

Linearly Converging Error Compensated SGD

Eduard Gorbunov

MIPT, Yandex and Sirius

Dmitry Kovalev
KAUST

Dmitry Makarenko
MIPT

Peter Richtárik
KAUST

Yandex



KAUST



Университет
Сириус



Dmitry Kovalev
PhD student
KAUST

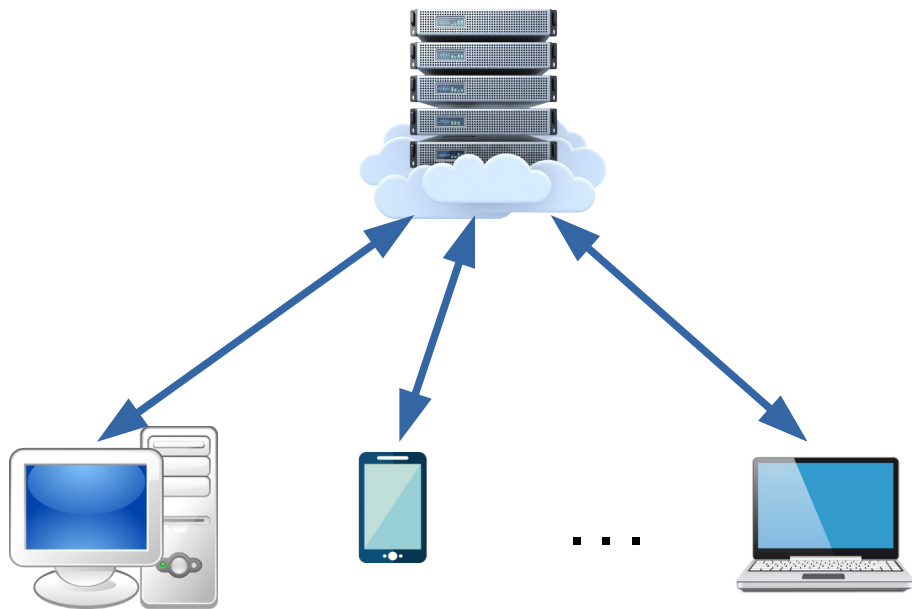


Dmitry Makarenko
PhD student
MIPT

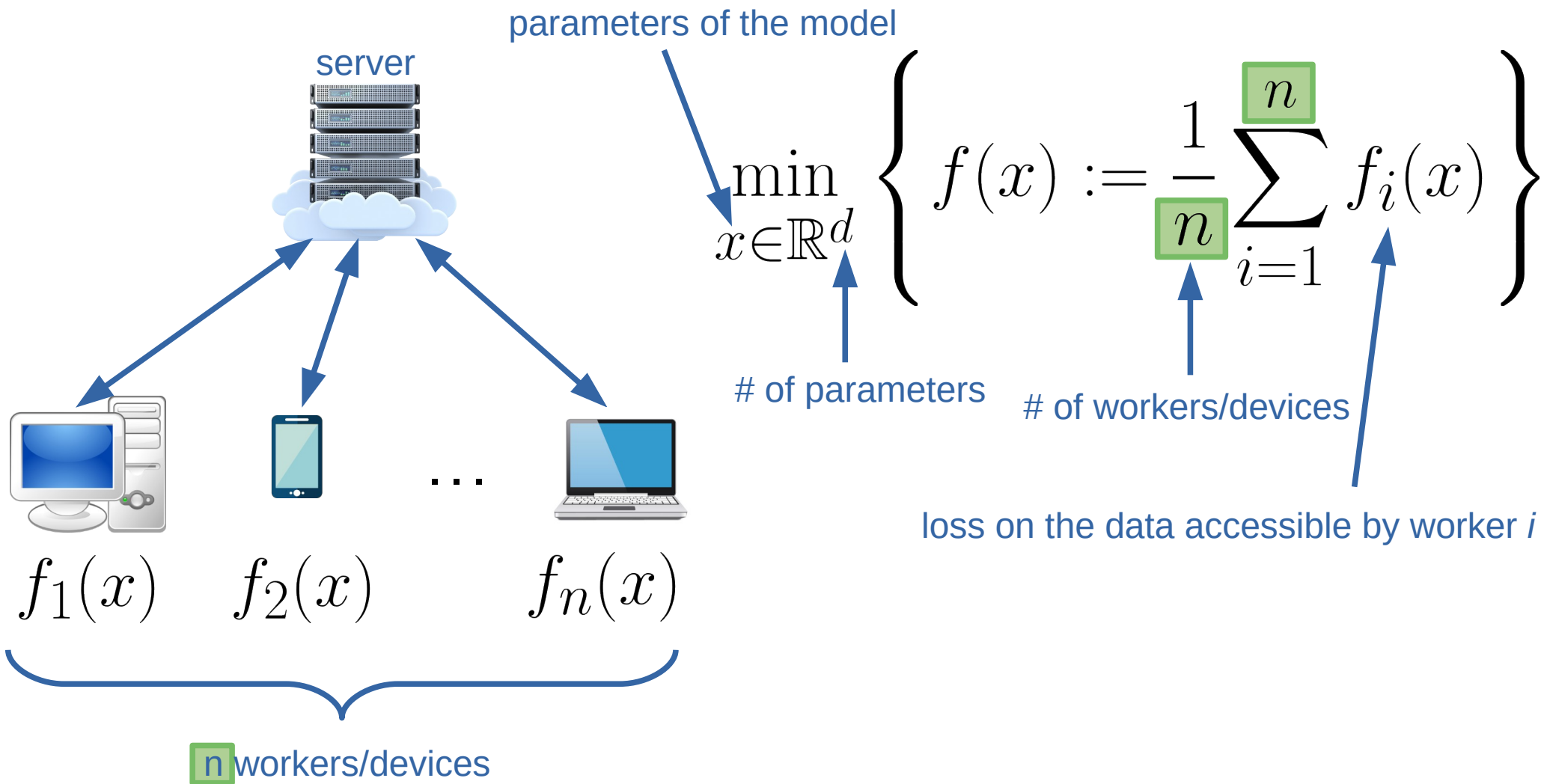


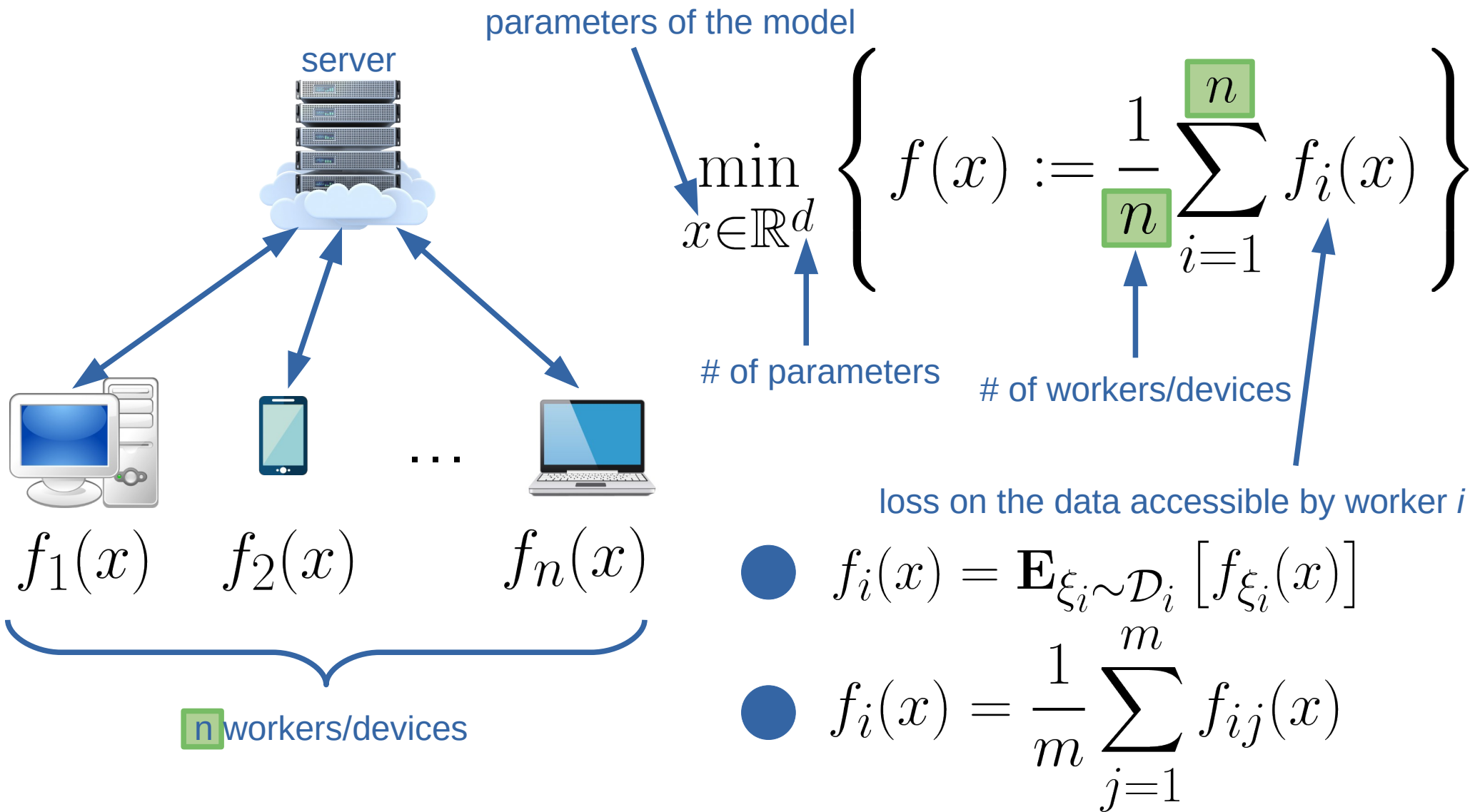
Peter Richtárik
Professor of Computer Science
KAUST

