

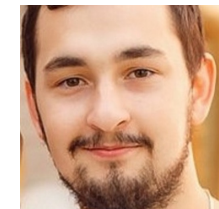
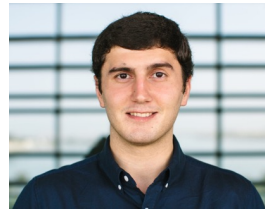
3PC: Three Point Compressors for Communication-Efficient Distributed Training

Peter Richtárik (KAUST)

ICML 2022, Baltimore, Maryland, USA

joint work with

Igor Sokolov (KAUST), Ilyas Fatkhullin (ETH), Elnur Gasanov (KAUST), Zhize Li (KAUST), Eduard Gorbunov (MIPT)



Distributed Nonconvex Optimization

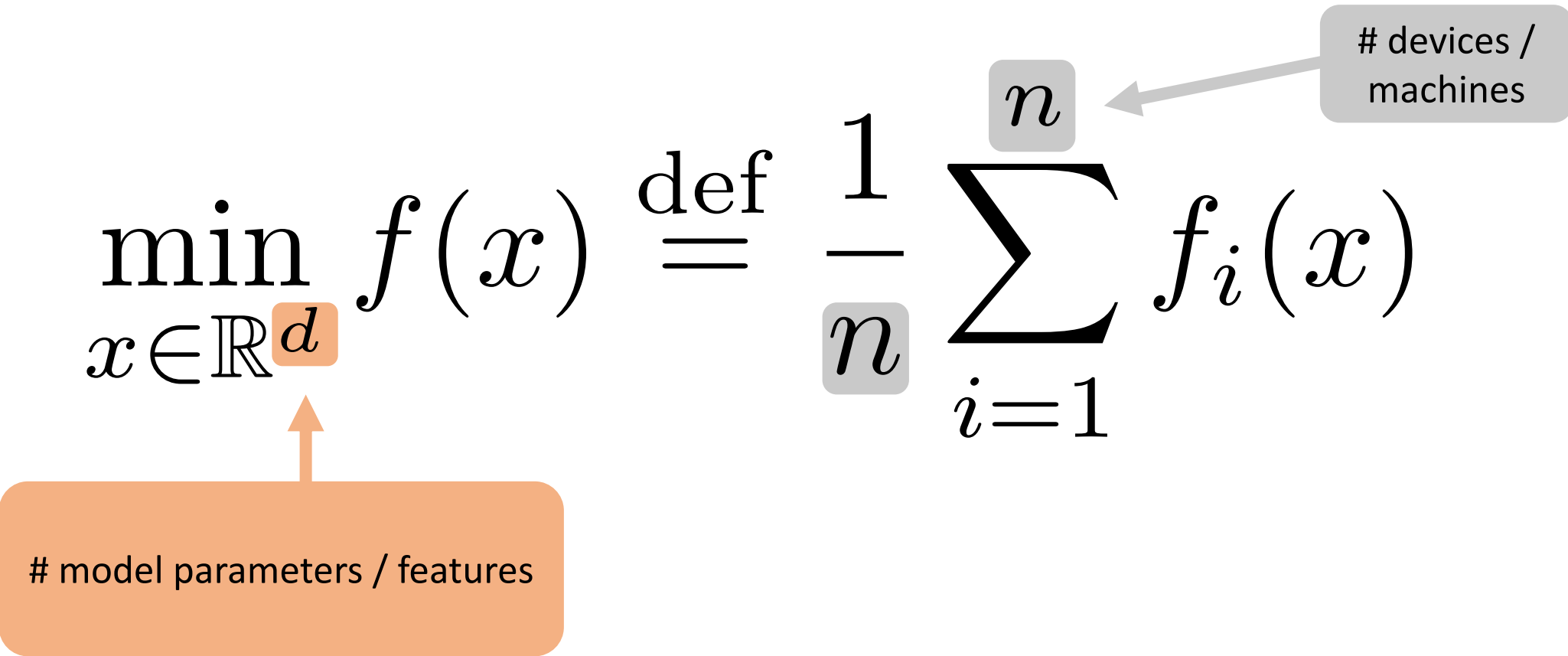
$$\min_{x \in \mathbb{R}^d} f(x) \stackrel{\text{def}}{=} \frac{1}{n} \sum_{i=1}^n f_i(x)$$

Distributed Nonconvex Optimization

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

model parameters / features

Distributed Nonconvex Optimization

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


The diagram illustrates the variables in the distributed optimization equation. The variable d in the domain $x \in \mathbb{R}^d$ is highlighted in an orange box. An orange arrow points from an orange rounded rectangle labeled "# model parameters / features" to this box. The variable n in the denominator and the summation limit is highlighted in a gray box. A gray arrow points from a gray rounded rectangle labeled "# devices / machines" to this box.

model parameters / features

devices /
machines

Distributed Nonconvex Optimization

$$\min_{x \in \mathbb{R}^d} f(x)$$

model parameters / features

$\stackrel{\text{def}}{=}$

$$\frac{1}{n} \sum_{i=1}^n$$

$$f_i(x)$$

devices /
machines

Lipschitz gradient:


$$\|\nabla f_i(x) - \nabla f_i(y)\| \leq L_i \|x - y\|, \quad \forall x, y \in \mathbb{R}^d$$

Distributed Compressed Gradient Descent

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

Distributed Compressed Gradient Descent

devices / machines


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

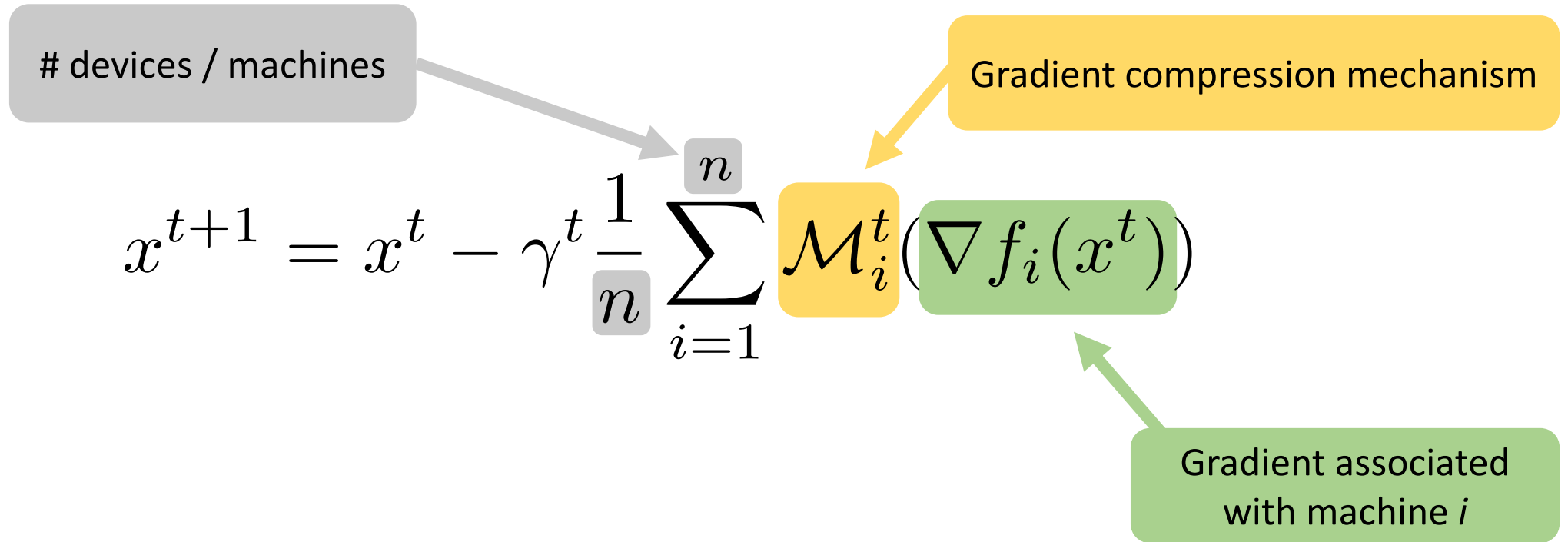
Distributed Compressed Gradient Descent

devices / machines

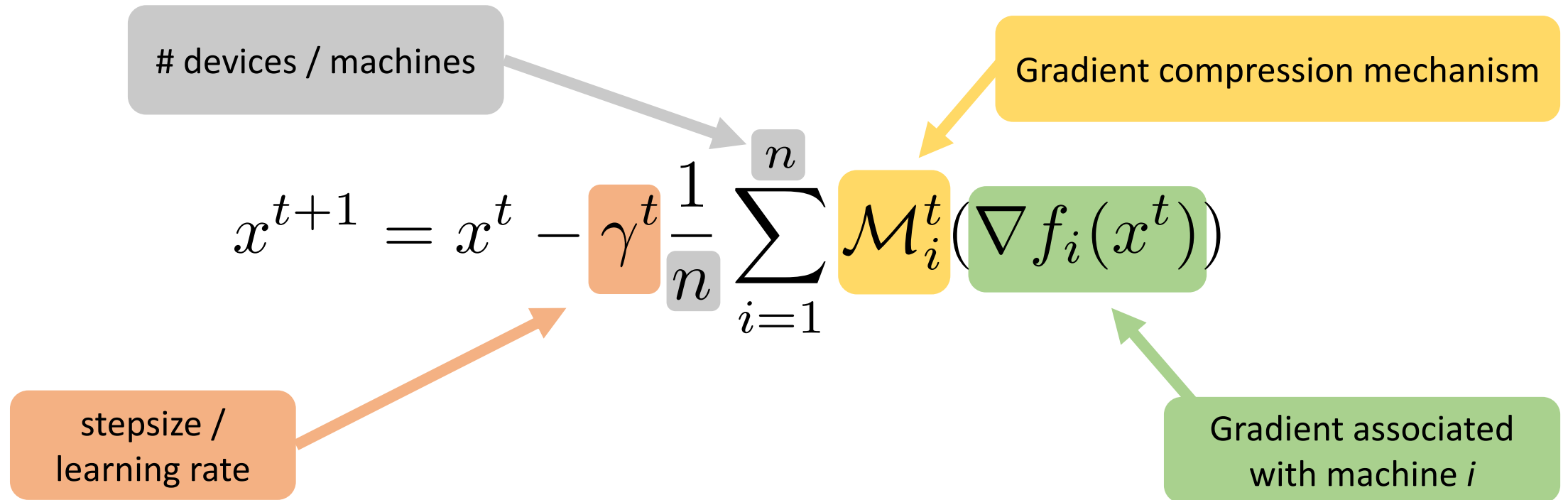
$$x^{t+1} = x^t - \gamma^t \frac{1}{n} \sum_{i=1}^n \mathcal{M}_i^t(\nabla f_i(x^t))$$

