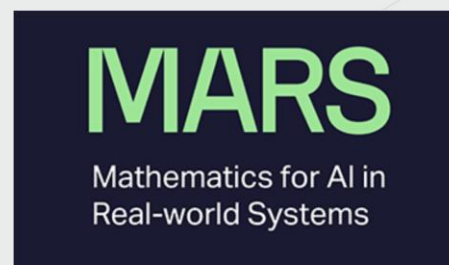


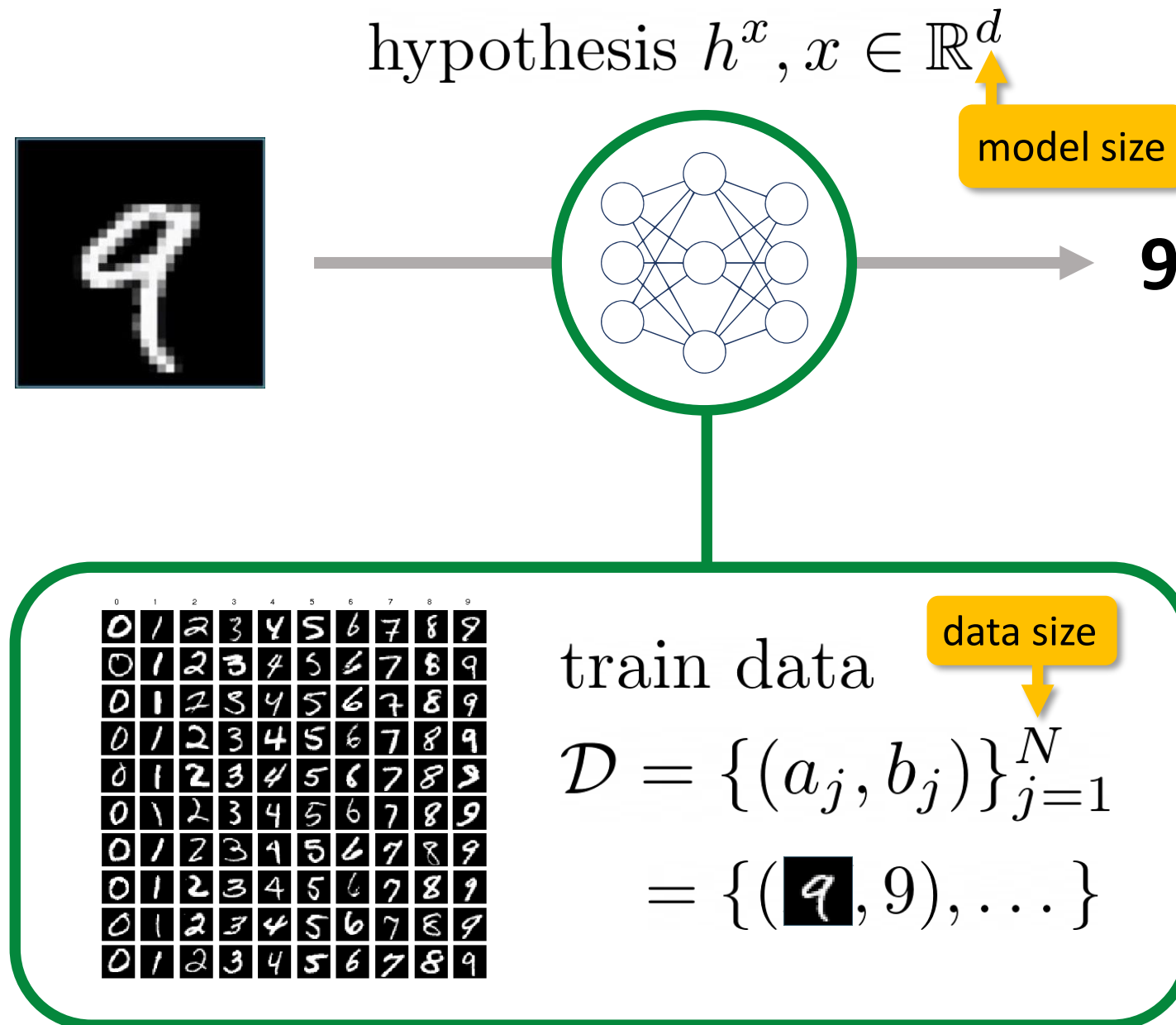
# Distributed Optimization and Error Feedback

Mher Safaryan

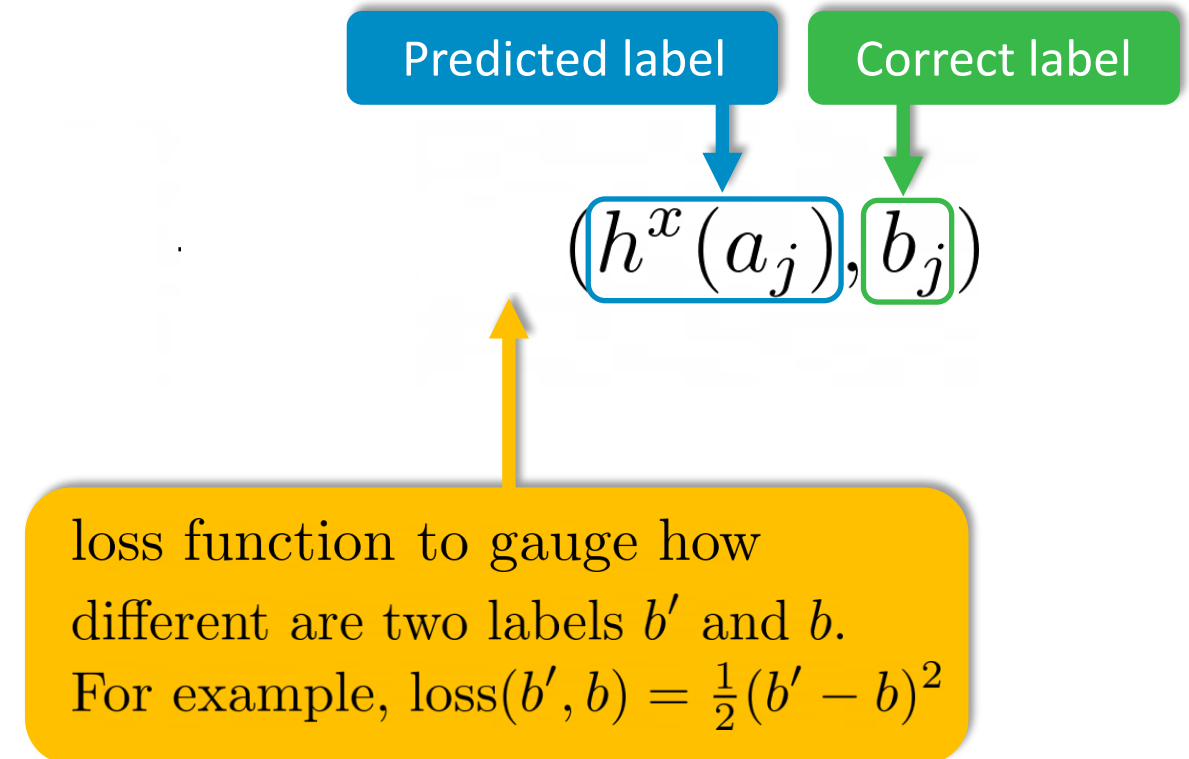
Assistant Professor (Lecturer)



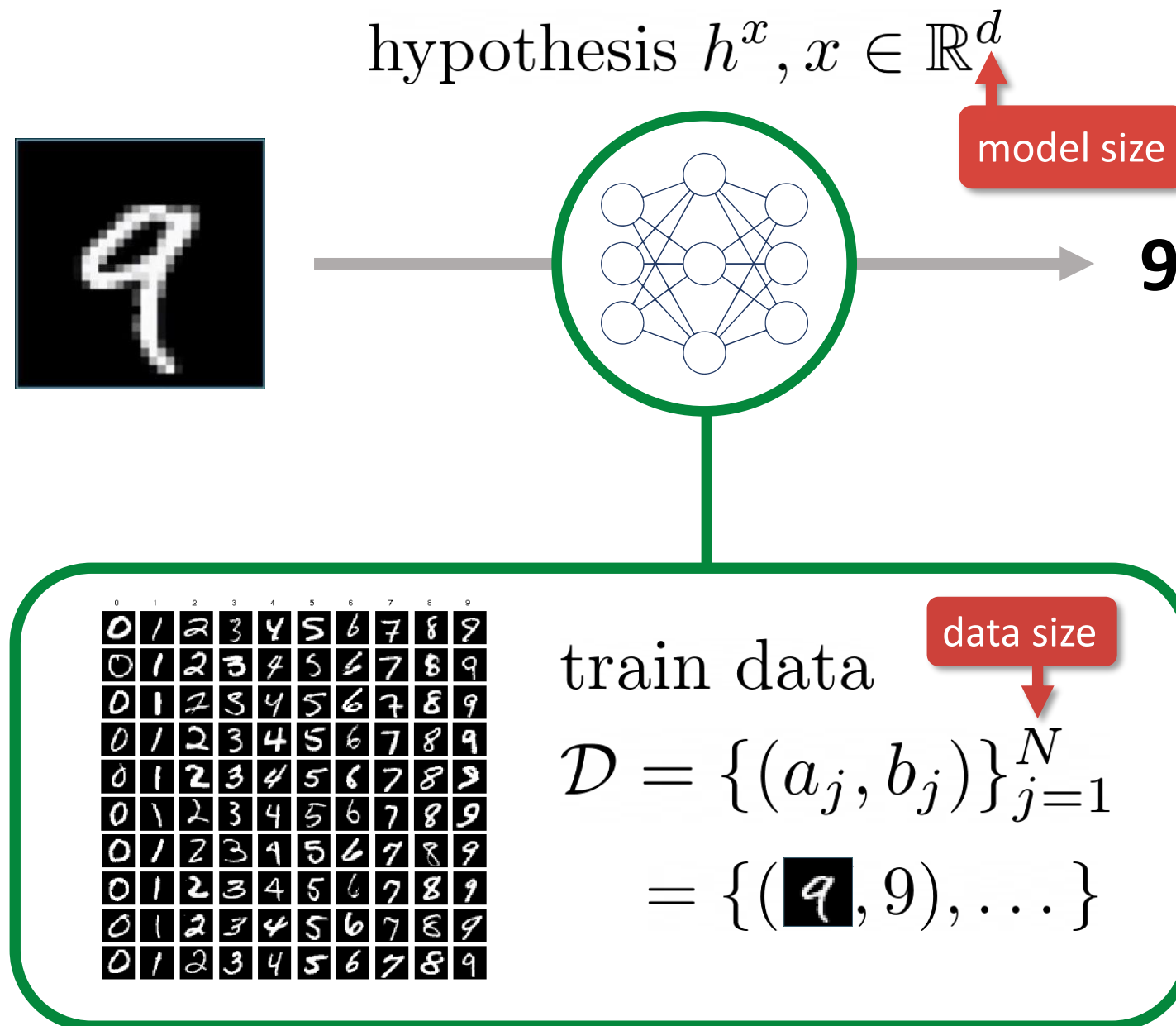
# Supervised Machine Learning



## Empirical Risk Minimization (ERM)



# Supervised Machine Learning



## Empirical Risk Minimization (ERM)

Predicted label

Correct label

$$\min_{x \in \mathbb{R}^d} \frac{1}{N} \sum_{j=1}^N \text{loss}(h^x(a_j), b_j)$$

loss function to gauge how different are two labels  $b'$  and  $b$ .  
For example,  $\text{loss}(b', b) = \frac{1}{2}(b' - b)^2$

