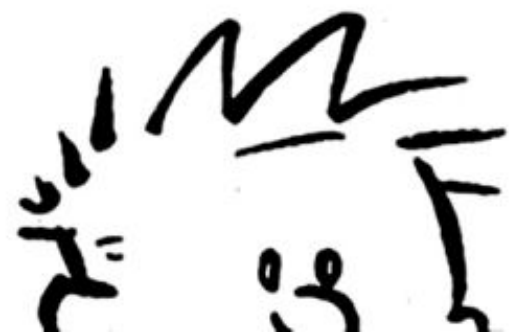# Role of Momentum in Byzantine Resilience

## of Distributed Learning

**Nirupam Gupta**

**EPFL**

"All models are wrong, but some are useful."

– George Box ?

# Machine Learning

# Machine Learning

**DATA**

# Machine Learning

**DATA**

# Machine Learning

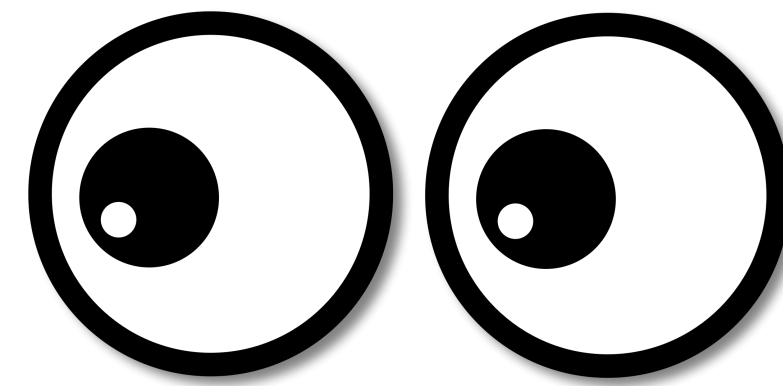**DATA**                                                                  **AMATEUR MACHINE**
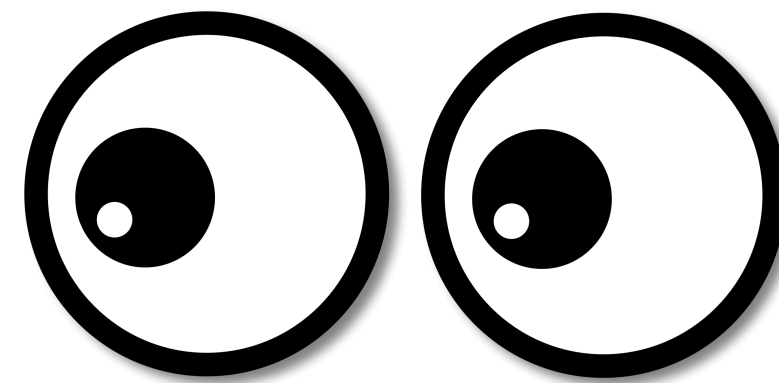
# Machine Learning

**DATA**



**AMATEUR MACHINE**

# Machine Learning

**DATA**

**AMATEUR MACHINE**

# Machine Learning

**DATA**



**WATCH & LEARN**

**AMATEUR MACHINE**

# Machine Learning

**DATA**

**AMATEUR MACHINE**

**WATCH & LEARN**

$\mathcal{D}$

# Machine Learning

**DATA**

**WATCH & LEARN**

**AMATEUR MACHINE**

$\mathcal{D}$

$\theta$

# Machine Learning

**DATA**



$\mathcal{D}$

**WATCH & LEARN**



**AMATEUR MACHINE**



$\theta$

**SOLVE**

$$\theta^* \leftarrow \arg\min_{\theta} Q(\theta) := \mathop{\mathbb{E}}_{x \sim \mathcal{D}} q(\theta, x)$$

# Machine Learning

**DATA**

**AMATEUR MACHINE**

**WATCH & LEARN**

$\mathcal{D}$

$\theta$

**SOLVE**

$$\theta^* \leftarrow \arg\min_\theta Q(\theta) := \mathbb{E}_{x \sim \mathcal{D}} \, q(\theta, x)$$

# Machine Learning

**DATA**

**AMATEUR MACHINE**

**WATCH & LEARN**

$\mathscr{D}$

$\theta$

**SOLVE**

$$\theta^* \leftarrow \arg\min_{\theta} Q(\theta) := \mathbb{E}_{x \sim \mathscr{D}} \, q\,(\theta, x)$$

# Machine Learning

**DATA**

**AMATEUR MACHINE**



$\mathscr{D}$

**WATCH & LEARN**

$\theta$

Loss per data-point

**SOLVE**

$$\theta^* \leftarrow \arg\min_{\theta} Q(\theta) := \mathbb{E}_{x \sim \mathscr{D}} \; q\,(\theta, x)$$

# Machine Learning

**DATA**

**EXPERT MACHINE**

$\mathscr{D}$

**WATCH & LEARN**

$\theta$

Loss per data-point

**SOLVE**

$$\theta^* \leftarrow \arg\min_{\theta} Q(\theta) := \mathbb{E}_{x \sim \mathscr{D}}\ q\,(\theta, x)$$

# Machine Learning

**DATA**

**EXPERT MACHINE**

$\mathscr{D}$

**WATCH & LEARN**

Loss per data-point

**SOLVE**

$$\theta* \leftarrow \arg \min_{\theta} Q(\theta) := \mathbb{E}_{x \sim \mathscr{D}} \, q\,(\theta, x)$$
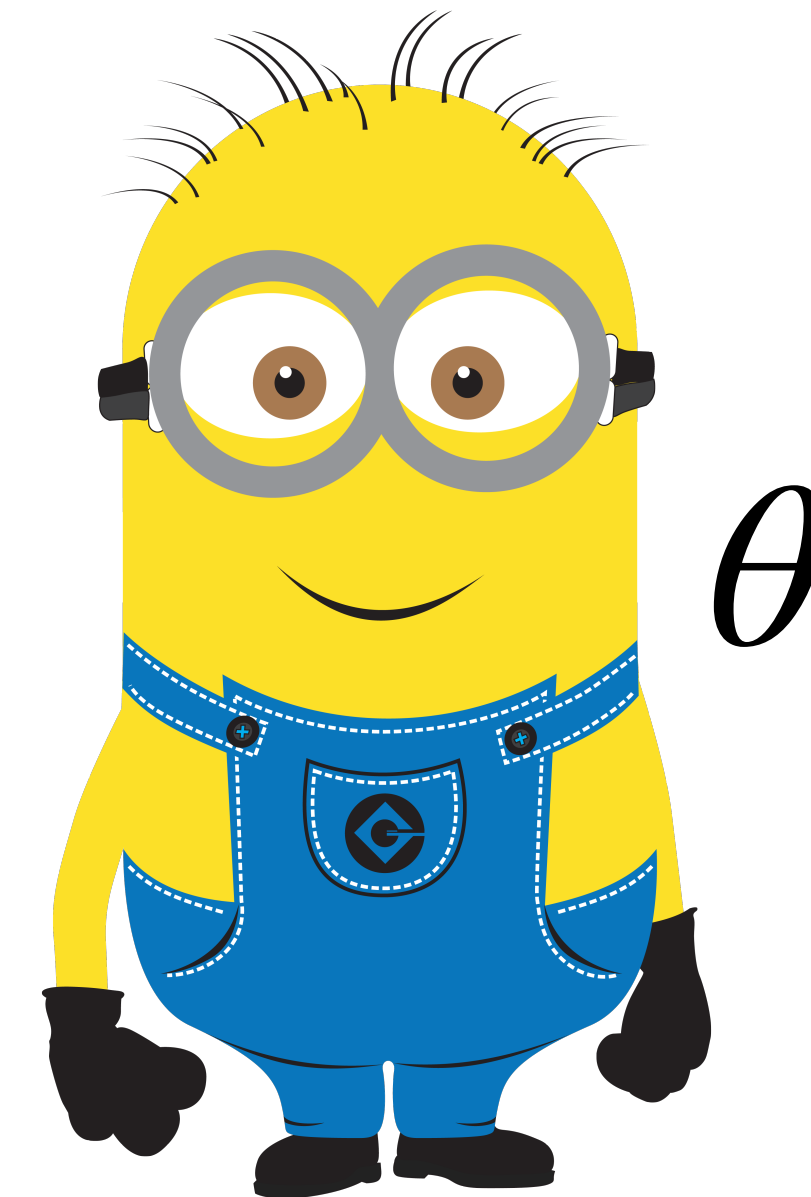
# Machine Learning

**DATA**

**EXPERT MACHINE**

$\mathscr{D}$

**WATCH & LEARN**

$\theta*$

Loss per data-point

**SOLVE**

$$\theta* \leftarrow \arg\min_{\theta} Q(\theta) := \mathbb{E}_{x \sim \mathscr{D}} \, q(\theta, x)$$

# Machine Learning

**DATA**

**EXPERT MACHINE**

$\mathscr{D}$

**WATCH & LEARN**

$\theta*$

Loss per data-point

**SOLVE**

$$\theta* \leftarrow \arg\min_{\theta} Q(\theta) := \mathbb{E}_{x \sim \mathscr{D}} \; q\,(\theta, x)$$
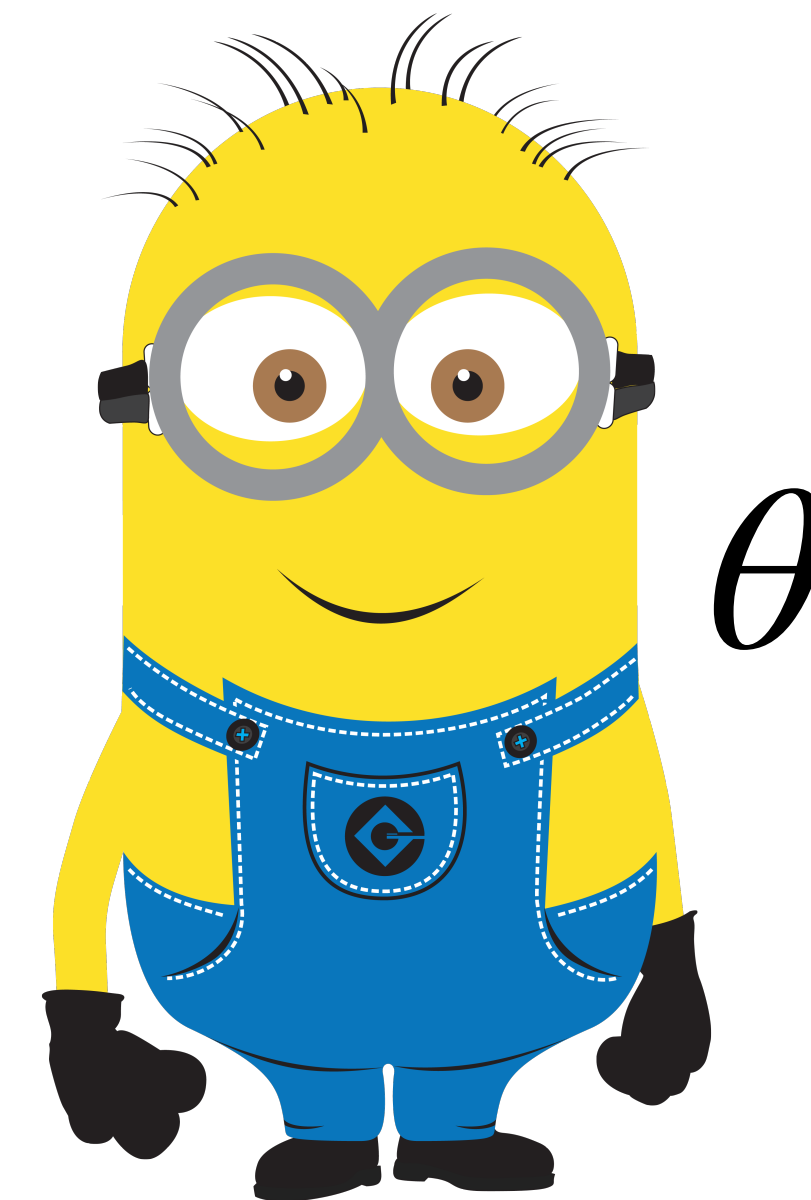
# Machine Learning

**DATA**

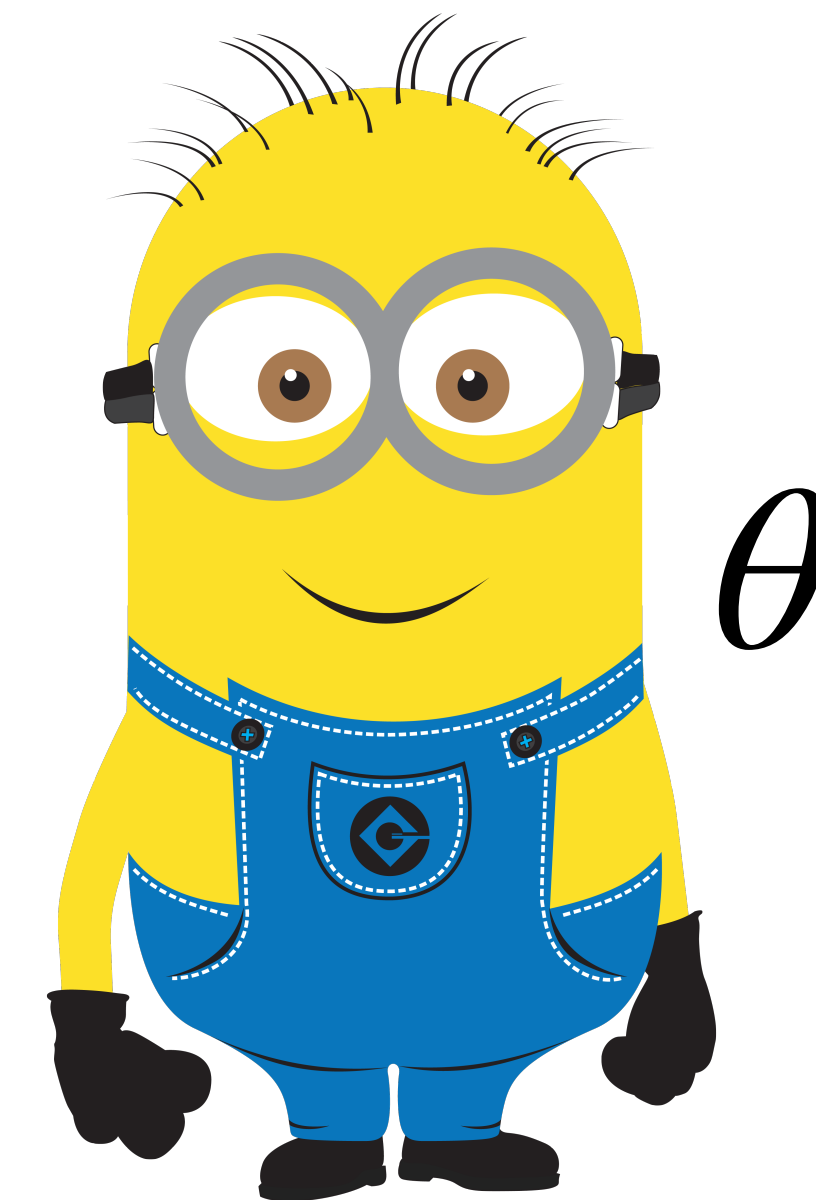**EXPERT MACHINE**

$\mathscr{D}$

**WATCH & LEARN**



$\theta*$

Loss per data-point

**SOLVE**

$$\theta* \leftarrow \arg\min_{\theta} Q(\theta) := \mathbb{E}_{x \sim \mathscr{D}} \, q\,(\theta, x)$$

Relax

# Machine Learning

**DATA**

**EXPERT MACHINE**

$$\mathscr{D}$$

**WATCH & LEARN**

$$\theta*$$

Loss per data-point

**SOLVE**

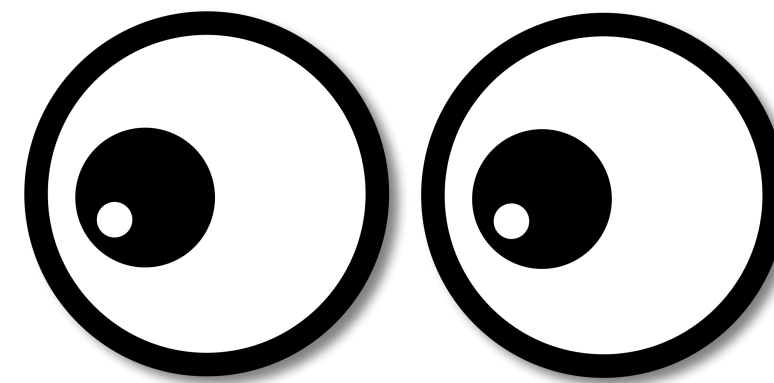$$\theta* \leftarrow \arg\min_{\theta} Q(\theta) := \mathbb{E}_{x \sim \mathscr{D}} \, q\,(\theta, x)$$

Relax

**MORE REASONABLE GOAL**

$$\theta* \in \big(\theta \,;\, \nabla Q(\theta) = 0\big)$$

# Machine Learning

**DATA**

**"EXPERT" MACHINE**

$\mathcal{D}$

**WATCH & LEARN**

$\theta*$

Loss per data-point

**SOLVE**

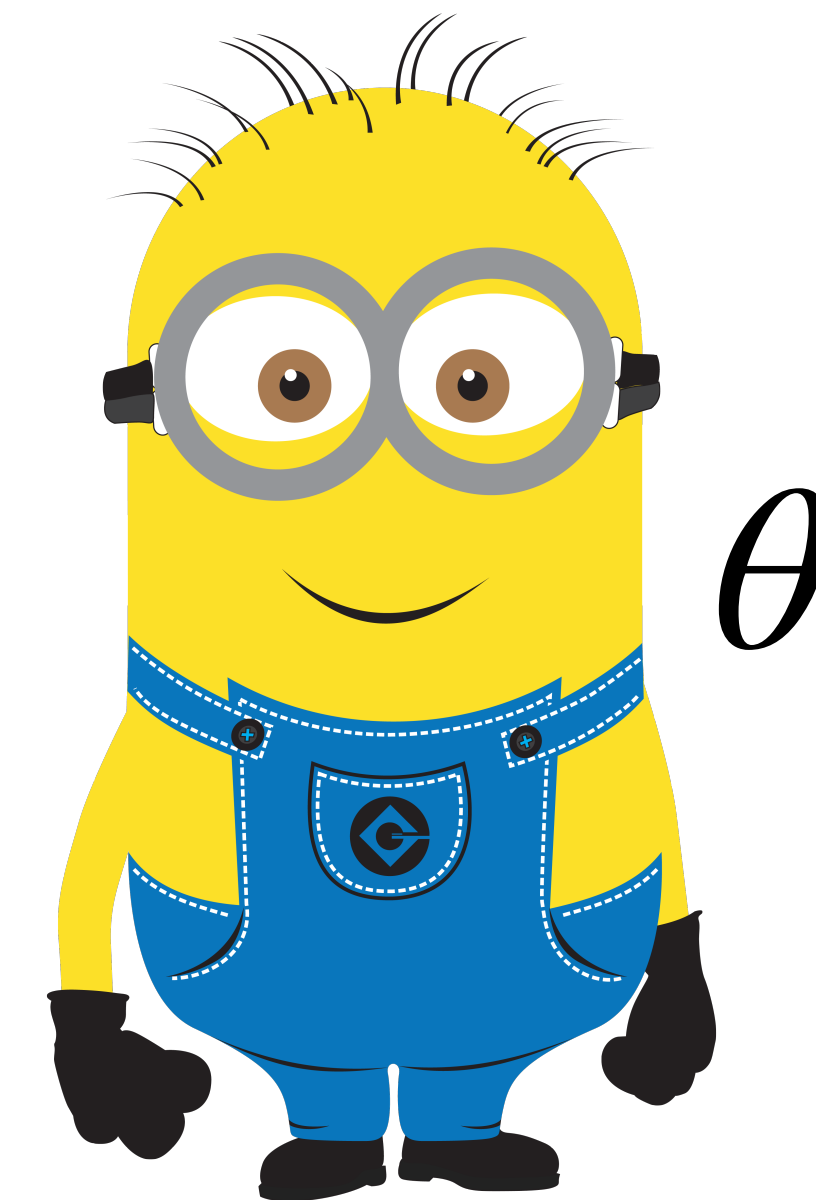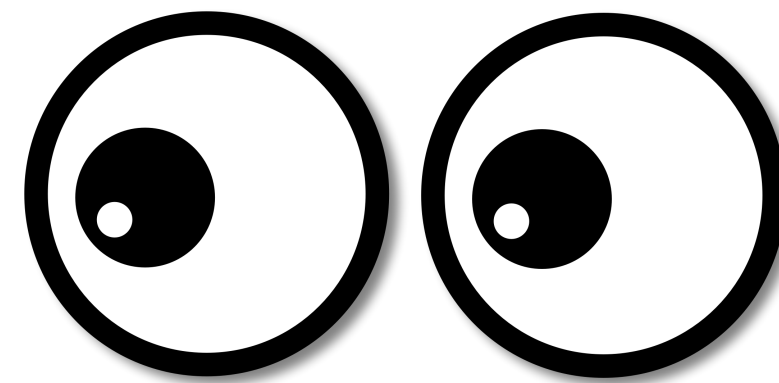$$\theta* \leftarrow \arg\min_{\theta} Q(\theta) := \mathbb{E}_{x \sim \mathcal{D}} q(\theta, x)$$

Relax

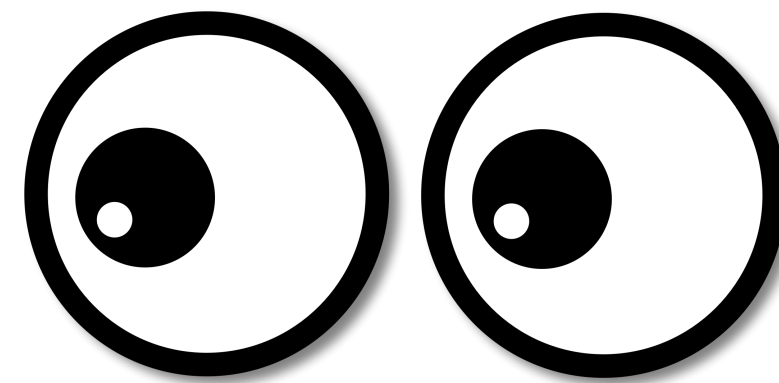**MORE REASONABLE GOAL**

$$\theta* \in (\theta \ ; \ \nabla Q(\theta) = 0)$$

# Distributed Learning

# Distributed Learning

**DATA IS GROWING**

# Distributed Learning

**DATA IS GROWING**

# Distributed Learning

**DATA IS GROWING**

**SO ARE THE MODELS**

# Distributed Learning

**DATA IS GROWING**

**SO ARE THE MODELS**

# Distributed Learning

**DATA IS GROWING**



**SO ARE THE MODELS**



Number of parameters

Switch-C
1.6t

Gshard
600b

OpenAI
GPT3
175b

T-NLG
17b

NVIDIA
MegatronLM
8.3b

OpenAI
GPT2
1.5b

UNIVERSITY of
WASHINGTON
Grover-Mega
1.5b

Ai2
Transformer
ELMo
465m

Google
BERT-Large
340m

MT-DNN
330m

XLM
665m

RoBERTa
355m

Ai2
ELMo
94m

OpenAI
GPT
100m

XLNET
340m
Carnegie
Mellon
University

DistilBERT
66m

1t

100b

10b

1b

.01b

January-18    April-18    July-18    October-18    February-19    May-19    August-19    December-19    March-20    June-20    September-20    January-...

# Distributed Learning



**DATA IS GROWING**

**SO ARE THE MODELS**

**NEED MANY MACHINES TO TRAIN**

# Distributed Learning

**DATA IS GROWING**

**SO ARE THE MODELS**

**NEED MANY MACHINES TO TRAIN**

# Distributed SGD

**Divides the workload per machine by the total size of the system**

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$\theta_t$$

SERVER

$\theta_t$

$\theta_t$

$\theta_t$

NODE 1

NODE 2

NODE 3

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ;$$

$$\theta_t$$

SERVER

$\theta_t$

$\theta_t$

$\theta_t$

NODE 1

NODE 2

NODE 3

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ;$$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ;$$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

Nodes query *stochastic gradients* with bounded variance

$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ;$$

True gradient

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

True gradient

$\theta_t$

SERVER

$\theta_t$ $\theta_t$ $\theta_t$

NODE 1 NODE 3

NODE 2

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \boxed{\nabla Q(\theta_t)} + u_t^i \quad ; \quad u_t^i \sim \mathscr{U}_t$$

True gradient

$$\mathbb{E}\left[\mathscr{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathscr{U}_t\|^2\right] \leq \sigma^2 < \infty$$

$$\theta_t$$

SERVER

$\theta_t$

$\theta_t$

$\theta_t$

NODE 1

NODE 3

NODE 2

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \boxed{\nabla Q(\theta_t)} + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

True gradient

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

$\theta_t$

SERVER

$\theta_t$

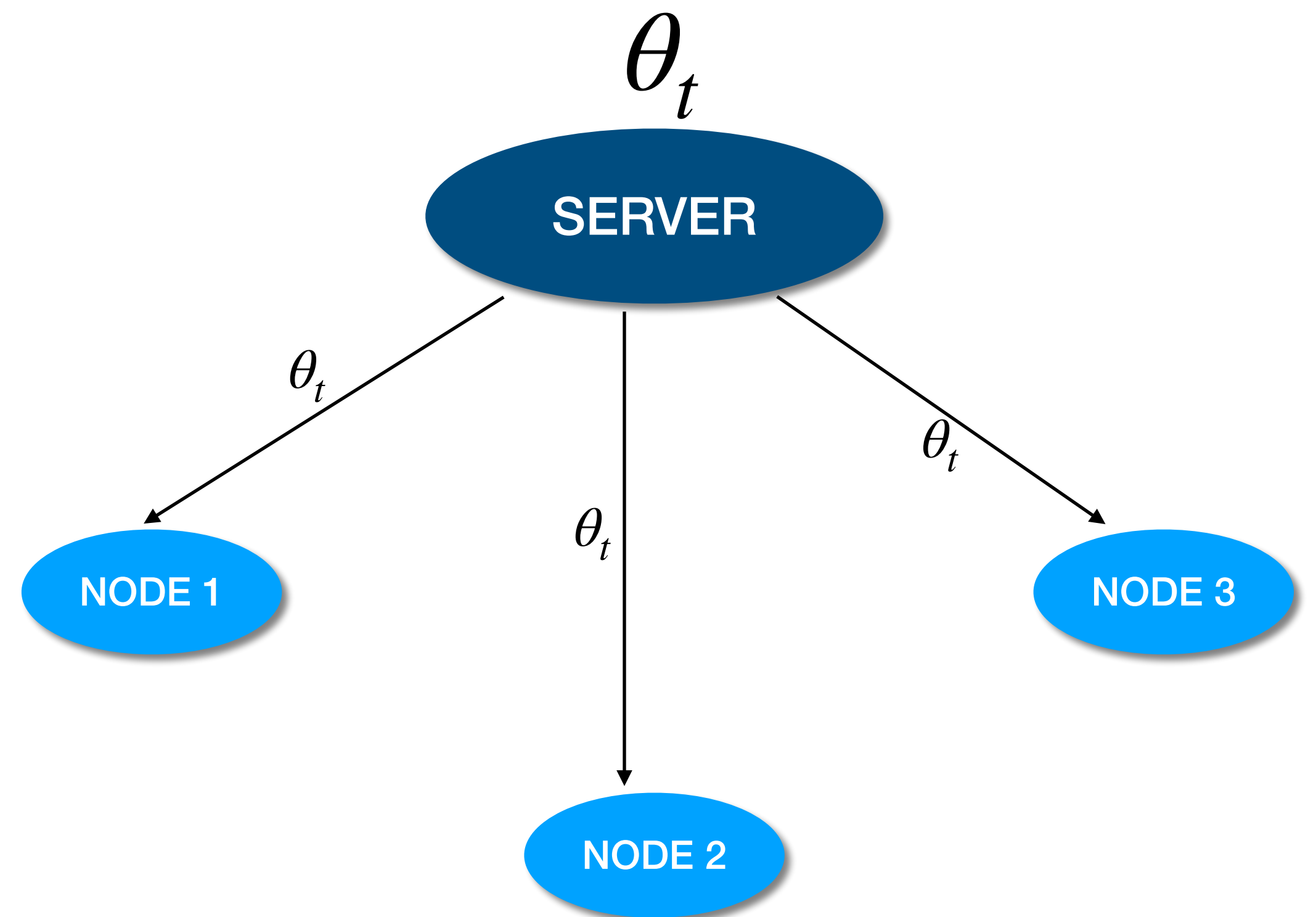$\theta_t$

$\theta_t$

NODE 1

NODE 3

NODE 2

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance
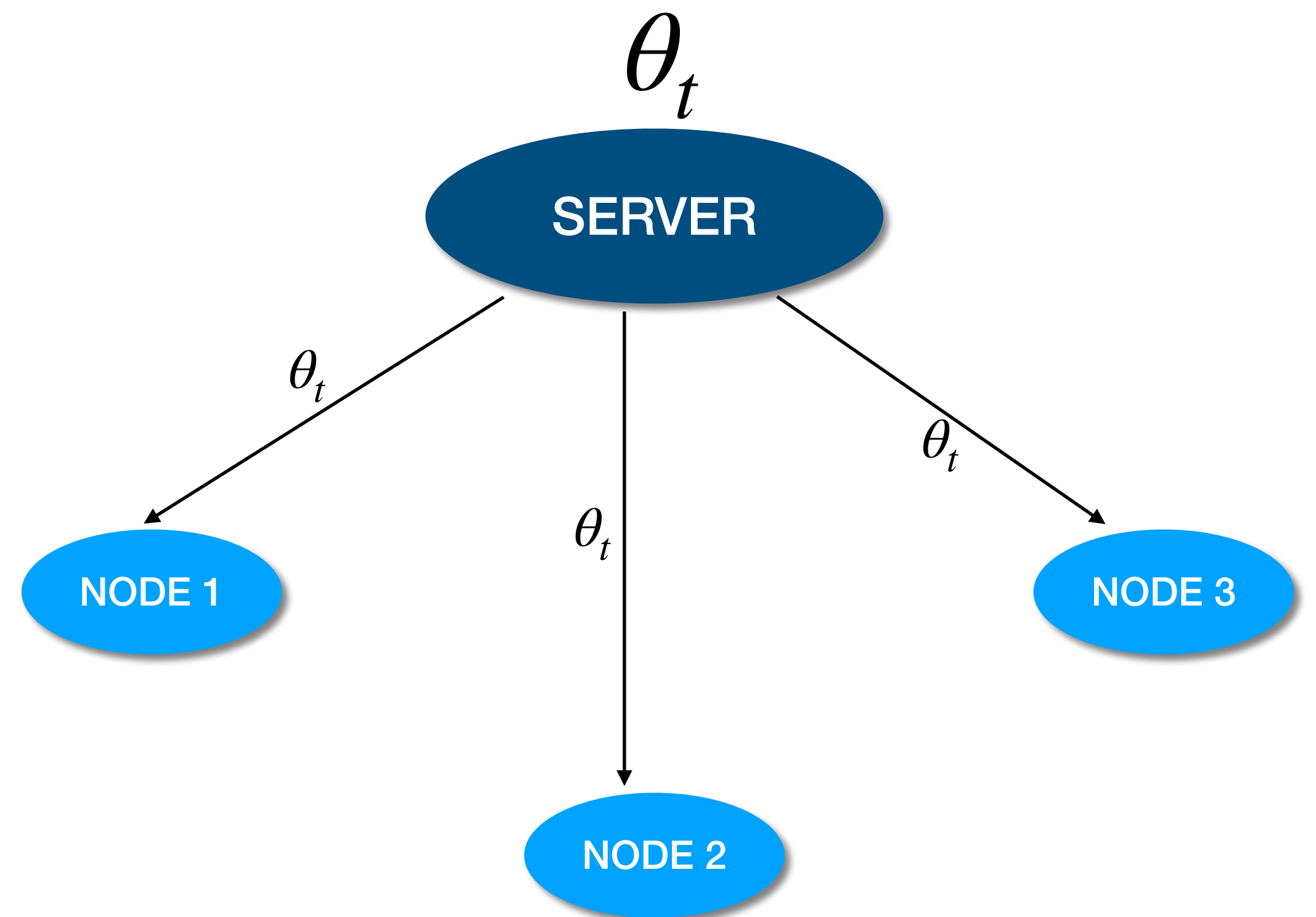
$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