1. The Problem



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





Assumptions

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

$$f_i(x) - f_i(y) \geq \langle \nabla f_i(y), x - y \rangle$$

● f_1, f_2, \dots, f_n – L-smooth and convex



Assumptions

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

$$f_i(x) - f_i(y) \geq \langle \nabla f_i(y), x - y \rangle$$

● f_1, f_2, \dots, f_n – L-smooth and convex

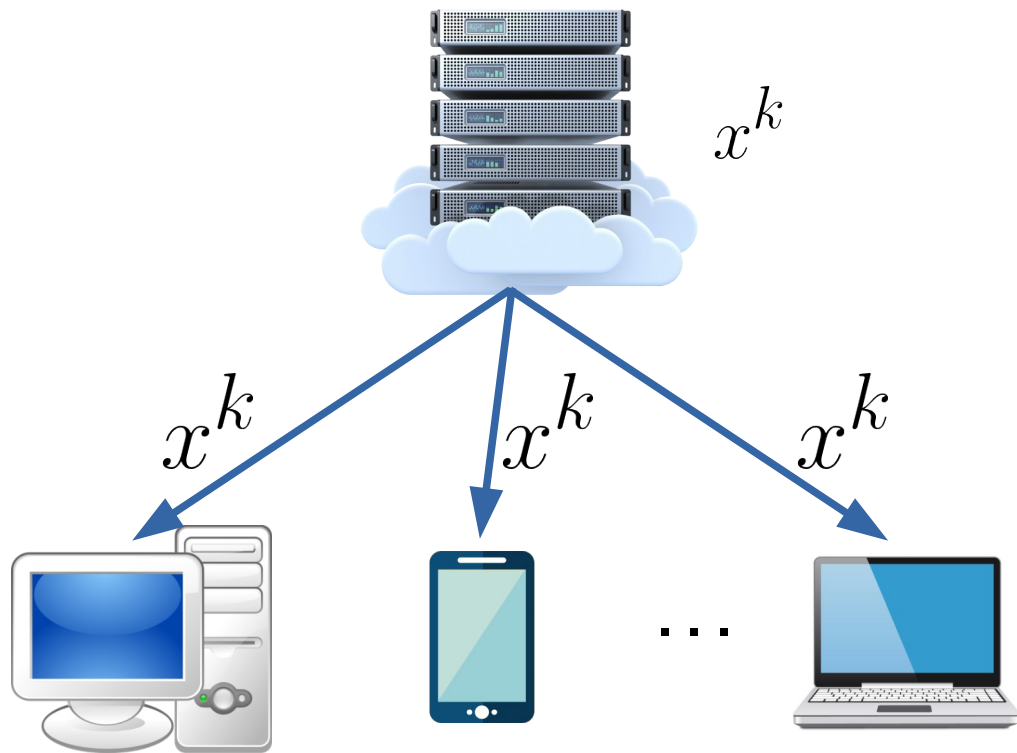
● f – strongly quasi-convex

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

the solution of the problem

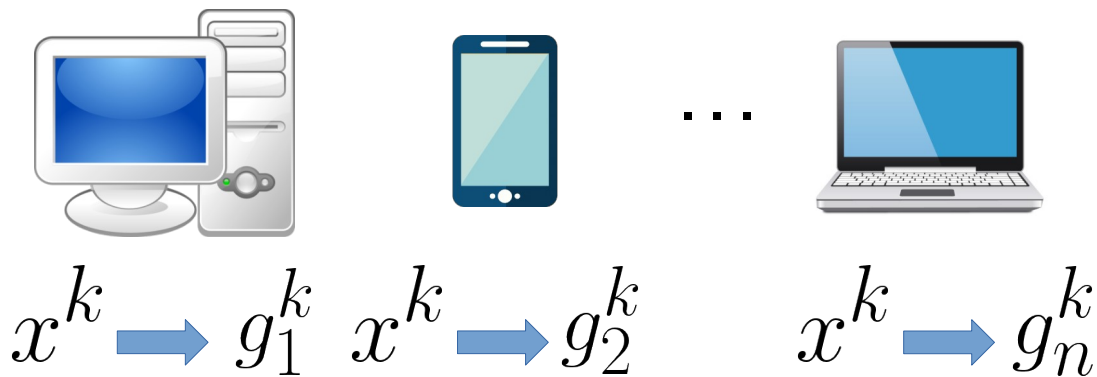
2. Parallel SGD

1 Server broadcasts the parameters



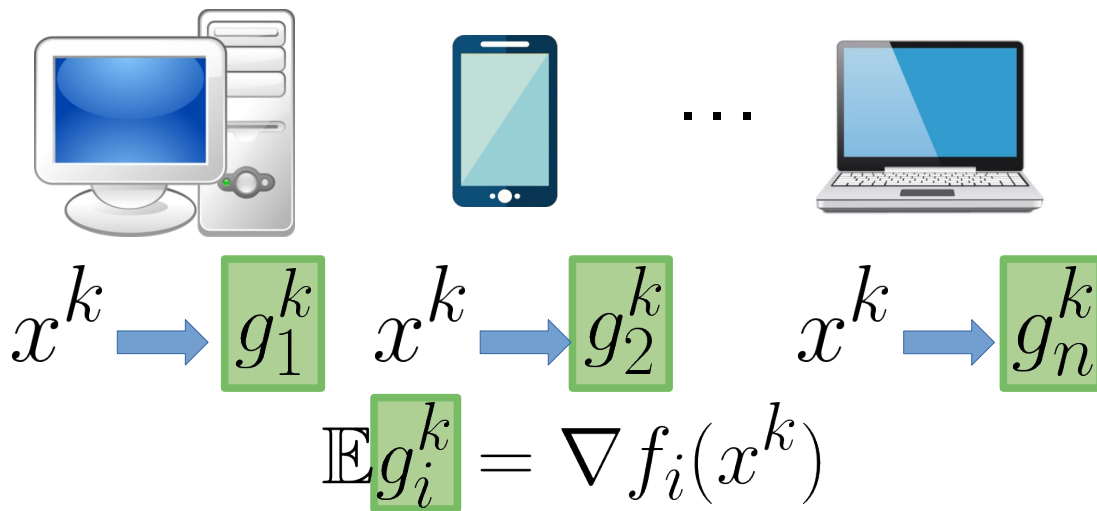
1 Server broadcasts the parameters

2 Devices compute **stochastic gradients** in parallel



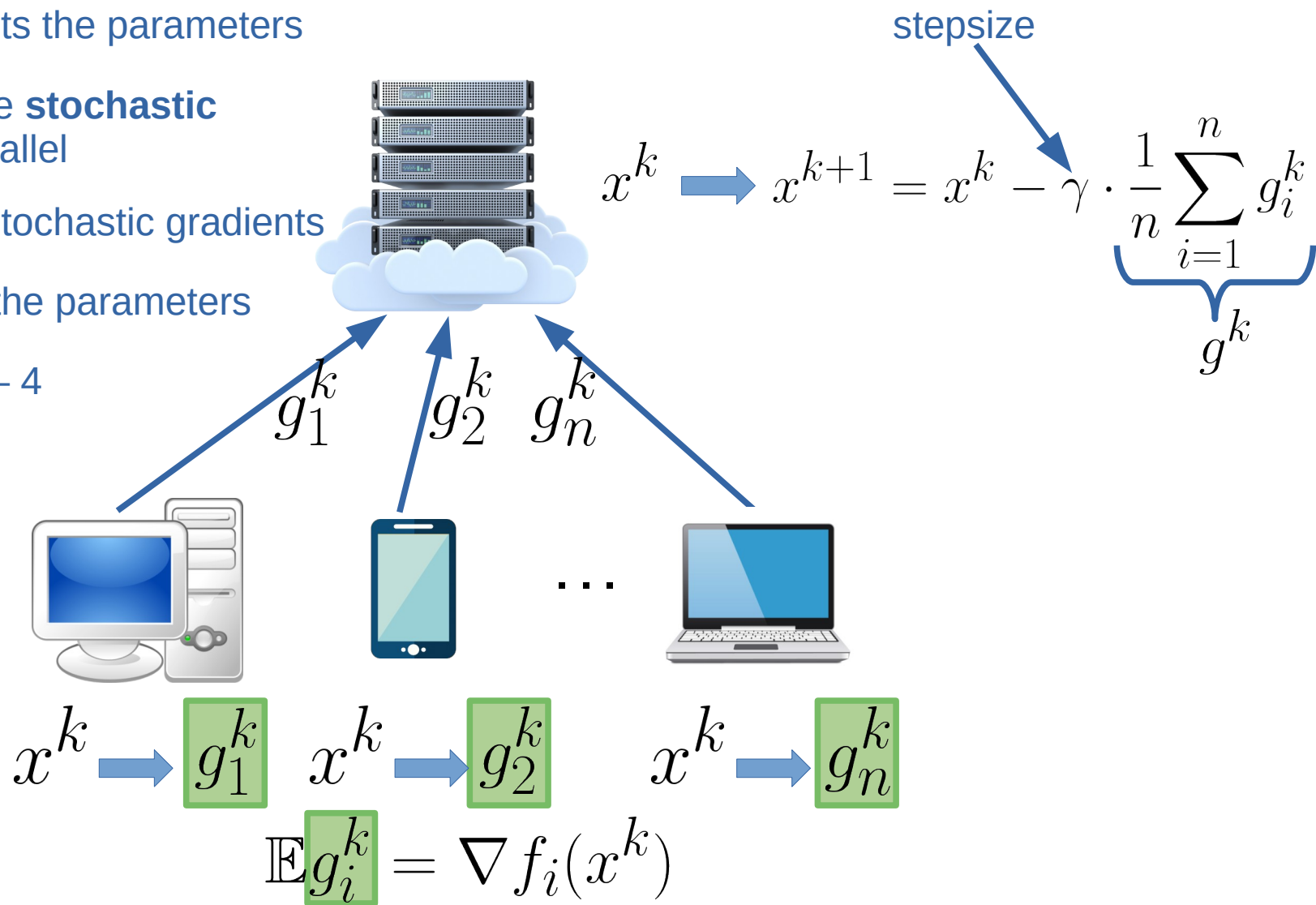
1 Server broadcasts the parameters

2 Devices compute **stochastic gradients** in parallel



13

- 1 Server broadcasts the parameters
- 2 Devices compute **stochastic gradients** in parallel
- 3 Server gathers stochastic gradients
- 4 Server updates the parameters
- 5 Repeat steps 1 – 4



14

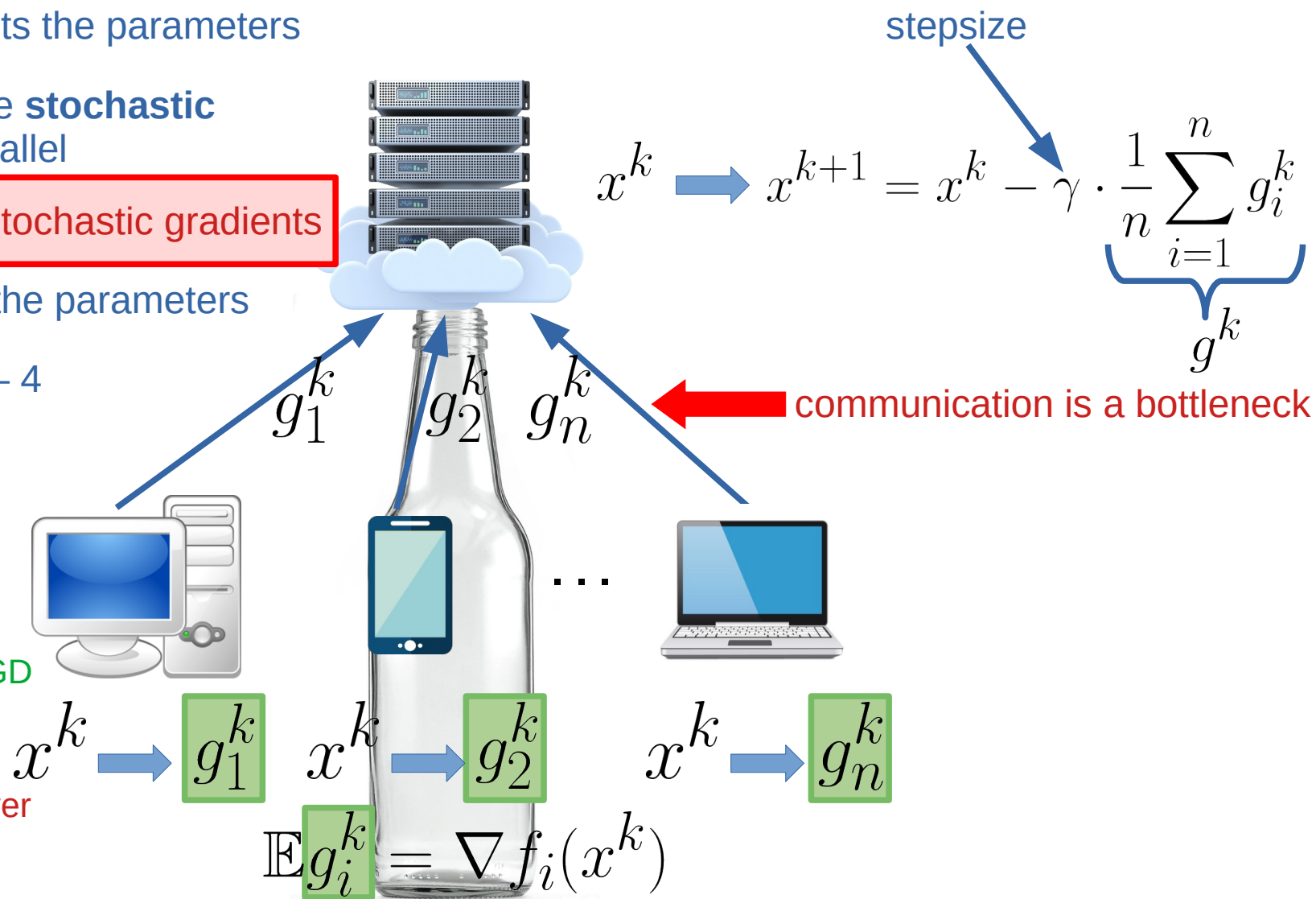
- 1 Server broadcasts the parameters
- 2 Devices compute **stochastic gradients** in parallel
- 3 **Server gathers stochastic gradients**
- 4 Server updates the parameters
- 5 Repeat steps 1 – 4

Good news:

- ✓ Very simple algorithm
- ✓ Can be much faster than non-parallel SGD

Issues:

- ✗ Overload of the server





3. Communication Bottleneck

How to Handle Communication Bottleneck?


● Change the topology of the network → Decentralized optimization

● Do more work on each worker and communicate less → Local-SGD/Federated Averaging

How to Handle Communication Bottleneck?

- Change the topology of the network  Decentralized optimization
- Do more work on each worker and communicate less  Local-SGD/Federated Averaging
- Send less information to reduce the communication cost

Example:

Send $g_i^k = \begin{pmatrix} 1 \\ -15 \\ 0.2 \\ -7 \\ 10 \end{pmatrix}$  Send $\mathcal{C}(g_i^k) = \begin{pmatrix} 0 \\ -15 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

How to Handle Communication Bottleneck?

- Change the topology of the network → Decentralized optimization
 - Do more work on each worker and communicate less → Local-SGD/Federated Averaging
 - Send less information to reduce the communication cost
- We focus on this approach

Example:

Send $g_i^k = \begin{pmatrix} 1 \\ -15 \\ 0.2 \\ -7 \\ 10 \end{pmatrix}$

Compression operator

Send $\mathcal{C}(g_i^k) = \begin{pmatrix} 0 \\ -15 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

How to Handle Communication Bottleneck?

- Change the topology of the network → Decentralized optimization
 - Do more work on each worker and communicate less → Local-SGD/Federated Averaging
 - Send less information to reduce the communication cost
- We focus on this approach

Example:

Send $g_i^k = \begin{pmatrix} 1 \\ -15 \\ 0.2 \\ -7 \\ 10 \end{pmatrix}$

Compression operator

Send $\mathcal{C}(g_i^k) = \begin{pmatrix} 0 \\ -15 \\ 0 \\ 0 \\ 0 \end{pmatrix}$

What are the options for choosing this?

Compression Operators



Unbiased compressors
(quantizations)

$$x \rightarrow \mathcal{Q}(x) \quad \mathbb{E}[\mathcal{Q}(x)] = x$$

Biased compressors

$$x \rightarrow \mathcal{C}(x)$$

Compression Operators



```
graph TD; A[Compression Operators] --> B[Unbiased compressors<br/>(quantizations)]; A --> C[Biased compressors];
```

Unbiased compressors
(quantizations)

$$x \rightarrow \mathcal{Q}(x) \quad \mathbb{E}[\mathcal{Q}(x)] = x$$

$$\mathbb{E} \|\mathcal{Q}(x) - x\|^2 \leq \omega \|x\|^2$$

Biased compressors

$$x \rightarrow \mathcal{C}(x)$$

$$\mathbb{E} \|\mathcal{C}(x) - x\|^2 \leq (1 - \delta) \|x\|^2$$

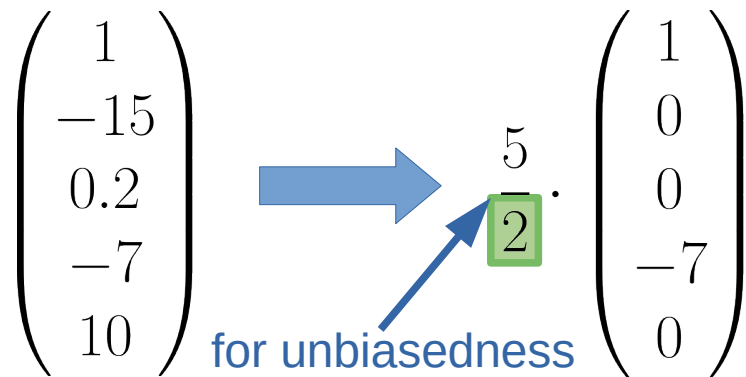
Compression Operators

Unbiased compressors
(quantizations)

$$x \rightarrow \mathcal{Q}(x) \quad \mathbb{E}[\mathcal{Q}(x)] = x$$

$$\mathbb{E}\|\mathcal{Q}(x) - x\|^2 \leq \omega \|x\|^2$$

Example: RandK (for $K = 2$)



$$\begin{pmatrix} 1 \\ -15 \\ 0.2 \\ -7 \\ 10 \end{pmatrix} \rightarrow \begin{pmatrix} 1 \\ 0 \\ 0 \\ -7 \\ 0 \end{pmatrix}$$

for unbiasedness

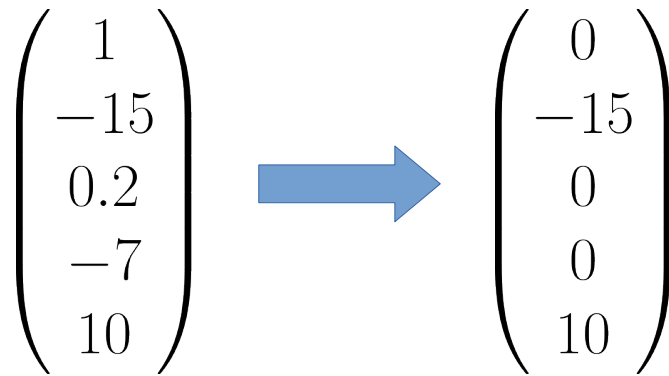
Pick $K = 2$ components uniformly at random

Biased compressors

$$x \rightarrow \mathcal{C}(x)$$

$$\mathbb{E}\|\mathcal{C}(x) - x\|^2 \leq (1 - \delta) \|x\|^2$$

Example: TopK (for $K = 2$)



$$\begin{pmatrix} 1 \\ -15 \\ 0.2 \\ -7 \\ 10 \end{pmatrix} \rightarrow \begin{pmatrix} 0 \\ -15 \\ 0 \\ 0 \\ 10 \end{pmatrix}$$

Pick $K = 2$ components with largest absolute value

Methods with Unbiased Compressors

QSGD



Alistarh, Dan, Demjan Grubic, Jerry Li, Ryota Tomioka, and Milan Vojnovic. **"QSGD: Communication-efficient SGD via gradient quantization and encoding."** In *Advances in Neural Information Processing Systems*, pp. 1709-1720. 2017.

TernGrad



Wen, Wei, Cong Xu, Feng Yan, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. **"Terngrad: Ternary gradients to reduce communication in distributed deep learning."** In *Advances in neural information processing systems*, pp. 1509-1519. 2017.

DQGD



Khirirat, Sarit, Hamid Reza Feyzmahdavian, and Mikael Johansson. **"Distributed learning with compressed gradients."** arXiv preprint arXiv:1806.06573 (2018).

DIANA



Mishchenko, Konstantin, Eduard Gorbunov, Martin Takáč, and Peter Richtárik. **"Distributed learning with compressed gradient differences."** arXiv preprint arXiv:1901.09269 (2019).



Horváth, Samuel, Dmitry Kovalev, Konstantin Mishchenko, Sebastian Stich, and Peter Richtárik. **"Stochastic distributed learning with gradient quantization and variance reduction."** arXiv preprint arXiv:1904.05115 (2019).