# Empirical Risk Minimization

$$\min_{x \in \mathbb{R}^d} \frac{1}{N} \sum_{j=1}^N \text{loss}(h^x(a_j), b_j)$$

$$\mathcal{D} = \{(a_j, b_j)\}_{j=1}^N$$



$$\min_{x \in \mathbb{R}^d} \mathbb{E}_{(a,b) \sim \mathcal{D}} [\text{loss}(h^x(a), b)]$$

$$\xi \stackrel{\text{def}}{=} (a, b)$$

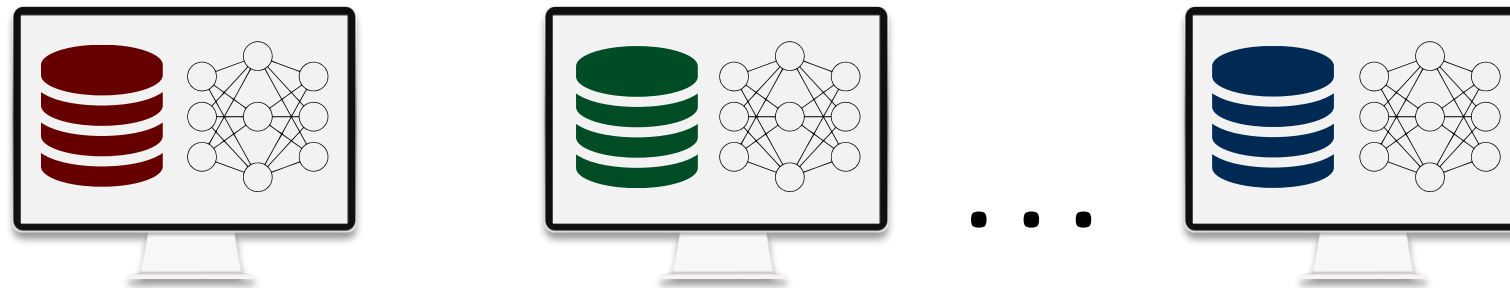
$$f_\xi(x) \stackrel{\text{def}}{=} \text{loss}(h^x(a), b)$$



$$\min_{x \in \mathbb{R}^d} \mathbb{E}_{\xi \sim \mathcal{D}} [f_\xi(x)]$$

# Distributed Machine Learning

## Why Distributed?



**Reason 1:** **BIG Data** does not fit into a single device

**Reason 2:** **Data Privacy** in Federated Learning

**Reason 3:** **Parallel Computation**

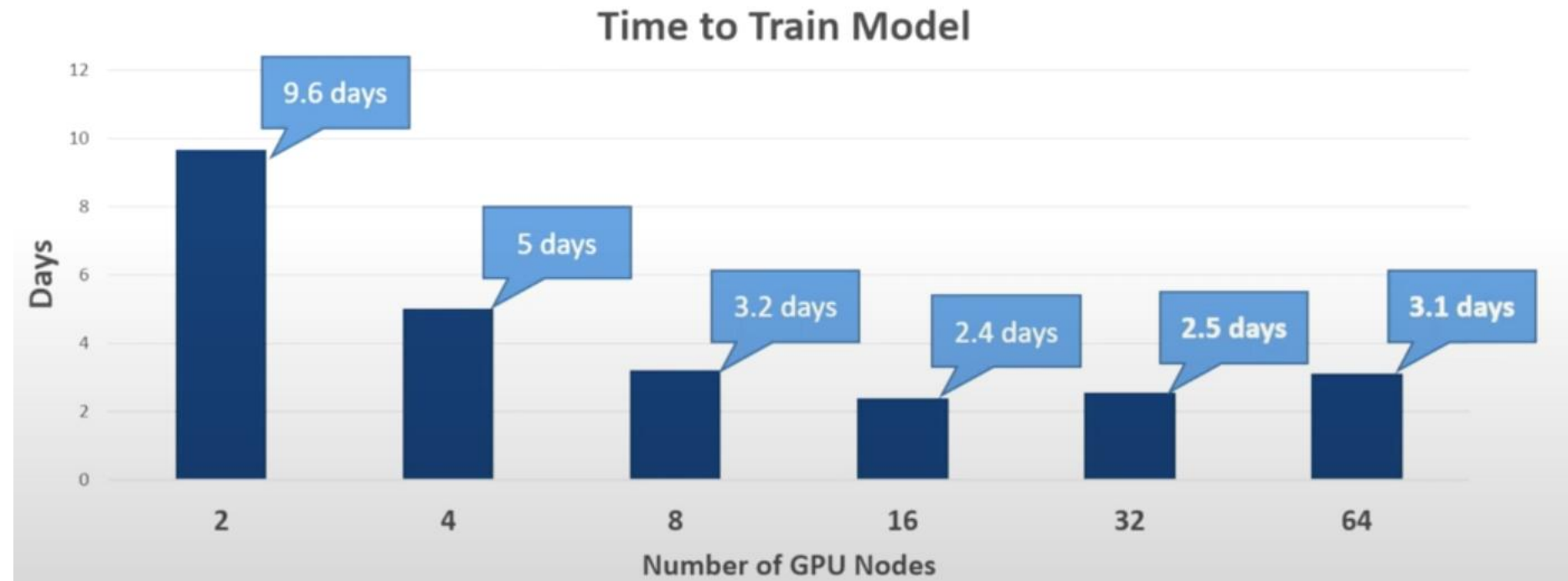
# Distributed Learning in Practice

CSCS: Europe's Top Supercomputer (World 4<sup>th</sup>)

- 4500+ GPU Nodes, state-of-the-art interconnect

Task:

- Image Classification (ResNet-152 on ImageNet)
- Single Node time (TensorFlow): **19 days**
- 1024 Nodes: **25 minutes** (*in theory*)



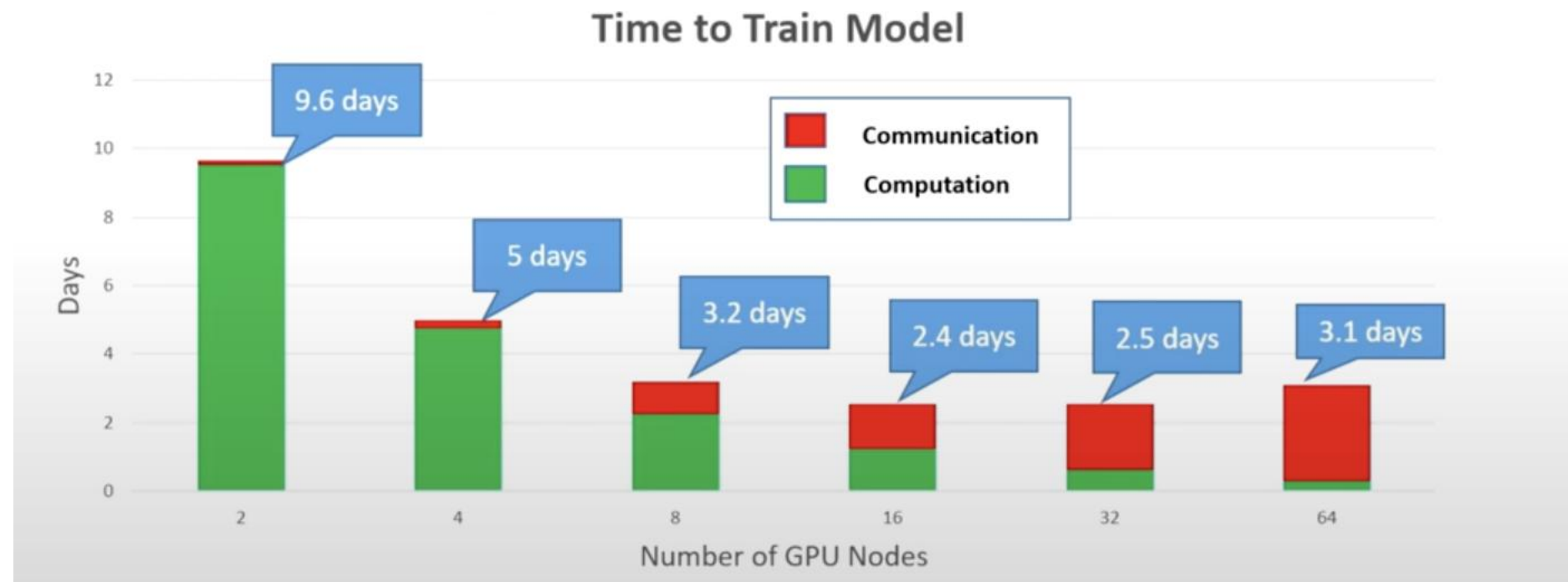
# Distributed Learning in Practice

CSCS: Europe's Top Supercomputer (World 4<sup>th</sup>)

- 4500+ GPU Nodes, state-of-the-art interconnect

Task:

- Image Classification (ResNet-152 on ImageNet)
- Single Node time (TensorFlow): **19 days**
- 1024 Nodes: **25 minutes** (*in theory*)

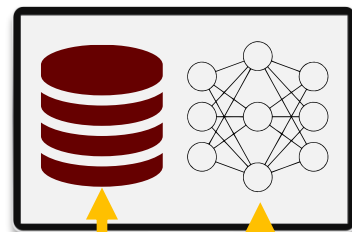




# Distributed Optimization Problem

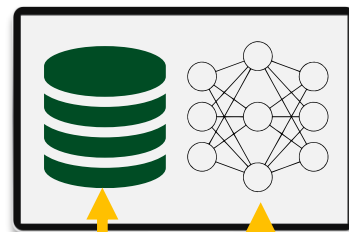
$$\min_{x \in \mathbb{R}^d} \left\{ f(x) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}$$

Overall risk/loss



$\mathcal{D}_1$

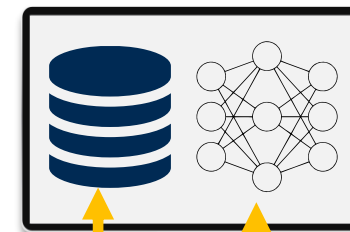
$x$



$\mathcal{D}_2$

$x$

...



