

# A Quadratic Synchronization Rule for Distributed Deep Learning



Xinran Gu<sup>\*,1</sup> Kaifeng Lyu<sup>\*,2</sup> Sanjeev Arora<sup>2</sup> Jingzhao Zhang<sup>1</sup> Longbo Huang<sup>1</sup>

<sup>1</sup>Tsinghua University <sup>2</sup>Princeton University

## Abstract

Local gradient methods, e.g., Local SGD, improve the communication efficiency of data parallel training by letting workers communicate **only every  $H$  steps**.

### • How to set the synchronization period $H$ ?

- Optimization: communication & convergence tradeoff
- Generalization: proper  $H \Rightarrow$  higher test acc. (Lin et al., 2020)
- We propose a theory-grounded strategy to set  $H$

### Quadratic Synchronization Rule (QSR)

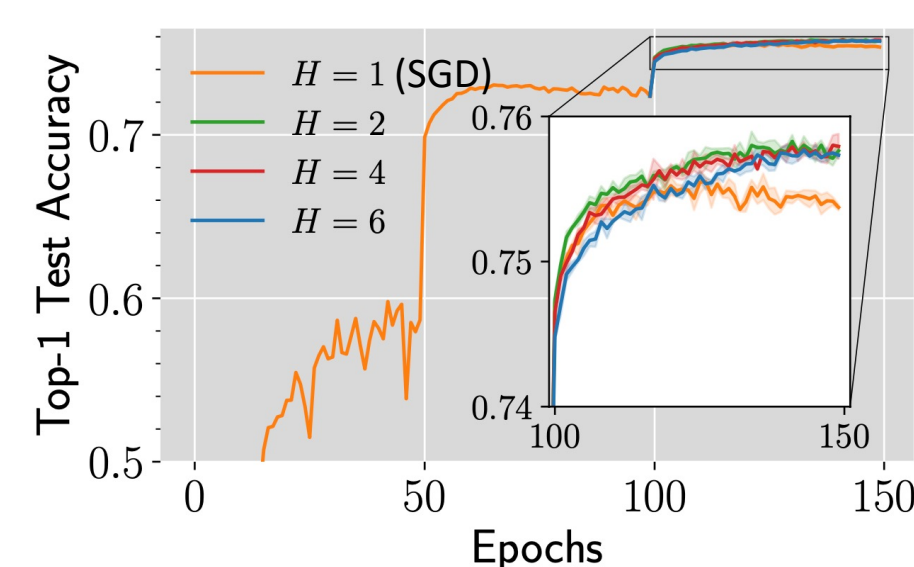
$$H \sim \eta^{-2} \quad (\eta: \text{learning rate})$$

Improve comm. efficiency & test acc. simultaneously!

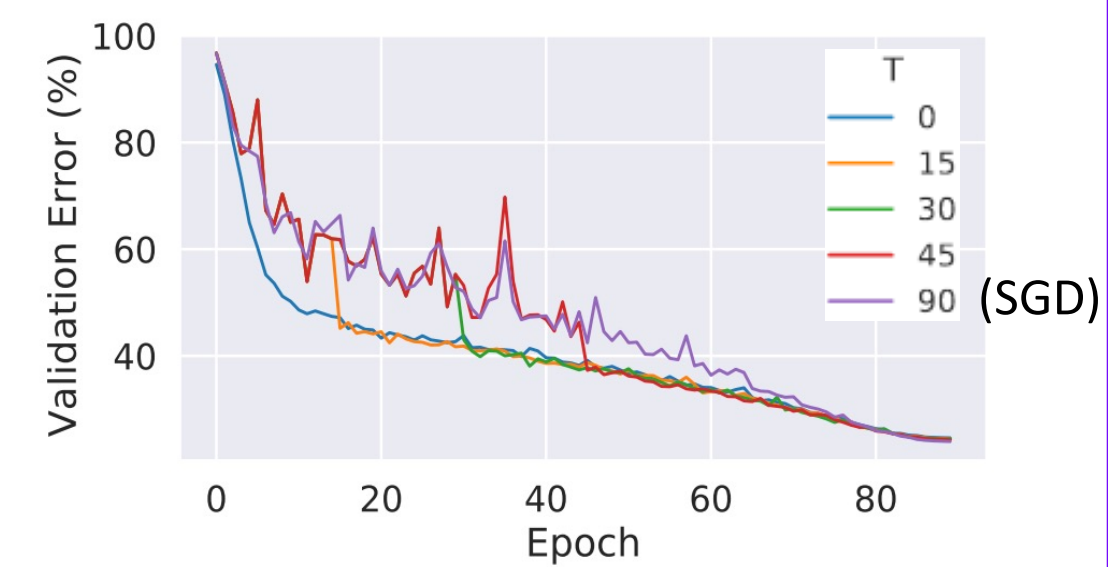
	time	val. acc.
data parallel	26.7h	79.86%
QSR	20.2h	80.98%