Sublinear convergence rates even in the case when workers quantize full gradients

Converges **linearly** when workers quantize full gradients

Parallel SGD with Biased Compressor Can Diverge at Exponential Rate



Beznosikov, Aleksandr, Samuel Horváth, Peter Richtárik, and Mher Safaryan. "On Biased Compression for Distributed Learning." arXiv preprint arXiv:2002.12410 (2020).

$$n = d = 3$$

$$f_1(x) = \langle a, x \rangle^2 + \frac{1}{4} \|x\|^2 \quad f_2(x) = \langle b, x \rangle^2 + \frac{1}{4} \|x\|^2 \quad f_3(x) = \langle c, x \rangle^2 + \frac{1}{4} \|x\|^2$$

$$a = (-3, 2, 2)^\top \quad b = (2, -3, 2)^\top \quad c = (2, 2, -3)^\top$$

$$x^0 = (t, t, t)^\top$$

In this case Parallel SGD with Top1 compression operator satisfies

$$x^k = \left(1 + \frac{11\gamma}{6}\right)^k x^0$$

Parallel SGD with Biased Compressor Can Diverge at Exponential Rate



Beznosikov, Aleksandr, Samuel Horváth, Peter Richtárik, and Mher Safaryan. "On Biased Compression for Distributed Learning." arXiv preprint arXiv:2002.12410 (2020).

$$n = d = 3$$

$$f_1(x) = \langle a, x \rangle^2 + \frac{1}{4} \|x\|^2 \quad f_2(x) = \langle b, x \rangle^2 + \frac{1}{4} \|x\|^2 \quad f_3(x) = \langle c, x \rangle^2 + \frac{1}{4} \|x\|^2$$

$$a = (-3, 2, 2)^\top \quad b = (2, -3, 2)^\top \quad c = (2, 2, -3)^\top$$

$$x^0 = (t, t, t)^\top$$

In this case Parallel SGD with Top1 compression operator satisfies

$$x^k = \left(1 + \frac{11\gamma}{6}\right)^k x^0$$

One can fix this using one special trick called **error-compensation**

4. Error-Compensated SGD

Papers on EC-SGD



Seide, Frank, Hao Fu, Jasha Droppo, Gang Li, and Dong Yu. **"1-bit stochastic gradient descent and its application to data-parallel distributed training of speech dnns."** *In Fifteenth Annual Conference of the International Speech Communication Association*. 2014.



Stich, Sebastian U., Jean-Baptiste Cordonnier, and Martin Jaggi. **"Sparsified SGD with memory."** *In Advances in Neural Information Processing Systems*, pp. 4447-4458. 2018.



Karimireddy, Sai Praneeth, Quentin Rebjock, Sebastian Stich, and Martin Jaggi. **"Error Feedback Fixes SignSGD and other Gradient Compression Schemes."** *In International Conference on Machine Learning*, pp. 3252-3261. 2019.



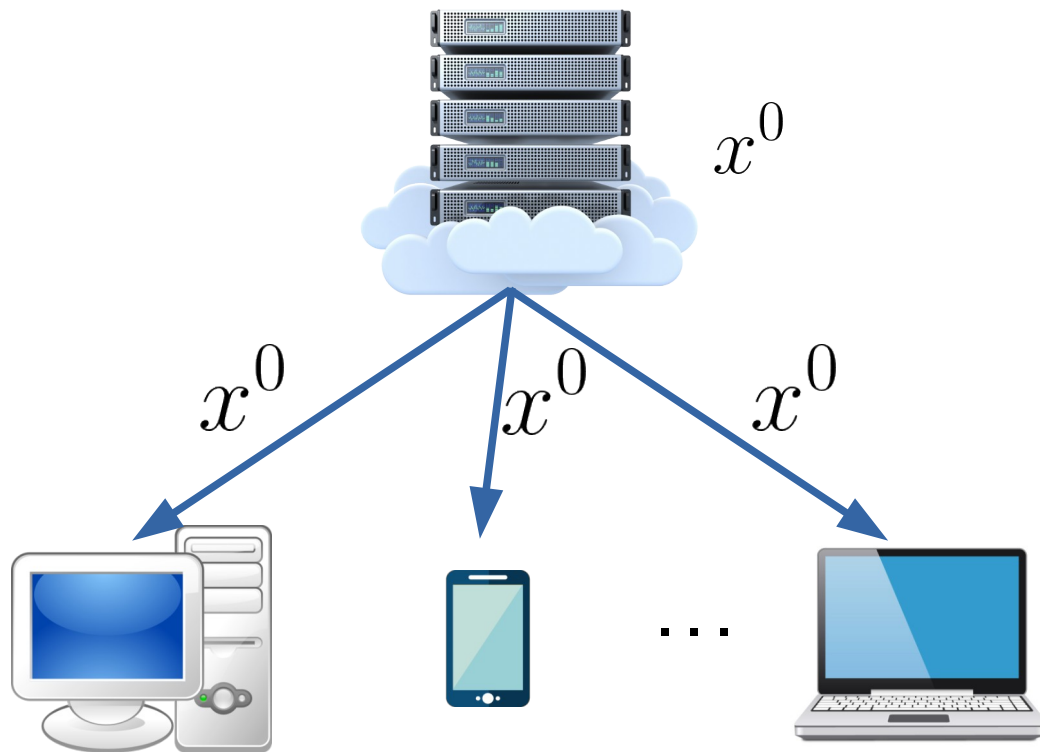
Stich, Sebastian U., and Sai Praneeth Karimireddy. **"The error-feedback framework: Better rates for SGD with delayed gradients and compressed communication."** arXiv preprint arXiv:1909.05350 (2019).



Beznosikov, Aleksandr, Samuel Horváth, Peter Richtárik, and Mher Safaryan. **"On Biased Compression for Distributed Learning."** arXiv preprint arXiv:2002.12410 (2020).

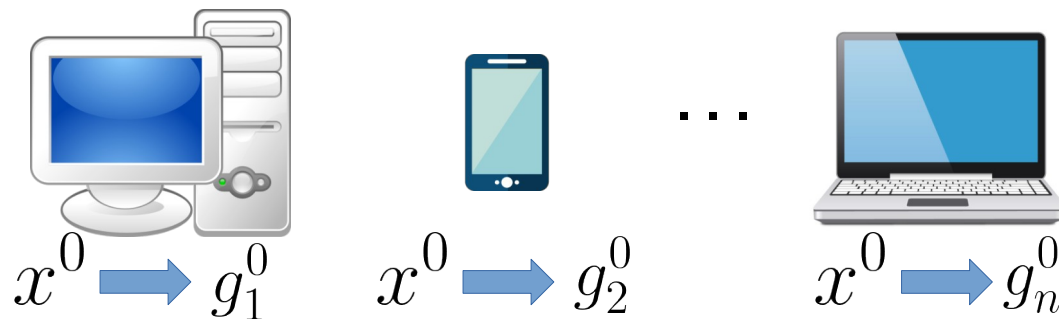
- 1 Server broadcasts the parameters

Step 1



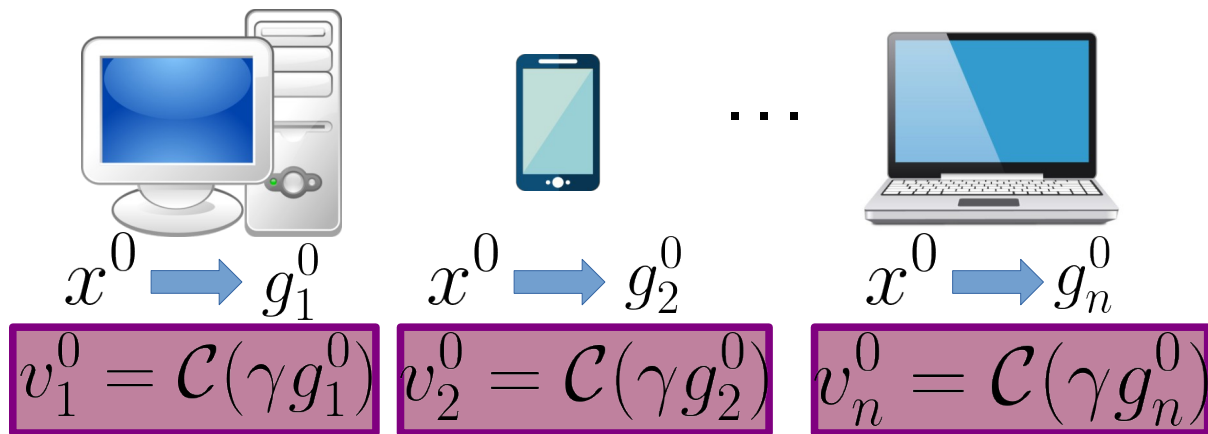
Step 1

- 1 Server broadcasts the parameters
- 2 Devices compute **stochastic gradients** in parallel



Step 1

- 1 Server broadcasts the parameters
- 2 Devices compute **stochastic gradients** in parallel
- 3 Compression



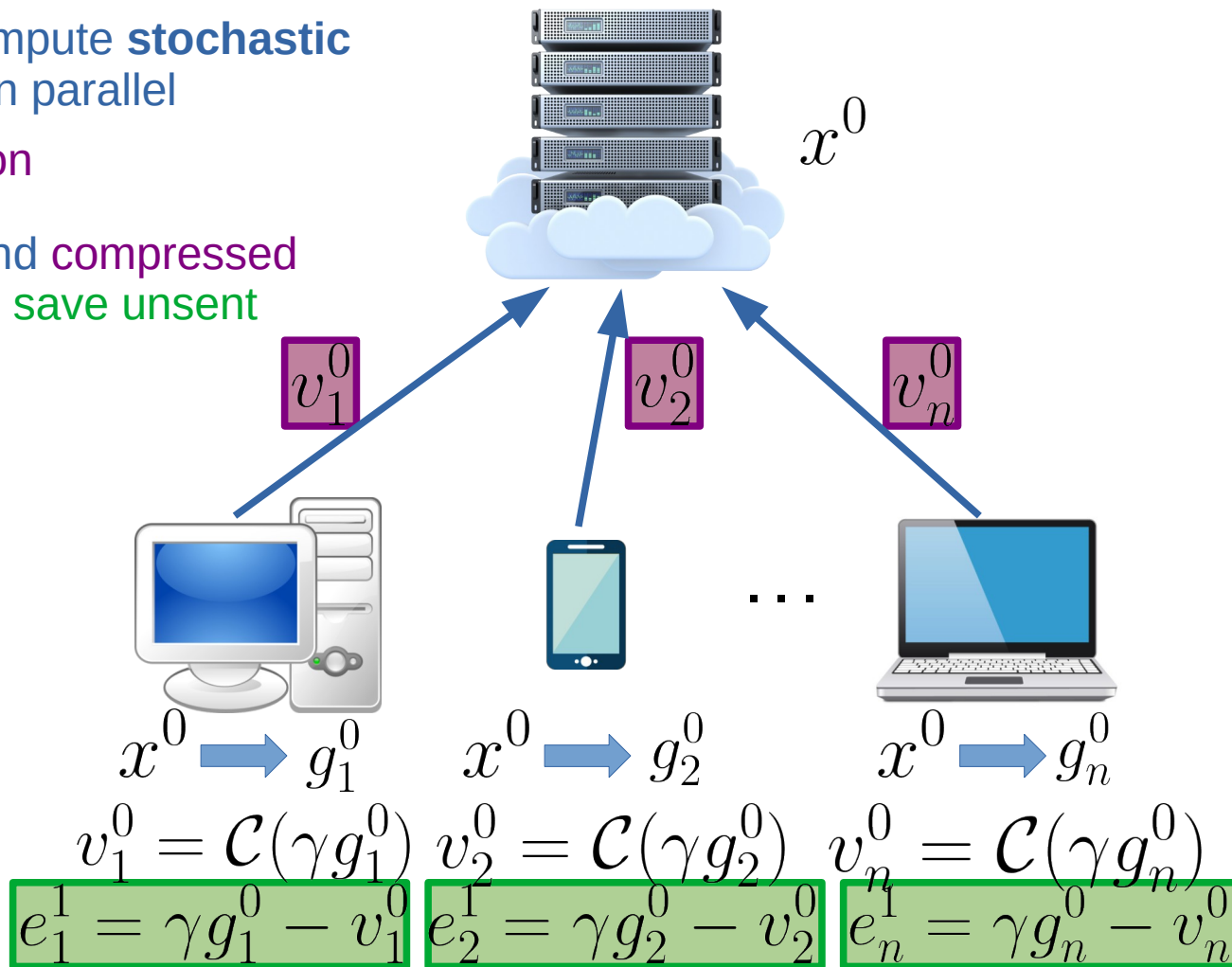
Step 1

1 Server broadcasts the parameters

2 Devices compute **stochastic gradients** in parallel

3 Compression

4 Devices send **compressed vectors** and **save unsent information**



Step 1

1 Server broadcasts the parameters

2 Devices compute **stochastic gradients** in parallel

3 **Compression**

4 Devices send **compressed vectors** and **save unsent information**

5 Server gathers the information and updates the parameters



$$x^0 \rightarrow x^1 = x^0 - \frac{1}{n} \sum_{i=1}^n v_i^0$$



$$x^0 \rightarrow g_1^0$$



$$x^0 \rightarrow g_2^0$$

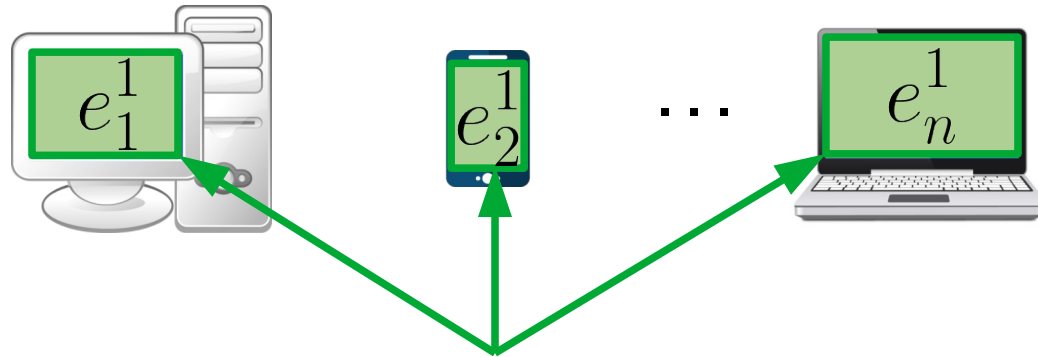
...



$$x^0 \rightarrow g_n^0$$

$$\begin{aligned} v_1^0 &= \mathcal{C}(\gamma g_1^0) & v_2^0 &= \mathcal{C}(\gamma g_2^0) & v_n^0 &= \mathcal{C}(\gamma g_n^0) \\ e_1^1 &= \gamma g_1^0 - v_1^0 & e_2^1 &= \gamma g_2^0 - v_2^0 & e_n^1 &= \gamma g_n^0 - v_n^0 \end{aligned}$$

