

# Convergence and Linear Speed-Up in Stochastic Federated Learning

Paul Mangold

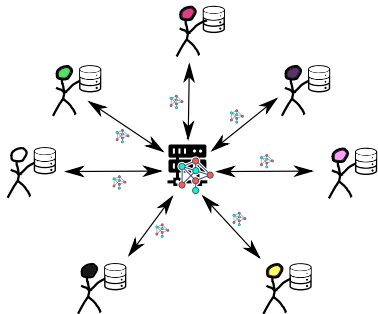
Workshop Fondation Mathématiques de l'IA

March 25th, 2025

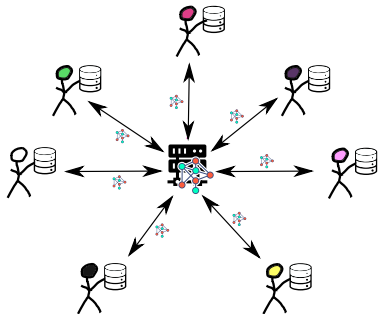
# Federated Learning



# Federated Learning



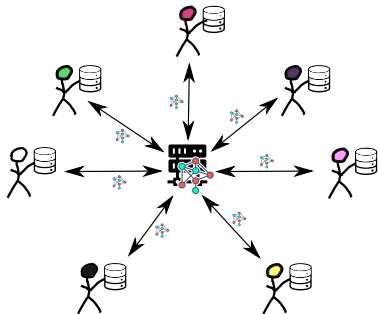
# Federated Learning



Collaborative optimization problem

$$\min_{x \in \mathbb{R}^d} \frac{1}{N} \sum_{c=1}^N f_c(x) \quad , \quad f_c(x) = \mathbb{E}_Z[F_c(x; Z)]$$

# Federated Learning



Collaborative optimization problem

$$\min_{x \in \mathbb{R}^d} \frac{1}{N} \sum_{c=1}^N f_c(x) \quad , \quad f_c(x) = \mathbb{E}_Z[F_c(x; Z)]$$

**Problem: data is heterogeneous, communication is expensive**

# I. Federated Averaging

# Federated Averaging<sup>1</sup>

$$x^* \in \arg \min_{x \in \mathbb{R}^d} \frac{1}{N} \sum_{c=1}^N \mathbb{E}_Z[F_c(x; Z)]$$

At each global iteration

- For  $c = 1$  à  $N$  in parallel

- Receive  $x^{(t)}$ , set  $x_c^{(t,0)} = x^{(t)}$

- For  $h = 0$  to  $H - 1$

$$x_c^{(t,h+1)} = x_c^{(t,h)} - \gamma \nabla F_c(x_c^{(t,h)}; Z_c^{(t,h+1)})$$

- Aggregate local models

$$x^{(t+1)} = \frac{1}{N} \sum_{c=1}^N x_c^{(t,H)}$$

---

<sup>1</sup>B. McMahan et al. "Communication-efficient learning of deep networks from decentralized data". In: **AISTATS**. 2017.

# Federated Averaging<sup>1</sup>

$$x^* \in \arg \min_{x \in \mathbb{R}^d} \frac{1}{N} \sum_{c=1}^N \mathbb{E}_Z[F_c(x; Z)]$$

At each global iteration

- For  $c = 1$  à  $N$  in parallel

- Receive  $x^{(t)}$ , set  $x_c^{(t,0)} = x^{(t)}$

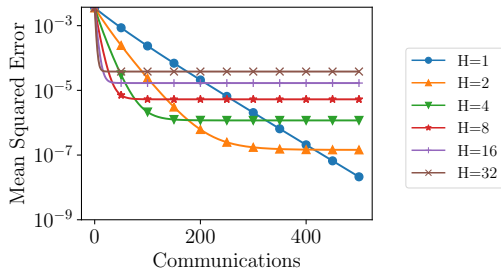
- For  $h = 0$  to  $H - 1$

$$x_c^{(t,h+1)} = x_c^{(t,h)} - \gamma \nabla F_c(x_c^{(t,h)}; Z_c^{(t,h+1)})$$

- Aggregate local models

$$x^{(t+1)} = \frac{1}{N} \sum_{c=1}^N x_c^{(t,H)}$$

With deterministic gradients:



<sup>1</sup>B. McMahan et al. "Communication-efficient learning of deep networks from decentralized data". In: **AISTATS**. 2017.



