



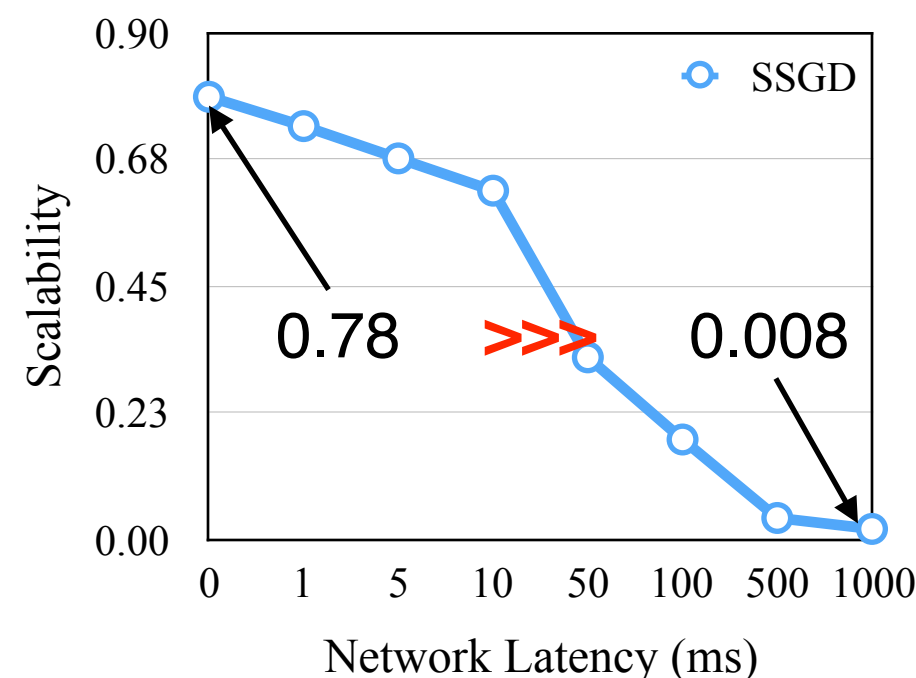
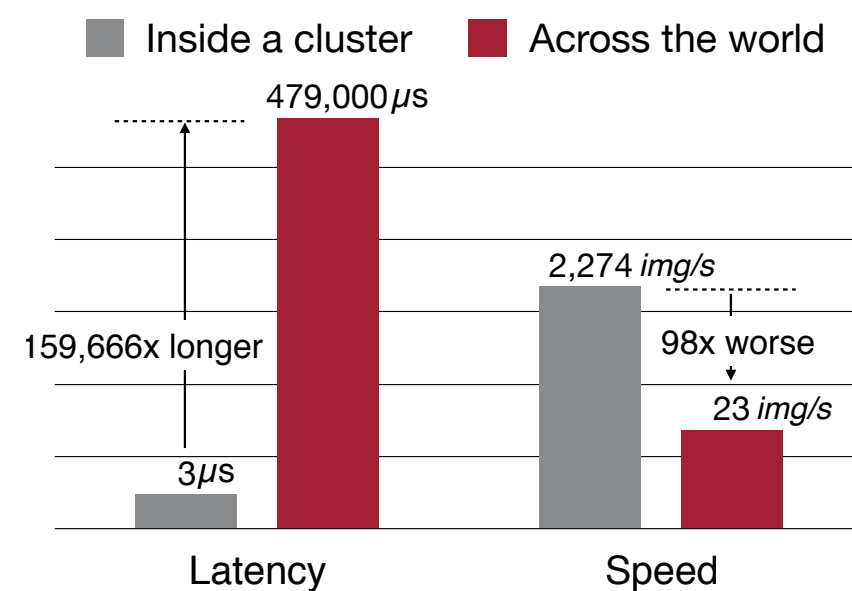
Training inside a cluster v.s. across world

Traditional distributed training is performed inside a cluster because it requires high end networking infrastructure.

