

Communication-efficient distributed optimization algorithms

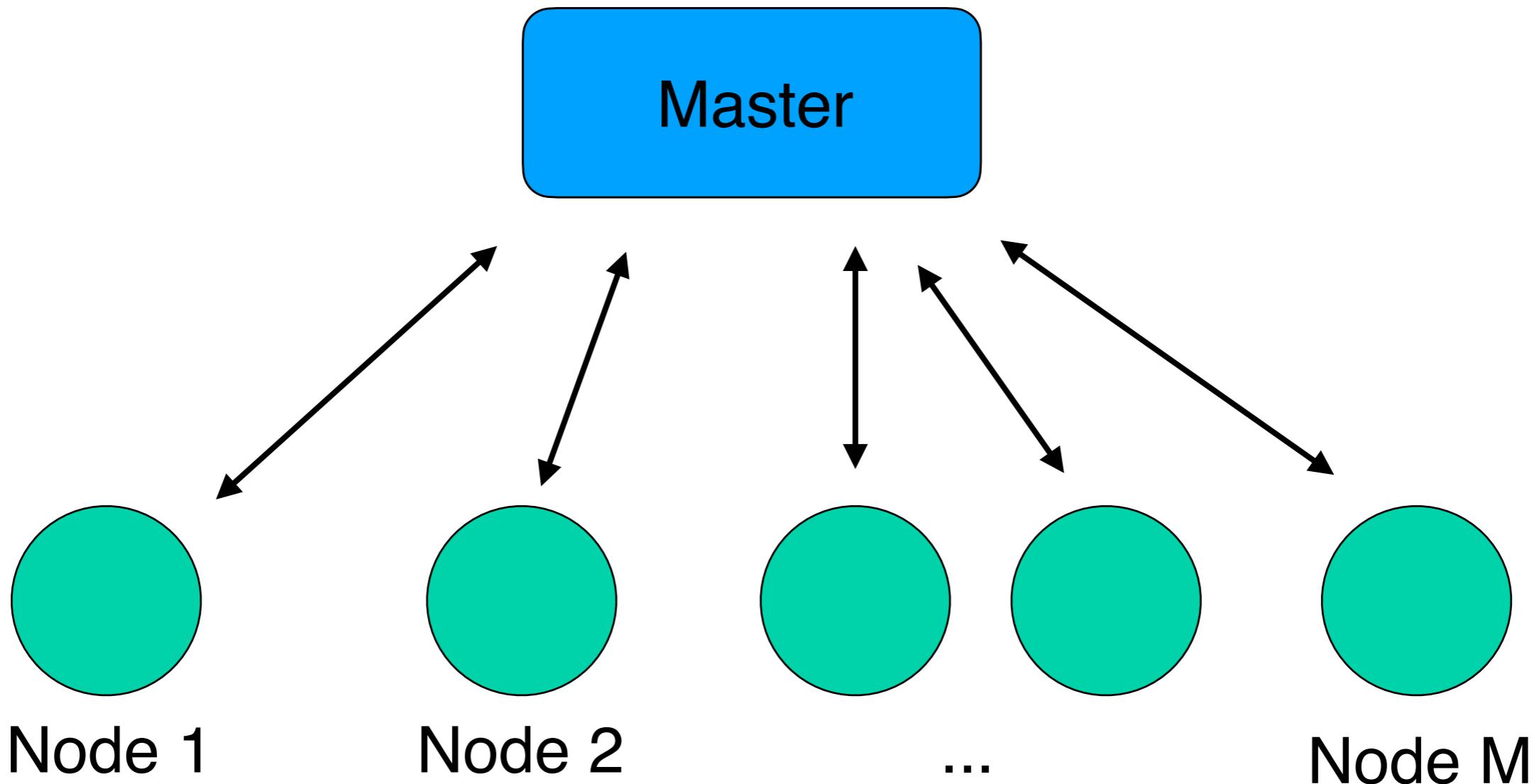
Laurent Condat, Peter Richtárik

King Abdullah University of Science and Technology
(KAUST)
Thuwal, Saudi Arabia

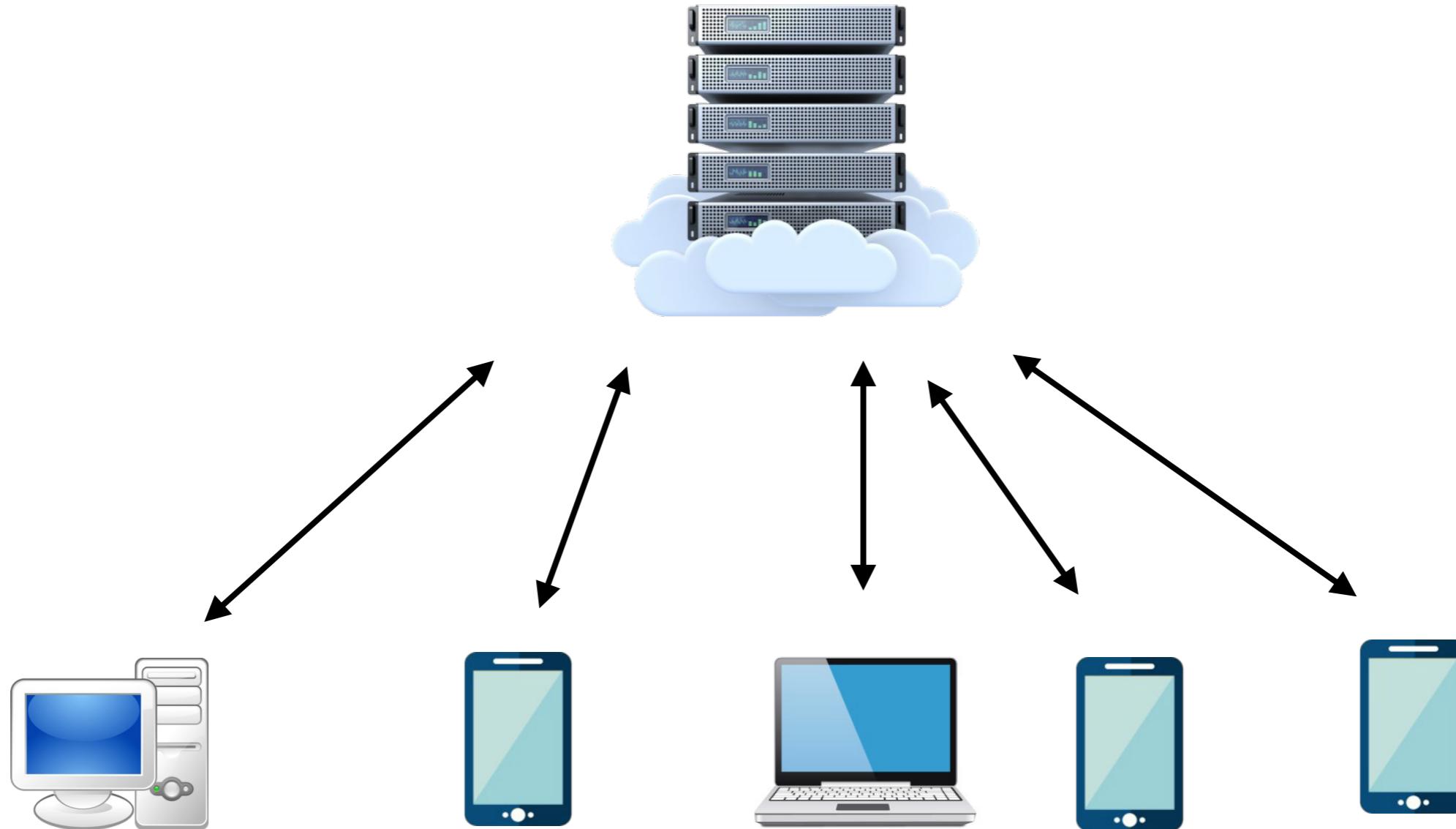


4th IMA Conf. Mathematical
Challenges of Big Data, Sept. 2022

Distributed computing



Federated learning



Federated learning



Convex optimization problem

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad R(x) + \frac{1}{M} \sum_{m=1}^M F_m(x)$$

- every function F_m is convex and L -smooth, for some $L > 0$, and μ -strongly convex, for some $\mu > 0$, i.e. $F - \frac{\mu}{2} \|\cdot\|^2$ is convex.

Convex optimization problem

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad R(x) + \frac{1}{M} \sum_{m=1}^M F_m(x)$$

- every function F_m is convex and L -smooth, for some $L > 0$, and μ -strongly convex, for some $\mu > 0$, i.e. $F - \frac{\mu}{2} \|\cdot\|^2$ is convex.
- $R : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ is a proper, closed, convex function whose proximity operator

$$\text{prox}_{\gamma R} : x \mapsto \arg \min_w \left(\gamma R(w) + \frac{1}{2} \|x - w\|^2 \right)$$

is easy to compute

Convex optimization problem

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad R(x) + \frac{1}{M} \sum_{m=1}^M F_m(x)$$

Prox-GD:

$$x^{k+1} := \text{prox}_{\gamma R} \left(x^k - \frac{\gamma}{M} \sum_{m=1}^M \nabla F_m(x^k) \right)$$

Convex optimization problem

$$\underset{x \in \mathbb{R}^d}{\text{minimize}} \quad R(x) + \frac{1}{M} \sum_{m=1}^M F_m(x)$$

Prox-GD:

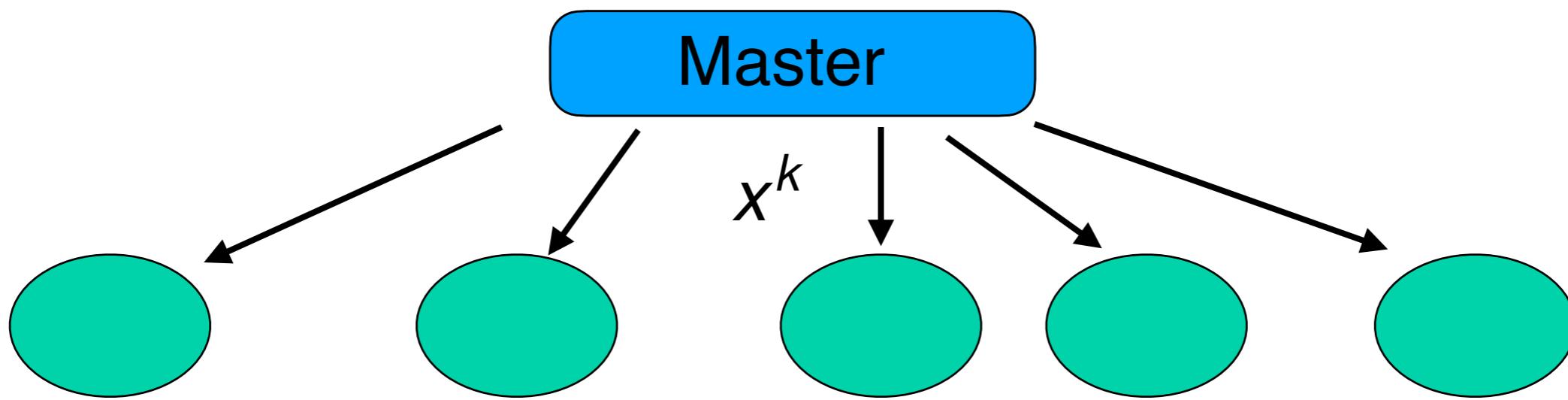
$$x^{k+1} := \text{prox}_{\gamma R} \left(x^k - \frac{\gamma}{M} \sum_{m=1}^M \nabla F_m(x^k) \right)$$

