

FedNL: Making Newton-Type Methods Applicable to Federated Learning

Mher Safaryan

Postdoctoral Research Fellow



جامعة الملك عبد الله
للعلوم والتقنية
King Abdullah University of
Science and Technology



ICML
International Conference
On Machine Learning



Rustem Islamov

Master's student



Xun Qian

Research Scientist

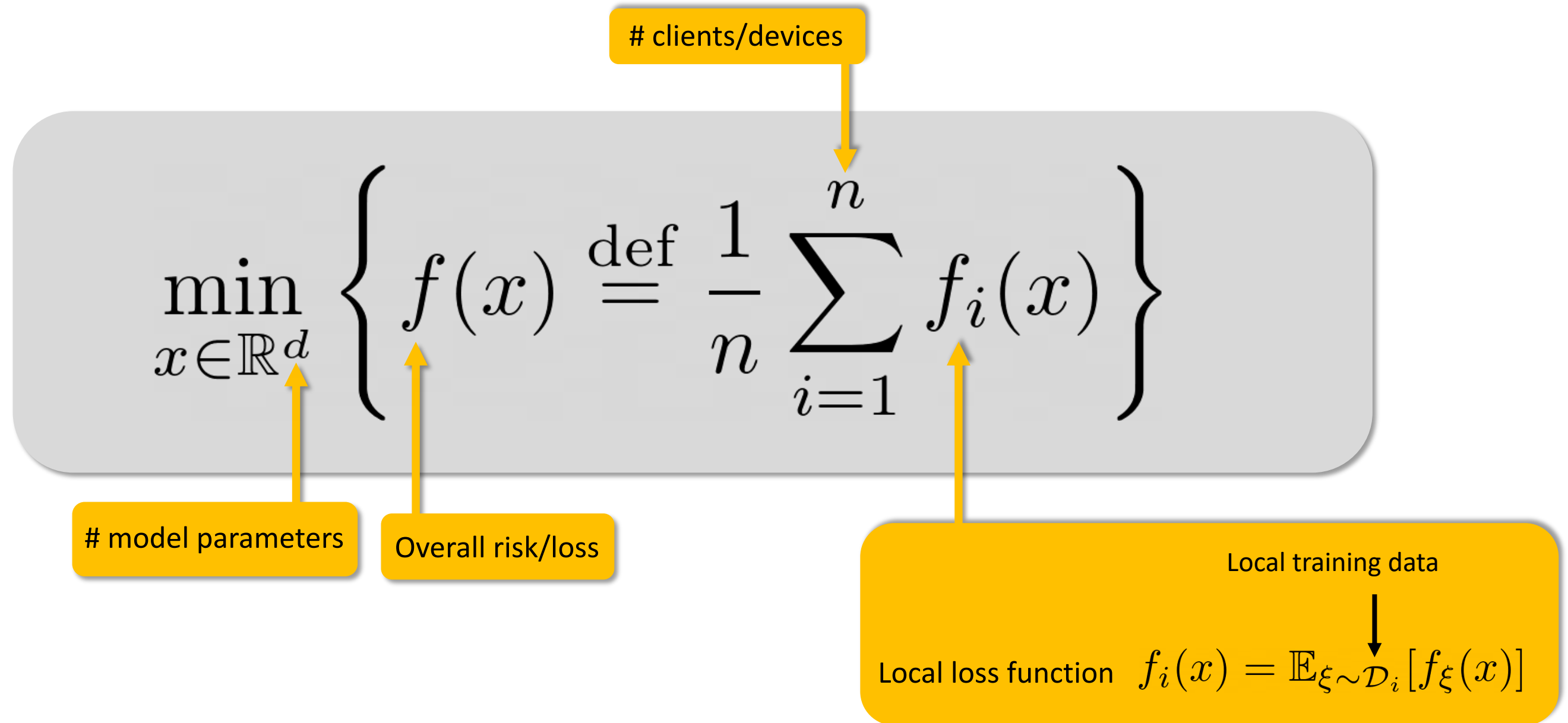