save 7h, improve 1%

Setting: 300 epoch on ViT-B, ImageNet



(a) Constant LR after switching



(b) Cosine LR decay (Ortiz et al., 2021)

ImageNet, ResNet-50

Issue: short-term generalization benefits on cos decay (Ortiz et al., 2021)

## Background: Local Gradient Methods

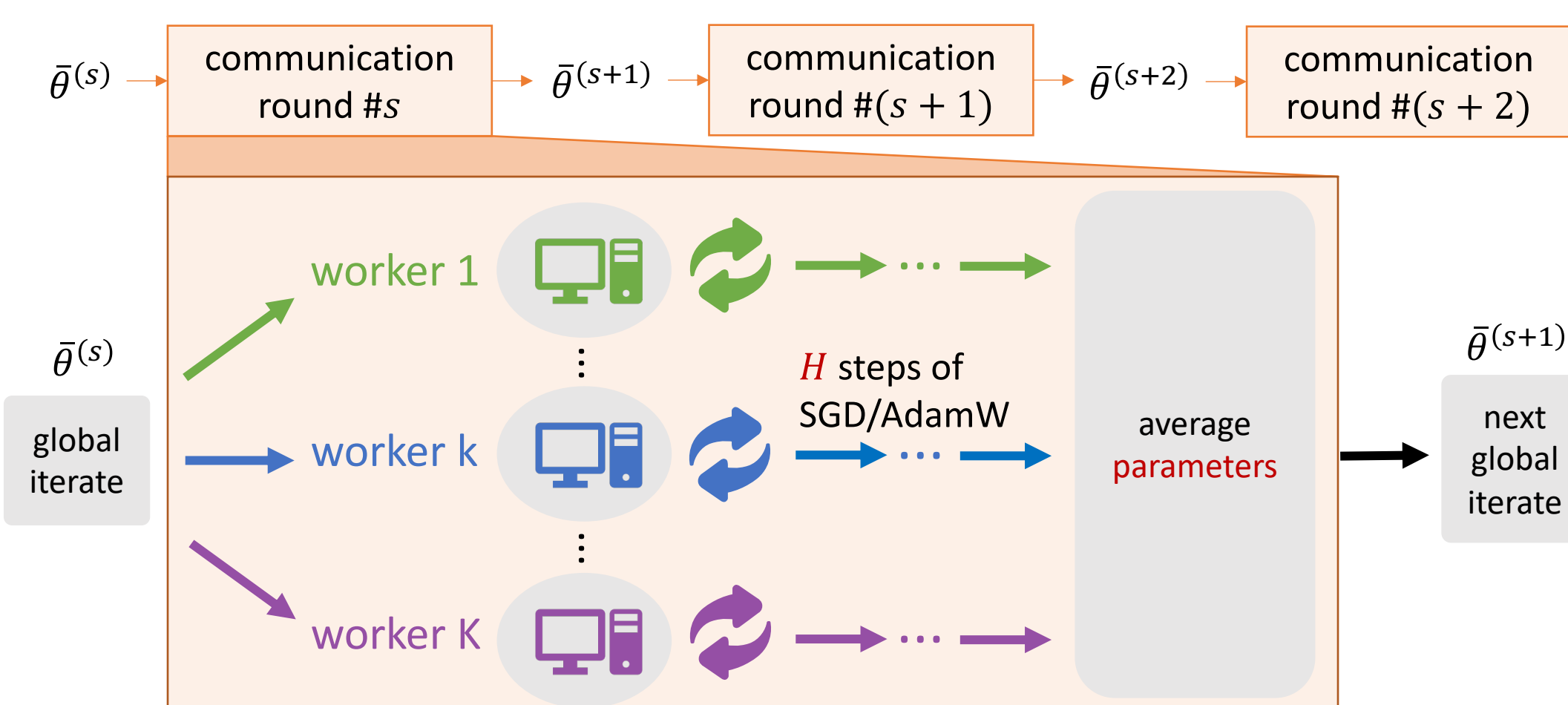
### • Data parallel approach

- Distribute gradient computation on  $B$  samples to  $K$  workers
- Each iteration, each worker:
  1. compute gradients on  $B/K$  samples
  2. average gradients via All-Reduce
  3. update using the averaged gradient & optimizer OPT

Issue: frequent sync.  $\Rightarrow$  high comm. cost

### • Local gradient methods

- Each worker locally updates its own replica with OPT
- Average model parameters every  $H$  steps