# Classical analyses of this algorithm

(For  $L$ -smooth,  $\mu$ -strongly convex functions)

Choose your favorite heterogeneity measure

- first-order<sup>1</sup>:  $\zeta = \frac{1}{N} \sum_{c=1}^N \|\nabla f_c(x^*) - \nabla f(x^*)\|^2$
- second-order<sup>2</sup>:  $\zeta = \frac{1}{N} \sum_{c=1}^N \|\nabla_c^2 f(x^*) - \nabla^2 f(x^*)\|^2$
- average drift<sup>3</sup>:  $\zeta = \left\| \frac{1}{NH} \sum_{c=1}^N \sum_{h=0}^{H-1} \nabla f(x_c^{(h)}) - \nabla f(x^*) \right\|^2$

---

<sup>1</sup>X. Lian et al. "Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel SGD". In: *NeurIPS* (2017).

<sup>2</sup>A. Khaled and C. Jin. "Faster federated optimization under second-order similarity". In: *arXiv* (2022).

<sup>3</sup>J. Wang et al. "On the Unreasonable Effectiveness of Federated Averaging with Heterogeneous Data". In: *TMLR* (2024).

# Classical analyses of this algorithm

(For  $L$ -smooth,  $\mu$ -strongly convex functions)

Choose your favorite heterogeneity measure

- first-order<sup>1</sup>:  $\zeta = \frac{1}{N} \sum_{c=1}^N \|\nabla f_c(x^*) - \nabla f(x^*)\|^2$
- second-order<sup>2</sup>:  $\zeta = \frac{1}{N} \sum_{c=1}^N \|\nabla_c^2 f(x^*) - \nabla^2 f(x^*)\|^2$
- average drift<sup>3</sup>:  $\zeta = \left\| \frac{1}{NH} \sum_{c=1}^N \sum_{h=0}^{H-1} \nabla f(x_c^{(h)}) - \nabla f(x^*) \right\|^2$

Show **convergence to a neighborhood** of  $x^*$

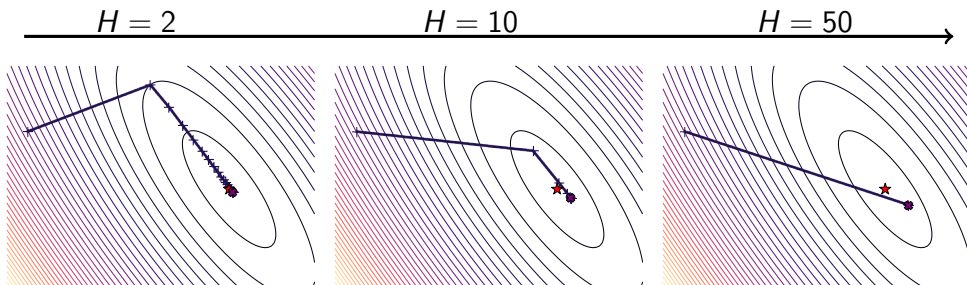
$$\|x^{(T)} - x^*\|^2 \lesssim (1 - \gamma\mu)^{HT} \|x^{(0)} - x^*\|^2 + \chi(\gamma, H, \zeta) \quad (\text{for some function } \chi)$$

---

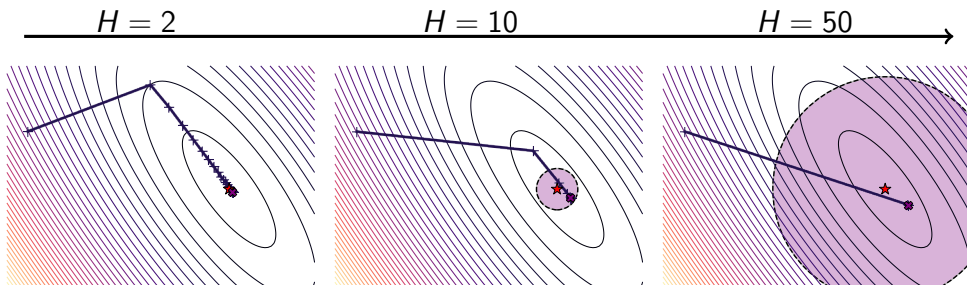
<sup>1</sup>X. Lian et al. "Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel SGD". In: *NeurIPS* (2017).

<sup>2</sup>A. Khaled and C. Jin. "Faster federated optimization under second-order similarity". In: *arXiv* (2022).