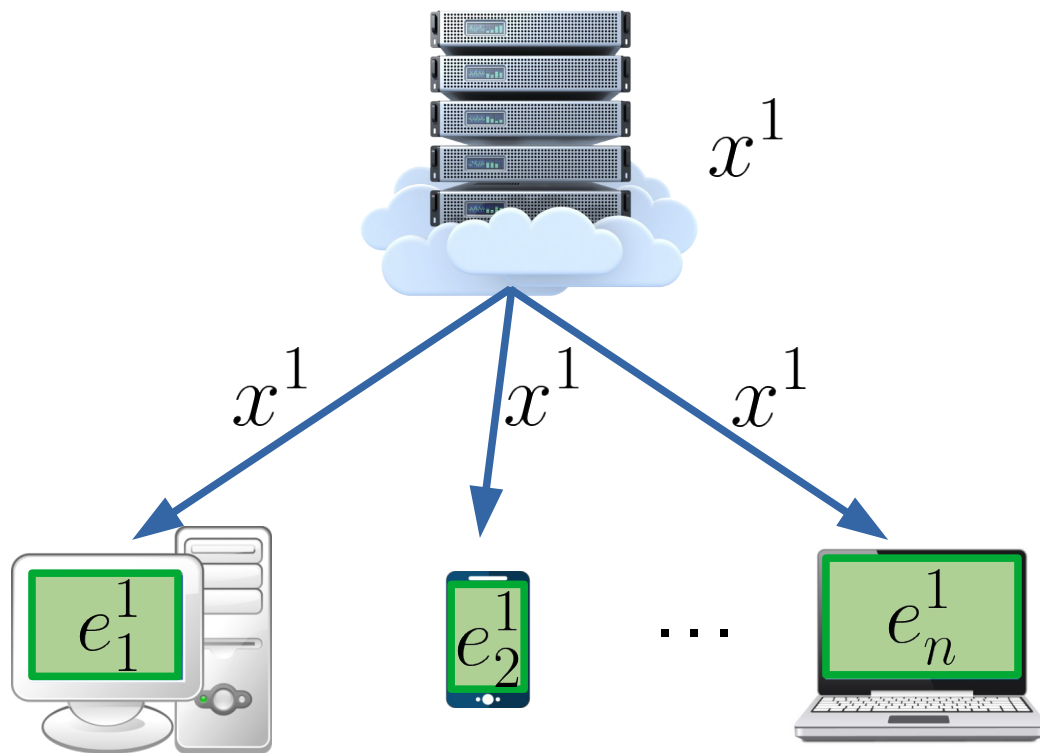
Step 1



devices keep these vectors for the next iterations
to *partially* send them later

- 1 Server broadcasts new parameters

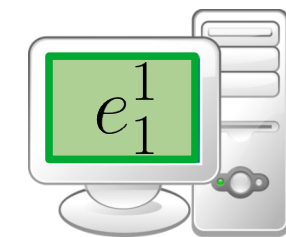
Step 2



Step 2

1 Server broadcasts new parameters

2 Devices compute **stochastic gradients** in parallel

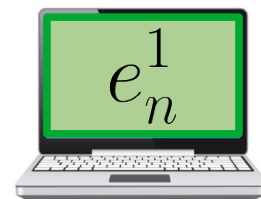


$$x^1 \rightarrow g_1^1$$



$$x^1 \rightarrow g_2^1$$

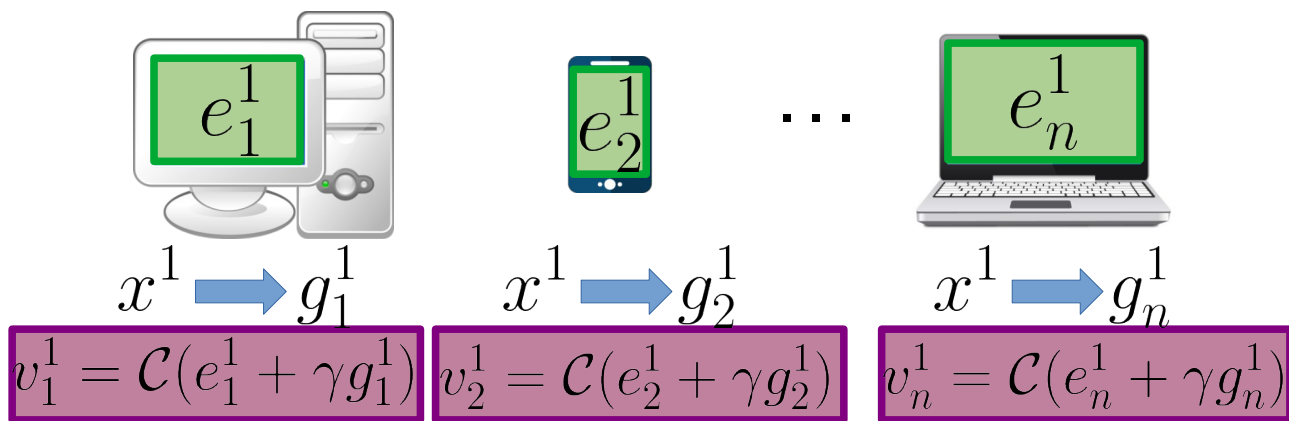
...



$$x^1 \rightarrow g_n^1$$

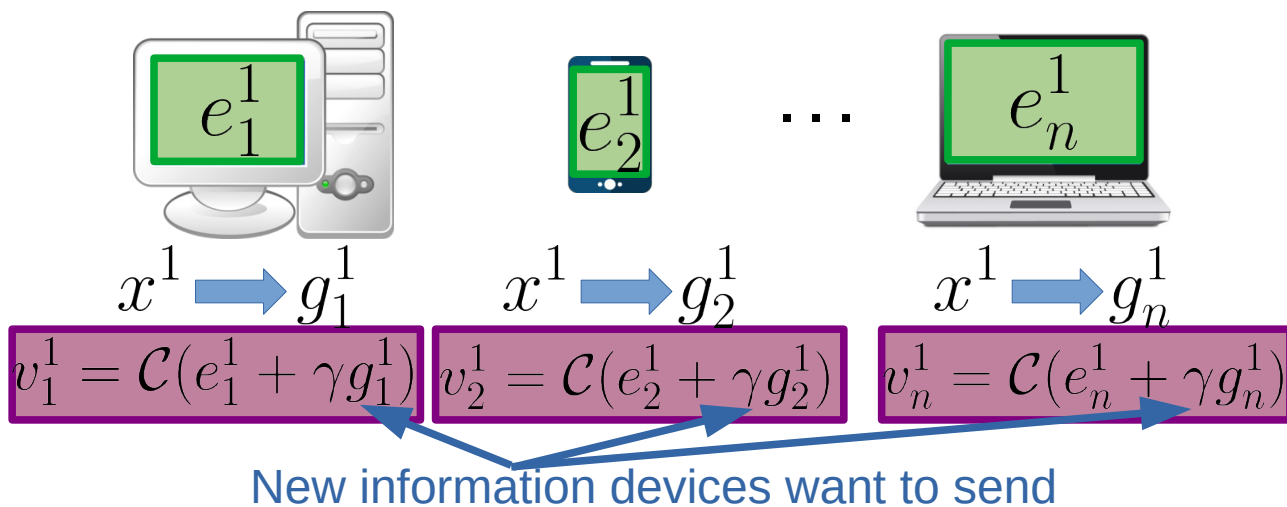
Step 2

- 1 Server broadcasts new parameters
- 2 Devices compute **stochastic gradients** in parallel
- 3 Compression



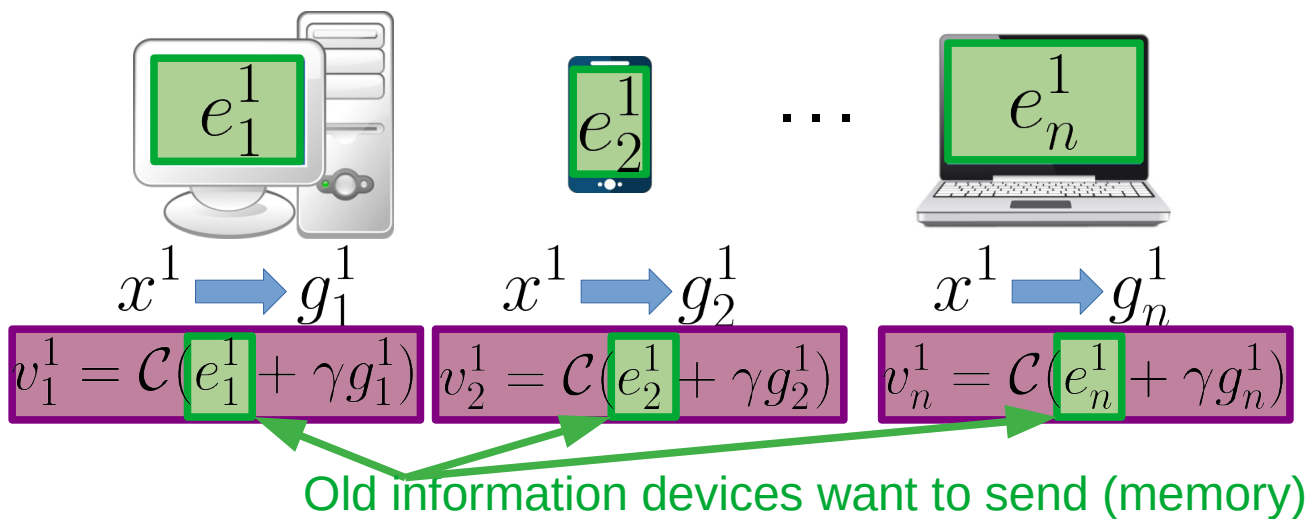
Step 2

- 1 Server broadcasts new parameters
- 2 Devices compute **stochastic gradients** in parallel
- 3 Compression



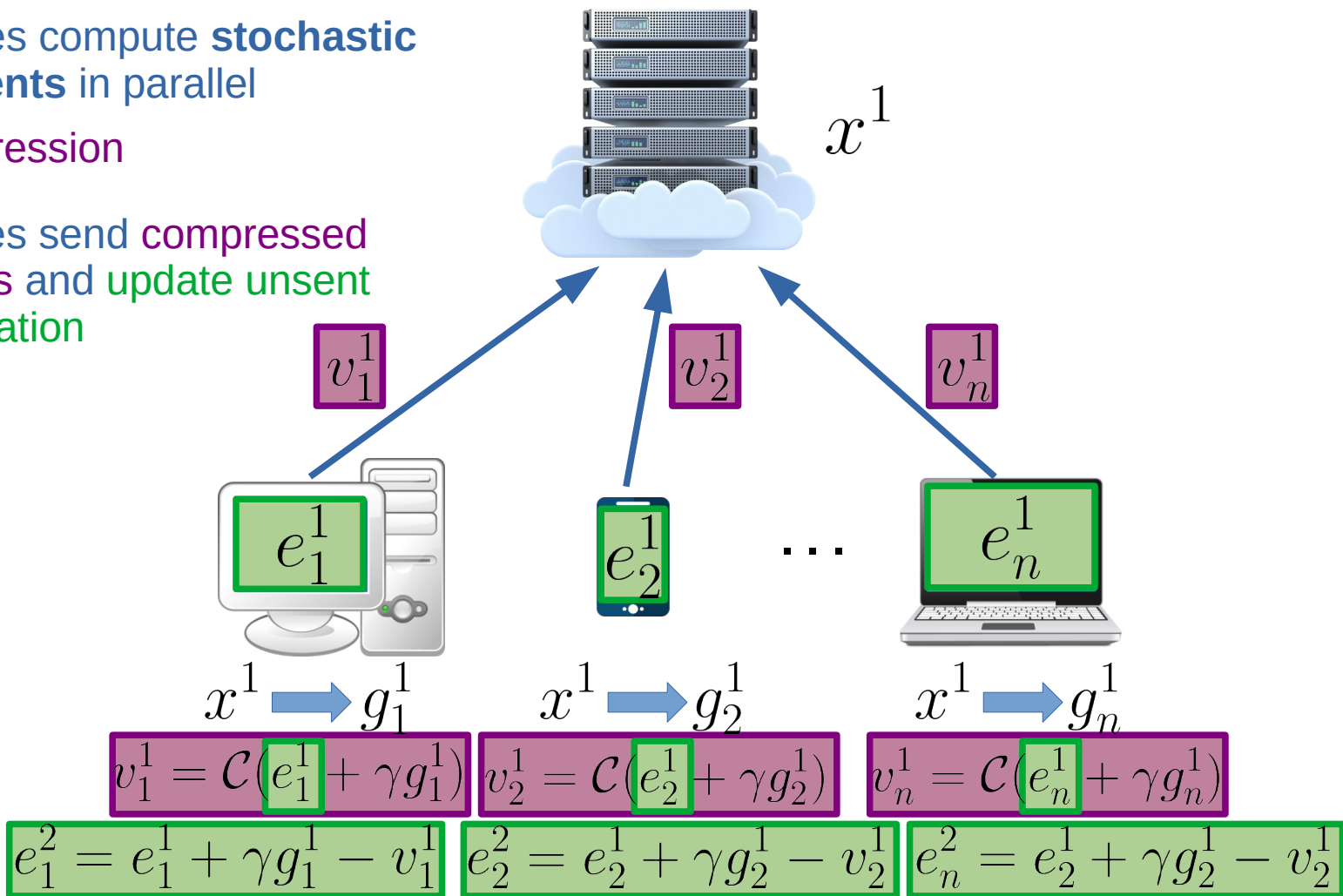
Step 2

- 1 Server broadcasts new parameters
- 2 Devices compute **stochastic gradients** in parallel
- 3 Compression



Step 2

- 1 Server broadcasts new parameters
- 2 Devices compute **stochastic gradients** in parallel
- 3 **Compression**
- 4 Devices send **compressed vectors** and **update unsent information**



Step 2

1 Server broadcasts new parameters

2 Devices compute **stochastic gradients** in parallel

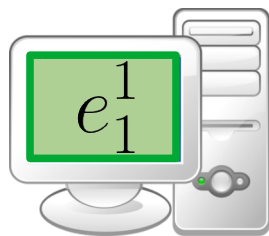
3 **Compression**

4 Devices send **compressed vectors** and **update unsent information**

5 Server gathers the information and updates the parameters



$$x^1 \rightarrow x^2 = x^1 - \frac{1}{n} \sum_{i=1}^n v_i^1$$

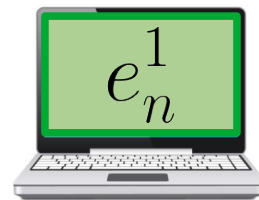


$$x^1 \rightarrow g_1^1$$



$$x^1 \rightarrow g_2^1$$

...



$$x^1 \rightarrow g_n^1$$

$$v_1^1 = \mathcal{C}(e_1^1 + \gamma g_1^1) \quad v_2^1 = \mathcal{C}(e_2^1 + \gamma g_2^1) \quad v_n^1 = \mathcal{C}(e_n^1 + \gamma g_n^1)$$

$$e_1^2 = e_1^1 + \gamma g_1^1 - v_1^1 \quad e_2^2 = e_2^1 + \gamma g_2^1 - v_2^1 \quad e_n^2 = e_n^1 + \gamma g_n^1 - v_n^1$$

Step 2

1 Server broadcasts new parameters

2 Devices compute **stochastic gradients** in parallel

3 **Compression**

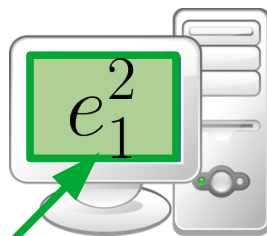
4 Devices send **compressed vectors** and **update unsent information**

5 Server gathers the information and updates the parameters

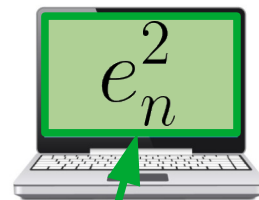
6 **Devices update their memory**



$$x^1 \rightarrow x^2 = x^1 - \frac{1}{n} \sum_{i=1}^n v_i^1$$



...