Peter Richtárik

Professor of Computer Science



The Problem



The Problem

Hessians are Lipschitz continuous

$$\|\nabla^2 f_i(x) - \nabla^2 f_i(y)\| \leq L\|x - y\|$$

2nd order smooth
non-convex

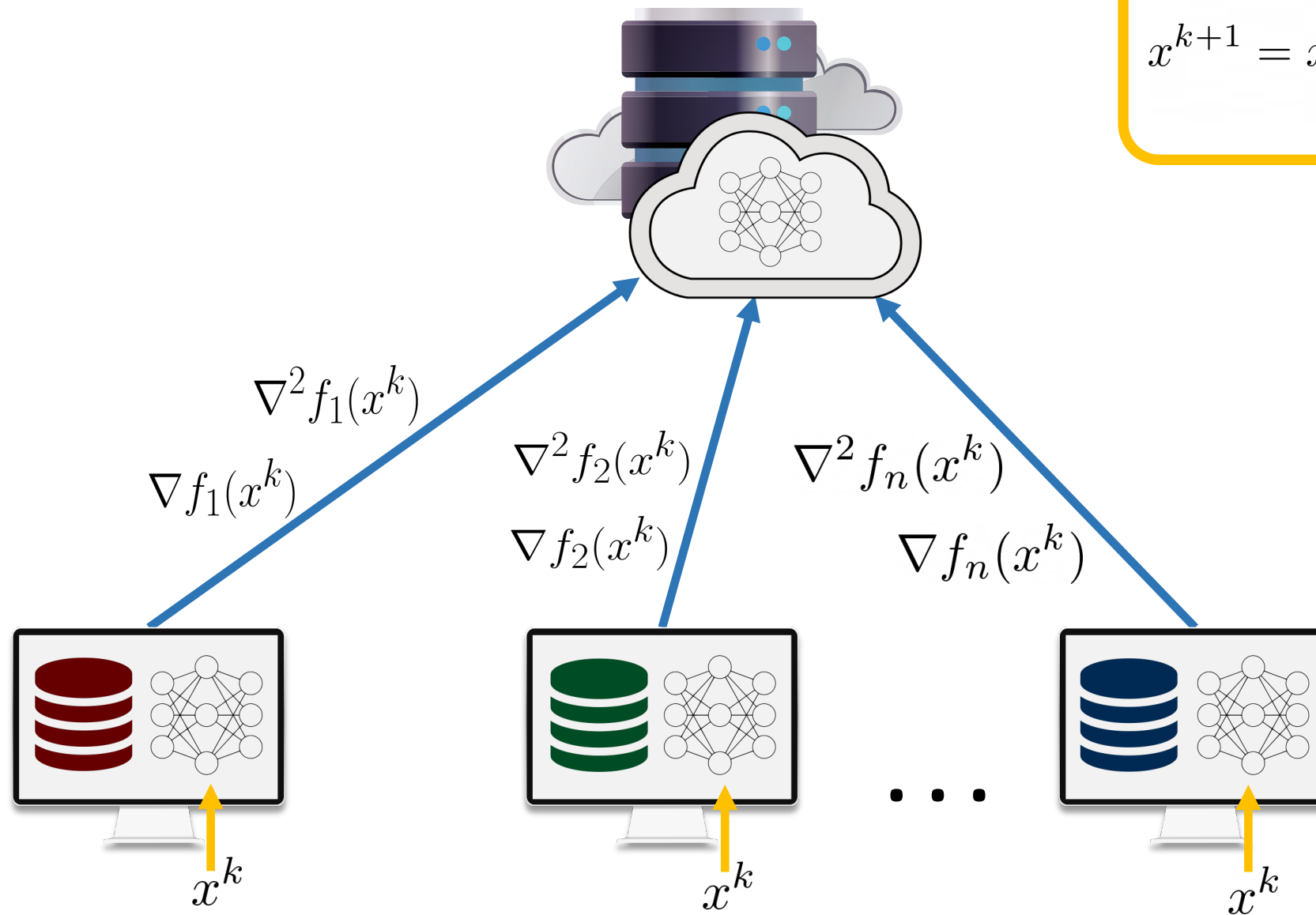
$$\min_{x \in \mathbb{R}^d} \left\{ f(x) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(x) \right\}$$