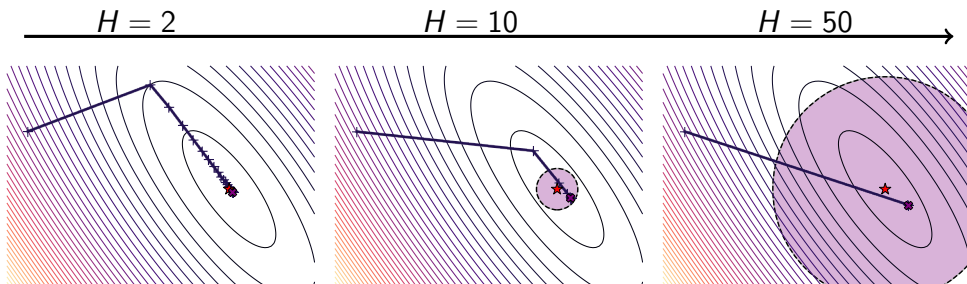
<sup>3</sup>J. Wang et al. "On the Unreasonable Effectiveness of Federated Averaging with Heterogeneous Data". In: *TMLR* (2024).



When the number of local iterations increases, bias increases



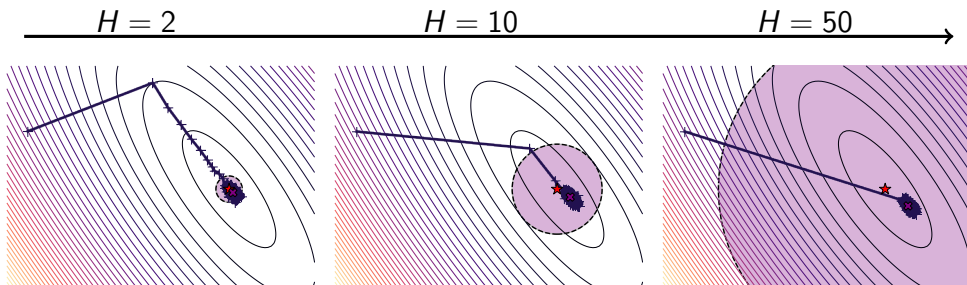
When the number of local iterations increases, bias increases  
... but the bound is oblivious to problem's geometry



When the number of local iterations increases, bias increases

... but the bound is oblivious to problem's geometry

**Remark:** It seems that iterates converge in some way?



When the number of local iterations increases, bias increases

... but the bound is oblivious to problem's geometry

**Remark:** It seems that iterates converge in some way?

# FedAvg (with stochastic gradients) converges!<sup>1</sup>

(For thrice derivable,  $L$ -smooth,  $\mu$ -strongly convex functions)

- FedAvg converges to a stationary distribution  $\pi^{(\gamma, H)}$ 
  - denoting  $x^{(t)} \sim \psi_{x^{(t)}}$ , we have

$$\mathcal{W}_2(\psi_{x^{(t)}}; \pi^{(\gamma, H)}) \leq (1 - \gamma\mu)^{Ht} \mathcal{W}_2(\psi_{x^{(0)}}; \pi^{(\gamma, H)})$$

- where  $\mathcal{W}_2$  is the second order Wasserstein distance

---

<sup>1</sup>P. Mangold et al. "Refined Analysis of Federated Averaging's Bias and Federated Richardson-Romberg Extrapolation". In: **AISTATS**. 2025.

# FedAvg (with stochastic gradients) converges!<sup>1</sup>

(For thrice derivable,  $L$ -smooth,  $\mu$ -strongly convex functions)

- FedAvg converges to a stationary distribution  $\pi^{(\gamma, H)}$
- FedAvg's iterates covariance is

$$\int (x - x^*)(x - x^*)^\top \pi^{(\gamma, H)}(dx) = \boxed{\frac{\gamma}{N} C(x^*)} + O(\gamma^{3/2} H)$$

---

<sup>1</sup>P. Mangold et al. "Refined Analysis of Federated Averaging's Bias and Federated Richardson-Romberg Extrapolation". In: **AISTATS**. 2025.



# FedAvg (with decreasing step-sizes) converges!<sup>1</sup>

(For  $\mu$ -strongly convex functions)

Linear speed-up !

variance decreases in  $1/N$   
variance scales in  $\gamma$

- FedAvg converges to the stationary distribution  $\pi^{(\gamma, H)}$
- FedAvg's iterates covariance is

$$\int (x - x^*)(x - x^*)^\top \pi^{(\gamma, H)}(dx) = \frac{\gamma}{N} C(x^*) + O(\gamma^{3/2} H)$$

<sup>1</sup>P. Mangold et al. "Refined Analysis of Federated Averaging's Bias and Federated Richardson-Romberg Extrapolation". In: **AISTATS**. 2025.