True gradient

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

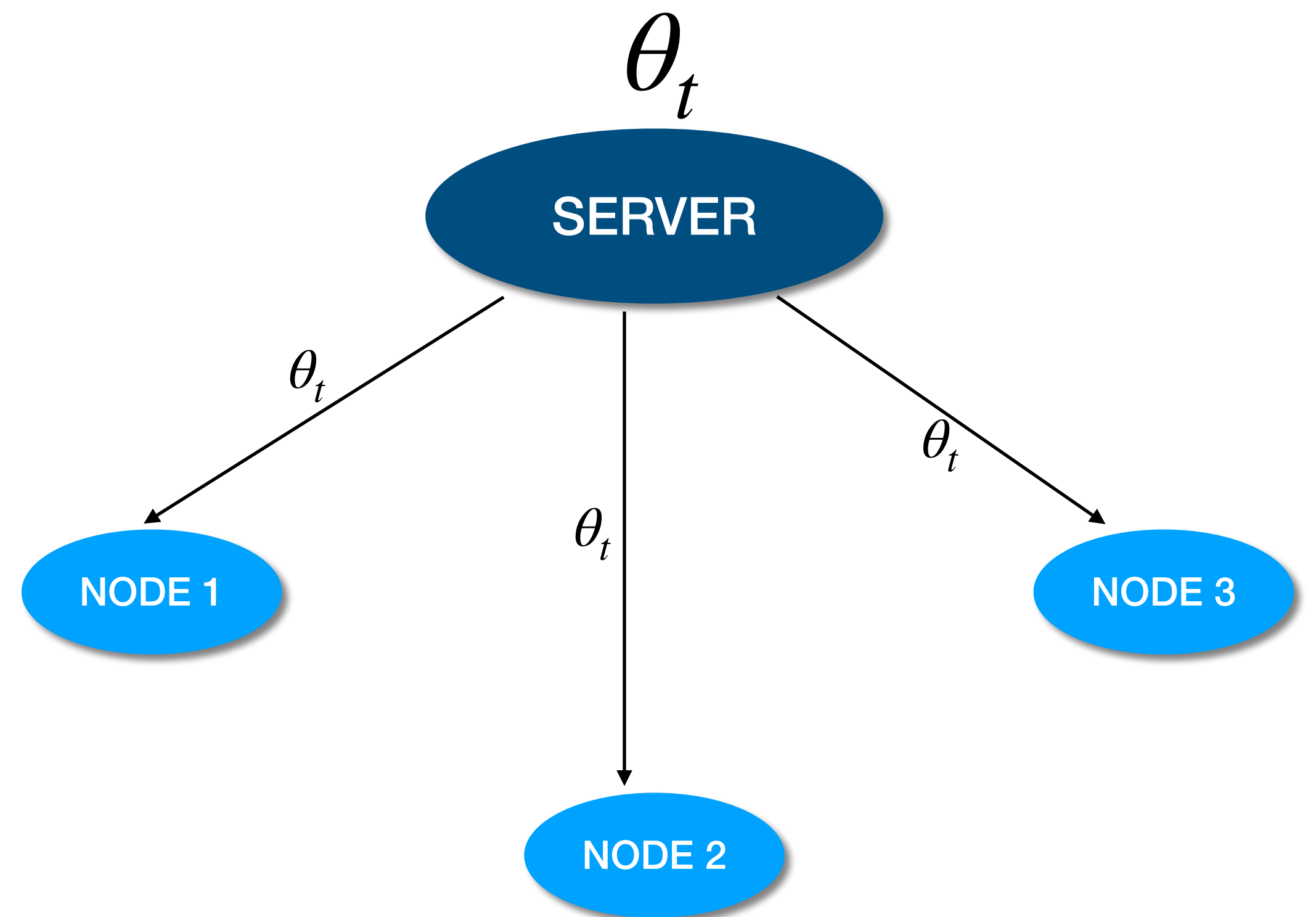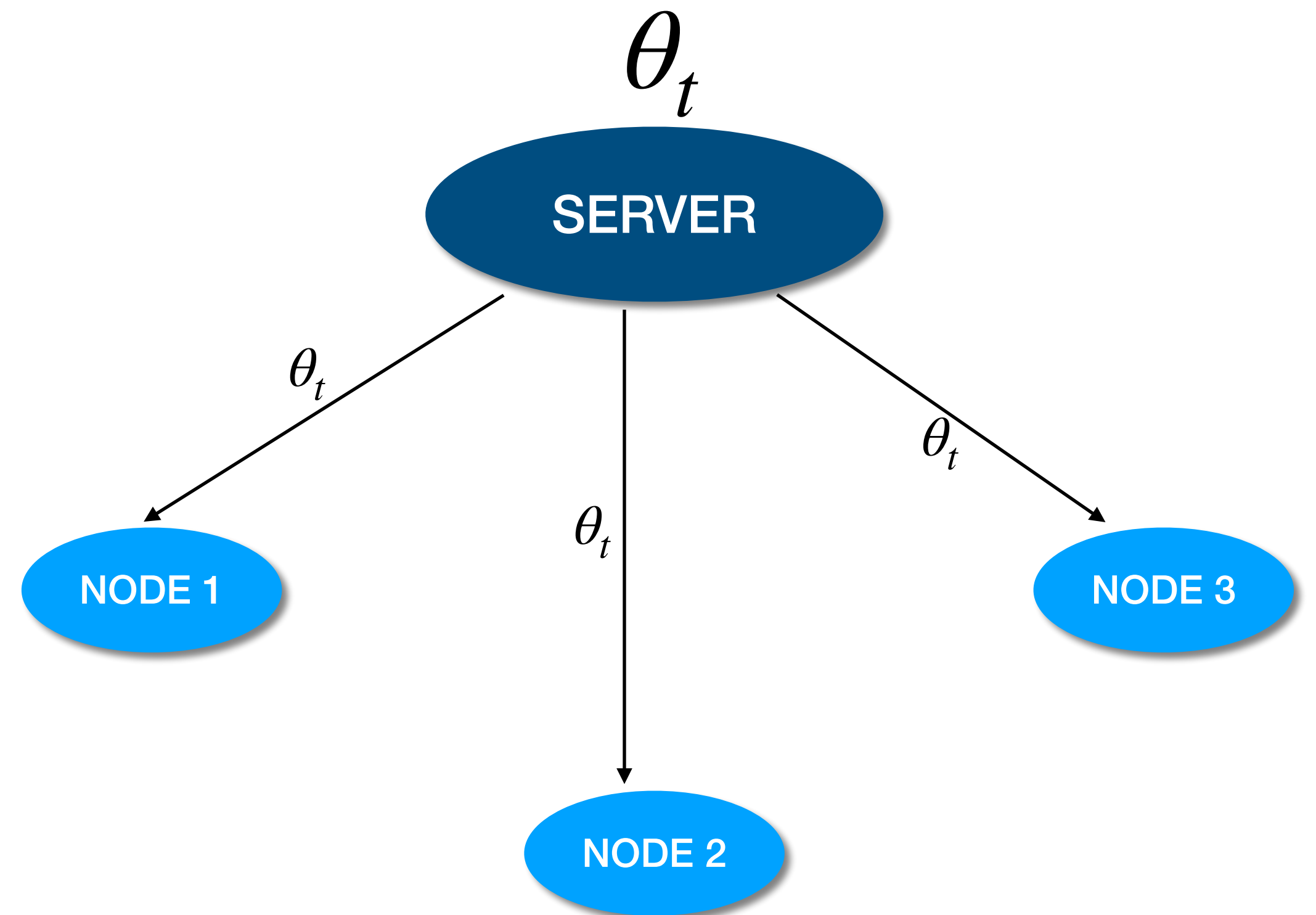$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathscr{U}_t$$

$$\mathbb{E}\left[\mathscr{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathscr{U}_t\|^2\right] \leq \sigma^2 < \infty$$
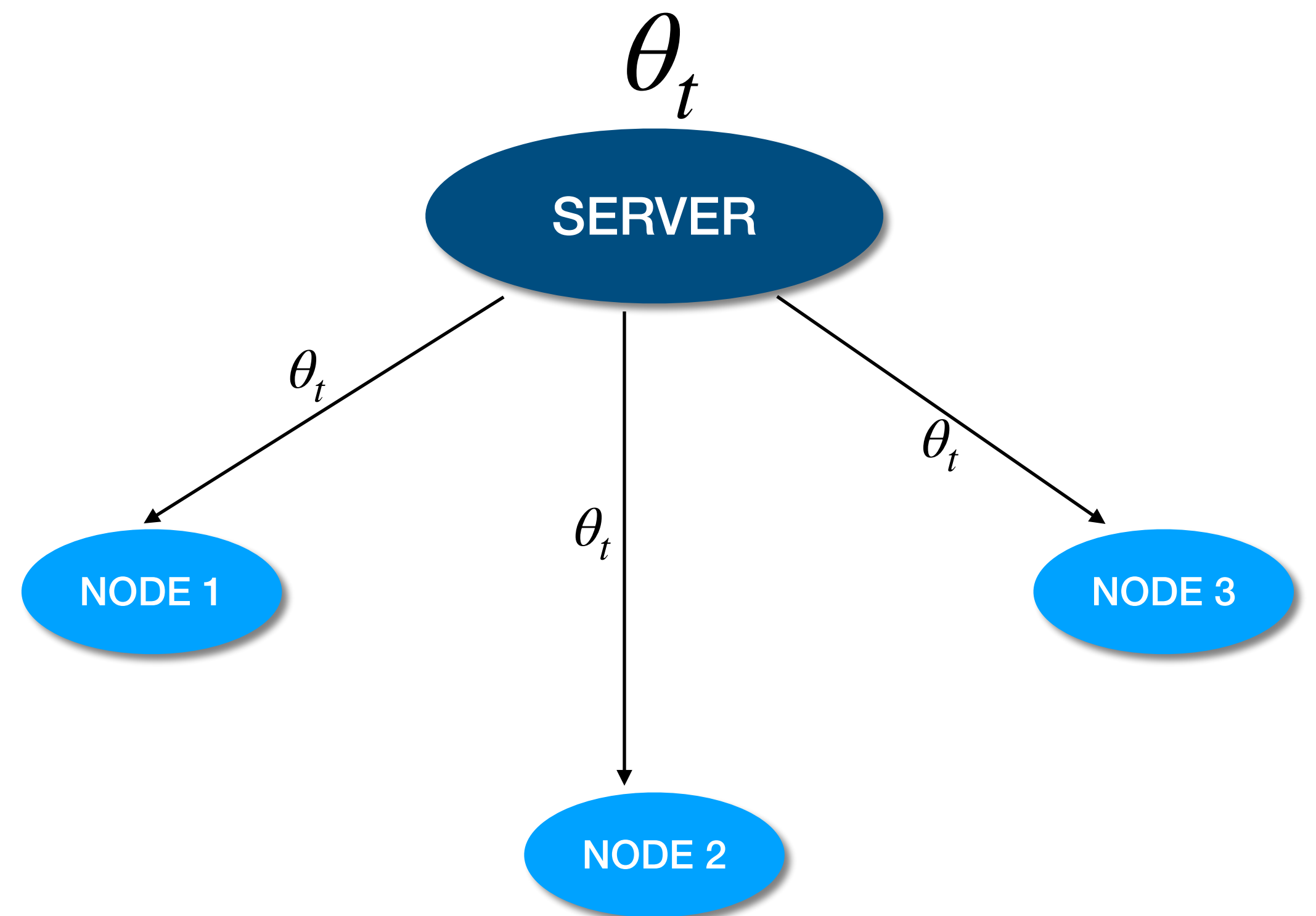
True gradient

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \boxed{\nabla Q(\theta_t)} + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

True gradient

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

$\theta_t$

SERVER

$\theta_t$

$g_t^1$

$\theta_t$ $g_t^2$

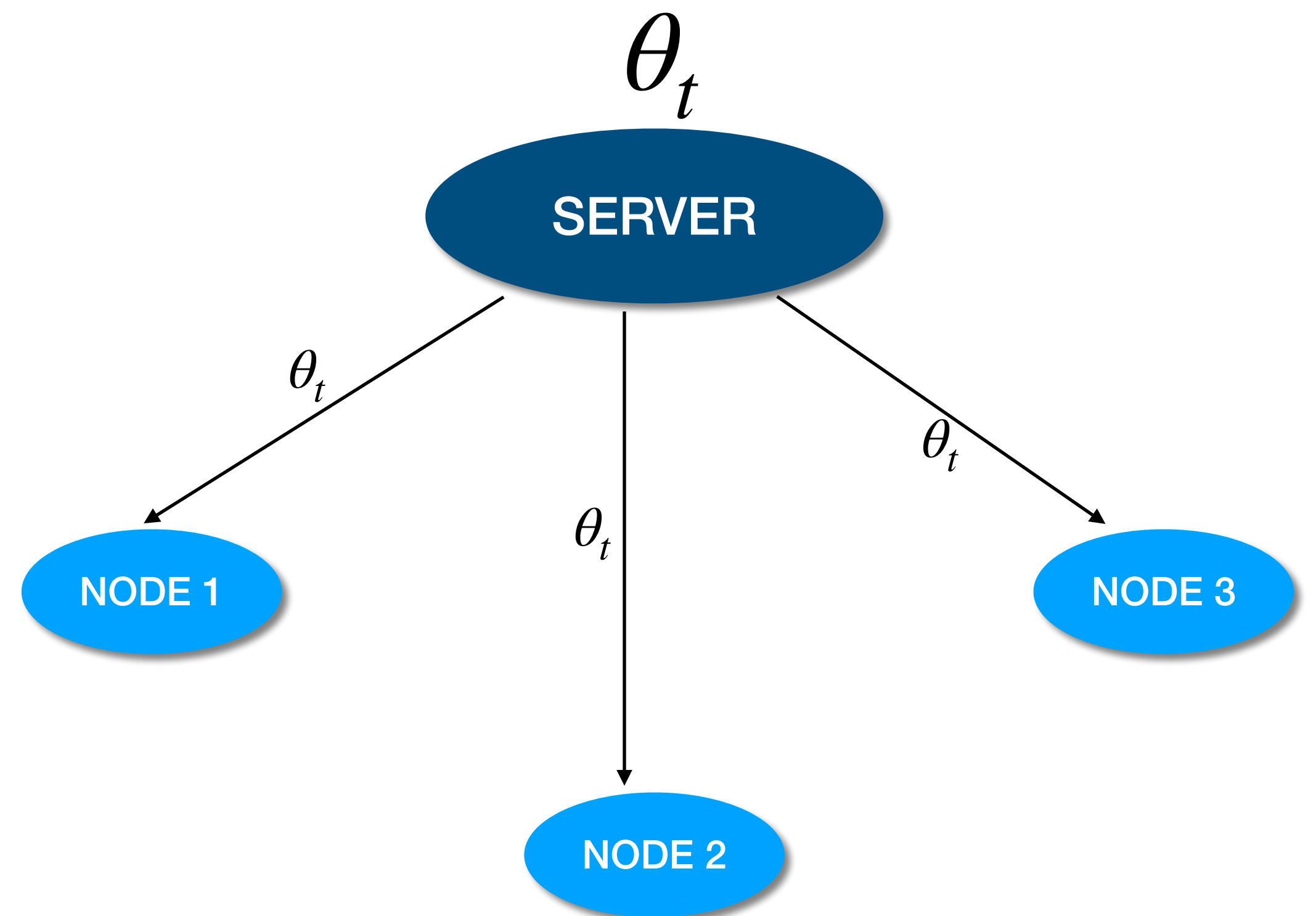$\theta_t$

NODE 1

NODE 3

NODE 2

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$
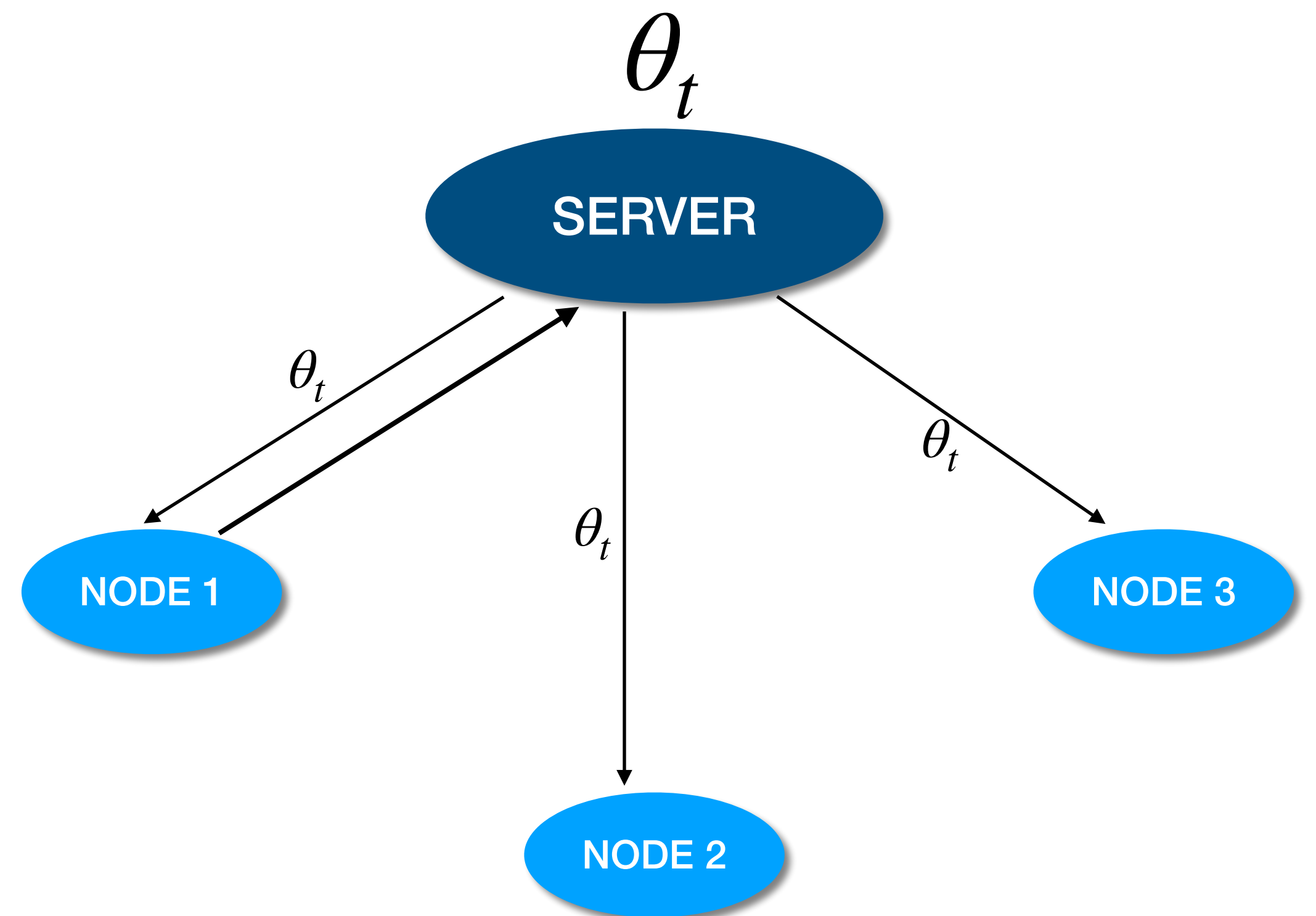
True gradient

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

True gradient

$\theta_t$

SERVER

$\theta_t$

$g_t^1$

$\theta_t$

$g_t^2$

$g_t^3$

$\theta_t$

NODE 1

NODE 2

NODE 3

# Distributed SGD

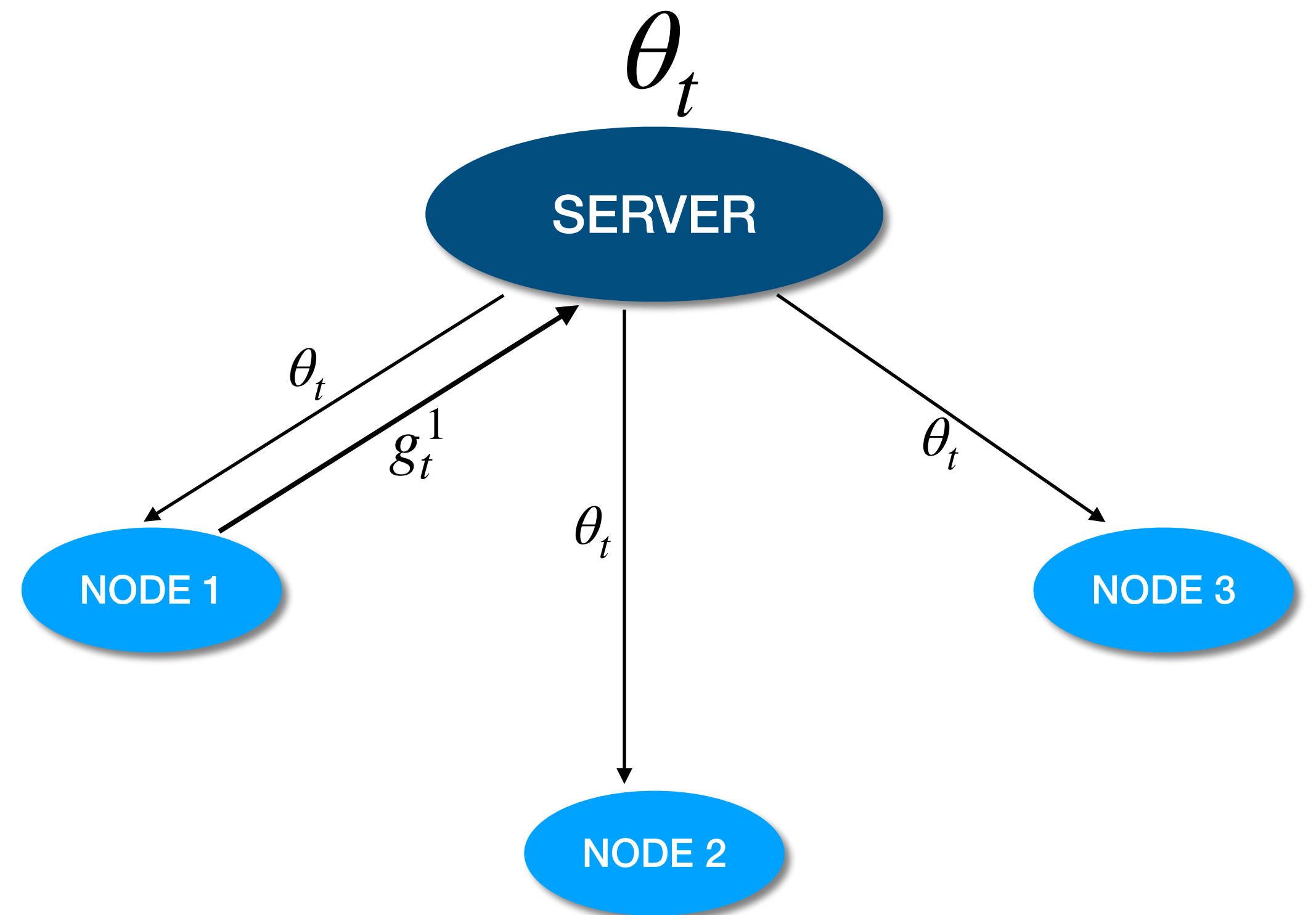**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \boxed{\nabla Q(\theta_t)} + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

True gradient

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

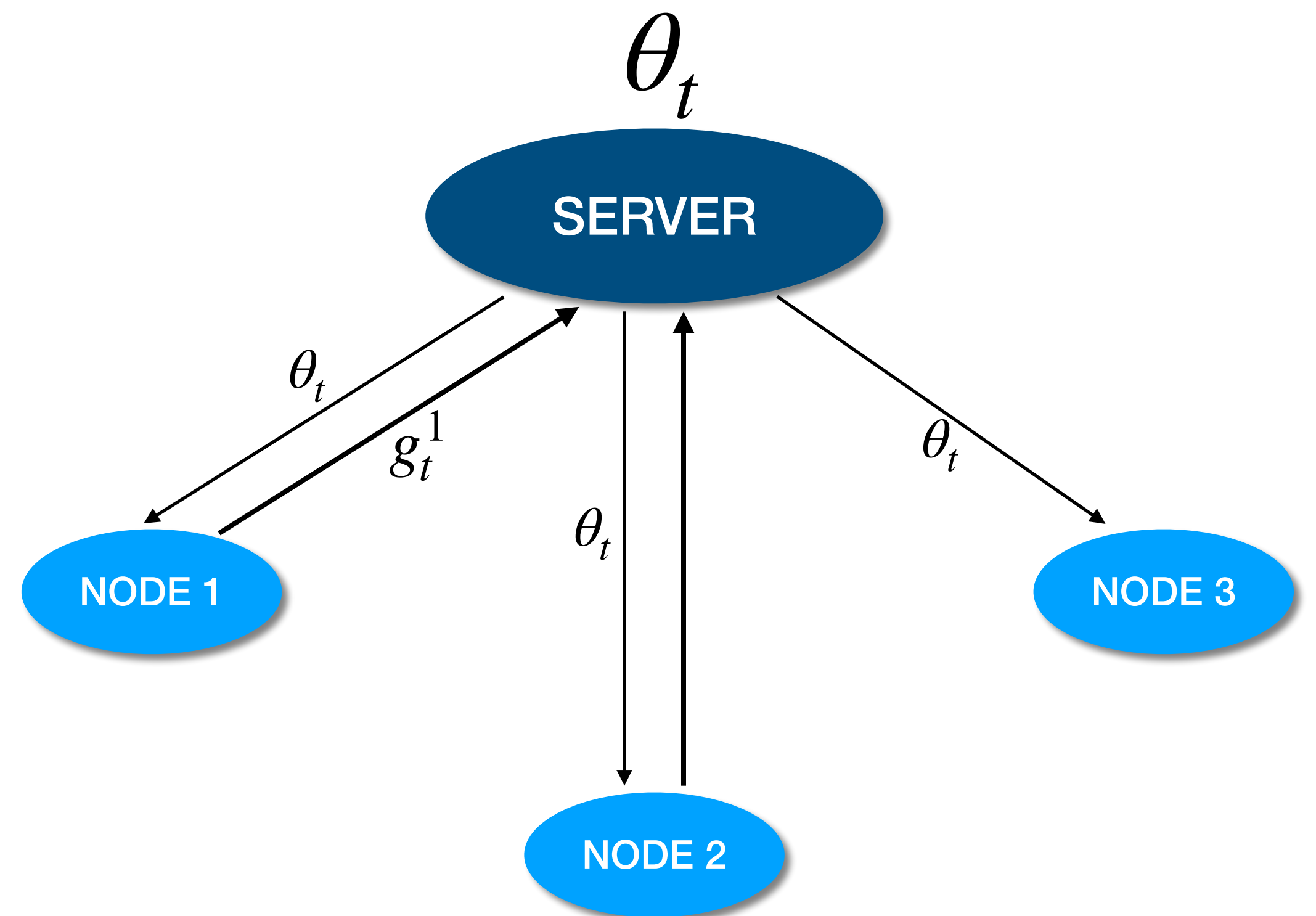**Server averages** the gradients, $\hat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \boxed{\nabla Q(\theta_t)} + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

True gradient

**Server averages** the gradients, $\displaystyle \widehat{g}_t = \frac{1}{n} \sum_i g_t^i$
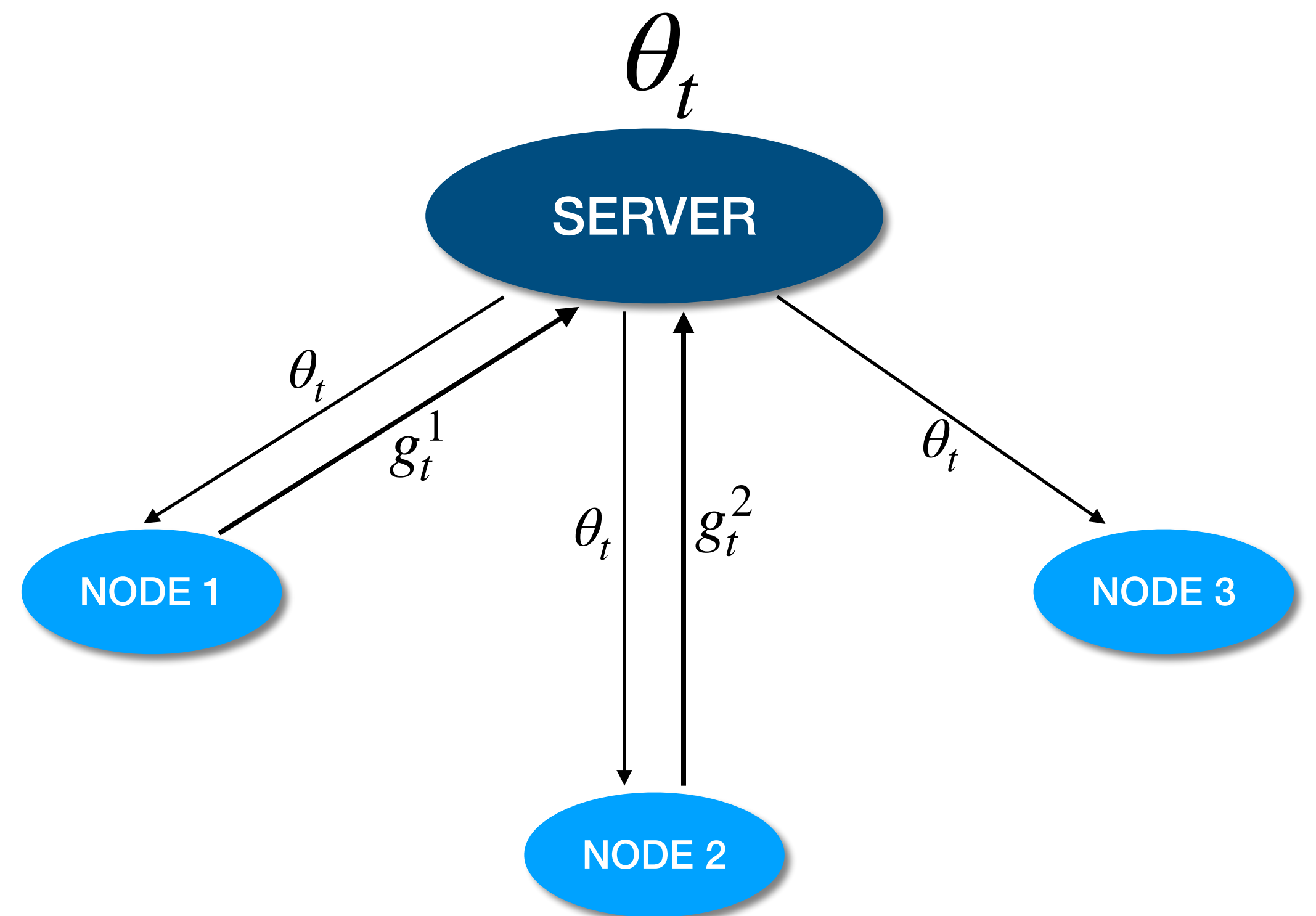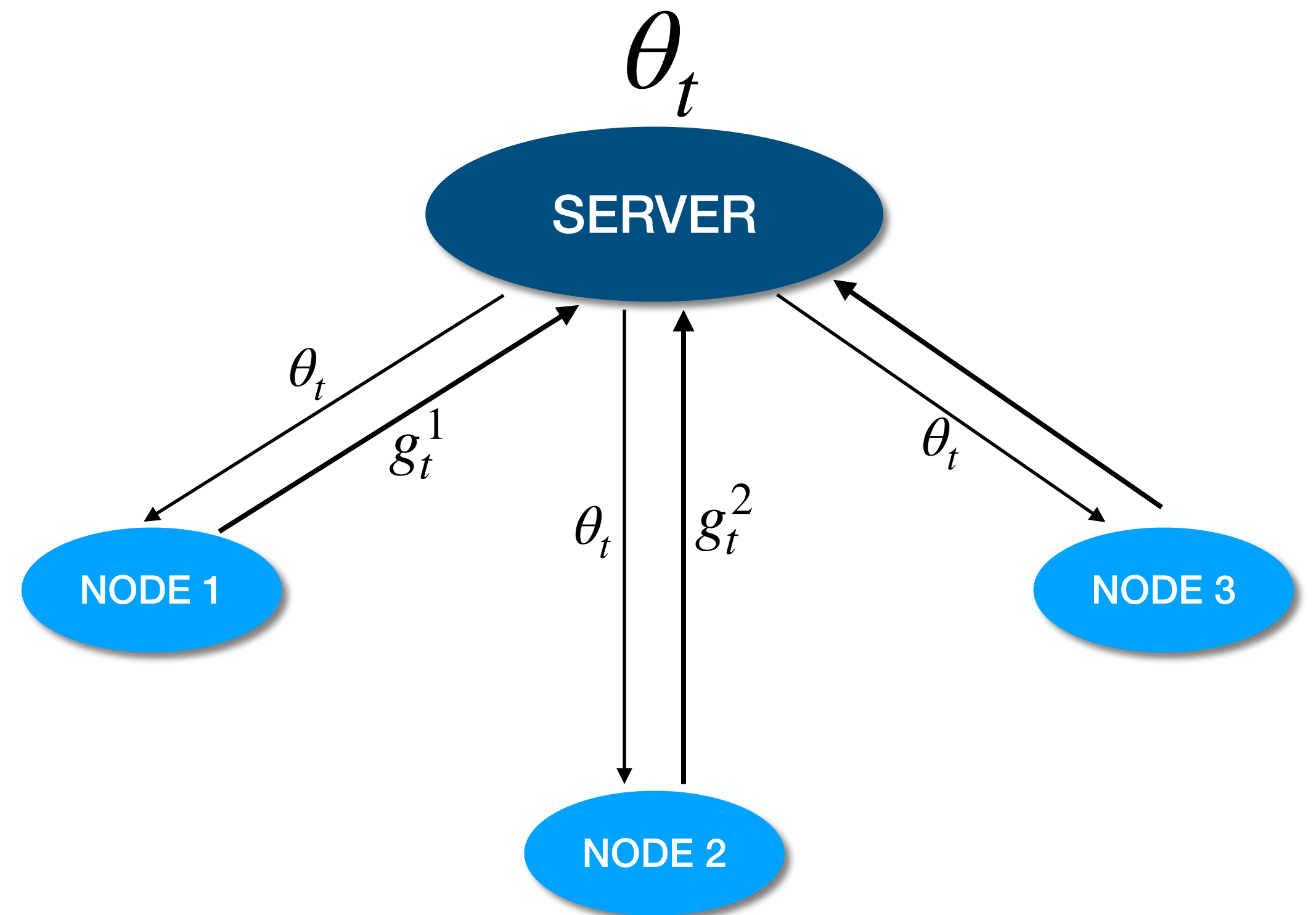
# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