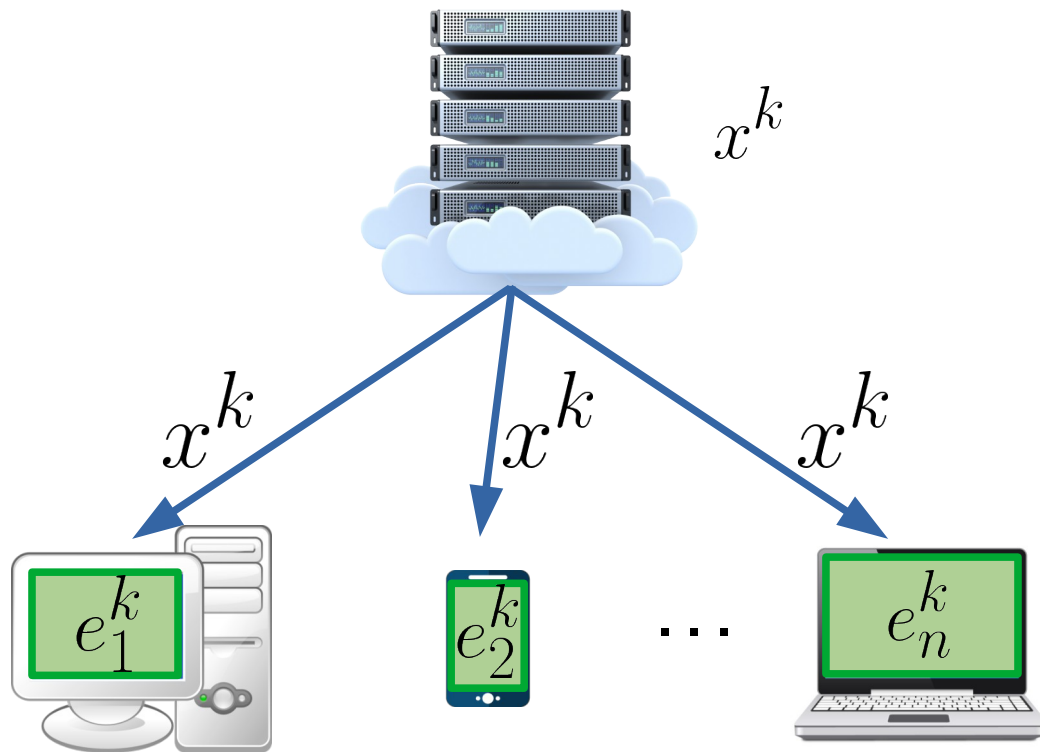
$$e_1^2 = e_1^1 + \gamma g_1^1 - v_1^1$$

$$e_2^2 = e_2^1 + \gamma g_2^1 - v_2^1$$

$$e_n^2 = e_n^1 + \gamma g_n^1 - v_n^1$$

1 Server broadcasts new parameters

Step $k+1$



Step $k+1$

1 Server broadcasts new parameters

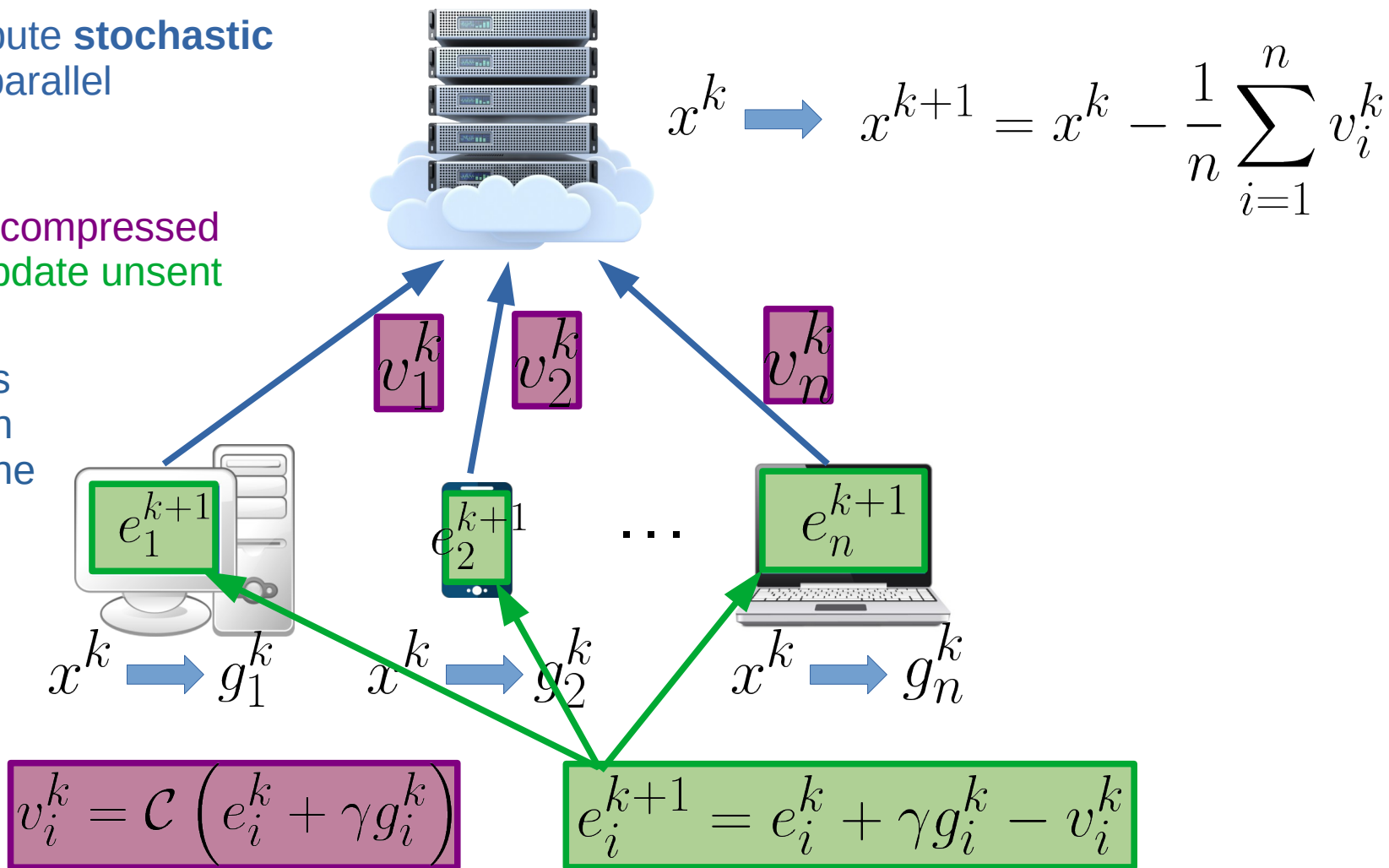
2 Workers compute **stochastic gradients** in parallel

3 **Compression**

4 Devices send **compressed vectors** and **update unsent information**

5 Server gathers the information and updates the parameters

6 Repeat steps 1 – 5



Error-Compensated SGD



Converges even with biased compression operators

EC-SGD finds such \hat{x} that $\mathbb{E}[f(\hat{x})] - f(x^*) \leq \varepsilon$ after

Hides logarithmical
factors

$$\rightarrow \tilde{\mathcal{O}} \left(\frac{L}{\delta \mu} + \frac{\sigma^2}{n \mu \varepsilon} + \frac{\sqrt{L(\sigma^2 + \zeta_*^2 / \delta)}}{\mu \sqrt{\delta \varepsilon}} \right) \text{ iterations}$$

$$\mathbb{E} \|\mathcal{C}(x) - x\|^2 \leq (1 - \delta) \|x\|^2$$

$$\mathbb{E} \left[\|g_i^k - \nabla f_i(x^k)\|^2 \mid x^k \right] \leq \sigma^2$$

$$\zeta_*^2 = \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*)\|^2$$

Error-Compensated SGD



Converges even with biased compression operators



Fails to converge with **linear rate** even when workers compute full gradients

EC-SGD finds such \hat{x} that $\mathbb{E}[f(\hat{x})] - f(x^*) \leq \varepsilon$ after

Hides logarithmical
factors

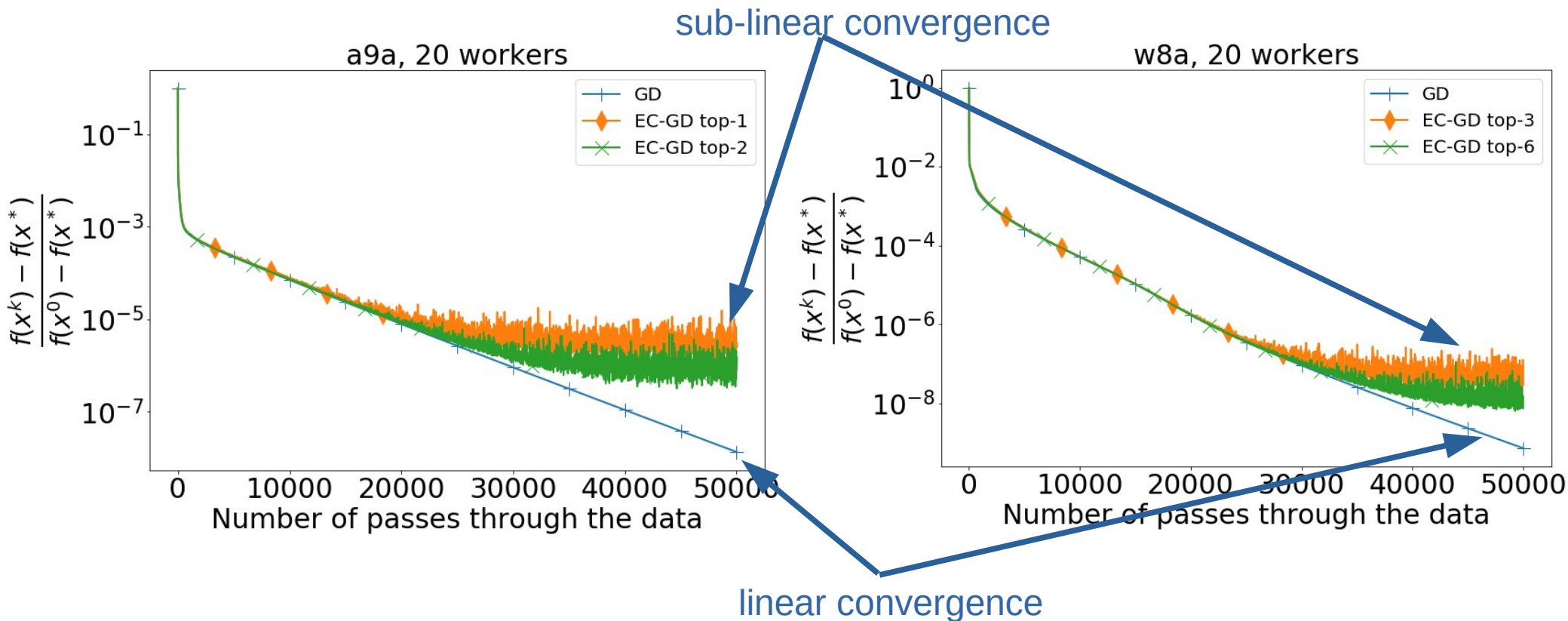
$$\rightarrow \tilde{\mathcal{O}} \left(\frac{L}{\delta\mu} + \frac{\sigma^2}{n\mu\varepsilon} + \frac{\sqrt{L(\sigma^2 + \zeta_*^2/\delta)}}{\mu\sqrt{\delta\varepsilon}} \right) \text{ iterations}$$

$$\mathbb{E}\|\mathcal{C}(x) - x\|^2 \leq (1 - \delta)\|x\|^2 \quad \mathbb{E} \left[\|g_i^k - \nabla f_i(x^k)\|^2 \mid x^k \right] \leq \sigma^2$$

$$\zeta_*^2 = \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*)\|^2$$

EC-GD and Logistic Regression

$$\min_{x \in \mathbb{R}^d} \left\{ f(x) = \frac{1}{N} \sum_{i=1}^N \log(1 + \exp(-y_i \cdot (Ax)_i)) + \frac{\mu}{2} \|x\|^2 \right\}$$



Error-Compensated SGD

- ✓ Converges even with biased compression operators
- ✗ Fails to converge with **linear rate** even when workers compute full gradients

Questions:

- 1 *Is it possible to design **linearly converging** SGD with error compensation when workers compute full gradients, i.e., linearly converging **EC-GD**?*
- 2 *Is it possible to design **linearly converging** SGD with error compensation when **the local loss functions have a finite-sum form**?*

The answer is Yes for both questions

5.1. New method: EC-GDstar

Error-Compensated GD

EC-GD finds such \hat{x} that $\mathbb{E}[f(\hat{x})] - f(x^*) \leq \varepsilon$ after

Hides logarithmical factors $\rightarrow \tilde{\mathcal{O}} \left(\frac{L}{\delta\mu} + \frac{\sqrt{L\zeta_*^2}}{\mu\delta\sqrt{\varepsilon}} \right)$ iterations

$$\zeta_*^2 = \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*)\|^2$$

Error-Compensated GD

EC-GD finds such \hat{x} that $\mathbb{E}[f(\hat{x})] - f(x^*) \leq \varepsilon$ after

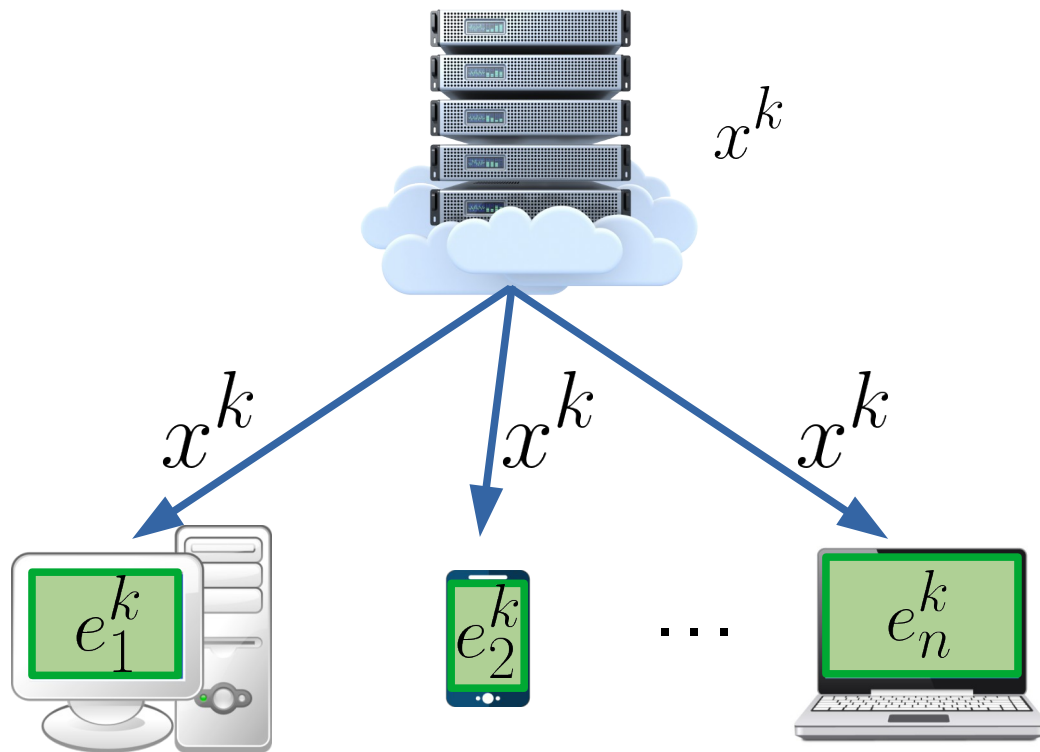
Hides logarithmical factors $\rightarrow \tilde{\mathcal{O}} \left(\frac{L}{\delta\mu} + \frac{\sqrt{L\zeta_*^2}}{\mu\delta\sqrt{\varepsilon}} \right)$ iterations

$$\zeta_*^2 = \frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x^*)\|^2$$

What if devices know these vectors from the beginning?