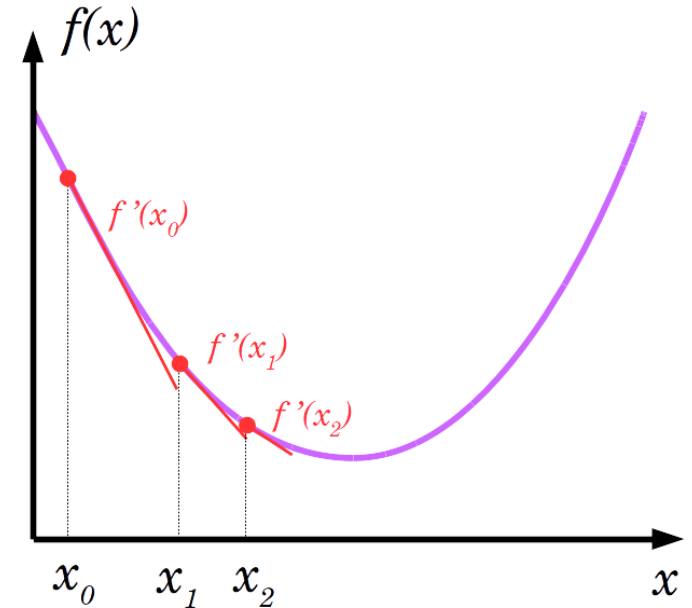
$\mathcal{D}_n$

$x$

$$f_1(x) = \mathbb{E}_{\xi \sim \mathcal{D}_1} [f_{\xi}(x)]$$

$$f_2(x) = \mathbb{E}_{\xi \sim \mathcal{D}_2} [f_{\xi}(x)]$$

$$f_n(x) = \mathbb{E}_{\xi \sim \mathcal{D}_n} [f_{\xi}(x)]$$

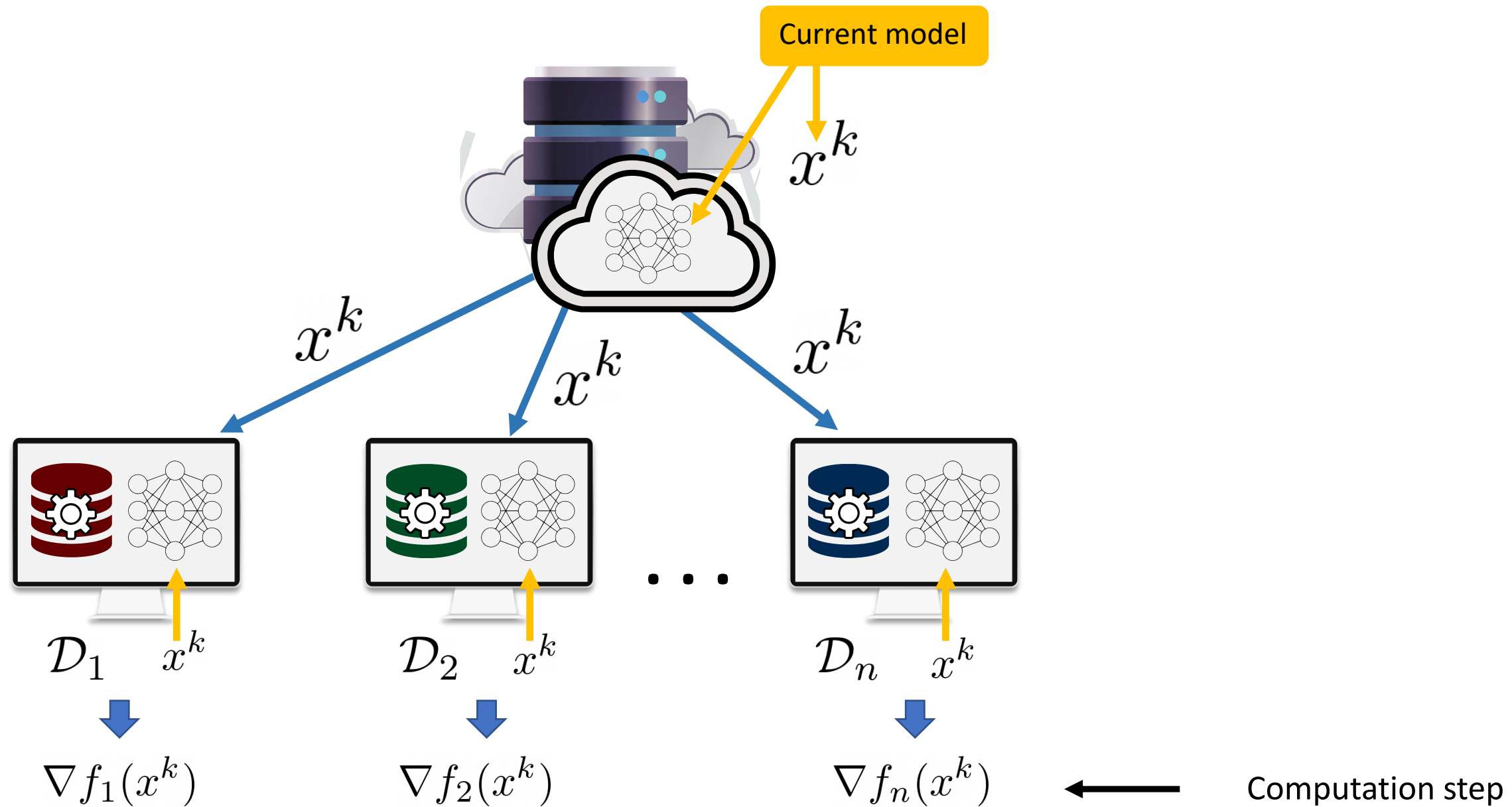


## Gradient Descent

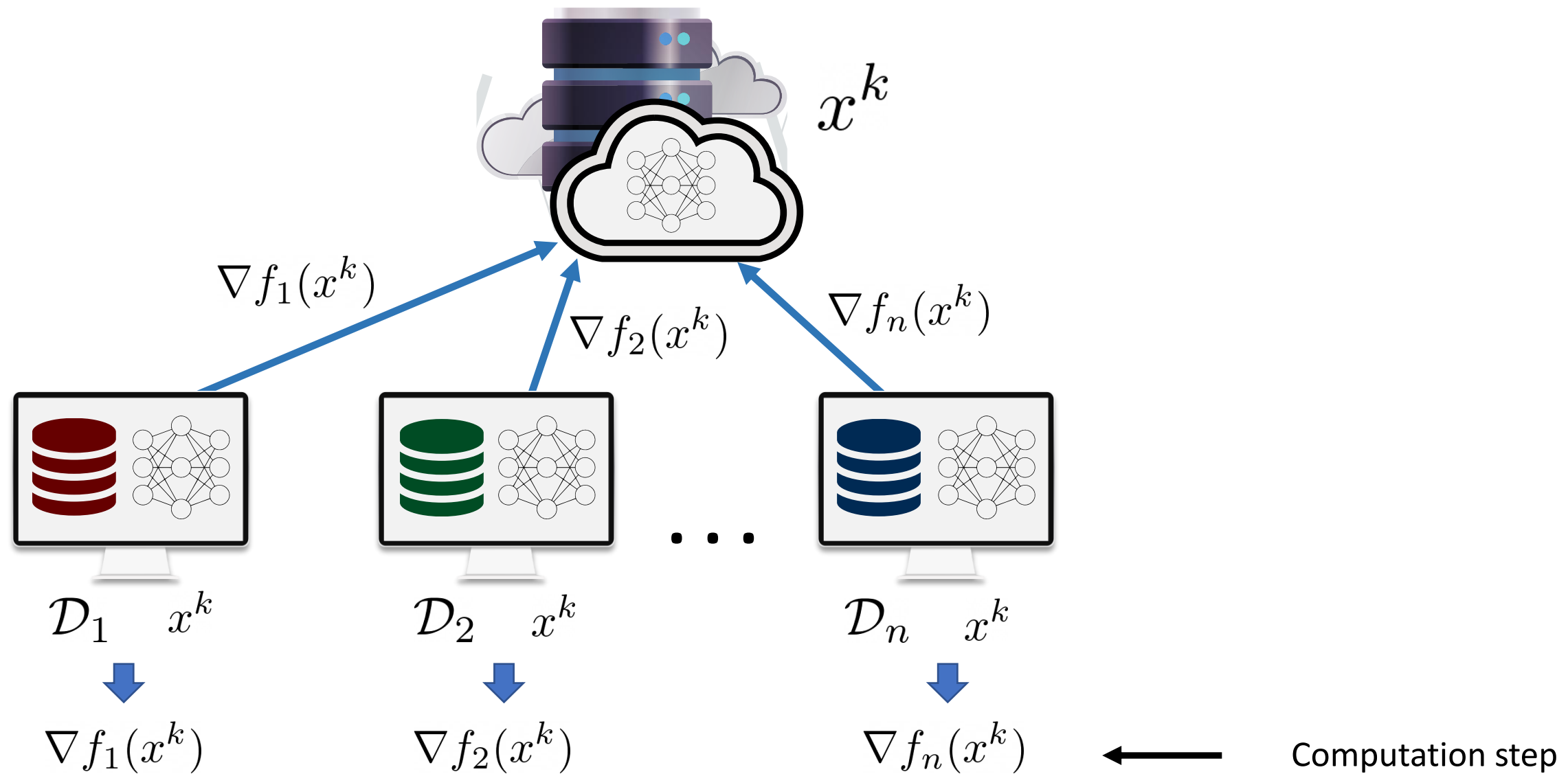
$$x^{k+1} = x^k - \gamma \nabla f(x^k)$$