True gradient

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\displaystyle\sum_i g_t^i$

$\theta_{t+1}$

SERVER

$\theta_t$

$g_t^1$

$\theta_t$ $g_t^2$

$g_t^3$

$\theta_t$

NODE 1

NODE 2
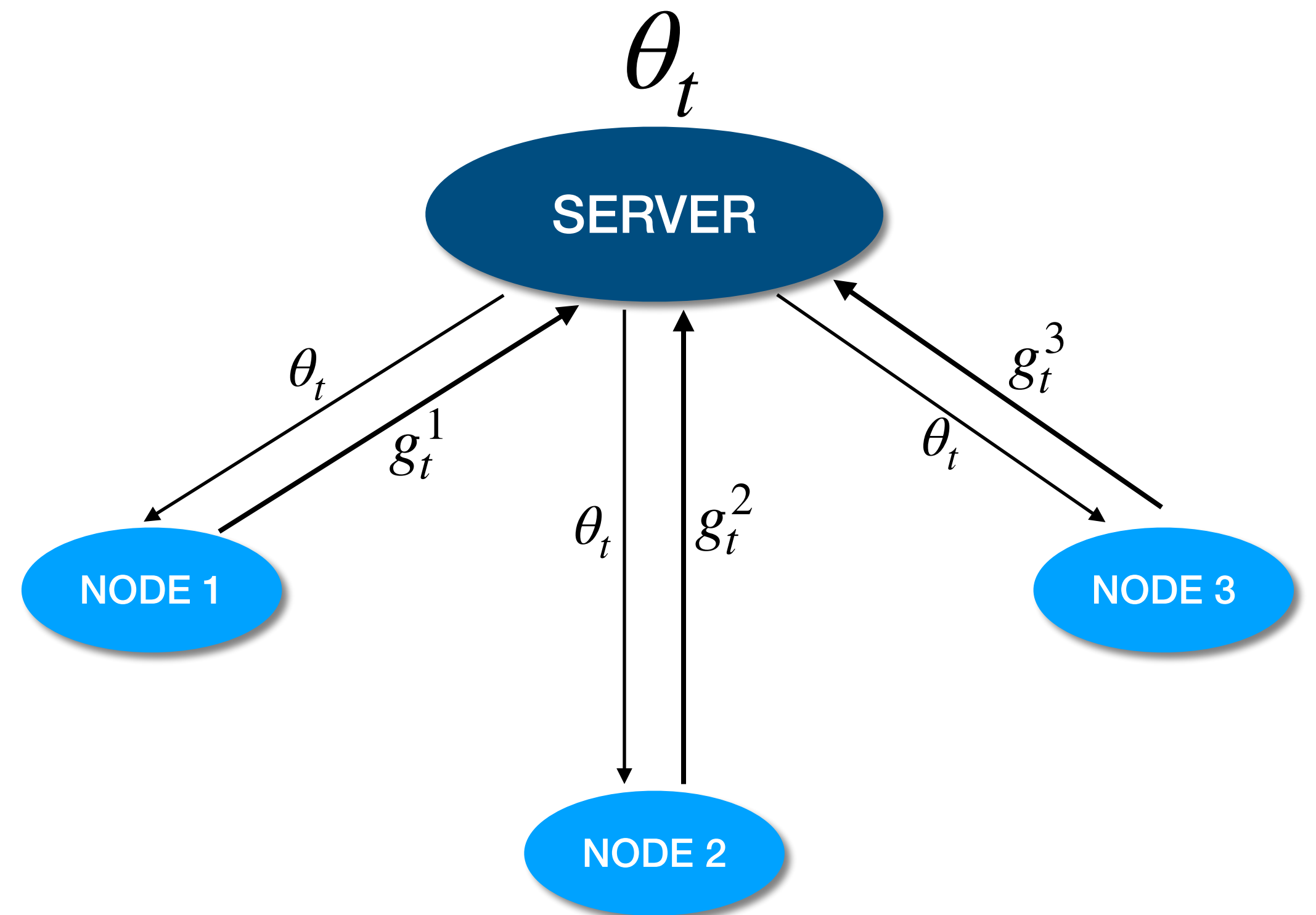
NODE 3

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

True gradient

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \, \widehat{g}_t$

# Distributed SGD

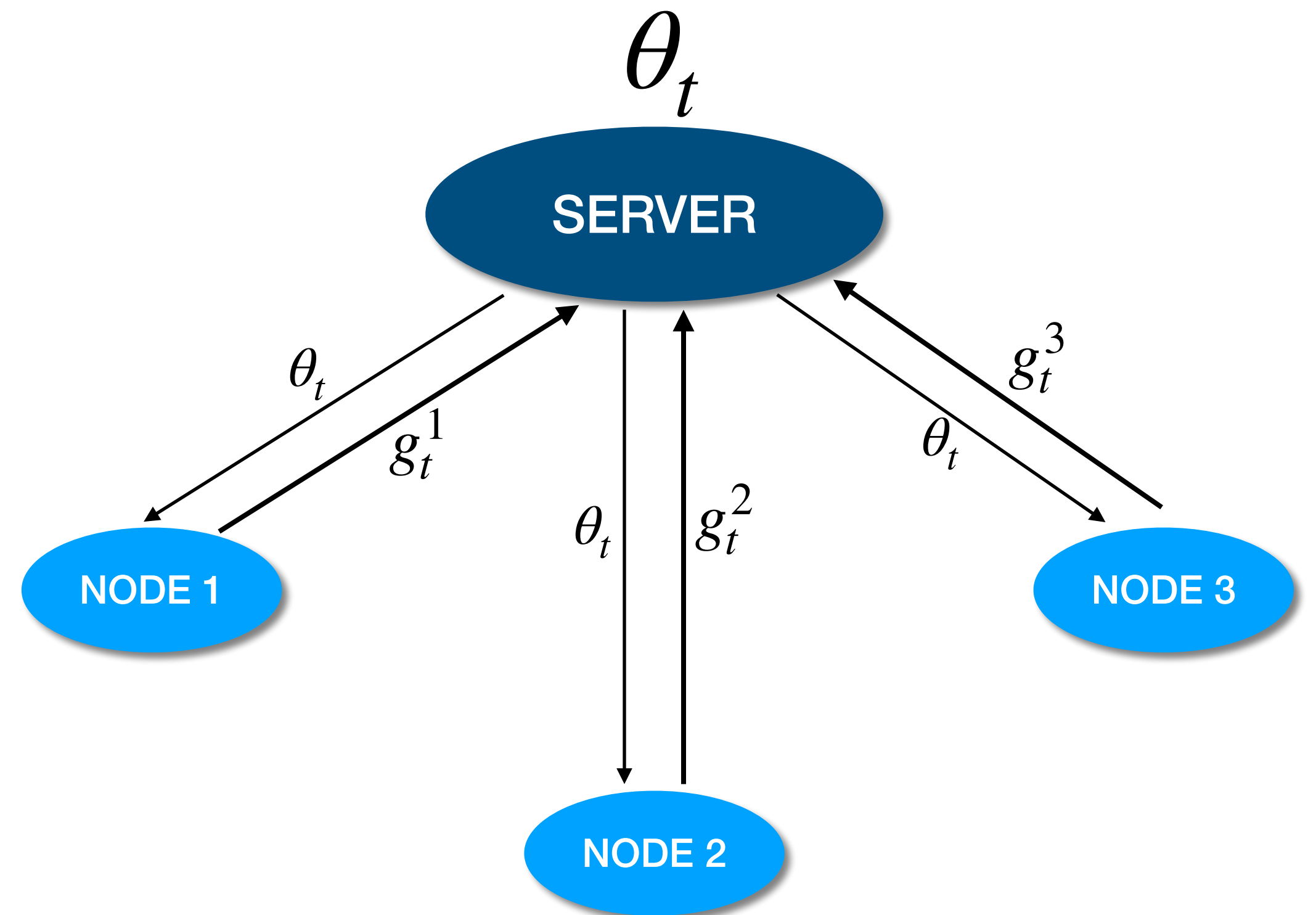**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

True gradient

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t\,\widehat{g}_t$

$\theta_{t+1}$

SERVER

$\theta_t$

$g_t^1$

$\theta_t$  $g_t^2$

$g_t^3$

$\theta_t$

NODE 1

NODE 3

NODE 2

# Distributed SGD

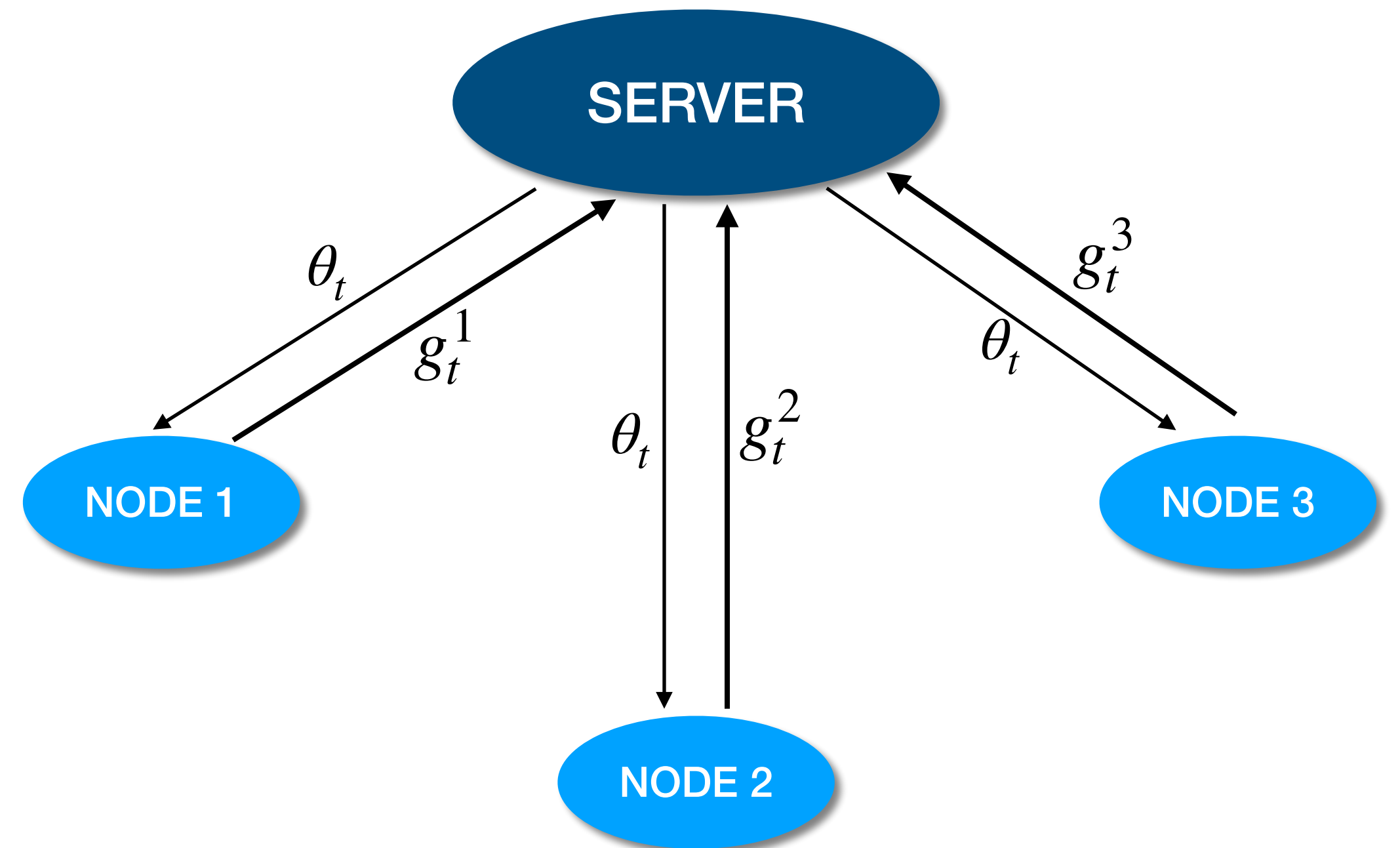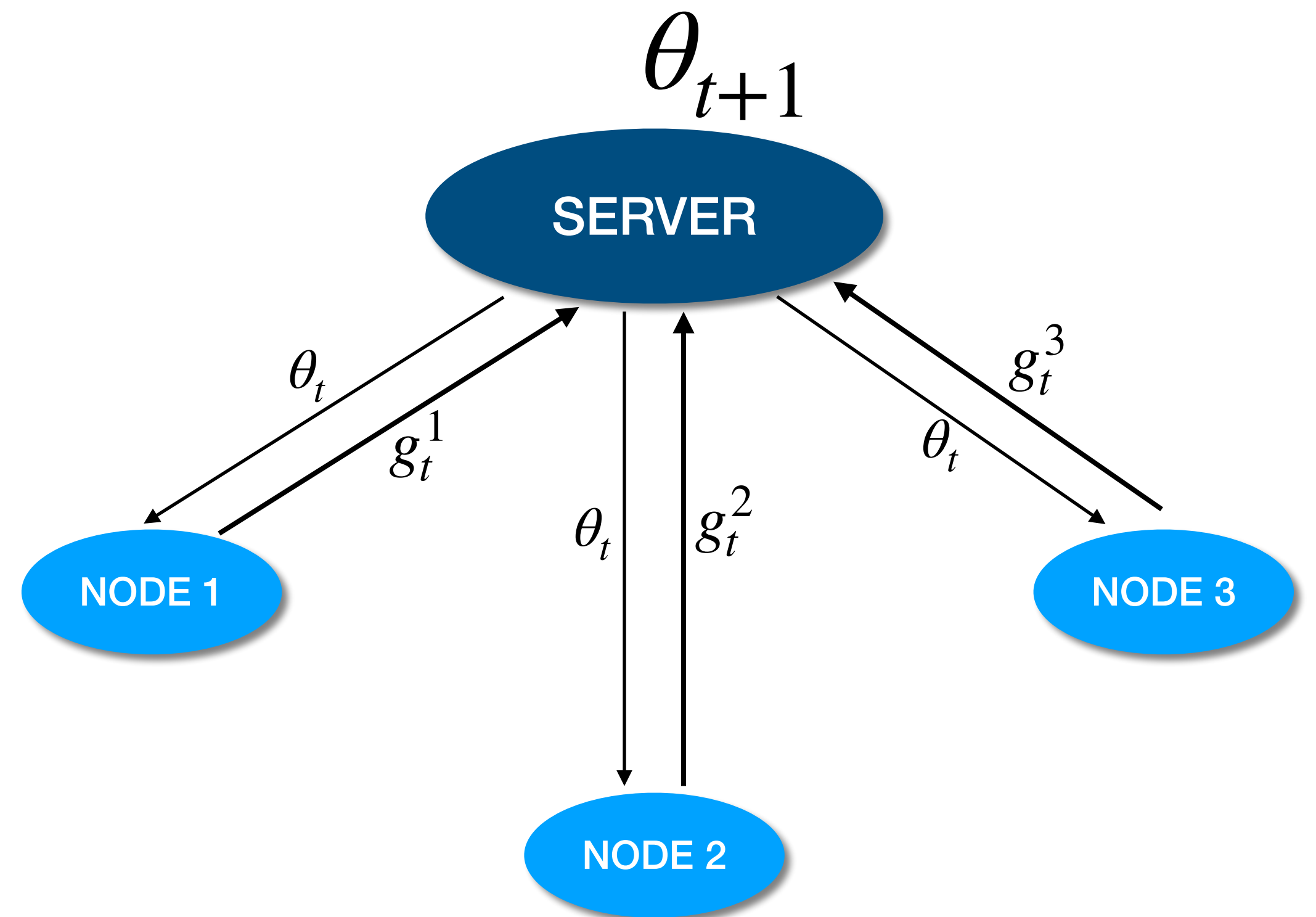**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathscr{U}_t$$

$$\mathbb{E}\left[\mathscr{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathscr{U}_t\|^2\right] \leq \sigma^2 < \infty$$

True gradient

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t\,\widehat{g}_t$

$\theta_{t+1}$

SERVER

$\theta_t$        $g_t^1$        $\theta_t$   $g_t^2$        $g_t^3$   $\theta_t$

NODE 1

NODE 3

NODE 2

# Distributed SGD

**Divides the workload per machine by the total size of the system**

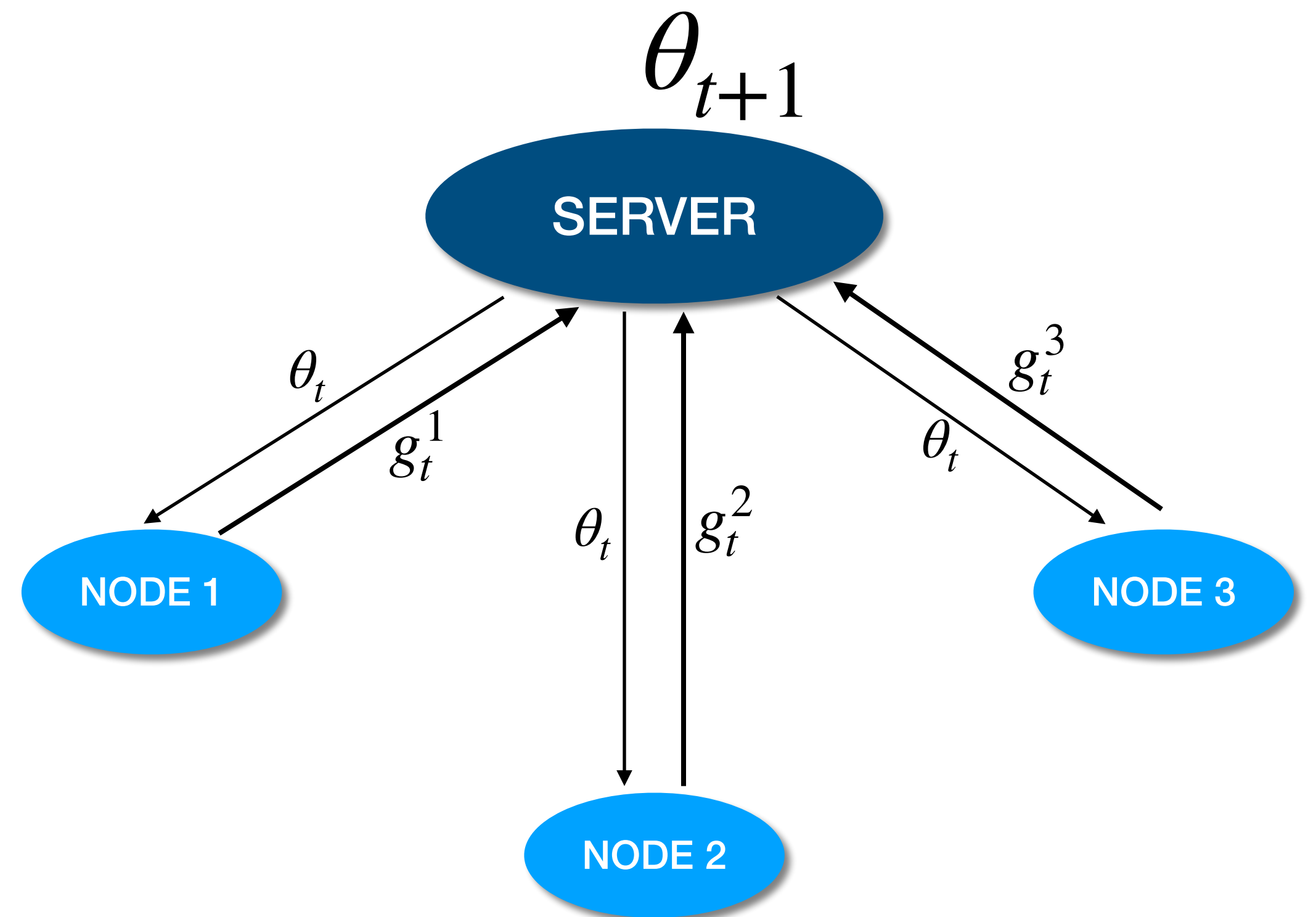**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

True gradient

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \widehat{g}_t$

$\theta_{t+1}$

SERVER

NODE 1

NODE 2

NODE 3

$\theta_t$

$g_t^1$

$\theta_t$

$g_t^2$

$g_t^3$

$\theta_t$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

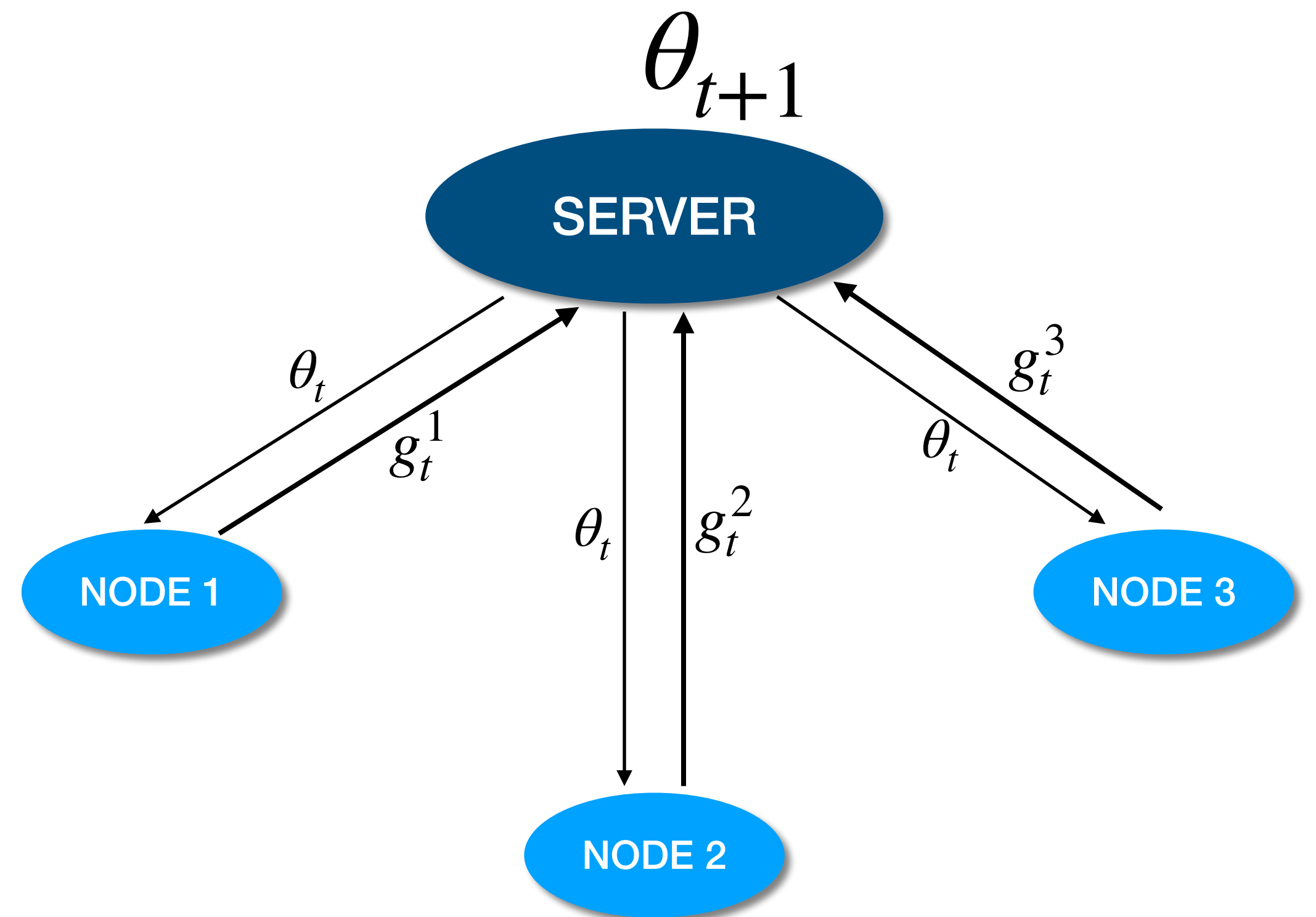**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathscr{U}_t$$

True gradient

$$\mathbb{E}\left[\mathscr{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathscr{U}_t\|^2\right] \leq \sigma^2 < \infty$$

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\displaystyle\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t\,\widehat{g}_t$

$\theta_{t+1}$

SERVER

$\theta_t$

$g_t^1$

$g_t^3$

$\theta_t$

$\theta_t$

$g_t^2$

NODE 1

NODE 3

NODE 2

# Distributed SGD

**Divides the workload per machine by the total size of the system**

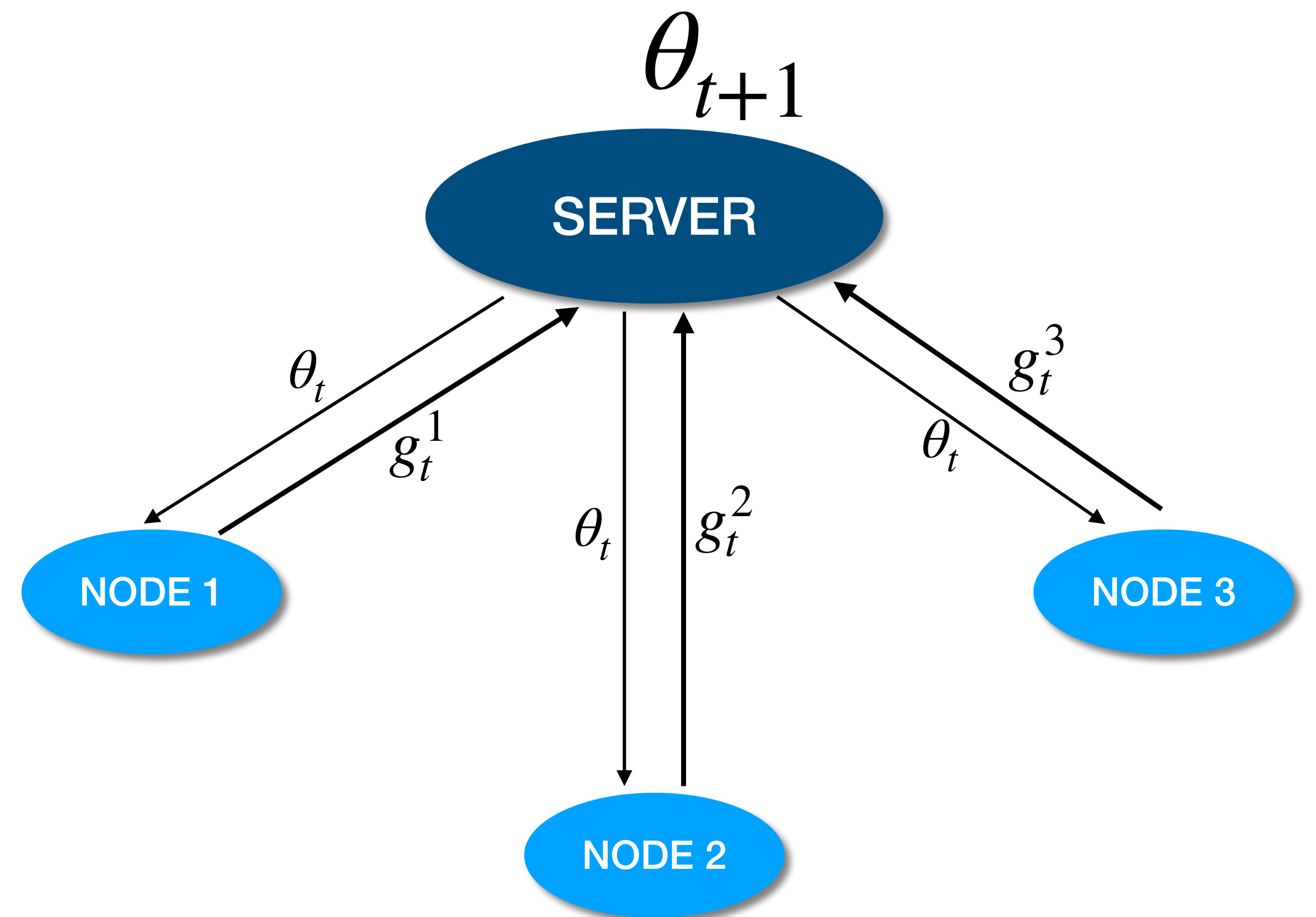**Nodes query** *stochastic gradients* with bounded variance

