

Moshpit SGD: Communication-Efficient Decentralized Training on Heterogeneous Unreliable Devices

Max Ryabinin*, Eduard Gorbunov*,
Vsevolod Plokhotnyuk, Gennady Pekhimenko

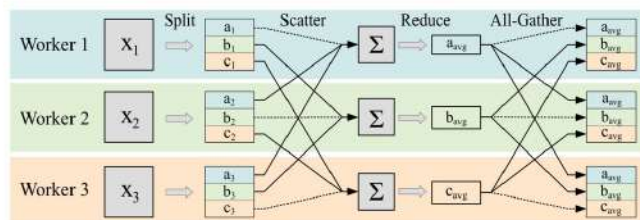


The Paper in Brief

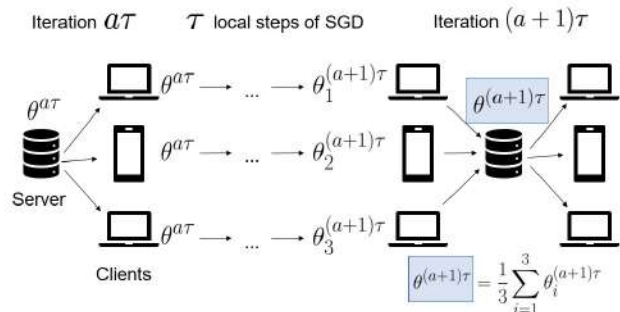
- We develop a simple method for distributed training with unstable (i.e. frequently joining and leaving) workers
- It combines the communication efficiency of All-Reduce with the fault tolerance of Gossip-based methods
- Has strong theoretical guarantees both for convergence to the actual average and stochastic optimization
- In practice, allows distributed training on preemptible instances and outperforms standard data-parallel training at a fraction of cost for ResNet and ALBERT

Background: Data-Parallel Training

- Most popular approach to distributed training: split batches across devices, average gradients, run the SGD step
- Naively sending all gradients is slow; more efficient versions (Ring, Butterfly) are used in practice
- However, standard All-Reduce fails in congested/high-latency networks
- Gossip fares better, but loses efficiency (sends all data to each neighbor)

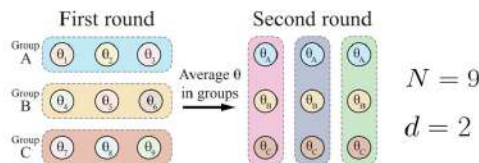


Background: Local SGD



Moshpit All-Reduce

Special case: exact averaging in d steps for a full "grid"



General case:

Algorithm 1 Moshpit All-Reduce (for i -th peer)

```

Input: parameters  $\{\theta_j\}_{j=1}^N$ , number of peers  $N$ ,  $d$ ,
 $M$ , number of iterations  $T$ , peer index  $i$ 
 $\theta_i^0 := \theta_i$ 
 $C_i^0 := \text{get\_initial\_index}(i)$ 
for  $t \in 1, \dots, T$  do
    DHT( $C_i^{t-1}, t$ ), add(address $_i$ )
    Matchmaking() // wait for peers to assemble
    peers $_t := \text{DHT.get}(C_i^{t-1}, t)$ 
     $\theta_i^t, c_i^t := \text{AllReduce}(\theta_i^{t-1}, \text{peers}_t)$ 
     $C_i^t := (C_i^{t-1} [1:], c_i^t)$  // same as eq. (1)
end for
Return  $\theta_i^T$ 
    
```

- Group key: $C_i^t := (c_i^{t-d+1}, c_i^{t-d+2}, \dots, c_i^t)$
- get_initial_index(i) = $(\lfloor i/M^{d-1} \rfloor \bmod M)_{j \in \{1, \dots, d\}}$
- Key property: if two peers were in the same group in round t , they choose different groups in round $t+1$

Theoretical guarantees:

