## Distributed Optimization for Machine Learning

Lecture 14 - Communication-efficient Distributed Training - Part I

Tianyi Chen

School of Electrical and Computer Engineering Cornell Tech, Cornell University

October 15, 2025





### Example: air pollution prediction in smart cities



There is a community of multiple houses, where each house has a smart sensor that records environmental information.

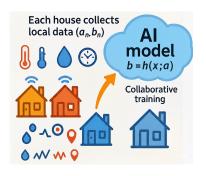
Each house collects data pairs  $\boldsymbol{\xi}_n = \{\boldsymbol{a}_n, \boldsymbol{b}_n\}$  over time, where:

 $a_n =$ (temperature, humidity, time, location, etc.)

 $\boldsymbol{b}_n = \text{concentration of a particular pollutant.}$ 



# Motivation of distributed training



**Goal:** All houses want to collaboratively train a machine learning model to predict future  $\boldsymbol{b}$  given  $\boldsymbol{a}$ :

$$\boldsymbol{b} = h(\boldsymbol{x}; \boldsymbol{a})$$

where x denotes the model parameters to be learned.



### Example: next-word prediction on smart keyboards

Each smartphone collects sequences of words typed by the user.

$$oldsymbol{a}_n = \mathrm{Vec} \left( egin{array}{c} \mathsf{current word} \\ dots \\ \mathsf{past word} \end{array} 
ight),$$

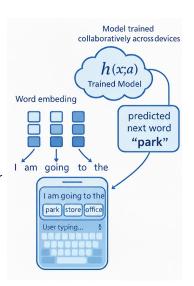
 $\boldsymbol{b}_n = \operatorname{Vec}(\operatorname{next} \operatorname{word})$ .

Here,  $\operatorname{Vec}(\cdot)$  denotes the Word-to-vector embedding operation.

**Goal:** Learn a model h(x; a) to predict the next word embedding:

$$\boldsymbol{b} = h(\boldsymbol{x}; \boldsymbol{a})$$





## Why perform distributed training?

**Key question:** Why *distributed* data?

#### Main reason: Privacy!

- Each house may not want to share its raw sensor data with others or with a central server.
- Instead, they exchange only model updates or gradients to preserve local data confidentiality.

#### Secondary reason: Bandwidth and latency!

- Reduce communication overhead of transferring large datasets.
- Enable real-time, edge-level learning across smart devices.

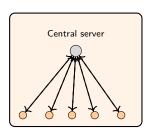


## Optimization formulation of data parallelism

A network of n nodes (such as mobile devices) collaborate to solve:

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}), \quad \text{where } \left[ f_i(\mathbf{x}) = \mathbb{E}_{\boldsymbol{\xi}_i \sim D_i} [F(\mathbf{x}; \boldsymbol{\xi}_i)] \right]$$

- Each component  $f_i : \mathbb{R}^d \to \mathbb{R}$  is local and private to node i.
- Random variable  $\xi_i$  denotes local data following distribution  $D_i$ .
- $D_i$  may be different  $\Rightarrow$  data heterogeneity.



Local data on nodes



## Parallel SGD: compute locally, communicate globally

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}), \quad \text{where } f_i(\mathbf{x}) = \mathbb{E}_{\xi_i \sim D_i}[F(\mathbf{x}; \xi_i)].$$

**PSGD** 

$$g_i^k = \nabla F(\mathbf{x}^k; \xi_i^k)$$
 (Local compt.)

$$g_i^k = \nabla F(\mathbf{x}^k; \xi_i^k)$$
 (Local compt.) 
$$\mathbf{x}^{k+1} = \mathbf{x}^k - \frac{\eta}{n} \sum_{i=1}^n g_i^k$$
 (Global comm.)

- Each node *i* samples mini-batch  $\xi_i^k$  and computes  $\nabla F(\mathbf{x}^k; \xi_i^k)$ .
- All nodes synchronize (i.e., globally average) to update x.



## Communication overhead of distributed training

#### Communication overhead:

- Each entry of a d-dimensional vector (model or gradient) requires 32 bits by default float32 (IEEE 754 single-precision floating-point).
- Each upload or download of the vector incurs:

Communication cost = 
$$32 \times d \times n$$

#### where

- 32: bits per entry,
- d: number of dimensions ( $10^6 \sim 10^{11}$ ), n: number of workers ( $10^3 \sim 10^4$ ).
- $\Rightarrow$  Total communication per round =  $\mathcal{O}(10^{10} \text{ to } 10^{16})$  bits.