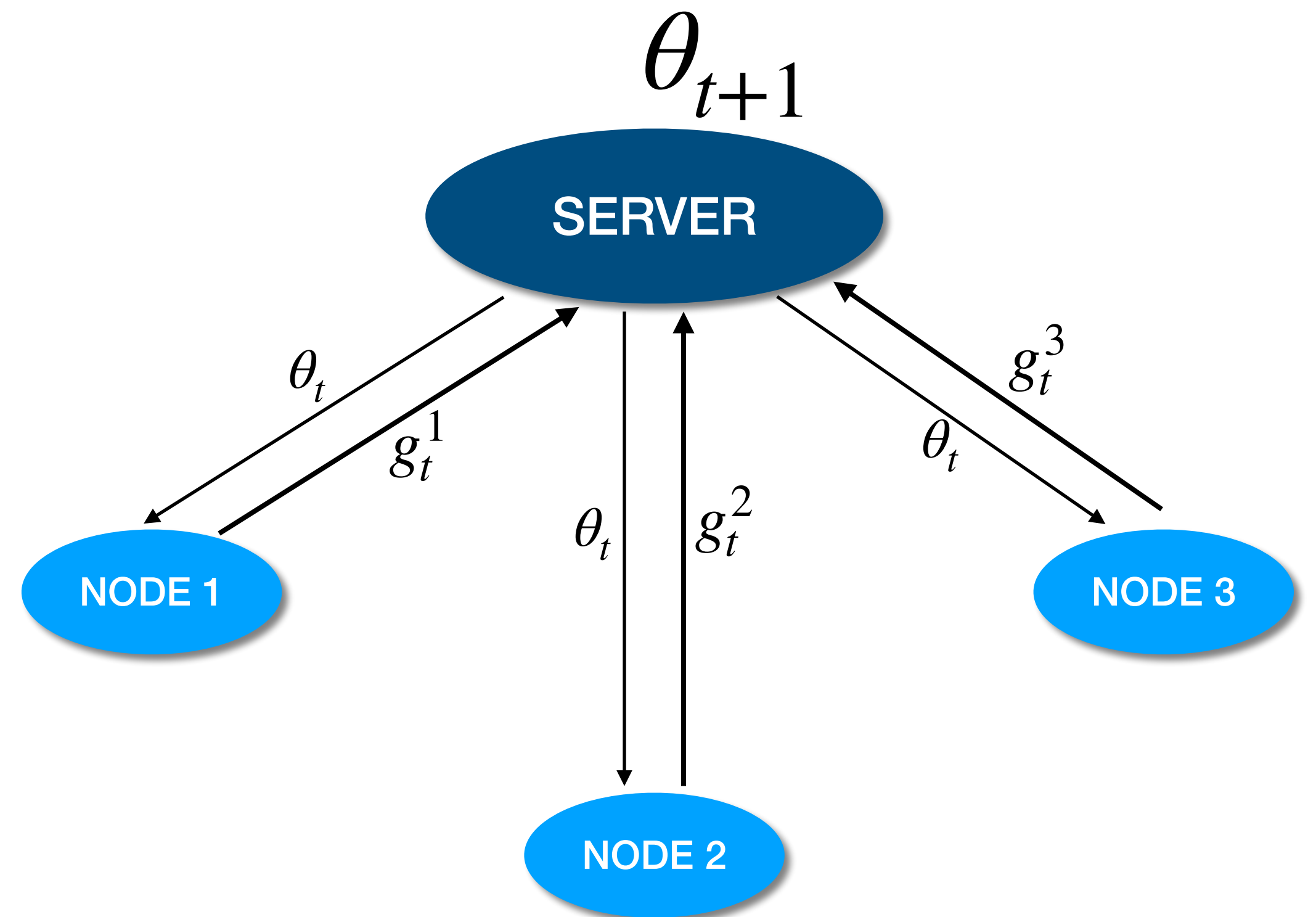$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

True gradient

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t\,\widehat{g}_t$

$\theta_{t+1}$

SERVER

$\theta_t$  $g_t^1$  $g_t^3$  $\theta_t$

NODE 1  $\theta_t$  $g_t^2$  NODE 3

NODE 2

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \boxed{\nabla Q(\theta_t)} + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

True gradient

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t\, \widehat{g}_t$

$\theta_{t+1}$

SERVER

$\theta_t$

$g_t^1$

$\theta_t$

$g_t^2$

$g_t^3$

$\theta_t$

NODE 1

NODE 2

NODE 3

$\theta_1$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

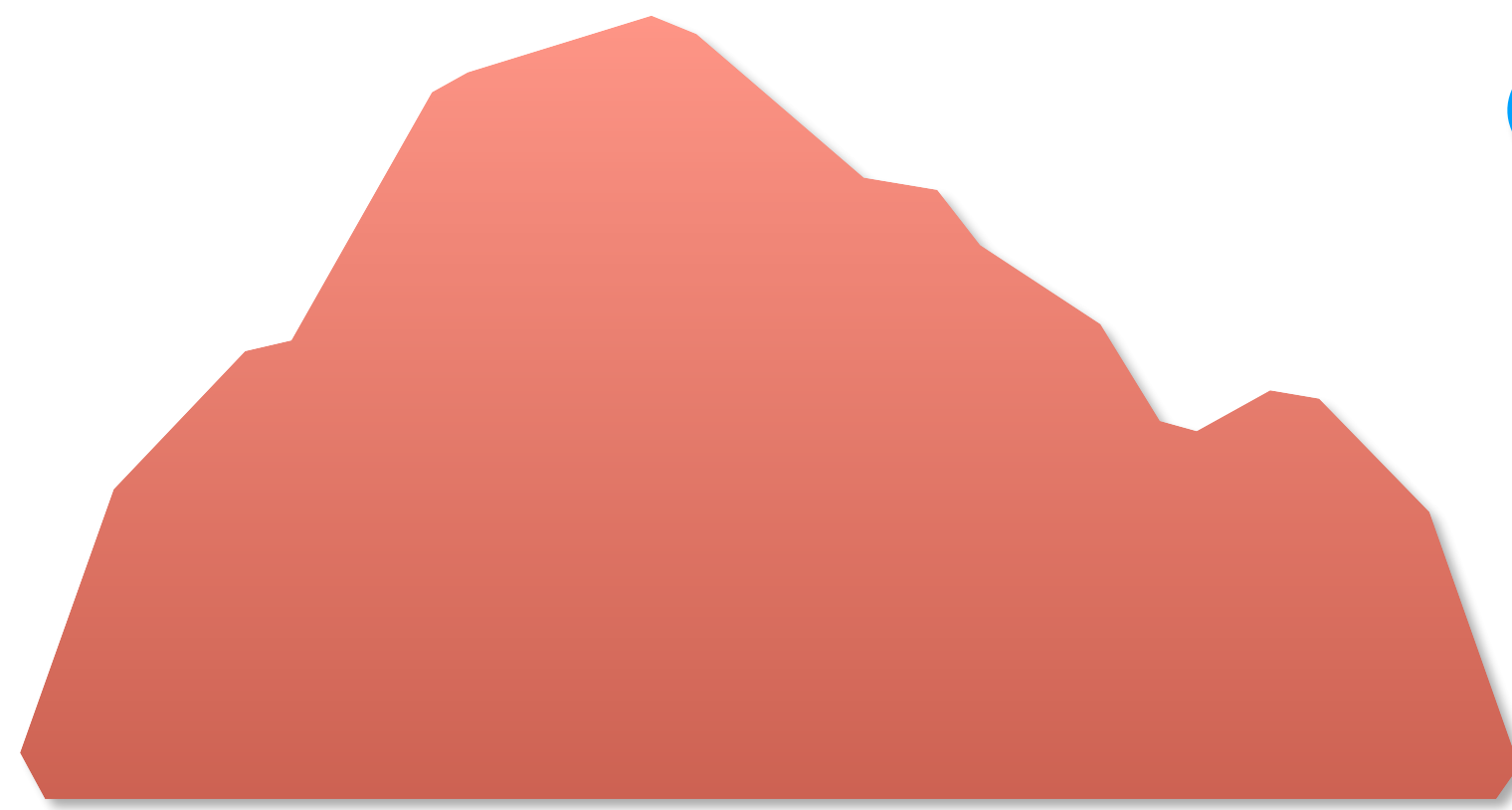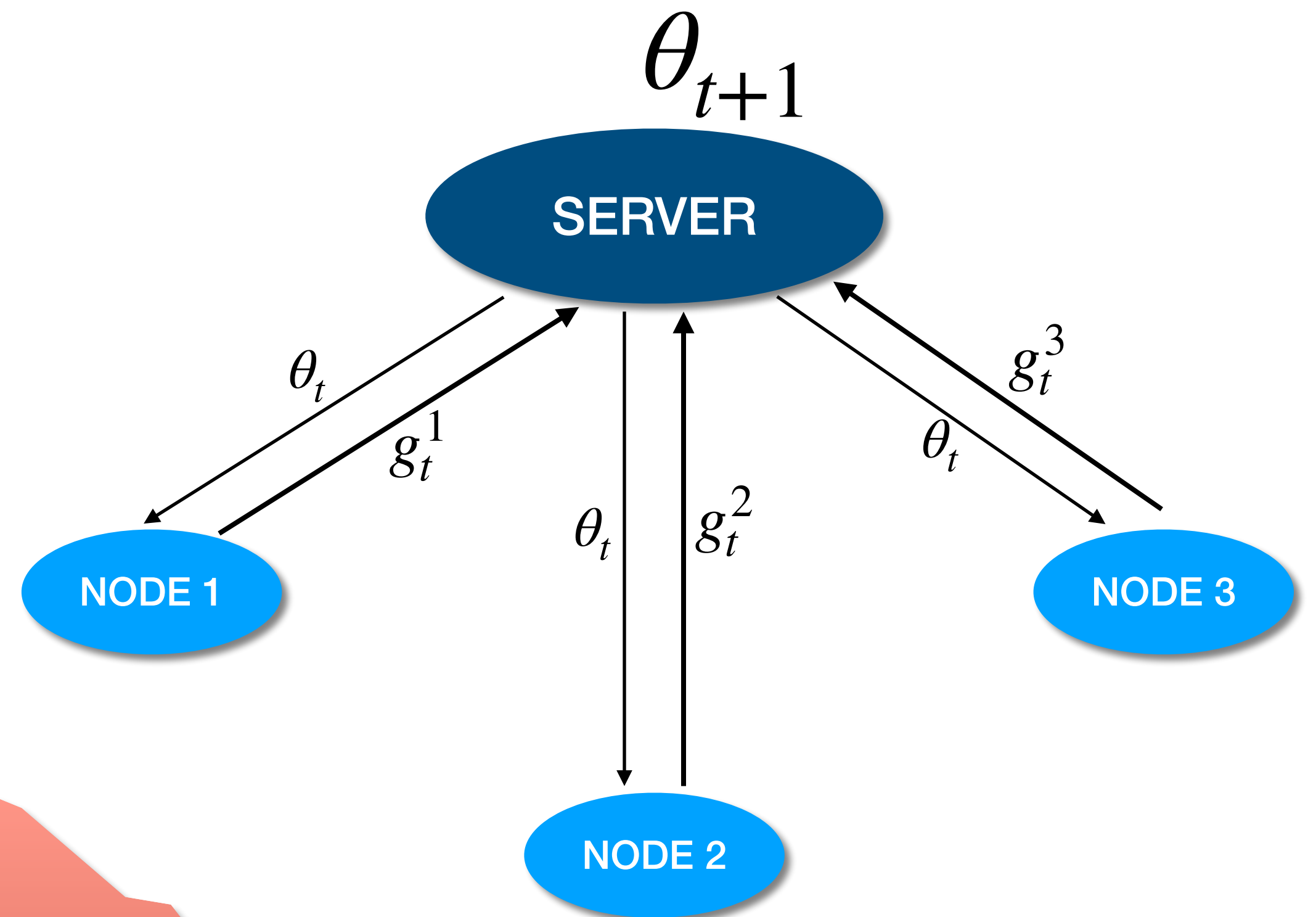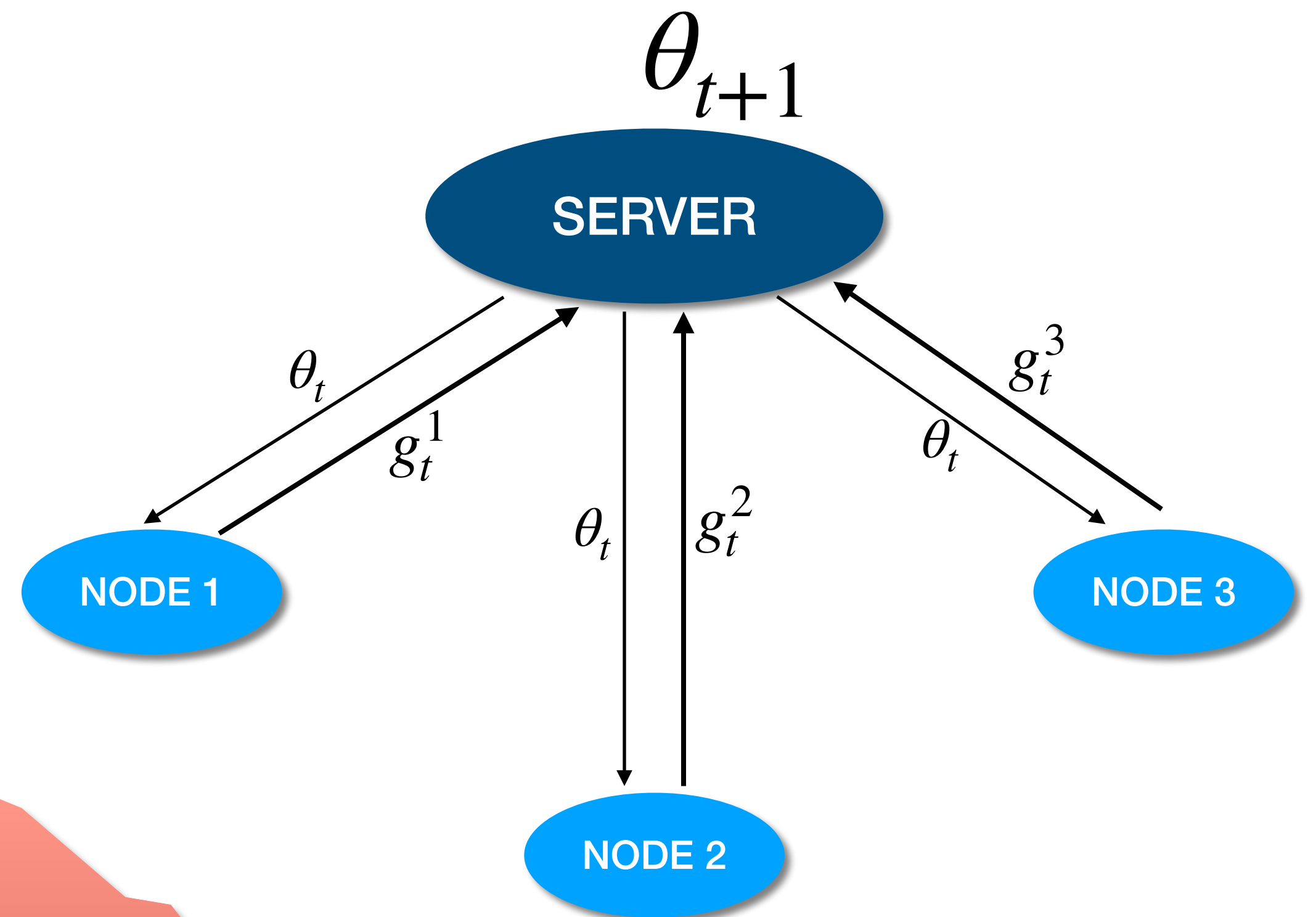$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

True gradient

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

**Server averages** the gradients, $\hat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \, \hat{g}_t$

$\theta_{t+1}$

SERVER

$\theta_t$   $g_t^1$   $\theta_t$   $g_t^2$   $g_t^3$   $\theta_t$

NODE 1

NODE 2

NODE 3

$\theta_1$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

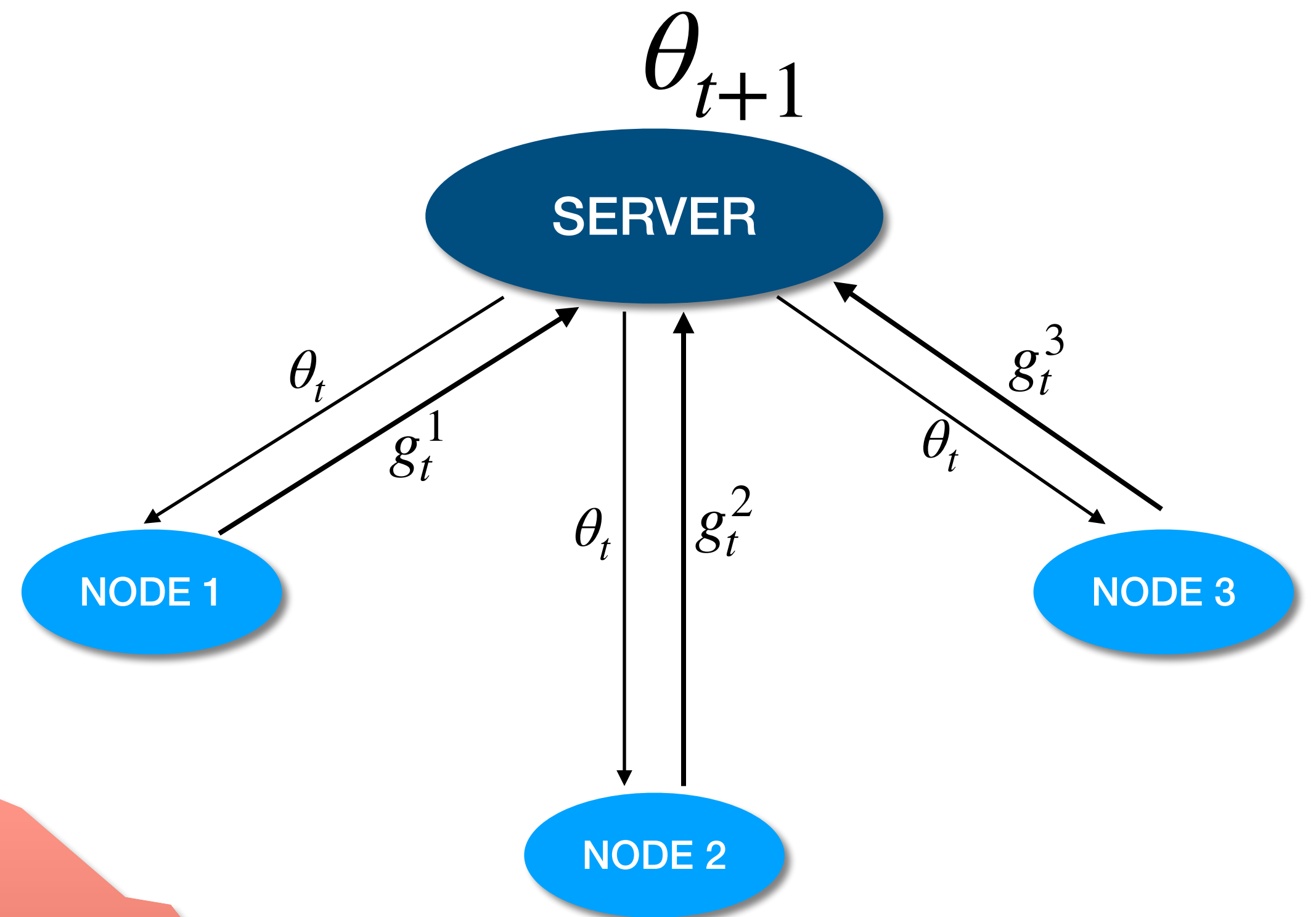$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

True gradient

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \, \widehat{g}_t$

$\theta_{t+1}$

SERVER

$\theta_t$   $g_t^1$    $\theta_t$   $g_t^2$    $g_t^3$   $\theta_t$

NODE 1    NODE 3

NODE 2

$\theta_1$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

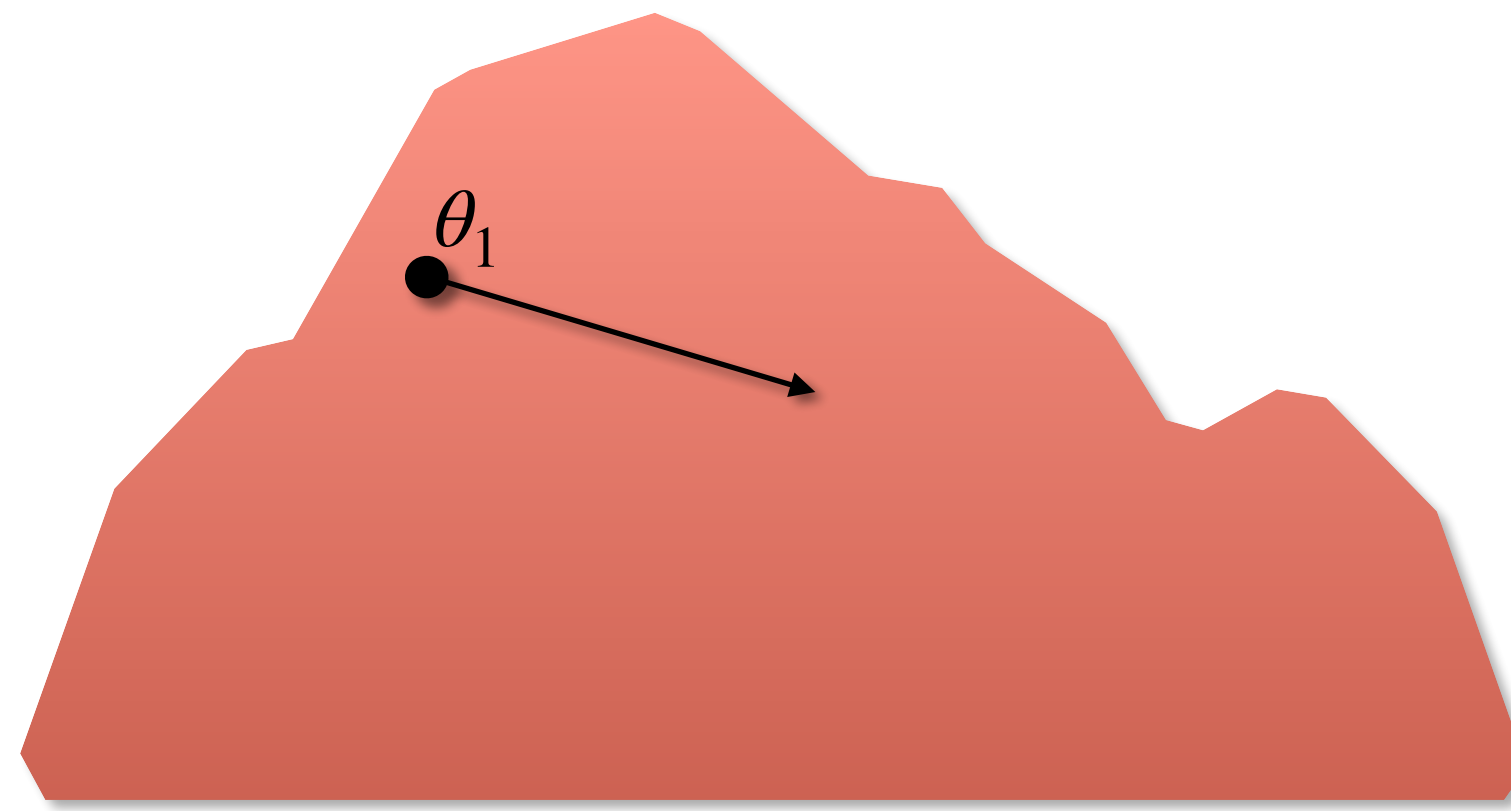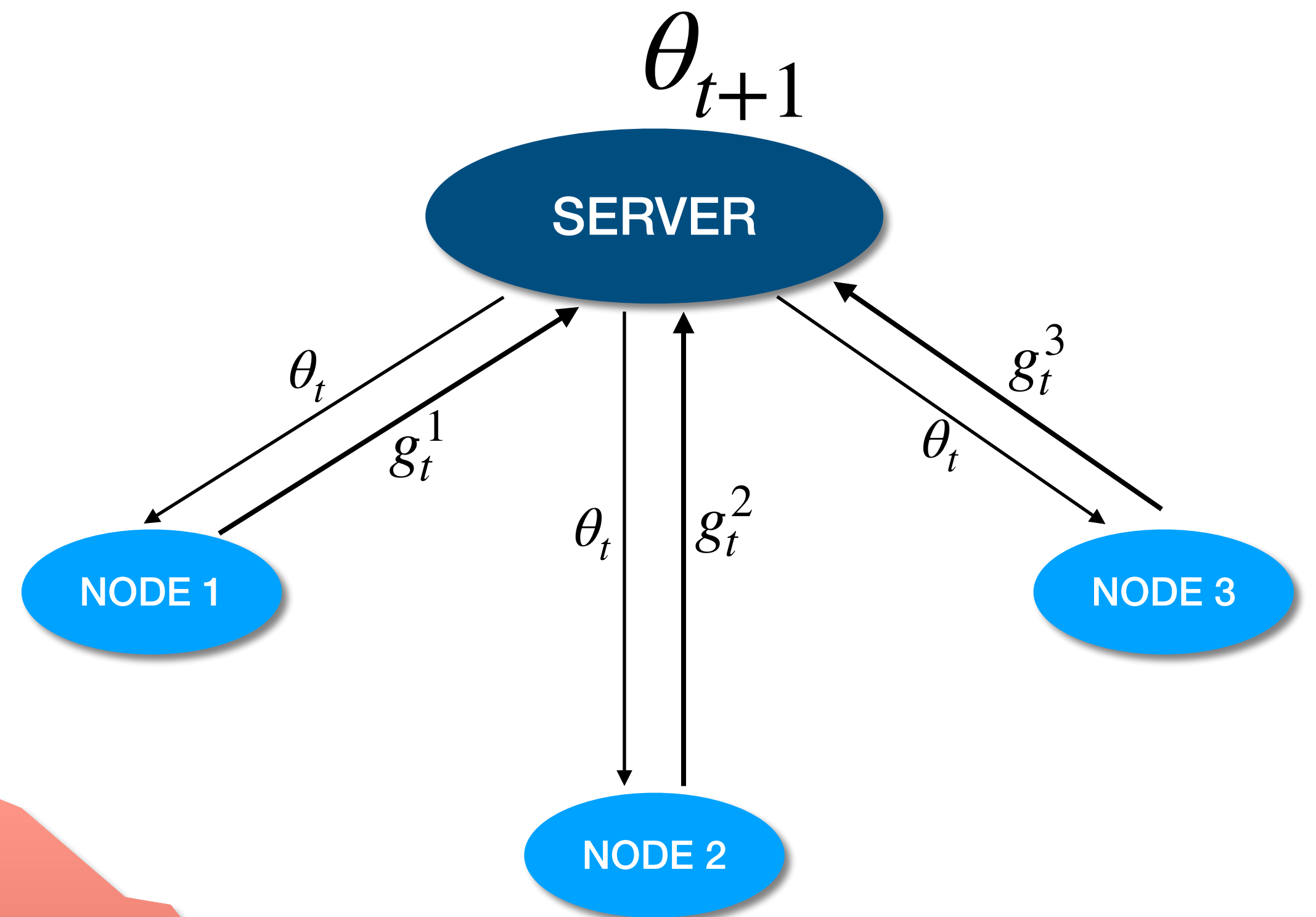$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

True gradient

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \widehat{g}_t$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$
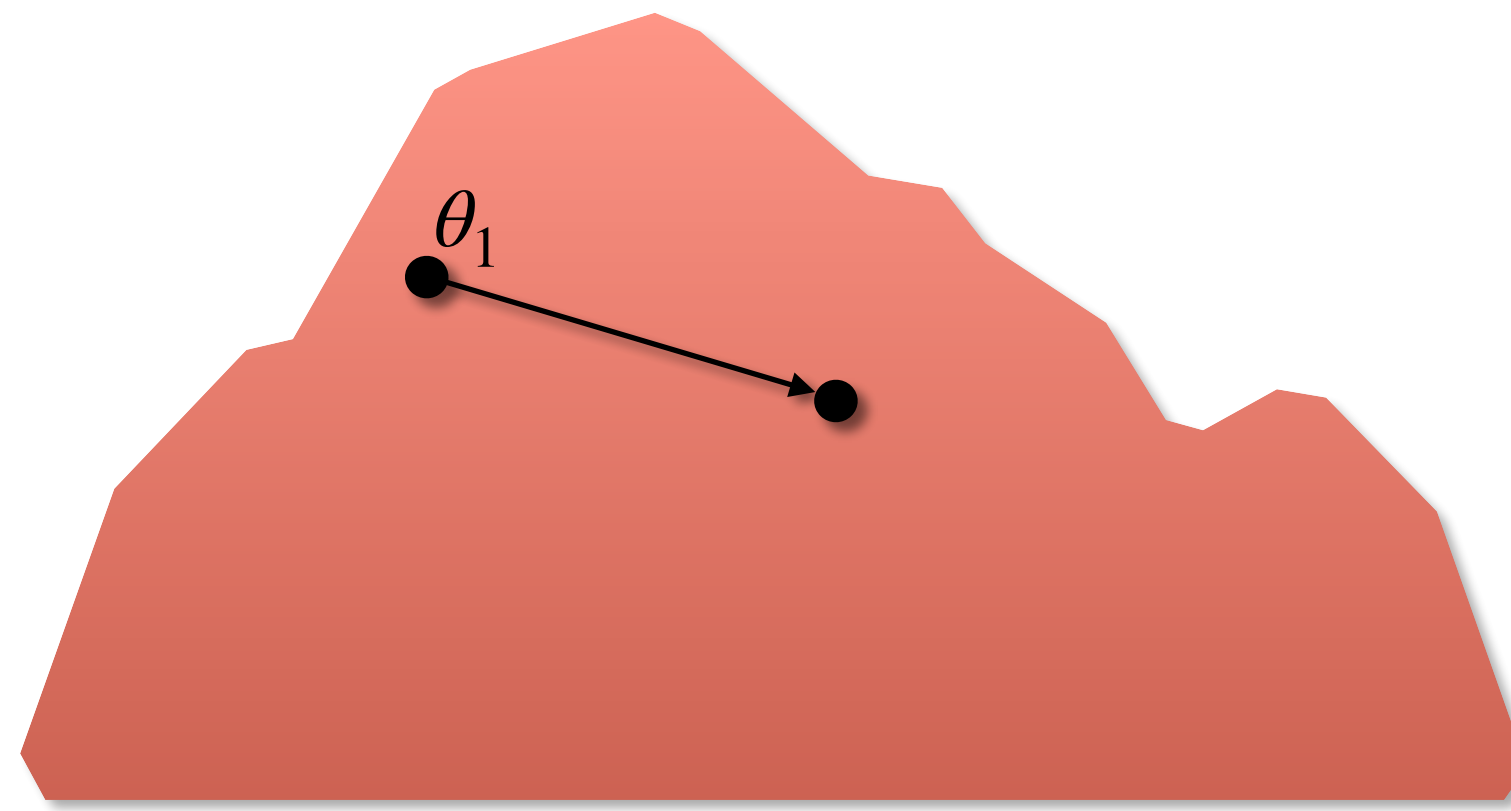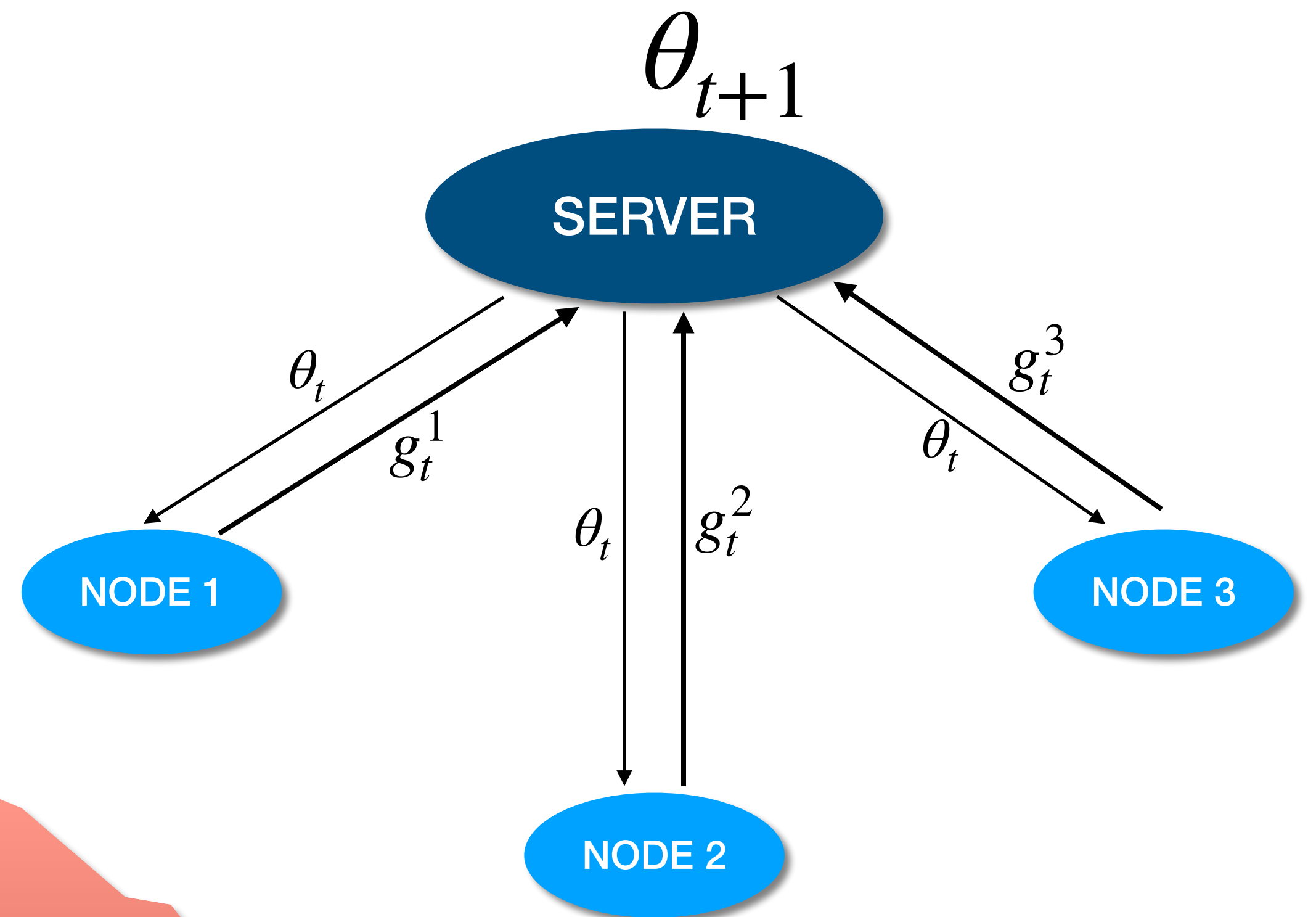
$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$
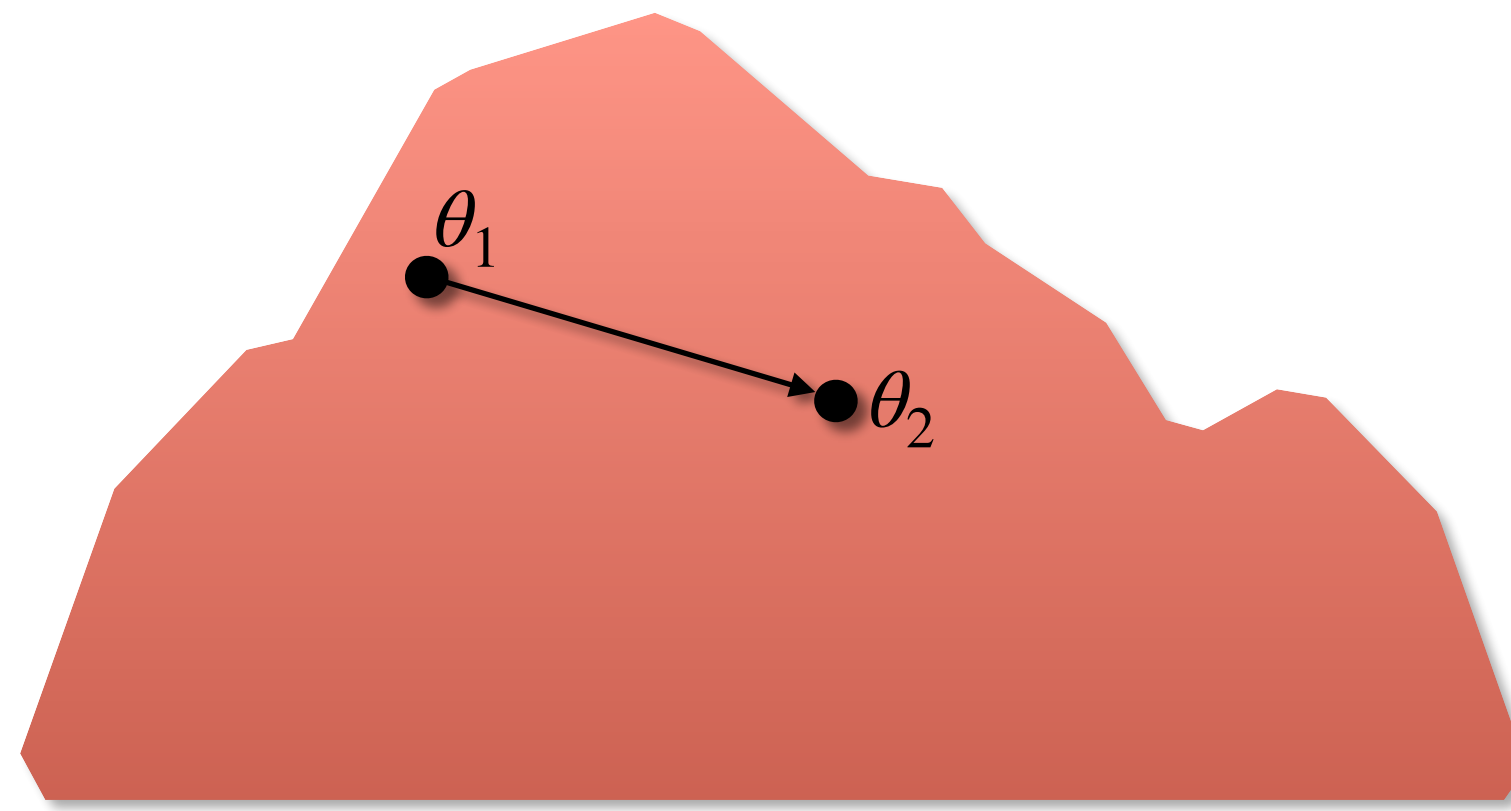
True gradient