# FedAvg (with stochastic gradients) converges!<sup>1</sup>

(For thrice derivable,  $L$ -smooth,  $\mu$ -strongly convex functions)

- FedAvg converges to a stationary distribution  $\pi^{(\gamma, H)}$
- FedAvg's iterates covariance is
- We can now give an **exact expansion of the bias**

$$\int x \pi^{(\gamma, H)}(dx) = x^* + \frac{\gamma(H-1)}{2N} \sum_{c=1}^N \nabla^2 f(x^*)^{-1} (\nabla^2 f_c(x^*) - \nabla^2 f(x^*)) \nabla f_c(x^*)$$
$$- \frac{\gamma}{2N} \nabla^2 f(x^*)^{-1} \nabla^3 f(x^*) A^{-1} C(x^*) + O(\gamma^{3/2} H)$$

<sup>1</sup>P. Mangold et al. "Refined Analysis of Federated Averaging's Bias and Federated Richardson-Romberg Extrapolation". In: **AISTATS**. 2025.

# FedAvg (with stochastic gradients) converges!<sup>1</sup>

## Heterogeneity bias

vanishes when  $\nabla^2 f_c(x^*) = \nabla^2 f(x^*)$   
or when  $\nabla f_c(x^*) = \nabla f(x^*)$

## Stochasticity bias

$A$  is some linear operator  
 $C(x^*)$  is  $\nabla f$ 's covariance at  $x^*$

- FedAvg's iterates covariance is
- We can now give an **exact expansion of the bias**

$$\int x \pi^{(\gamma, H)}(dx) = x^* + \frac{\gamma(H-1)}{2N} \sum_{c=1}^N \nabla^2 f(x^*)^{-1} (\nabla^2 f_c(x^*) - \nabla^2 f(x^*)) \nabla f_c(x^*)$$

$$- \frac{\gamma}{2N} \nabla^2 f(x^*)^{-1} \nabla^3 f(x^*) A^{-1} C(x^*) + O(\gamma^{3/2} H)$$

<sup>1</sup>P. Mangold et al. "Refined Analysis of Federated Averaging's Bias and Federated Richardson-Romberg Extrapolation". In: **AISTATS**. 2025.

## II. Correcting heterogeneity: Scaffold

# Scaffold<sup>1</sup>

$$x^* \in \arg \min_{x \in \mathbb{R}^d} \frac{1}{N} \sum_{c=1}^N \mathbb{E}_Z[F_c(x; Z)]$$

At each global iteration

- For  $c = 1$  to  $N$  in parallel

- Receive  $x^{(t)}$ , set  $x_c^{(t,0)} = x^{(t)}$

- For  $h = 0$  to  $H - 1$

$$x_c^{(t,h+1)} = x_c^{(t,h)} - \gamma (\nabla F_c(x_c^{(t,h)}; Z_c^{(t,h+1)}) + \xi_c^{(t)})$$

- Aggregate models, update control variates

$$x^{(t+1)} = \frac{1}{N} \sum_{c=1}^N x_c^{(t,H)}$$

$$\xi_c^{(t+1)} = \xi_c^{(t)} + \frac{1}{\gamma H} (\theta_c^{t,H} - \theta^{(t+1)})$$

---

<sup>1</sup>S. P. Karimireddy et al. "Scaffold: Stochastic controlled averaging for federated learning". In: **International conference on machine learning**. PMLR, 2020, pp. 5132–5143.

# Scaffold<sup>1</sup>

$$x^* \in \arg \min_{x \in \mathbb{R}^d} \frac{1}{N} \sum_{c=1}^N \mathbb{E}_Z[F_c(x; Z)]$$

At each global iteration

- For  $c = 1$  to  $N$  in parallel

- Receive  $x^{(t)}$ , set  $x_c^{(t,0)} = x^{(t)}$

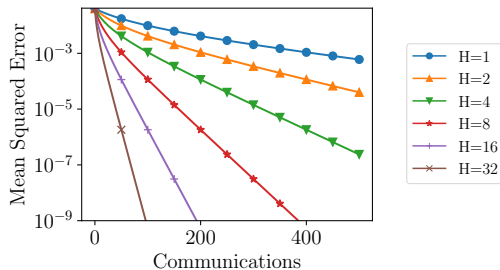
- For  $h = 0$  to  $H - 1$

$$x_c^{(t,h+1)} = x_c^{(t,h)} - \gamma (\nabla F_c(x_c^{(t,h)}; Z_c^{(t,h+1)}) + \xi_c^{(t)})$$