#### Solutions to overcome communication overhead

**Goal:** Reduce the total communication cost per iteration:

$$32 \times d \times n$$

#### Possible solutions:

- S1: Reduce the communication rounds: e.g., Local SGD: perform  $\tau$  local updates before synchronization to reduce communication frequency while maintaining accuracy.
- S2: Reduce the number of bits via quantization/sparsification: e.g., Stochastic or deterministic quantization, threshold-based or Top-k sparsification.
- S3: Reduce the number of workers: e.g., Randomized / cyclic and adaptive worker selection (LAG).



#### Table of Contents

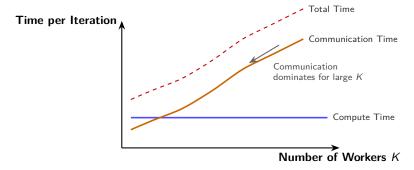
Reduce the Communication Rounds via Local SGD

Reduce the Number of Bits via Quantization



#### Communication bottleneck in Parallel SGD

- In Parallel SGD, workers synchronize after every step.
- $\blacksquare$  Comm. dominates runtime when n is large or network is slow.





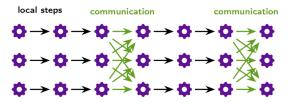
## Idea: Local updates before synchronization

**Key idea:** Each node *i* performs several SGD steps before averaging.

$$\mathbf{x}_{i}^{(s+1)} = \mathbf{x}_{i}^{(s)} - \eta \nabla F(\mathbf{x}_{i}^{(s)}; \boldsymbol{\xi}_{i}^{(s)}), \quad s = 0, \dots, \tau - 1$$

where  $\mathbf{x}_{i}^{(0)} = \mathbf{x}^{k}$ . After every  $\tau$  steps:

$$\mathbf{x}^{k+1} = \frac{1}{n} \sum_{i=1}^{n} \mathbf{x}_i^{(\tau)}$$



**Benefit:** Reduces communication by a factor of  $\tau$ .



#### Mini-batch SGD vs. Local SGD

#### Mini-batch or Parallel SGD:

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \frac{1}{n\tau} \sum_{i=1}^n \sum_{s=1}^\tau \nabla F(\mathbf{x}^t; \boldsymbol{\xi}_i^{t,s})$$

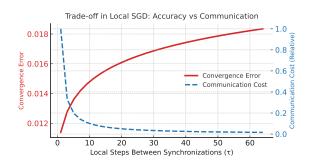
#### Local SGD:

$$\mathbf{x}_{i}^{t+1} = \begin{cases} \mathbf{x}_{i}^{t} - \eta_{t} \nabla F(\mathbf{x}_{i}^{t}; \boldsymbol{\xi}_{i}^{t}), & t \mod \tau \neq 0 \\ \frac{1}{n} \sum_{i=1}^{n} \left( \mathbf{x}_{i}^{t} - \eta_{t} \nabla F(\mathbf{x}_{i}^{t}; \boldsymbol{\xi}_{i}^{t}) \right), & t \mod \tau = 0 \end{cases}$$

Method	Mini-batch SGD	Local SGD
# Comm. rounds	K	K
Batch size	$n\tau$	n
# Model updates	K	$\tau K$
# Gradient calcs	nτK	nτK



### Communication vs. computation trade-off



Runtime per iteration = Compute 
$$+\frac{1}{\tau}$$
Comm.

**Insight:** Increasing  $\tau$  improves efficiency but risks model drift.



# Why Local SGD works under homogeneous data?

If  $f(\mathbf{x})$  is convex and all workers start synchronized  $(\mathbf{x}_i^t = \bar{\mathbf{x}}^t)$ :

$$f_i(\mathbf{x}_i^{t+\tau}) \leq f_i(\bar{\mathbf{x}}^t) - (\text{descent term for worker } i).$$

Thus, synchronization preserves global descent.

$$f(\bar{\mathbf{x}}^{t+\tau}) \leq \frac{1}{n} \sum_{i=1}^{n} f_i(\mathbf{x}_i^{t+\tau}) \leq f(\bar{\mathbf{x}}^t) - (\text{averaged progress}).$$

Not true in general if  $f_i$  differ across workers (non-i.i.d.).



