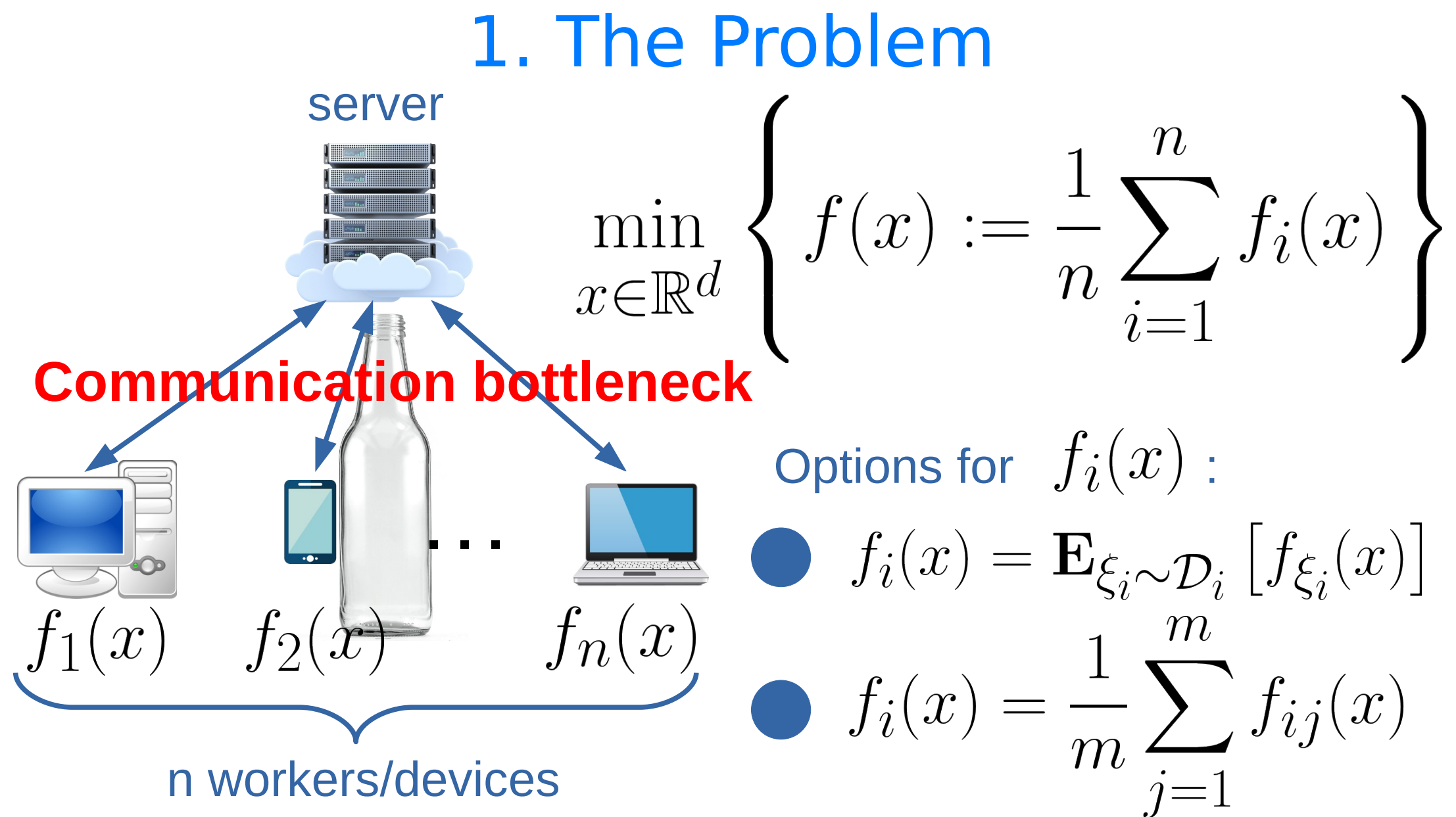


Local SGD: Unified Theory and New Efficient Methods



Eduard Gorbunov^{1,2,3}, Filip Hanzely⁴, Peter Richtárik³
¹MIPT (Russia), ²Yandex (Russia), ³KAUST (Saudi Arabia), ⁴TTIC (United States)