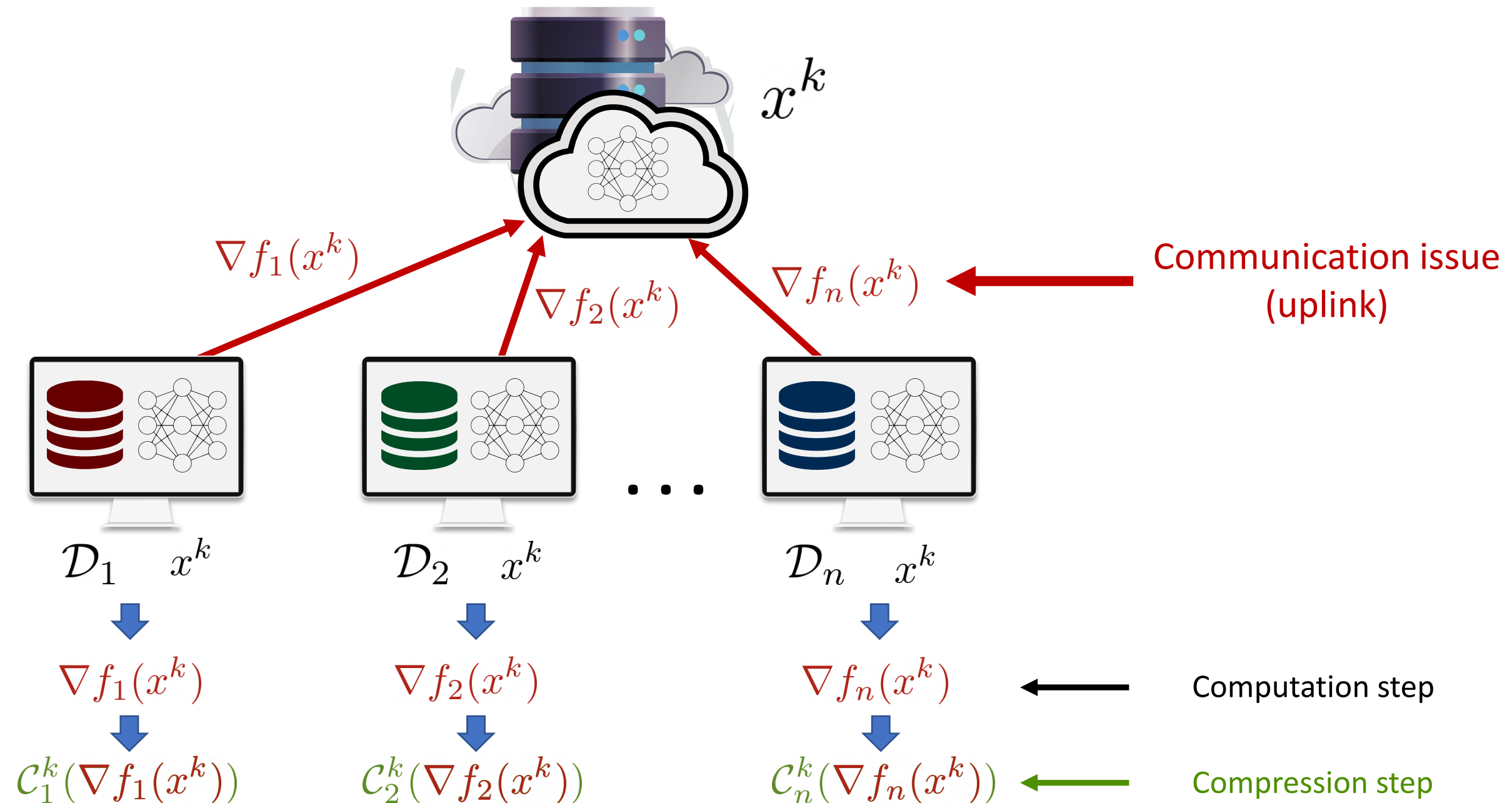
# Distributed Gradient Descent



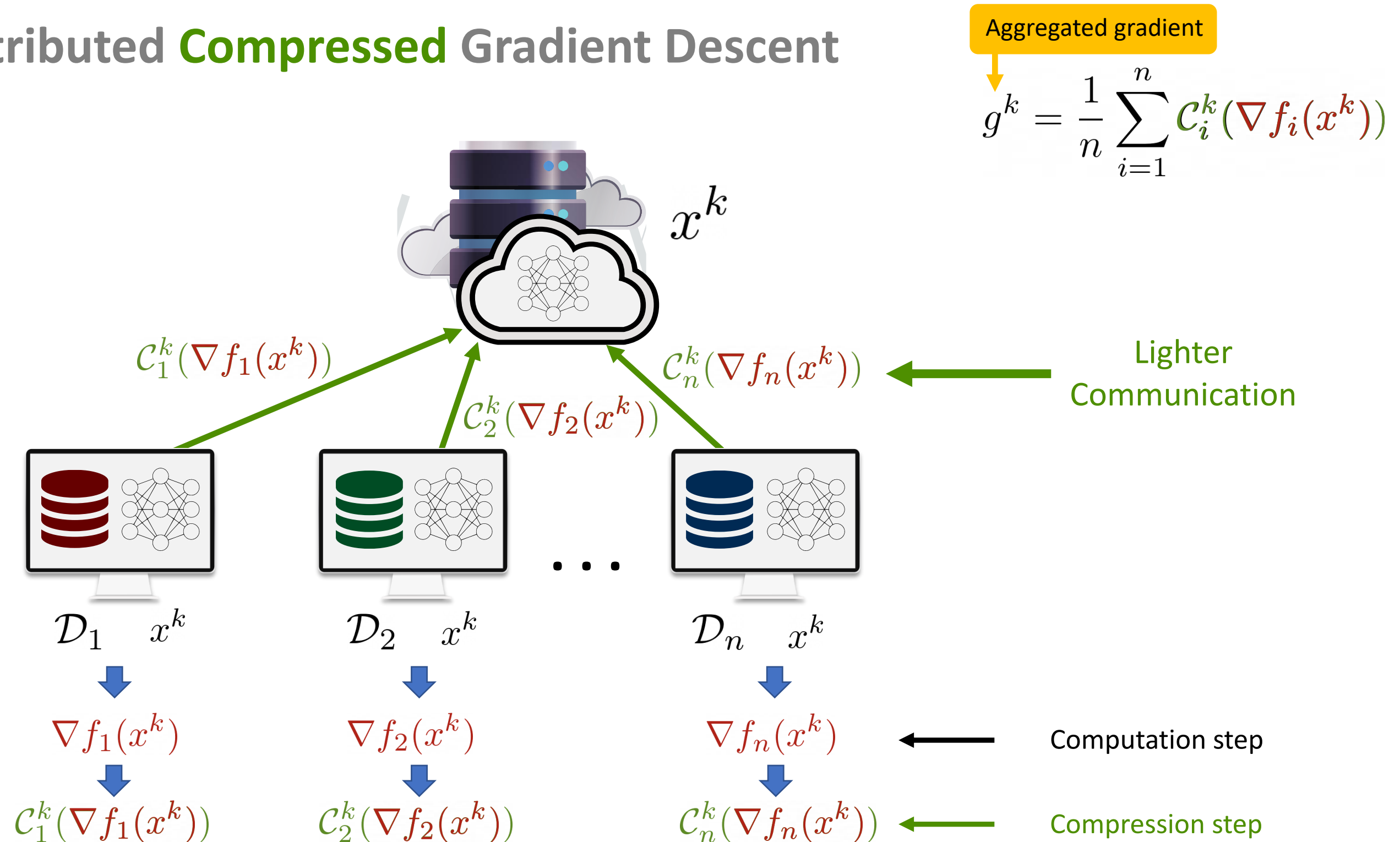
# Distributed Gradient Descent



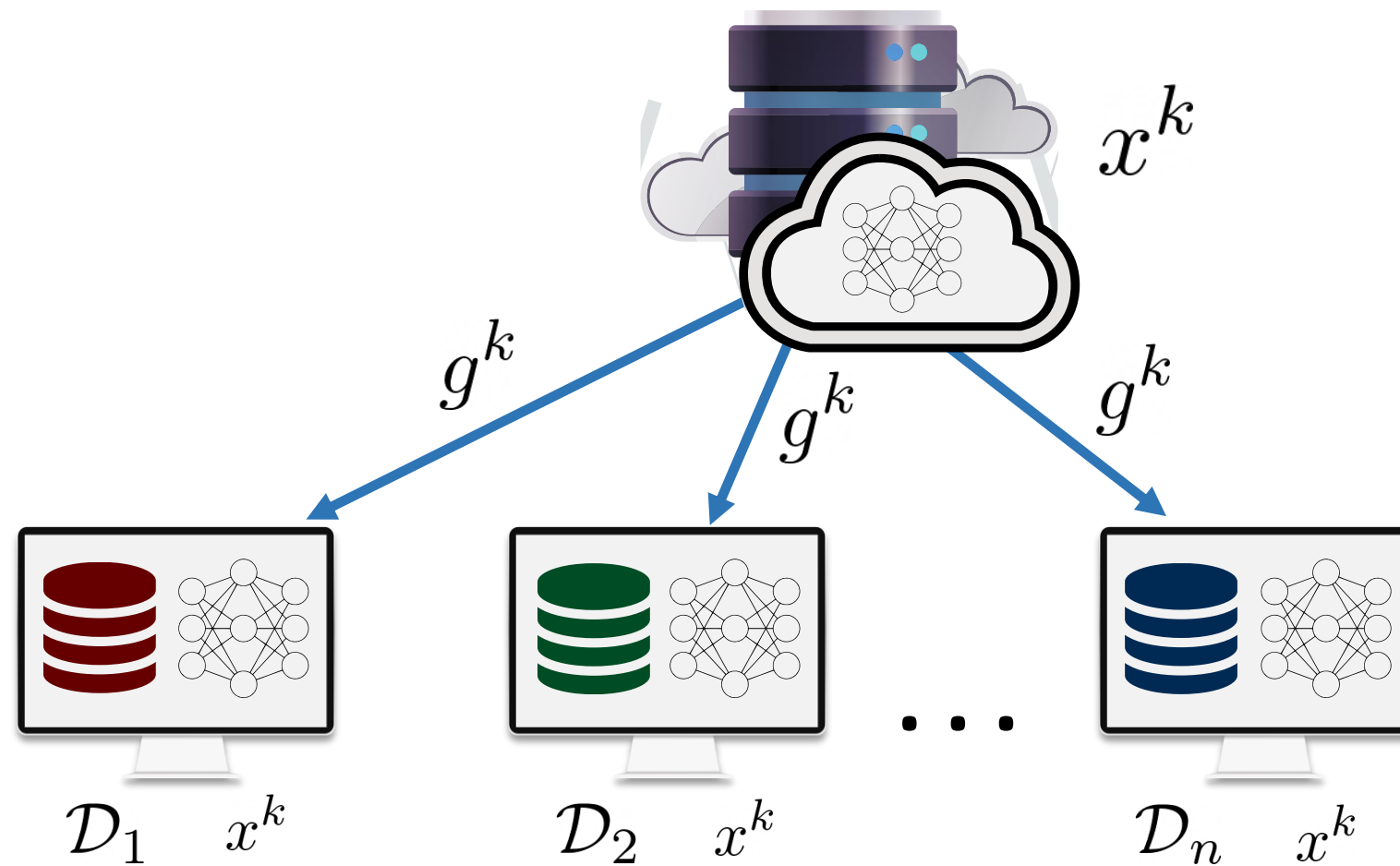
# Distributed Gradient Descent



# Distributed **Compressed** Gradient Descent



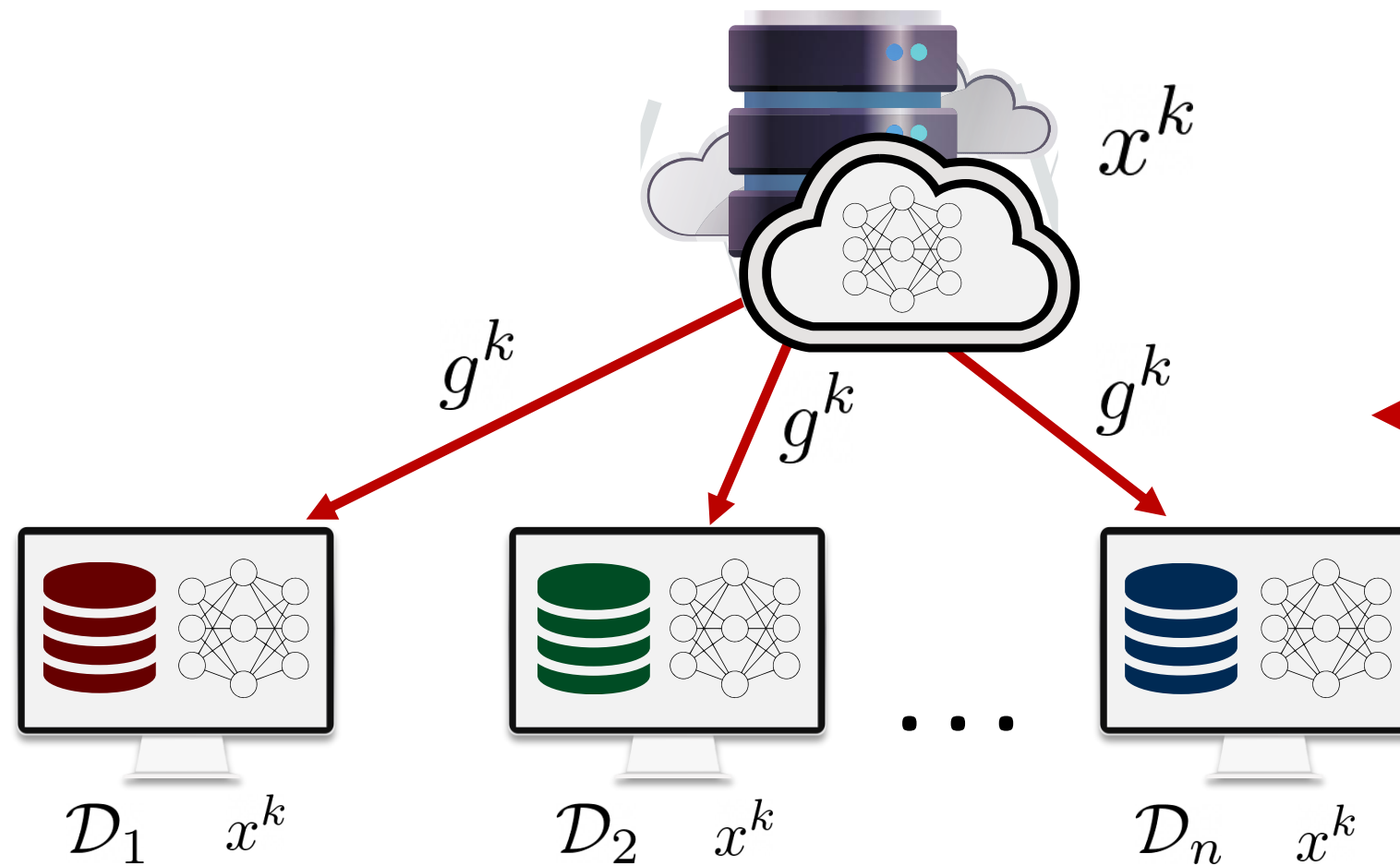
# Distributed **Compressed** Gradient Descent