Gradient associated
with machine i

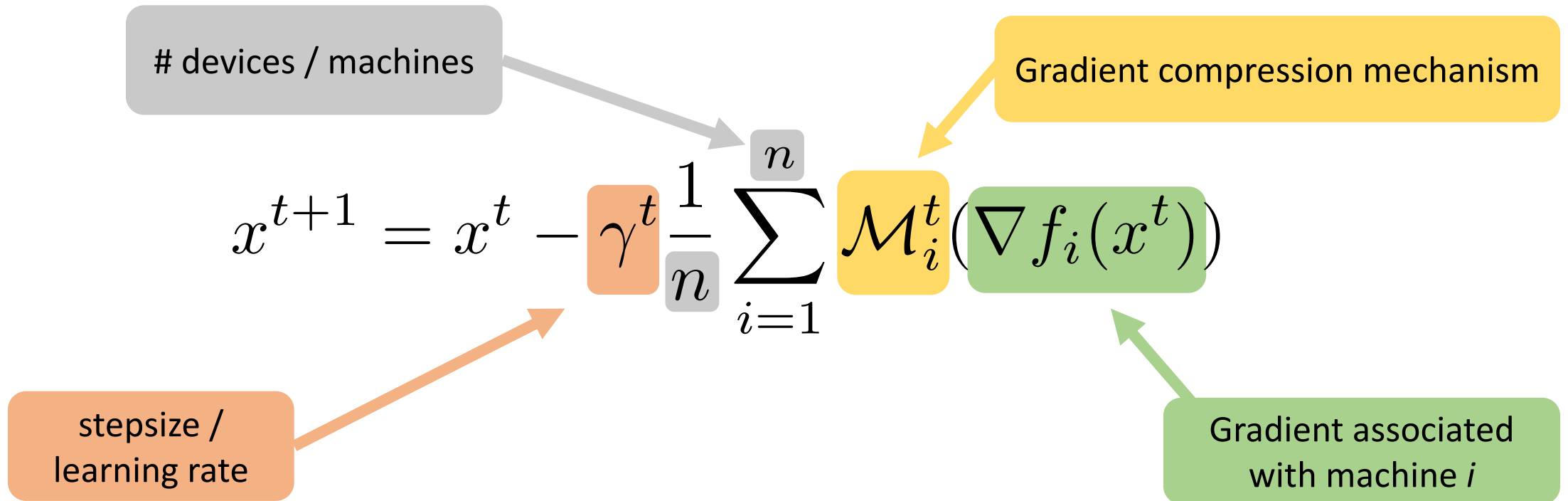
Distributed Compressed Gradient Descent



Distributed Compressed Gradient Descent



Distributed Compressed Gradient Descent



How to design a good compression mechanism?

3PC: Three Point Compressor

$$\begin{aligned} g_i^t &= \mathcal{M}_i^t \left(\nabla f_i(x^t) \right) \\ &= \mathcal{M} \left(g_i^{t-1}, \nabla f_i(x^{t-1}), \nabla f_i(x^t) \right) \\ &= \mathcal{M}_{g_i^{t-1}, \nabla f_i(x^{t-1})} \left(\nabla f_i(x^t) \right) \end{aligned}$$

Current
compressed
gradient

3PC: Three Point Compressor

$$\begin{aligned} g_i^t &= \mathcal{M}_i^t \left(\nabla f_i(x^t) \right) \\ &= \mathcal{M} \left(g_i^{t-1}, \nabla f_i(x^{t-1}), \nabla f_i(x^t) \right) \\ &= \mathcal{M}_{g_i^{t-1}, \nabla f_i(x^{t-1})} \left(\nabla f_i(x^t) \right) \end{aligned}$$

