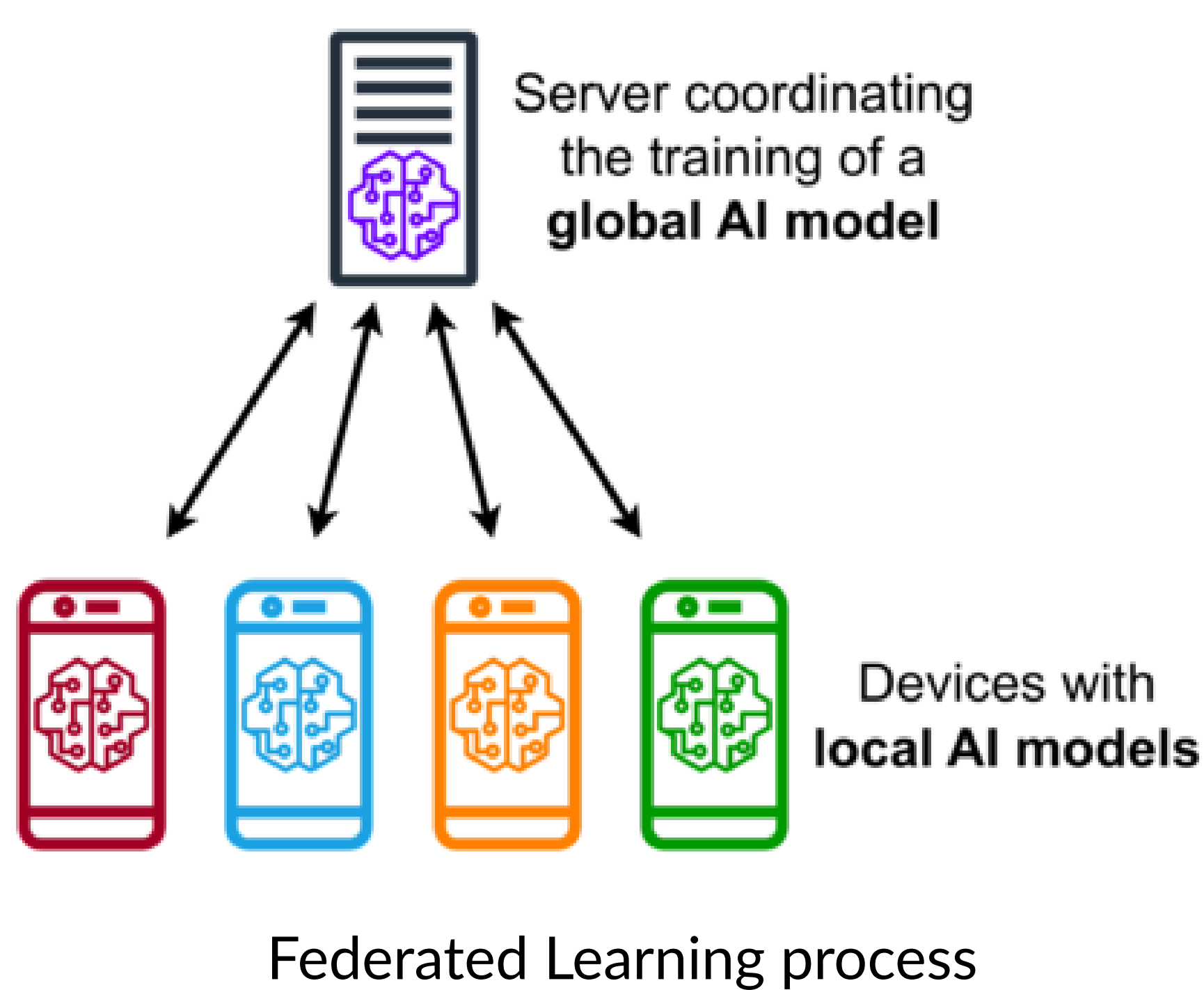


# Local SGD converges faster for quadratic-like objectives and requires less communication.

## Motivation and Challenges

- ▶ Larger models need data and tasks to be shared across many ( $M$ ) devices.
- ▶ Devices calculate local stochastic gradients and transmit them to a central server.
- ▶ Transmitting large amounts of data is costly.
- ▶ Our goal: **Reduce** the number of **communication** rounds.