Aggregated gradient

$$g^k = \frac{1}{n} \sum_{i=1}^n c_i^k (\nabla f_i(x^k))$$

# Distributed **Compressed** Gradient Descent

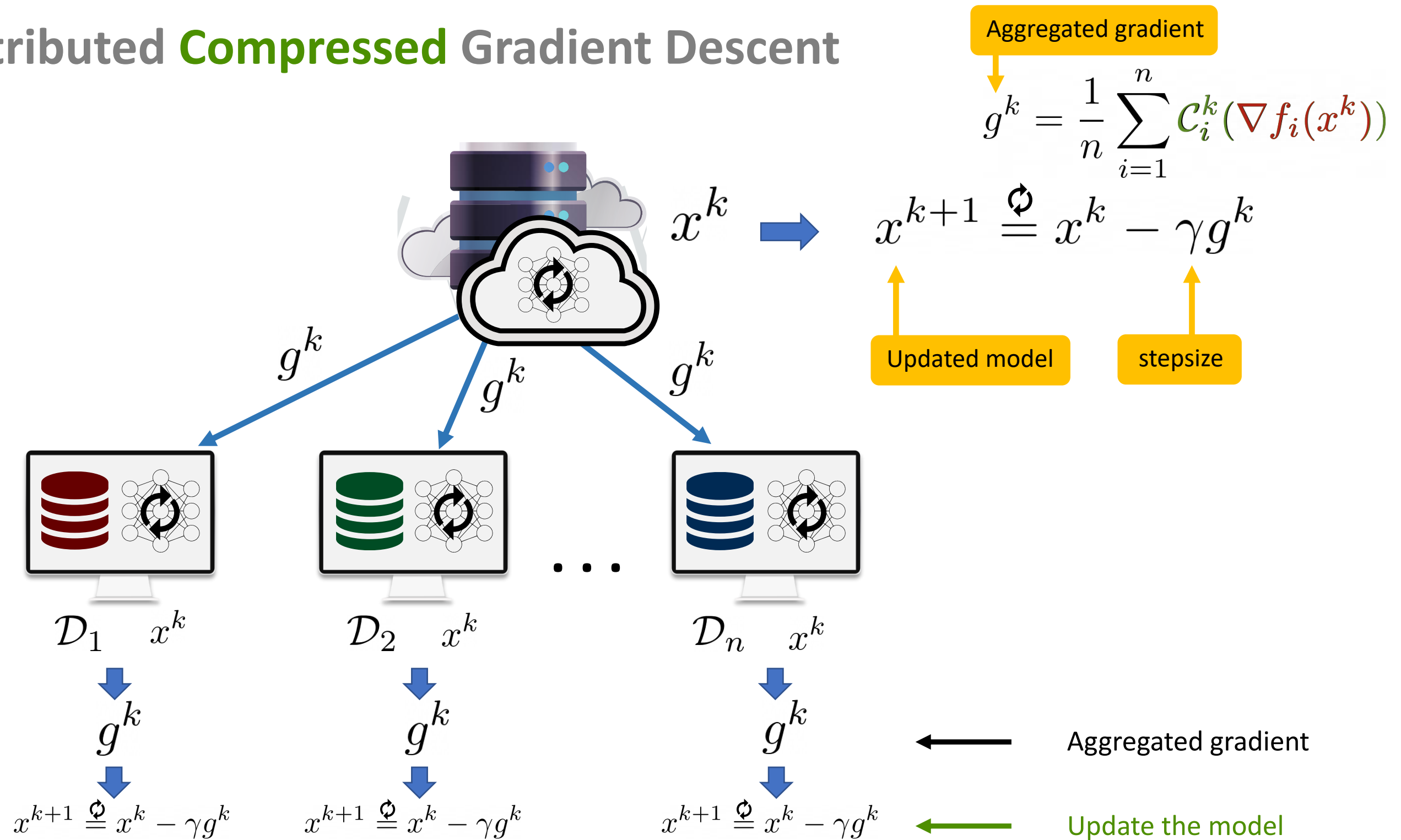


Aggregated gradient

$$g^k = \frac{1}{n} \sum_{i=1}^n c_i^k (\nabla f_i(x^k))$$

Communication issue  
(downlink)

# Distributed **Compressed** Gradient Descent





# Distributed **Compressed** Gradient Descent

Compression operator

$$\mathcal{C}_i^k : \mathbb{R}^d \rightarrow \mathbb{R}^d$$

$$x^{k+1} = x^k - \frac{\gamma}{n} \sum_{i=1}^n \mathcal{C}_i^k (\nabla f_i(x^k))$$

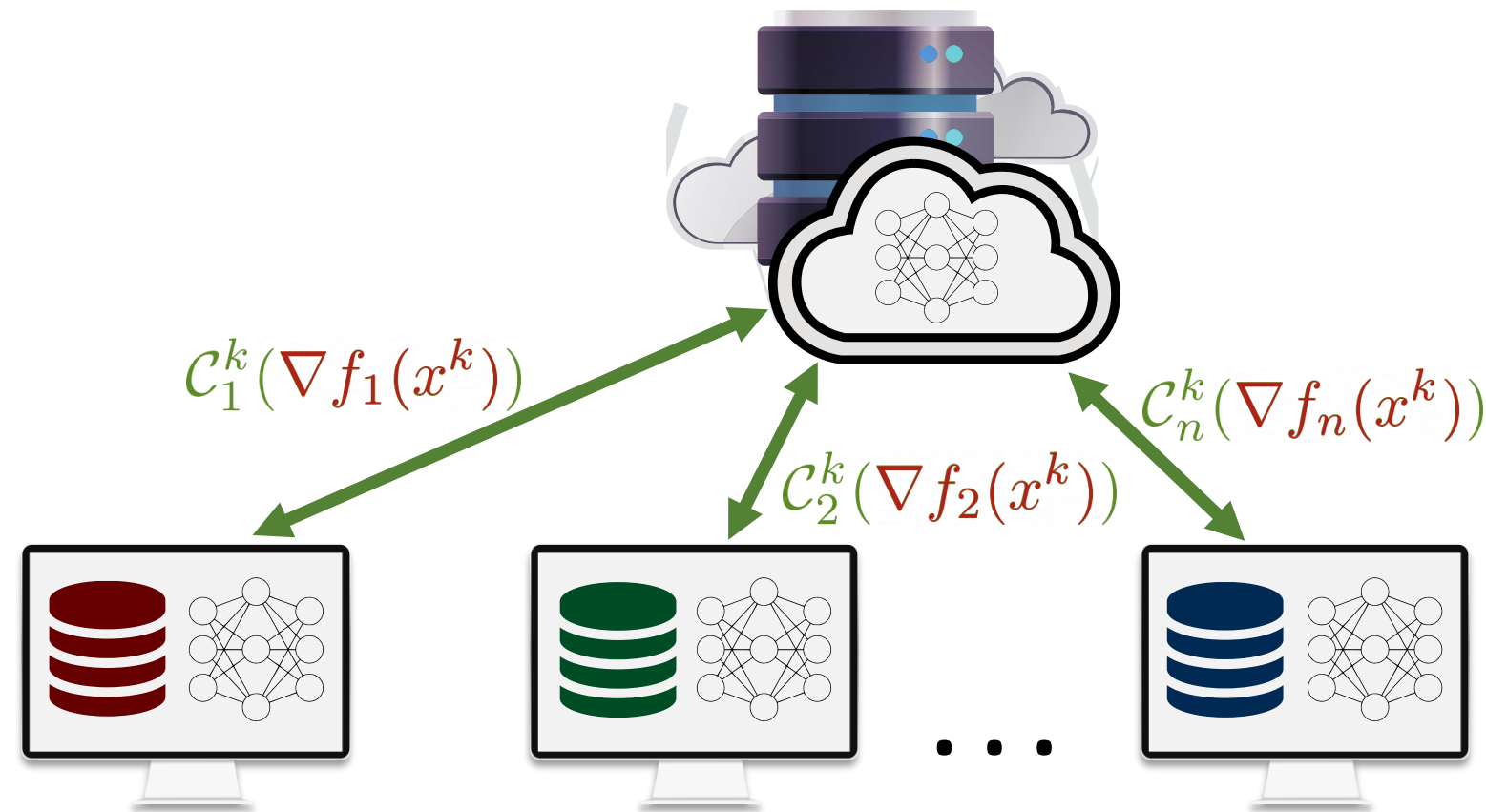
# Contributions

## ➤ Compressed communication

- Sign (1-bit) compression



Mher Safaryan, Peter Richtárik  
Stochastic Sign Descent Methods: New  
Algorithms and Better Theory, *ICML 2021*

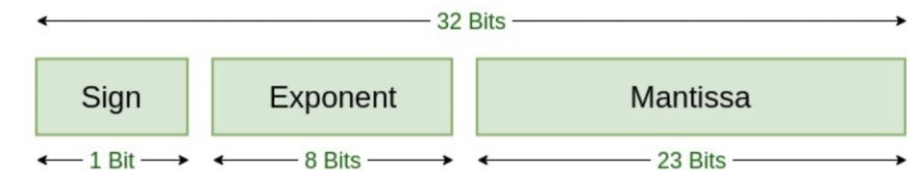


vector  $g \in \mathbb{R}^d$

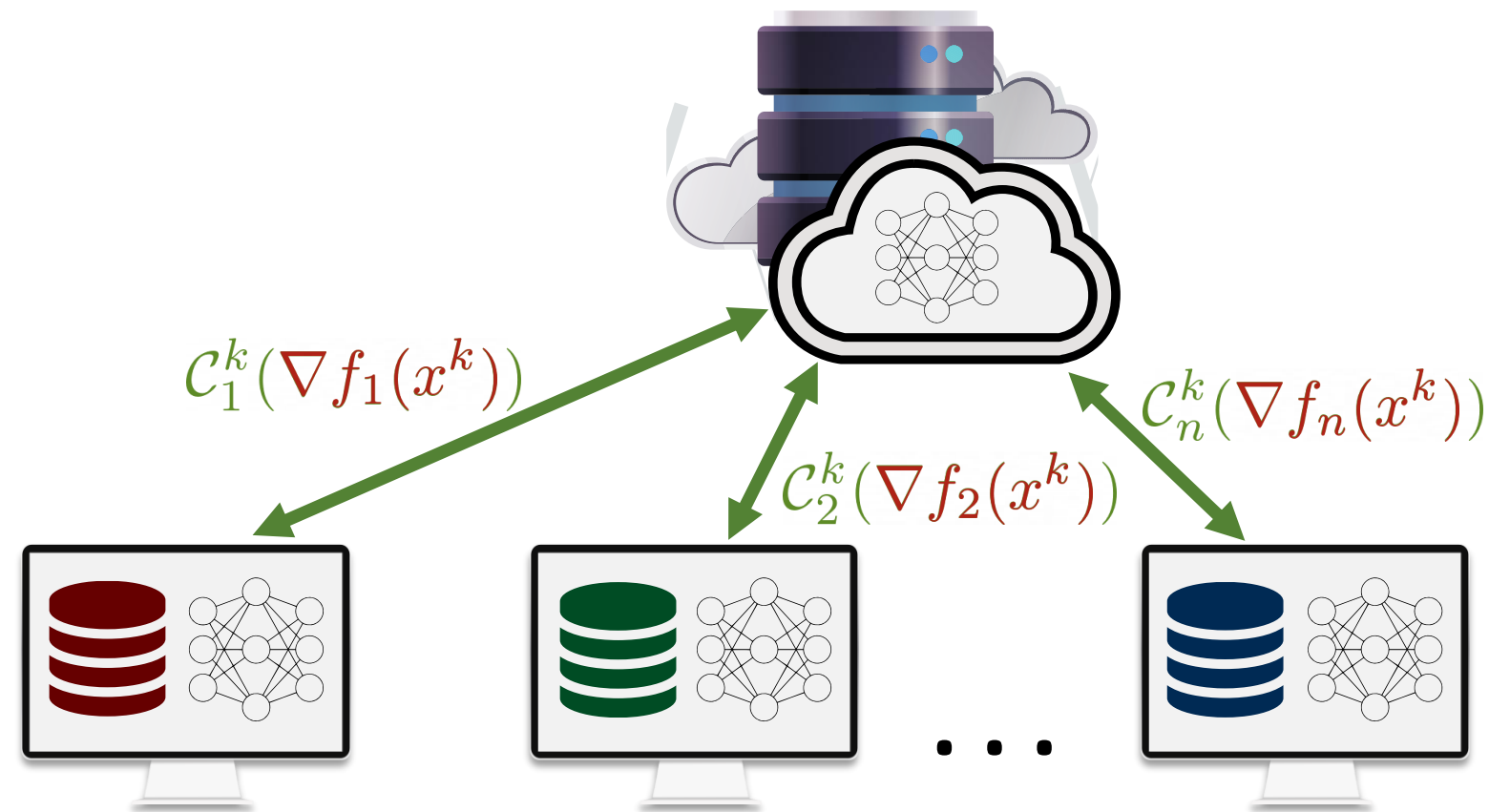
$$(\text{sign } g)_j = \begin{cases} +1 & \text{if } g_j > 0 \\ -1 & \text{if } g_j < 0 \\ 0 & \text{if } g_j = 0 \end{cases}$$

entry  $j \in \{1, 2, \dots, d\}$

$$\text{sign} \begin{bmatrix} 0.4 \\ 0 \\ -0.3 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ -1 \end{bmatrix}$$