Strongly convex

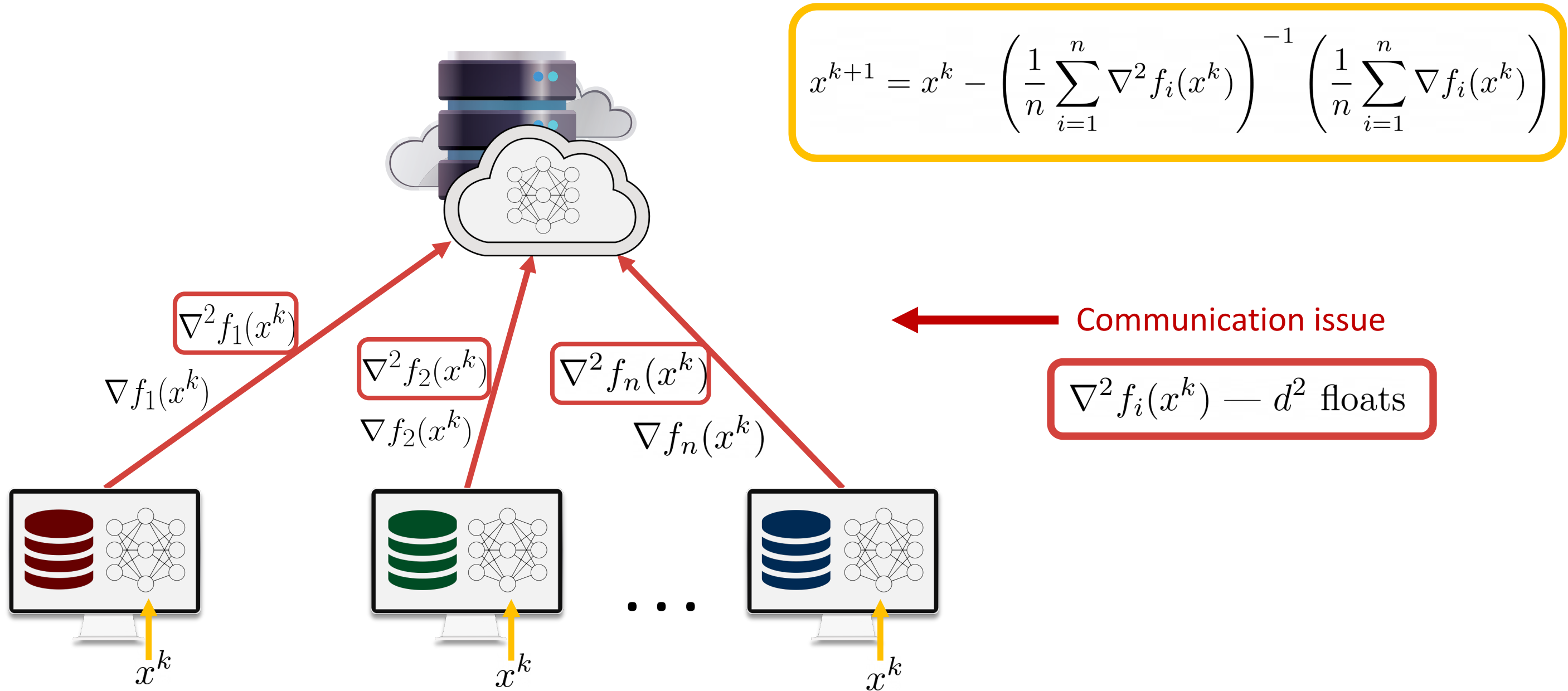
for some $\mu > 0$ and for any $x, y \in \mathbb{R}^d$
$$f(x) \geq f(y) + \langle \nabla f(y), x - y \rangle + \frac{\mu}{2} \|x - y\|^2$$

Distributed Implementation of Newton's method

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



Distributed Implementation of Newton's method



Newton's Method

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

Newton's Method

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

Diagram illustrating the Newton's Method update formula. The formula shows the next iterate x^{k+1} is calculated by subtracting the product of the inverse of the average Hessian and the average gradient from the current iterate x^k .

Annotations for the Hessian term $\nabla^2 f_i(x^k)$:

- Can be computed locally
- Expensive to communicate: $\mathcal{O}(d^2)$

Annotations for the Gradient term $\nabla f_i(x^k)$:

- Can be computed locally
- Easy to communicate: $\mathcal{O}(d)$

- ✗ $\mathcal{O}(d)$ communication cost per round
- ✓ Implementability in practice
- ✓ Local quadratic convergence rate independent of the condition number

$$\|x^{k+1} - x^*\| \leq \frac{L}{2\mu} \|x^k - x^*\|^2$$

Diagram illustrating the local quadratic convergence rate. The inequality shows the distance from the next iterate to the optimal point x^* is bounded by a quadratic term of the current distance, scaled by the ratio of the Hessian Lipschitz constant L to the strong convexity constant 2μ .

Annotations for the inequality:

- Hessian Lipschitz constant (points to L)
- Strong convexity constant (points to 2μ)
- Local quadratic rate (points to the squared term $\|x^k - x^*\|^2$)

Newton **Star** Method

The unique minimizer of $f(x)$

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



Rustem Islamov, Xun Qian and Peter Richtárik
**Distributed second order methods with fast
rates and compressed communication,**
ICML 2021.

Newton **Star** Method

The unique minimizer of $f(x)$

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

Can NOT be computed locally

Single communication of $\mathcal{O}(d^2)$

Can be computed locally

Easy to communicate: $\mathcal{O}(d)$



Rustem Islamov, Xun Qian and Peter Richtárik
Distributed second order methods with fast rates and compressed communication,
ICML 2021.