## Generalization Benefits of Local SGD

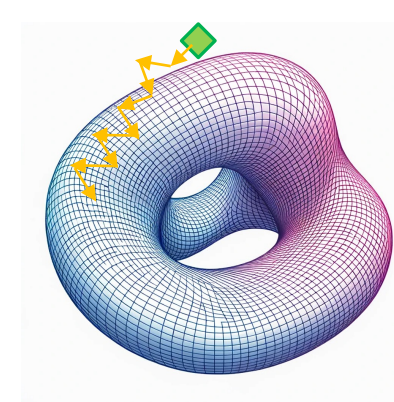
### • Local steps improve generalization (Lin et al., 2020)

- Run #1: Parallel SGD ( $\equiv$  Local SGD with  $H = 1$ )
- Run #2: Same as #1 but switch to Local SGD with  $H > 1$  at some epoch  $t_0$ , named "Post-local SGD"
- Result: test acc. #2 > #1

## Theory: Why does Local SGD Generalize Better?

### • Setting (Follow Blanc et al., 2020; Damian et al., 2021; Li et al., 2022)

- Assume a minimizer manifold  $\Gamma$
- Assume a smaller LR  $\eta$
- Analyze dynamics of (Local) SGD near  $\Gamma$

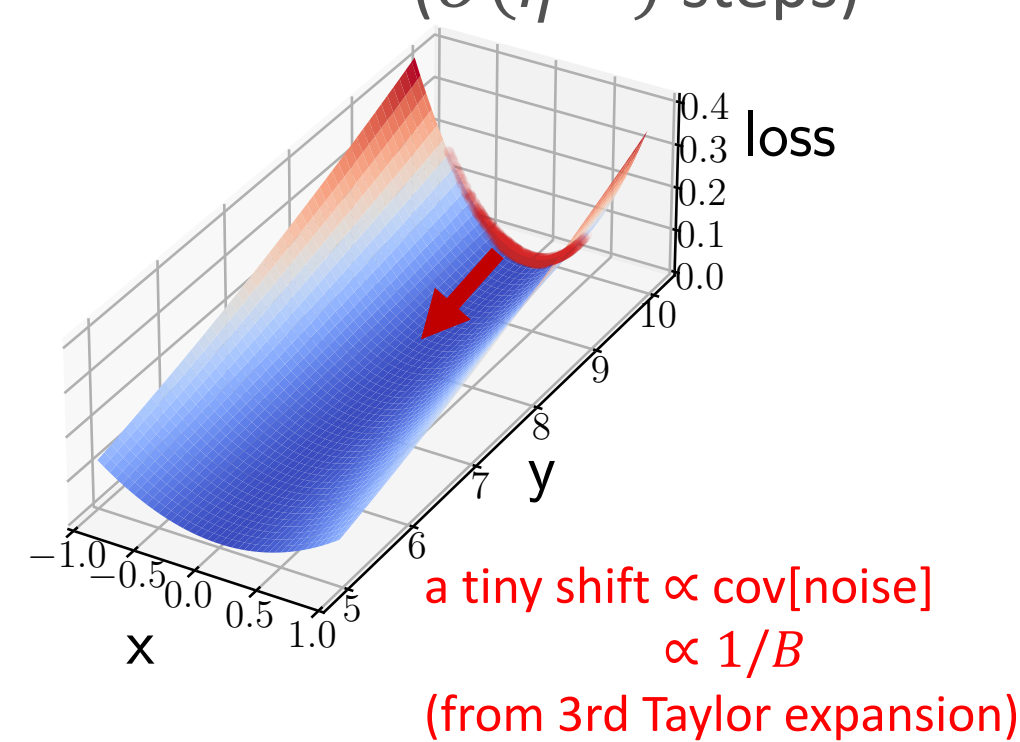


### • Fast and slow dynamics in SGD

(Blanc et al., 2020; Damian et al., 2021; Li et al., 2022)