EC-GDstar

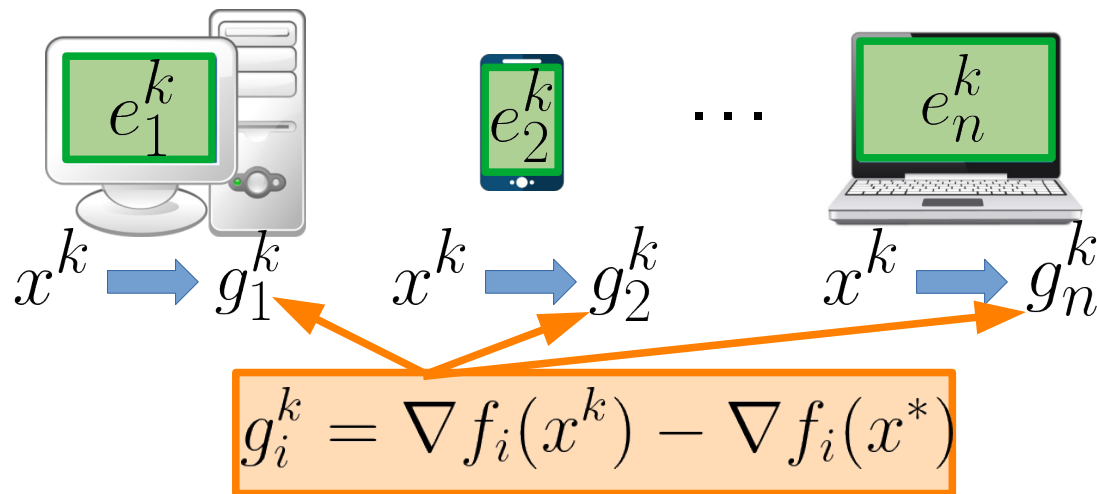
- 1 Server broadcasts new parameters



EC-GDstar

1 Server broadcasts new parameters

2 Workers compute **shifted gradients** in parallel



EC-GDstar

1 Server broadcasts new parameters

2 Workers compute **shifted gradients** in parallel

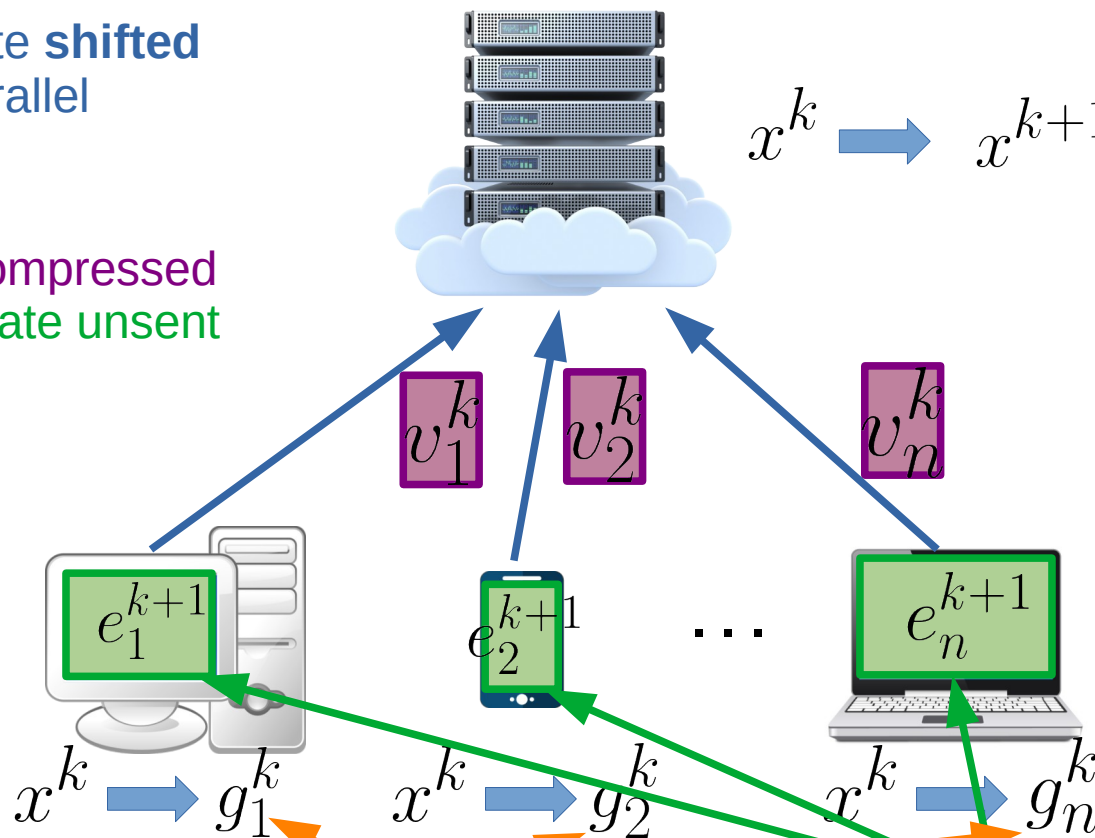
3 **Compression**

4 Devices send **compressed vectors** and **update unsent information**

5 Server gathers the information and updates the parameters

6 Repeat steps 1 – 5

$$x^k \rightarrow x^{k+1} = x^k - \frac{1}{n} \sum_{i=1}^n v_i^k$$



$$v_i^k = \mathcal{C} \left(e_i^k + \gamma g_i^k \right)$$

$$g_i^k = \nabla f_i(x^k) - \nabla f_i(x^*)$$

$$e_i^{k+1} = e_i^k + \gamma g_i^k - v_i^k$$

EC-GDstar: Rate of Convergence

EC-GDstar finds such \hat{x} that $\mathbb{E}[f(\hat{x})] - f(x^*) \leq \varepsilon$ after

$$\mathcal{O}\left(\frac{L}{\delta\mu} \ln \frac{1}{\varepsilon}\right) \text{ iterations}$$



Linear rate



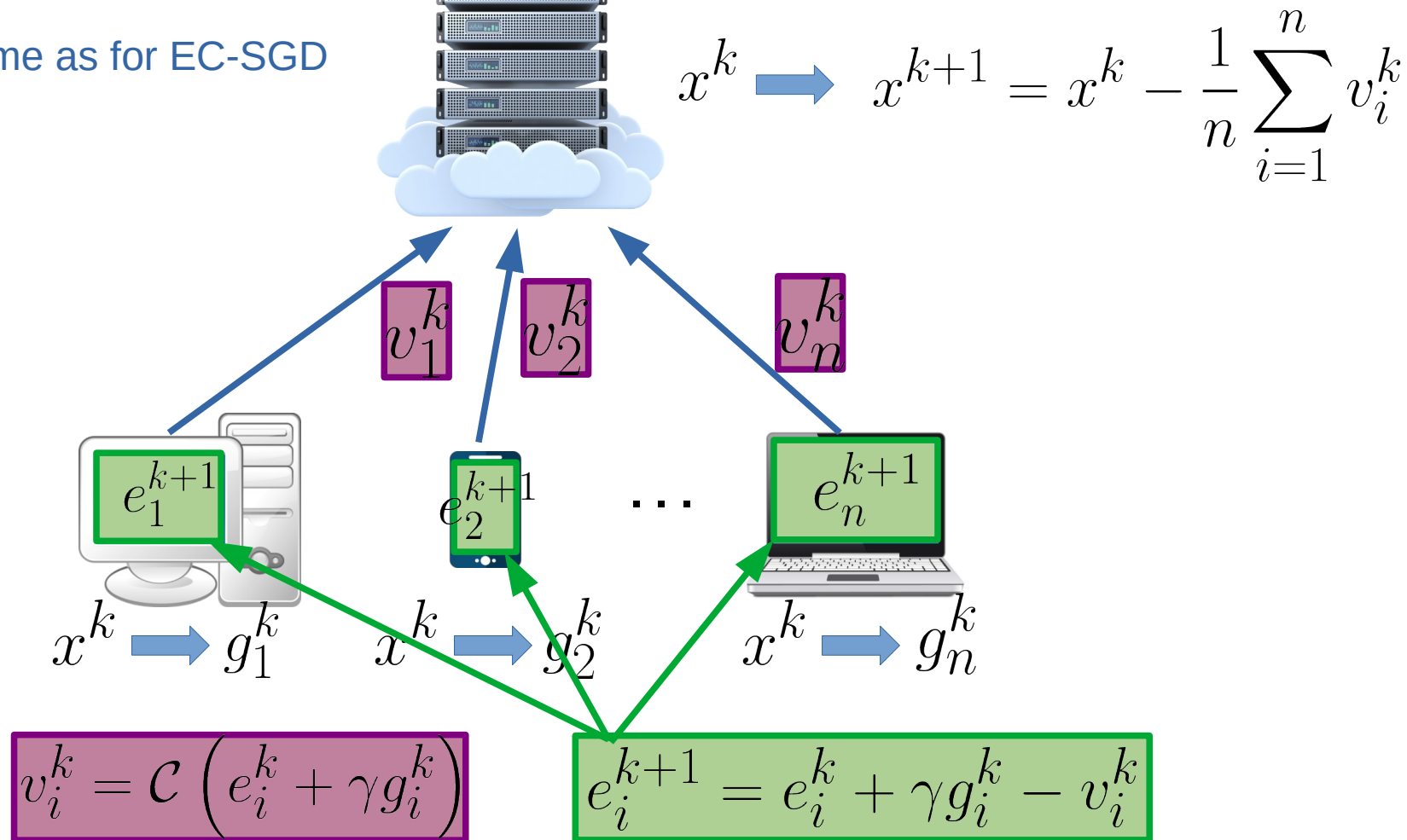
The method is impractical: it uses the gradients at the solution

Can we develop a practical analog?

5.2. New method: EC-SGD-DIANA

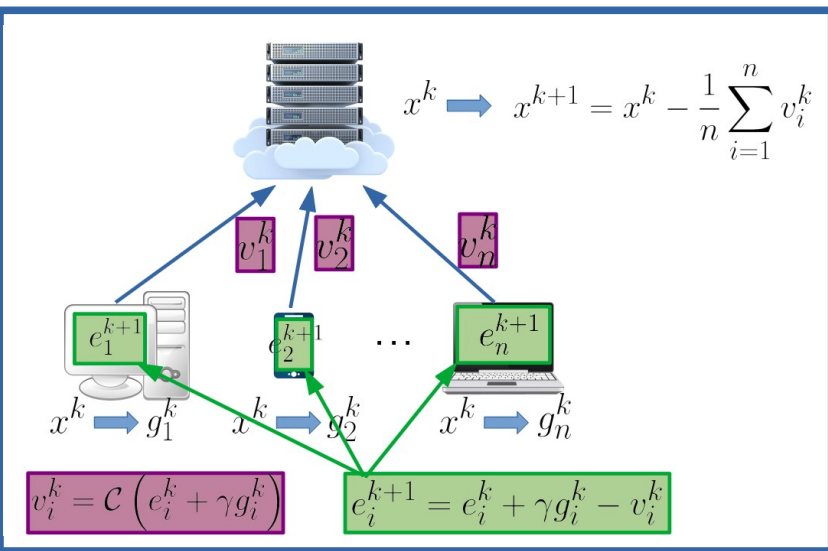
EC-SGD-DIANA

The same scheme as for EC-SGD



EC-SGD-DIANA

$$g_i^k = \hat{g}_i^k - h_i^k + h^k$$



EC-SGD-DIANA

The key insight why do we need $\{h_i^k\}_{i=1}^n$:

it reduces the variance coming from compressions via learning the gradients at the solution!

$$\mathbb{E} \left[\hat{g}_i^k \mid x^k \right] = \nabla f_i(x^k)$$

$$g_i^k = \boxed{\hat{g}_i^k} - h_i^k + h^k \leftarrow h^k = \frac{1}{n} \sum_{i=1}^n h_i^k$$

stepsize

$$h_i^{k+1} = h_i^k + \alpha \mathcal{Q} \left(\hat{g}_i^k - h_i^k \right)$$

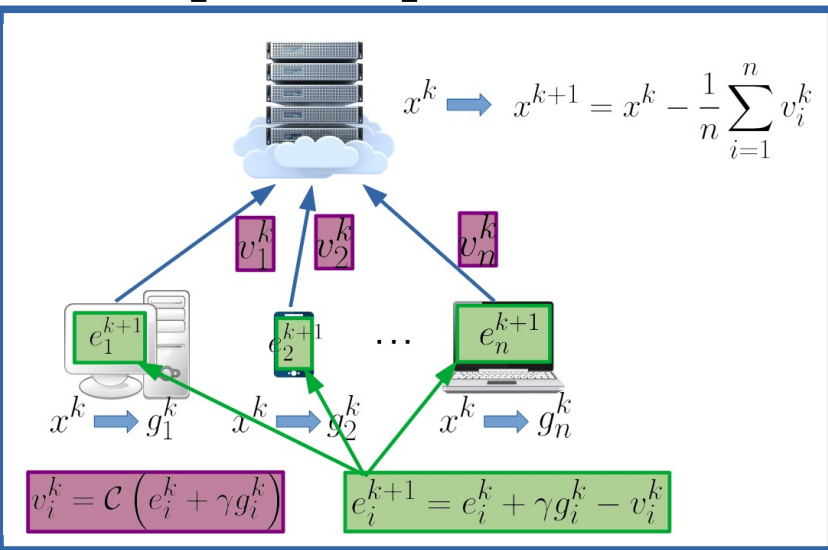
Server broadcasts this vector to the workers

Workers send these vectors to the server

Works for both cases:

● $f_i(x) = \mathbf{E}_{\xi_i \sim \mathcal{D}_i} [f_{\xi_i}(x)] \quad \boxed{\hat{g}_i^k} = \nabla f_{\xi_i}(x^k)$

● $f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x) \quad \boxed{\hat{g}_i^k} = \nabla f_{il}(x^k)$
 $l \sim [m]$ uniformly at random



EC-SGD-DIANA: Rate of Convergence

EC-SGD-DIANA finds such \hat{x} that $\mathbb{E}[f(\hat{x})] - f(x^*) \leq \varepsilon$ after

Hides logarithmical factors

Option I: $\tilde{\mathcal{O}} \left(\omega + \frac{L}{\delta\mu} + \frac{\sigma^2}{n\mu\varepsilon} + \frac{\sqrt{L\sigma^2}}{\delta\mu\sqrt{\varepsilon}} \right)$ iterations

Option II: $\tilde{\mathcal{O}} \left(\frac{1 + \omega}{\delta} + \frac{L}{\delta\mu} + \frac{\sigma^2}{n\mu\varepsilon} + \frac{\sqrt{L\sigma^2}}{\mu\sqrt{\delta\varepsilon}} \right)$