- ✓ $\mathcal{O}(d)$ communication cost per round
- ✗ Implementability in practice
- ✓ Local quadratic convergence rate independent of the condition number

Hessian Lipschitz constant

$$\|x^{k+1} - x^*\| \leq \frac{L}{2\mu} \|x^k - x^*\|^2$$

Strong convexity constant

Local quadratic rate

Learning the Optimal Hessian Matrices

Newton Star

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

Learning the Optimal Hessian Matrices

Newton Star

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

Idea! Learn the optimal Hessians $\nabla^2 f_i(x^*)$ in communication efficient manner:

(i) $\mathbf{H}_i^k \rightarrow \nabla^2 f_i(x^*)$ as $k \rightarrow \infty$ (ii) $\mathbf{H}_i^{k+1} - \mathbf{H}_i^k$ is compressed

$$\begin{aligned} x^{k+1} &= x^k - \left(\frac{1}{n} \sum_{i=1}^n \mathbf{H}_i^k \right)^{-1} \nabla f(x^k) \\ &= x^k - \left(\mathbf{H}^k \right)^{-1} \nabla f(x^k) \end{aligned}$$

FedNL: Two Options for Updating the Global Model

Option 1

$$x^{k+1} = x^k - \left(\begin{bmatrix} \mathbf{H}^k \\ \mu \end{bmatrix} \right)^{-1} \nabla f(x^k)$$

Projection onto the cone
of positive definite
matrices

Option 2

$$x^{k+1} = x^k - \left(\mathbf{H}^k + l^k \mathbf{I} \right)^{-1} \nabla f(x^k)$$

$$l^k = \frac{1}{n} \sum_{i=1}^n \|\mathbf{H}_i^k - \nabla^2 f_i(x^k)\|_F$$

FedNL: Hessian Learning Rate Options

$$\mathbf{H}_i^{k+1} = \mathbf{H}_i^k + \alpha \mathcal{C}_i^k (\nabla^2 f_i(x^k) - \mathbf{H}_i^k)$$

Stepsize depends only on
the compression, e.g.,

$$\alpha = 1$$

$$\alpha = 1 - \sqrt{1 - \delta}$$

$$\alpha = \frac{1}{\omega + 1}$$

FedNL: New Hessian Learning Technique

$$\mathbf{H}_i^{k+1} = \mathbf{H}_i^k + \alpha \mathcal{C}_i^k (\nabla^2 f_i(x^k) - \mathbf{H}_i^k)$$

Compression operator

Contractive compressor $\mathbb{C}(\delta)$, $\delta \in [0, 1)$

$$\|\mathcal{C}(\mathbf{M})\|_F \leq \|\mathbf{M}\|_F$$

$$\|\mathcal{C}(\mathbf{M}) - \mathbf{M}\|_F^2 \leq (1 - \delta) \|\mathbf{M}\|_F^2 \quad \forall \mathbf{M} \in \mathbb{R}^{d \times d}$$

Unbiased compressor $\mathbb{B}(\omega)$, $\omega \geq 0$

$$\mathbb{E}[\mathcal{C}(\mathbf{M})] = \mathbf{M}$$

$$\mathbb{E} \left[\|\mathcal{C}(\mathbf{M}) - \mathbf{M}\|_F^2 \right] \leq \omega \|\mathbf{M}\|_F^2 \quad \forall \mathbf{M} \in \mathbb{R}^{d \times d}$$

Greedy sparsification (Top- K)

$$\begin{bmatrix} -0.4 & 12.1 & 0.76 \\ 2.8 & -9.7 & -1.1 \\ 0.24 & 4.5 & 0.9 \end{bmatrix} \xrightarrow{K=2} \begin{bmatrix} 0 & 12.1 & 0 \\ 0 & -9.7 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Random sparsification (Rand- K)

$$\begin{bmatrix} -0.4 & 12.1 & 0.76 \\ 2.8 & -9.7 & -1.1 \\ 0.24 & 4.5 & 0.9 \end{bmatrix} \xrightarrow{K=2} \frac{9}{2} \begin{bmatrix} -0.4 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 4.5 & 0 \end{bmatrix}$$