# Contributions



$$\left(\widetilde{\text{sign}} g\right)_j = \begin{cases} +1 & \text{with prob. } \frac{1}{2} + \frac{1}{2} \frac{g_j}{\|g\|} \\ -1 & \text{with prob. } \frac{1}{2} - \frac{1}{2} \frac{g_j}{\|g\|} \end{cases}$$

Stochastic sign

Bernoulli random variable  $B\left(\frac{1}{2} + \frac{1}{2} \frac{g_j}{\|g\|}\right)$

## ➤ Compressed communication

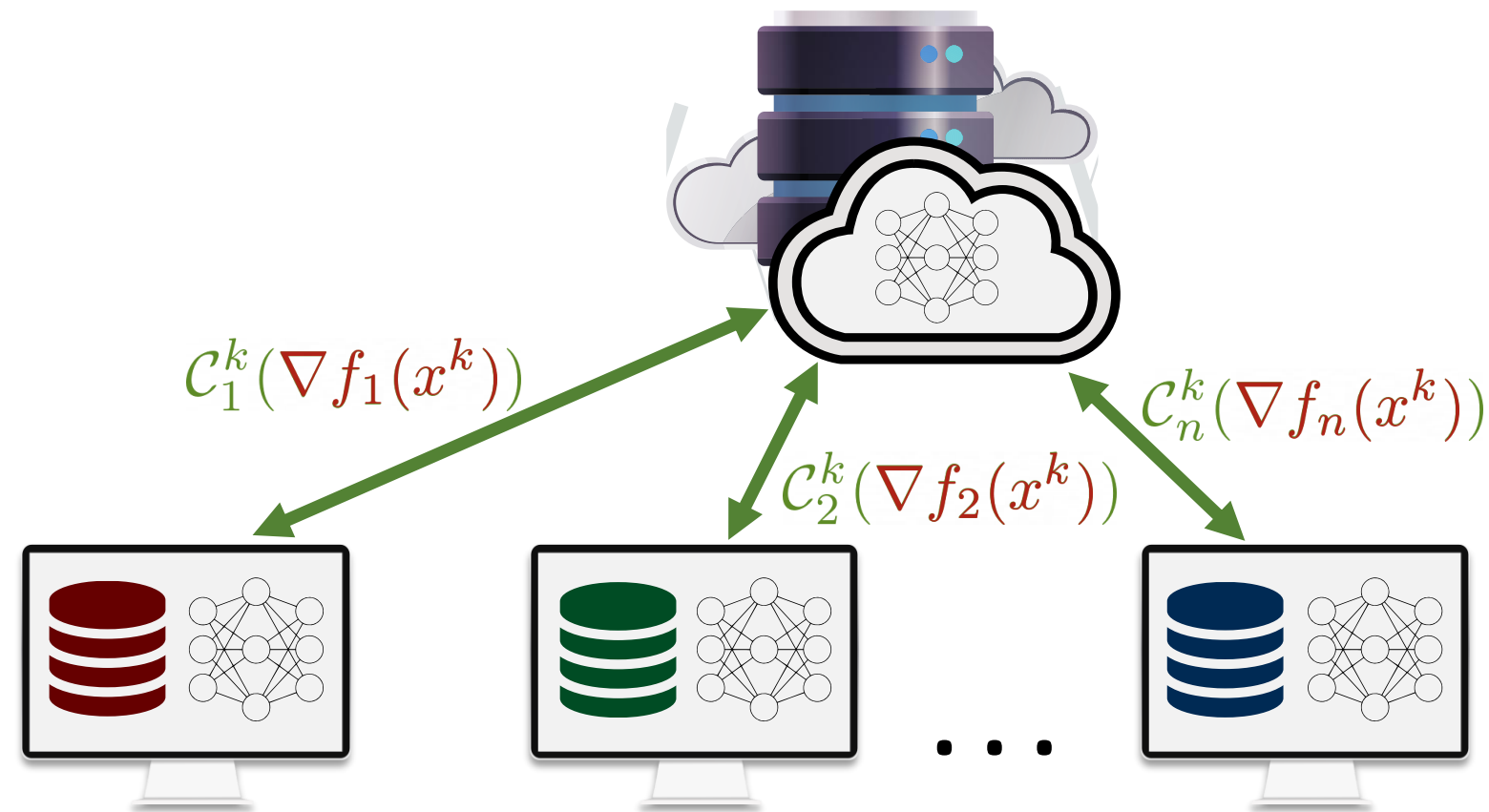
- Sign (1-bit) compression



Mher Safaryan, Peter Richtárik  
Stochastic Sign Descent Methods: New Algorithms and Better Theory, *ICML 2021*

$$\widetilde{\text{sign}} \begin{bmatrix} 0.4 \\ 0 \\ -0.3 \end{bmatrix} = \begin{bmatrix} B(0.9) \\ B(0.5) \\ B(0.2) \end{bmatrix}$$

# Contributions



$$\begin{array}{ccc}
 \text{Top}_k & & \text{Rand}_k \\
 \begin{bmatrix} -0.4 \\ 12.1 \\ 0.76 \\ 2.8 \\ -9.7 \end{bmatrix} & \xrightarrow{k=2} & \begin{bmatrix} 0 \\ 12.1 \\ 0 \\ 0 \\ -9.7 \end{bmatrix} \\
 & & \begin{bmatrix} -0.4 \\ 12.1 \\ 0.76 \\ 2.8 \\ -9.7 \end{bmatrix} \xrightarrow{k=2} \frac{5}{2} \cdot \begin{bmatrix} -0.4 \\ 0 \\ 0 \\ 2.8 \\ 0 \end{bmatrix}
 \end{array}$$

## ➤ Compressed communication

- Sign (1-bit) compression
- Contractive compression



Aleksandr Beznosikov, Samuel Horváth, Peter Richtárik, Mher Safaryan  
**On Biased Compression for Distributed Learning**, *JMLR 2023*

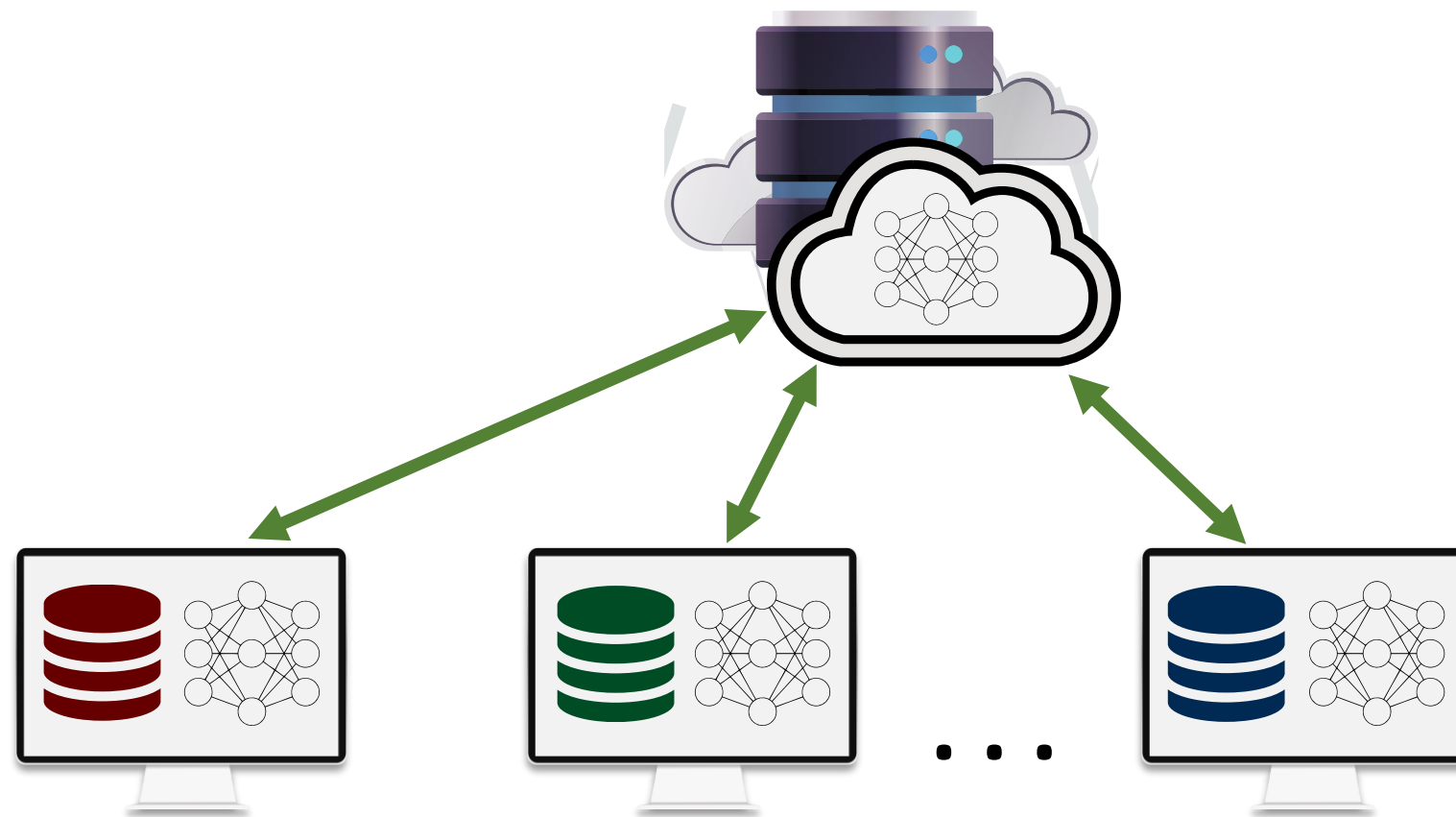


Mher Safaryan, Filip Hanzely, Peter Richtárik  
**Smoothness Matrices Beat Smoothness Constants: Better Communication Compression Techniques for Distributed Optimization**, *NeurIPS 2021*



Bokun Wang, Mher Safaryan, Peter Richtárik  
**Theoretically Better and Numerically Faster Distributed Optimization with Smoothness-Aware Quantization Techniques**, *NeurIPS 2022*

# Contributions



## ➤ Compressed communication

- Sign (1-bit) compression
- Contractive compression
- Second-order optimization



Mher Safaryan, Rustem Islamov,  
Xun Qian, Peter Richtárik  
**FedNL: Making Newton-type methods  
applicable to federated learning, *ICML 2022***