$$0 < \gamma \leq \frac{2}{L + \mu}$$

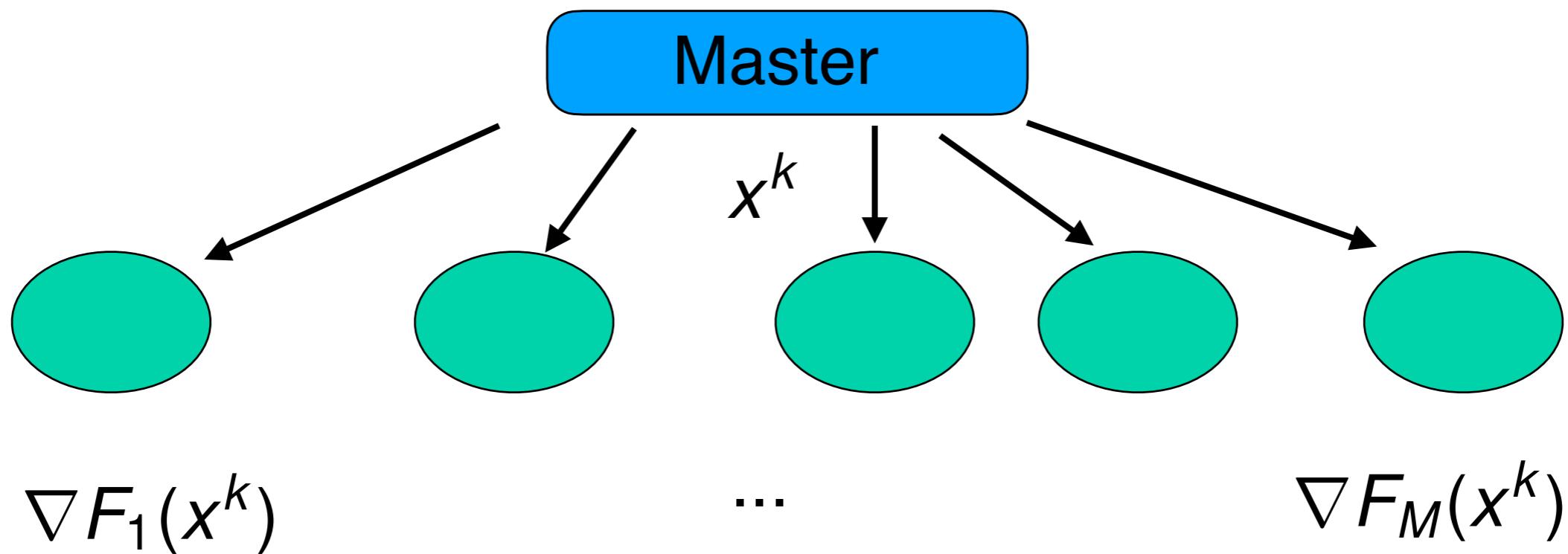


$$\|x^k - x^*\| \leq (1 - \gamma\mu)^k \|x^0 - x^*\|$$

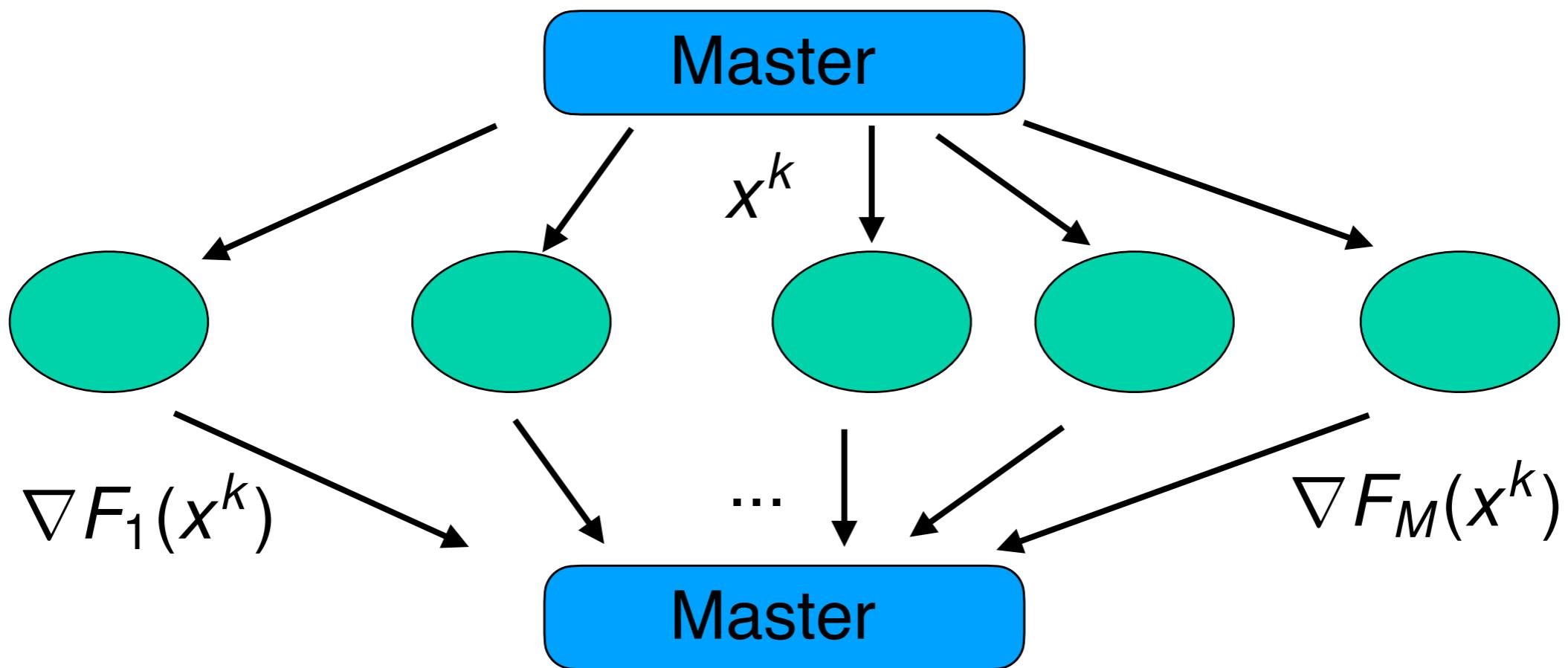
Distributed prox. GD



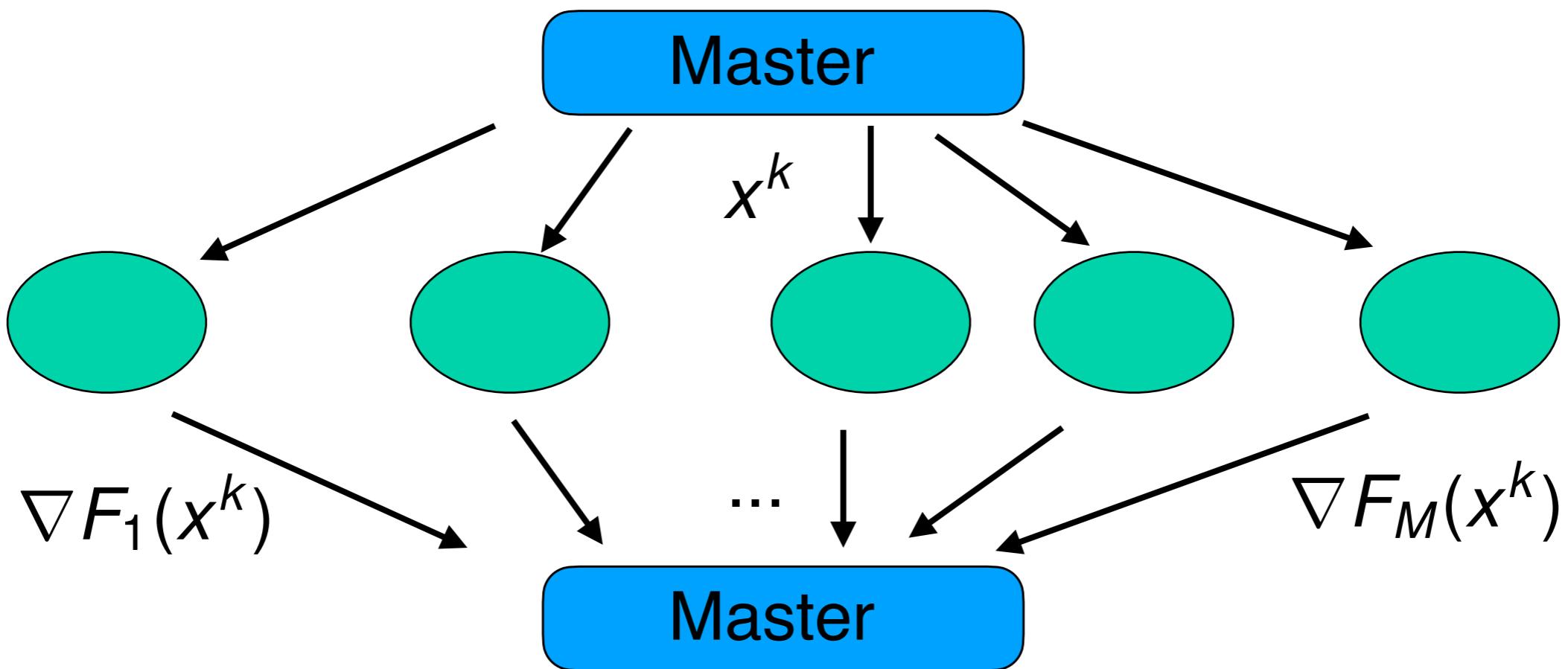
Distributed prox. GD



Distributed prox. GD

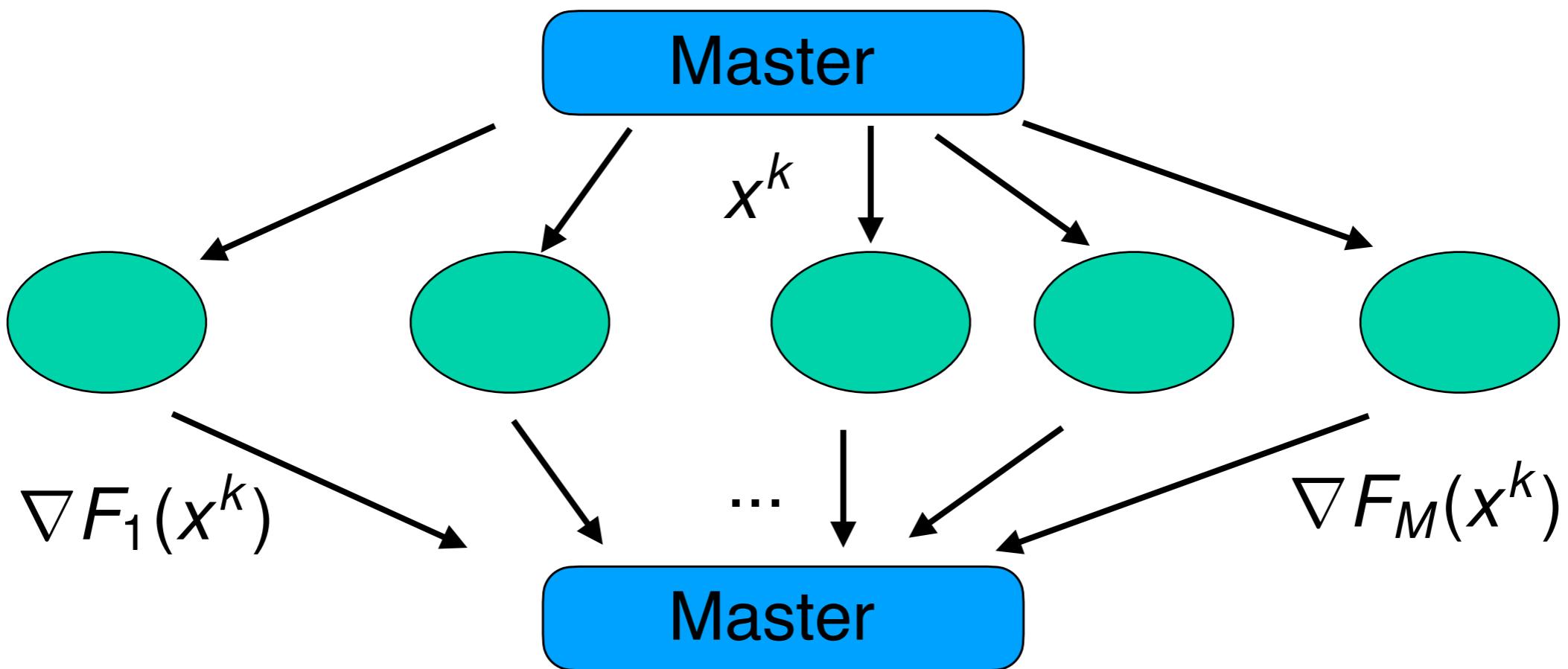


Distributed prox. GD



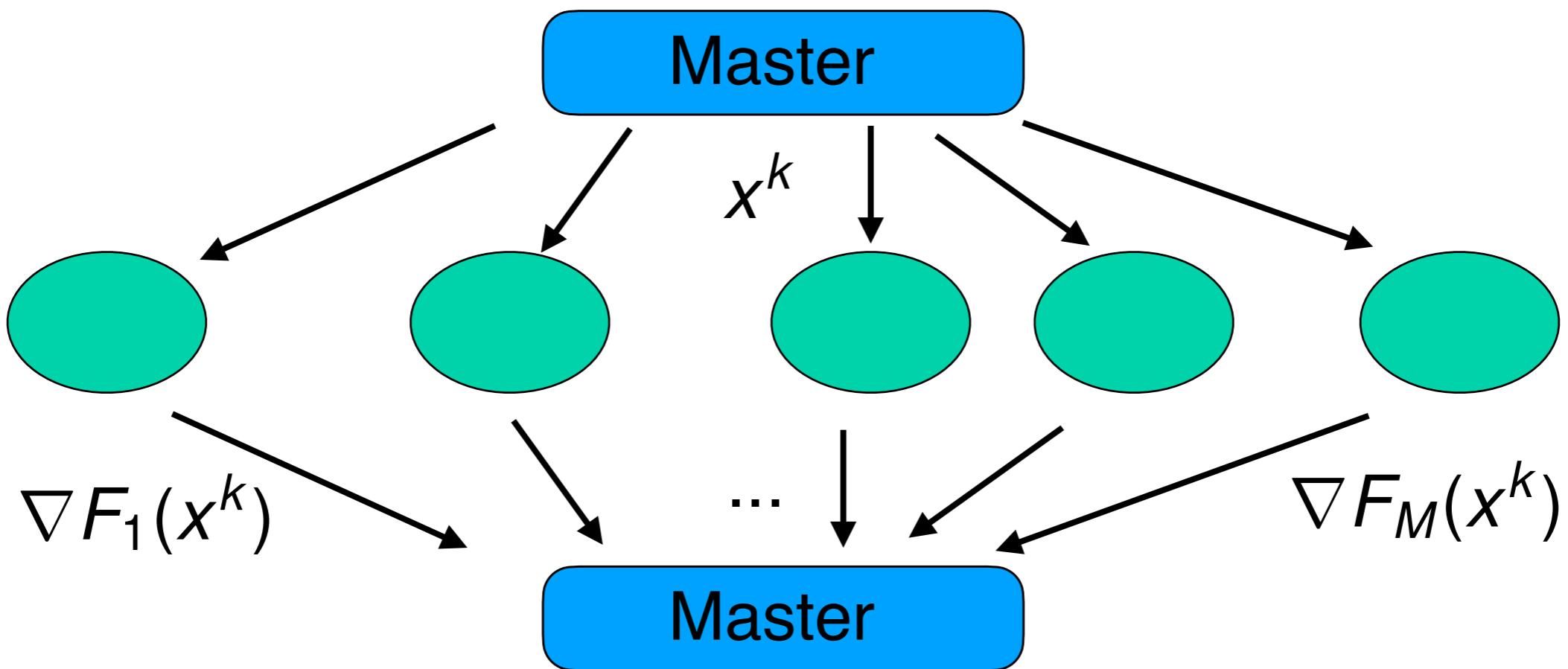
$$\frac{1}{M} \sum_{m=1}^M \nabla F_m(x^k)$$

Distributed prox. GD



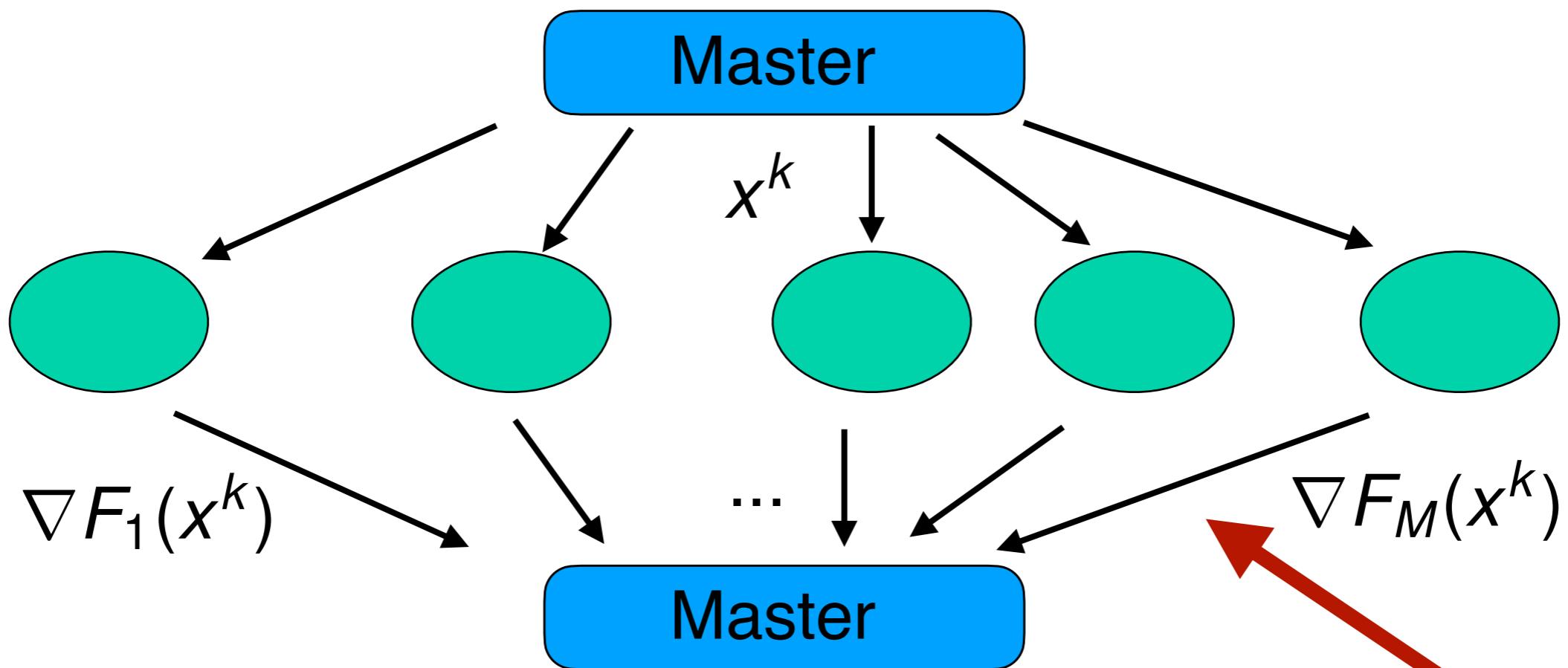
$$x^k - \frac{\gamma}{M} \sum_{m=1}^M \nabla F_m(x^k)$$