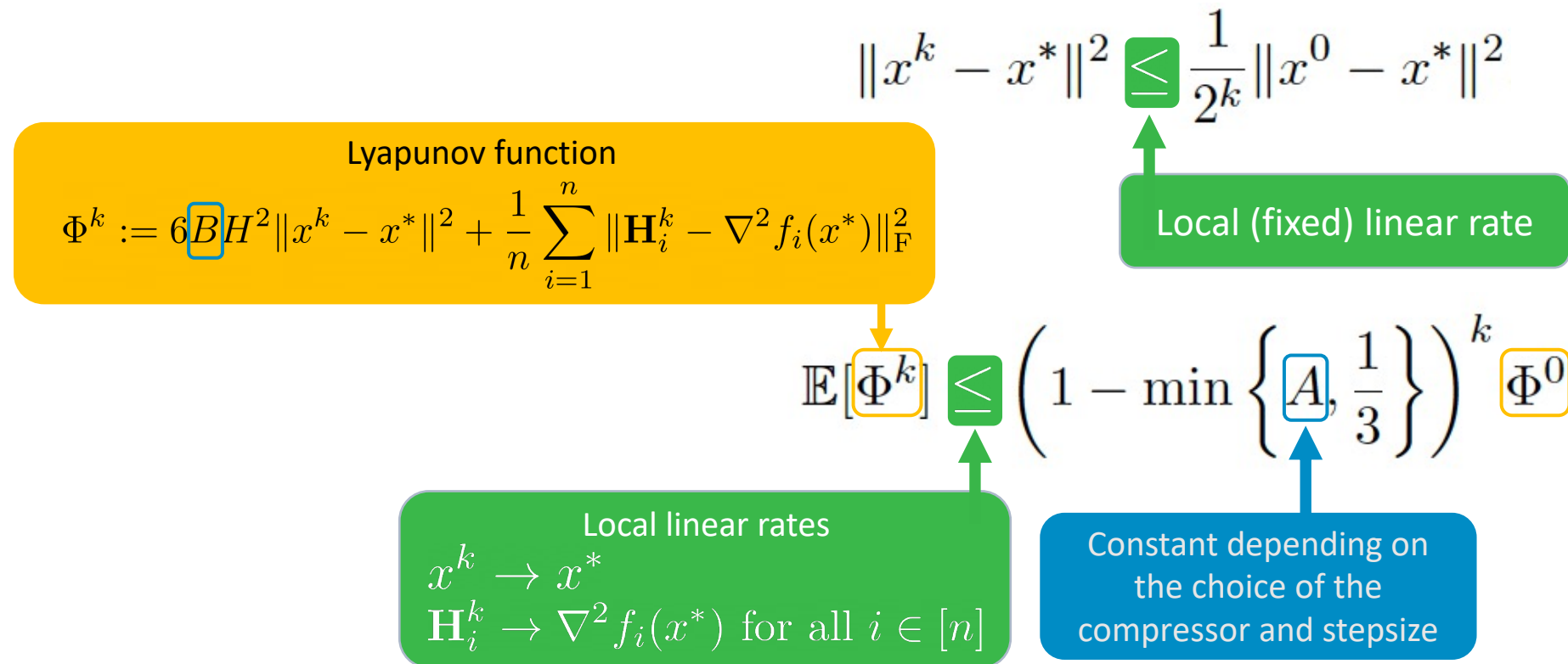
factor preserving
unbiasedness

FedNL: Local Convergence Theory

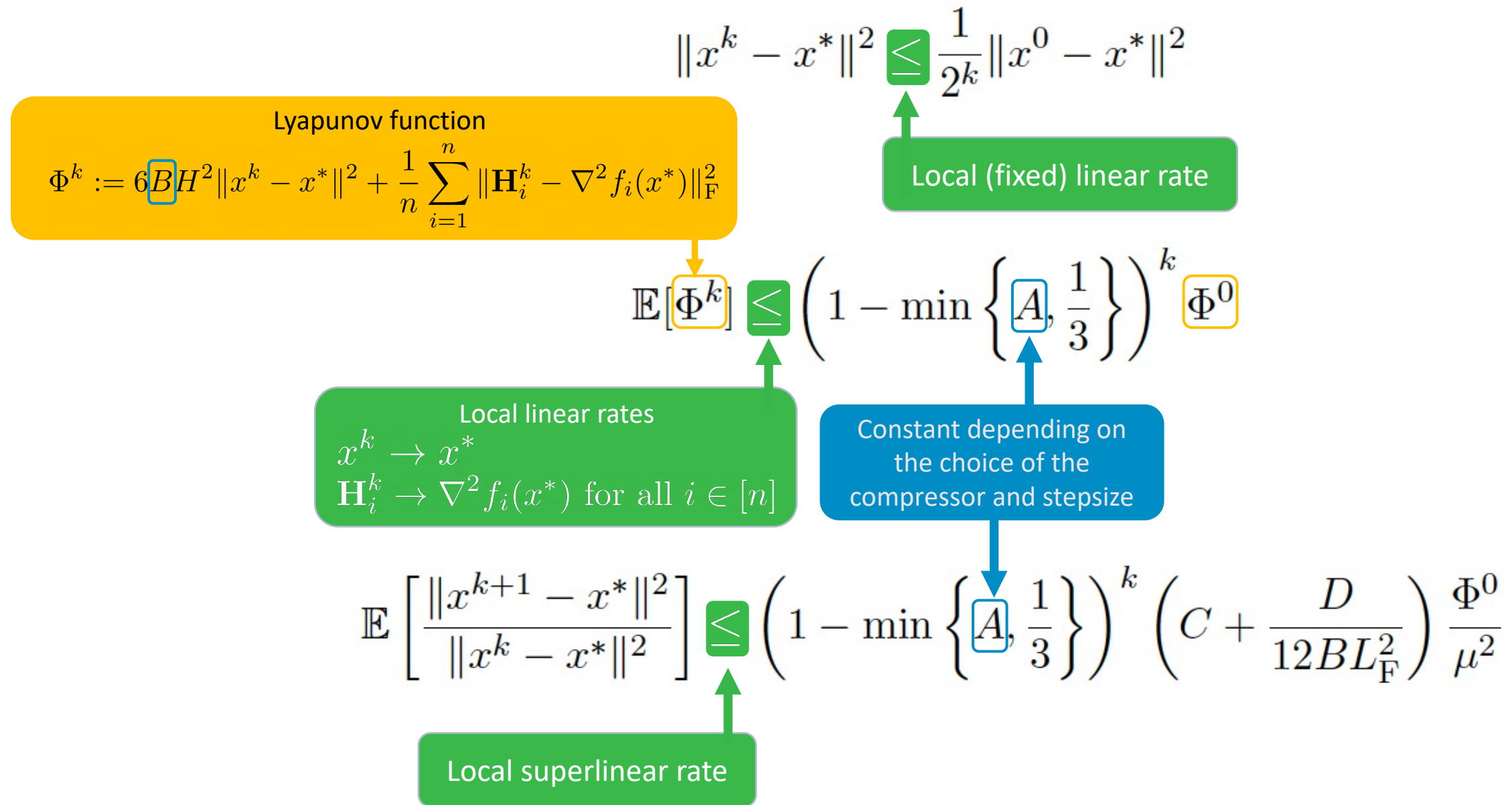
$$\|x^k - x^*\|^2 \leq \frac{1}{2^k} \|x^0 - x^*\|^2$$