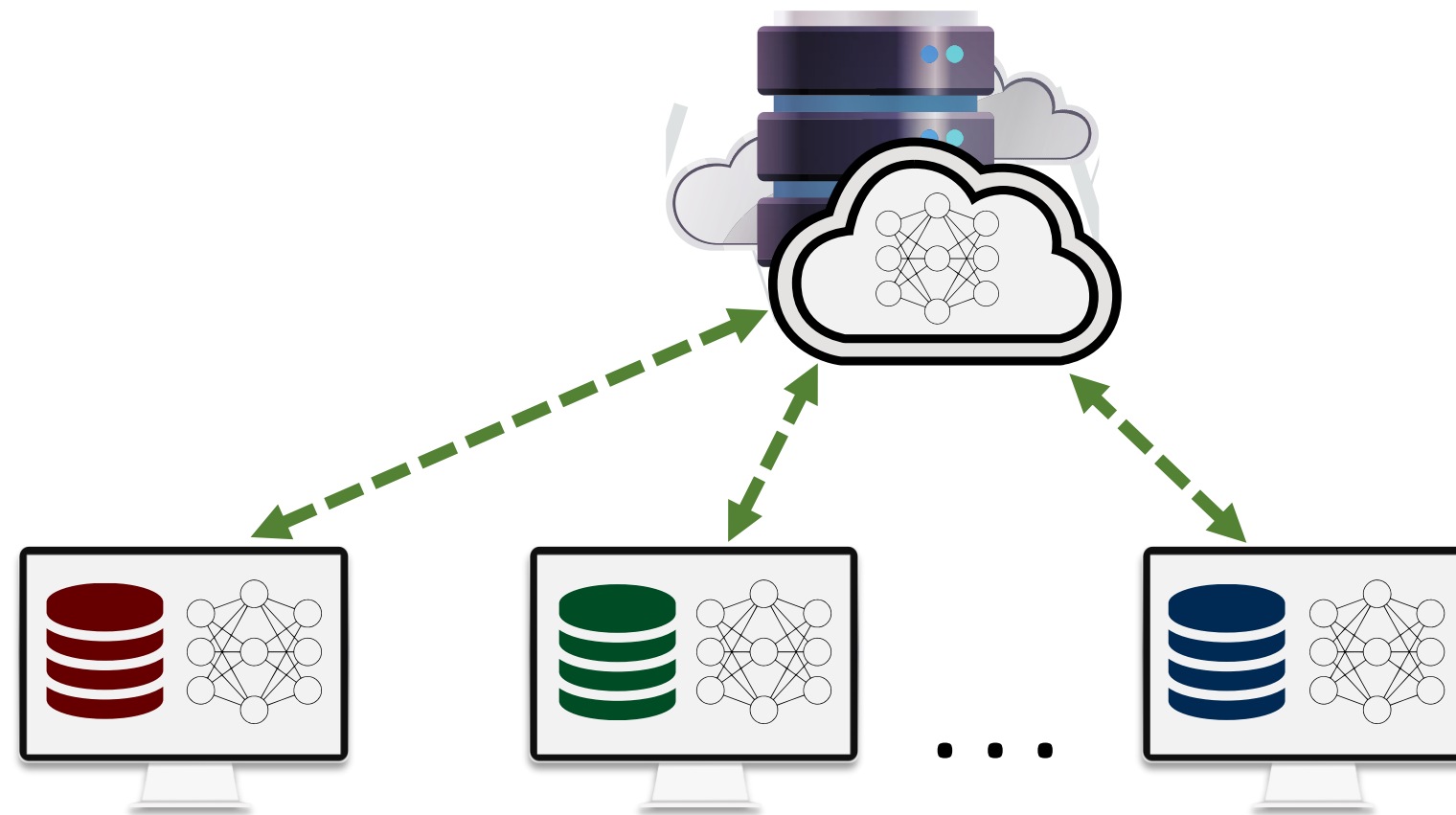
Xun Qian, Rustem Islamov,  
Mher Safaryan, Peter Richtárik  
**Basis Matters: Better Communication-Efficient  
Second Order Methods for Federated Learning,  
*AISTATS 2022***



Rustem Islamov, Xun Qian, Slavomír Hanzely,  
Mher Safaryan, Peter Richtárik  
**Distributed Newton-type methods with  
communication compression and Bernoulli  
aggregation, *TMLR 2023***

$$x^{k+1} = x^k - \left( \frac{1}{n} \sum_{i=1}^n \nabla^2 f_i(x^k) \right)^{-1} \left( \frac{1}{n} \sum_{i=1}^n \nabla f_i(x^k) \right)$$

# Contributions



## ➤ Compressed communication

- Sign (1-bit) compression
- Contractive compression
- Second-order optimization

## ➤ Infrequent communication



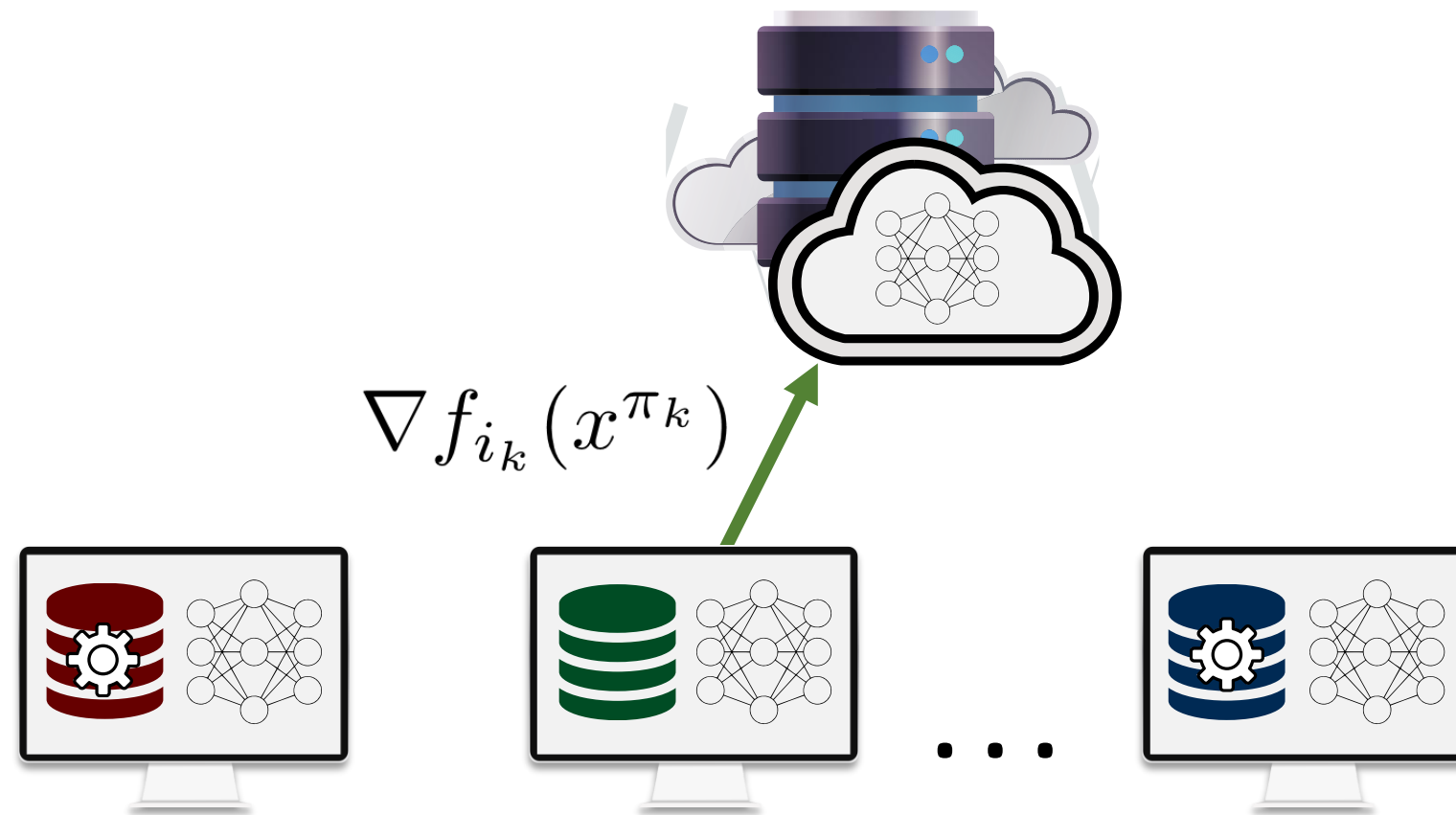
Artavazd Maranjyan, Mher Safaryan, Peter Richtárik  
**Gradskip: Communication-accelerated local  
gradient methods with better computational  
complexity, *Master's thesis, YSU, 2022***

$$x_i^{k+1} = \begin{cases} x_i^k - \gamma g_i^k, & \text{if } c_{k+1} = 0, \\ \frac{1}{n} \sum_{i=1}^n (x_i^k - \gamma g_i^k), & \text{if } c_{k+1} = 1, \end{cases}$$

← No communication!

← Communication step

# Contributions



$$x^{k+1} = x^k - \gamma \nabla f_{i_k}(x^{\pi_k})$$

Index of a machine

Potentially outdated model

## ➤ Compressed communication

- Sign (1-bit) compression
- Contractive compression
- Second-order optimization

## ➤ Infrequent communication

## ➤ Asynchronous communication



Rustem Islamov, Mher Safaryan, Dan Alistarh  
A Sharp Unified Analysis of Asynchronous-SGD  
Algorithms, *Master's thesis, IP Paris, 2023*,  
*AISTATS 2024*



# Distributed **Compressed** Gradient Descent

Compression operator

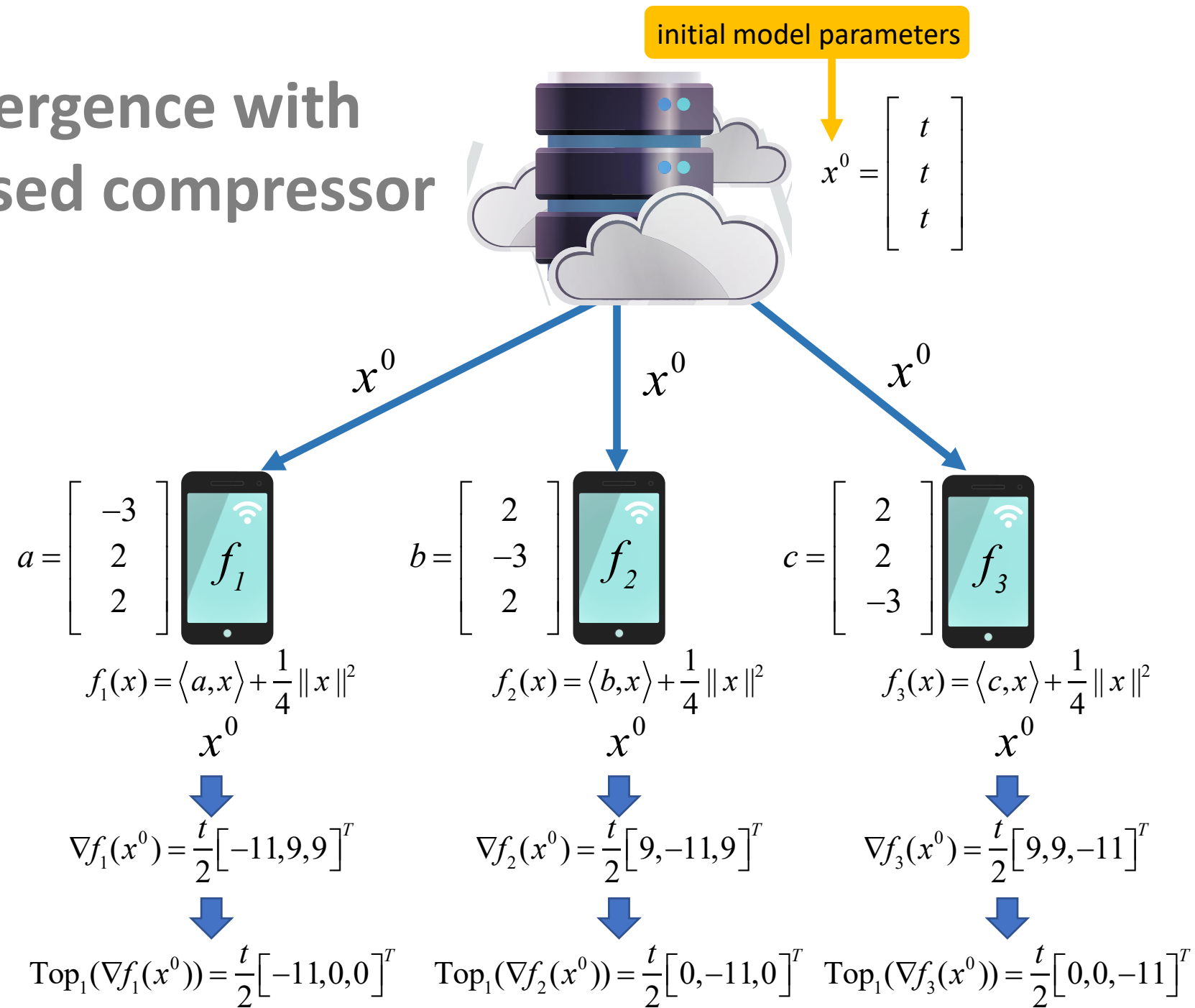
$$\mathcal{C}_i^k : \mathbb{R}^d \rightarrow \mathbb{R}^d$$

$$x^{k+1} = x^k - \frac{\gamma}{n} \sum_{i=1}^n \mathcal{C}_i^k (\nabla f_i(x^k))$$

