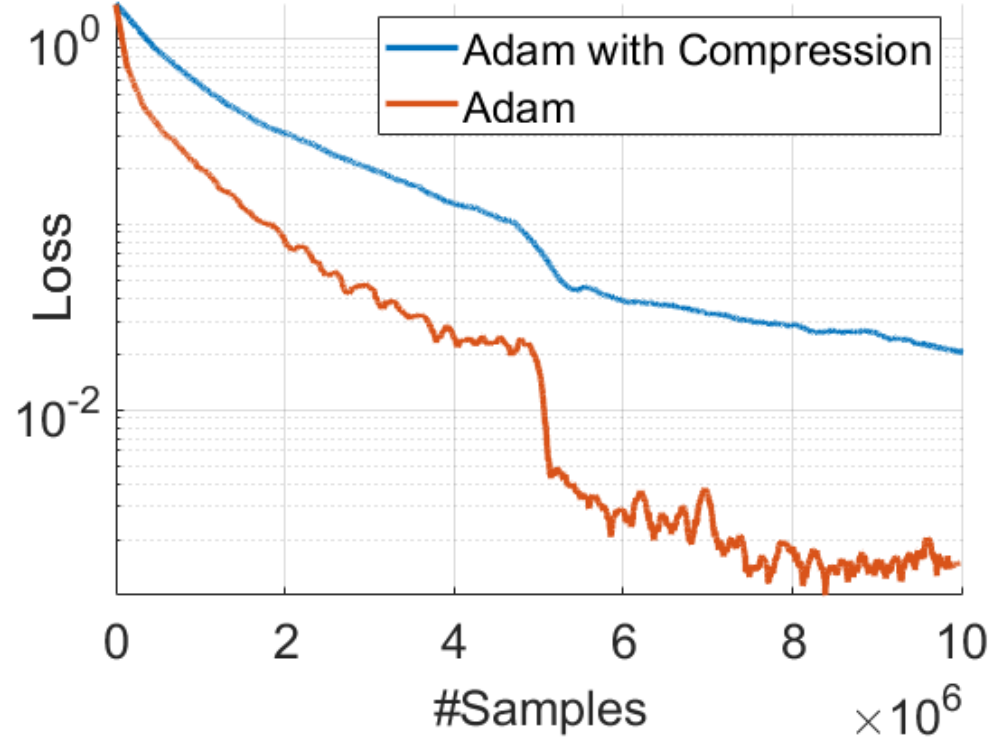
$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2$$

$$= \sum_{s=0}^t \beta_2^{t-s} \mathbf{g}_s^2$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma \frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_t}}$$

Nonlinear term

Incompatible with error compensation



Adam:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma \frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_t}}$$

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$$

$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t^2$$

Constant v

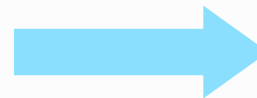


Quasi Adam:

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma \frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_t}}$$

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$$

Define: $\gamma_v = \frac{\gamma}{\sqrt{\mathbf{v}}}$

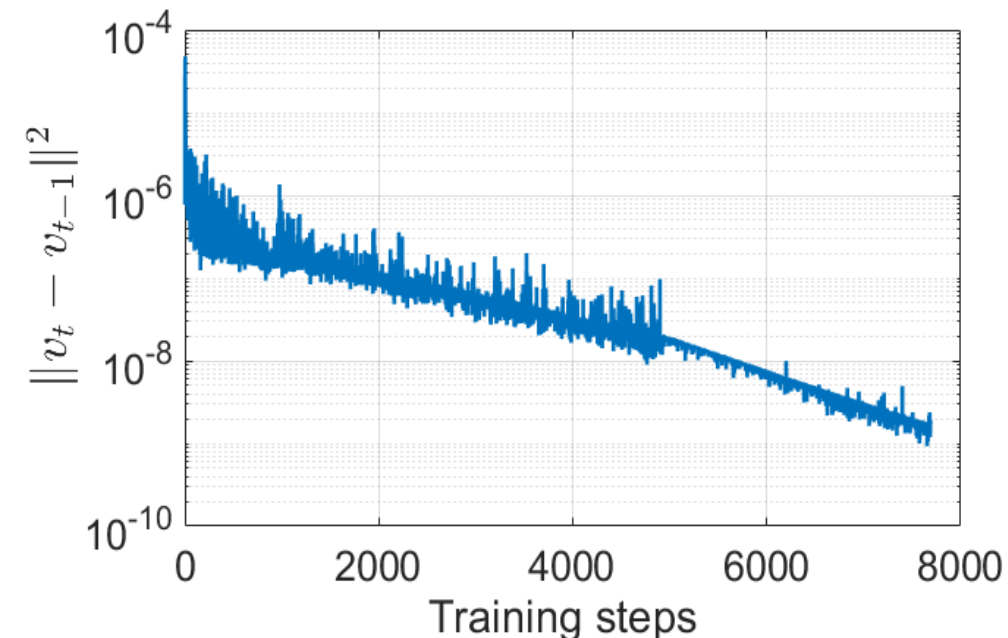


Momentum SGD:

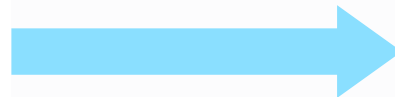
$$\mathbf{x}_{t+1} = \mathbf{x}_t - \gamma_v \mathbf{m}_t;$$

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t;$$

\mathbf{v}_t becomes more stable along with training.