Local (fixed) linear rate

FedNL: Local Convergence Theory



FedNL: Local Convergence Theory



Experiments: Regularized Logistic Regression

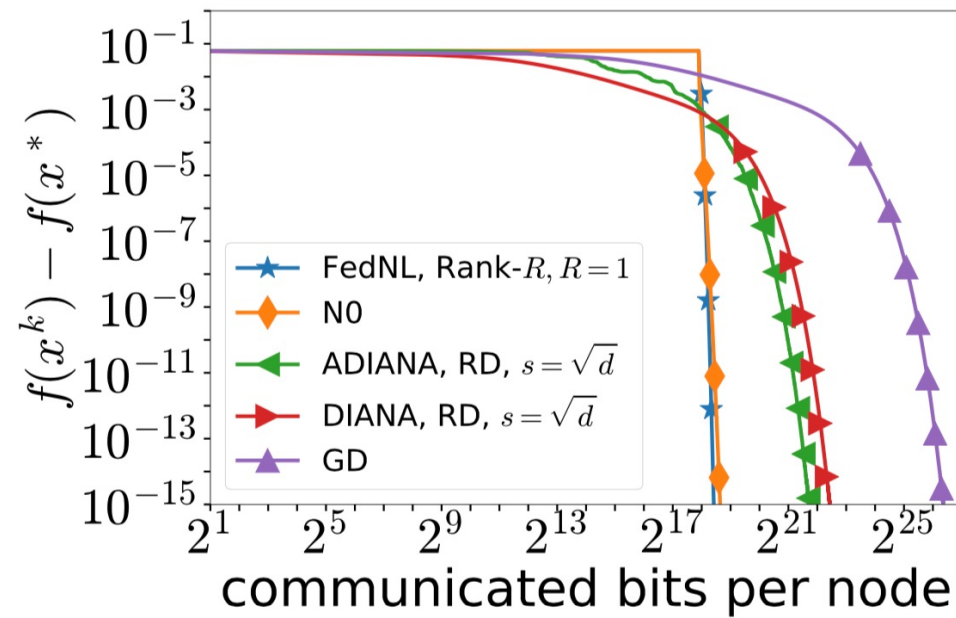
$$\min_{x \in \mathbb{R}^d} \left\{ f(x) := \frac{1}{n} \sum_{i=1}^n f_i(x) + \frac{\lambda}{2} \|x\|^2 \right\}, \quad f_i(x) = \frac{1}{m} \sum_{j=1}^m \log \left(1 + \exp(-b_{ij} a_{ij}^\top x) \right),$$

Regularization parameter λ points to the regularization term $\frac{\lambda}{2} \|x\|^2$.

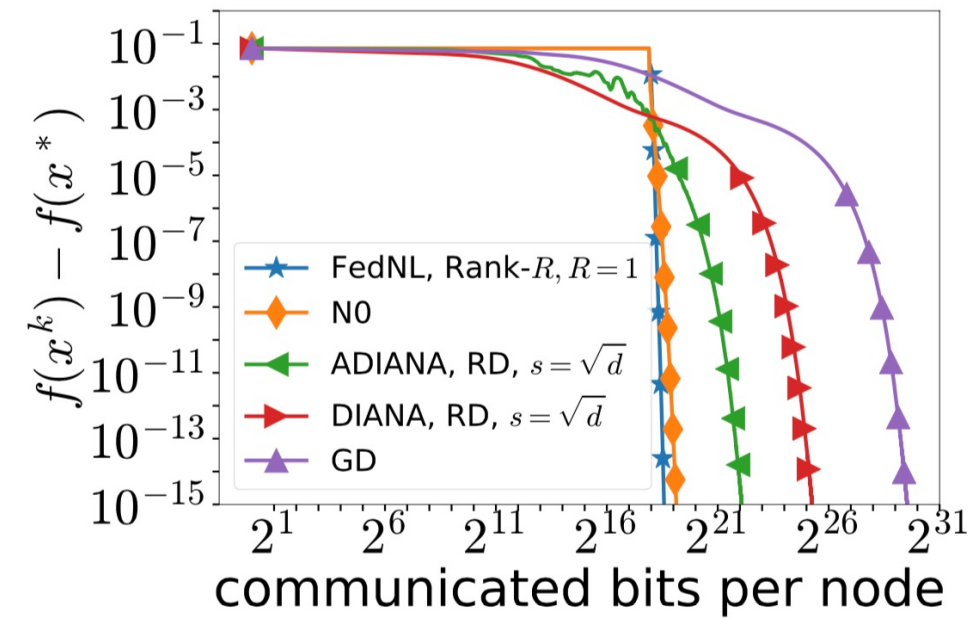
Training data points a_{ij} and b_{ij} point to the data points in the loss function $f_i(x)$.