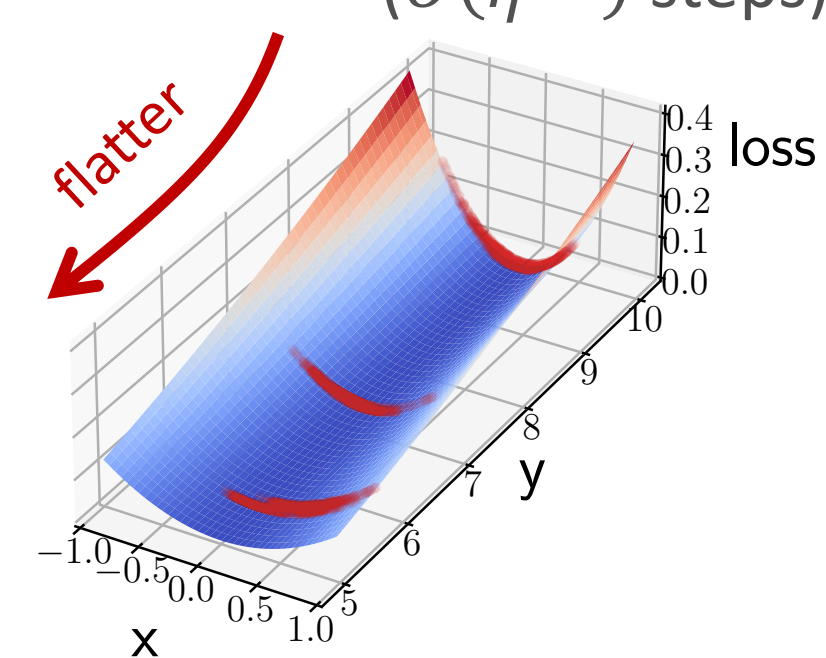
#### Fast Dynamics (short term)

Diffuse locally near a minimizer ( $O(\eta^{-1})$  steps)



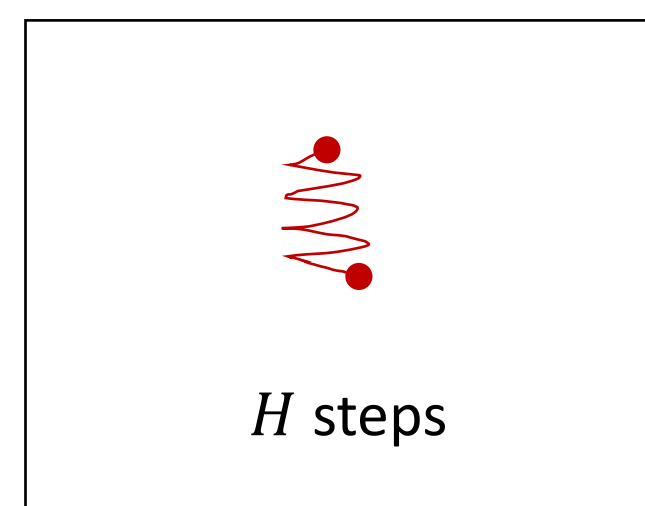
#### Slow Dynamics (long term)

"Center" of the diffusion shifts ( $O(\eta^{-2})$  steps)