Current
compressed
gradient

3PC: Three Point Compressor

g_i^t

$$= \mathcal{M}_i^t \left(\nabla f_i(x^t) \right)$$

Current gradient

$$= \mathcal{M} \left(g_i^{t-1}, \nabla f_i(x^{t-1}), \nabla f_i(x^t) \right)$$

$$= \mathcal{M}_{g_i^{t-1}, \nabla f_i(x^{t-1})} \left(\nabla f_i(x^t) \right)$$

Current
compressed
gradient

3PC: Three Point Compressor

g_i^t

=

$$\mathcal{M}_i^t \left(\nabla f_i(x^t) \right)$$

Current gradient

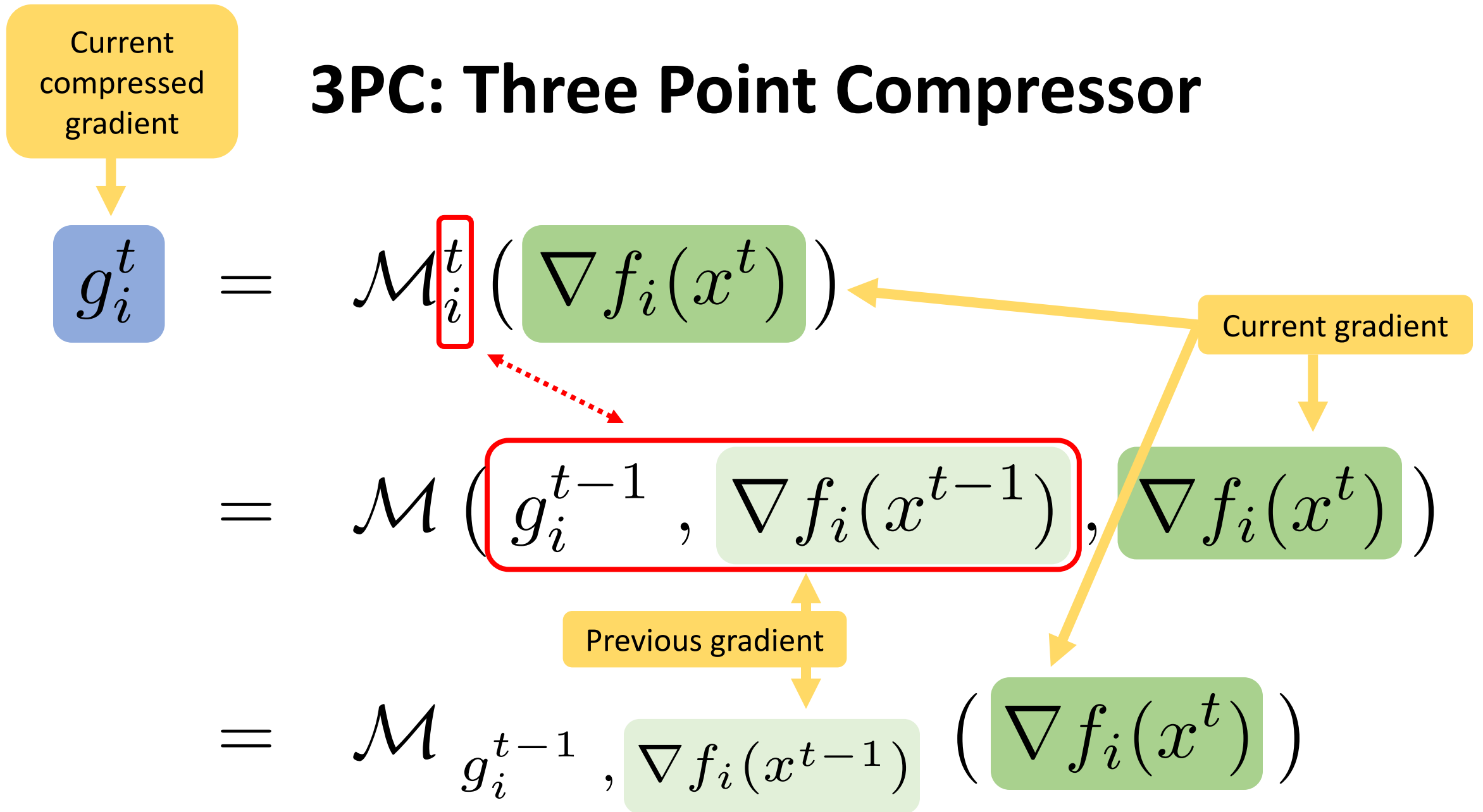
=

$$\mathcal{M} \left(g_i^{t-1}, \nabla f_i(x^{t-1}), \nabla f_i(x^t) \right)$$

=

$$\mathcal{M}_{g_i^{t-1}, \nabla f_i(x^{t-1})} \left(\nabla f_i(x^t) \right)$$

3PC: Three Point Compressor



3PC: Three Point Compressor

Current
compressed
gradient

$$g_i^t$$

=

$$\mathcal{M}_i^t \left(\nabla f_i(x^t) \right)$$

Current gradient

=

$$\mathcal{M} \left(g_i^{t-1}, \nabla f_i(x^{t-1}), \nabla f_i(x^t) \right)$$

Previous compressed gradient

Previous gradient

=

$$\mathcal{M} \left(g_i^{t-1}, \nabla f_i(x^{t-1}), \nabla f_i(x^t) \right)$$

3PC: Three Point Compressor

Current
compressed
gradient

$$g_i^t$$

=

$$\mathcal{M}_i^t \left(\nabla f_i(x^t) \right)$$

Current gradient

=

$$\mathcal{M} \left(g_i^{t-1}, \nabla f_i(x^{t-1}), \nabla f_i(x^t) \right)$$

Previous compressed gradient

Previous gradient

=

$$\mathcal{M} \left(g_i^{t-1}, \nabla f_i(x^{t-1}), \nabla f_i(x^t) \right)$$

t $t - 1$ Exact
gradientCompressed
gradient

3PC: Three Point Compressor

Current
compressed
gradient

$$g_i^t$$

=

$$\mathcal{M}_i^t \left(\nabla f_i(x^t) \right)$$

Current gradient

=

$$\mathcal{M} \left(g_i^{t-1}, \nabla f_i(x^{t-1}), \nabla f_i(x^t) \right)$$

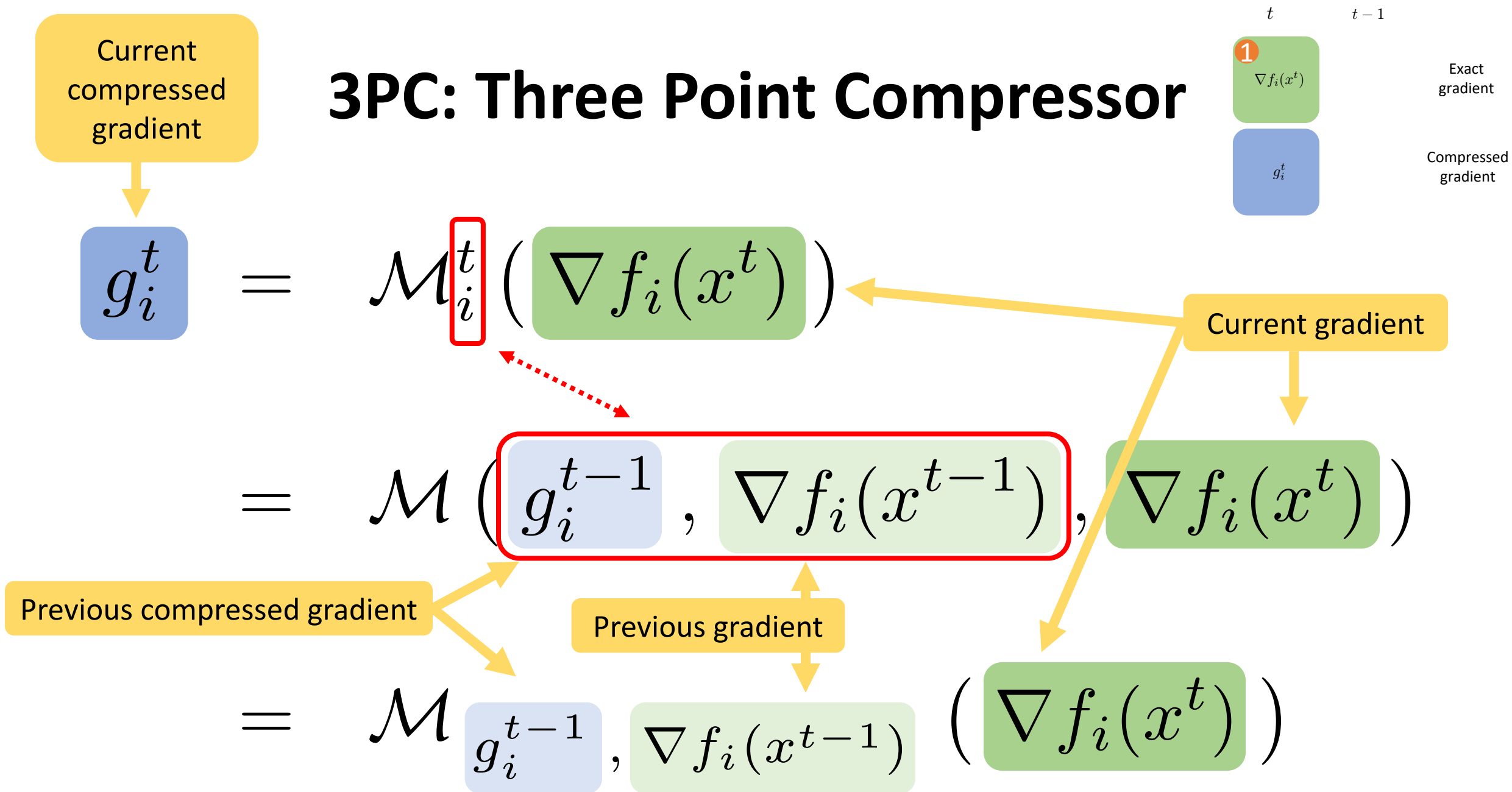
Previous compressed gradient

Previous gradient

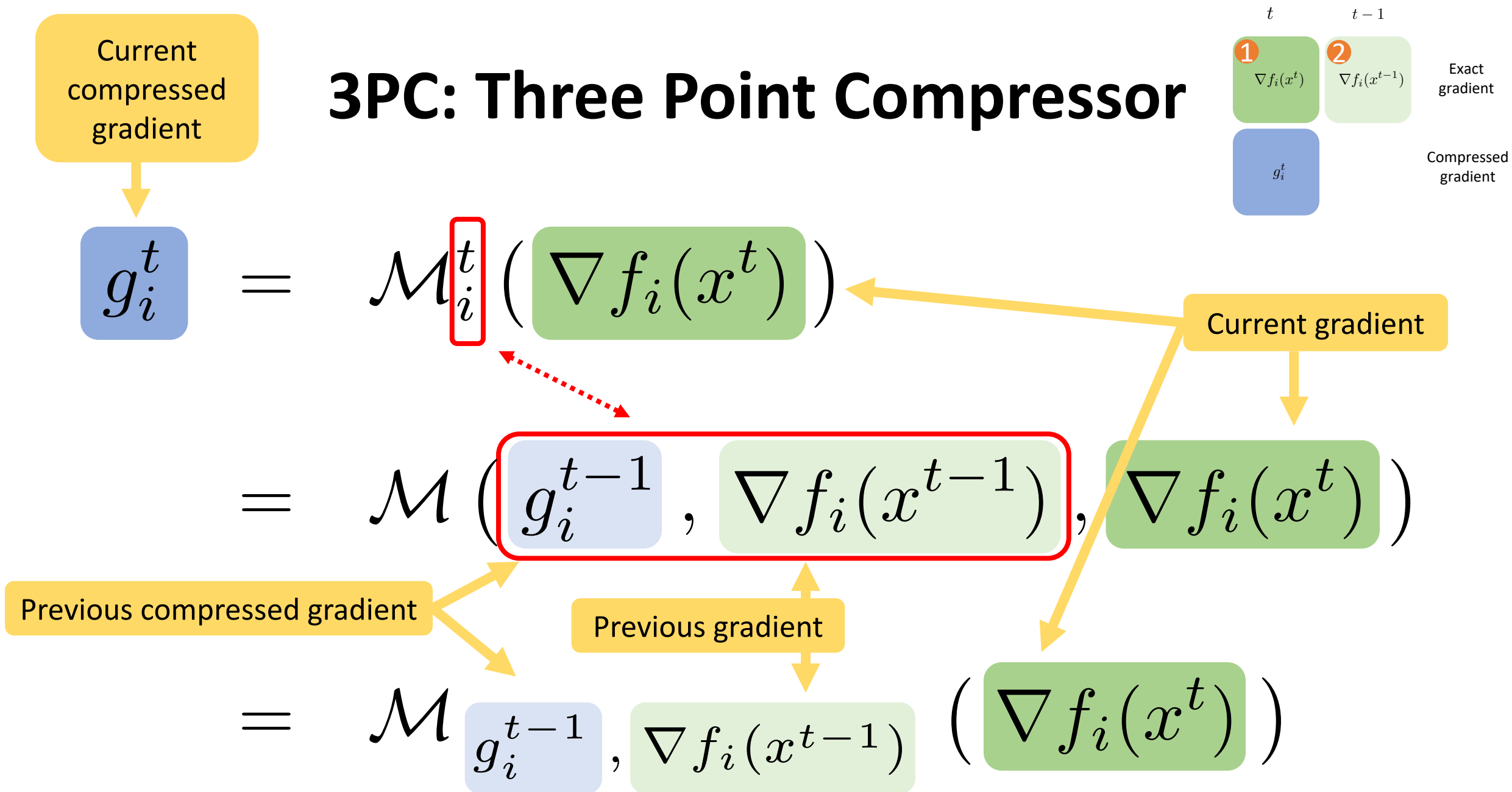
=

$$\mathcal{M} \left(g_i^{t-1}, \nabla f_i(x^{t-1}), \nabla f_i(x^t) \right)$$

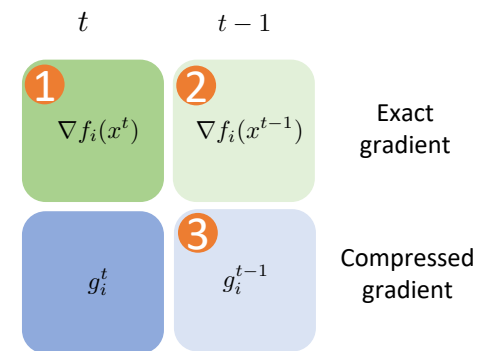
3PC: Three Point Compressor



3PC: Three Point Compressor



3PC: Three Point Compressor



Current compressed gradient

$$g_i^t$$

$$= \mathcal{M}_i^t \left(\nabla f_i(x^t) \right)$$

Current gradient

$$= \mathcal{M} \left(g_i^{t-1}, \nabla f_i(x^{t-1}), \nabla f_i(x^t) \right)$$

Previous compressed gradient

Previous gradient

$$= \mathcal{M} \left(g_i^{t-1}, \nabla f_i(x^{t-1}), \nabla f_i(x^t) \right)$$

3PC Inequality

3PC Inequality

Inequality characterizing a Contractive Compressor:

3PC Inequality

Inequality characterizing a Contractive Compressor:

$$\mathbb{E} [\|g_i^t - \nabla f_i(x^t)\|^2] \leq (1 - \lambda) \|\nabla f_i(x^t)\|^2$$

3PC Inequality

Inequality characterizing a Contractive Compressor:

$$\mathbb{E} [\|g_i^t - \nabla f_i(x^t)\|^2] \leq (1 - \lambda) \|\nabla f_i(x^t)\|^2$$

3PC Inequality

Inequality characterizing a Contractive Compressor:

$$\mathbb{E} [\|g_i^t - \nabla f_i(x^t)\|^2] \leq (1 - \lambda) \|\nabla f_i(x^t)\|^2$$

3PC Inequality

Inequality characterizing a Contractive Compressor:

$$\mathbb{E} [\|g_i^t - \nabla f_i(x^t)\|^2] \leq (1 - \lambda) \|\nabla f_i(x^t)\|^2$$