where $\{a_{ij}, b_{ij}\}_{j \in [m]}$ are data points at the i -th device. The datasets were taken from LibSVM library [Chang and Lin, 2011]: [a1a](#), [a9a](#), [w7a](#), [w8a](#), and [phishing](#).

Experiments: FedNL vs Gradient-Type Methods

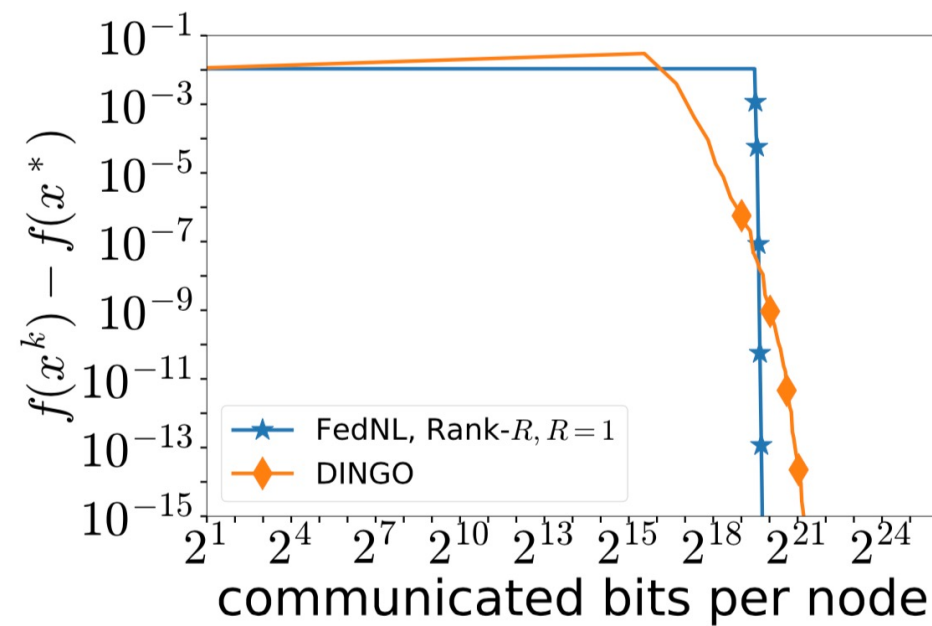


(a) **a1a**, $\lambda = 10^{-3}$

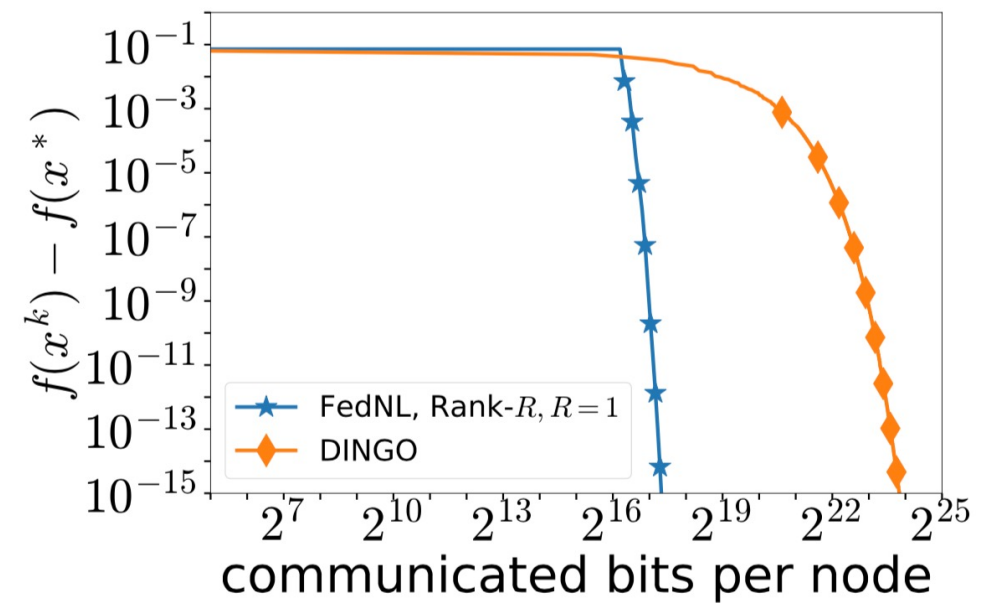