But in many cases, *the data are distributed in different geographical locations and cannot be centralized*. For example, medical history, keyboard inputs.

Due to the long physical distance, high latency cannot be avoided and traditional algorithm scales poorly under such conditions.

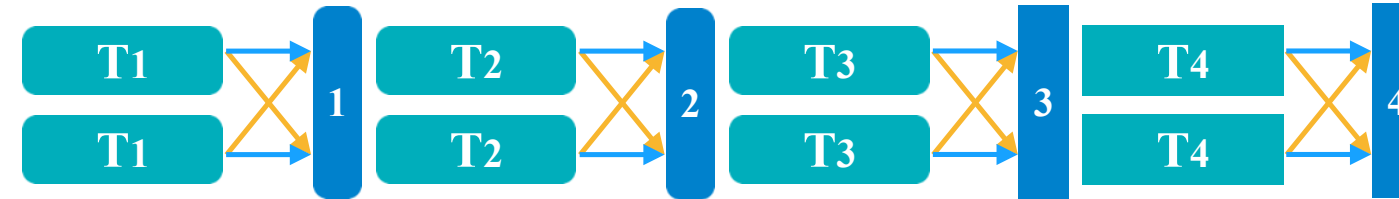
In this work, we aim to **scale Synchronous SGD across the world without loss of speed and accuracy!**



Vanilla SGD

For iteration $n = 0, 1, \dots$, on j^{th} worker

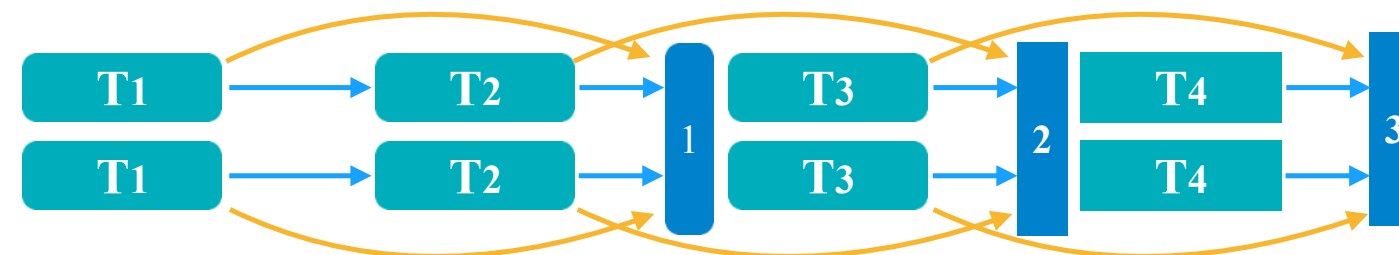
$$\overline{w}_{(n,j)} = \overline{w}_{(n-1)} - \gamma \overline{\nabla w}_{(n-1)}$$



Delayed Update

For iteration $n = 0, 1, \dots$, on j^{th} worker, when $n > t$

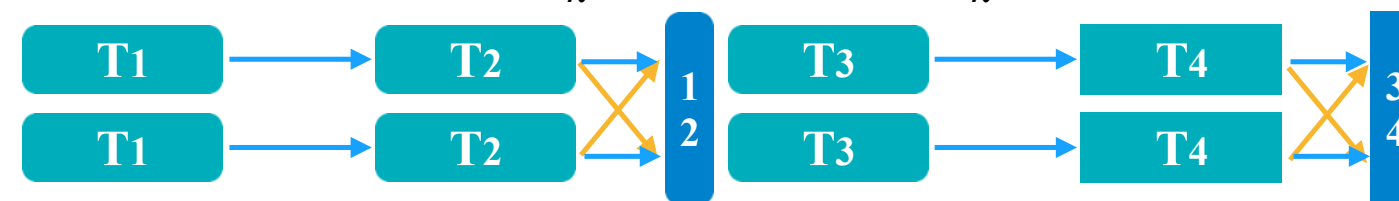
$$w_{(n,j)} = w_{(n-1,j)} - \gamma (\nabla w_{(n-1,j)} + \overline{\nabla w}_{(n-t)} - \nabla w_{(n-t,j)})$$