We denote the concept above as "Federated Learning"



## Local SGD

- ▶ The most popular Federated Learning method is called **Local SGD**.
- ▶ It performs multiple local SGD steps between communications.
- ▶ Problem: if we reduce the number of communications and increase the number of local steps ( $H$ ), the performance **degrades**.

Woodworth et al., 2020 noted the following: For **quadratic** objectives, Local SGD convergence rate **is not affected** by the number of local steps, making it **highly efficient** for such problems.

## But what happens when we diverge from pure quadratic setting?

Thus, our aim was to establish better communication complexity rates for objectives somehow **close to the quadratic form**.

In order to measure the proximity of an objective  $F$  to the quadratic form, we decompose  $F$  into the sum:  $F = Q + R$ , where  $Q$  is a convex quadratic function, and  $R$  is some convex residue.

Then we introduce was *quadraticity* parameter  $\epsilon := \frac{L^2 R}{L^2} \leq 1$ .

- ▶ For quadratic objectives, where  $F$  is equal to  $Q$ , **the value of  $\epsilon$  is zero**
- ▶ For quadratic-like objectives, i.e. cases where  $Q$  is somewhat larger than  $R$ ,  $\epsilon$  is **small**.

## Quadraticity concept allows us to improve over the previous lower bounds for Local SGD

### Breaking existing bounds

Under the assumption of *uniformly bounded variance*, when  $E \|\nabla F(x) - \nabla F(x, z)\|^2 \leq \sigma^2$ :

$$E[F(x_T) - F(x_*)] = O\left(\frac{LD^2}{T} + \frac{\sigma D}{\sqrt{MT}} + \left(\frac{HL\sigma^2 D^4}{T^2}\right)^{1/3}\right) \text{ [Woodworth et al., 2020]}$$

↓

$$O\left(\frac{LD^2}{T} + \frac{\sigma D}{\sqrt{MT}} + \left(\frac{\epsilon HL\sigma^2 D^4}{T^2}\right)^{1/3}\right) \text{ [This work]}$$

If we denote  $\lambda = \mu_Q + \mu_R$  we can also get an estimate for the case  $\lambda > 0$ :

$$\tilde{O}\left(\text{exp.decay} + \frac{\sigma^2}{\mu MT} + \frac{HL\sigma^2}{\mu^2 T^2}\right) \text{ [Woodworth et al., 2020]}$$

↓

$$\tilde{O}\left(\text{exp.decay} + \frac{\sigma^2}{\lambda MT} + \frac{\epsilon HL\sigma^2}{\lambda^2 T^2}\right) \text{ [This work]}$$

### Abandoning restrictive assumption

If we replace uniformly bounded variance assumption with more **general** one, i.e.

$$E \|\nabla F(x) - \nabla F(x, z)\|^2 \leq \sigma^2 + \rho \|\nabla F(x)\|^2$$

the acceleration given by quadraticity **persists**.

Case  $\lambda > 0$ :

$$E[F(x_T) - F(x_*)] = O\left(\text{exp.decay} + \frac{\sigma^2}{\lambda MT} + \frac{\rho HL^2 \sigma^2}{\lambda^3 MT^3} + \frac{\epsilon HL\sigma^2}{\lambda^2 T^2}\right)$$

Variance reduction term

Represents the impact of  $\rho \|\nabla F(x)\|^2$

Represents the **drift** caused by rare communication

In all the estimates above, the last term represents **drift** that appears due to many local steps (or rare communication, which is equivalent)

So, when it is multiplied by the  $\epsilon$  factor it shows that **the influence of rare communications weakens for quadratic functions**.