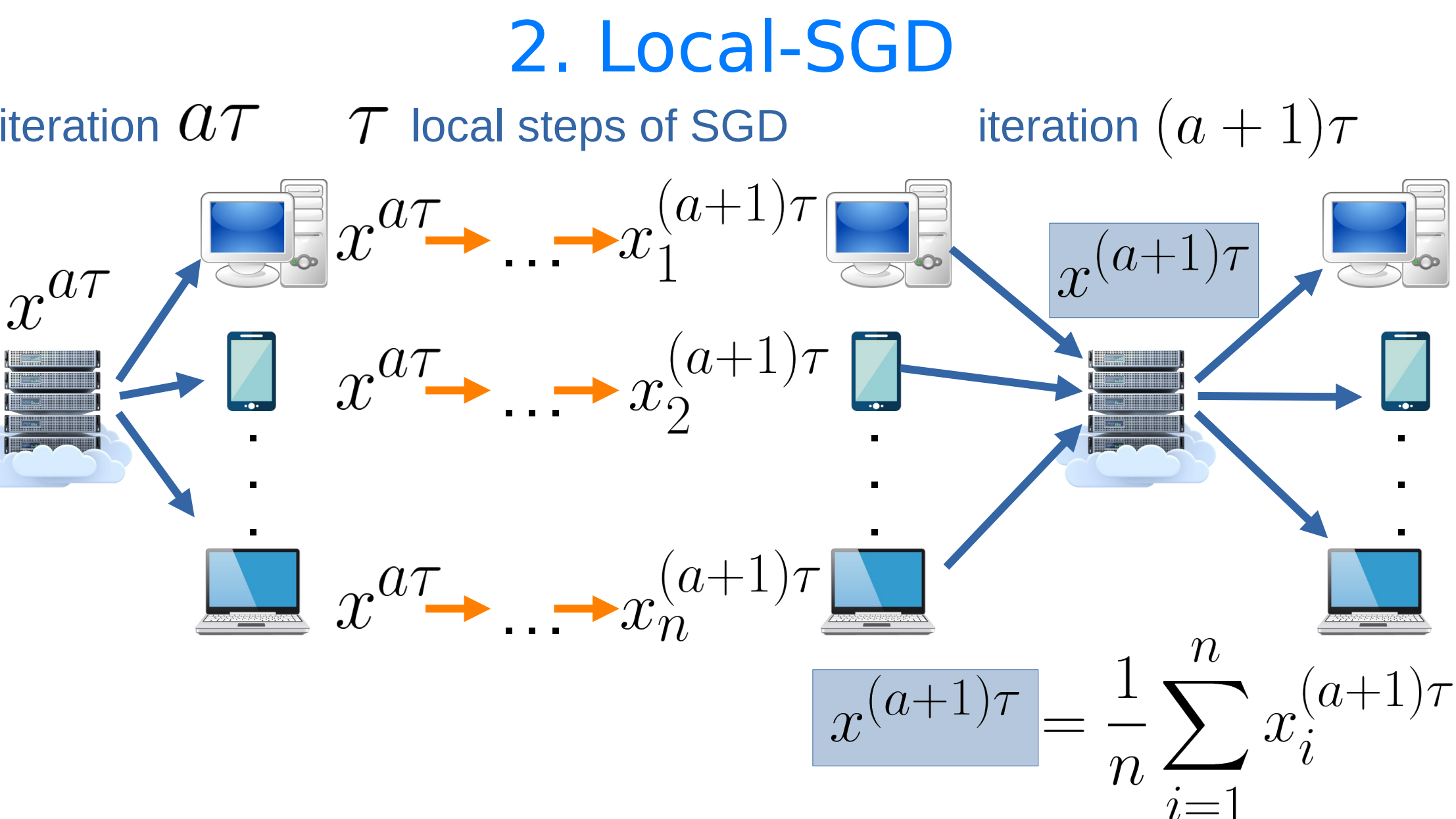


Problem: Distributed optimization / training, where n workers (devices/clients) jointly solve a problem by communicating with a central server.