(b) **a9a**, $\lambda = 10^{-4}$

Experiments: FedNL vs DINGO

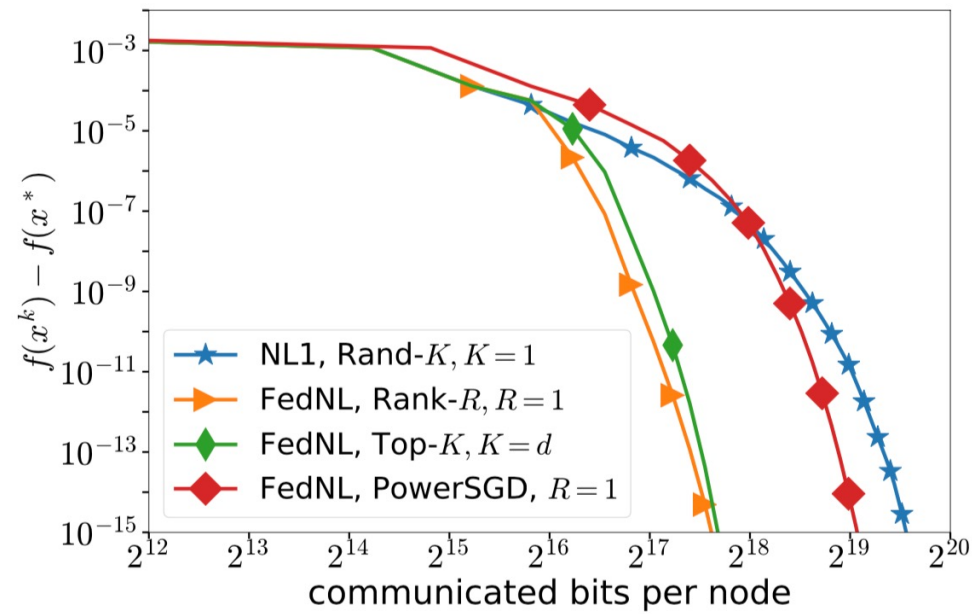


(c) **w8a**, $\lambda = 10^{-3}$

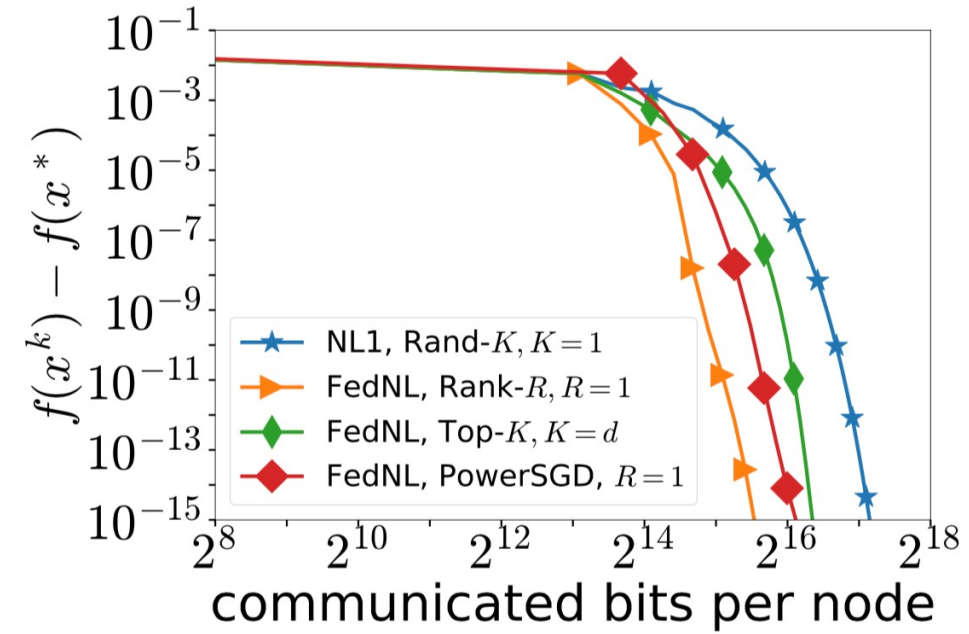


(d) **phishing**, $\lambda = 10^{-4}$

Experiments: FedNL vs NEWTON-LEARN (NL)



(a) **w8a**, $\lambda = 10^{-3}$



(b) **phishing**, $\lambda = 10^{-3}$

Thank you

