

Distributed Learning over Unreliable Networks

Chen Yu, Hanlin Tang, Cedric Renggli, Simon Kassing, Ankit Singla, Dan Alistarh, Ce Zhang, Ji Liu

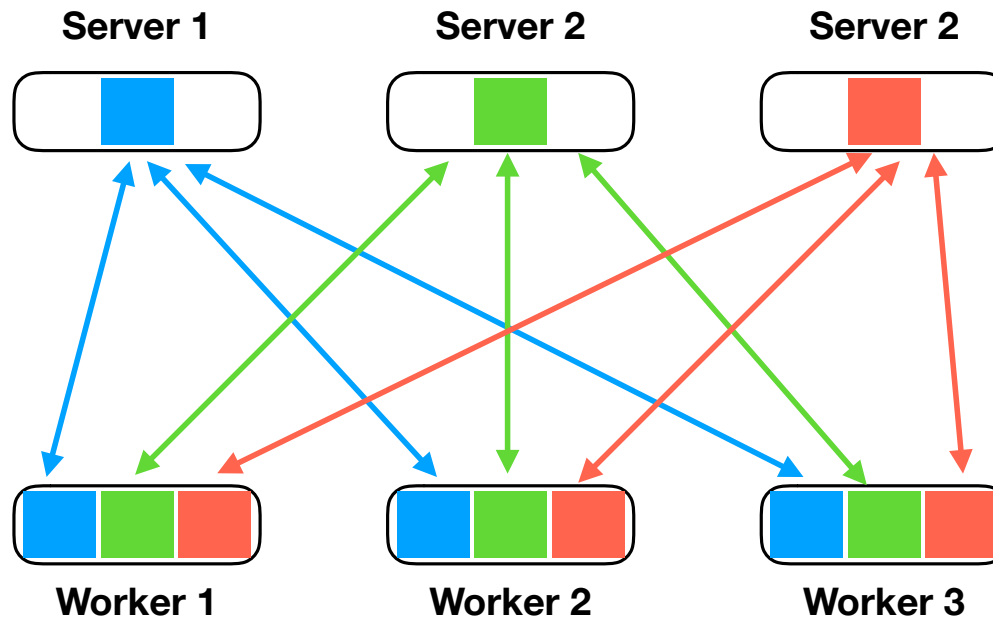
Presenter: Chen Yu



ETH zürich



AllReduce SGD

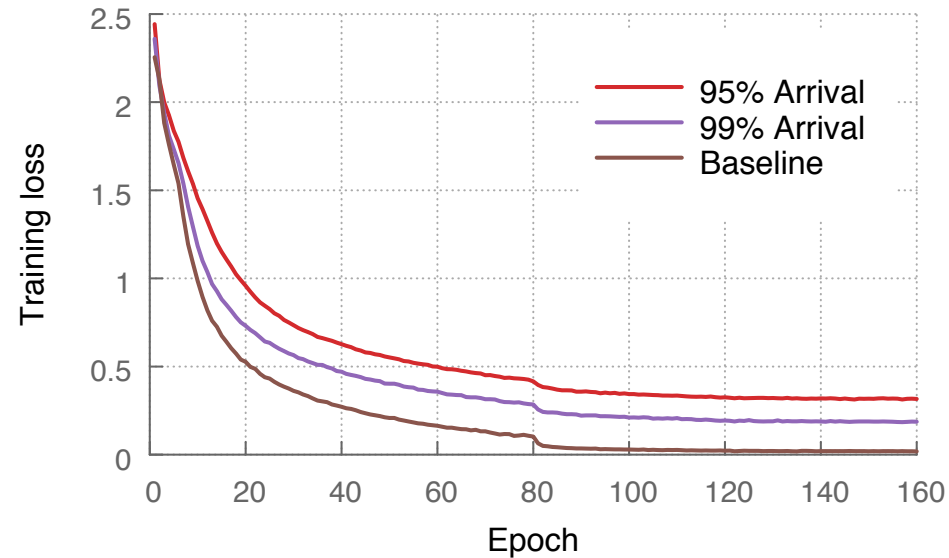
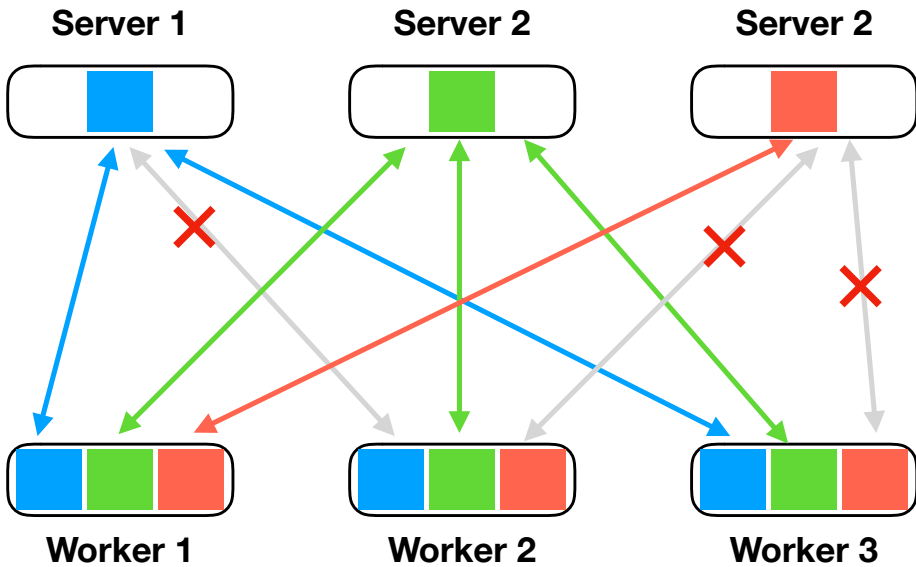


$$x_{t+1} = x_t - \frac{1}{n} \sum_{i=1}^n \nabla F(x_t; \xi_t^{(i)})$$



AllReduce

Unreliable Network



Sharing Gradients **Won't** Work

Reliable Parameter Server (**RPS**)

High Level: Share **Models**

Local Partition:

$$v_t^{(i)} = x_t^{(i)} - \gamma g_t^{(i)}, \quad v_t^{(i)} = \left((v_t^{(i,1)})^\top, (v_t^{(i,2)})^\top, \dots, (v_t^{(i,n)})^\top \right)^\top.$$

Robust Averaging:

$$\tilde{v}_t^{(i)} = \frac{1}{|\mathcal{N}_t^{(i)}|} \sum_{j \in \mathcal{N}_t^{(i)}} v_t^{(i,j)}$$

Model Update:

$$x_{t+1}^{(i,j)} = \begin{cases} \tilde{v}_t^{(j)} & j \in \tilde{\mathcal{N}}_t^{(i)} \\ v_t^{(i,j)} & j \notin \tilde{\mathcal{N}}_t^{(i)} \end{cases}.$$

Convergence Rate

Assumptions:

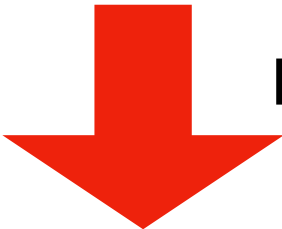
$f(x)$ **Non Convex, with L-Lipschitz Gradient;**

$\mathbb{E}_{\xi \sim \mathcal{D}_i} \|\nabla F(x; \xi) - \nabla f_i(x)\|^2 \leq \sigma^2, \forall i, \forall x;$ **Bounded Data Variance**

$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x) - \nabla f(x)\|^2 \leq \zeta^2, \forall i, \forall x.$ **Bounded Dataset Difference**