# Quadratic objectives: analytical insight

For local quadratic objectives  $f_i(\mathbf{x}) = \frac{1}{2}\mathbf{x}^{\top}\mathbf{A}_i\mathbf{x} - \mathbf{b}_i^{\top}\mathbf{x}$ :

$$\mathbf{x}_i^{k+1} = \mathbf{x}_i^k - \eta_k \nabla f_i(\mathbf{x}_i^k).$$

Averaging yields:

$$\mathbb{E}[\bar{\boldsymbol{x}}^{k+1}|\mathcal{F}^k] = \bar{\boldsymbol{x}}^k - \eta_k \frac{1}{n} \sum_{i=1}^n \nabla f_i(\boldsymbol{x}_i^k) \approx \bar{\boldsymbol{x}}^k - \eta_k \nabla f(\bar{\boldsymbol{x}}^k).$$

Hence, Local SGD mimics global descent dynamics.



# Local SGD improves efficiency in quadratic setting

#### **Theorem 1** (Local SGD under smooth and convex loss)

The error bound for Local SGD with  $\tau$  local updates equals the bound for Mini-batch SGD with batch size n and  $K\tau$  rounds:

$$\epsilon_{ ext{L-SGD}} := rac{1}{K} \sum_{k=1}^{K} \mathbb{E}[||
abla f(oldsymbol{x}^k)||_2^2] = oldsymbol{\Theta}igg(rac{1}{K au} + rac{\sigma}{\sqrt{nK au}}igg)$$

- lacktriangle More local updates au always help convergence.
- Mini-batch SGD:  $\epsilon_{\mathrm{MB-SGD}} = \Theta\left(\frac{1}{K} + \frac{\sigma}{\sqrt{nK\tau}}\right)$
- Local SGD can be better given the same computation budget.



## Performance for general convex objectives\*

Upper and lower bounds for Local SGD:

$$\text{Upper:} \hspace{0.5cm} \epsilon_{\textit{L-SGD}} = \mathcal{O}\bigg(\frac{\sigma^{2/3}}{\textit{K}^{2/3}\tau^{1/3}} + \frac{\sigma}{\sqrt{\textit{nK}\tau}}\bigg)$$

Lower: 
$$\Omega\left(\frac{\sigma^{2/3}}{K^{2/3}\tau^{2/3}} + \frac{\sigma}{\sqrt{nK\tau}}\right)$$

Mini-batch SGD: 
$$\Theta\left(\frac{1}{K} + \frac{\sigma}{\sqrt{nK\tau}}\right)$$

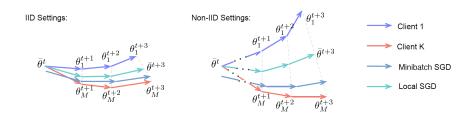
Local SGD better when  $K \lesssim \tau$ , worse when  $K \gtrsim \tau$ .

Woodworth et al. "Is Local SGD Better than Mini-batch SGD?", ICML 2020



# Why local SGD may fail under heterogeneous data?

In heterogeneous (non-i.i.d.) data settings, local gradients are misaligned



- Local updates diverge due to heterogeneous data (Γ<sup>2</sup>).
- Need additional assumptions to control gradient dissimilarity.
- Larger  $\tau \Rightarrow$  greater deviation from the global model.

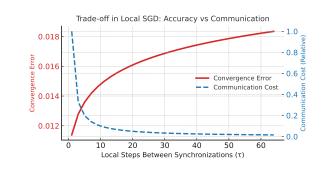


# Summary: Parallel SGD vs local SGD

Aspect	Parallel SGD	Local SGD
Communication	Every iteration	Every $ au$ iterations
Local computation	1 gradient step	au local steps
Speed	Communication-limited	Compute-efficient
Convergence rate	Stable	Slower $((\eta^2 \tau \Gamma^2)$ bias)
Best for	Data centers, i.i.d. data	Federated settings



### Takeaway: Communication - accuracy trade-off



- au au = 1: fully synchronized (Parallel SGD)
- $lue{ au}$  au>1: fewer syncs  $\Rightarrow$  faster but drift grows
- $lue{}$  Choose au based on network bandwidth and data heterogeneity



Rule of thumb:  $au^* \propto \sqrt{rac{c_{
m comm}}{c_{
m comp}}}$ 

#### When to use local SGD?

#### Recommended if:

- Communication cost  $c_{\text{comm}} \gg c_{\text{comp}}$
- Data across workers are relatively homogeneous
- Occasional synchronization suffices for convergence