- Aggregate models, update control variates

$$x^{(t+1)} = \frac{1}{N} \sum_{c=1}^N x_c^{(t,H)}$$

$$\xi_c^{(t+1)} = \xi_c^{(t)} + \frac{1}{\gamma H} (\theta_c^{t,H} - \theta^{(t+1)})$$



→ No more heterogeneity bias!

<sup>1</sup>S. P. Karimireddy et al. "Scaffold: Stochastic controlled averaging for federated learning". In: *International conference on machine learning*. PMLR, 2020, pp. 5132–5143.

# Scaffold also converges !<sup>1</sup>

(For  $L$ -smooth,  $\mu$ -strongly convex functions with  $\nabla^3 f(x)$  bounded by  $Q$ )

- Scaffold converges if  $\gamma HL \leq 1$ , towards a distribution  $\pi^{(\gamma, H)}$ 
  - denoting  $x^{(t)} \sim \psi_{x^{(t)}}$ , we have

$$\mathcal{W}_2(\psi_{x^{(t)}}; \pi^{(\gamma, H)}) \leq (1 - \gamma\mu)^{Ht} \mathcal{W}_2(\psi_{x^{(0)}}; \pi^{(\gamma, H)})$$

- where  $\mathcal{W}_2$  is the second order Wasserstein distance

---

<sup>1</sup>P. Mangold et al. "Scaffold with Stochastic Gradients: New Analysis with Linear Speed-Up". In: [arxiv preprint](#). 2025.

# Scaffold also converges !<sup>1</sup>

(For  $L$ -smooth,  $\mu$ -strongly convex functions with  $\nabla^3 f(x)$  bounded by  $Q$ )

- Scaffold converges if  $\gamma HL \leq 1$ , towards a distribution  $\pi^{(\gamma, H)}$
- Scaffold's variance is close to FedAvg's variance

$$\int (x - x^*)(x - x^*)^\top \pi^{(\gamma, H)}(dx) = \boxed{\frac{\gamma}{N} C(x^*)} + O(\gamma^{3/2})$$

---

<sup>1</sup>P. Mangold et al. "Scaffold with Stochastic Gradients: New Analysis with Linear Speed-Up". In: [arxiv preprint](#). 2025.



# Scaffold also converges !<sup>1</sup>

(For  $L$ -smooth functions with  $\nabla^3 f(x)$  bounded by  $Q$ )

Linear speed-up !

variance decreases in  $1/N$

variance scales in  $\gamma$

- Scaffold converges to distribution  $\pi^{(\gamma, H)}$
- Scaffold's variance is close to FedAvg's variance

$$\int (x - x^*)(x - x^*)^\top \pi^{(\gamma, H)}(dx) = \frac{\gamma}{N} C(x^*) + O(\gamma^{3/2})$$

<sup>1</sup>P. Mangold et al. "Scaffold with Stochastic Gradients: New Analysis with Linear Speed-Up". In: [arxiv preprint](#). 2025.

# Scaffold also converges !<sup>1</sup>

(For  $L$ -smooth,  $\mu$ -strongly convex functions with  $\nabla^3 f(x)$  bounded by  $Q$ )

- Scaffold converges if  $\gamma HL \leq 1$ , towards a distribution  $\pi^{(\gamma, H)}$
- Scaffold's variance is close to FedAvg's variance
- Scaffold still has some bias

$$\int x \pi^{(\gamma, H)}(dx) = x^* - \frac{\gamma}{2N} \nabla^2 f(x^*)^{-1} \nabla^3 f(x^*) A^{-1} C(x^*) + O(\gamma^{3/2})$$

---

<sup>1</sup>P. Mangold et al. "Scaffold with Stochastic Gradients: New Analysis with Linear Speed-Up". In: [arxiv preprint](#). 2025.