EC-SGD-DIANA: Rate of Convergence

EC-SGD-DIANA finds such \hat{x} that $\mathbb{E}[f(\hat{x})] - f(x^*) \leq \varepsilon$ after

Hides logarithmical factors

Option I: $\tilde{\mathcal{O}} \left(\omega + \frac{L}{\delta\mu} + \frac{\sigma^2}{n\mu\varepsilon} + \frac{\sqrt{L\sigma^2}}{\delta\mu\sqrt{\varepsilon}} \right)$ iterations

Option II: $\tilde{\mathcal{O}} \left(\frac{1+\omega}{\delta} + \frac{L}{\delta\mu} + \frac{\sigma^2}{n\mu\varepsilon} + \frac{\sqrt{L\sigma^2}}{\mu\sqrt{\delta\varepsilon}} \right)$

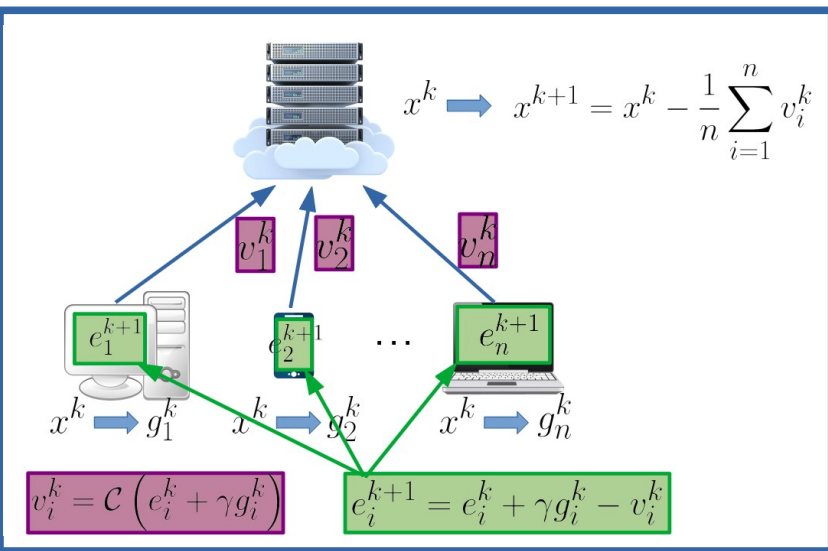
Moreover, if workers compute **full gradients**, then the rate of convergence is linear

$$\mathcal{O} \left(\left(\omega + \frac{L}{\delta\mu} \right) \log \frac{1}{\varepsilon} \right)$$

5.3. New method: EC-LSVRG-DIANA

EC-LSVRG-DIANA

$$g_i^k = \hat{g}_i^k - h_i^k + h^k$$



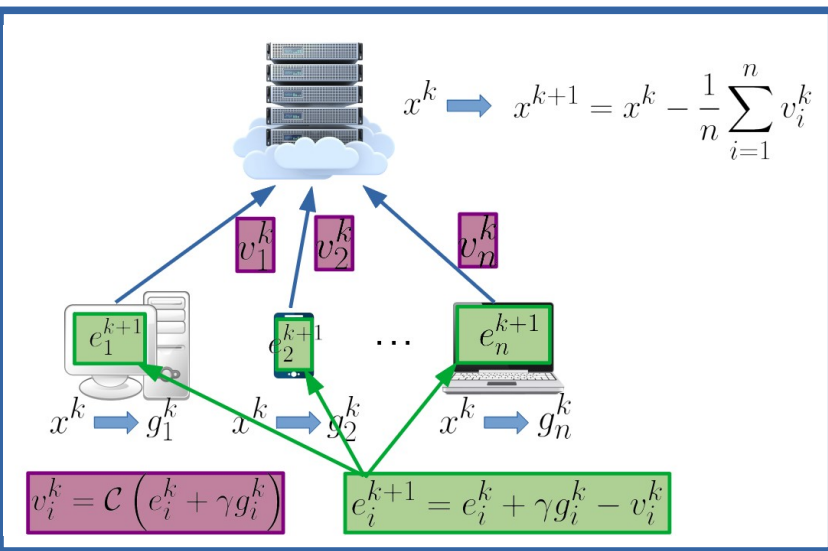
Works for the case:

$$\bullet f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x)$$

EC-LSVRG-DIANA

$$g_i^k = \hat{g}_i^k - h_i^k + h^k \quad l \sim [m] \text{ uniformly at random}$$

$$\hat{g}_i^k = \nabla f_{il} \left(x^k \right) - \nabla f_{il} \left(w_i^k \right) + \nabla f_i \left(w_i^k \right)$$



Works for the case:

$$\bullet \quad f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x)$$

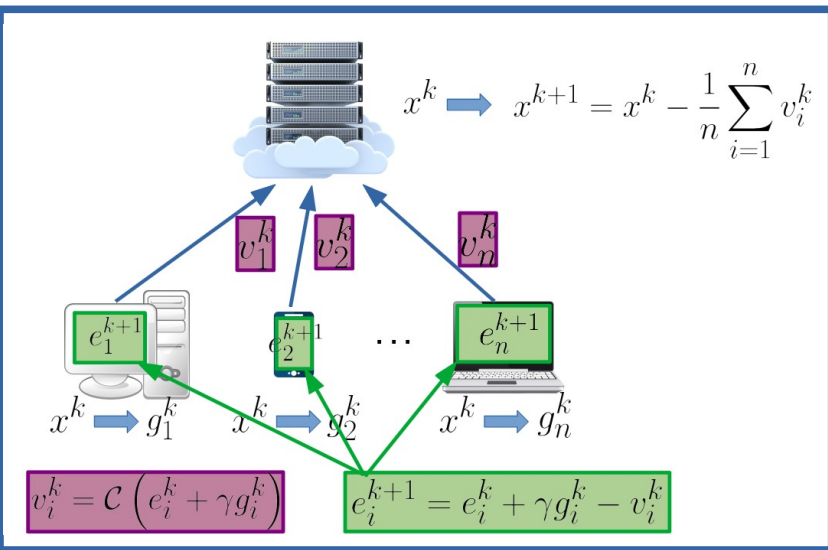
EC-LSVRG-DIANA

$$g_i^k = \boxed{\hat{g}_i^k} - h_i^k + h^k \quad \boxed{l} \sim [m] \text{ uniformly at random}$$

$$\boxed{\hat{g}_i^k} = \nabla f_{i\boxed{l}}(x^k) - \nabla f_{i\boxed{l}}(\boxed{w_i^k}) + \nabla f_i(\boxed{w_i^k})$$

Reduction of the variance introduced due to the stochasticity of the gradients

$$\boxed{w_i^{k+1}} = \begin{cases} x^k, & \text{with probability } p \\ \boxed{w_i^k}, & \text{with probability } 1 - p \end{cases}$$



Works for the case:

$$\bullet f_i(x) = \frac{1}{m} \sum_{j=1}^m f_{ij}(x)$$

EC-LSVRG-DIANA: Rate of Convergence

EC-LSVRG-DIANA finds such \hat{x} that $\mathbb{E}[f(\hat{x})] - f(x^*) \leq \varepsilon$ after

$$\mathcal{O} \left(\left(\omega + m + \frac{L}{\delta\mu} \right) \log \frac{1}{\epsilon} \right) \text{ iterations}$$

6. Unified Convergence Analysis of Methods with Error Compensation

Key Assumption

$$g^k = \frac{1}{n} \sum_{i=1}^n g_i^k, \quad \mathbb{E} [g^k \mid x^k] = \nabla f (x^k) \quad \bar{g}_i^k = \mathbb{E} [g_i^k \mid x^k]$$

$$\frac{1}{n} \sum_{i=1}^n \|\bar{g}_i^k\|^2 \leq 2A (f (x^k) - f (x^*)) + B_1 \sigma_{1,k}^2 + B_2 \sigma_{2,k}^2 + D_1$$

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|g_i^k - \bar{g}_i^k\|^2 \mid x^k] \leq 2\tilde{A} (f (x^k) - f (x^*)) + \tilde{B}_1 \sigma_{1,k}^2 + \tilde{B}_2 \sigma_{2,k}^2 + \tilde{D}_1$$

$$\mathbb{E} [\|g^k\|^2 \mid x^k] \leq 2A' (f (x^k) - f (x^*)) + B'_1 \sigma_{1,k}^2 + B'_2 \sigma_{2,k}^2 + D'_1$$

$$\mathbb{E} [\sigma_{1,k+1}^2 \mid \sigma_{1,k}^2, \sigma_{2,k}^2] \leq (1 - \rho_1) \sigma_{1,k}^2 + 2C_1 (f (x^k) - f (x^*)) + G\rho_1 \sigma_{2,k}^2 + D_2$$

$$\mathbb{E} [\sigma_{2,k+1}^2 \mid \sigma_{2,k}^2] \leq (1 - \rho_2) \sigma_{2,k}^2 + 2C_2 (f (x^k) - f (x^*))$$


Key Assumption

$$g^k = \frac{1}{n} \sum_{i=1}^n g_i^k, \quad \mathbb{E} [g^k \mid x^k] = \nabla f(x^k) \quad \bar{g}_i^k = \mathbb{E} [g_i^k \mid x^k]$$

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \|\bar{g}_i^k\|^2 &\leq 2A (f(x^k) - f(x^*)) + B_1 \sigma_{1,k}^2 + B_2 \sigma_{2,k}^2 + D_1 \\ \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|g_i^k - \bar{g}_i^k\|^2 \mid x^k] &\leq 2\tilde{A} (f(x^k) - f(x^*)) + \tilde{B}_1 \sigma_{1,k}^2 + \tilde{B}_2 \sigma_{2,k}^2 + \tilde{D}_1 \\ \mathbb{E} [\|g^k\|^2 \mid x^k] &\leq 2A' (f(x^k) - f(x^*)) + B'_1 \sigma_{1,k}^2 + B'_2 \sigma_{2,k}^2 + D'_1 \end{aligned}$$

$$\mathbb{E} [\sigma_{1,k+1}^2 \mid \sigma_{1,k}^2, \sigma_{2,k}^2] \leq (1 - \rho_1) \sigma_{1,k}^2 + 2C_1 (f(x^k) - f(x^*)) + G\rho_1 \sigma_{2,k}^2 + D_2$$

$$\mathbb{E} [\sigma_{2,k+1}^2 \mid \sigma_{2,k}^2] \leq (1 - \rho_2) \sigma_{2,k}^2 + 2C_2 (f(x^k) - f(x^*))$$

 Reflects smoothness properties of the problem and noises introduced by compressions and stochastic gradients

Key Assumption

$$g^k = \frac{1}{n} \sum_{i=1}^n g_i^k, \quad \mathbb{E} [g^k \mid x^k] = \nabla f(x^k) \quad \bar{g}_i^k = \mathbb{E} [g_i^k \mid x^k]$$

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \|\bar{g}_i^k\|^2 &\leq 2A (f(x^k) - f(x^*)) + B_1 \sigma_{1,k}^2 + B_2 \sigma_{2,k}^2 + D_1 \\ \frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|g_i^k - \bar{g}_i^k\|^2 \mid x^k] &\leq 2\tilde{A} (f(x^k) - f(x^*)) + \tilde{B}_1 \sigma_{1,k}^2 + \tilde{B}_2 \sigma_{2,k}^2 + \tilde{D}_1 \\ \mathbb{E} [\|g^k\|^2 \mid x^k] &\leq 2A' (f(x^k) - f(x^*)) + B'_1 \sigma_{1,k}^2 + B'_2 \sigma_{2,k}^2 + D'_1 \end{aligned}$$

$$\mathbb{E} [\sigma_{1,k+1}^2 \mid \sigma_{1,k}^2, \sigma_{2,k}^2] \leq (1 - \rho_1) \sigma_{1,k}^2 + 2C_1 (f(x^k) - f(x^*)) + G\rho_1 \sigma_{2,k}^2 + D_2$$

$$\mathbb{E} [\sigma_{2,k+1}^2 \mid \sigma_{2,k}^2] \leq (1 - \rho_2) \sigma_{2,k}^2 + 2C_2 (f(x^k) - f(x^*))$$

Reflects smoothness properties of the problem and noises introduced by compressions and stochastic gradients

Describes the process of variance reduction of the variance coming from compressions

Key Assumption

$$g^k = \frac{1}{n} \sum_{i=1}^n g_i^k, \quad \mathbb{E} [g^k \mid x^k] = \nabla f(x^k) \quad \bar{g}_i^k = \mathbb{E} [g_i^k \mid x^k]$$

$$\frac{1}{n} \sum_{i=1}^n \|\bar{g}_i^k\|^2 \leq 2A (f(x^k) - f(x^*)) + B_1 \sigma_{1,k}^2 + B_2 \sigma_{2,k}^2 + D_1$$

$$\frac{1}{n} \sum_{i=1}^n \mathbb{E} [\|g_i^k - \bar{g}_i^k\|^2 \mid x^k] \leq 2\tilde{A} (f(x^k) - f(x^*)) + \tilde{B}_1 \sigma_{1,k}^2 + \tilde{B}_2 \sigma_{2,k}^2 + \tilde{D}_1$$

$$\mathbb{E} [\|g^k\|^2 \mid x^k] \leq 2A' (f(x^k) - f(x^*)) + B'_1 \sigma_{1,k}^2 + B'_2 \sigma_{2,k}^2 + D'_1$$

$$\mathbb{E} [\sigma_{1,k+1}^2 \mid \sigma_{1,k}^2, \sigma_{2,k}^2] \leq (1 - \rho_1) \sigma_{1,k}^2 + 2C_1 (f(x^k) - f(x^*)) + G\rho_1 \sigma_{2,k}^2 + D_2$$

$$\mathbb{E} [\sigma_{2,k+1}^2 \mid \sigma_{2,k}^2] \leq (1 - \rho_2) \sigma_{2,k}^2 + 2C_2 (f(x^k) - f(x^*))$$

Reflects smoothness properties of the problem and noises introduced by compressions and stochastic gradients

Describes the process of variance reduction of the variance coming from compressions

Describes the process of variance reduction of the variance coming from stochastic gradients

Main Theorem

Some quantity depending only on the
starting point and stepsize

$$\mathbb{E} \left[f \left(\bar{x}^K \right) - f \left(x^* \right) \right] \leq (1 - \eta)^K \frac{\Psi \left(x^0, \gamma \right)}{\gamma} + \gamma \Phi \left(D_1, \tilde{D}_1, D'_1, D_2 \right)$$

Linear function

$$\eta = \min \left\{ \frac{\gamma \mu}{2}, \frac{\rho_1}{4}, \frac{\rho_2}{4} \right\}$$