#### Avoid if:

- Highly non-i.i.d. data (strong gradient heterogeneity)
- Models are unstable to small parameter changes



#### Table of Contents

Reduce the Communication Rounds via Local SGD

Reduce the Number of Bits via Quantization



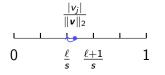
### Deterministic quantization

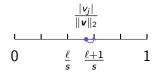
**Definition:** For any vector  $\mathbf{v} = [v_1, v_2, \dots, v_d]^{\top} \in \mathbb{R}^d$ , the j-th entry of the s-level quantized vector  $Q_s(\mathbf{v})$  is defined as:

$$[Q_s(\mathbf{v})]_j := \|\mathbf{v}\|_2 \cdot \operatorname{sign}(v_j) \cdot \zeta_j(\mathbf{v}, s),$$

Let  $0 \le \ell < s$  be an integer such that  $\frac{|v_j|}{\|v\|_2} \in \left[\frac{\ell}{s}, \frac{\ell+1}{s}\right]$ . Then:

$$\zeta_j(\mathbf{v},s) = \begin{cases} \frac{\ell}{s}, & \text{if } \frac{|\mathbf{v}_j|}{\|\mathbf{v}\|_2} - \frac{\ell}{s} \leq \frac{1}{2s}, \\ \frac{\ell+1}{s}, & \text{otherwise} \end{cases}$$





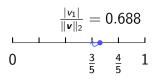


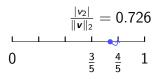
# Example: deterministic quantization (s = 5)

**Example:** Consider a 2-D vector  $\mathbf{v} = [0.36, \ 0.38]$ . Let s = 5. Its  $\ell_2$ -norm is  $\|\mathbf{v}\|_2 = \sqrt{0.36^2 + 0.38^2} \approx 0.523$ . Thus,

$$\frac{|v_1|}{\|\mathbf{v}\|_2} = 0.688, \qquad \frac{|v_2|}{\|\mathbf{v}\|_2} = 0.726.$$

Both values fall into the same quantization interval  $\left[\frac{3}{5}, \frac{4}{5}\right] = [0.6, 0.8].$ 





According to the rule:  $[Q_s(\mathbf{v})]_j = \|\mathbf{v}\|_2 \cdot \text{sign}(v_j) \cdot \zeta_j(\mathbf{v}, s)$ , we obtain:  $Q_5(\mathbf{v}) = 0.523 [0.6, 0.8] = [0.314, 0.418]$ .



## Loss of deterministic quantization

**Problem with this strategy:** Higher quantization error for values that are further away from the center of the interval.

**Lemma.** For any vector  $\mathbf{v} \in \mathbb{R}^d$ , we have:

(i) 
$$\|Q_s(oldsymbol{v}) - oldsymbol{v}\|_{\infty} \leq rac{1}{s} \|oldsymbol{v}\|_2$$
 (bias)

(ii) 
$$||Q_s(\mathbf{v}) - \mathbf{v}||_2^2 \le \frac{d^2}{s} ||\mathbf{v}||_2^2$$



## Stochastic quantization

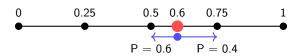
For any vector  $\mathbf{v} = [v_1, \dots, v_d]^{\top} \in \mathbb{R}^d$ , the j-th entry of the s-level quantized vector  $Q_s(\mathbf{v})$  is:

$$[Q_s(\mathbf{v})]_j := \|\mathbf{v}\|_2 \operatorname{sign}(v_j) \zeta_j(\mathbf{v}, s),$$

where the random variable  $\zeta_j(\mathbf{v},s)$  is:

$$\zeta_j(oldsymbol{v},s) = egin{dcases} rac{\ell+1}{s}, & ext{with probability } sigg(rac{|v_j|}{\|oldsymbol{v}\|_2} - rac{\ell}{s}igg)\,, \ rac{\ell}{s}, & ext{otherwise} \end{cases}$$

See example for s = 4 levels below:





### Lemma: properties of stochastic quantization

**Lemma:** For any vector  $\mathbf{v} \in \mathbb{R}^d$ , if we apply stochastic quantization  $Q_s(\mathbf{v})$ , then we have:

(i) Unbiasedness:

$$\mathbb{E}[Q_s(\mathbf{v})] = \mathbf{v}$$

(ii) Bounded variance:

$$\mathbb{E}\big[\|Q_s(\boldsymbol{\nu})-\boldsymbol{\nu}\|_2^2\big] \leq \min\!\left(\frac{d}{s^2},\frac{\sqrt{d}}{s}\right)\|\boldsymbol{\nu}\|_2^2$$

The proof of the second property is given in Appendix 1 of the QSGD paper https://arxiv.org/pdf/1610.02132.



