# **Delay Gradient Averaging: Tolerate the Communication Latency in Federated Learning**

Ligeng Zhu<sup>1</sup>, Hangzhou Lin<sup>2</sup>, Yao Lu<sup>3</sup>, Yujun Lin<sup>1</sup>, Song Han<sup>1</sup> <sup>3</sup>Google <sup>2</sup>Amazon <sup>1</sup>MIT

### **Federated Learning Allows Training** without Sharing



- Security: Data never leaves devices thus promises security and regularization.
- Customization: Models continually adapt to new data from the sensors.

## Challenge: Network Communication Bottleneck



Wired ethernet or infinity band Bandwidth as high as 100Gb/s, Latency as low as **1us** 



WiFi or Cellular network Bandwidth up to 1Gb/s, Latency ~200ms.

### **Existing Methods Improves the Bandwidth,** but not the Latency

- Bandwidth can be always improved by
  - Hardware upgrade
  - Gradient compression and quantization

### • Latency is hard to improve because

- Physical limits: Shanghai to Boston, even with the speed of light, still takes 162ms.
- Signal congestion: Urban office and home creates a lot of signal contention.











$$\nabla f_j(x) - \nabla f_j(y) || \le L ||x - y||. \quad \forall x, y \in \mathbb{R}^d$$

$$\|\nabla F_{j}(w;\zeta_{i})\|^{2} \leq G^{2}, \forall w, \forall j, \mathbb{E}_{\zeta_{j}}\|\nabla F_{j}(w;\zeta_{j}) - \nabla f_{j}(w)\|^{2} \leq \sigma^{2}, \forall w, \forall j.$$

ts	Partition	FedAvg (K=5)		FedAvg (K=10)		FedAvg (K=20)		DGA (K=5, D=20)	
		Acc	Speedup	Acc	Speedup	Acc	Speedup	Acc	Speedup
R	i.i.d non-i.i.d	88.7 48.2	$1 \times$	88.5 47.2	1.51×	88.1 43.9	2.05  imes	88.6 48.0	3.16×
let	i.i.d non-i.i.d	76.6 55.4	$1 \times$	76.5 52.5	1.43×	76.2 48.6	1.81×	76.4 54.9	2.55×
eare	i.i.d non-i.i.d	47.6 36.9	$1 \times$	47.3 34.3	1.66×	47.4 30.1	2.51×	47.1 36.3	4.07×