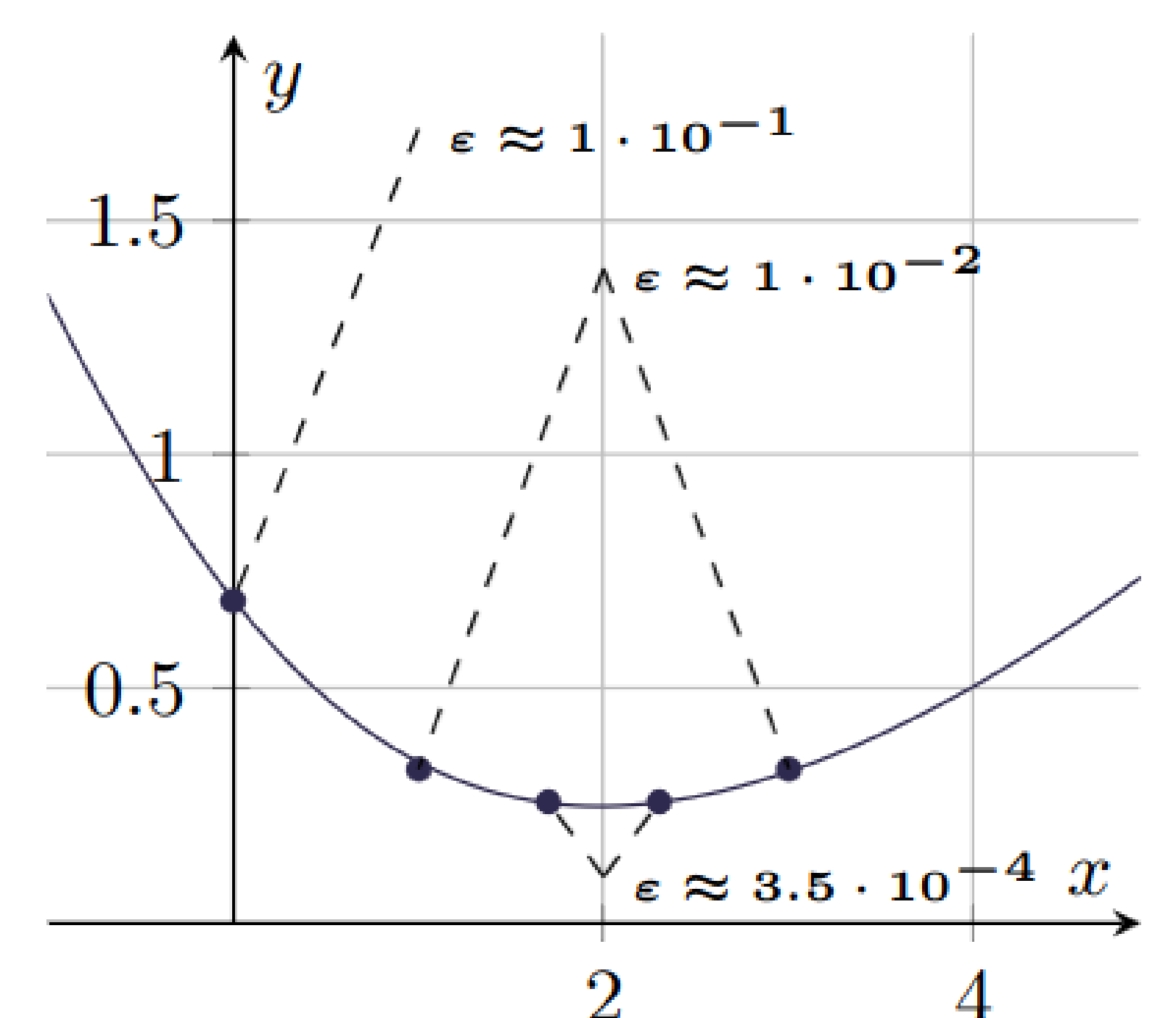
## Notation

The following symbols and definitions are used throughout this work:

Symbol	Definition
$M$	Number of devices
$H$	Number of local SGD steps
$T$	Total number of iterations on a given device
$D$	Initial distance to the optimum, $\ x_0 - x_*\ $
$\mu$	Strong convexity constant
$L$	Lipschitz gradient constant

## Discussion

An important observation about quadraticity is that for functions with a **Lipschitz Hessian**,  $\epsilon$  decreases rapidly, as illustrated in the graph below.



Decrease of  $\epsilon$  for LogLoss with  $l_2$  regularization

## References

Woodworth, B., Patel, K. K., Stich, S. U., Dai, Z., Bullins, B., McMahan, H. B., Shamir, O., & Srebro, N. (2020). Is local sgd better than minibatch sgd?