1-bit Adam



Warmup phase:

- Original Adam (No compression)

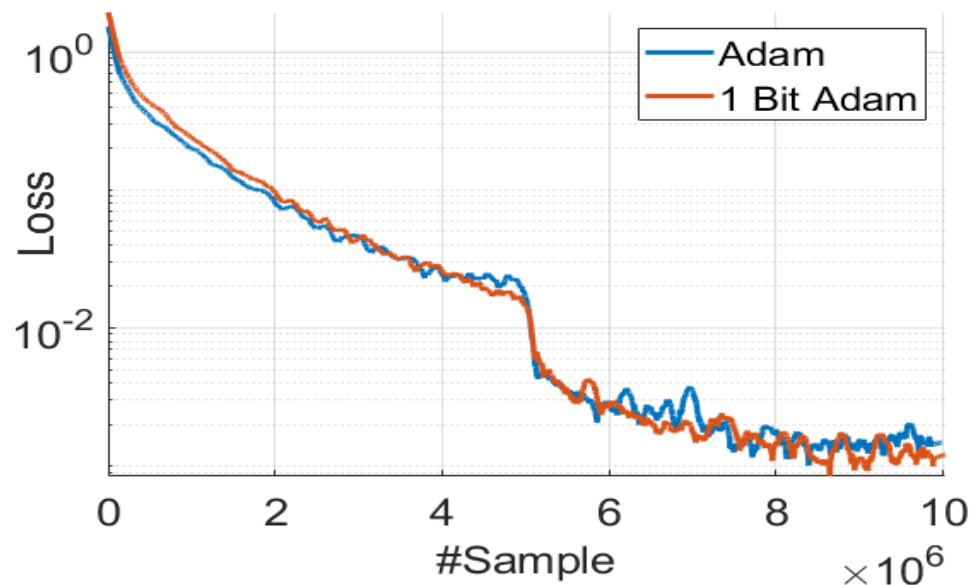
Compression phase:

- Keep \mathbf{v}_t unchanged
- Compress \mathbf{m}_t only

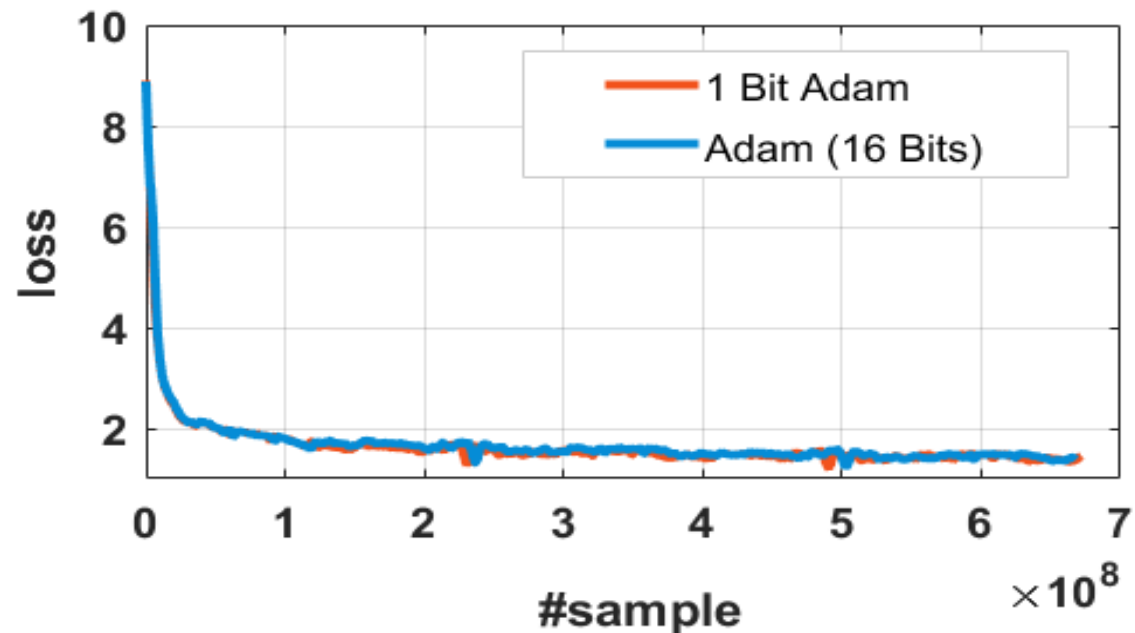
Convergence Results

32x speedup after warmup, **5x** speedup in total (10%~15% warmup needed)