Methods with Error Compensation Covered by Our Framework

Problem	Method	Alg #	Citation	Sec #	Rate (constants ignored)
(1)+(3)	EC-SGDsr	Alg 3	new	H.1	$\tilde{\mathcal{O}} \left(\frac{\mathcal{L}}{\mu} + \frac{\sqrt{L\mathcal{L}}}{\delta\mu} + \frac{\text{Var}}{n\mu\varepsilon} + \frac{\sqrt{L(\text{Var}+\zeta_*^2/\delta)}}{\mu\sqrt{\delta\varepsilon}} \right)$
(1)+(2)	EC-SGD	Alg 4	[45]	H.2	$\tilde{\mathcal{O}} \left(\frac{\kappa}{\delta} + \frac{\text{Var}}{n\mu\varepsilon} + \frac{\sqrt{L(\text{Var}+\zeta_*^2/\delta)}}{\delta\mu\sqrt{\varepsilon}} \right)$
(1)+(3)	EC-GDstar	Alg 5	new	H.3	$\mathcal{O} \left(\frac{\kappa}{\delta} \log \frac{1}{\varepsilon} \right)$
(1)+(2)	EC-SGD-DIANA	Alg 6	new	H.4	Option I: $\tilde{\mathcal{O}} \left(\omega + \frac{\kappa}{\delta} + \frac{\sigma^2}{n\mu\varepsilon} + \frac{\sqrt{L\sigma^2}}{\delta\mu\sqrt{\varepsilon}} \right)$ Option II: $\tilde{\mathcal{O}} \left(\frac{1+\omega}{\delta} + \frac{\kappa}{\delta} + \frac{\sigma^2}{n\mu\varepsilon} + \frac{\sqrt{L\sigma^2}}{\mu\sqrt{\delta\varepsilon}} \right)$
(1)+(3)	EC-SGDsr-DIANA	Alg 7	new	H.5	Option I: $\tilde{\mathcal{O}} \left(\omega + \frac{\mathcal{L}}{\mu} + \frac{\sqrt{L\mathcal{L}}}{\delta\mu} + \frac{\text{Var}}{n\mu\varepsilon} + \frac{\sqrt{L\text{Var}}}{\delta\mu\sqrt{\varepsilon}} \right)$ Option II: $\tilde{\mathcal{O}} \left(\frac{1+\omega}{\delta} + \frac{\mathcal{L}}{\mu} + \frac{\sqrt{L\mathcal{L}}}{\delta\mu} + \frac{\text{Var}}{n\mu\varepsilon} + \frac{\sqrt{L\text{Var}}}{\mu\sqrt{\delta\varepsilon}} \right)$
(1)+(2)	EC-GD-DIANA [†]	Alg 6	new	H.4	$\mathcal{O} \left(\left(\omega + \frac{\kappa}{\delta} \right) \log \frac{1}{\varepsilon} \right)$
(1)+(3)	EC-LSVRG	Alg 8	new	H.6	$\tilde{\mathcal{O}} \left(m + \frac{\kappa}{\delta} + \frac{\sqrt{L\zeta_*^2}}{\delta\mu\sqrt{\varepsilon}} \right)$
(1)+(3)	EC-LSVRGstar	Alg 9	new	H.7	$\mathcal{O} \left(\left(m + \frac{\kappa}{\delta} \right) \log \frac{1}{\varepsilon} \right)$
(1)+(3)	EC-LSVRG-DIANA	Alg 10	new	H.8	$\mathcal{O} \left(\left(\omega + m + \frac{\kappa}{\delta} \right) \log \frac{1}{\varepsilon} \right)$

Methods with Error Compensation Covered by Our Framework

Problem	Method	Alg #	Citation	Sec #	Rate (constants ignored)
(1)+(3)	EC-SGDsr	Alg 3	new	H.1	$\tilde{\mathcal{O}} \left(\frac{\mathcal{L}}{\mu} + \frac{\sqrt{L\mathcal{L}}}{\delta\mu} + \frac{\text{Var}}{n\mu\varepsilon} + \frac{\sqrt{L(\text{Var}+\zeta_*^2/\delta)}}{\mu\sqrt{\delta\varepsilon}} \right)$
(1)+(2)	EC-SGD	Alg 4	[45]	H.2	$\tilde{\mathcal{O}} \left(\frac{\kappa}{\delta} + \frac{\text{Var}}{n\mu\varepsilon} + \frac{\sqrt{L(\text{Var}+\zeta_*^2/\delta)}}{\delta\mu\sqrt{\varepsilon}} \right)$
(1)+(3)	EC-GDstar	Alg 5	new	H.3	$\mathcal{O} \left(\frac{\kappa}{\delta} \log \frac{1}{\varepsilon} \right)$
(1)+(2)	EC-SGD-DIANA	Alg 6	new	H.4	Option I: $\tilde{\mathcal{O}} \left(\omega + \frac{\kappa}{\delta} + \frac{\sigma^2}{n\mu\varepsilon} + \frac{\sqrt{L\sigma^2}}{\delta\mu\sqrt{\varepsilon}} \right)$ Option II: $\tilde{\mathcal{O}} \left(\frac{1+\omega}{\delta} + \frac{\kappa}{\delta} + \frac{\sigma^2}{n\mu\varepsilon} + \frac{\sqrt{L\sigma^2}}{\mu\sqrt{\delta\varepsilon}} \right)$
(1)+(3)	EC-SGDsr-DIANA	Alg 7	new	H.5	Option I: $\tilde{\mathcal{O}} \left(\omega + \frac{\mathcal{L}}{\mu} + \frac{\sqrt{L\mathcal{L}}}{\delta\mu} + \frac{\text{Var}}{n\mu\varepsilon} + \frac{\sqrt{L\text{Var}}}{\delta\mu\sqrt{\varepsilon}} \right)$ Option II: $\tilde{\mathcal{O}} \left(\frac{1+\omega}{\delta} + \frac{\mathcal{L}}{\mu} + \frac{\sqrt{L\mathcal{L}}}{\delta\mu} + \frac{\text{Var}}{n\mu\varepsilon} + \frac{\sqrt{L\text{Var}}}{\mu\sqrt{\delta\varepsilon}} \right)$
(1)+(2)	EC-GD-DIANA [†]	Alg 6	new	H.4	$\mathcal{O} \left(\left(\omega + \frac{\kappa}{\delta} \right) \log \frac{1}{\varepsilon} \right)$
(1)+(3)	EC-LSVRG	Alg 8	new	H.6	$\tilde{\mathcal{O}} \left(m + \frac{\kappa}{\delta} + \frac{\sqrt{L\zeta_*^2}}{\delta\mu\sqrt{\varepsilon}} \right)$
(1)+(3)	EC-LSVRGstar	Alg 9	new	H.7	$\mathcal{O} \left(\left(m + \frac{\kappa}{\delta} \right) \log \frac{1}{\varepsilon} \right)$
(1)+(3)	EC-LSVRG-DIANA	Alg 10	new	H.8	$\mathcal{O} \left(\left(\omega + m + \frac{\kappa}{\delta} \right) \log \frac{1}{\varepsilon} \right)$

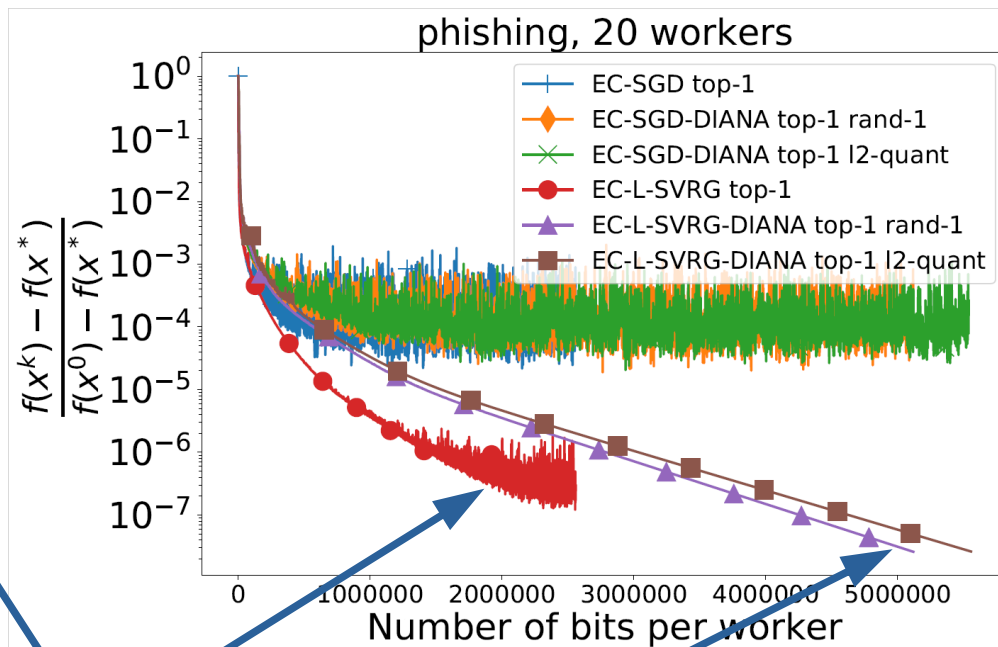
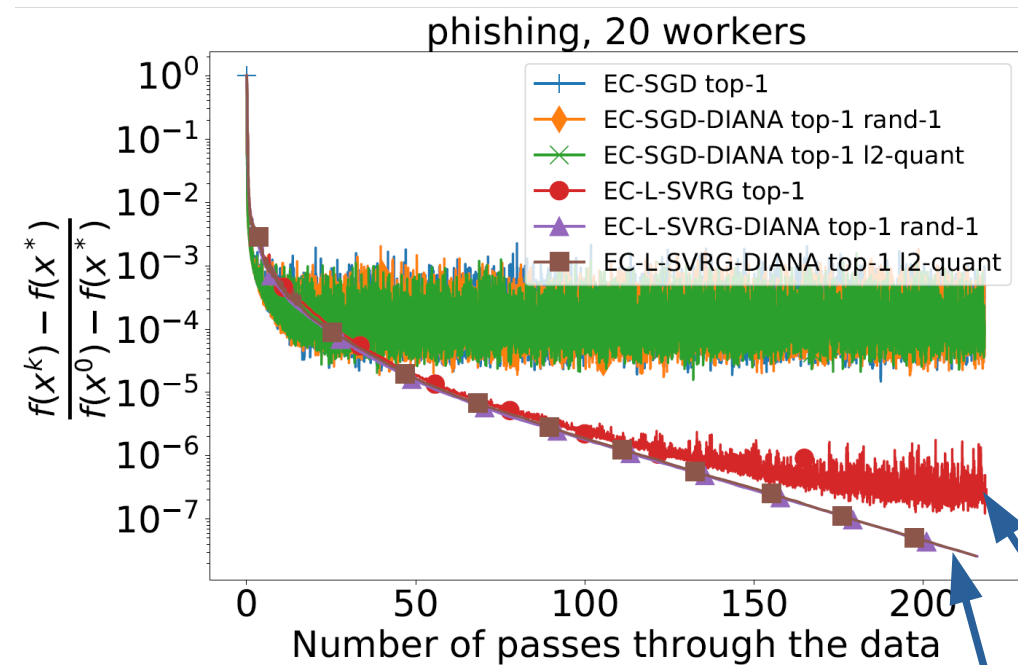
Methods with Error Compensation Covered by Our Framework

Problem	Method	Alg #	Citation	Sec #	Rate (constants ignored)
(1)+(3)	EC-SGDsr	Alg 3	new	H.1	$\tilde{\mathcal{O}} \left(\frac{\mathcal{L}}{\mu} + \frac{\sqrt{L\mathcal{L}}}{\delta\mu} + \frac{\text{Var}}{n\mu\varepsilon} + \frac{\sqrt{L(\text{Var}+\zeta_*^2/\delta)}}{\mu\sqrt{\delta\varepsilon}} \right)$
(1)+(2)	EC-SGD	Alg 4	[45]	H.2	$\tilde{\mathcal{O}} \left(\frac{\kappa}{\delta} + \frac{\text{Var}}{n\mu\varepsilon} + \frac{\sqrt{L(\text{Var}+\zeta_*^2/\delta)}}{\delta\mu\sqrt{\varepsilon}} \right)$
(1)+(3)	EC-GDstar	Alg 5	new	H.3	$\mathcal{O} \left(\frac{\kappa}{\delta} \log \frac{1}{\varepsilon} \right)$
(1)+(2)	EC-SGD-DIANA	Alg 6	new	H.4	Option I: $\tilde{\mathcal{O}} \left(\omega + \frac{\kappa}{\delta} + \frac{\sigma^2}{n\mu\varepsilon} + \frac{\sqrt{L\sigma^2}}{\delta\mu\sqrt{\varepsilon}} \right)$ Option II: $\tilde{\mathcal{O}} \left(\frac{1+\omega}{\delta} + \frac{\kappa}{\delta} + \frac{\sigma^2}{n\mu\varepsilon} + \frac{\sqrt{L\sigma^2}}{\mu\sqrt{\delta\varepsilon}} \right)$
(1)+(3)	EC-SGDsr-DIANA	Alg 7	new	H.5	Option I: $\tilde{\mathcal{O}} \left(\omega + \frac{\mathcal{L}}{\mu} + \frac{\sqrt{L\mathcal{L}}}{\delta\mu} + \frac{\text{Var}}{n\mu\varepsilon} + \frac{\sqrt{L\text{Var}}}{\delta\mu\sqrt{\varepsilon}} \right)$ Option II: $\tilde{\mathcal{O}} \left(\frac{1+\omega}{\delta} + \frac{\mathcal{L}}{\mu} + \frac{\sqrt{L\mathcal{L}}}{\delta\mu} + \frac{\text{Var}}{n\mu\varepsilon} + \frac{\sqrt{L\text{Var}}}{\mu\sqrt{\delta\varepsilon}} \right)$
(1)+(2)	EC-GD-DIANA [†]	Alg 6	new	H.4	$\mathcal{O} \left(\left(\omega + \frac{\kappa}{\delta} \right) \log \frac{1}{\varepsilon} \right)$
(1)+(3)	EC-LSVRG	Alg 8	new	H.6	$\tilde{\mathcal{O}} \left(m + \frac{\kappa}{\delta} + \frac{\sqrt{L\zeta_*^2}}{\delta\mu\sqrt{\varepsilon}} \right)$
(1)+(3)	EC-LSVRGstar	Alg 9	new	H.7	$\mathcal{O} \left(\left(m + \frac{\kappa}{\delta} \right) \log \frac{1}{\varepsilon} \right)$
(1)+(3)	EC-LSVRG-DIANA	Alg 10	new	H.8	$\mathcal{O} \left(\left(\omega + m + \frac{\kappa}{\delta} \right) \log \frac{1}{\varepsilon} \right)$

Our framework covers even methods without error compensation and methods with delayed updates

7. Experiments

Logistic Regression with l2-regularization



partial variance reduction

full variance reduction

8. More Methods

More Methods Fitting our Framework

The generality of our approach helps to obtain convergence guarantees for a big number of different stochastic methods (even without error compensation). Here are some examples.

- Methods without error feedback: SGD, SGD-SR (arbitrary sampling), SAGA, SVRG, L-SVRG, QSGD, TernGrad, DQGD, DIANA, **DIANAsr-DQ**, VR-DIANA, JacSketch, SEGA
- Methods with delayed updates: D-SGD, **D-SGD-SR** (arbitrary sampling), **D-QSGD**, **D-SGD-DIANA**, **D-LSVRG**, **D-QLSVRG**, **D-LSVRG-DIANA**

More Methods Fitting our Framework

The generality of our approach helps to obtain convergence guarantees for a big number of different stochastic methods (even without error compensation). Here are some examples.

- Methods without error feedback: SGD, SGD-SR (arbitrary sampling), SAGA, SVRG, L-SVRG, QSGD, TernGrad, DQGD, DIANA, **DIANAsr-DQ**, VR-DIANA, JacSketch, SEGA
- Methods with delayed updates: D-SGD, **D-SGD-SR** (arbitrary sampling), **D-QSGD**, **D-SGD-DIANA**, **D-LSVRG**, **D-QLSVRG**, **D-LSVRG-DIANA**
- ✓ In one theorem, we recover the sharpest rates for all known special cases
- ✓ 16 new methods
- ✓ Our analysis works for weakly convex objectives as well

Thank you for watching!

bold font = new method