3PC Inequality

Inequality characterizing a Contractive Compressor:

$$\mathbb{E} [\|g_i^t - \nabla f_i(x^t)\|^2] \leq (1 - \lambda) \|\nabla f_i(x^t)\|^2$$

$$\mathbb{E} [\|\mathcal{C}(x) - x\|^2] \leq (1 - \lambda) \|x\|^2 \quad \forall x \in \mathbb{R}^d$$

3PC Inequality

Inequality characterizing a 3PC:

Inequality characterizing a Contractive Compressor:

$$\mathbb{E} [\|g_i^t - \nabla f_i(x^t)\|^2] \leq (1 - \lambda) \|\nabla f_i(x^t)\|^2$$

$$\mathbb{E} [\|\mathcal{C}(x) - x\|^2] \leq (1 - \lambda) \|x\|^2 \quad \forall x \in \mathbb{R}^d$$

3PC Inequality

Inequality characterizing a 3PC:

$$\mathbb{E} [\|g_i^t - \nabla f_i(x^t)\|^2] \leq (1 - A)\|g_i^{t-1} - \nabla f_i(x^{t-1})\|^2 + B\|\nabla f_i(x^t) - \nabla f_i(x^{t-1})\|^2$$

Inequality characterizing a Contractive Compressor:

$$\mathbb{E} [\|g_i^t - \nabla f_i(x^t)\|^2] \leq (1 - \lambda)\|\nabla f_i(x^t)\|^2$$

$$\mathbb{E} [\|\mathcal{C}(x) - x\|^2] \leq (1 - \lambda)\|x\|^2 \quad \forall x \in \mathbb{R}^d$$

3PC Inequality

Inequality characterizing a 3PC:

$$\mathbb{E} [\|g_i^t - \overset{\textcircled{1}}{\nabla f_i(x^t)}\|^2] \leq (1 - A) \|g_i^{t-1} - \nabla f_i(x^{t-1})\|^2 + B \|\overset{\textcircled{1}}{\nabla f_i(x^t)} - \nabla f_i(x^{t-1})\|^2$$

Inequality characterizing a Contractive Compressor:

$$\mathbb{E} [\|\overset{\textcircled{1}}{g_i^t} - \nabla f_i(x^t)\|^2] \leq (1 - \lambda) \|\nabla f_i(x^t)\|^2 \qquad \mathbb{E} [\|\mathcal{C}(x) - \overset{\textcircled{1}}{x}\|^2] \leq (1 - \lambda) \|\overset{\textcircled{1}}{x}\|^2 \quad \forall \overset{\textcircled{1}}{x} \in \mathbb{R}^d$$

3PC Inequality

Inequality characterizing a 3PC:

$$\mathbb{E} [\|g_i^t - \overset{\textcircled{1}}{\nabla f_i(x^t)}\|^2] \leq (1 - A) \|g_i^{t-1} - \overset{\textcircled{2}}{\nabla f_i(x^{t-1})}\|^2 + B \|\overset{\textcircled{1}}{\nabla f_i(x^t)} - \overset{\textcircled{2}}{\nabla f_i(x^{t-1})}\|^2$$

Inequality characterizing a Contractive Compressor:

$$\mathbb{E} [\|\overset{\text{blue}}{g_i^t} - \nabla f_i(x^t)\|^2] \leq (1 - \overset{\text{yellow}}{\lambda}) \|\nabla f_i(x^t)\|^2 \qquad \mathbb{E} [\|\overset{\text{blue}}{\mathcal{C}(x)} - \overset{\text{green}}{x}\|^2] \leq (1 - \overset{\text{yellow}}{\lambda}) \|\overset{\text{green}}{x}\|^2 \qquad \forall \overset{\text{green}}{x} \in \mathbb{R}^d$$

3PC Inequality

Inequality characterizing a 3PC:

$$\mathbb{E} [\|g_i^t - \overset{\textcircled{1}}{\nabla f_i(x^t)}\|^2] \leq (1 - A) \|\overset{\textcircled{3}}{g_i^{t-1}} - \overset{\textcircled{2}}{\nabla f_i(x^{t-1})}\|^2 + B \|\overset{\textcircled{1}}{\nabla f_i(x^t)} - \overset{\textcircled{2}}{\nabla f_i(x^{t-1})}\|^2$$

Inequality characterizing a Contractive Compressor:

$$\mathbb{E} [\|\overset{\textcircled{3}}{g_i^t} - \nabla f_i(x^t)\|^2] \leq (1 - \lambda) \|\nabla f_i(x^t)\|^2 \qquad \mathbb{E} [\|\mathcal{C}(x) - \overset{\textcircled{1}}{x}\|^2] \leq (1 - \lambda) \|\overset{\textcircled{1}}{x}\|^2 \qquad \forall \overset{\textcircled{1}}{x} \in \mathbb{R}^d$$

3PC Inequality

Inequality characterizing a 3PC:

$$\mathbb{E} [\| \overset{\textcircled{1}}{g_i^t} - \overset{\textcircled{1}}{\nabla f_i(x^t)} \|^2] \leq (1 - A) \| \overset{\textcircled{3}}{g_i^{t-1}} - \overset{\textcircled{2}}{\nabla f_i(x^{t-1})} \|^2 + B \| \overset{\textcircled{1}}{\nabla f_i(x^t)} - \overset{\textcircled{2}}{\nabla f_i(x^{t-1})} \|^2$$

Inequality characterizing a Contractive Compressor:

$$\mathbb{E} [\| \overset{\textcircled{1}}{g_i^t} - \overset{\textcircled{1}}{\nabla f_i(x^t)} \|^2] \leq (1 - \overset{\textcircled{1}}{\lambda}) \| \overset{\textcircled{1}}{\nabla f_i(x^t)} \|^2 \qquad \mathbb{E} [\| \overset{\textcircled{1}}{\mathcal{C}(x)} - \overset{\textcircled{1}}{x} \|^2] \leq (1 - \overset{\textcircled{1}}{\lambda}) \| \overset{\textcircled{1}}{x} \|^2 \quad \forall \overset{\textcircled{1}}{x} \in \mathbb{R}^d$$

3PC Inequality

Inequality characterizing a 3PC:

A, B : Parameters characterizing a 3PC

$$\mathbb{E} [\| \overset{\textcircled{1}}{g_i^t} - \overset{\textcircled{1}}{\nabla f_i(x^t)} \|^2] \leq (1 - \overset{\textcircled{3}}{A}) \overset{\textcircled{2}}{\| g_i^{t-1} - \nabla f_i(x^{t-1}) \|^2} + \overset{\textcircled{1}}{B} \overset{\textcircled{2}}{\| \nabla f_i(x^t) - \nabla f_i(x^{t-1}) \|^2}$$

Inequality characterizing a Contractive Compressor:

$$\mathbb{E} [\| \overset{\textcircled{1}}{g_i^t} - \overset{\textcircled{1}}{\nabla f_i(x^t)} \|^2] \leq (1 - \overset{\textcircled{2}}{\lambda}) \overset{\textcircled{2}}{\| \nabla f_i(x^t) \|^2}$$

$$\mathbb{E} [\| \overset{\textcircled{1}}{\mathcal{C}(x)} - \overset{\textcircled{1}}{x} \|^2] \leq (1 - \overset{\textcircled{2}}{\lambda}) \overset{\textcircled{2}}{\| x \|^2} \quad \forall \overset{\textcircled{1}}{x} \in \mathbb{R}^d$$

3PC Inequality

Inequality characterizing a 3PC:

A, B : Parameters characterizing a 3PC

$$\mathbb{E} [\underbrace{\|g_i^t - \nabla f_i(x^t)\|^2}_{\text{Compression error at iteration } t}] \leq (1 - A) \|g_i^{t-1} - \nabla f_i(x^{t-1})\|^2 + B \|\nabla f_i(x^t) - \nabla f_i(x^{t-1})\|^2$$

Diagram annotations: A yellow box labeled "A, B : Parameters characterizing a 3PC" has arrows pointing to the yellow boxes containing A and B in the inequality. Orange numbers 1, 2, and 3 are placed above the terms g_i^t , $\nabla f_i(x^t)$, and g_i^{t-1} respectively.

Inequality characterizing a Contractive Compressor:

$$\mathbb{E} [\|g_i^t - \nabla f_i(x^t)\|^2] \leq (1 - \lambda) \|\nabla f_i(x^t)\|^2$$

$$\mathbb{E} [\|\mathcal{C}(x) - x\|^2] \leq (1 - \lambda) \|x\|^2 \quad \forall x \in \mathbb{R}^d$$

3PC Inequality

Inequality characterizing a 3PC:

A, B : Parameters characterizing a 3PC

$$\mathbb{E} [\underbrace{\|g_i^t - \nabla f_i(x^t)\|^2}_{\text{Compression error at iteration } t}] \leq (1 - \underbrace{A}_{\text{Compression error at iteration } t-1}) \underbrace{\|g_i^{t-1} - \nabla f_i(x^{t-1})\|^2}_{\text{Compression error at iteration } t-1} + B \underbrace{\|\nabla f_i(x^t) - \nabla f_i(x^{t-1})\|^2}_{\text{Compression error at iteration } t-1}$$

Diagram illustrating the 3PC inequality. The inequality relates the expected squared norm of the compression error at iteration t to the squared norm of the compression error at iteration $t-1$ and the squared norm of the difference between the gradients at iterations t and $t-1$. The parameters A and B characterize the 3PC. The terms are labeled with numbers 1, 2, and 3, and the compression errors are labeled with red brackets.