Distributed prox. GD



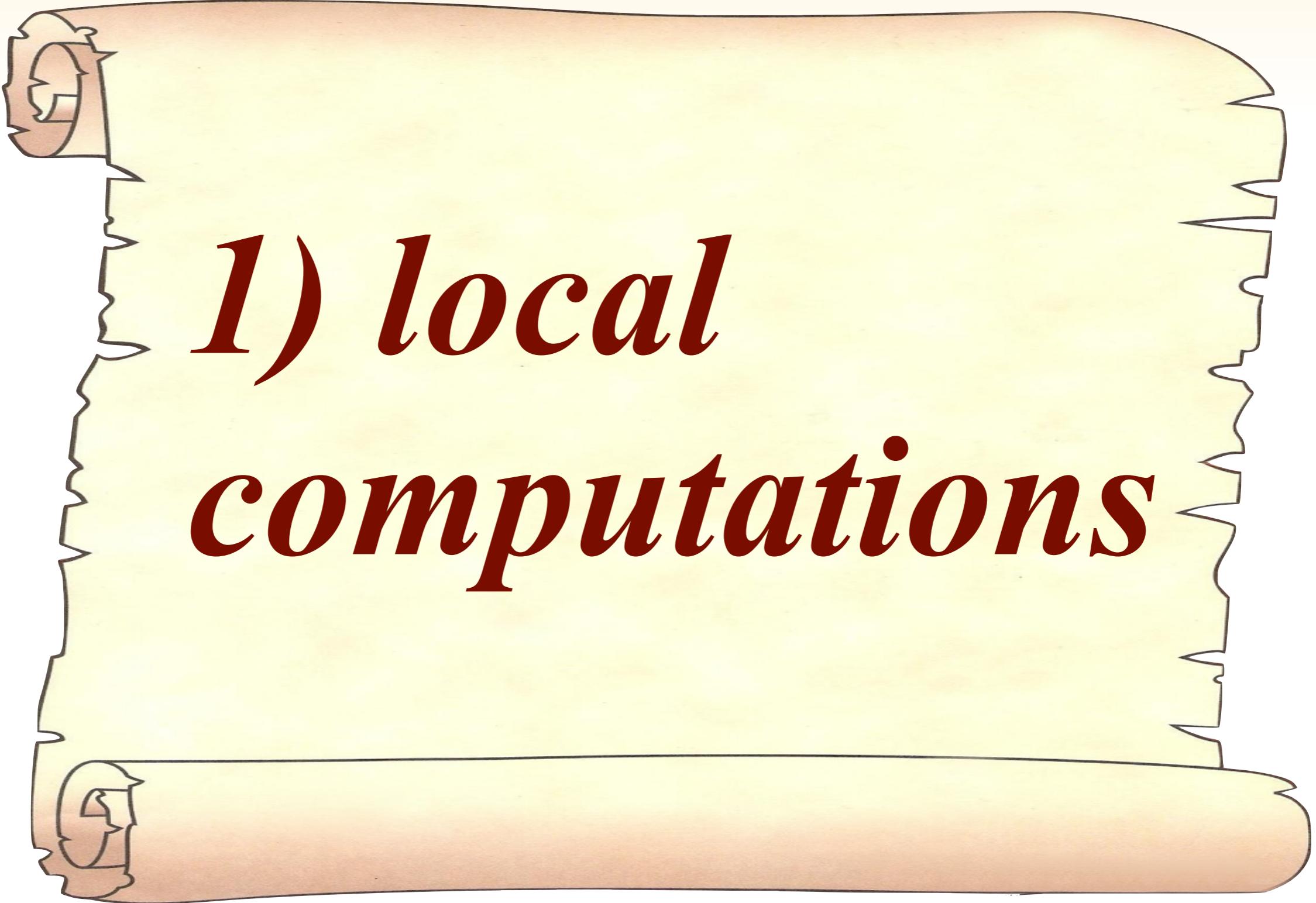
$$x^{k+1} := \text{prox}_{\gamma R} \left(x^k - \frac{\gamma}{M} \sum_{m=1}^M \nabla F_m(x^k) \right)$$

Distributed prox. GD



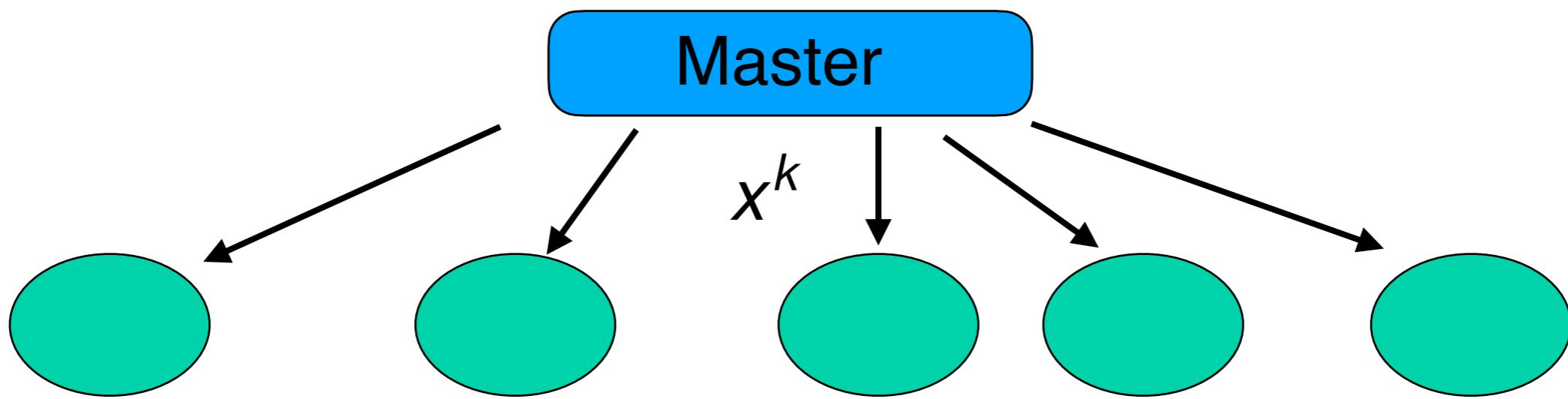
$$x^{k+1} := \text{prox}_{\gamma R} \left(x^k - \frac{\gamma}{M} \sum_{m=1}^M \nabla F_m(x^k) \right)$$





*1) local
computations*

Distributed prox. GD

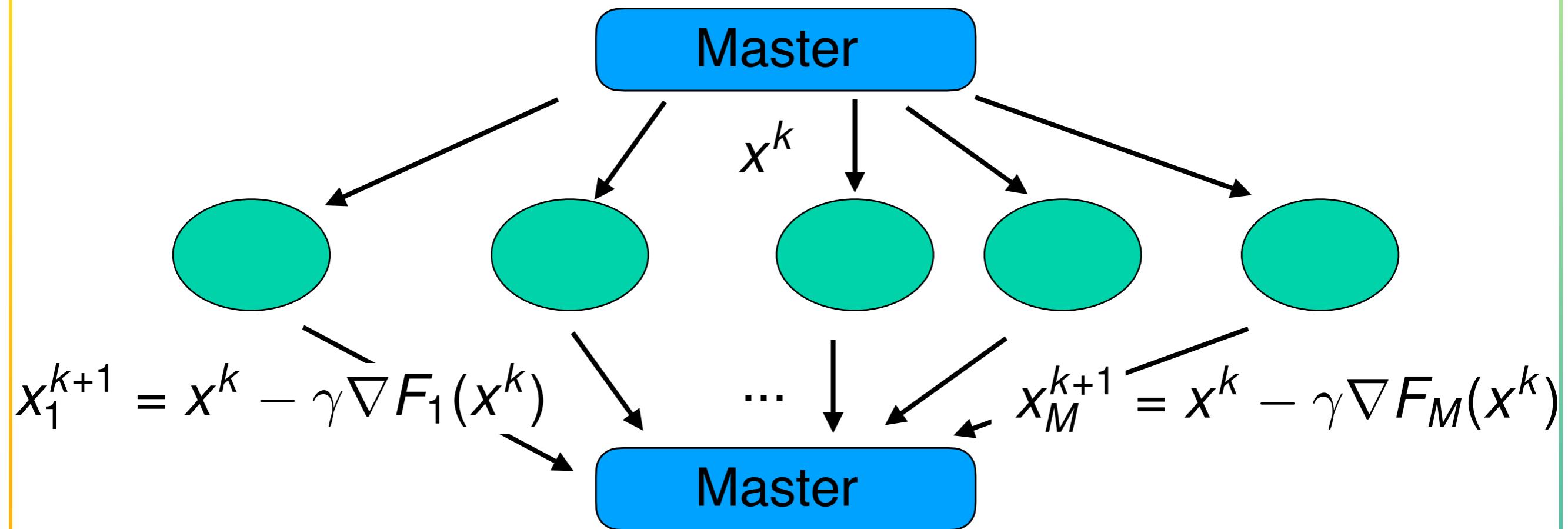


$$x_1^{k+1} = x^k - \gamma \nabla F_1(x^k)$$

...

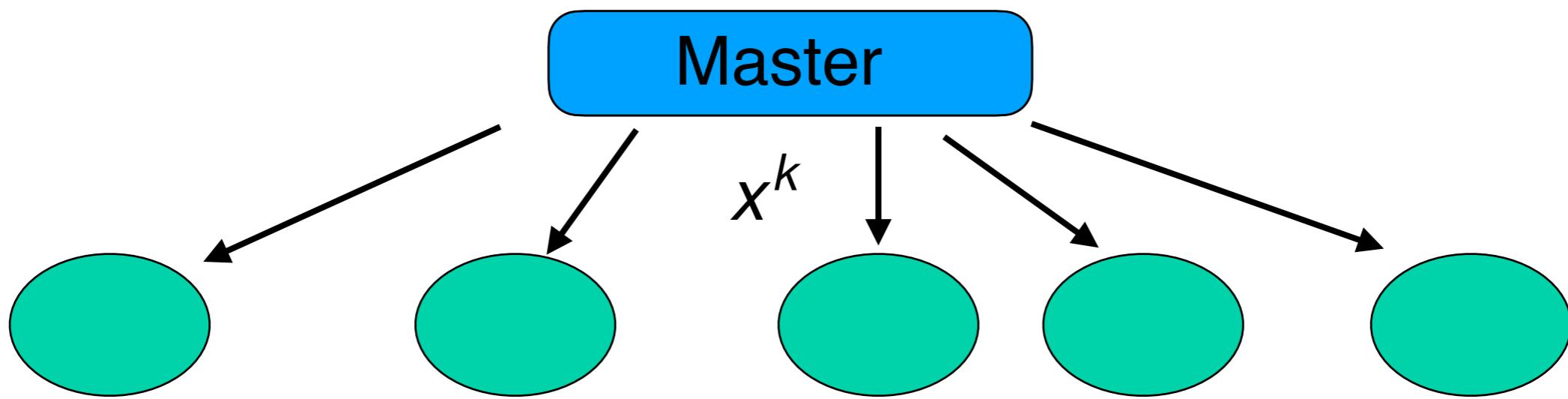
$$x_M^{k+1} = x^k - \gamma \nabla F_M(x^k)$$

Distributed prox. GD



$$x^{k+1} := \text{prox}_{\gamma R} \left(\frac{1}{M} \sum_{m=1}^M x_m^{k+1} \right)$$

Distributed prox. Local GD



$$x_1^{k+1} = x^k - \gamma \nabla F_1(x^k)$$

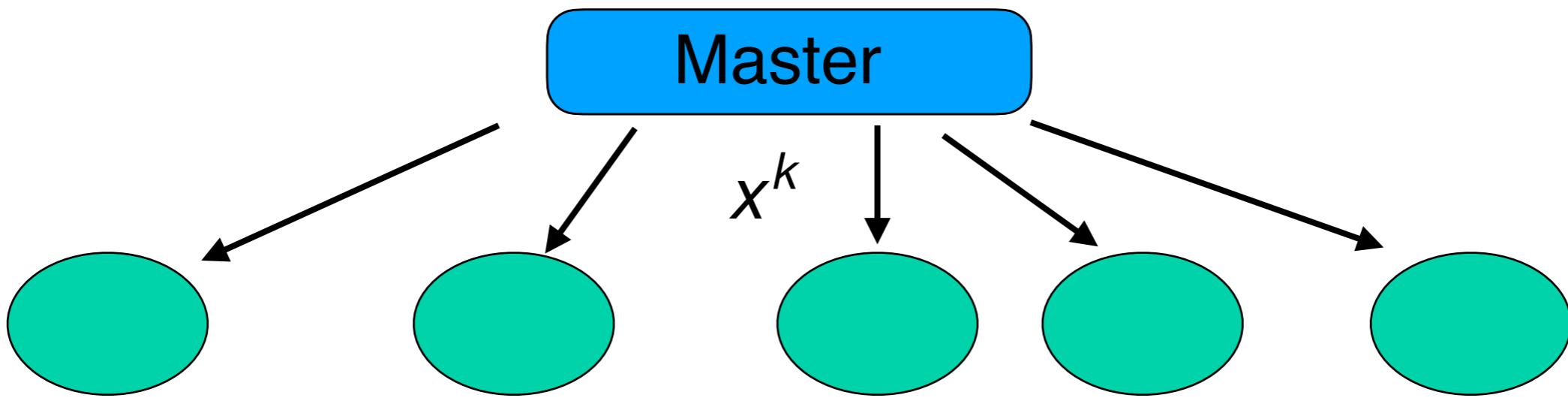
...