### • Local SGD drifts faster to flatter minima

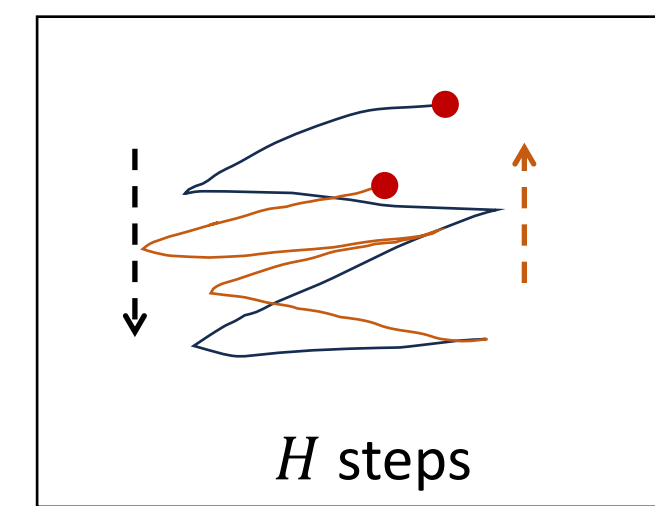
SGD, batch size  $B$



$H$  steps

drift slowly

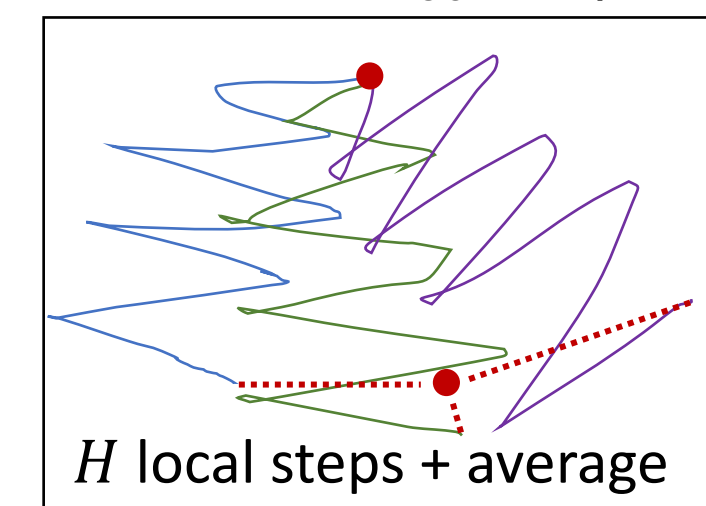
SGD, batch size  $B/K$



$H$  steps

drift fast in expectation, but go back and forth (large var.)

Local SGD,  $B_{loc} = B/K$



$H$  local steps + average

drift fast in expectation, averaging reduces var.

### • SDE approximations for different scalings of $H$

**Theorem (informal).** For  $O(\eta^{-2})$  steps, Local SGD with different scalings of  $H$  can be approximated by the following SDEs on  $\Gamma$ :

1.  $H = \beta/\eta$  (Gu et al., 2023)

$$d\zeta(t) = P_{\zeta} \left( \underbrace{\frac{1}{\sqrt{B}} \Sigma_{\parallel}^{1/2}(\zeta) dW_t - \frac{1}{2B} \nabla^3 \mathcal{L}(\zeta) [\hat{\Sigma}_{\diamond}(\zeta)] dt}_{\text{Same as SGD (Li et al., 2022)}} - \underbrace{\frac{K-1}{2B} \nabla^3 \mathcal{L}(\zeta) [\hat{\Psi}(\zeta)] dt}_{\text{Unique drift term of Local SGD}} \right)$$

-  $\hat{\Psi}(\zeta)$  increases with  $H$ , goes to 0 as  $H\eta \rightarrow 0$  and goes to  $\hat{\Sigma}_{\diamond}(\zeta)$  as  $H\eta \rightarrow \infty$

2.  $H = (\alpha/\eta)^2$  (our new result)

$$d\zeta(t) = P_{\zeta} \left( \frac{1}{\sqrt{B}} \Sigma_{\parallel}^{1/2}(\zeta) dW(t) - \underbrace{\frac{K}{2B} \nabla^3 \mathcal{L}(\zeta) [\hat{\Sigma}_{\diamond}(\zeta)] dt}_{K \text{ times of SGD; Local SGD with } H = \beta/\eta \text{ when } \beta \rightarrow \infty} \right)$$

$H \sim \eta^{-1}$  to see the benefit,  $H \sim \eta^{-2}$  to maximize it!

Cannot find valid SDE approximation on the manifold for more aggressive scalings.