Inequality characterizing a Contractive Compressor:

$$\mathbb{E} [\|g_i^t - \nabla f_i(x^t)\|^2] \leq (1 - \lambda) \|\nabla f_i(x^t)\|^2$$

$$\mathbb{E} [\|\mathcal{C}(x) - x\|^2] \leq (1 - \lambda) \|x\|^2 \quad \forall x \in \mathbb{R}^d$$

3PC Inequality

Inequality characterizing a 3PC:

A, B : Parameters characterizing a 3PC

$$\mathbb{E} [\| \overset{\textcircled{1}}{g_i^t} - \overset{\textcircled{2}}{\nabla f_i(x^t)} \|^2] \leq (1 - \overset{\textcircled{3}}{A}) \| \overset{\textcircled{3}}{g_i^{t-1}} - \overset{\textcircled{2}}{\nabla f_i(x^{t-1})} \|^2 + \overset{\textcircled{1}}{B} \| \overset{\textcircled{1}}{\nabla f_i(x^t)} - \overset{\textcircled{2}}{\nabla f_i(x^{t-1})} \|^2$$

Compression error
at iteration t

Compression error
at iteration $t - 1$

$$\mathbb{E} [\| \mathcal{M}_{h,y}(x) - x \|^2] \leq (1 - A) \| h - y \|^2 + B \| x - y \|^2$$

$\forall h, y, x \in \mathbb{R}^d$

Inequality characterizing a Contractive Compressor:

$$\mathbb{E} [\| g_i^t - \nabla f_i(x^t) \|^2] \leq (1 - \lambda) \| \nabla f_i(x^t) \|^2$$

$$\mathbb{E} [\| \mathcal{C}(x) - x \|^2] \leq (1 - \lambda) \| x \|^2 \quad \forall x \in \mathbb{R}^d$$

$$g_i^t \equiv \mathcal{M}_i^t(\nabla f_i(x^t))$$

Examples of Compression Mechanisms

$$g_i^t \equiv \mathcal{M}_i^t(\nabla f_i(x^t))$$

Examples of Compression Mechanisms

$g_i^t \equiv \nabla f_i(x^t)$	Gradient Descent (GD)

$$g_i^t \equiv \mathcal{M}_i^t(\nabla f_i(x^t))$$

Examples of Compression Mechanisms

$g_i^t \equiv \nabla f_i(x^t)$	+ fast $O(1/t)$ convergence - communicates a lot	Gradient Descent (GD)

$$g_i^t \equiv \mathcal{M}_i^t(\nabla f_i(x^t))$$

Examples of Compression Mechanisms

$$g_i^t \equiv \nabla f_i(x^t)$$

+ fast $O(1/t)$ convergence
- communicates a lot

Gradient Descent (GD)

$$g_i^t \equiv \mathcal{C}(\nabla f_i(x^t))$$

**Compressed
Gradient Descent (CGD)**

$$\mathbb{E} \|\mathcal{C}(u) - u\|^2 \leq (1 - \lambda)\|u\|^2, \quad \forall u \in \mathbb{R}^d$$

$$g_i^t \equiv \mathcal{M}_i^t(\nabla f_i(x^t))$$

Examples of Compression Mechanisms

$g_i^t \equiv \nabla f_i(x^t)$	<p>+ fast $O(1/t)$ convergence - communicates a lot</p>	<p>Gradient Descent (GD)</p>
$g_i^t \equiv \mathcal{C}(\nabla f_i(x^t))$	<p>- diverges! + communicates little</p>	<p>Compressed Gradient Descent (CGD)</p>
<p>$\mathbb{E} \ \mathcal{C}(u) - u\ ^2 \leq (1 - \lambda)\ u\ ^2, \quad \forall u \in \mathbb{R}^d$</p>		

$$g_i^t \equiv \mathcal{M}_i^t(\nabla f_i(x^t))$$

Examples of Compression Mechanisms

$g_i^t \equiv \nabla f_i(x^t)$	<p>+ fast $O(1/t)$ convergence - communicates a lot</p>	<p>Gradient Descent (GD)</p>
$g_i^t \equiv \mathcal{C}(\nabla f_i(x^t))$	<p>- diverges! + communicates little</p>	<p>Compressed Gradient Descent (CGD)</p>
$g_i^t \equiv g_i^{t-1} + \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1})$	<p>$\mathbb{E} \ \mathcal{C}(u) - u\ ^2 \leq (1 - \lambda)\ u\ ^2, \quad \forall u \in \mathbb{R}^d$</p>	<p>Error Feedback 2021 (EF21) [R., Sokolov, Fatkhullin @ NeurIPS 2021]</p>

$$g_i^t \equiv \mathcal{M}_i^t(\nabla f_i(x^t))$$

Examples of Compression Mechanisms

$g_i^t \equiv \nabla f_i(x^t)$	<p>+ fast $O(1/t)$ convergence - communicates a lot</p>	<p>Gradient Descent (GD)</p>
$g_i^t \equiv \mathcal{C}(\nabla f_i(x^t))$	<p>- diverges! + communicates little</p>	<p>Compressed Gradient Descent (CGD)</p>
$g_i^t \equiv g_i^{t-1} + \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1})$	<p>+ fast $O(1/t)$ convergence + communicates little</p>	<p>Error Feedback 2021 (EF21) [R., Sokolov, Fatkhullin @ NeurIPS 2021]</p>
<p>$\mathbb{E} \ \mathcal{C}(u) - u\ ^2 \leq (1 - \lambda)\ u\ ^2, \quad \forall u \in \mathbb{R}^d$</p>		

$$g_i^t \equiv \mathcal{M}_i^t(\nabla f_i(x^t))$$

Examples of Compression Mechanisms

$g_i^t \equiv \nabla f_i(x^t)$	<p>+ fast $O(1/t)$ convergence - communicates a lot</p>	<p>Gradient Descent (GD)</p>
$g_i^t \equiv \mathcal{C}(\nabla f_i(x^t))$	<p>- diverges! + communicates little</p>	<p>Compressed Gradient Descent (CGD)</p>
$g_i^t \equiv g_i^{t-1} + \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1})$	<p>+ fast $O(1/t)$ convergence + communicates little</p>	<p>Error Feedback 2021 (EF21) [R., Sokolov, Fatkhullin @ NeurIPS 2021]</p>
$g_i^t \equiv \begin{cases} \nabla f_i(x^t) & \text{if } \ \nabla f_i(x^t) - g_i^{t-1}\ ^2 > \zeta \ \nabla f_i(x^t) - \nabla f_i(x^{t-1})\ ^2 \\ g_i^{t-1} & \text{otherwise} \end{cases}$		<p>Lazily Aggregated Gradient (LAG) [Chen, Giannakis, Sun, Yin @ NeurIPS 2018]</p>

$$g_i^t \equiv \mathcal{M}_i^t(\nabla f_i(x^t))$$

Examples of Compression Mechanisms

$g_i^t \equiv \nabla f_i(x^t)$	<p>+ fast $O(1/t)$ convergence - communicates a lot</p>	<p>Gradient Descent (GD)</p>
$g_i^t \equiv \mathcal{C}(\nabla f_i(x^t))$	<p>- diverges! + communicates little</p>	<p>Compressed Gradient Descent (CGD)</p>
$g_i^t \equiv g_i^{t-1} + \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1})$	<p>+ fast $O(1/t)$ convergence + communicates little</p>	<p>Error Feedback 2021 (EF21) [R., Sokolov, Fatkhullin @ NeurIPS 2021]</p>
$g_i^t \equiv \begin{cases} \nabla f_i(x^t) & \text{if } \ \nabla f_i(x^t) - g_i^{t-1}\ ^2 > \zeta \ \nabla f_i(x^t) - \nabla f_i(x^{t-1})\ ^2 \\ g_i^{t-1} & \text{otherwise} \end{cases}$ <p>Communicates when a trigger is fired!</p>		<p>Lazily Aggregated Gradient (LAG) [Chen, Giannakis, Sun, Yin @ NeurIPS 2018]</p>

$$g_i^t \equiv \mathcal{M}_i^t(\nabla f_i(x^t))$$

Examples of Compression Mechanisms

$g_i^t \equiv \nabla f_i(x^t)$	<p>+ fast $O(1/t)$ convergence - communicates a lot</p>	<p>Gradient Descent (GD)</p>
$g_i^t \equiv \mathcal{C}(\nabla f_i(x^t))$	<p>- diverges! + communicates little</p>	<p>Compressed Gradient Descent (CGD)</p>
$g_i^t \equiv g_i^{t-1} + \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1})$	<p>+ fast $O(1/t)$ convergence + communicates little</p>	<p>Error Feedback 2021 (EF21) [R., Sokolov, Fatkhullin @ NeurIPS 2021]</p>
$g_i^t \equiv \begin{cases} \nabla f_i(x^t) & \text{if } \ \nabla f_i(x^t) - g_i^{t-1}\ ^2 > \zeta \ \nabla f_i(x^t) - \nabla f_i(x^{t-1})\ ^2 \\ g_i^{t-1} & \text{otherwise} \end{cases}$ <p>Communicates when a trigger is fired!</p>	<p>- convergence not understood + communication savings unclear</p>	<p>Lazily Aggregated Gradient (LAG) [Chen, Giannakis, Sun, Yin @ NeurIPS 2018]</p>

$$g_i^t \equiv \mathcal{M}_i^t(\nabla f_i(x^t))$$

Examples of Compression Mechanisms

$g_i^t \equiv \nabla f_i(x^t)$ <p>+ fast $O(1/t)$ convergence - communicates a lot</p>	Gradient Descent (GD)
$g_i^t \equiv \mathcal{C}(\nabla f_i(x^t))$ <p>- diverges! + communicates little</p>	Compressed Gradient Descent (CGD)
$g_i^t \equiv g_i^{t-1} + \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1})$ <p>+ fast $O(1/t)$ convergence + communicates little</p>	Error Feedback 2021 (EF21) [R., Sokolov, Fatkhullin @ NeurIPS 2021]
$g_i^t \equiv \begin{cases} \nabla f_i(x^t) & \text{if } \ \nabla f_i(x^t) - g_i^{t-1}\ ^2 > \zeta \ \nabla f_i(x^t) - \nabla f_i(x^{t-1})\ ^2 \\ g_i^{t-1} & \text{otherwise} \end{cases}$ <p>Communicates when a trigger is fired!</p> <p>- convergence not understood + communication savings unclear</p>	Lazily Aggregated Gradient (LAG) [Chen, Giannakis, Sun, Yin @ NeurIPS 2018]
$g_i^t \equiv \begin{cases} g_i^{t-1} + \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1}) & \text{if } \ \nabla f_i(x^t) - g_i^{t-1}\ ^2 > \zeta \ \nabla f_i(x^t) - \nabla f_i(x^{t-1})\ ^2 \\ g_i^{t-1} & \text{otherwise} \end{cases}$	CLAG = EF21 + LAG (Lazily Aggregated EF21) [NEW @ ICML 2022]