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \leftarrow \theta_t - \gamma_t \, \widehat{g}_t$

$\theta_{t+1}$

SERVER

NODE 1

NODE 2

NODE 3

$\theta_t$

$g_t^1$

$\theta_t$

$g_t^2$

$g_t^3$

$\theta_t$

$\theta_1$

$\theta_2$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

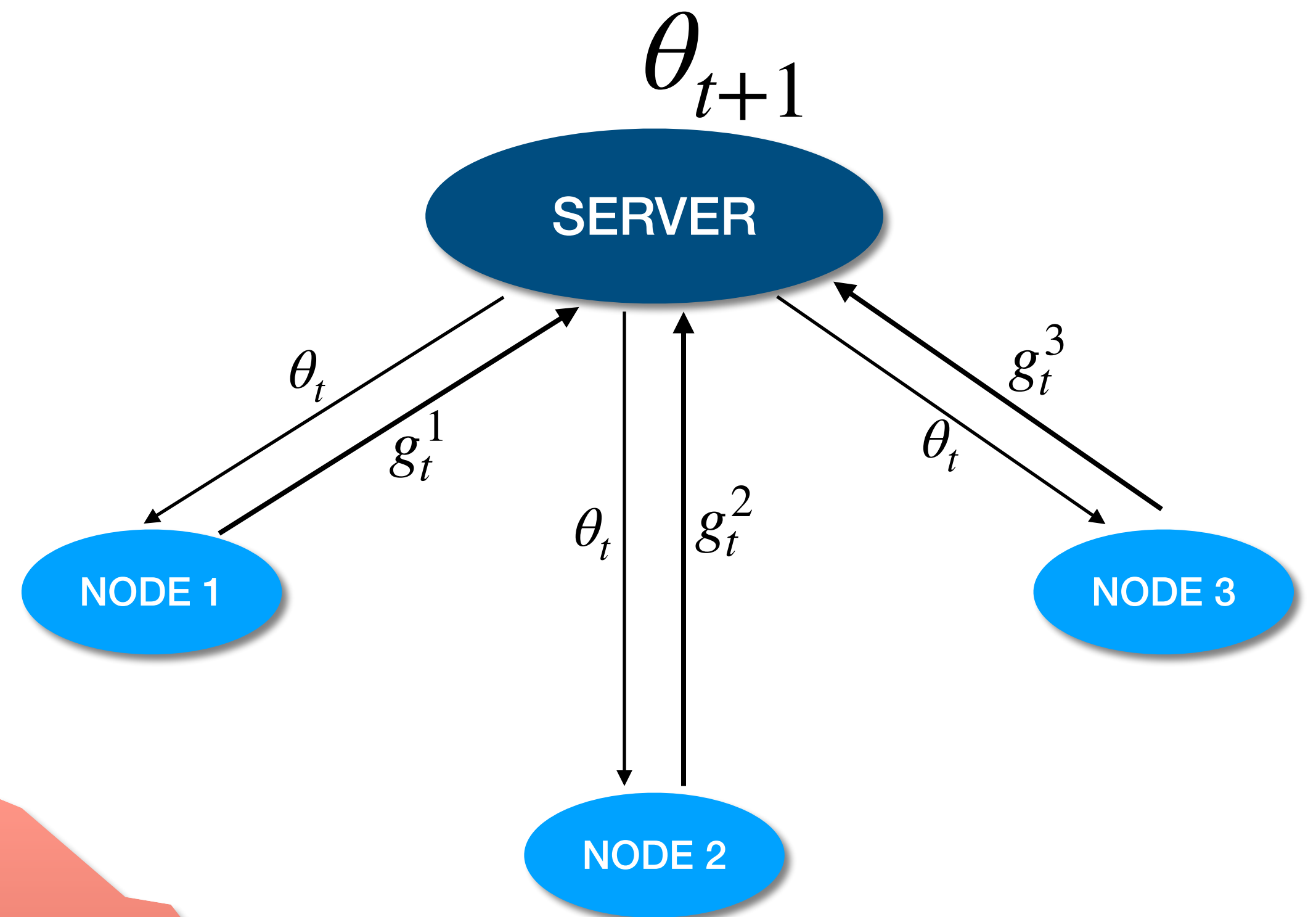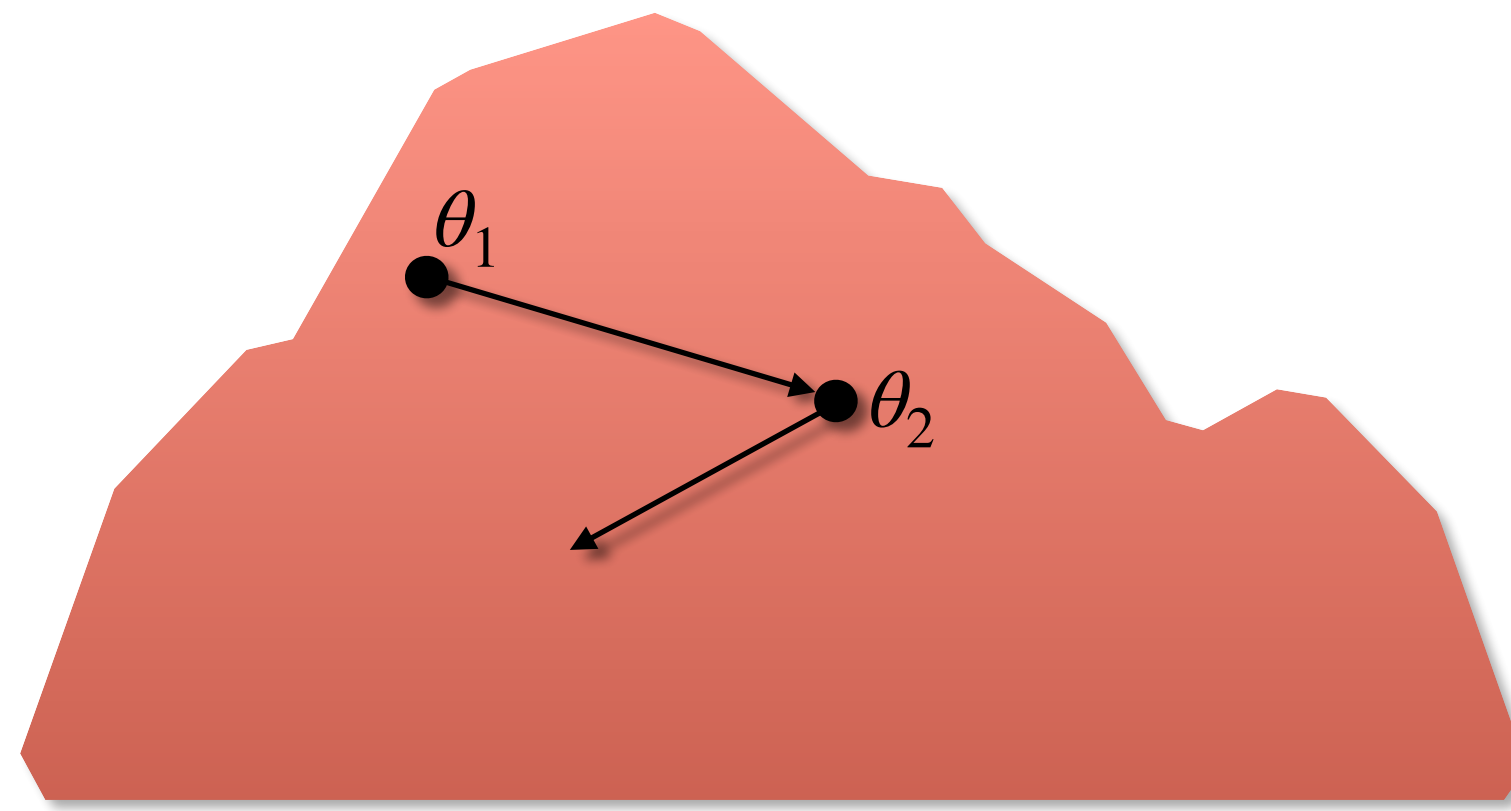$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \le \sigma^2 < \infty$$

True gradient

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \, \widehat{g}_t$

$\theta_{t+1}$

SERVER

$\theta_t$   $g_t^1$   $\theta_t$   $g_t^2$   $g_t^3$   $\theta_t$

NODE 1     NODE 3

NODE 2

$\theta_1$

$\theta_2$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

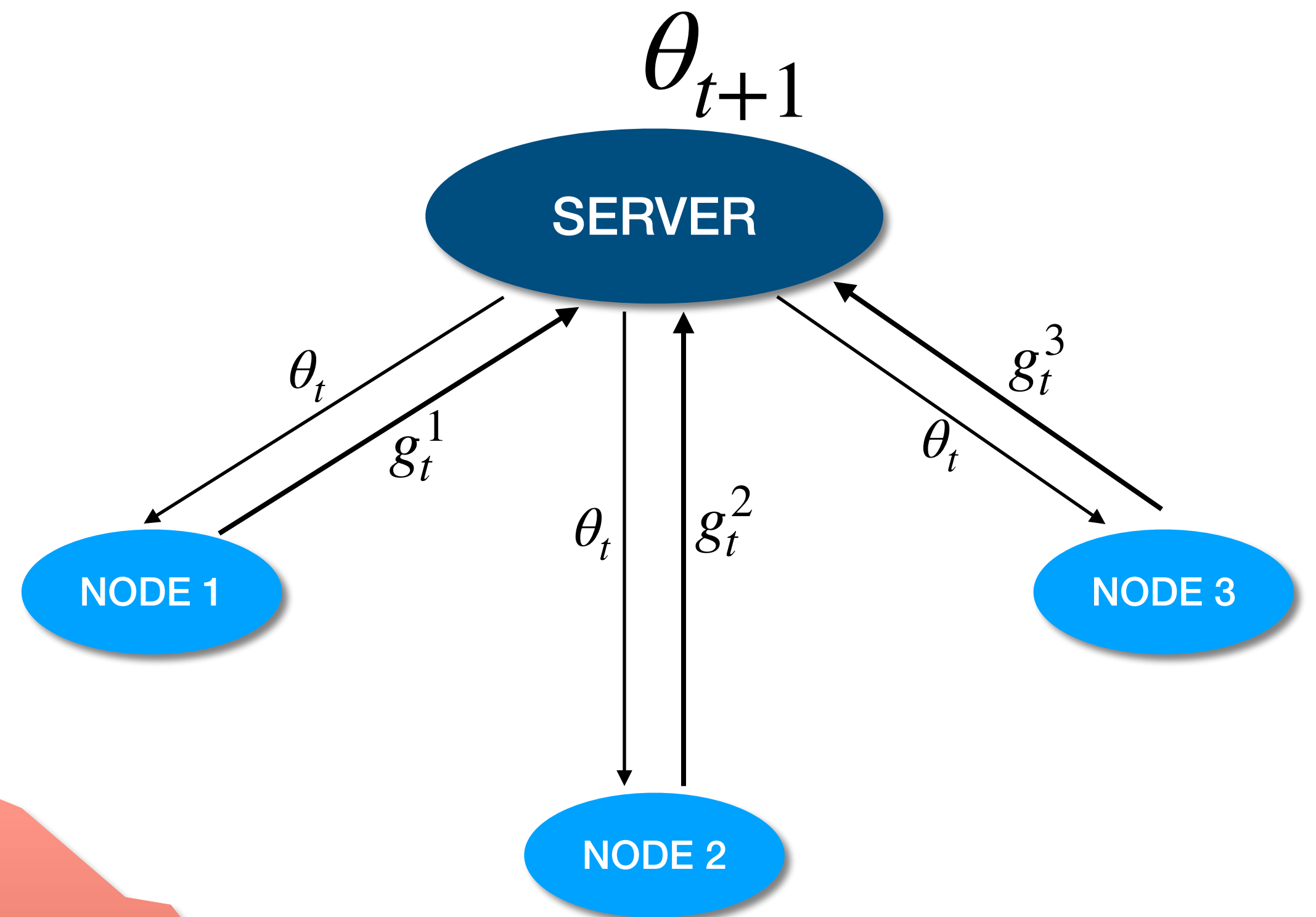$$g_t^i = \boxed{\nabla Q(\theta_t)} + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

True gradient

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\displaystyle\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t\, \widehat{g}_t$

$\theta_{t+1}$

SERVER

NODE 1

NODE 2

NODE 3

$\theta_t$ $g_t^1$ $\theta_t$ $g_t^2$ $g_t^3$ $\theta_t$

$\theta_1$ $\theta_2$ $\theta_3$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

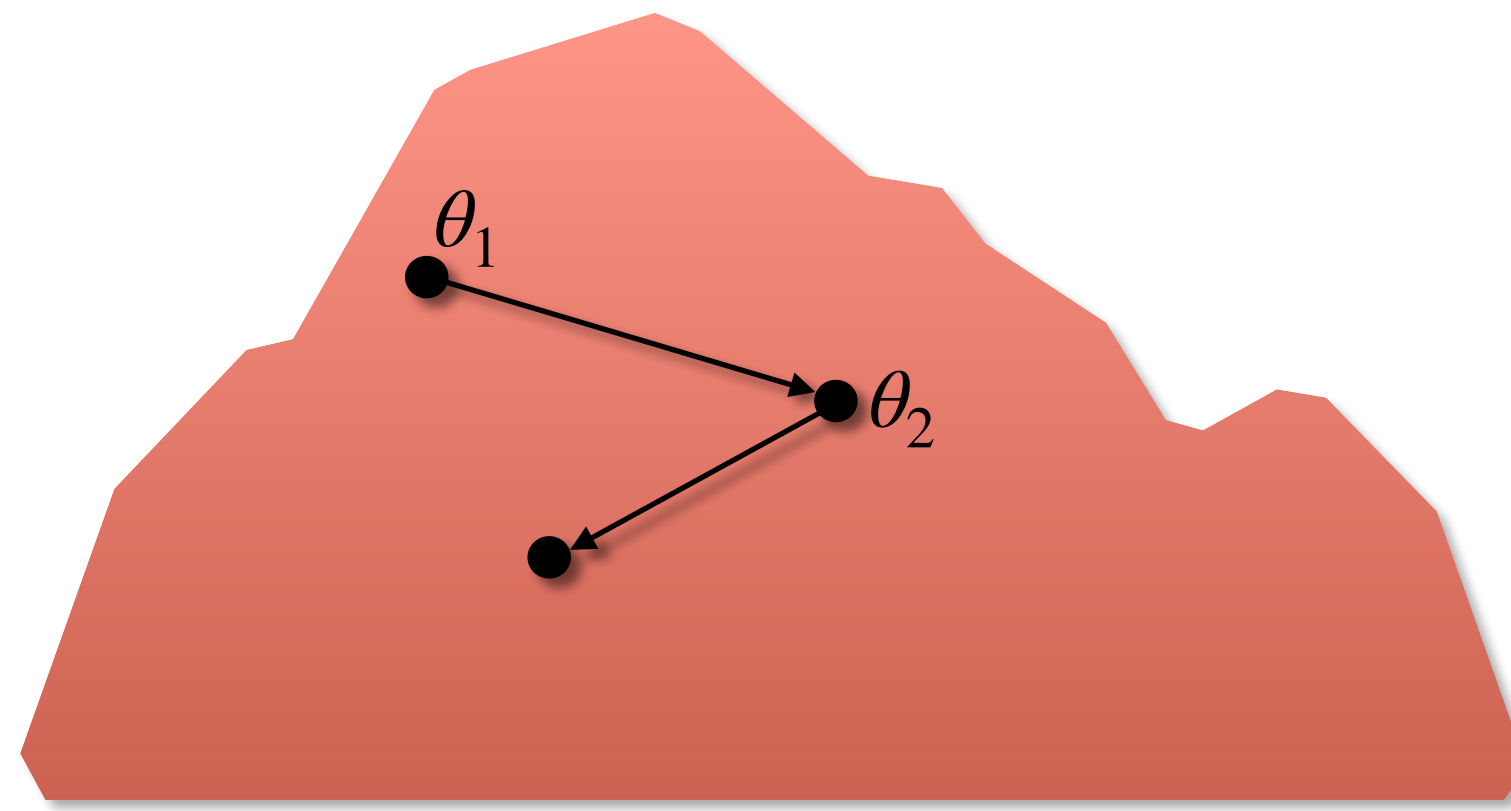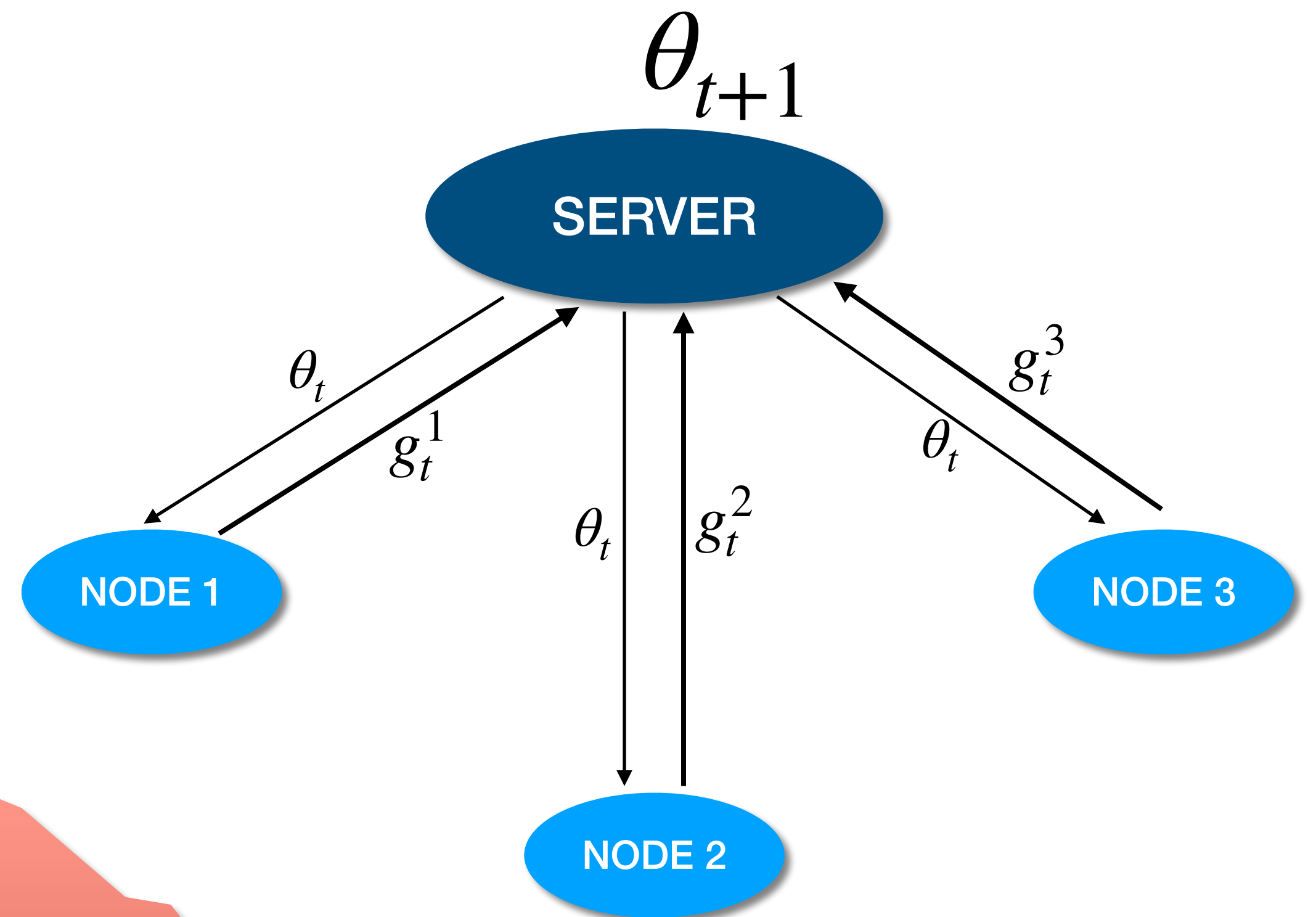$$g_t^i = \boxed{\nabla Q(\theta_t)} + u_t^i \quad ; \quad u_t^i \sim \mathscr{U}_t$$

$$\mathbb{E}\left[\mathscr{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathscr{U}_t\|^2\right] \le \sigma^2 < \infty$$

True gradient

**Server averages** the gradients, $\quad \widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t\, \widehat{g}_t$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

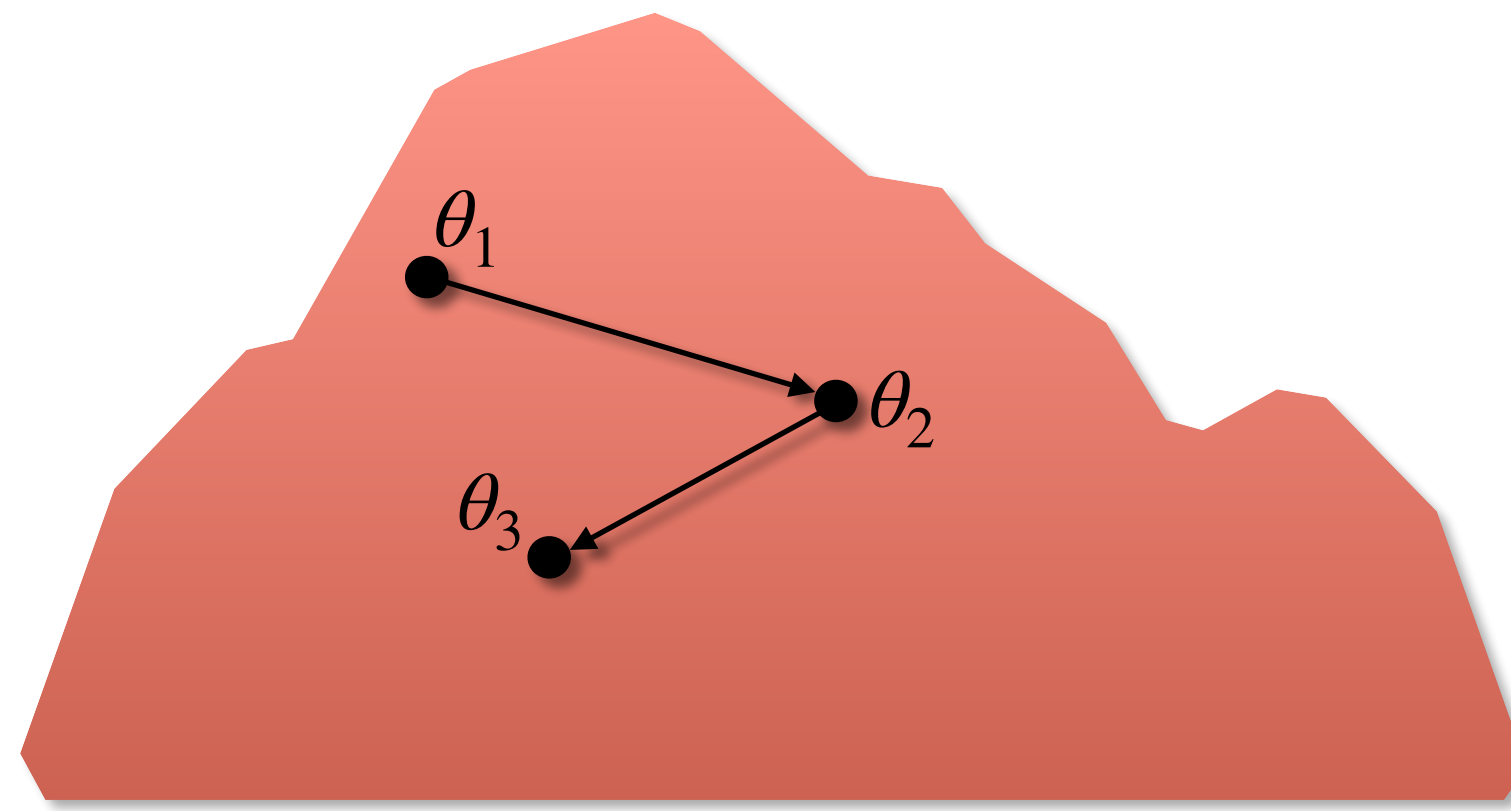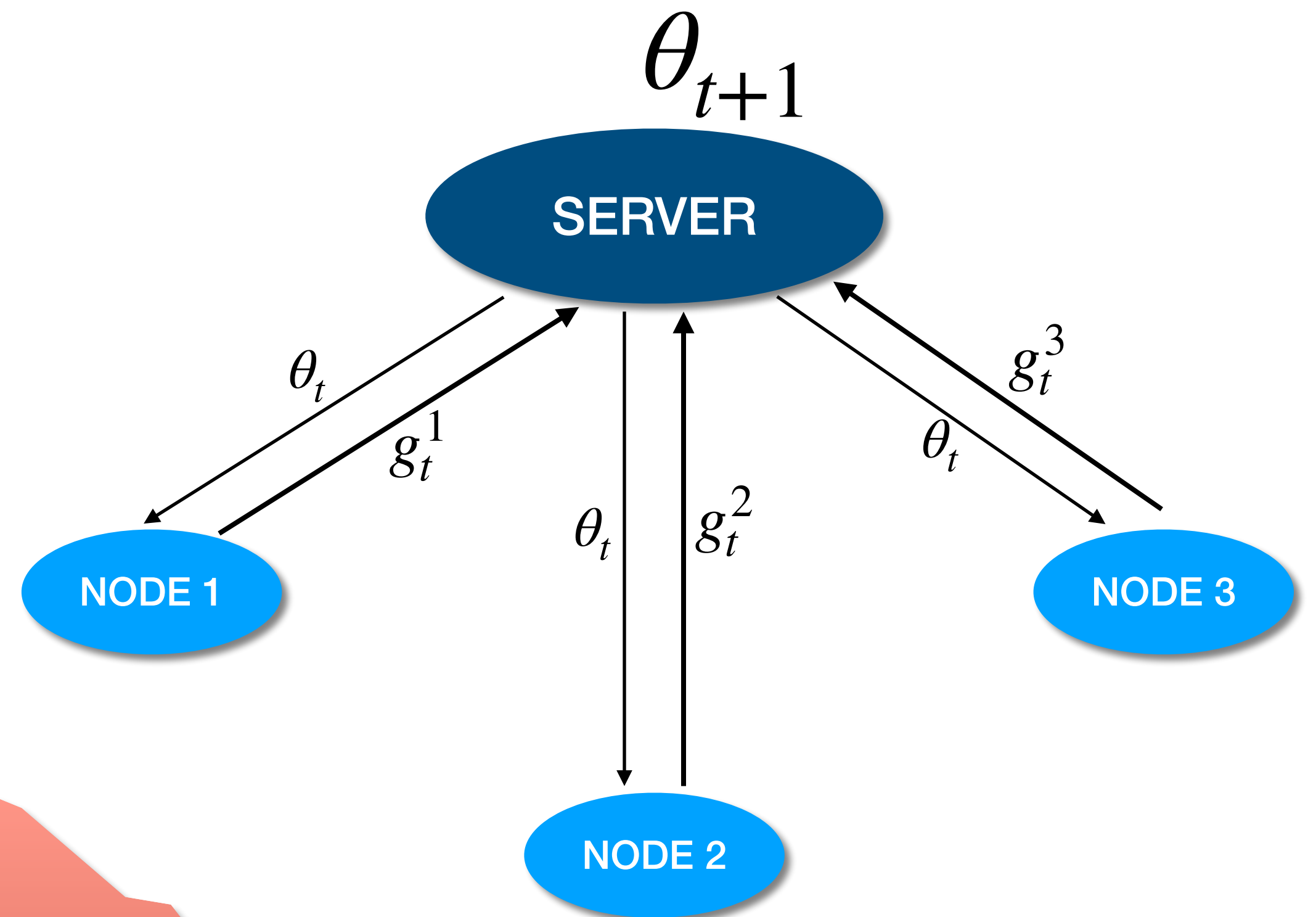$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathscr{U}_t$$

$$\mathbb{E}\left[\mathscr{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathscr{U}_t\|^2\right] \leq \sigma^2 < \infty$$