- Correctness:** If all workers have a non-zero probability of successfully running a communication round and the order of peers, is random, then all local vectors converge to the global average with probability 1:

$$\forall i, \left\| \theta_i^t - \frac{1}{N} \sum_i \theta_i^0 \right\|_2 \xrightarrow[t \rightarrow \infty]{} 0$$

- Exponential convergence to the average:** for a version of Moshpit All-Reduce with random splitting into r groups at each step, we have

$$\mathbb{E} \left[\frac{1}{N} \sum_{i=1}^N \|\theta_i^t - \bar{\theta}\|^2 \right] = \left(\frac{r-1}{N} + \frac{r}{N^2} \right)^T \frac{1}{N} \sum_{i=1}^N \|\theta_i - \bar{\theta}\|^2$$

Moshpit SGD & Its Convergence

Moshpit SGD = Local SGD + Moshpit All-Reduce

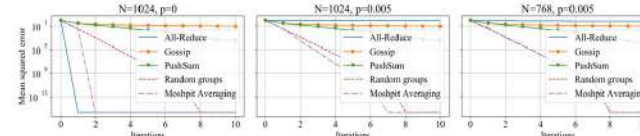
Theoretical guarantees: under the standard assumptions of bounded variance of stochastic gradients, reasonable number of iterations of Moshpit All-Reduce, and the bounded effect of peers' vanishing we recover:

- The best known rates from (Khaled et al., 2020; Woodworth et al., 2020) in **convex and strongly convex** cases
- The best known rates from (Koloskova et al., 2020; Li et al., 2019) in the **non-convex** case

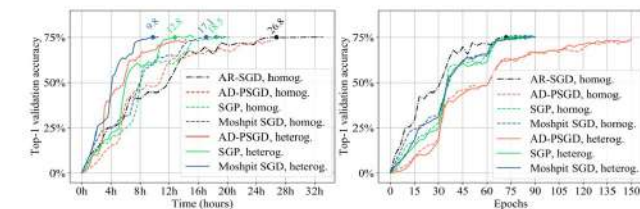
Experiments

Averaging performance

- We compare per-iteration convergence in a simulated setup
- All-Reduce takes too long to average with non-zero failure probability
- Gossip/SGP converge much slower (>10 iterations for target precision)



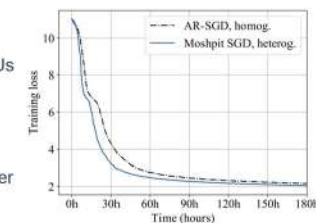
ResNet-50 on ImageNet



- We evaluate Moshpit-SGD and several baselines in two environments
- (16 nodes with 1xV100 and 64 workers with 81 different GPUs)
- Comparable to All-Reduce in terms of iterations, faster in terms of time
- Decentralized methods run faster, but achieve worse results

ALBERT on BookCorpus

- Baseline: All-Reduce on 8 V100
- Moshpit SGD: 66 preemptible GPUs
- For standard DDP, latency and failures make it impossible to run
- Costs of spot instances are much smaller, yet we converge 1.5x faster



References

Ahmed Khaled, Konstantin Mishchenko, and Peter Richtárik. Tighter Theory for Local SGD on Identical and Heterogeneous Data. In IAISTATS, pages 4519–4529. PMLR, 2020.

Blake Woodworth, Kumar Kshitij Patel, Sebastian Stich, Zhen Dai, Brian Bullins, Brendan McMahan, Ohad Shamir, and Nathan Srebro. Is Local SGD Better Than Minibatch SGD? In ICML, pages 10334–10343. PMLR, 2020.

Anastasia Koloskova, Nicolas Loizou, Sadra Boreiri, Martin Jaggi, and Sebastian Stich. A Unified Theory of Decentralized SGD with Changing Topology and Local Updates. In ICML, pages 5381–5393. PMLR, 2020.

Xiang Li, Wenhao Yang, Shusen Wang, and Zhihua Zhang. Communication Efficient Decentralized Learning with Multiple Local Updates. arXiv:1910.09126, 2019.

Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, Radu Soricut. ALBERT: A Lite BERT for Self-supervised Learning of Language Representations. ICML, 2020.