$$g_i^t \equiv \mathcal{M}_i^t(\nabla f_i(x^t))$$

Examples of Compression Mechanisms

$g_i^t \equiv \nabla f_i(x^t)$ <p>+ fast $O(1/t)$ convergence - communicates a lot</p>	Gradient Descent (GD)
$g_i^t \equiv \mathcal{C}(\nabla f_i(x^t))$ <p>- diverges! + communicates little</p>	Compressed Gradient Descent (CGD)
$g_i^t \equiv g_i^{t-1} + \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1})$ <p>+ fast $O(1/t)$ convergence + communicates little</p>	Error Feedback 2021 (EF21) [R., Sokolov, Fatkhullin @ NeurIPS 2021]
$g_i^t \equiv \begin{cases} \nabla f_i(x^t) & \text{if } \ \nabla f_i(x^t) - g_i^{t-1}\ ^2 > \zeta \ \nabla f_i(x^t) - \nabla f_i(x^{t-1})\ ^2 \\ g_i^{t-1} & \text{otherwise} \end{cases}$ <p>Communicates when a trigger is fired!</p> <p>- convergence not understood + communication savings unclear</p>	Lazily Aggregated Gradient (LAG) [Chen, Giannakis, Sun, Yin @ NeurIPS 2018]
$g_i^t \equiv \begin{cases} g_i^{t-1} + \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1}) & \text{if } \ \nabla f_i(x^t) - g_i^{t-1}\ ^2 > \zeta \ \nabla f_i(x^t) - \nabla f_i(x^{t-1})\ ^2 \\ g_i^{t-1} & \text{otherwise} \end{cases}$	CLAG = EF21 + LAG (Lazily Aggregated EF21) [NEW @ ICML 2022]

$$g_i^t \equiv \mathcal{M}_i^t(\nabla f_i(x^t))$$

Examples of Compression Mechanisms

$g_i^t \equiv \nabla f_i(x^t)$ <p>+ fast $O(1/t)$ convergence - communicates a lot</p>	Gradient Descent (GD)
$g_i^t \equiv \mathcal{C}(\nabla f_i(x^t))$ <p>- diverges! + communicates little</p>	Compressed Gradient Descent (CGD)
$g_i^t \equiv g_i^{t-1} + \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1})$ <p>+ fast $O(1/t)$ convergence + communicates little</p>	Error Feedback 2021 (EF21) [R., Sokolov, Fatkhullin @ NeurIPS 2021]
$g_i^t \equiv \begin{cases} \nabla f_i(x^t) & \text{if } \ \nabla f_i(x^t) - g_i^{t-1}\ ^2 > \zeta \ \nabla f_i(x^t) - \nabla f_i(x^{t-1})\ ^2 \\ g_i^{t-1} & \text{otherwise} \end{cases}$ <p>Communicates when a trigger is fired!</p> <p>- convergence not understood + communication savings unclear</p>	Lazily Aggregated Gradient (LAG) [Chen, Giannakis, Sun, Yin @ NeurIPS 2018]
$g_i^t \equiv \begin{cases} g_i^{t-1} + \mathcal{C}(\nabla f_i(x^t) - g_i^{t-1}) & \text{if } \ \nabla f_i(x^t) - g_i^{t-1}\ ^2 > \zeta \ \nabla f_i(x^t) - \nabla f_i(x^{t-1})\ ^2 \\ g_i^{t-1} & \text{otherwise} \end{cases}$ <p>+ combines benefits of EF21 and LAG</p>	CLAG = EF21 + LAG (Lazily Aggregated EF21) [NEW @ ICML 2022]

Three Point Compressors: Theory

Three Point Compressors: Theory

Theorem

If \mathcal{M}_i^t is a 3PC with parameters A and B , then the method

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

satisfies

Three Point Compressors: Theory

Theorem

If \mathcal{M}_i^t is a 3PC with parameters A and B , then the method

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

satisfies

$$\mathbb{E} \|\nabla f(\hat{x}^t)\|^2 \leq \frac{2(f(x^0) - \inf_x f(x))}{\gamma t} + \frac{\frac{1}{n} \sum_{i=1}^n \mathbb{E} \|g_i^0 - \nabla f_i(x^0)\|^2}{At}$$

Three Point Compressors: Theory

Theorem

If \mathcal{M}_i^t is a 3PC with parameters A and B , then the method

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

satisfies

$$\mathbb{E} \|\nabla f(\hat{x}^t)\|^2 \leq \frac{2(f(x^0) - \inf_x f(x))}{\gamma t} + \frac{\frac{1}{n} \sum_{i=1}^n \mathbb{E} \|g_i^0 - \nabla f_i(x^0)\|^2}{At}$$

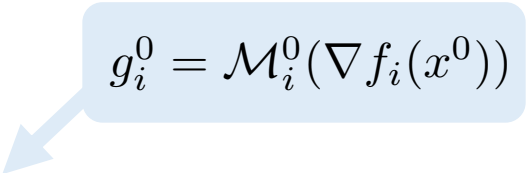
Three Point Compressors: Theory

Theorem

If \mathcal{M}_i^t is a 3PC with parameters A and B , then the method

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

satisfies

$$\mathbb{E} \|\nabla f(\hat{x}^t)\|^2 \leq \frac{2(f(x^0) - \inf_x f(x))}{\gamma t} + \frac{\frac{1}{n} \sum_{i=1}^n \mathbb{E} \|g_i^0 - \nabla f_i(x^0)\|^2}{At}$$


$g_i^0 = \mathcal{M}_i^0(\nabla f_i(x^0))$

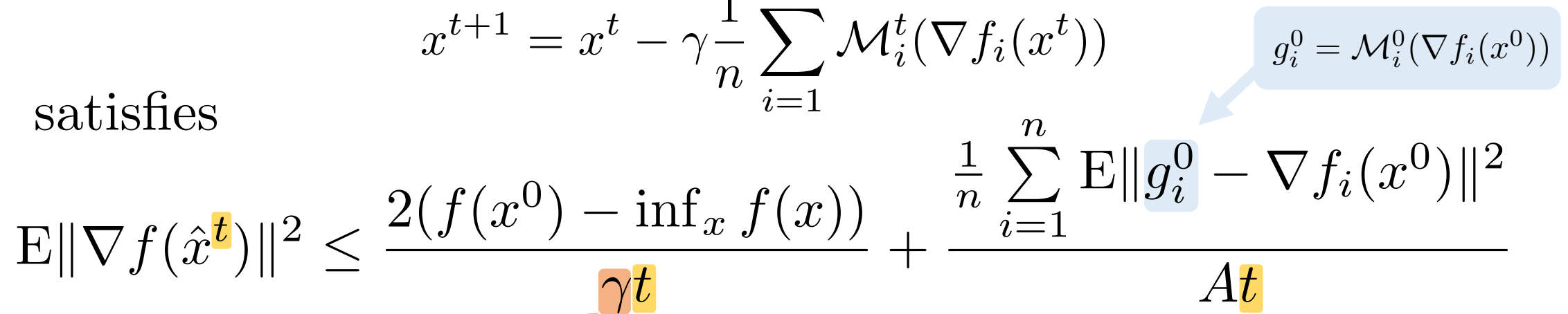
Three Point Compressors: Theory

Theorem

If \mathcal{M}_i^t is a 3PC with parameters A and B , then the method

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

satisfies

$$\mathbb{E} \|\nabla f(\hat{x}^t)\|^2 \leq \frac{2(f(x^0) - \inf_x f(x))}{\gamma t} + \frac{\frac{1}{n} \sum_{i=1}^n \mathbb{E} \|g_i^0 - \nabla f_i(x^0)\|^2}{At}$$


stepsize $0 < \gamma \leq \left(L_- + L_+ \sqrt{\frac{B}{A}}\right)^{-1}$

$$\|\nabla f(x) - \nabla f(y)\| \leq L_- \|x - y\|, \quad \forall x, y \in \mathbb{R}^d$$

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x) - \nabla f_i(y)\|^2 \leq L_+^2 \|x - y\|^2, \quad \forall x, y \in \mathbb{R}^d$$

Three Point Compressors: Theory

Theorem

If \mathcal{M}_i^t is a 3PC with parameters A and B , then the method

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

satisfies

$$\mathbb{E} \|\nabla f(\hat{x}^t)\|^2 \leq \frac{2(f(x^0) - \inf_x f(x))}{\gamma t} + \frac{\frac{1}{n} \sum_{i=1}^n \mathbb{E} \|g_i^0 - \nabla f_i(x^0)\|^2}{At}$$

$$g_i^0 = \mathcal{M}_i^0(\nabla f_i(x^0))$$

stepsize $0 < \gamma \leq \left(L_- + L_+ \sqrt{\frac{B}{A}}\right)^{-1}$

$$\|\nabla f(x) - \nabla f(y)\| \leq L_- \|x - y\|, \quad \forall x, y \in \mathbb{R}^d$$

$$\frac{1}{n} \sum_{i=1}^n \|\nabla f_i(x) - \nabla f_i(y)\|^2 \leq L_+^2 \|x - y\|^2, \quad \forall x, y \in \mathbb{R}^d$$

Corollary

$$t = \mathcal{O} \left(\frac{(f(x^0) - \inf_x f(x)) \left(L_- + L_+ \sqrt{\frac{B}{A}}\right)}{\varepsilon} \right) \Rightarrow \mathbb{E} \|\nabla f(\hat{x}^t)\|^2 \leq \varepsilon$$

Better Rates for Lazy Aggregation

Table 2 Comparison of existing and proposed theoretically-supported methods employing lazy aggregation. In the rates for our methods, $M_1 = L_- + L_+ \sqrt{B/A}$ and $M_2 = \max \left\{ L_- + L_+ \sqrt{2B/A}, A/2\mu \right\}$.

Method	Simple method?	Uses a contractive compressor \mathcal{C} ?	Strongly convex rate	PŁ nonconvex rate	General nonconvex rate
LAG (Chen et al., 2018)	✓	✗	linear ⁽⁹⁾	✗	✗
LAQ (Sun et al., 2019)	✗	✓ ⁽¹⁾	linear ⁽³⁾	✗	✗
LENA (Ghadikolaie et al., 2021) ⁽⁷⁾	✓ ⁽⁴⁾	✓ ⁽⁸⁾	$\mathcal{O}(G^4/T^2\mu^2)$ ^{(5), (6)}	$\mathcal{O}(G^4/T^2\mu^2)$ ^{(5), (6)}	$\mathcal{O}(G^{4/3}/T^{2/3})$ ⁽⁶⁾
LAG (NEW, 2022)	✓	✗	$\mathcal{O}(\exp(-T\mu/M_2))$	$\mathcal{O}(\exp(-T\mu/M_2))$	$\mathcal{O}(M_1/T)$
CLAG (NEW, 2022)	✓	✓ ⁽²⁾	$\mathcal{O}(\exp(-T\mu/M_2))$	$\mathcal{O}(\exp(-T\mu/M_2))$	$\mathcal{O}(M_1/T)$

⁽¹⁾ They consider a specific form of quantization only.