True gradient

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \, \widehat{g}_t$

$$\theta_{t+1}$$

SERVER

$\theta_t$

$g_t^1$

$\theta_t$   $g_t^2$

$g_t^3$

$\theta_t$

NODE 1

NODE 2

NODE 3

$\theta_1$

$\theta_2$

$\theta_3$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

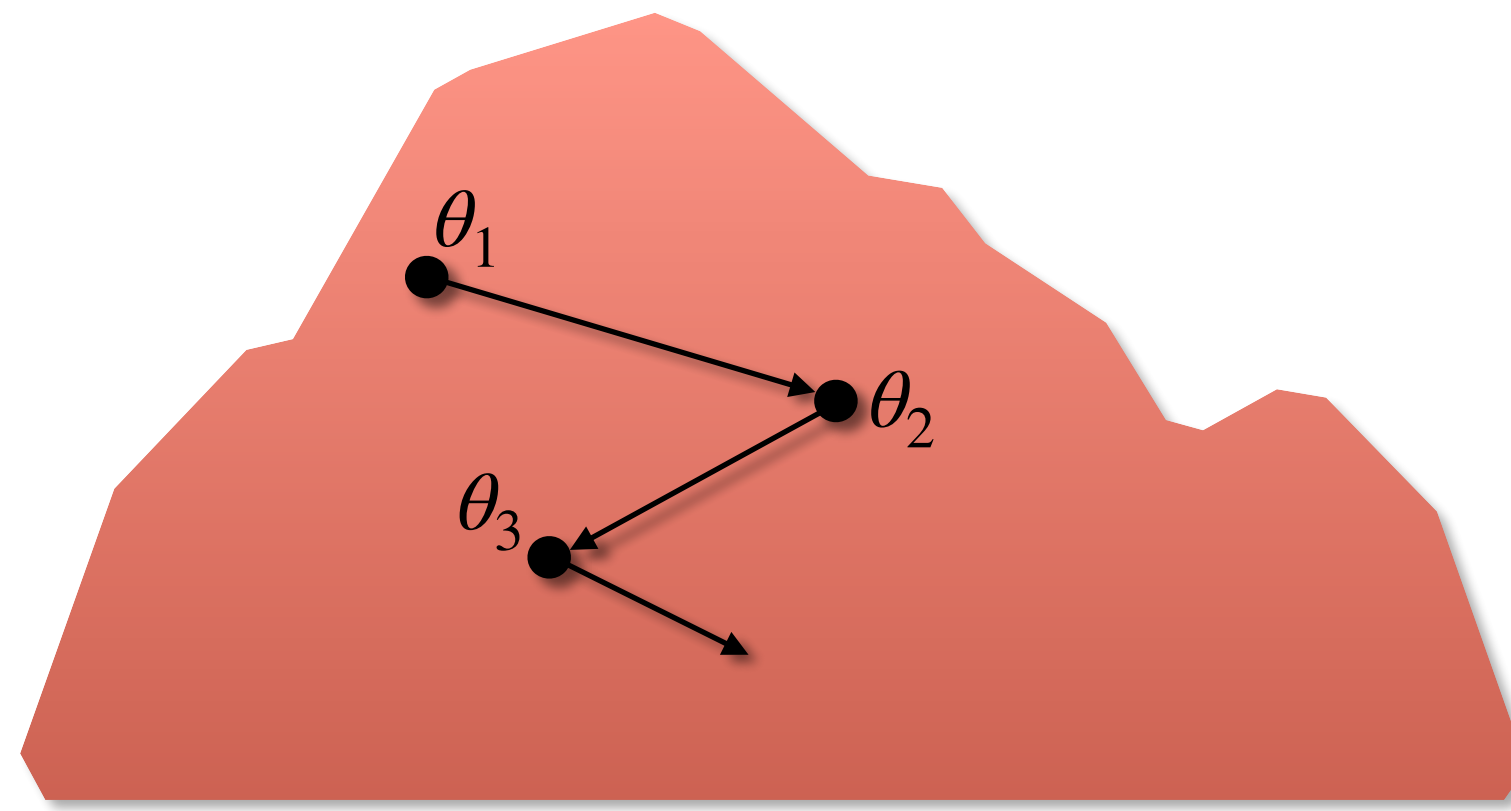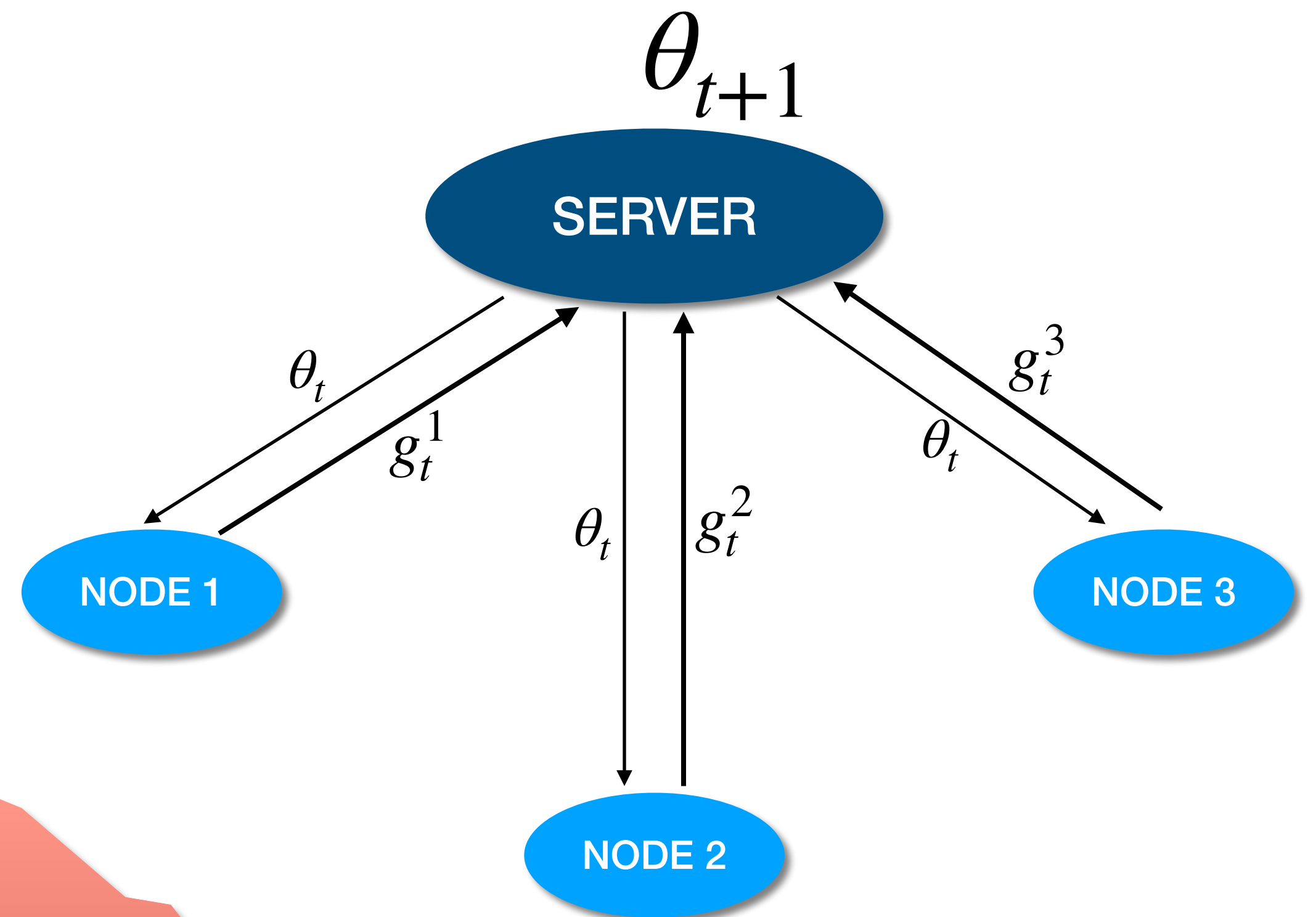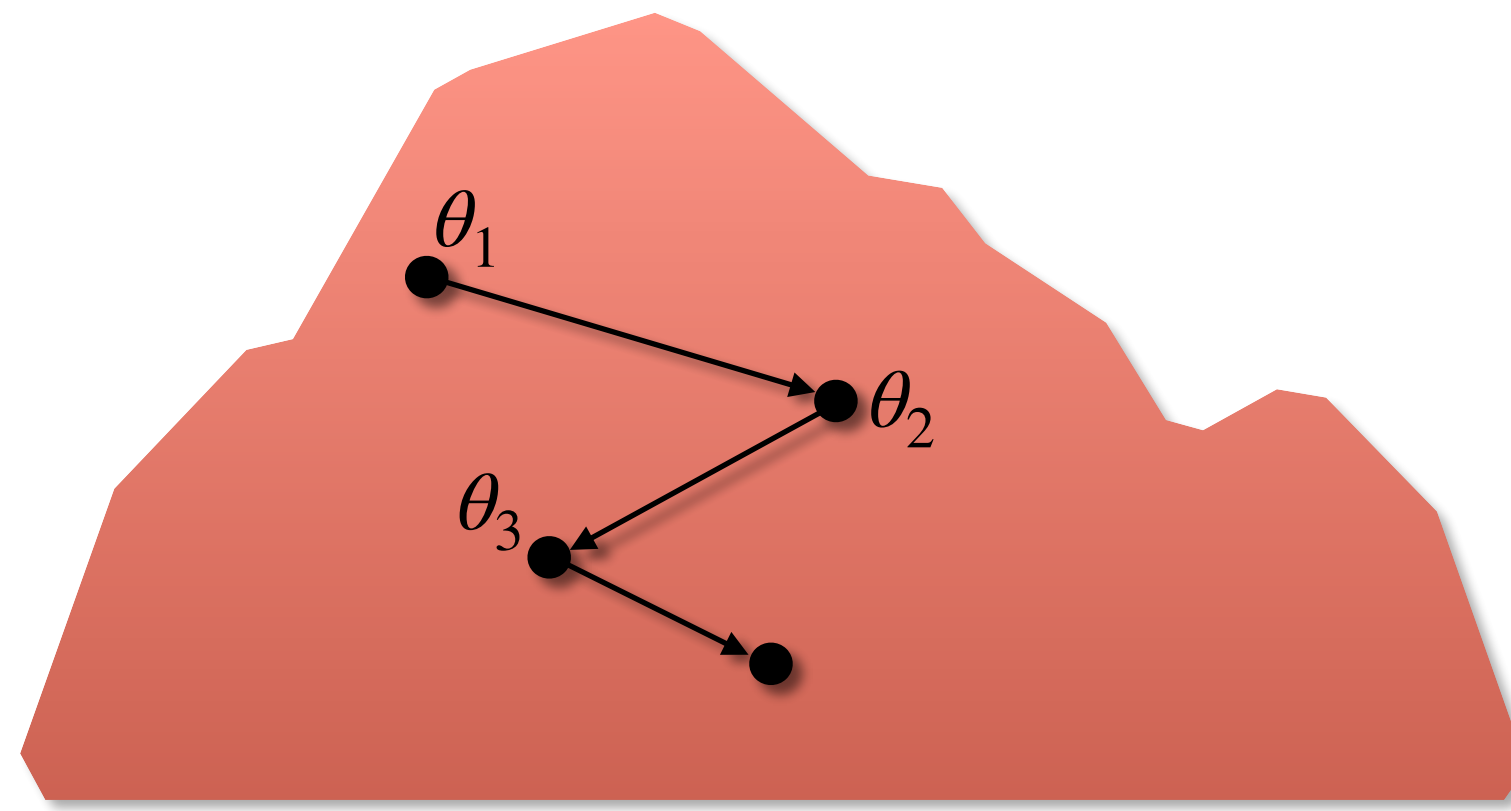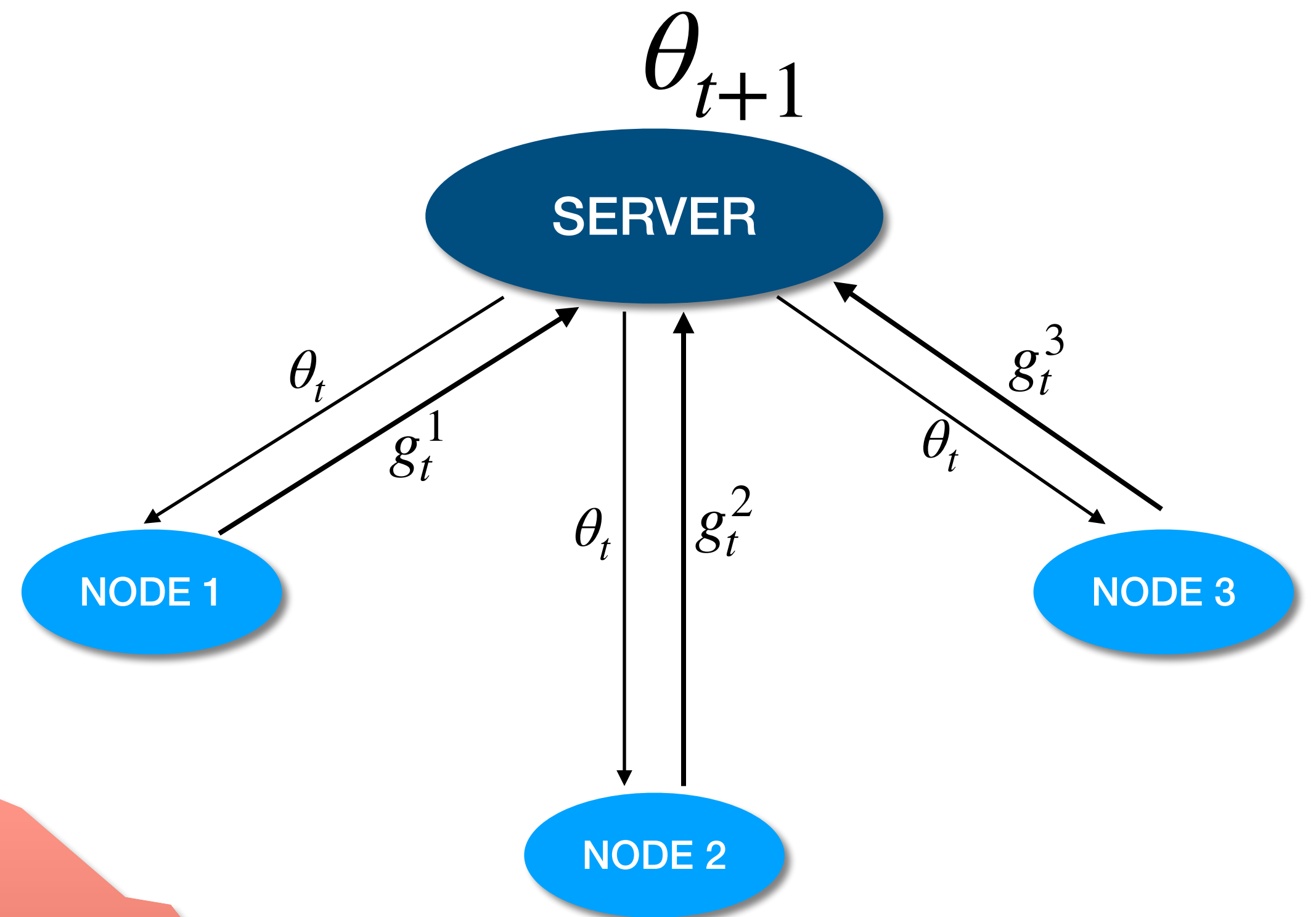$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

True gradient

**Server averages** the gradients, $\hat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \leftarrow \theta_t - \gamma_t \hat{g}_t$

$\theta_{t+1}$

SERVER

$\theta_t$    $g_t^1$      $g_t^3$   $\theta_t$

NODE 1

$\theta_t$   $g_t^2$

NODE 3

NODE 2

$\theta_1$

$\theta_2$

$\theta_3$

$\theta_4$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \boxed{\nabla Q(\theta_t)} + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

True gradient

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \, \widehat{g}_t$

$\theta_{t+1}$

SERVER

$\theta_t$   $g_t^1$    $\theta_t$   $g_t^2$    $g_t^3$   $\theta_t$

NODE 1     NODE 2     NODE 3

$\theta_1$

$\theta_2$

$\theta_3$

$\theta_4$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

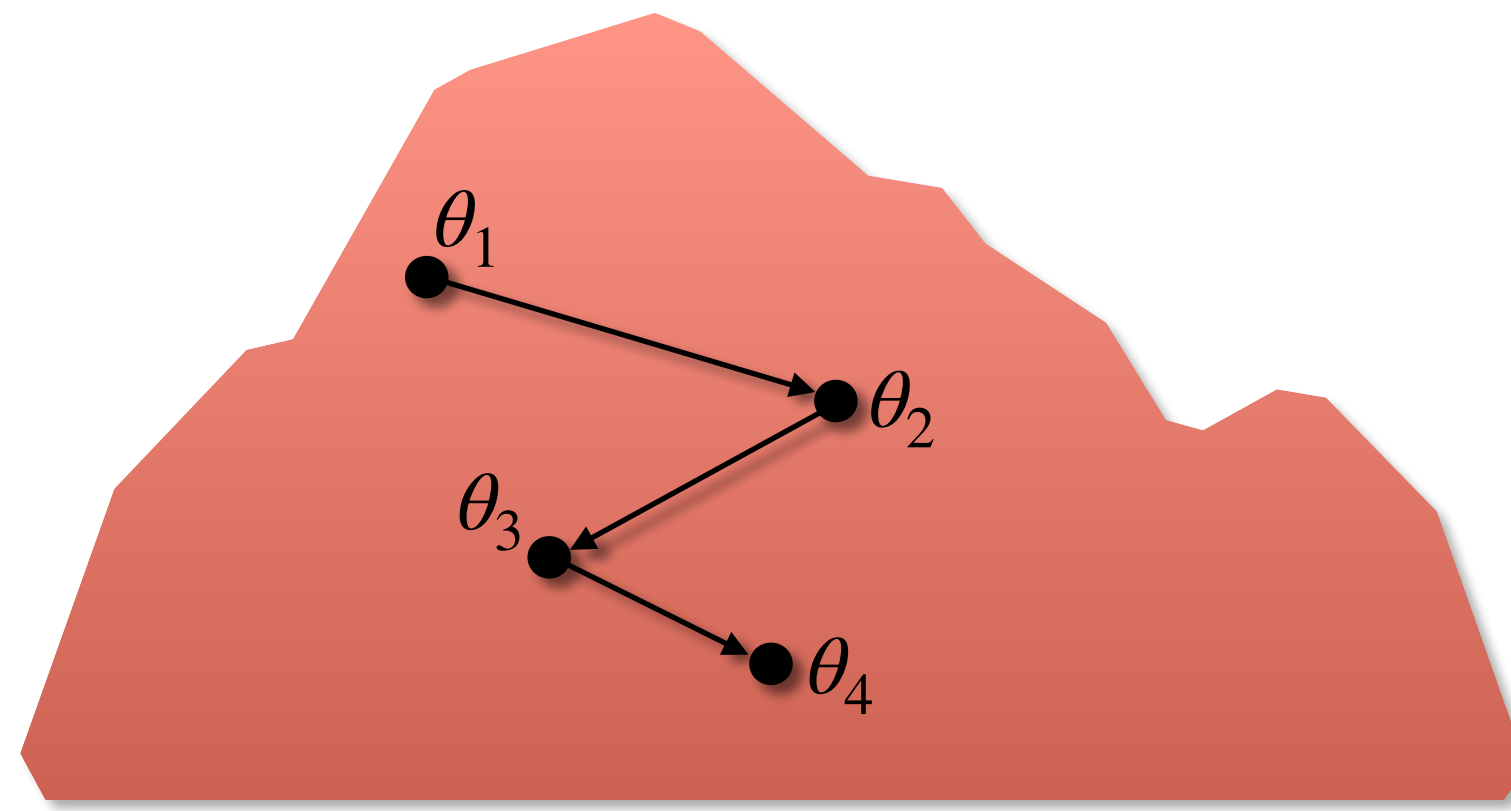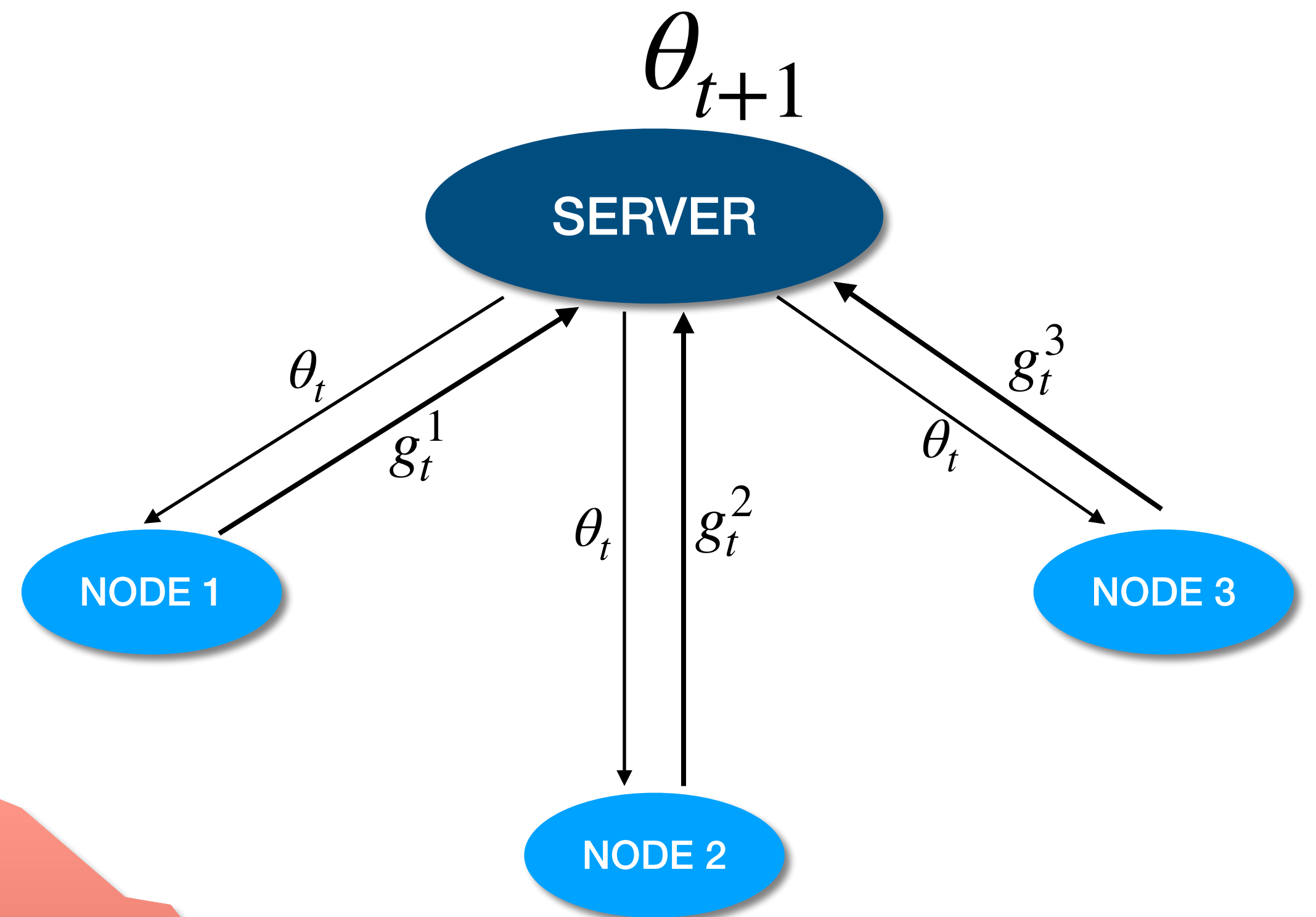$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$
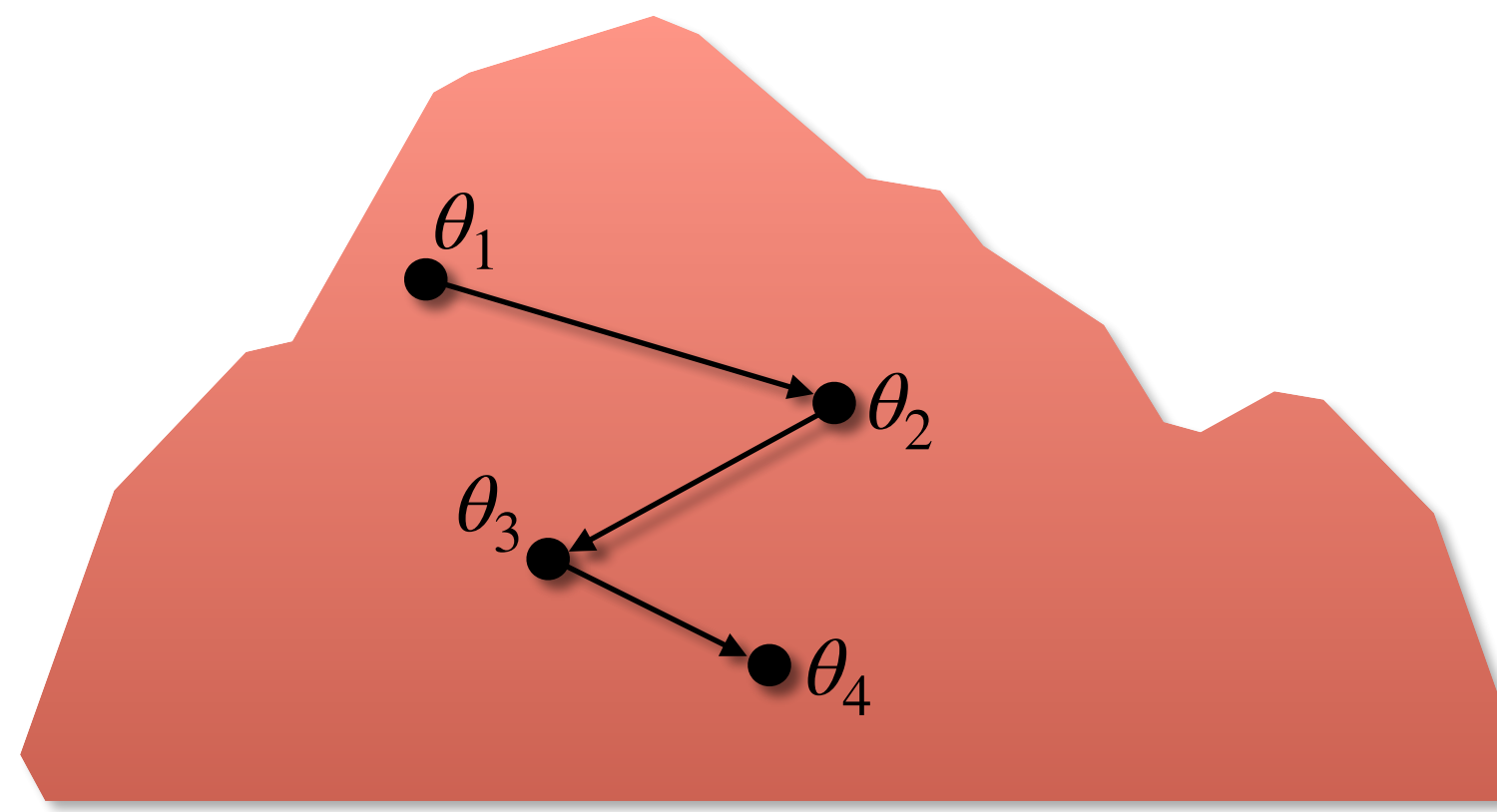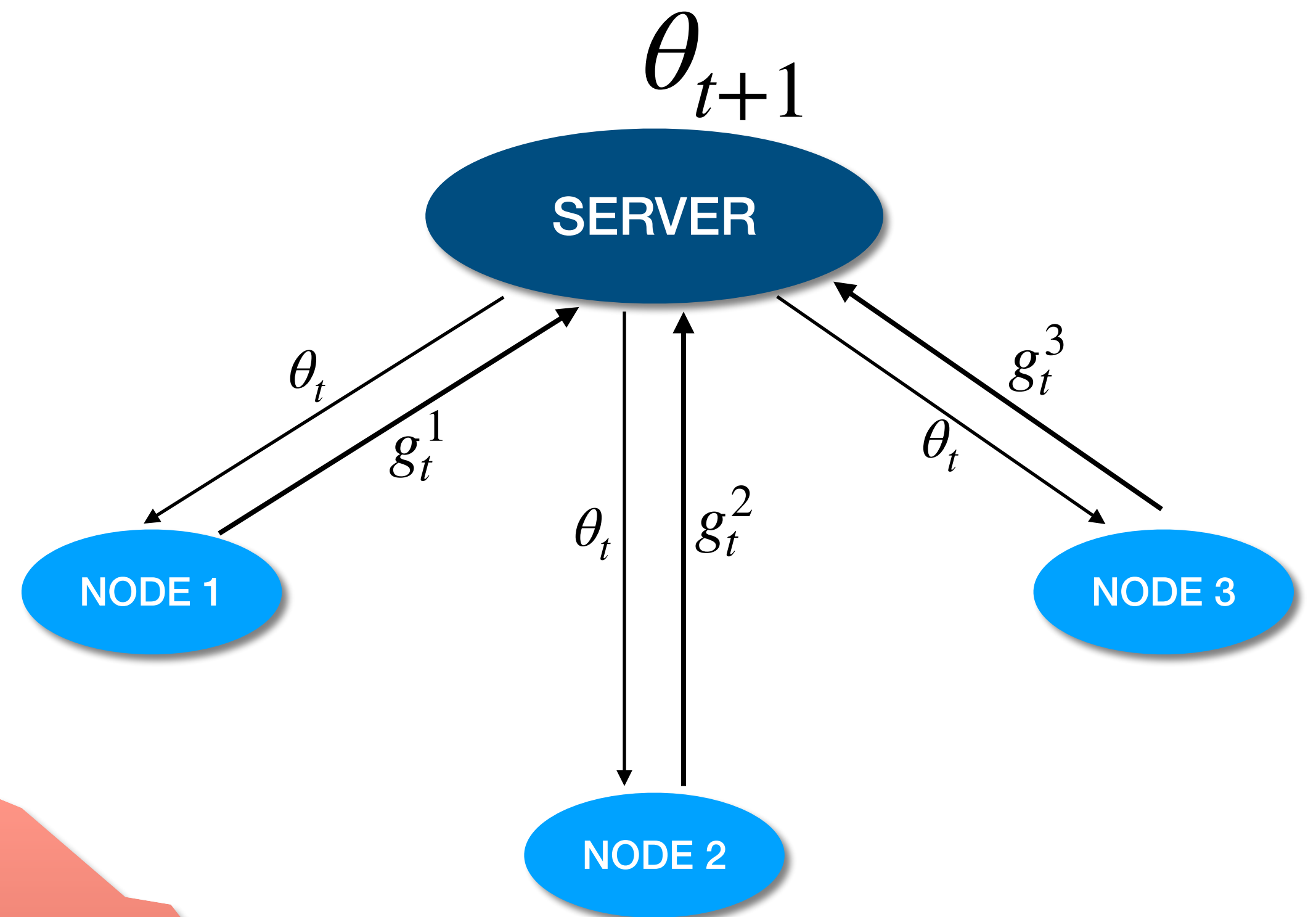
True gradient