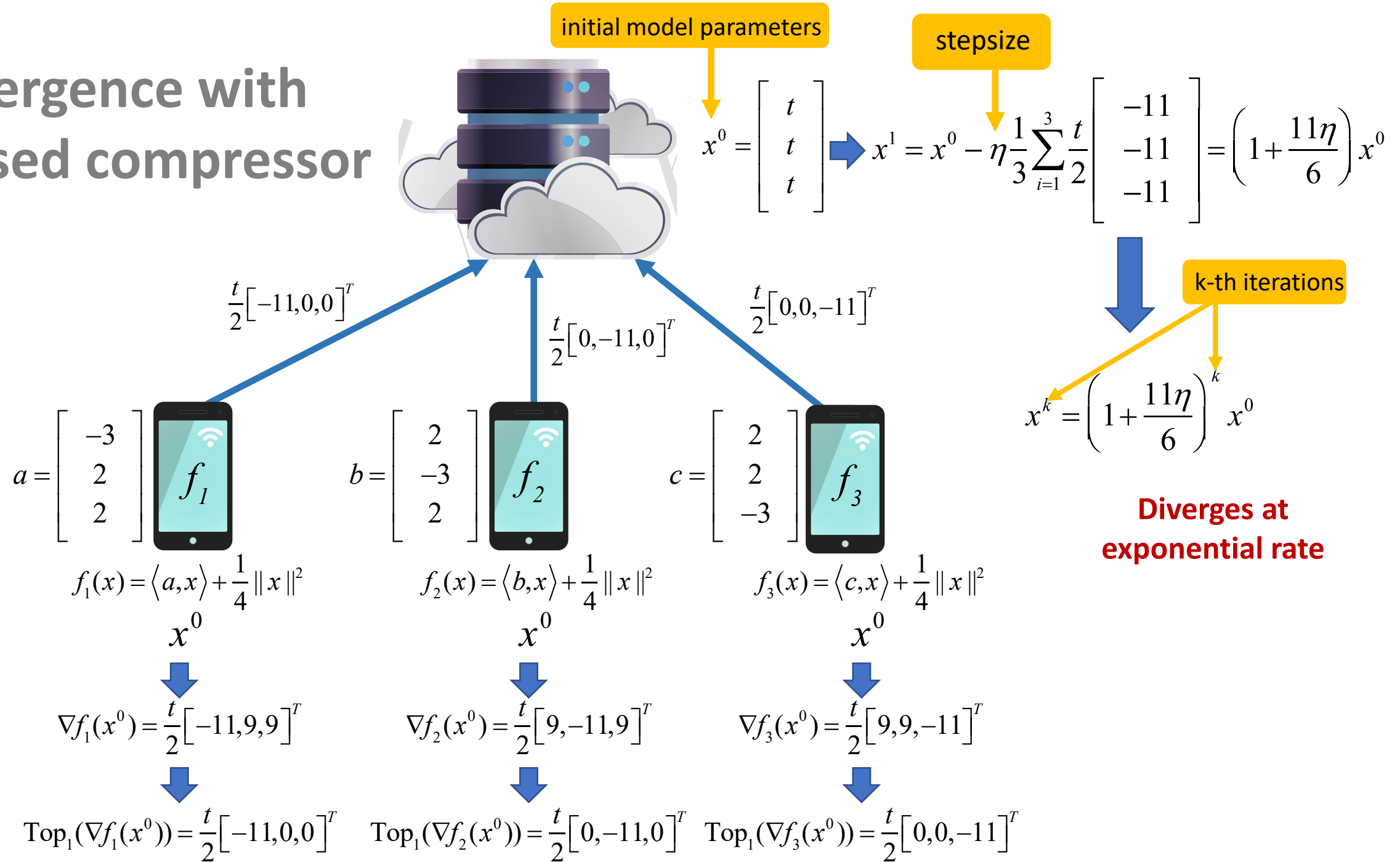


Do we have convergence ?

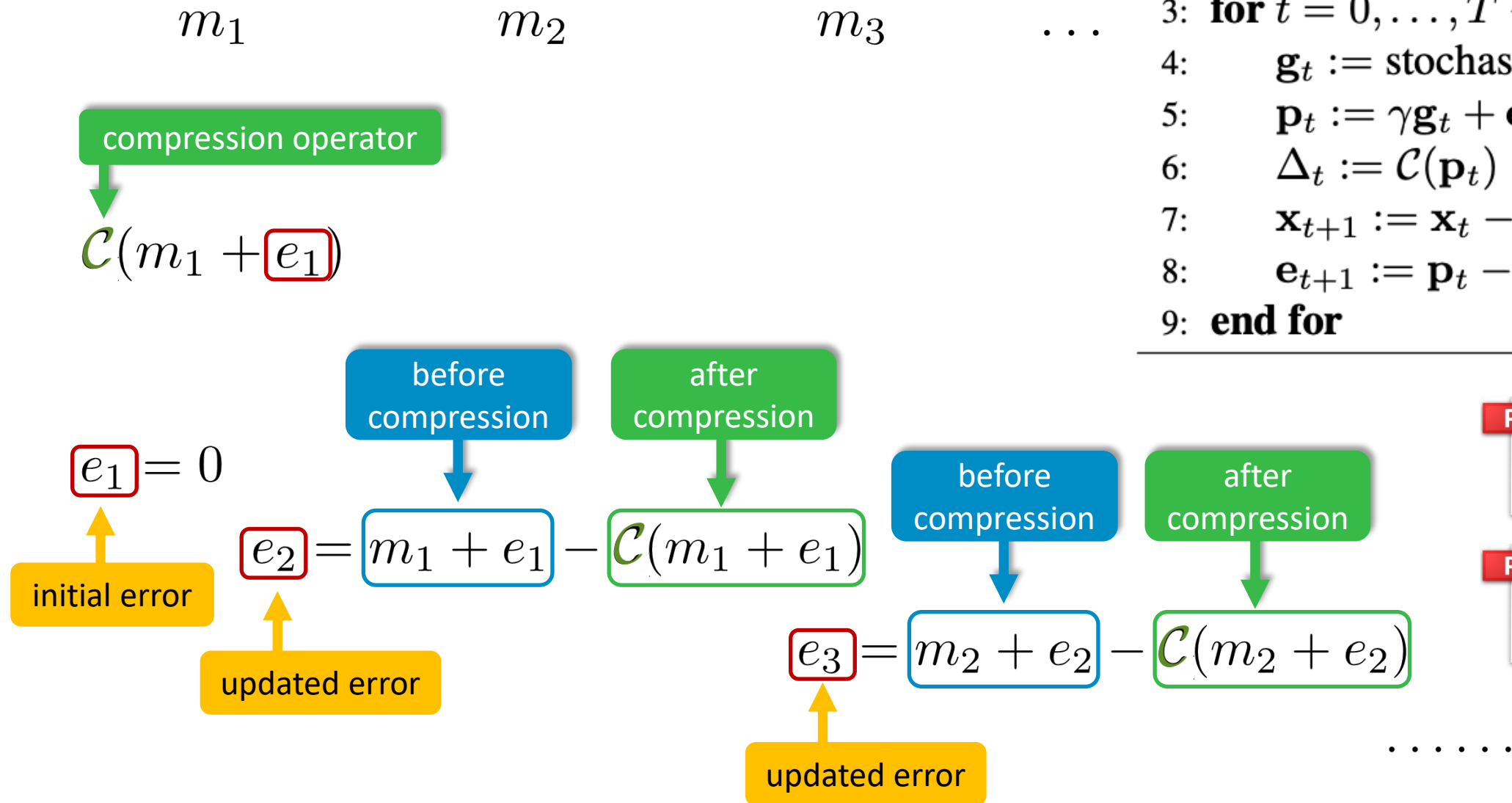
# Divergence with biased compressor



# Divergence with biased compressor



# Error Feedback



## Algorithm 2 EF-SGD (Compr. SGD with Error-Feedback)

- 1: **Input:** learning rate  $\gamma$ , compressor  $\mathcal{C}(\cdot)$ ,  $\mathbf{x}_0 \in \mathbb{R}^d$
- 2: **Initialize:**  $\mathbf{e}_0 = \mathbf{0} \in \mathbb{R}^d$
- 3: **for**  $t = 0, \dots, T - 1$  **do**
- 4:      $\mathbf{g}_t := \text{stochasticGradient}(\mathbf{x}_t)$
- 5:      $\mathbf{p}_t := \gamma \mathbf{g}_t + \mathbf{e}_t$   $\triangleright$  error correction
- 6:      $\Delta_t := \mathcal{C}(\mathbf{p}_t)$   $\triangleright$  compression
- 7:      $\mathbf{x}_{t+1} := \mathbf{x}_t - \Delta_t$   $\triangleright$  update iterate
- 8:      $\mathbf{e}_{t+1} := \mathbf{p}_t - \Delta_t$   $\triangleright$  update residual error
- 9: **end for**



Sebastian U. Stich and Sai Praneeth Karimireddy  
**The error-feedback framework: Better rates for SGD with delayed gradients and compressed communication.**  
 arXiv preprint arXiv:1909.05350, 2019.



Frank Seide, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu  
**1-bit stochastic gradient descent and application to data-parallel distributed training of speech DNNs.**  
 In Interspeech 2014, September 2014.

# Distributed SGD with Biased Compression and Error Feedback

**Parameters:** Compressors  $\mathcal{C}_i^k \in \mathbb{B}^3(\delta)$ ; Stepsizes  $\{\eta^k\}_{k \geq 0}$ ; Iteration count  $K$

**Initialization:** Choose  $x^0 \in \mathbb{R}^d$  and  $e_i^0 = 0$  for all  $i$

**for**  $k = 0, 1, 2, \dots, K$  **do**

Server sends  $x^k$  to all  $n$  machines

All machines in parallel perform these updates:

$$\begin{aligned}\tilde{g}_i^k &= \mathcal{C}_i^k(e_i^k + \eta^k g_i^k) \\ e_i^{k+1} &= e_i^k + \eta^k g_i^k - \tilde{g}_i^k\end{aligned}$$

Each machine  $i$  sends  $\tilde{g}_i^k$  to the server

Server performs aggregation:

$$x^{k+1} = x^k - \frac{1}{n} \sum_{i=1}^n \tilde{g}_i^k$$

**end for**

# Thank you