Assumptions: Smoothness and quasi-strong convexity of local loss functions:

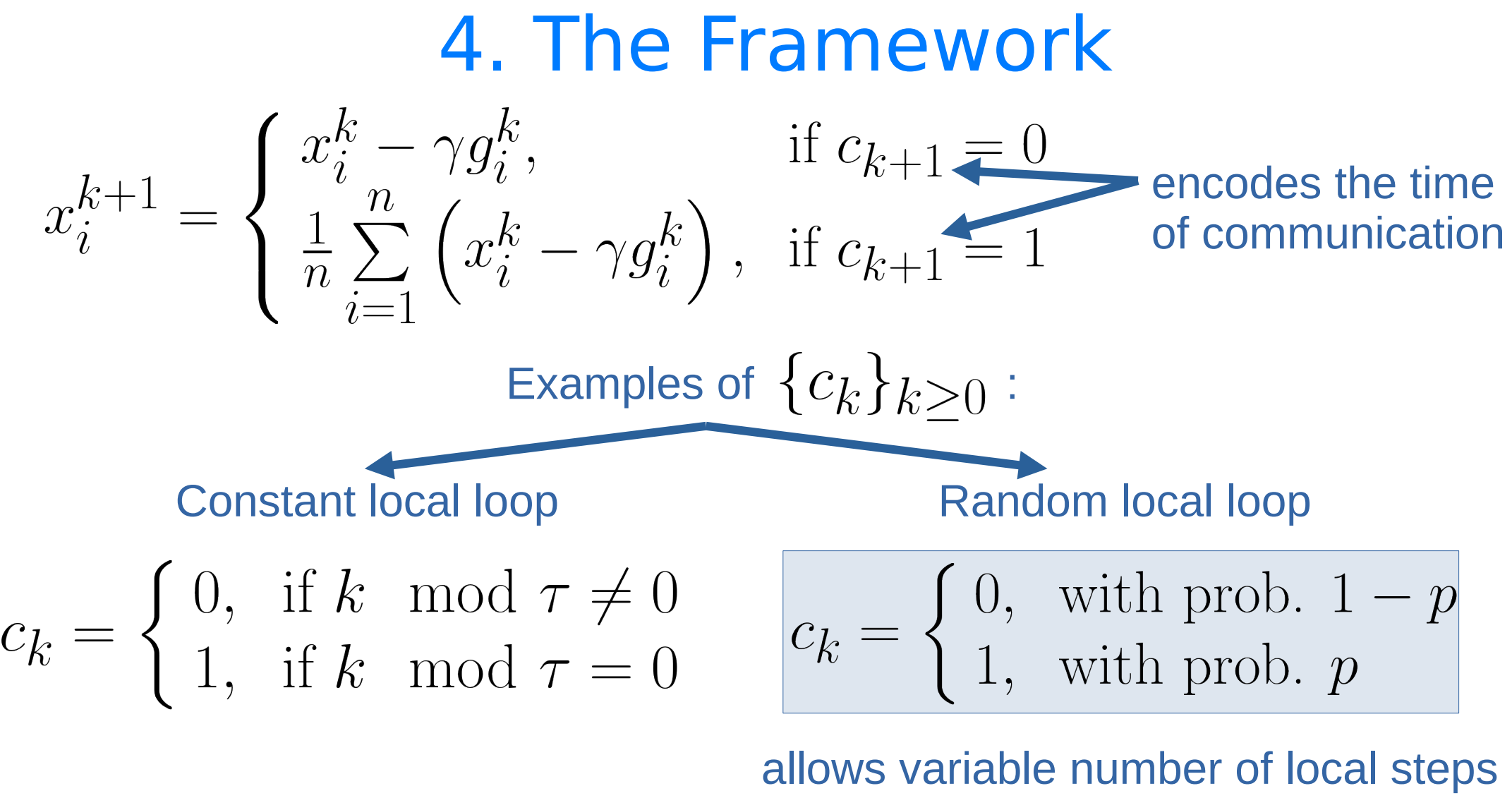
$$\|\nabla f_i(x) - \nabla f_i(y)\| \leq L\|x - y\|$$

$$f_i(x^*) \geq f_i(x) + \langle \nabla f_i(x), x^* - x \rangle + \frac{\mu}{2} \|x - x^*\|^2$$