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t\, \widehat{g}_t$

$\theta_{t+1}$

SERVER

$\theta_t$ $g_t^1$ $\theta_t$ $g_t^2$ $g_t^3$ $\theta_t$

NODE 1

NODE 3

NODE 2

$\theta_1$

$\theta_2$

$\theta_3$

$\theta_4$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

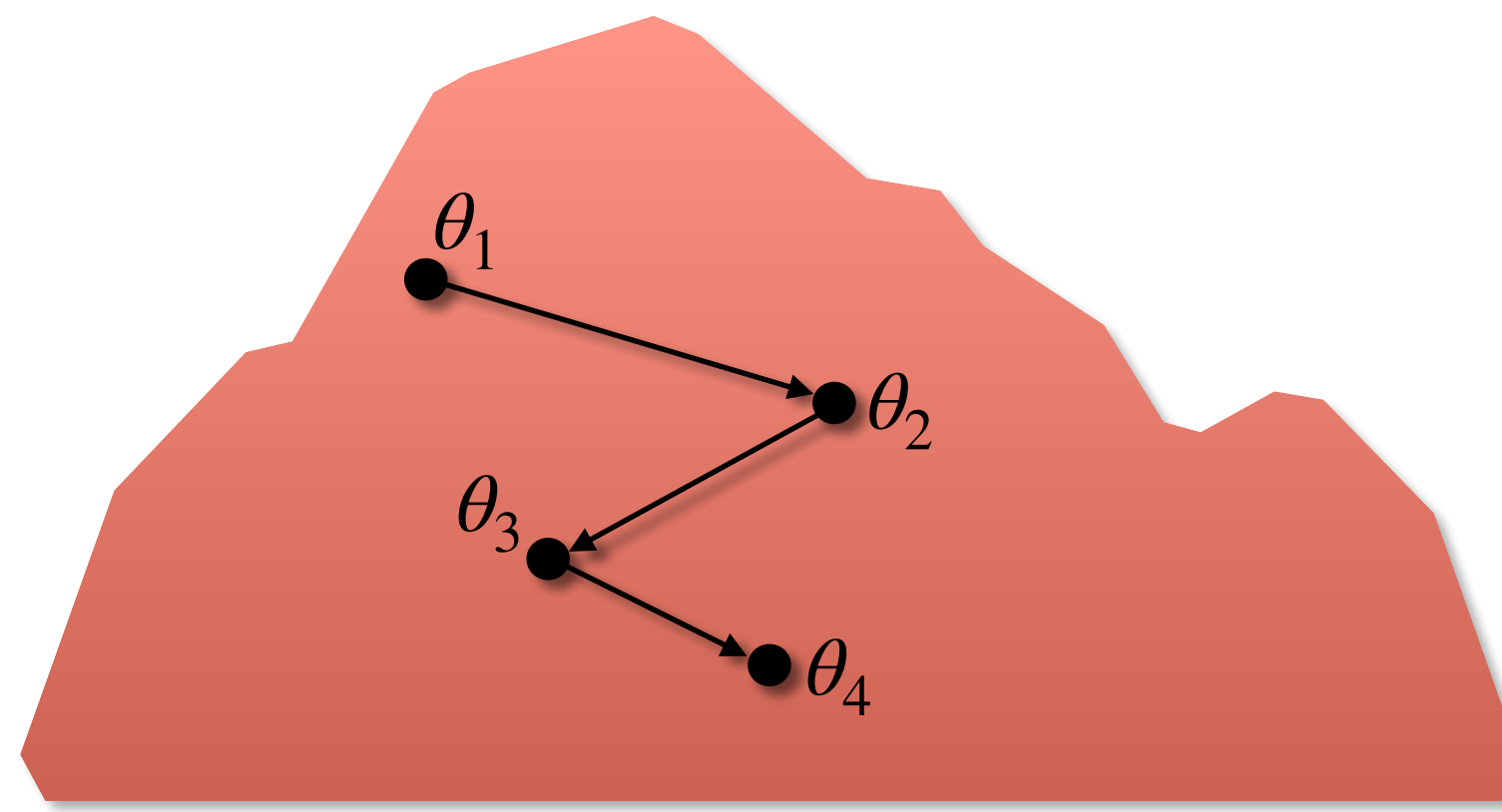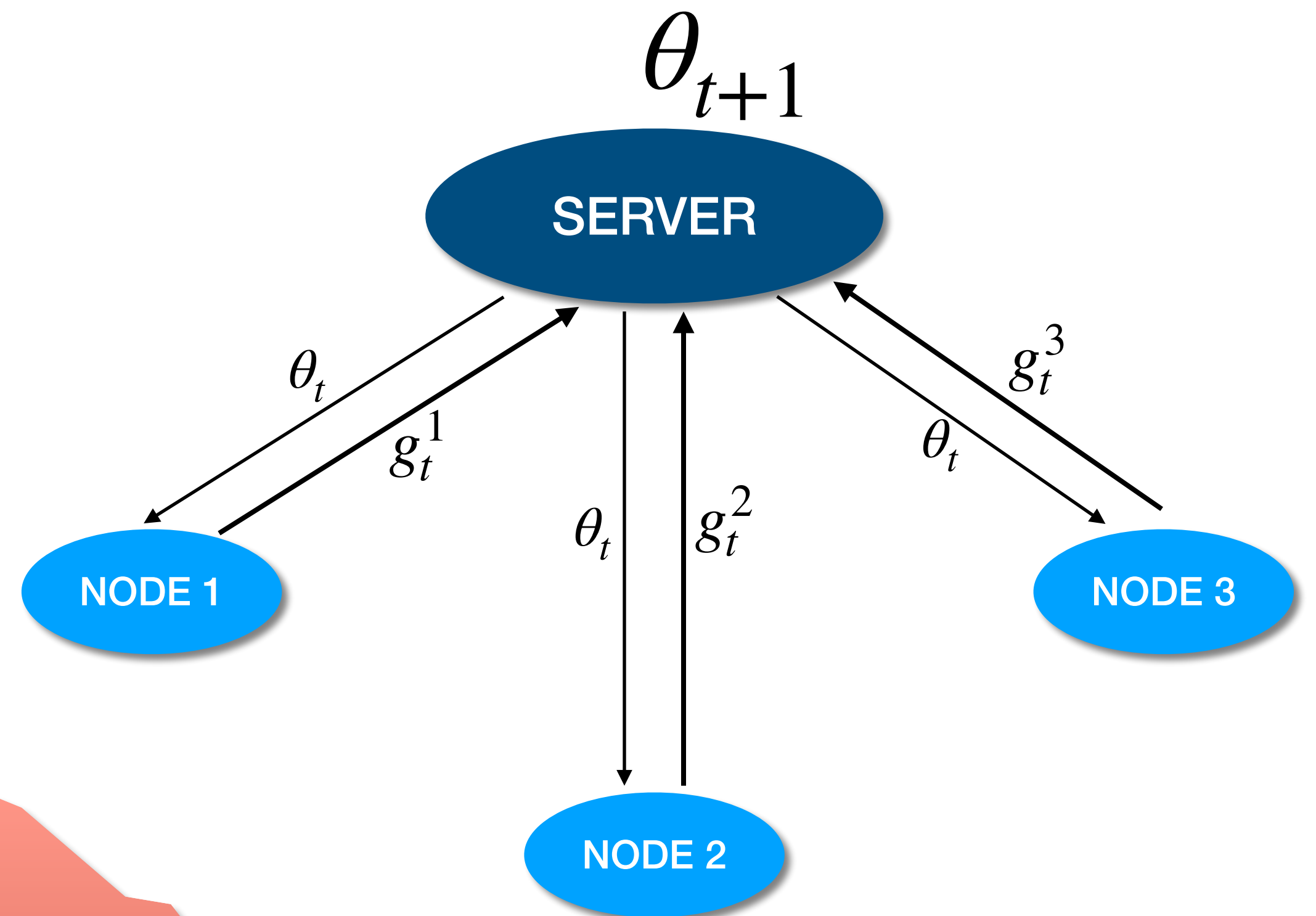$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathscr{U}_t$$

True gradient

$$\mathbb{E}\left[\mathscr{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathscr{U}_t\|^2\right] \leq \sigma^2 < \infty$$

reduces variance
by factor of n

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t\,\widehat{g}_t$

$\theta_{t+1}$

SERVER

$\theta_t$

$g_t^1$

$\theta_t$ $g_t^2$

$g_t^3$

$\theta_t$

NODE 1

NODE 3

NODE 2

$\theta_1$

$\theta_2$

$\theta_3$

$\theta_4$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \boxed{\nabla Q(\theta_t)} + u_t^i \quad ; \quad u_t^i \sim \mathscr{U}_t$$

True gradient

$$\mathbb{E}\left[\mathscr{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathscr{U}_t\|^2\right] \leq \sigma^2 < \infty$$
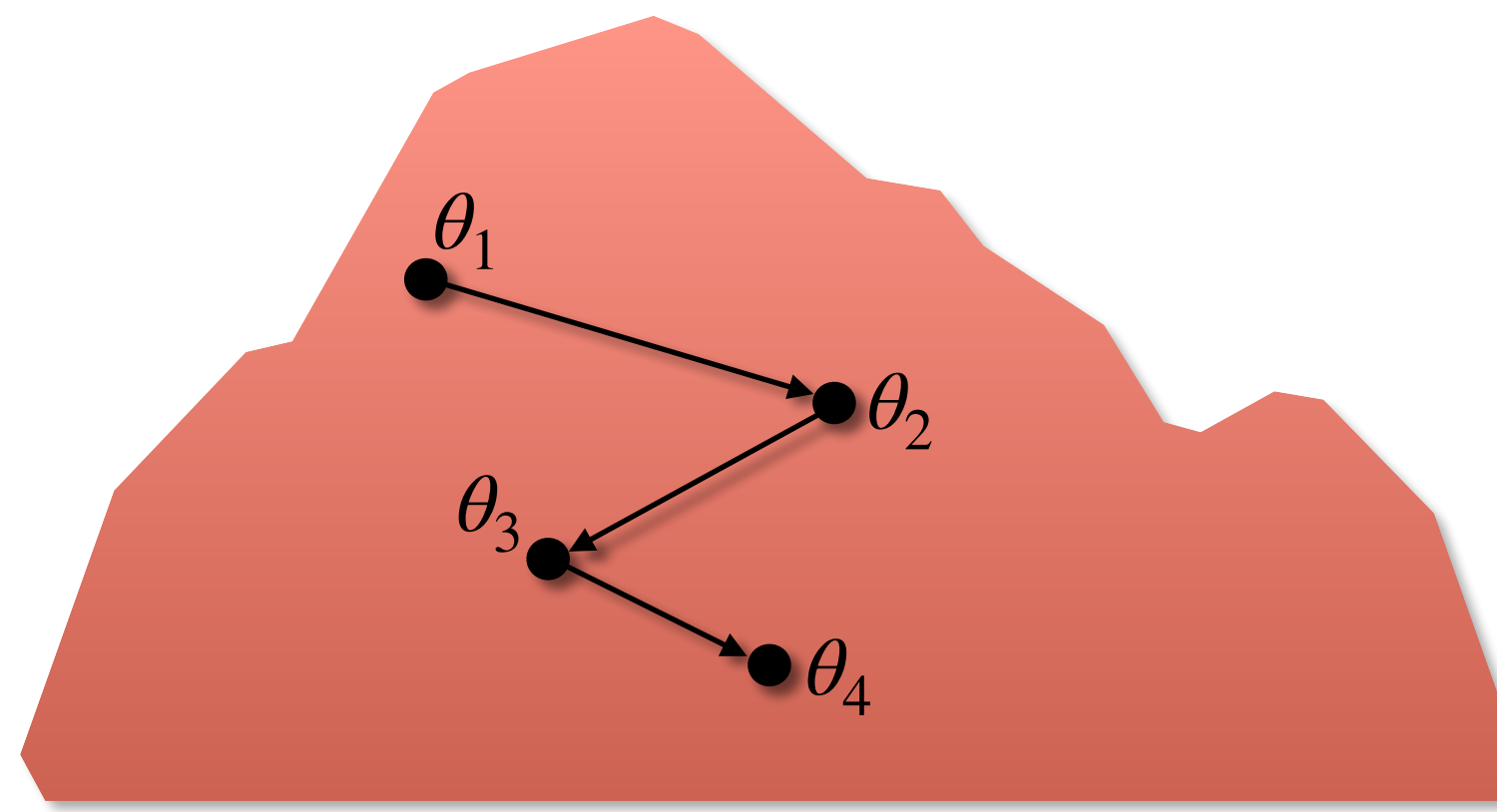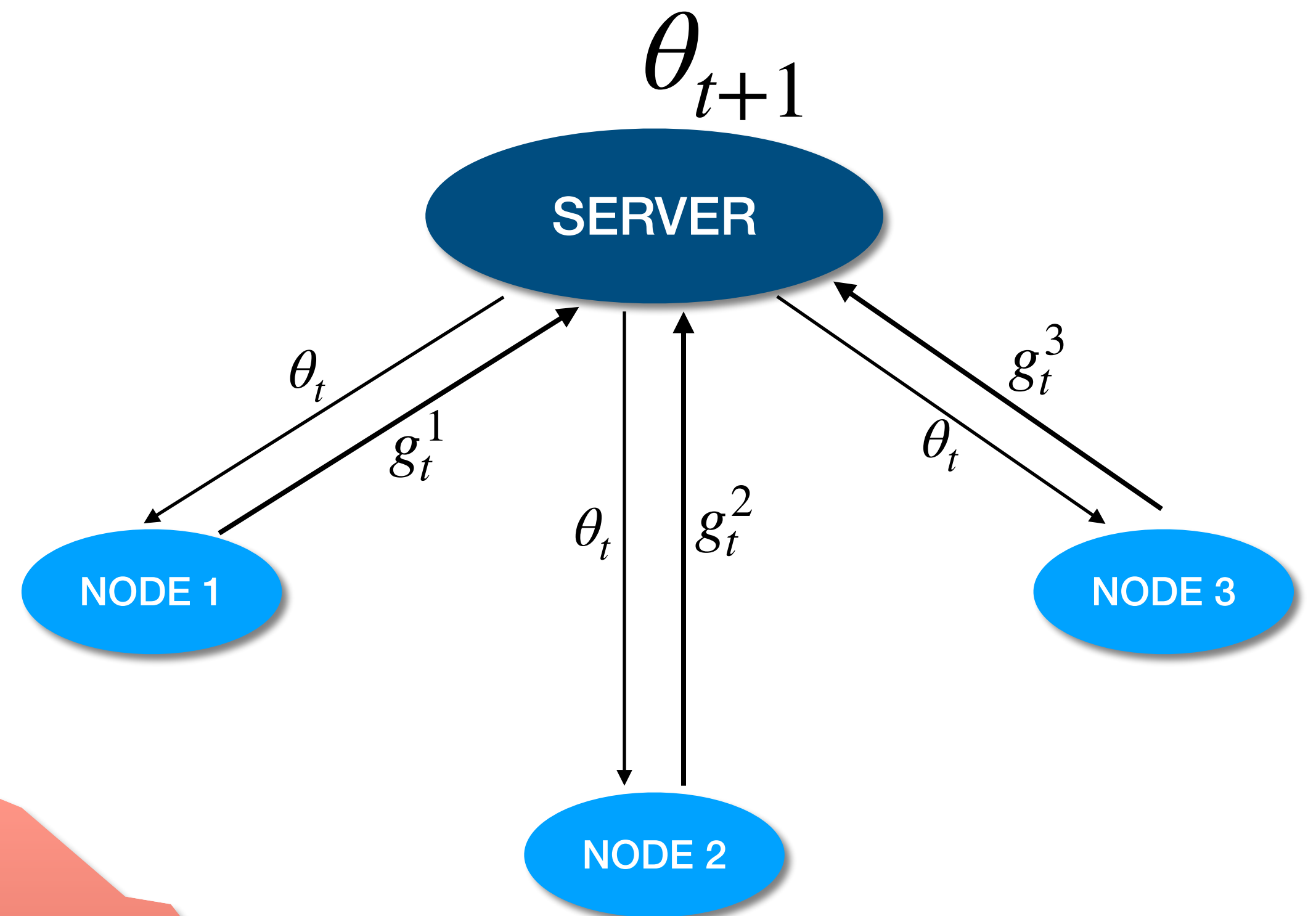
reduces variance
by factor of n

**Server averages** the gradients, $\widehat{g}_t = \frac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \widehat{g}_t$

$\theta_{t+1}$

SERVER

NODE 1

NODE 2

NODE 3

$\theta_t$   $g_t^1$   $\theta_t$   $g_t^2$   $g_t^3$   $\theta_t$

$\theta_1$

$\theta_2$

$\theta_3$

$\theta_4$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \nabla Q(\theta_t) + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \leq \sigma^2 < \infty$$
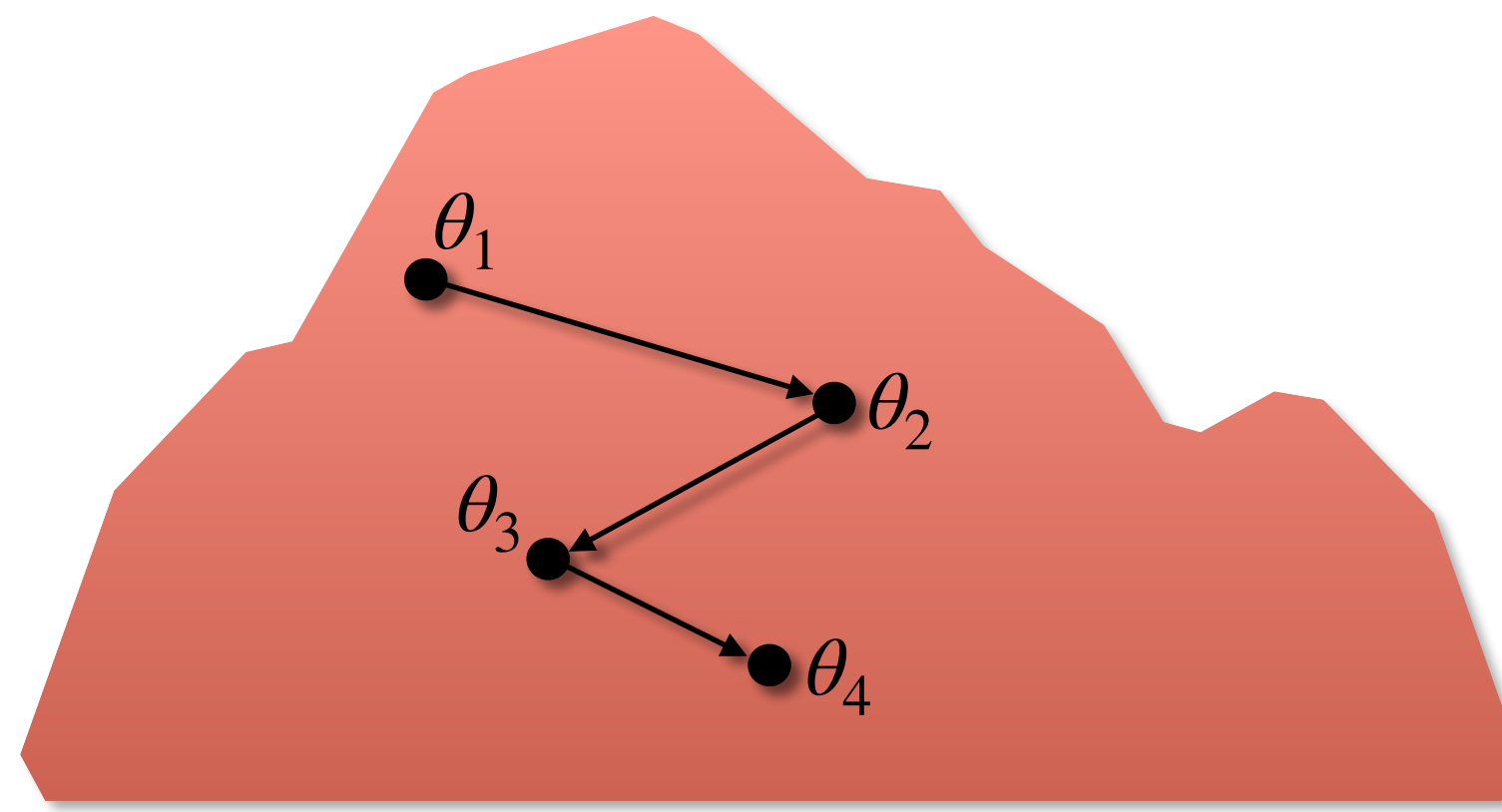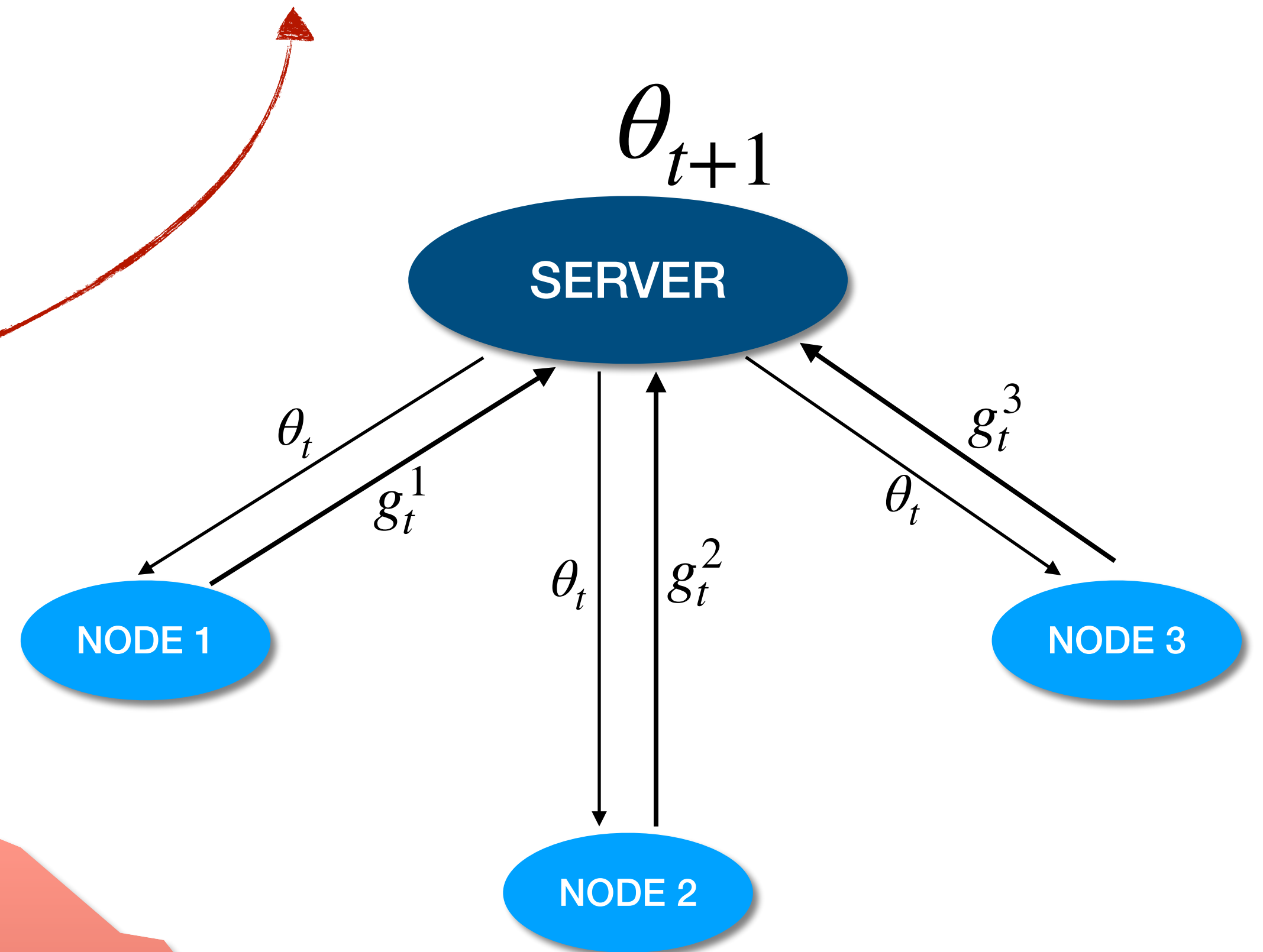
True gradient

reduces variance
by factor of n

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \widehat{g}_t$

Upon $T$ iterations

$\theta_{t+1}$

SERVER

$\theta_t$

$g_t^1$

$\theta_t$

$g_t^2$

$g_t^3$

$\theta_t$

NODE 1

NODE 2

NODE 3

$\theta_1$

$\theta_2$

$\theta_3$

$\theta_4$

# Distributed SGD

**Divides the workload per machine by the total size of the system**

**Nodes query** *stochastic gradients* with bounded variance

$$g_t^i = \boxed{\nabla Q(\theta_t)} + u_t^i \quad ; \quad u_t^i \sim \mathcal{U}_t$$

True gradient

$$\mathbb{E}\left[\mathcal{U}_t\right] = 0 \quad ; \quad \mathbb{E}\left[\|\mathcal{U}_t\|^2\right] \le \sigma^2 < \infty$$
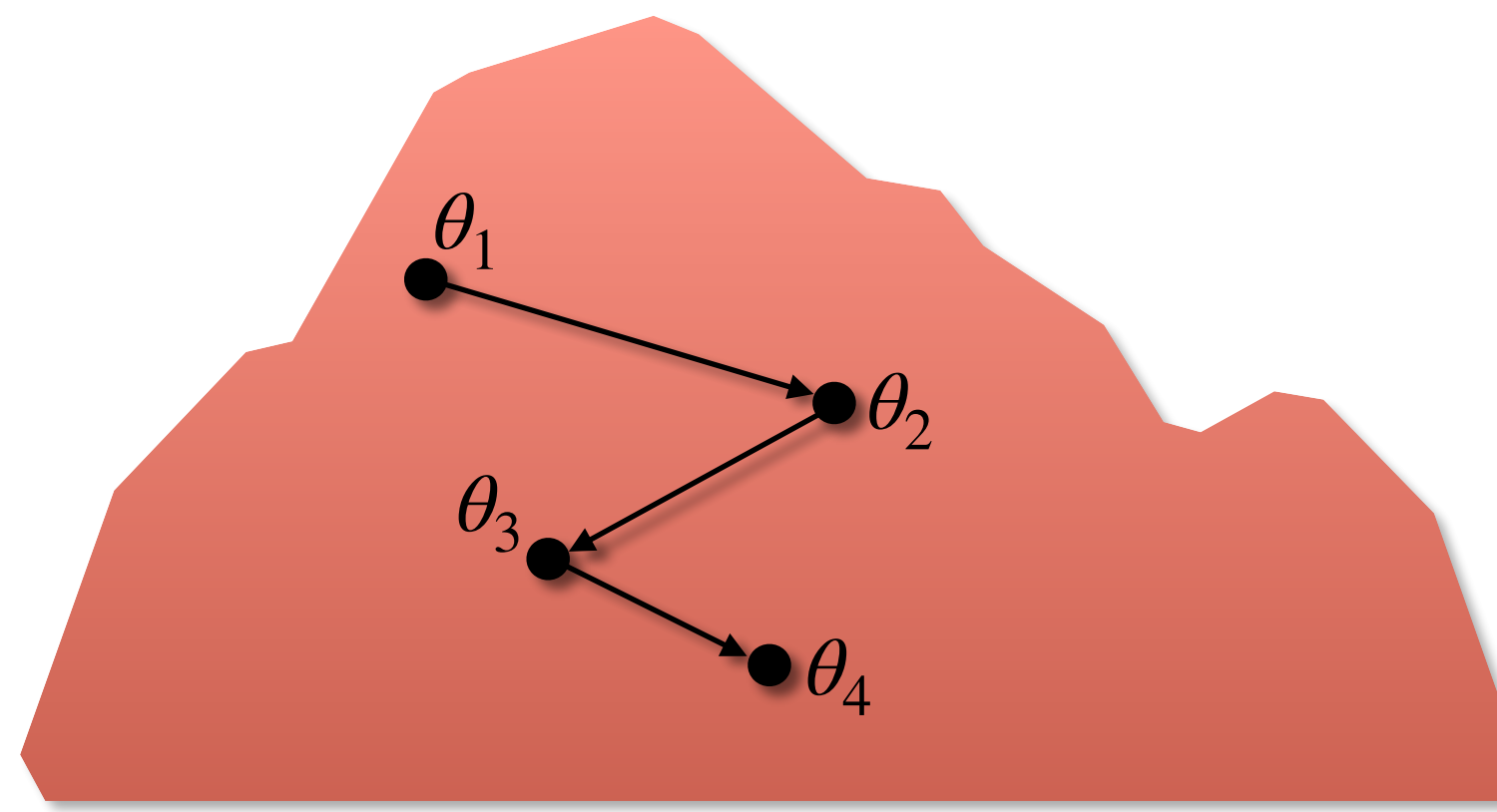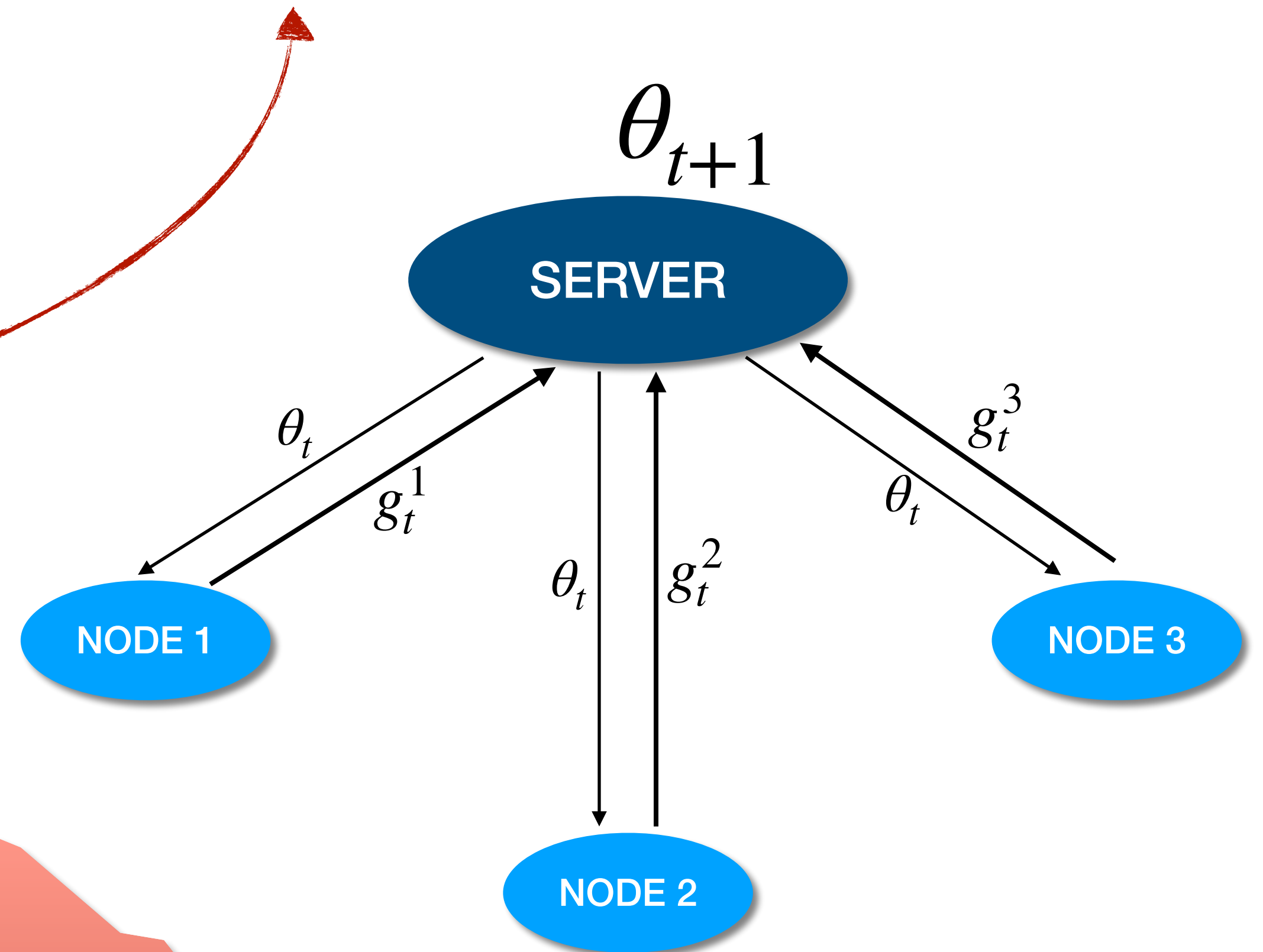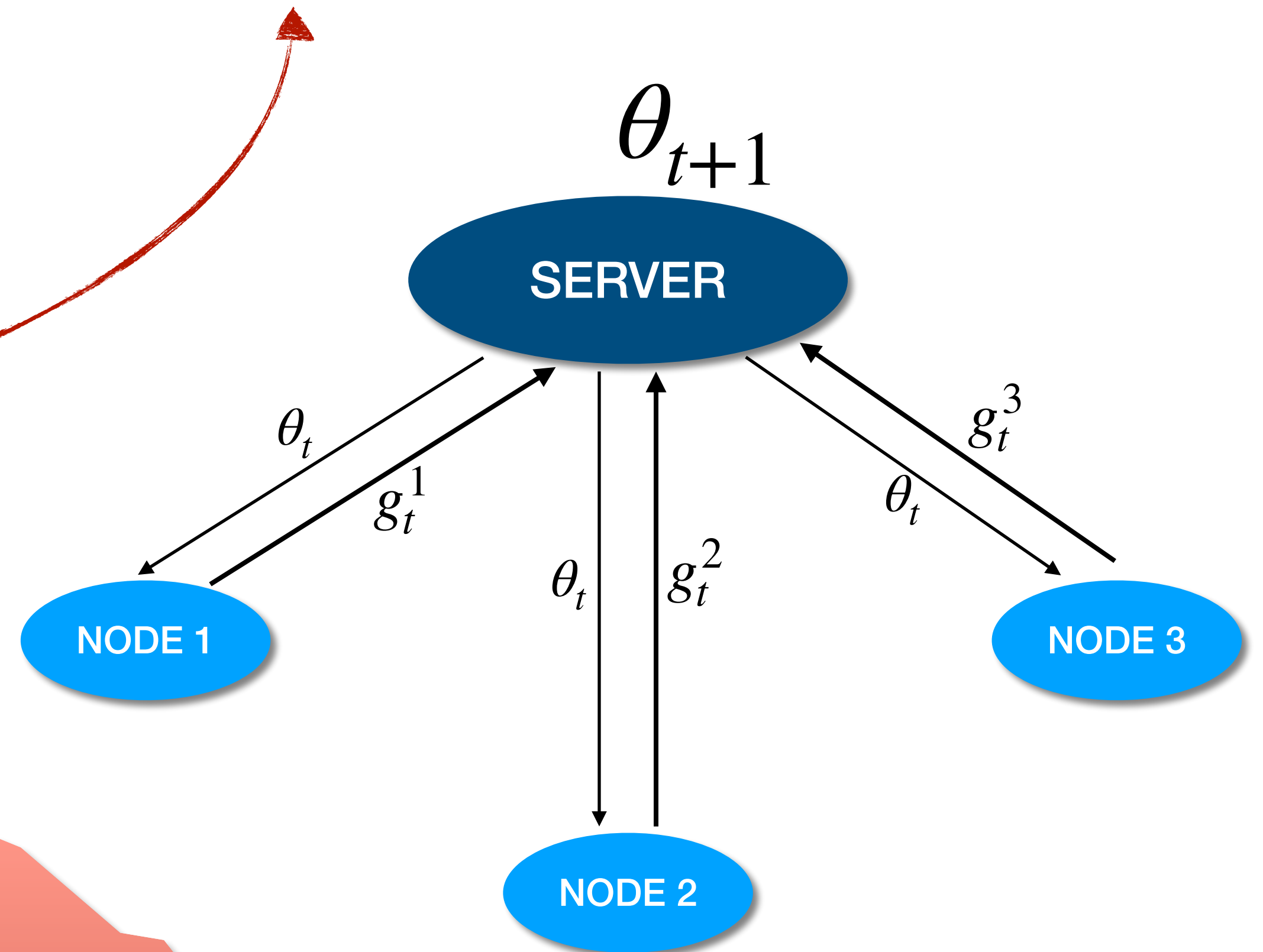
reduces variance
by factor of n