T: Total Iterations

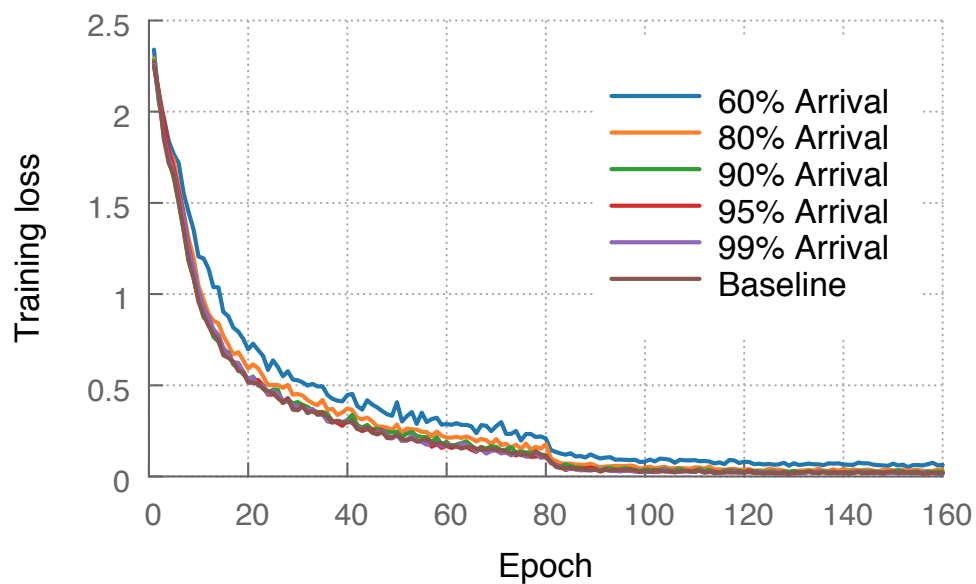
p: Package Dropping Rate


$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla f(\bar{x}_t)\|^2 \lesssim \frac{(\sigma + \zeta) \left(1 + \sqrt{p(1-p)}\right)}{(1 - \sqrt{p})\sqrt{nT}} + \frac{1}{T}$$

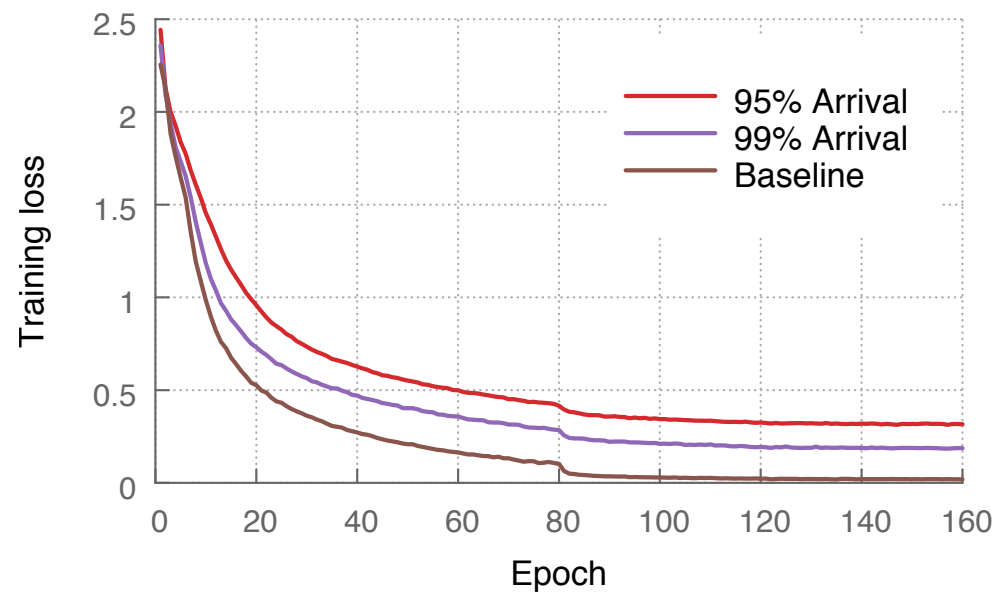
Experiments

16 NVIDIA TITAN Xp GPUs, ResNet-110 on CIFAR-10

RPS is **Robust**



Standard SGD is **Vulnerable**



Thanks

Welcome to Pacific Ballroom #97 to see the poster for
more detail