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

$$x_1^{k+2} = x_1^{k+1} - \gamma \nabla F_1(x_1^{k+1})$$

$$x_M^{k+2} = x_M^{k+1} - \gamma \nabla F_M(x_M^{k+1})$$

Distributed prox. Local GD



$$x_1^{k+1} = x^k - \gamma \nabla F_1(x^k)$$

...

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

$$x_1^{k+2} = x_1^{k+1} - \gamma \nabla F_1(x_1^{k+1})$$

$$x_M^{k+2} = x_M^{k+1} - \gamma \nabla F_M(x_M^{k+1})$$

...

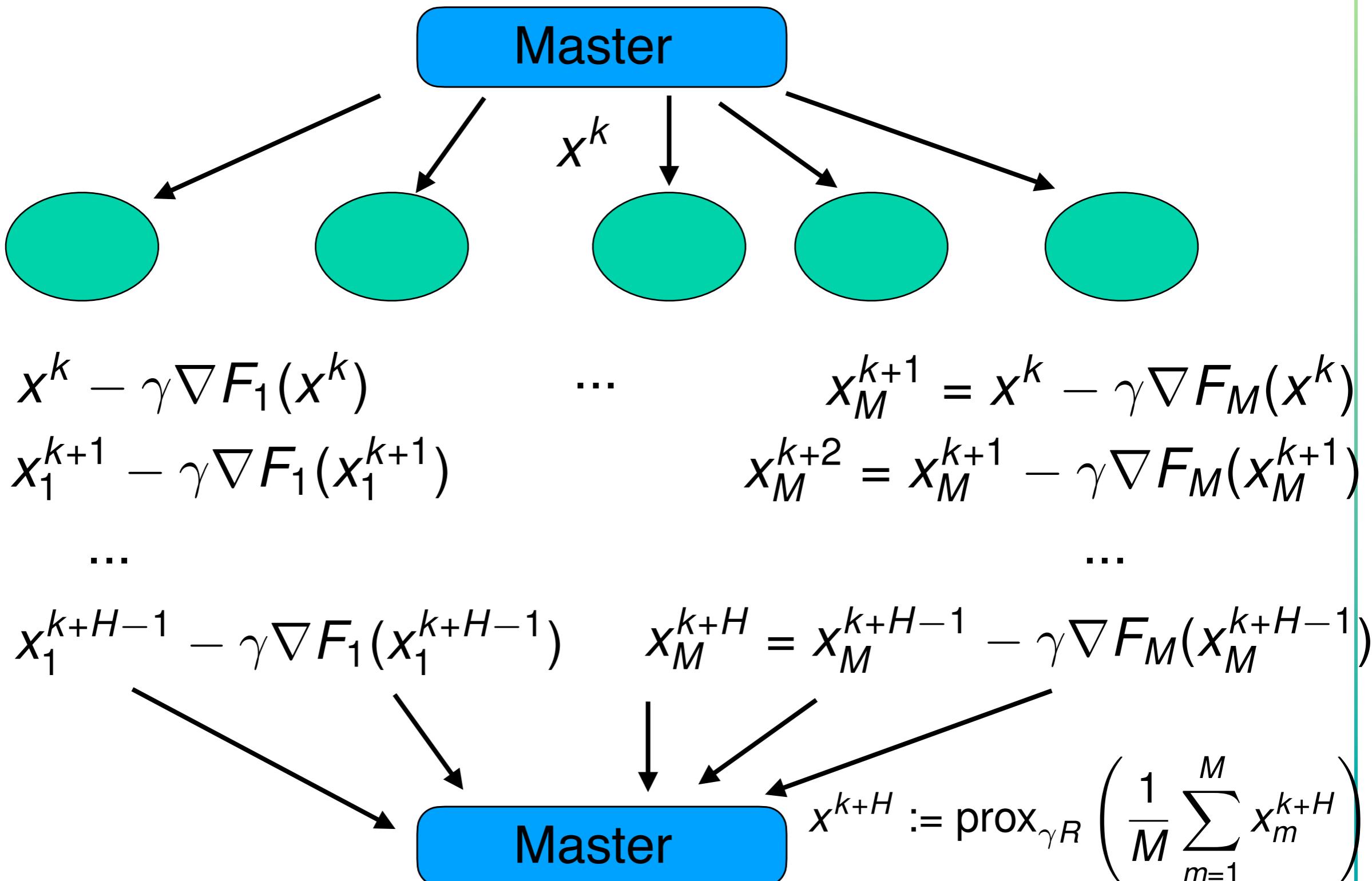
...

$$x_1^{k+H} = x_1^{k+H-1} - \gamma \nabla F_1(x_1^{k+H-1})$$

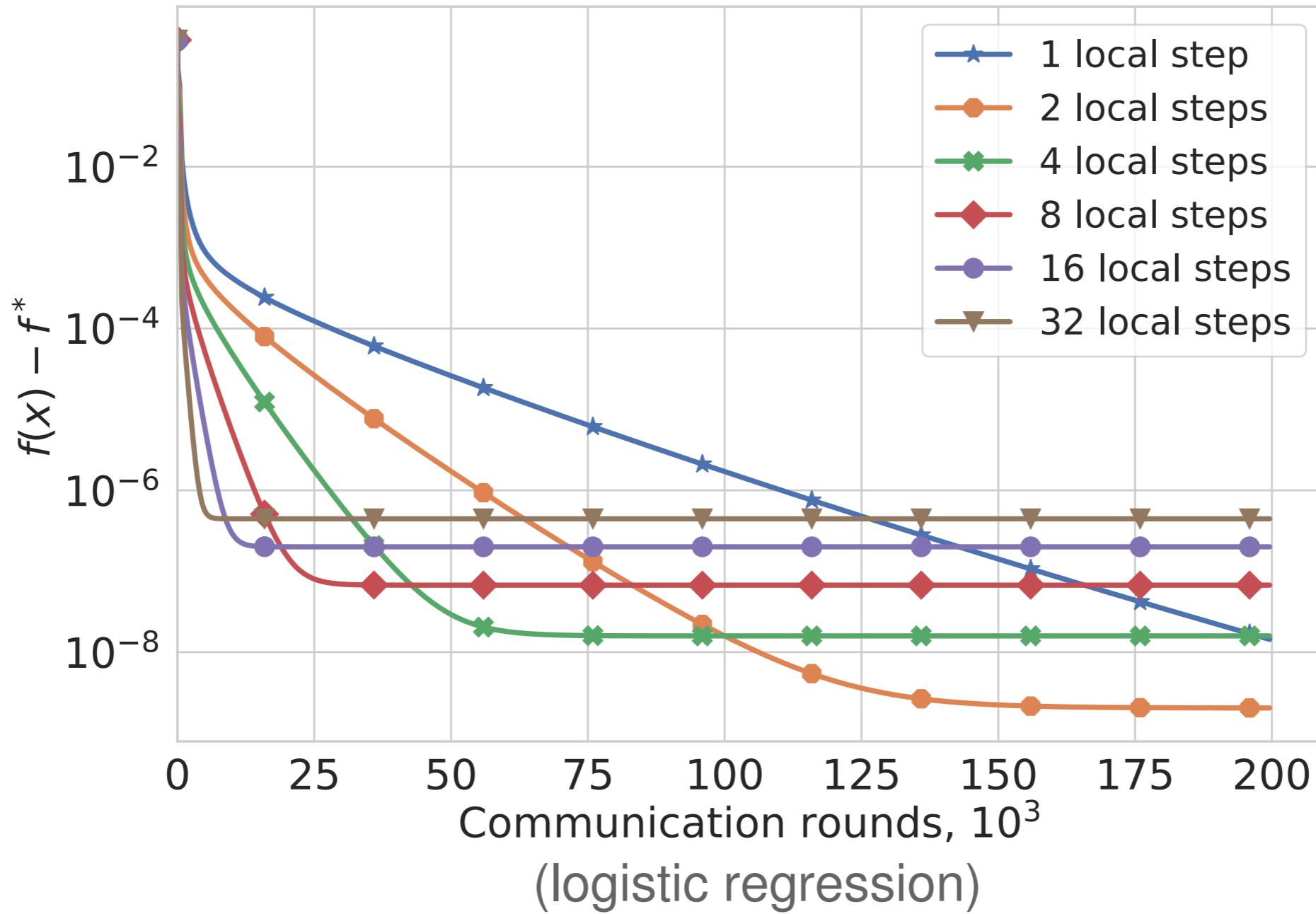
$$x_M^{k+H} = x_M^{k+H-1} - \gamma \nabla F_M(x_M^{k+H-1})$$