Temporally Sparse Update

For iteration $n = 0, 1, \dots$, on j^{th} worker, when $n \equiv 0 \pmod{n}$

$$w_{(n,j)} = w_{(n-1,j)} - \gamma (\nabla w_{(n-1,j)} + \sum_k^d \nabla w_{(n-t+k)} - \sum_k^d \nabla w_{(n-t+k,j)})$$

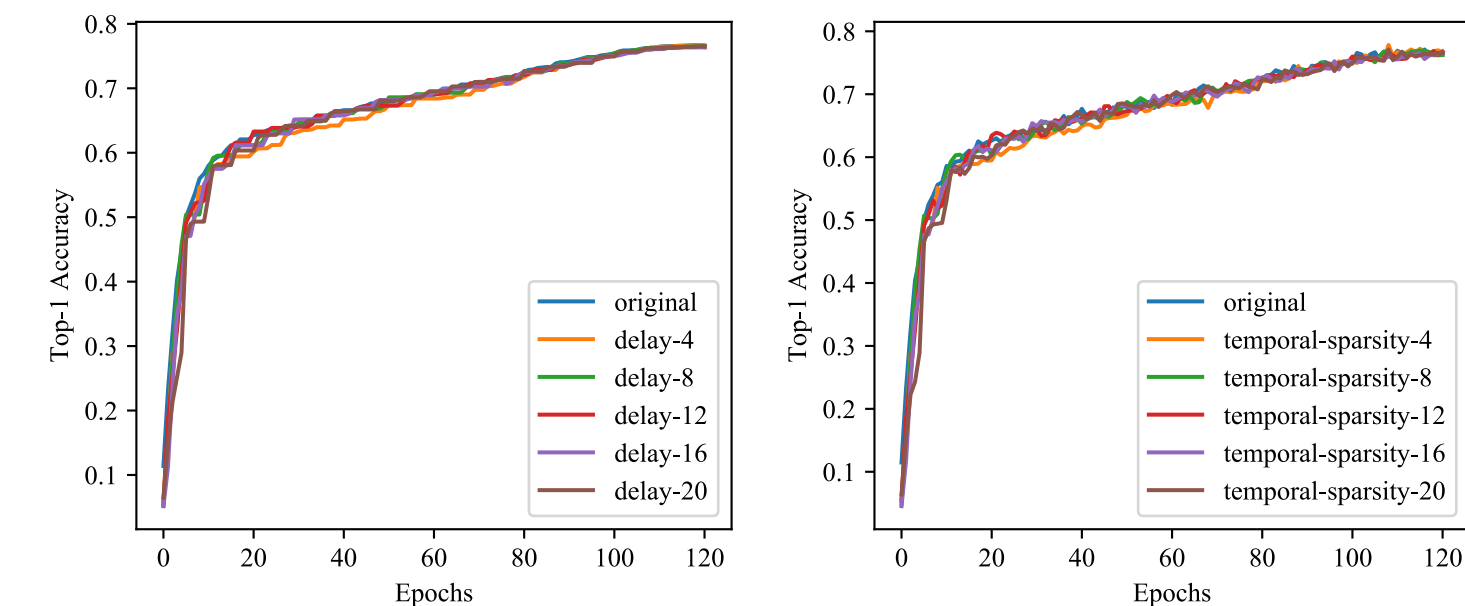
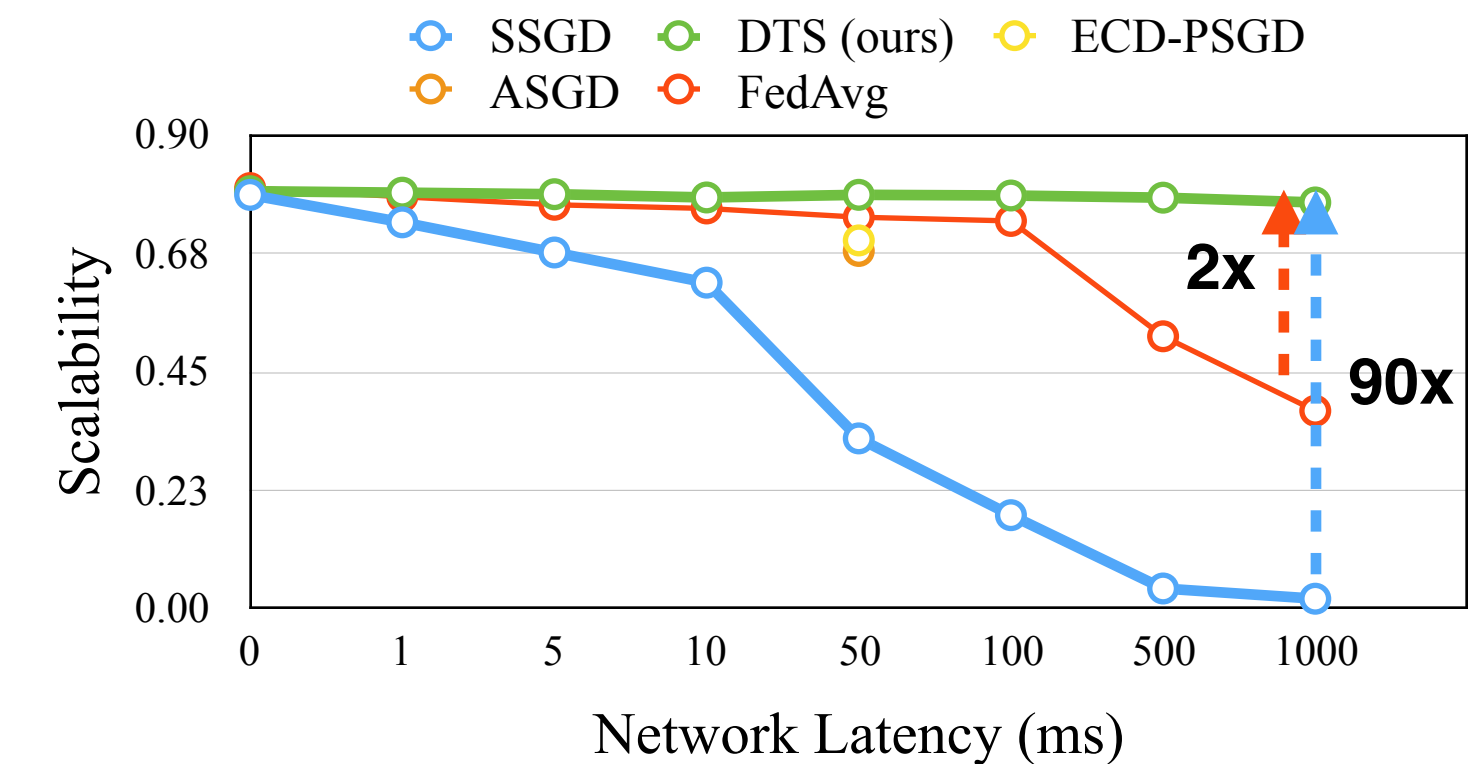


Theoretical Results

Convergence: $O(\frac{1}{\sqrt{NJ}}) + O(\frac{(t+d)^2 J}{N})$

$c < O(N^{\frac{1}{4}} J^{-\frac{3}{4}}) \rightarrow$ no slower than SGD

Empirical Results



On AWS, eight 8-V100 p3.16x instance located in *Oregon, Ohio, London* and *Tokyo*. Our algorithm demonstrates scalability of 0.72 while original SGD suffers at scalability of 0.008 (**90x**).