⁽²⁾ Works with any contractive compressor, including low rank approximation, Top- K , Rand- K , quantization, and more.

⁽³⁾ Their Theorem 1 does not present any *explicit* linear rate.

⁽⁴⁾ LENA employs the classical EF mechanism, but it is not clear what is this mechanism supposed to do.

⁽⁵⁾ They consider an assumption (μ -quasi-strong convexity) that is slightly stronger than our PŁ assumption. Both are weaker than strong convexity.

⁽⁶⁾ They assume the local gradients to be bounded by G ($\|\nabla f_i(x)\| \leq G$ for all x). We do not need such a strong assumption.

⁽⁷⁾ They also consider the 0-quasi-strong convex case (slight generalization of convexity); we do not consider the convex case. Moreover, they consider the stochastic case as well, we do not. We specialized all their results to the deterministic (i.e., full gradient) case for the purposes of this table.

⁽⁸⁾ Their contractive compressor depends on the trigger.

⁽⁹⁾ It is possible to specialize their method and proof so as to recover LAG as presented in our work, and to recover a rate similar to ours.

Better Rates for Lazy Aggregation

Table 2 Comparison of existing and proposed theoretically-supported methods employing lazy aggregation. In the rates for our methods, $M_1 = L_- + L_+ \sqrt{B/A}$ and $M_2 = \max \left\{ L_- + L_+ \sqrt{2B/A}, A/2\mu \right\}$.

Method	Simple method?	Uses a contractive compressor \mathcal{C} ?	Strongly convex rate	PŁ nonconvex rate	General nonconvex rate
LAG (Chen et al., 2018)	✓	✗	linear ⁽⁹⁾	✗	✗
LAQ (Sun et al., 2019)	✗	✓ ⁽¹⁾	linear ⁽³⁾	✗	✗
LENA (Ghadikolaei et al., 2021) ⁽⁷⁾	✓ ⁽⁴⁾	✓ ⁽⁸⁾	$\mathcal{O}(G^4/T^2\mu^2)$ ^{(5), (6)}	$\mathcal{O}(G^4/T^2\mu^2)$ ^{(5), (6)}	$\mathcal{O}(G^{4/3}/T^{2/3})$ ⁽⁶⁾
LAG (NEW, 2022)	✓	✗	$\mathcal{O}(\exp(-T\mu/M_2))$	$\mathcal{O}(\exp(-T\mu/M_2))$	$\mathcal{O}(M_1/T)$
CLAG (NEW, 2022)	✓	✓ ⁽²⁾	$\mathcal{O}(\exp(-T\mu/M_2))$	$\mathcal{O}(\exp(-T\mu/M_2))$	$\mathcal{O}(M_1/T)$

⁽¹⁾ They consider a specific form of quantization only.

⁽²⁾ Works with any contractive compressor, including low rank approximation, Top- K , Rand- K , quantization, and more.

⁽³⁾ Their Theorem 1 does not present any *explicit* linear rate.

⁽⁴⁾ LENA employs the classical EF mechanism, but it is not clear what is this mechanism supposed to do.

⁽⁵⁾ They consider an assumption (μ -quasi-strong convexity) that is slightly stronger than our PŁ assumption. Both are weaker than strong convexity.

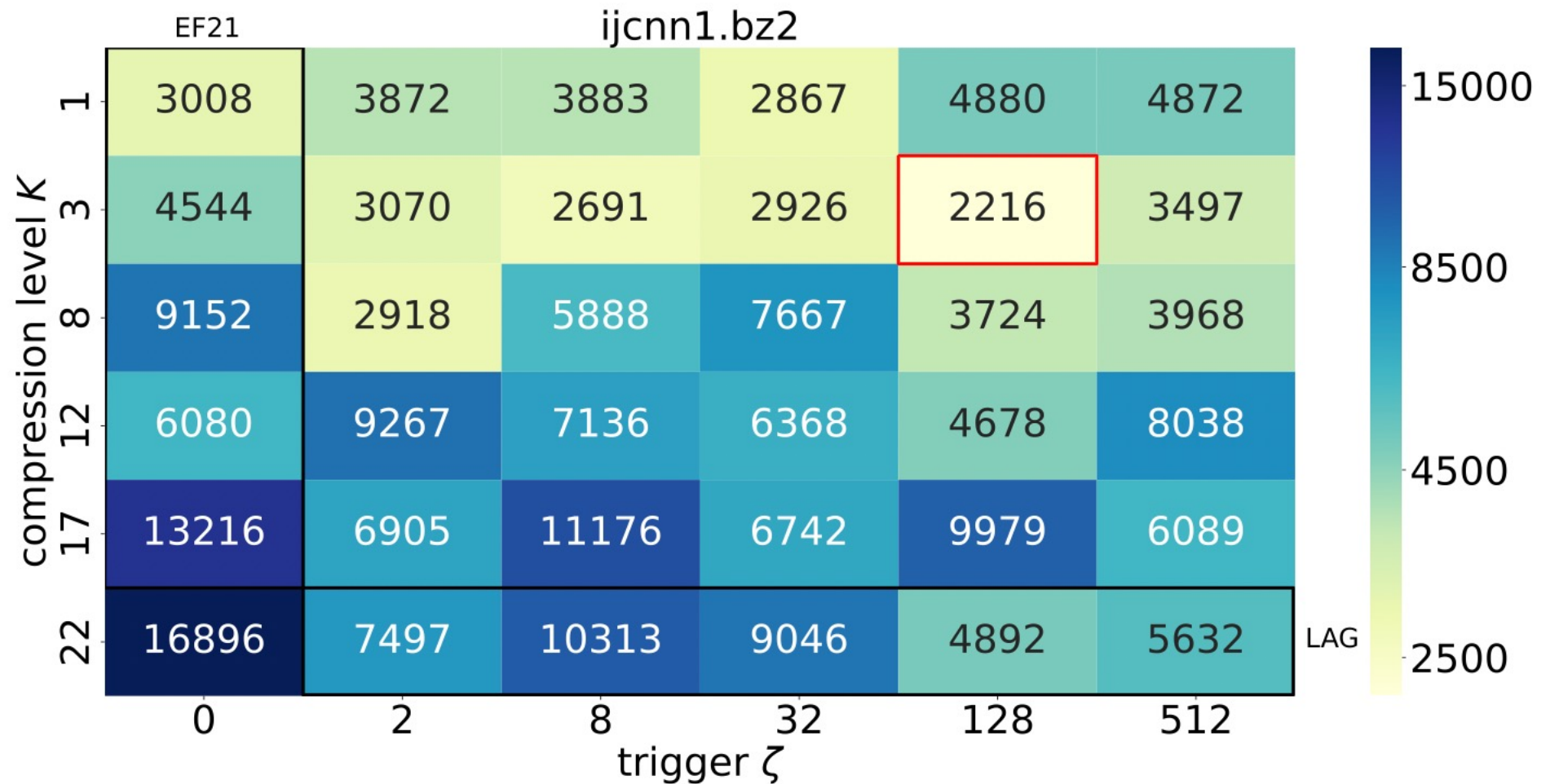
⁽⁶⁾ They assume the local gradients to be bounded by G ($\|\nabla f_i(x)\| \leq G$ for all x). We do not need such a strong assumption.

⁽⁷⁾ They also consider the 0-quasi-strong convex case (slight generalization of convexity); we do not consider the convex case. Moreover, they consider the stochastic case as well, we do not. We specialized all their results to the deterministic (i.e., full gradient) case for the purposes of this table.

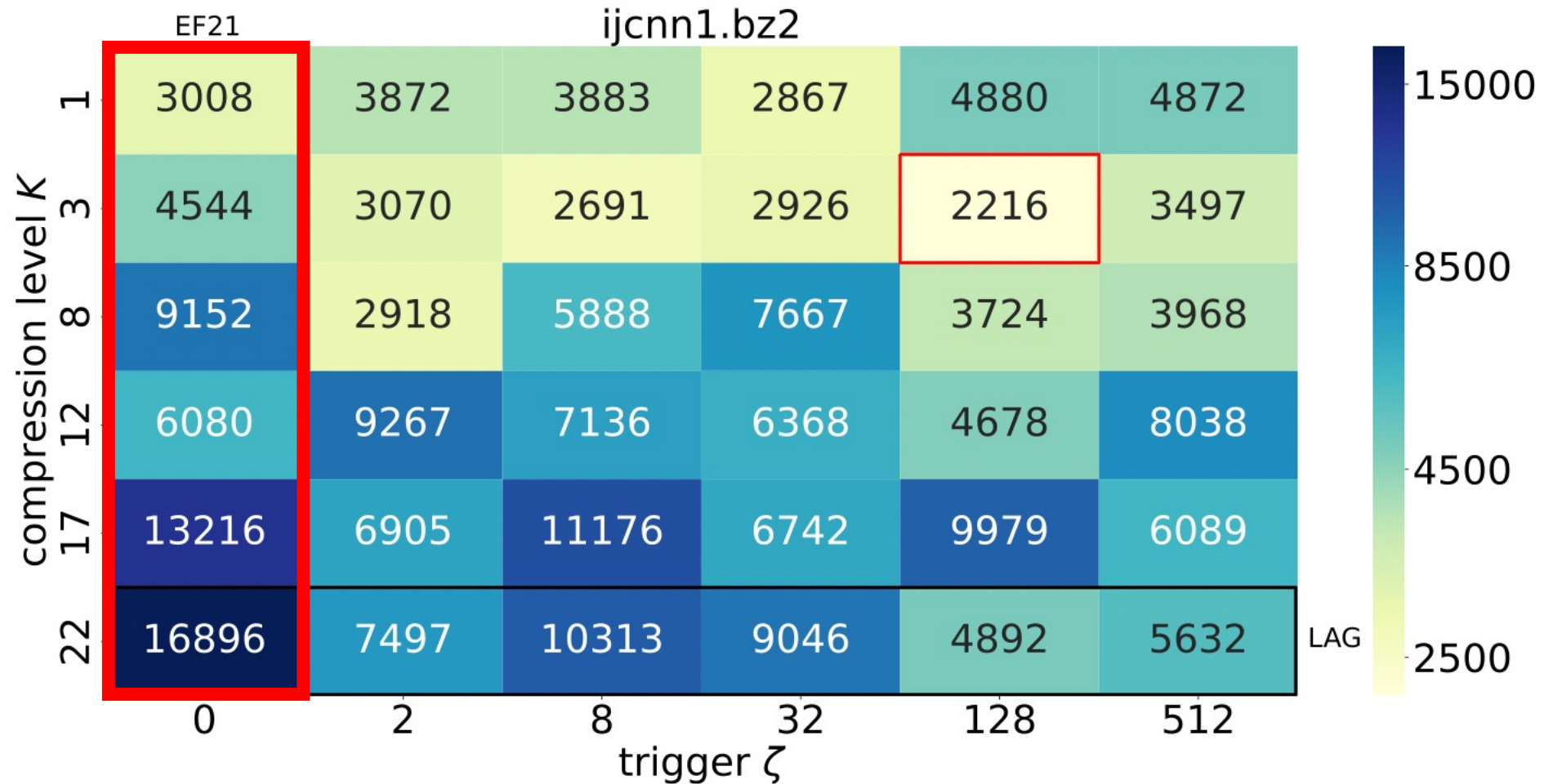
⁽⁸⁾ Their contractive compressor depends on the trigger.