$$H \geq 1$$

Distributed prox. Local GD



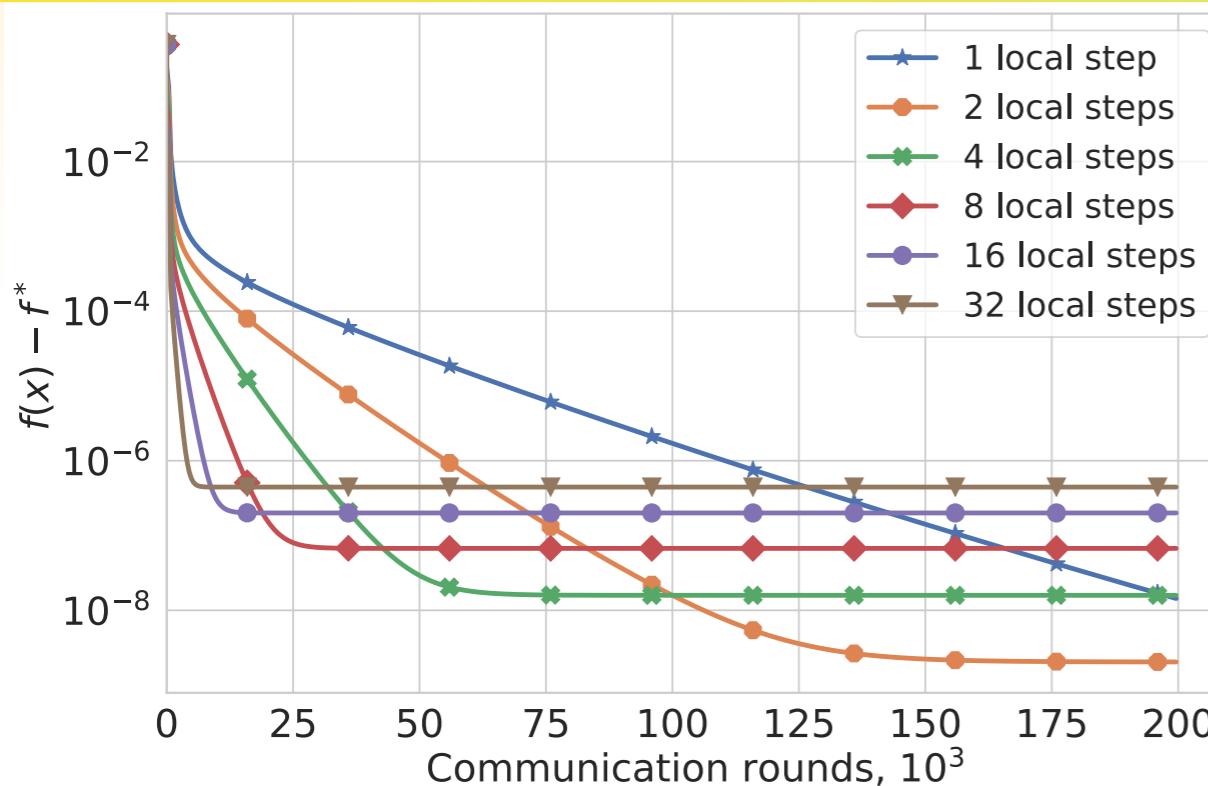
Local GD: performance



Local GD: analysis

- * Stich, Local SGD Converges Fast and Communicates Little. *ICLR* 2019.
- * Khaled, Mishchenko, and Richtárik, First analysis of local GD on heterogeneous data. *NeurIPS Workshop on Federated Learning for Data Privacy and Confidentiality*, 2019.
- * Khaled, Mishchenko, and Richtárik, Tighter theory for local SGD on identical and heterogeneous data. *AISTATS* 2020.
- * Ma, Konecny, Jaggi, Smith, Jordan, Richtárik, and Takáć, Distributed optimization with arbitrary local solvers. *Optimization Methods and Software*, 2017.
- * Haddadpour and Mahdavi, On the convergence of local descent methods in federated learning, *arXiv:1910.14425*, 2019.

Local GD: analysis



Malinovsky, Kovalev, Gasanov, Condat, Richtárik, “From local SGD to local fixed point methods for federated learning,”
ICML 2020

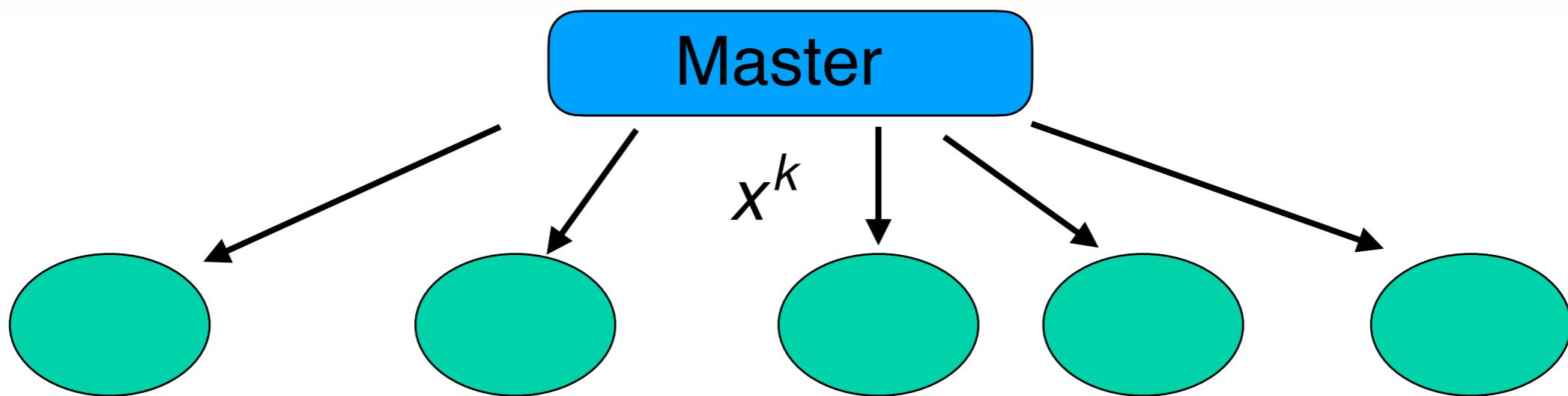
Theorem 2.11 (linear convergence) With $\gamma \in (0, \frac{2}{L+\mu}]$, $(x^{nH})_{n \geq 0}$ converges linearly to x^\dagger with rate ξ^H , where $\xi = 1 - \gamma\mu$, and

$$\|x^\dagger - x^*\| \leq S,$$

where

$$S = \frac{\xi}{1 - \xi} \frac{1 - \xi^{H-1}}{1 - \xi^H} \frac{1}{M} \sum_{m=1}^M \|\nabla F_m(x^*)\|.$$

Variance-reduced local GD



$$x_1^{k+1} = x^k - \gamma \nabla F_1(x^k) + h_1^k$$

...

$$x_M^{k+1} = x^k - \gamma \nabla F_M(x^k) + h_M^k$$

$$x_1^{k+2} = x_1^{k+1} - \gamma \nabla F_1(x_1^{k+1}) + h_1^k$$

$$x_M^{k+2} = x_M^{k+1} - \gamma \nabla F_M(x_M^{k+1}) + h_M^k$$

...

$$x_1^{k+H} = x_1^{k+H-1} - \gamma \nabla F_1(x_1^{k+H-1}) + h_1^k$$

$$x_M^{k+H} = x_M^{k+H-1} - \gamma \nabla F_M(x_M^{k+H-1}) + h_M^k$$

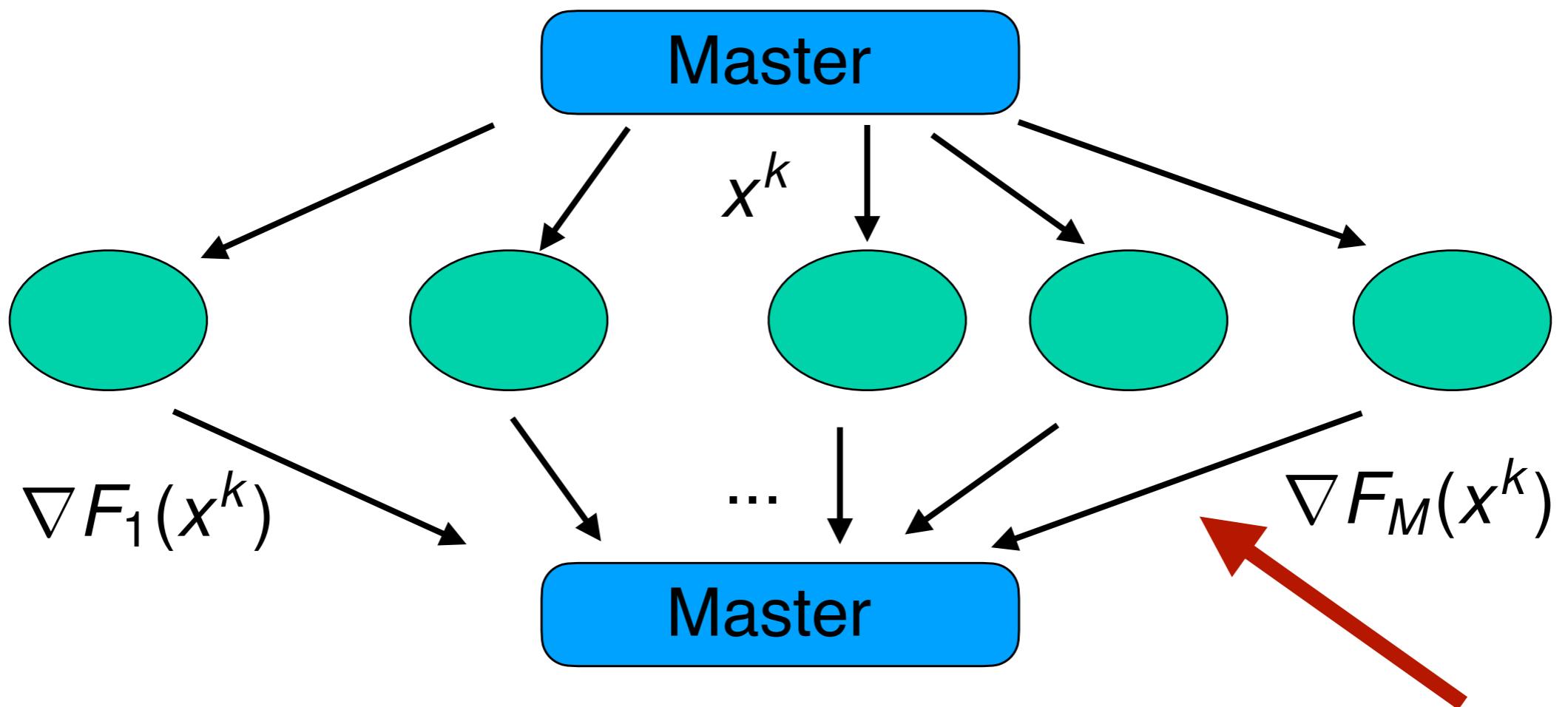
Mishchenko, Malinovsky, Stich, Richtárik, “ProxSkip: Yes! Local Gradient Steps Provably Lead to Communication Acceleration!” ICML 2022

Condat and Richtárik, “RandProx: Primal-Dual Optimization Algorithms with Randomized Proximal Updates,” arXiv:2207.12891, 2022