ResNet-18



BERT-large



System Implementation

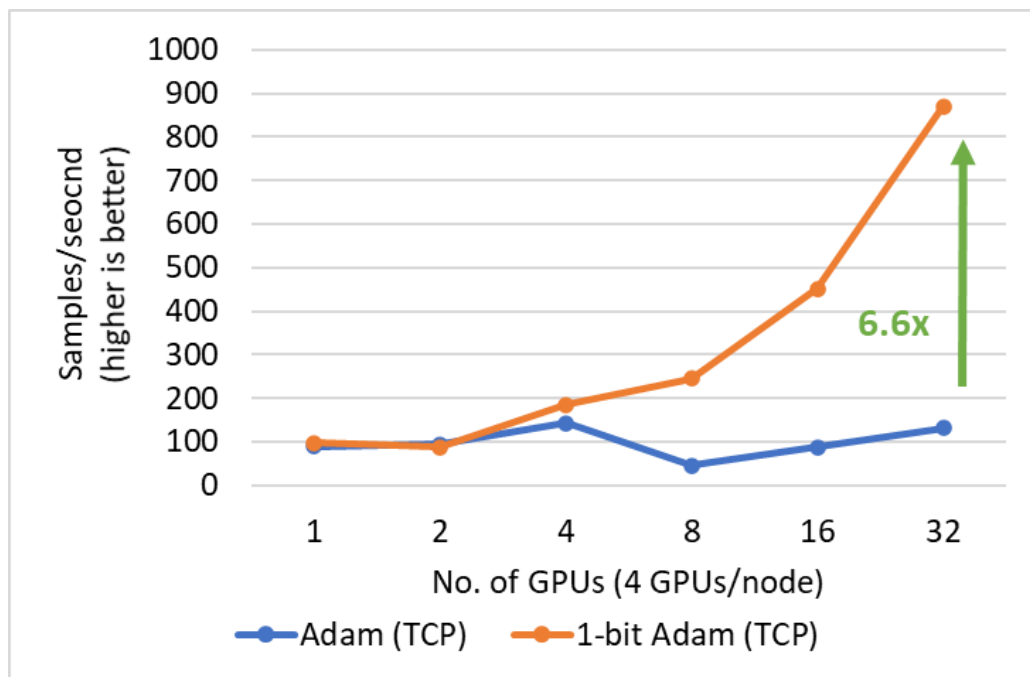
System Implementation

- An efficient implementation of **CUDA-Aware Collective Communication** backend that works for both InfiniBand and Ethernet network.
- Efficient compression kernel for encoding gradients (**32**-bits) into **1**-bit.
- Compatible with half-precision training (16-bits).
- Released in **Micorsoft DeepSpeed**.

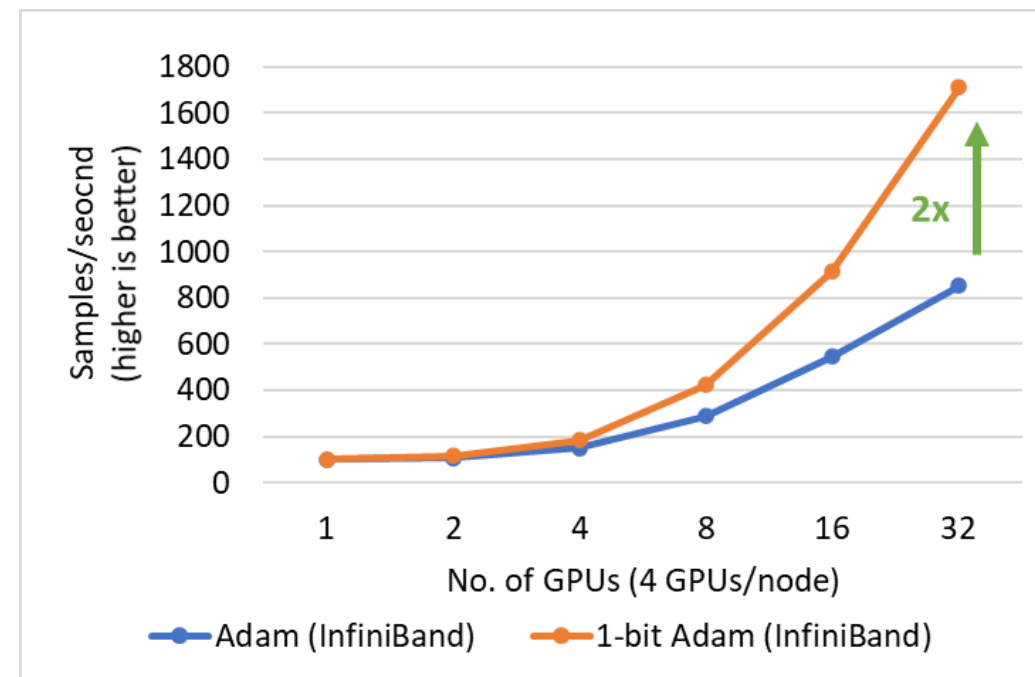
System Implementation

Results: BERT-Large pretraining. **End-to-end** training speed comparison (including both computation and communication)

4 V100/node, 40 Gbps Ethernet



8 V100/node, 100 Gbps InfiniBand



Thank You