- ✔ A lot of results are already known...
- ✘ ... but many fruitful directions are **unexplored**
- better understanding of the local shifts
- importance sampling
- variance reduction
- variable number of local steps
- general theory for multiple data similarity types

- ### 3. Our Contributions
- General theoretical framework for local first-order methods covering
 - local shifts
 - importance sampling
 - variance reduction
 - variable number of local steps
 - Our framework
 - recovers tight rates for known optimizers
 - fills missing gaps for known methods
 - extends the established optimizers
 - New efficient methods
 - S-Local-SVRG – the first linearly converging stochastic method with local updates in **heterogeneous case**



The assumption below covers a very broad class of methods.

$$\frac{1}{n} \sum_{i=1}^n \mathbf{E} [g_i^k | x_1^k, \dots, x_n^k] = \frac{1}{n} \sum_{i=1}^n \nabla f_i(x_i^k) \quad V_k = \frac{1}{n} \sum_{i=1}^n \|x_i^k - x^k\|^2$$

$$\frac{1}{n} \sum_{i=1}^n \mathbf{E} [\|g_i^k\|^2] \leq 2A \mathbf{E} [f(x^k) - f(x^*)] + B \mathbf{E} [\sigma_k^2] + F \mathbf{E} [V_k] + D_1$$

$$\mathbf{E} \left[\left\| \frac{1}{n} \sum_{i=1}^n g_i^k \right\|^2 \right] \leq 2A' \mathbf{E} [f(x^k) - f(x^*)] + B' \mathbf{E} [\sigma_k^2] + F' \mathbf{E} [V_k] + D_1'$$

$$\mathbf{E} [\sigma_{k+1}^2] \leq (1 - \rho) \mathbf{E} [\sigma_k^2] + 2C \mathbf{E} [f(x^k) - f(x^*)] + G \mathbf{E} [V_k] + D_2$$

$$2L \sum_{k=0}^K w_k \mathbf{E} [V_k] \leq \frac{1}{2} \sum_{k=0}^K w_k \mathbf{E} [f(x^k) - f(x^*)] + 2LH \mathbf{E} \sigma_0^2 + 2LD_3 \gamma^2 W_K$$

- Reflects smoothness properties of the problem and noises introduced by stochastic gradients and functions dissimilarity
 - Describes the process of local shifts' learning and variance reduction
 - Bounds the workers iterates' discrepancy
- weights: $w_k = \frac{1}{(1 - \min\{\gamma\mu, \frac{\rho}{4}\})^{k+1}} \quad W_K = \sum_{k=0}^K w_k$

5. Main Theorem: Simplified Version

$$\bar{x}^K = \frac{1}{W_K} \sum_{k=0}^K w_k x^k$$

depends only on the starting point and stepsize

linear function

$$\mathbf{E} [f(\bar{x}^K)] - f(x^*) \leq \left(1 - \min\left\{\gamma\mu, \frac{\rho}{4}\right\}\right)^K \frac{\Phi^0(x^0, \gamma)}{\gamma} + \gamma \Psi^0(D_1, D_2, D_3)$$

6. New Method: Shifted Local-SVRG

Finite-sum case: $f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x)$

$$x_i^{k+1} = \begin{cases} x_i^k - \gamma g_i^k, & \text{with prob. } 1-p \\ \frac{1}{n} \sum_{i=1}^n (x_i^k - \gamma g_i^k), & \text{with prob. } p \end{cases}$$

$$g_i^k = \nabla f_{i, \bar{j}_i}(x_i^k) - \nabla f_{i, \bar{j}_i}(y^k) + \nabla f(y^k) \quad \bar{j}_i \sim \{1, \dots, m\} \text{ uniformly at random}$$

$$y^{k+1} = \begin{cases} x^k, & \text{with prob. } q \\ y^k, & \text{with prob. } 1-q \end{cases} \quad q = \frac{1}{m}$$

Iteration complexity: $\mathcal{O} \left(\left(m + \frac{L}{p\mu} + \frac{\max L_{ij}}{n\mu} + \frac{\sqrt{(1-p)L \max L_{ij}}}{p\mu} \right) \log \frac{1}{\epsilon} \right)$

The first linearly converging local method for heterogeneous data!
 It is just an example. In fact, our approach covers a lot of different setups, methods and even the algorithms without local updates.