# Convergence guarantees for QSGD: error bound

If  $Q(\mathbf{v}_i^{(t)})$  is an unbiased stochastic estimator of  $\nabla F_i(\mathbf{x}^t)$ , then the quantized update is equivalent to a stochastic gradient update, and the standard SGD analysis can be applied.

#### Theorem 2 (Convergence of QSGD)

Let f be L-smooth and  $\eta_k \equiv \eta = 1/\sqrt{K}$ . Then the following holds:

$$\frac{1}{K} \sum_{k=1}^K \mathbb{E} \big[ \| \nabla f(\boldsymbol{x}^k) \|_2^2 \big] \leq \mathcal{O} \left( \frac{\sigma}{\sqrt{nK}} \sqrt{1 + \min \left( \frac{d}{s^2}, \frac{\sqrt{d}}{s} \right)} \right).$$

■ The error versus iterations convergence becomes worse if we use fewer quantization levels *s*.



# Proof of QSGD error convergence bound

By combining the variance upper bound with the bounded estimation error property of the stochastic quantizer, we have:

$$\mathbb{E}_{Q}\big[\|Q(g(\boldsymbol{x};\boldsymbol{\xi})) - g(\boldsymbol{x};\boldsymbol{\xi})\|_{2}^{2}\big] \leq \min\left(\frac{d}{s^{2}}, \frac{\sqrt{d}}{s}\right)\|g(\boldsymbol{x};\boldsymbol{\xi})\|_{2}^{2}$$

$$\Rightarrow \quad \mathbb{E}_{Q}\big[\|Q(g(\boldsymbol{x};\boldsymbol{\xi}))\|_{2}^{2}\big] \leq \|g(\boldsymbol{x};\boldsymbol{\xi})\|_{2}^{2} + \min\left(\frac{d}{s^{2}},\frac{\sqrt{d}}{s}\right)\|g(\boldsymbol{x};\boldsymbol{\xi})\|_{2}^{2}$$

$$\mathbb{E}_{\xi}\big[\mathbb{E}_{Q}\big[\|Q(g(\boldsymbol{x};\boldsymbol{\xi}))\|_{2}^{2}\big]\big] \leq \mathbb{E}_{\xi}\big[\|g(\boldsymbol{x};\boldsymbol{\xi})\|_{2}^{2}\big] + \min\left(\frac{d}{s^{2}},\frac{\sqrt{d}}{s}\right)\mathbb{E}_{\xi}\big[\|g(\boldsymbol{x};\boldsymbol{\xi})\|_{2}^{2}\big]$$

$$\mathbb{E}\big[\|Q(g(\boldsymbol{x};\boldsymbol{\xi}))\|_2^2\big] \leq \left(1 + \min\left(\frac{d}{s^2}, \frac{\sqrt{d}}{s}\right)\right) \left(\|\nabla F(\boldsymbol{x})\|_2^2 + \sigma^2\right)$$



## Implementation of stochastic quantization

After quantization, we transmit  $Q_s(\mathbf{v})$  instead of the full vector  $\mathbf{v}$ .

The quantized vector  $Q_s(\mathbf{v})$  can be represented by the tuple:

$$Q_s(\mathbf{v}) = \left(\underbrace{\|\mathbf{v}\|_2}_{32 \text{ bits}}, \underbrace{\operatorname{sign}(v_j)_{j=1}^d}_{d \text{ bits}}, \underbrace{\zeta_j(\mathbf{v}, s)_{j=1}^d}_{d \log_2 s \text{ bits}}\right)$$

Total:

$$32 + d(1 + \log_2 s)$$
 vs.  $32d$  (full precision)

**Conclusion:** Quantization effectively reduces communication cost while introducing a moderate increase in the error versus during convergence.



### Recap and fine-tuning

- What we have talked about today?
  - ⇒ **Local SGD** reduces communication rounds by allowing each worker to perform multiple local updates before synchronization.
  - $\Rightarrow$  **Quantization** reduces communication cost by transmitting gradients with fewer bits.



Welcome anonymous survey!