⁽⁹⁾ It is possible to specialize their method and proof so as to recover LAG as presented in our work, and to recover a rate similar to ours.

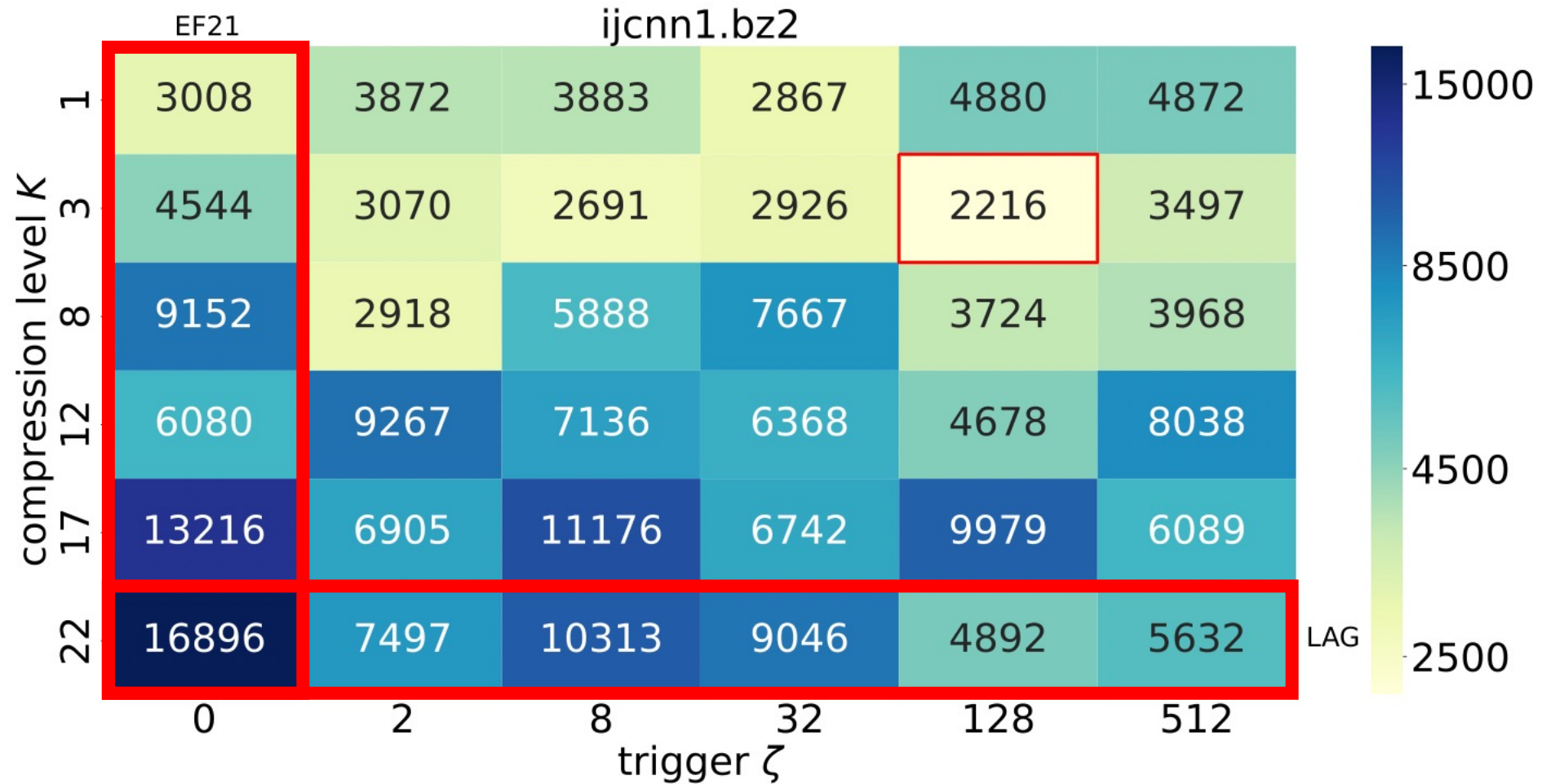
CLAG: Combining the Benefits of EF21 and LAG



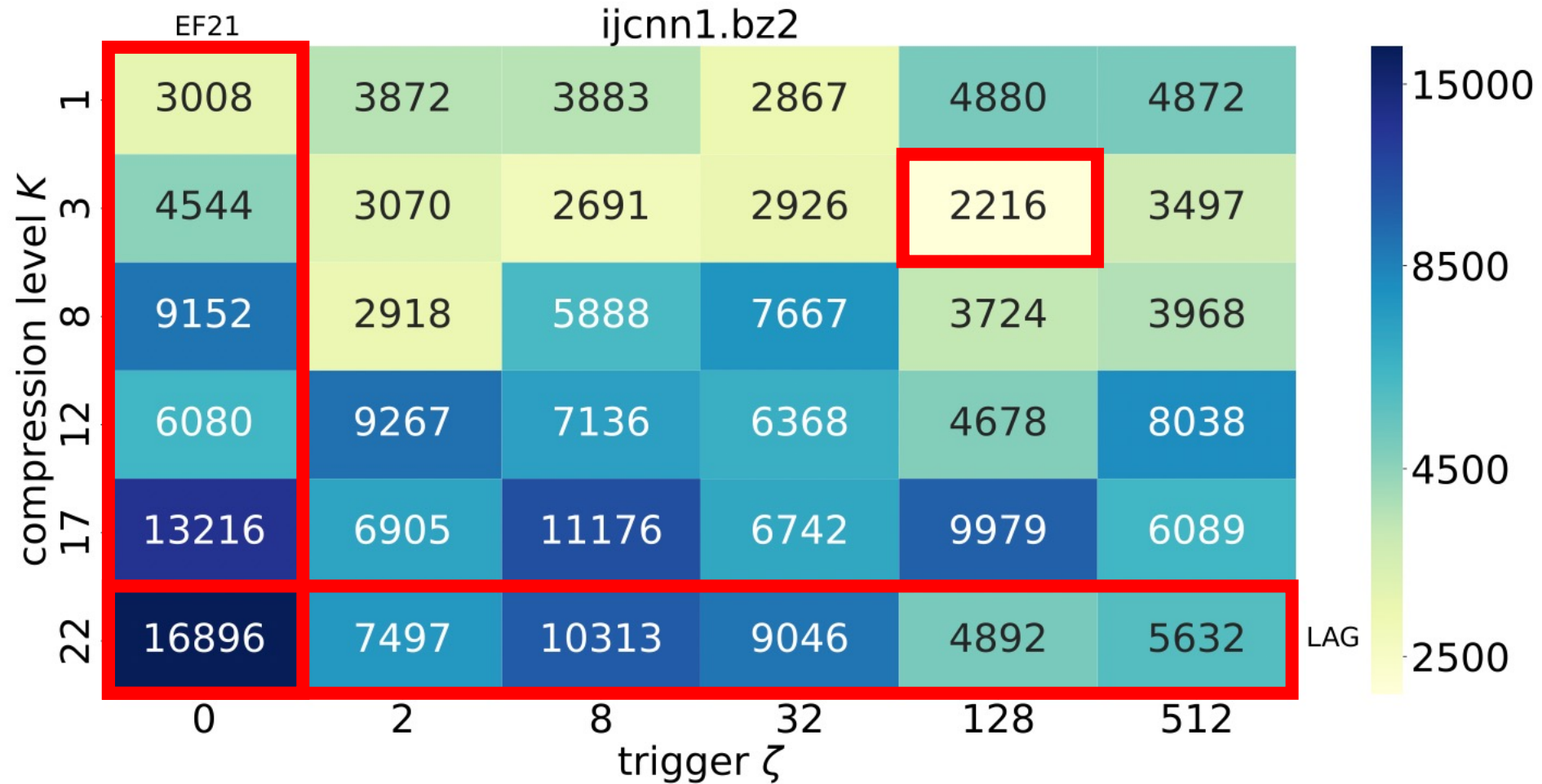
CLAG: Combining the Benefits of EF21 and LAG



CLAG: Combining the Benefits of EF21 and LAG



CLAG: Combining the Benefits of EF21 and LAG



Summary & Extensions

- 3PC recovers
 - **gradient descent** (GD) and its rate
 - **error feedback** (EF21) and its rate
 - **lazy aggregation** (LAG) & gets a **better** rate
- 3PC uncovers a hidden link between error feedback and lazy aggregation mechanisms & literature
- 3PC includes
 - **new method (CLAG)** which combines the benefits of EF21 and LAG
 - several additional new methods (not mentioned in this talk)
- We prove
 - **Fast $O(1/t)$ rate** for smooth nonconvex functions
 - **linear rate** under the Polyak-Łojasiewicz condition (not mentioned in this talk)