2) compression

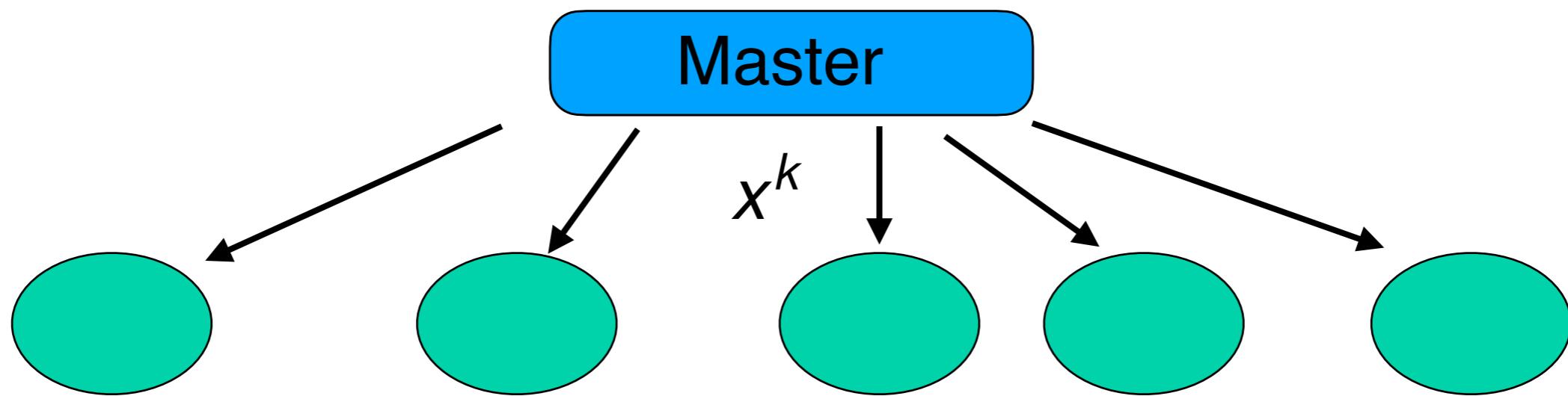
Dist. prox. GD



$$x^{k+1} := \text{prox}_{\gamma R} \left(x^k - \frac{\gamma}{M} \sum_{m=1}^M \nabla F_m(x^k) \right)$$

compression

Dist. prox. GD with compression

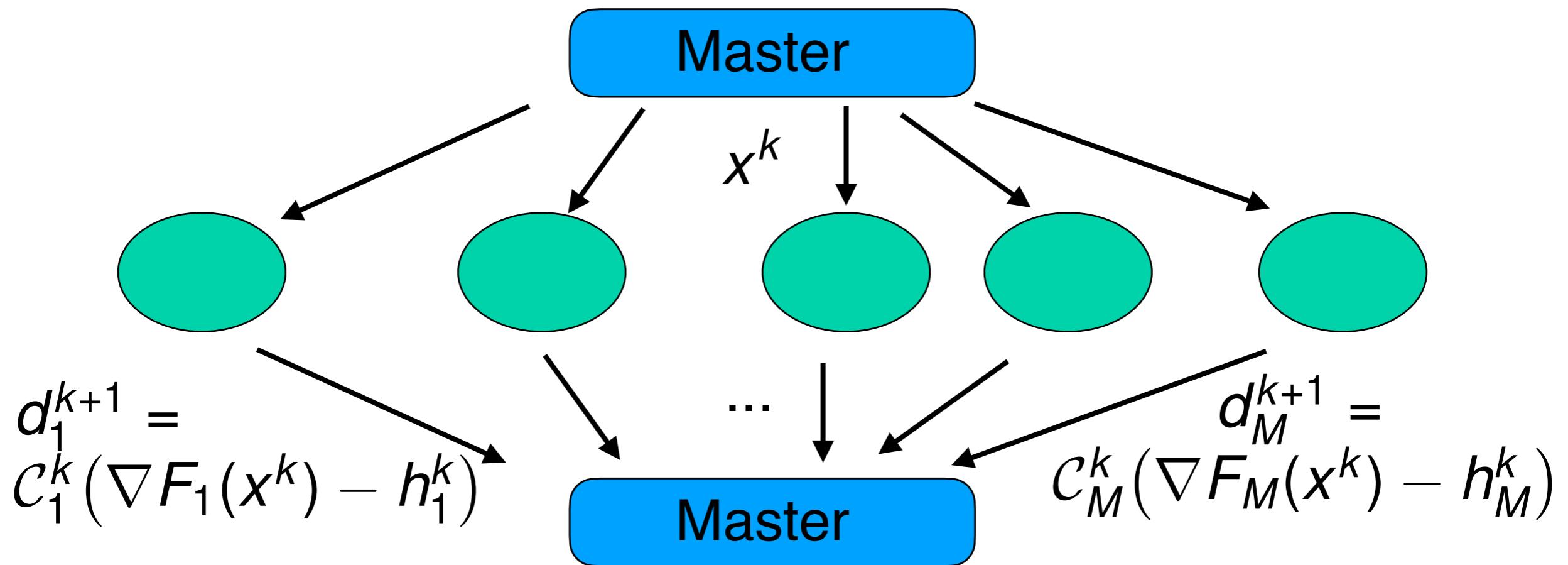


$$d_1^{k+1} = \mathcal{C}_1^k (\nabla F_1(x^k) - h_1^k)$$

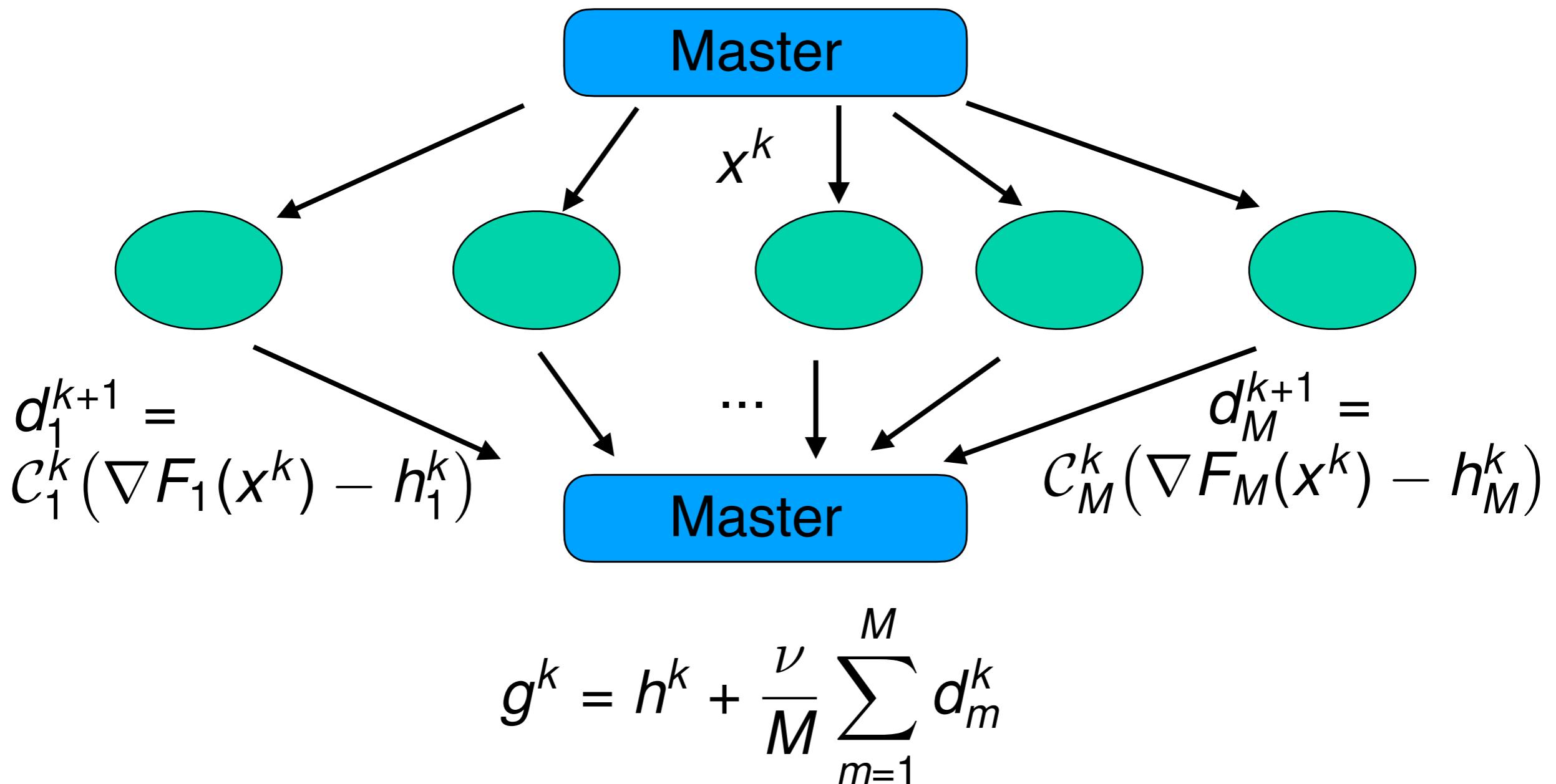
...