# Scaffold also converges !<sup>1</sup>

(For  $L$ -smooth,  $\mu$ -strongly convex function)

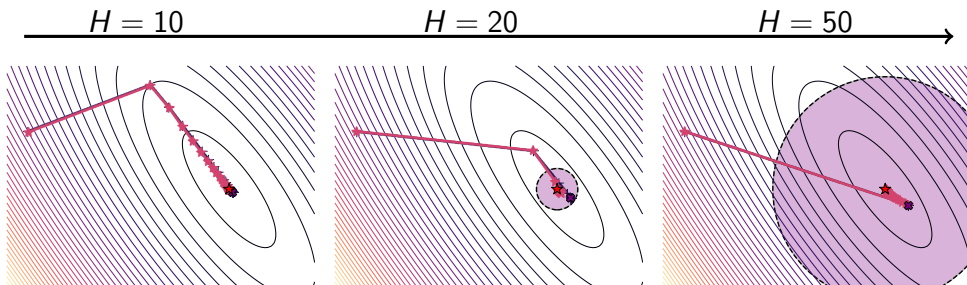
Stochasticity bias remains!

$A$  is some linear operator  
 $C(x^*)$  is  $\nabla f$ 's covariance at  $x^*$

- Scaffold converges if  $\gamma HL \leq 1$ , towards a distribution
- Scaffold's variance is close to FedAvg's variance
- Scaffold still has some bias

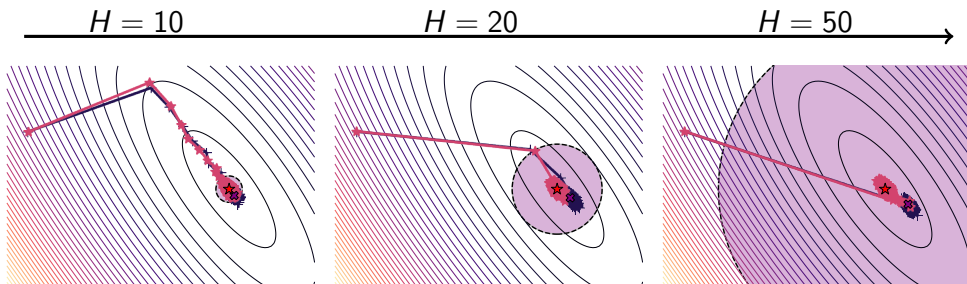
$$\int x \pi^{(\gamma, H)}(dx) = x^* - \frac{\gamma}{2N} \nabla^2 f(x^*)^{-1} \nabla^3 f(x^*) A^{-1} C(x^*) + O(\gamma^{3/2})$$

<sup>1</sup>P. Mangold et al. "Scaffold with Stochastic Gradients: New Analysis with Linear Speed-Up". In: [arxiv preprint](#). 2025.



Scaffold converges to the right point

... and its variance is similar to FedAvg!



Scaffold converges to the right point

... and its variance is similar to FedAvg!

# New Convergence Rate for Scaffold

(For  $L$ -smooth,  $\mu$ -strongly convex functions with  $\nabla^3 f(x)$  bounded by  $Q$ )

$$\mathbb{E} [\|x^{(T)} - x^*\|^2] \lesssim \left(1 - \frac{\gamma\mu}{4}\right)^{HT} \left\{ \|x^{(0)} - x^*\|^2 + 2\gamma^2 H^2 \zeta^2 + \frac{\sigma_\star^2}{L\mu} \right\} \\ + \frac{\gamma}{\mathbf{N}\mu} \sigma_\star^2 + \frac{\gamma^{3/2} Q}{\mu^{5/2}} \sigma_\star^3 + \frac{\gamma^3 H Q^2}{\mu^3} \sigma_\star^4$$

where

- $\sigma_\star^2 = \mathbb{E}[\frac{1}{N} \sum_{c=1}^N \|\nabla F_c^Z(x^*) - \nabla f_c(x^*)\|^2]$  is the variance at  $x^*$
- $\zeta^2 = \frac{1}{N} \sum_{c=1}^N \|\nabla f_c^Z(x^*)\|^2$  measures gradient heterogeneity

# Linear Speed-Up!

As long as  $N$  is not too large, one can obtain  $\mathbb{E} [\|x^{(T)} - x^*\|^2] \leq \epsilon^2$  with

$$\text{\#grad per client} = \tilde{O}\left(\frac{\sigma_*^2}{N\mu^2\epsilon^2} \log\left(\frac{1}{\epsilon}\right)\right)$$

# Conclusion

- FedAvg and Scaffold converge (even with stochastic gradients)
- This allows to derive new analyses for these problems, with exact first-order expression for bias
- And we proved that Scaffold has:
  - variance similar to FedAvg's variance
  - *linear speed-up* in the number of clients!!