$$d_M^{k+1} = \mathcal{C}_M^k (\nabla F_M(x^k) - h_M^k)$$

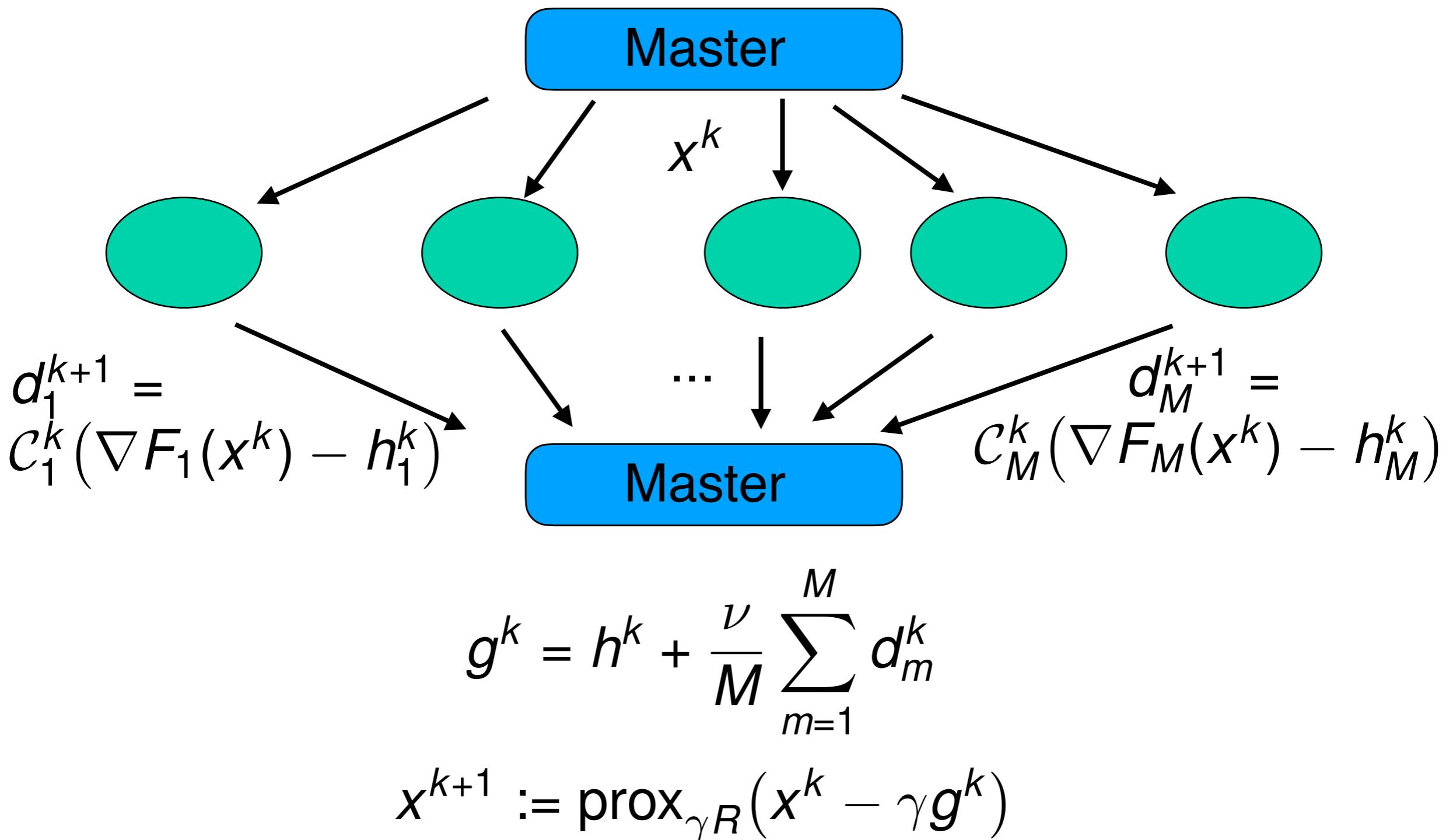
Dist. prox. GD with compression



Dist. prox. GD with compression



Dist. prox. GD with compression



Dist. prox. GD with compression

Algorithm 1 (EF-BV)

```
1: input: parameters  $\gamma > 0$ ,  $\lambda > 0$ ,  $\nu > 0$ ,  
2: initial vectors  $x^0 \in \mathbb{R}^d$  and  $h_m^0 \in \mathbb{R}^d$   
3:  $h^0 := \frac{1}{M} \sum_{m=1}^M h_m^0$   
4: for  $k = 0, 1, \dots$  do  
5:   for  $m = 1, \dots, M$  in parallel do  
6:      $d_m^{k+1} := C_m^k (\nabla F_m(x^k) - h_m^k)$   
7:      $h_m^{k+1} := h_m^k + \lambda d_m^{k+1}$   
8:   end for  
9:   // at master:  
10:   $d^{k+1} := \frac{1}{M} \sum_{m=1}^M d_m^{k+1}$   
11:   $x^{k+1} := \text{prox}_{\gamma R}(x^k - \gamma(h^k + \nu d^{k+1}))$   
12:   $h^{k+1} := h^k + \lambda d^{k+1}$   
13: end for
```

Condat, Yi, Richtárik,
“EF-BV: A unified theory
of error feedback and
variance reduction...”,
NeurIPS 2022

Dist. prox. GD with compression

Algorithm 1 (EF-BV)

```
1: input: parameters  $\gamma > 0$ ,  $\lambda > 0$ ,  $\nu > 0$ ,  
2: initial vectors  $x^0 \in \mathbb{R}^d$  and  $h_m^0 \in \mathbb{R}^d$   
3:  $h^0 := \frac{1}{M} \sum_{m=1}^M h_m^0$   
4: for  $k = 0, 1, \dots$  do  
5:   for  $m = 1, \dots, M$  in parallel do  
6:      $d_m^{k+1} := C_m^k (\nabla F_m(x^k) - h_m^k)$   
7:      $h_m^{k+1} := h_m^k + \lambda d_m^{k+1}$   
8:   end for  
9:   // at master:  
10:   $d^{k+1} := \frac{1}{M} \sum_{m=1}^M d_m^{k+1}$   
11:   $x^{k+1} := \text{prox}_{\gamma R}(x^k - \gamma(h^k + \nu d^{k+1}))$   
12:   $h^{k+1} := h^k + \lambda d^{k+1}$   
13: end for
```

Condat, Yi, Richtárik,
“EF-BV: A unified theory
of error feedback and
variance reduction...”,
NeurIPS 2022

$\nu = 1$ and unbiased
compressors: DIANA
[Mishchenko et al. 2019]
generalized in:

Condat and Richtárik,
“MURANA: A Generic
Framework for
Stochastic Variance-
Reduced Optimization,”
MSML 2022

DIANA

$$x^{k+1} := \text{prox}_{\gamma R} \left(x^k - \frac{\gamma}{M} \sum_{m=1}^M g_m^k \right)$$

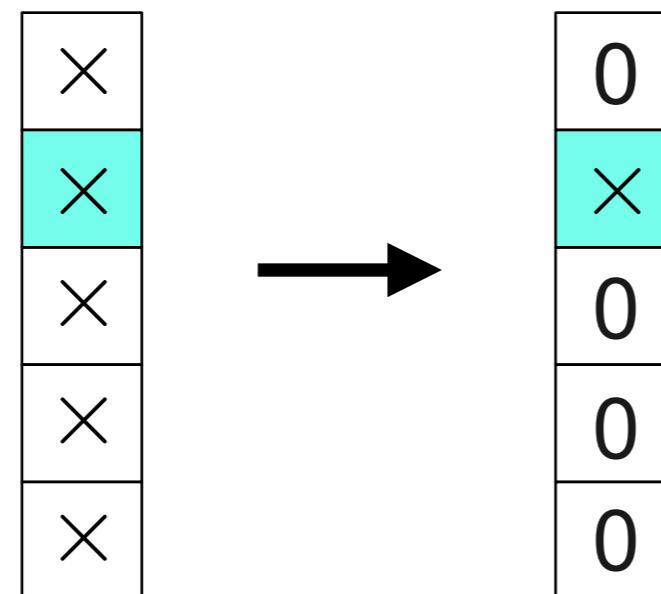
with stochastic gradients

$$g_m^k = h_m^k + C_m^k (\nabla F_m(x^k) - h_m^k) \approx \nabla F_m(x^k)$$

which are unbiased: $\mathbb{E}[g_m^k] = \nabla F_m(x^k)$

Examples of unbiased compressors

- rand- s : s elements out of d chosen unif. at random and scaled by $\frac{d}{s}$, other ones set to 0.



Examples of unbiased compressors

- rand- s : s elements out of d chosen uniformly at random and scaled by $\frac{d}{s}$, other ones set to 0.
- quantization of the real values:
Example: 0.2 represented by
$$\begin{cases} 0 \text{ with probability } \frac{4}{5} \\ 1 \text{ with probability } \frac{1}{5} \end{cases}$$

Examples of unbiased compressors

- rand- s : s elements out of d chosen uniformly at random and scaled by $\frac{d}{s}$, other ones set to 0.

- quantization of the real values:

Example: 0.2 represented by

$$\begin{cases} 0 \text{ with probability } \frac{4}{5} \\ 1 \text{ with probability } \frac{1}{5} \end{cases}$$

Albasyoni, Safaryan,
Condat, Richtárik “Optimal
Gradient Compression for
Distributed and Federated
Learning,” 2020

Variance of unbiased compressors

$\exists \omega \geq 0$ such that for any $v \in \mathbb{R}^d$,

- $\mathbb{E}[\mathcal{C}_m^k(v)] = v$
- $\mathbb{E}[||\mathcal{C}_m^k(v) - v||^2] \leq \omega \|v\|^2$

Variance of unbiased compressors

$\exists \omega \geq 0$ such that for any $v \in \mathbb{R}^d$,

- $\mathbb{E}[\mathcal{C}_m^k(v)] = v$
- $\mathbb{E}[||\mathcal{C}_m^k(v) - v||^2] \leq \omega \|v\|^2$

We also define $\omega_{av} \in [0, \omega]$ and $\zeta \in [0, \omega_{av}]$ such that, for any $(v_m)_{m=1}^M$,

$$\mathbb{E} \left[\left\| \frac{1}{M} \sum_{m=1}^M (\mathcal{C}_m^k(v_m) - v_m) \right\|^2 \right] \leq \frac{\omega_{av}}{M} \sum_{m=1}^M \|v_m\|^2 - \zeta \left\| \frac{1}{M} \sum_{m=1}^M v_m \right\|^2$$



If the $(\mathcal{C}_m^k)_{m=1}^M$ are independent, $\omega_{av} = \frac{\omega}{M}$

DIANA: convergence

Theorem [MURANA] In DIANA, set $\lambda := \frac{1}{1+\omega}$ and suppose that

$$0 < \gamma < \frac{2}{L} \frac{1}{a + 4\omega_{av}},$$

where $a := \max(1 - 2\zeta, 0)$. Choose $b > 1$ s.t. $\eta := 1 - \gamma \left(\frac{2}{L} \frac{1}{a + (1+b)^2 \omega_{av}} \right)^{-1} \in (0, 1)$. Define the Lyapunov function, for every $k \geq 0$,

$$\Psi^k := \|x^k - x^*\|^2 + (b^2 + b)\gamma^2 \omega_{av}(1 + \omega) \frac{1}{M} \sum_{m=1}^M \|h_m^k - h_m^*\|^2.$$

Then, for every $k \geq 0$, we have $\mathbb{E}[\Psi^k] \leq c^k \Psi^0$, where

$$c := 1 - \min \left\{ 2\gamma\eta\mu, \frac{1 - b^{-2}}{1 + \omega} \right\} < 1.$$

DIANA: convergence

Theorem [MURANA] In DIANA, set $\lambda := \frac{1}{1+\omega}$ and suppose that

$$0 < \gamma < \frac{2}{L} \frac{1}{a + 4\omega_{av}},$$



DIANA achieves ϵ -accuracy
with iteration complexity

$$\mathcal{O}\left(\left(\frac{L}{\mu}(1 + \omega_{av}) + \omega\right) \log(\epsilon^{-1})\right)$$

DIANA: convergence

Theorem [MURANA] In DIANA, set $\lambda := \frac{1}{1+\omega}$ and suppose that

$$0 < \gamma < \frac{2}{L} \frac{1}{a + 4\omega_{av}},$$



Typically, the communication complexity can be reduced from

$$\mathcal{O}\left(d \frac{L}{\mu} \log(\epsilon^{-1})\right) \text{ to } \mathcal{O}\left(\left(\frac{L}{\mu} + d\right) \log(\epsilon^{-1})\right)$$

Conclusion

2 ideas to reduce communication:

- 1) use **local computations**: communicate less frequently.
- 2) use **compression**: communicate compressed vectors.

Conclusion

2 ideas to reduce communication:

- 1) use **local computations**: communicate less frequently.
- 2) use **compression**: communicate compressed vectors.

Combining the 2 ideas: work in progress!