**Server averages** the gradients, $\widehat{g}_t = \dfrac{1}{n}\sum_i g_t^i$

Learning rate

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \, \widehat{g}_t$

Upon $T$ iterations

$$\frac{1}{T}\sum_{t=1}^{T} \mathbb{E}\left[\|\nabla Q\left(\theta_t\right)\|^2\right] \le \epsilon \in \mathcal{O}\left(\sqrt{\frac{\sigma^2}{nT}}\right)$$

$\theta_{t+1}$

SERVER

$\theta_t$

$g_t^1$

$g_t^3$

$\theta_t$

$\theta_t$

$g_t^2$

NODE 1

NODE 3

NODE 2

$\theta_1$

$\theta_2$

$\theta_3$

$\theta_4$

# Byzantine Resilience in Distributed Learning

# Byzantine Resilience in Distributed Learning

$f$ **out of** $n$ **nodes are Byzantine faulty**

# Byzantine Resilience in Distributed Learning

$f$ **out of** $n$ **nodes are Byzantine faulty**

**CAN THE SERVER STILL OBTAIN**

$$\theta* \in \left( \theta \ ; \ \nabla Q(\theta) = 0 \right)$$

$\theta_t$

SERVER

$\theta_t$

$g_t^1$

NODE 1

Honest

$\theta_t$

$g_t^2$

NODE 2

$\widetilde{g_t^3}$

$\theta_t$

NODE 3

Byzantine

# Byzantine Resilience in Distributed Learning

$f$ **out of** $n$ **nodes are Byzantine faulty**

**CAN THE SERVER STILL OBTAIN**

$$\theta^* \in \left( \theta \ ; \ \nabla Q(\theta) = 0 \right)$$

**Feasible** as nodes are assumed
to sample data from the same distribution*

$\theta_t$

SERVER

$\theta_t$

$g_t^1$

$\widetilde{g}_t^3$

$\theta_t$

$\theta_t$ $g_t^2$

NODE 1
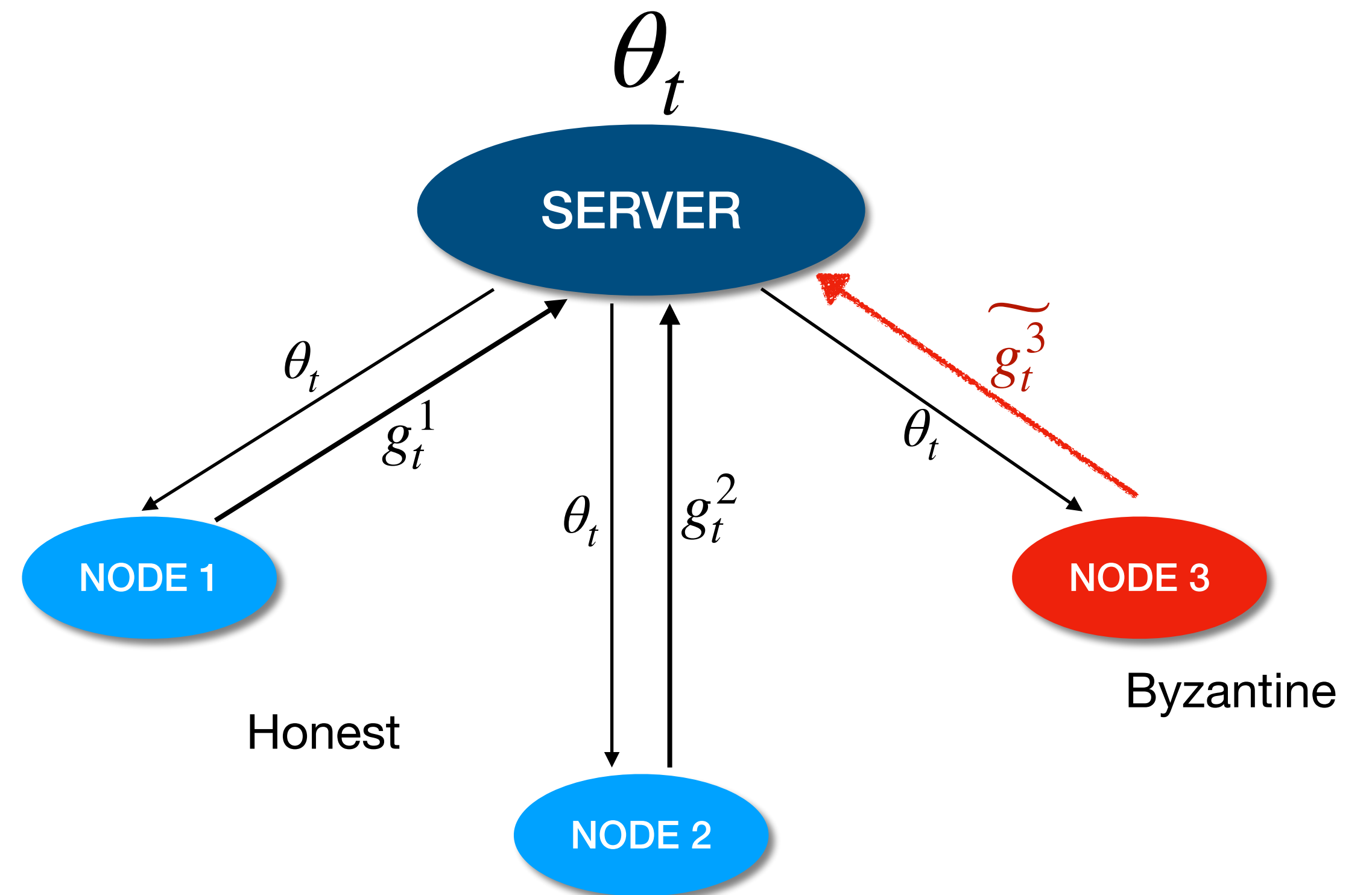
NODE 3

NODE 2

Honest

Byzantine

# Byzantine Resilience in Distributed Learning

$f$ **out of** $n$ **nodes are Byzantine faulty**

CAN THE SERVER STILL OBTAIN

$$\theta^* \in \left( \theta \ ; \ \nabla Q(\theta) = 0 \right)$$

**Feasible** as nodes are assumed
to sample data from the same distribution*

$$\theta_t$$

SERVER

$\theta_t$

$g_t^1$

$\widetilde{g}_t^3$

$\theta_t$

$\theta_t$ $g_t^2$

NODE 1

NODE 3

Byzantine

Honest

NODE 2

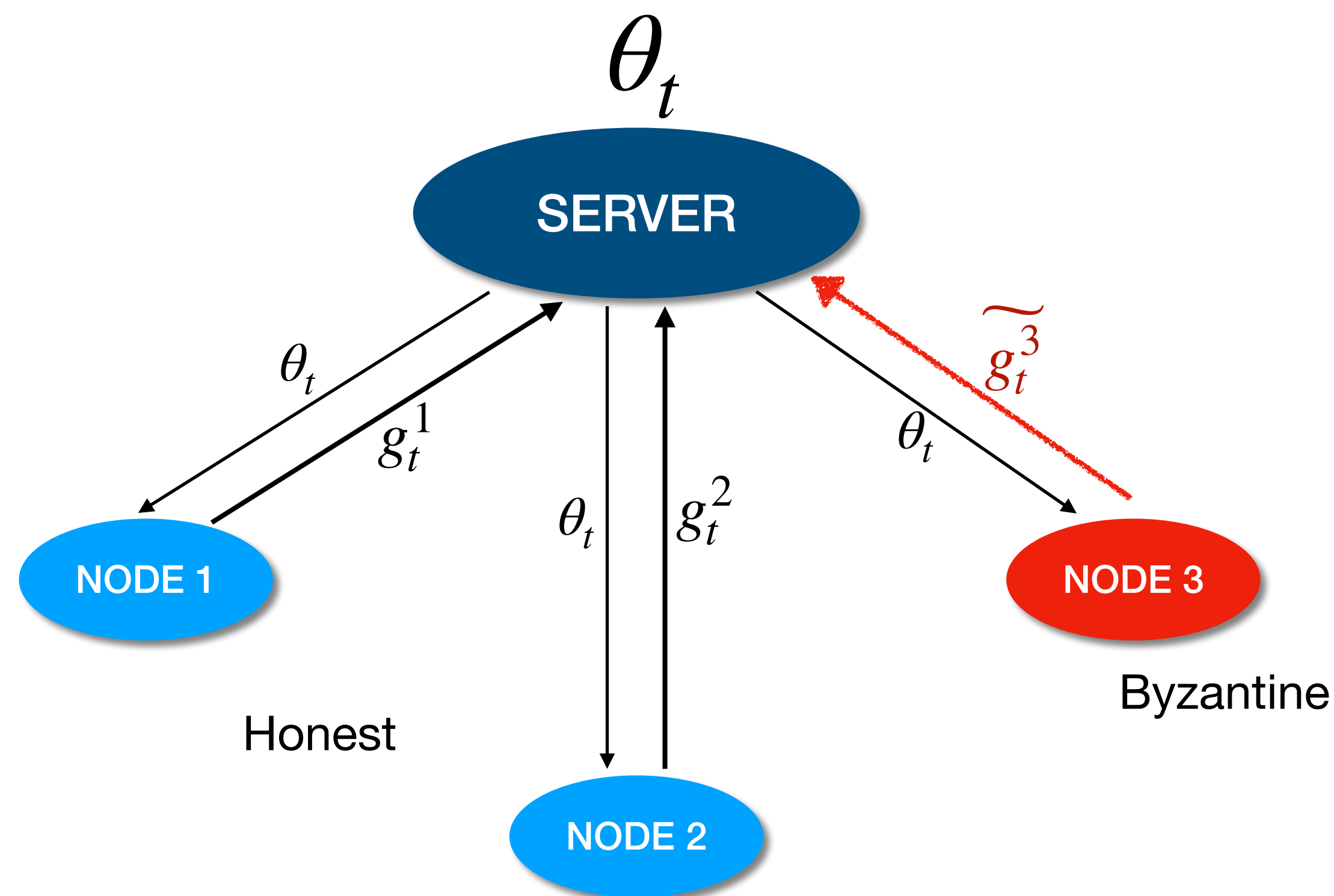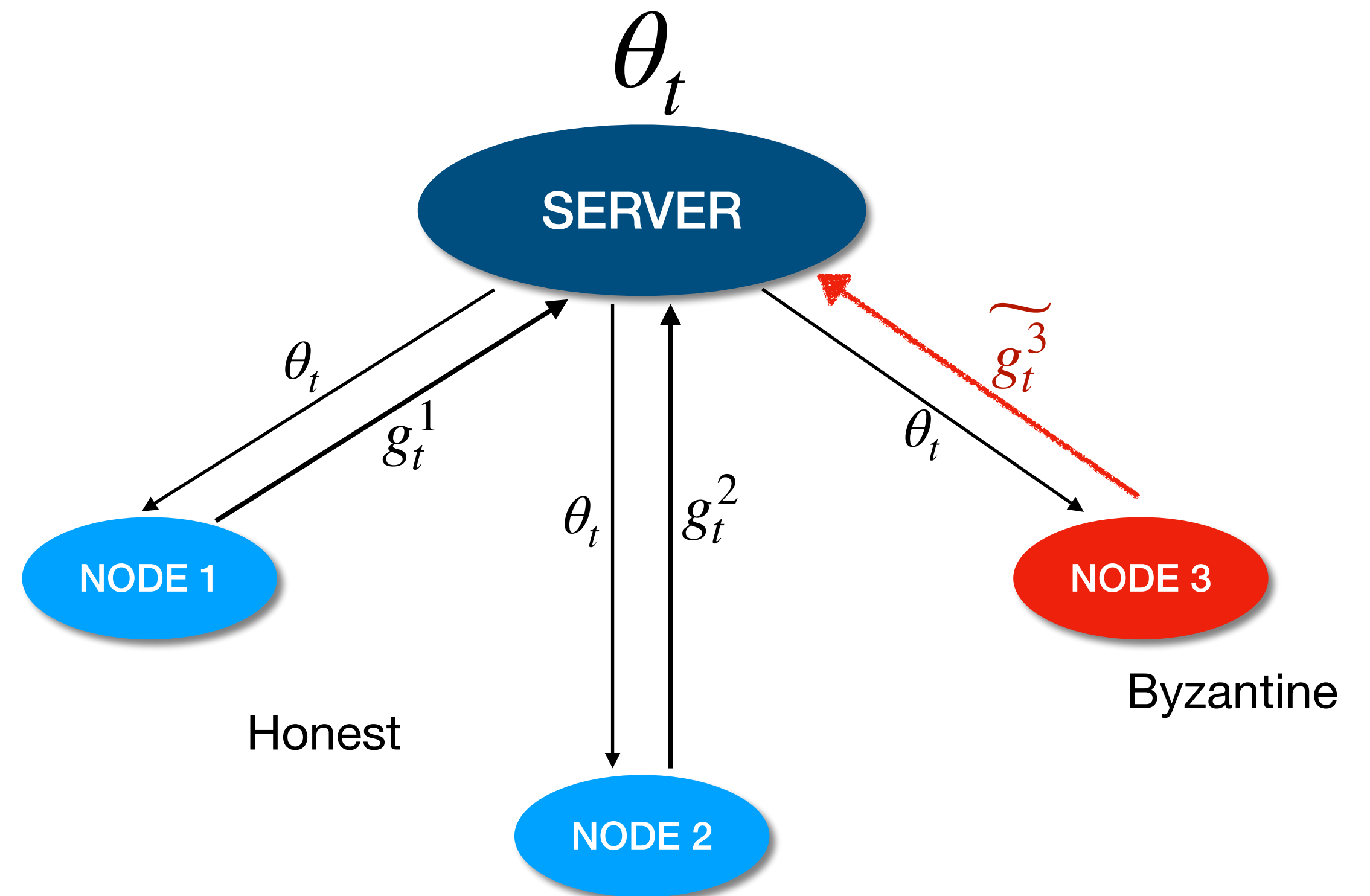* In general, when data distributions differ, we require $2f$-*redundancy* in data (Gupta and Vaidya, 2021)

# Byzantine Resilience in Distributed Learning

**$f$ out of $n$ nodes are Byzantine faulty**

CAN THE SERVER STILL OBTAIN

$$\theta^* \in \left( \theta \ ; \ \nabla Q(\theta) = 0 \right)$$

**Feasible** as nodes are assumed
to sample data from the same distribution*

$(f, \epsilon)$**-Resilience**

$\theta_t$

SERVER

$\theta_t$

$g_t^1$

$\theta_t$

$g_t^2$

$\theta_t$

$\widetilde{g_t^3}$

NODE 1

Honest

NODE 2

NODE 3

Byzantine

* In general, when data distributions differ, we require 2*f-redundancy* in data (Gupta and Vaidya, 2021)

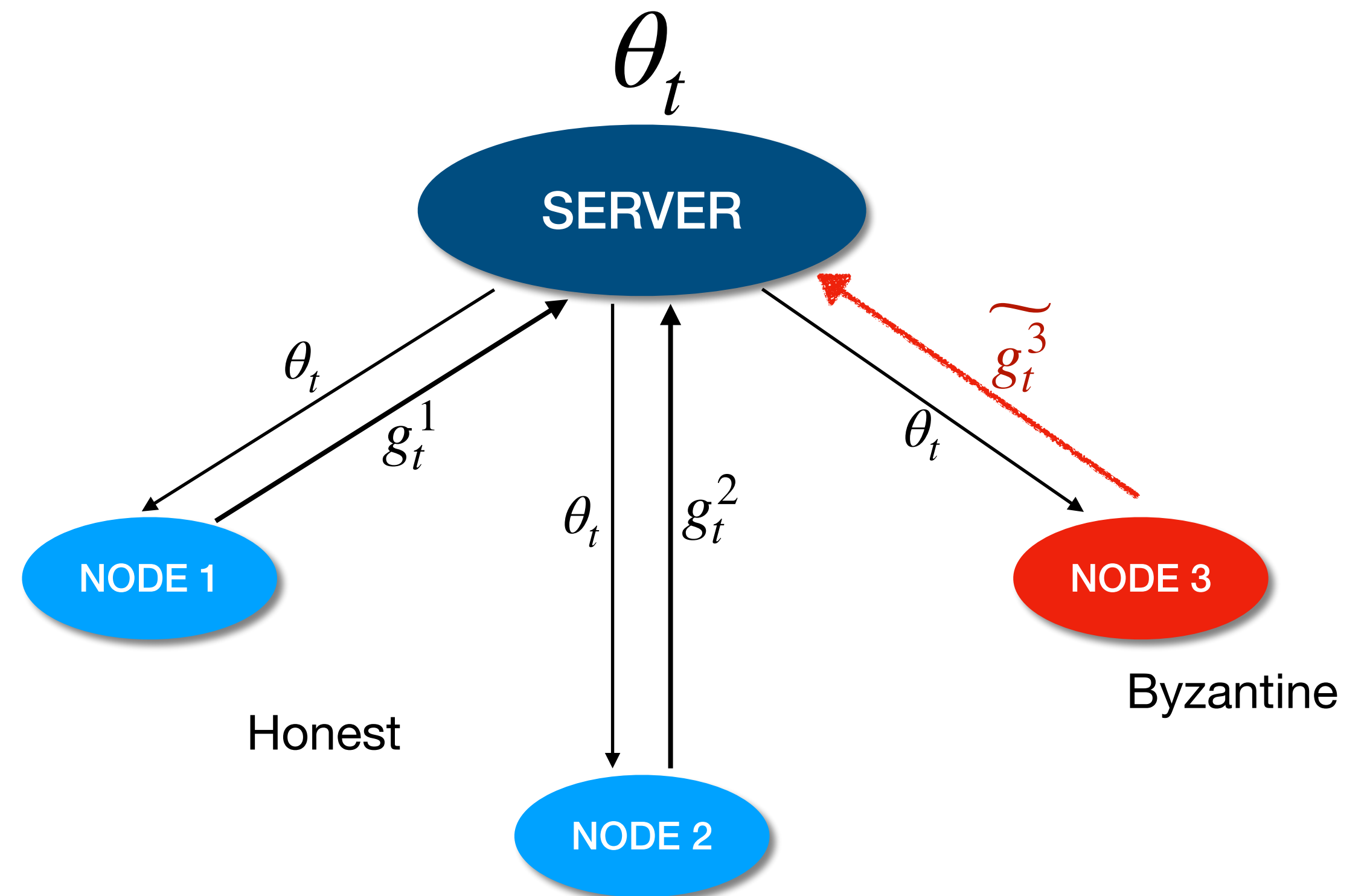# Byzantine Resilience in Distributed Learning

$f$ **out of** $n$ **nodes are Byzantine faulty**

**CAN THE SERVER STILL OBTAIN**

$$\theta* \in \left(\theta \; ; \; \nabla Q(\theta) = 0\right)$$
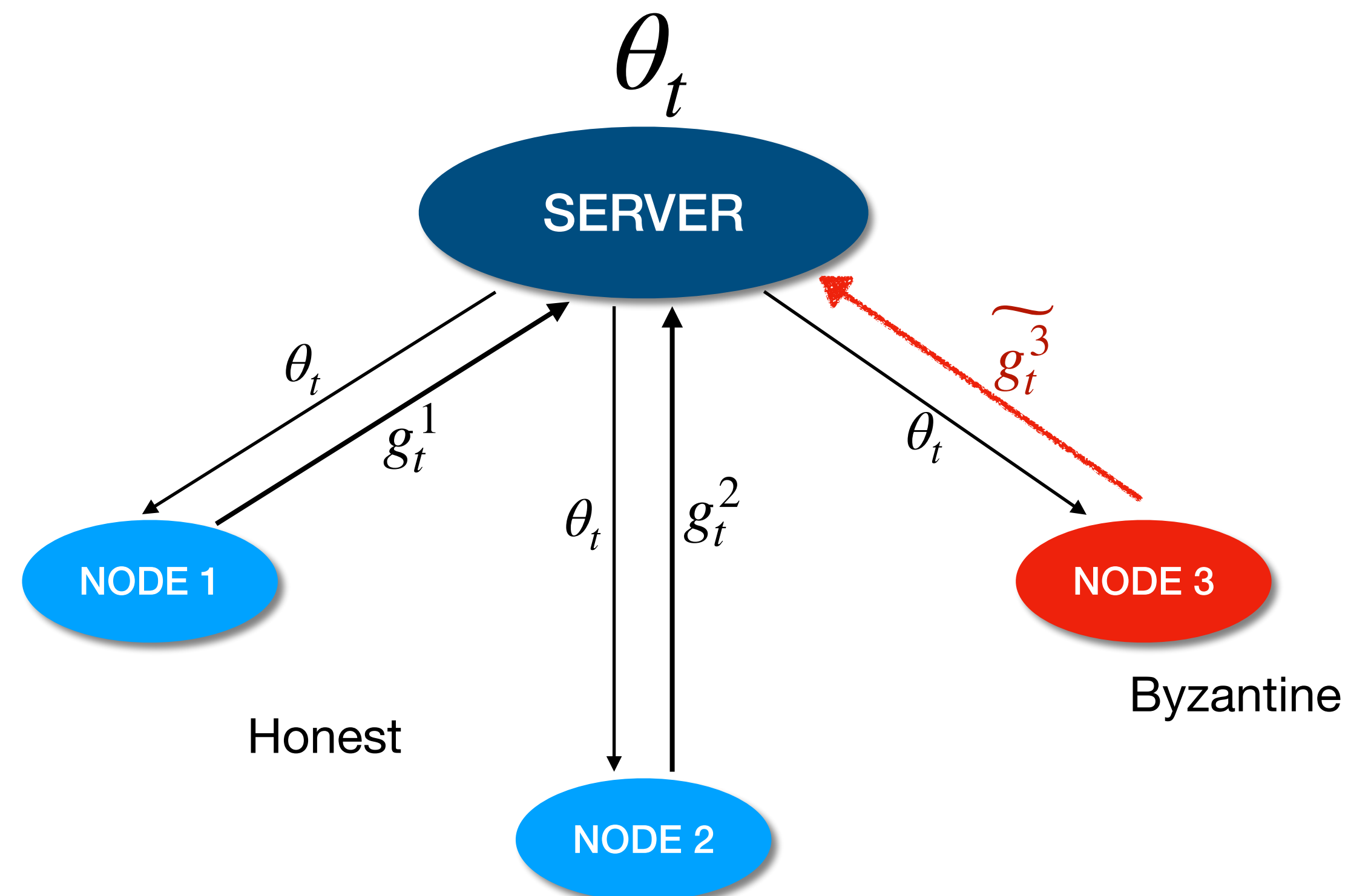
**Feasible** as nodes are assumed
to sample data from the same distribution*

$(f, \epsilon)$**-Resilience**

An algorithm that allows the honest

nodes to compute $\widehat{\theta}$ such that

$$\|\nabla Q\left(\widehat{\theta}\right)\|^2 \leq \epsilon$$

$\theta_t$

SERVER

$\theta_t$

$g_t^1$

$\widetilde{g_t^3}$

$\theta_t$

$\theta_t$ $g_t^2$

$\theta_t$

NODE 1

NODE 3

NODE 2

Honest

Byzantine

* In general, when data distributions differ, we require 2*f-redundancy* in data (Gupta and Vaidya, 2021)

# Byzantine Resilience in Distributed Learning

$f$ **out of** $n$ **nodes are Byzantine faulty**

CAN THE SERVER STILL OBTAIN

$$\theta^* \in \left(\theta \ ; \ \nabla Q(\theta) = 0\right)$$
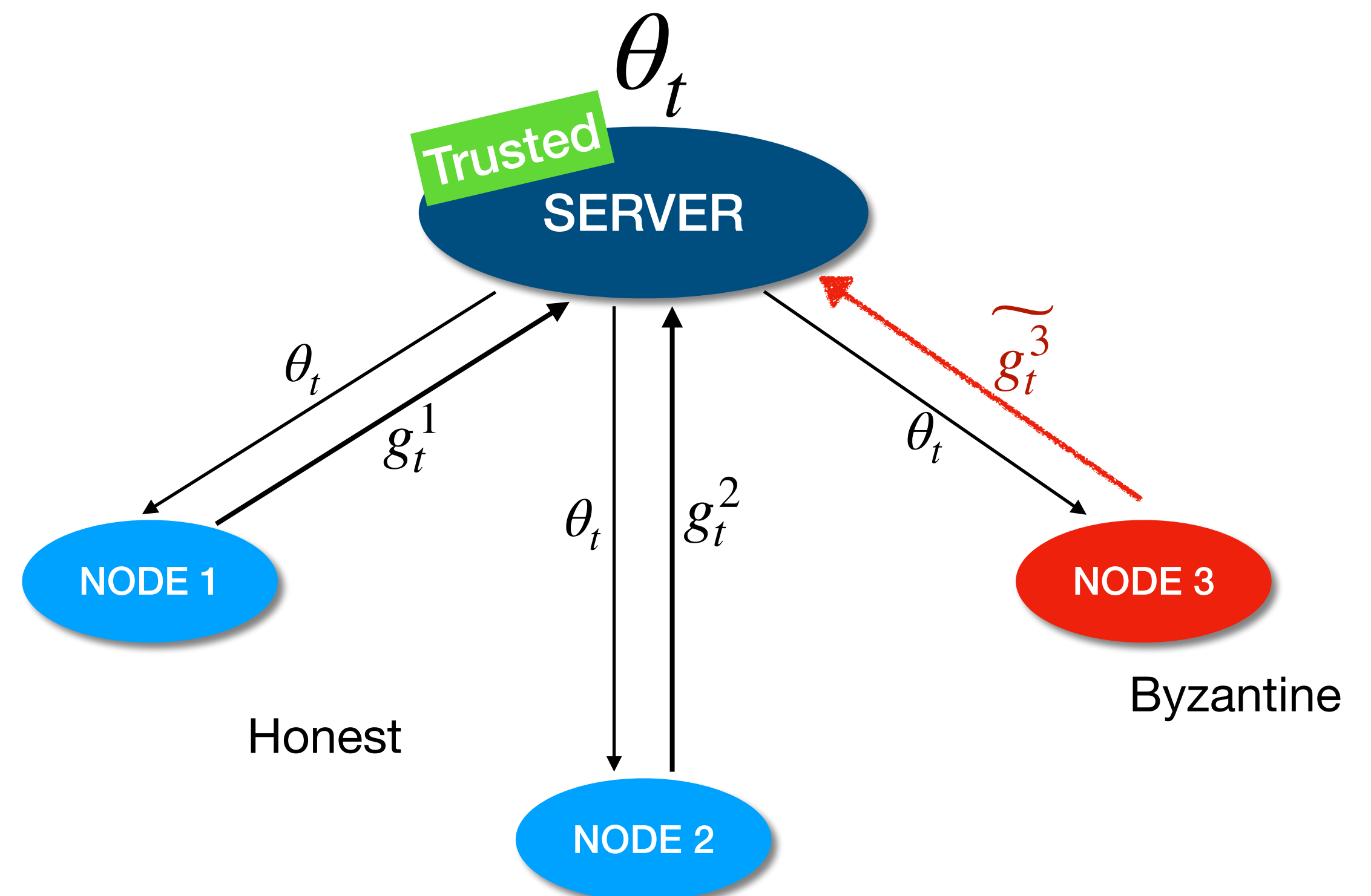
**Feasible** as nodes are assumed
to sample data from the same distribution*

$(f, \epsilon)$**-Resilience**

An algorithm that allows the honest

nodes to compute $\widehat{\theta}$ such that

$$\|\nabla Q\left(\widehat{\theta}\right)\|^2 \leq \epsilon$$



$\theta_t$

Trusted

SERVER

$\theta_t$

$g_t^1$

$\widetilde{g_t^3}$

$\theta_t$

$\theta_t$

$g_t^2$

NODE 1

Honest

NODE 2

NODE 3

Byzantine

\* In general, when data distributions differ, we require 2*f-redundancy* in data (Gupta and Vaidya, 2021)

# Byzantine Resilience in Distributed Learning

$f$ **out of** $n$ **nodes are Byzantine faulty**

CAN THE SERVER STILL OBTAIN

$$\theta* \in \left( \theta \; ; \; \nabla Q(\theta) = 0 \right)$$
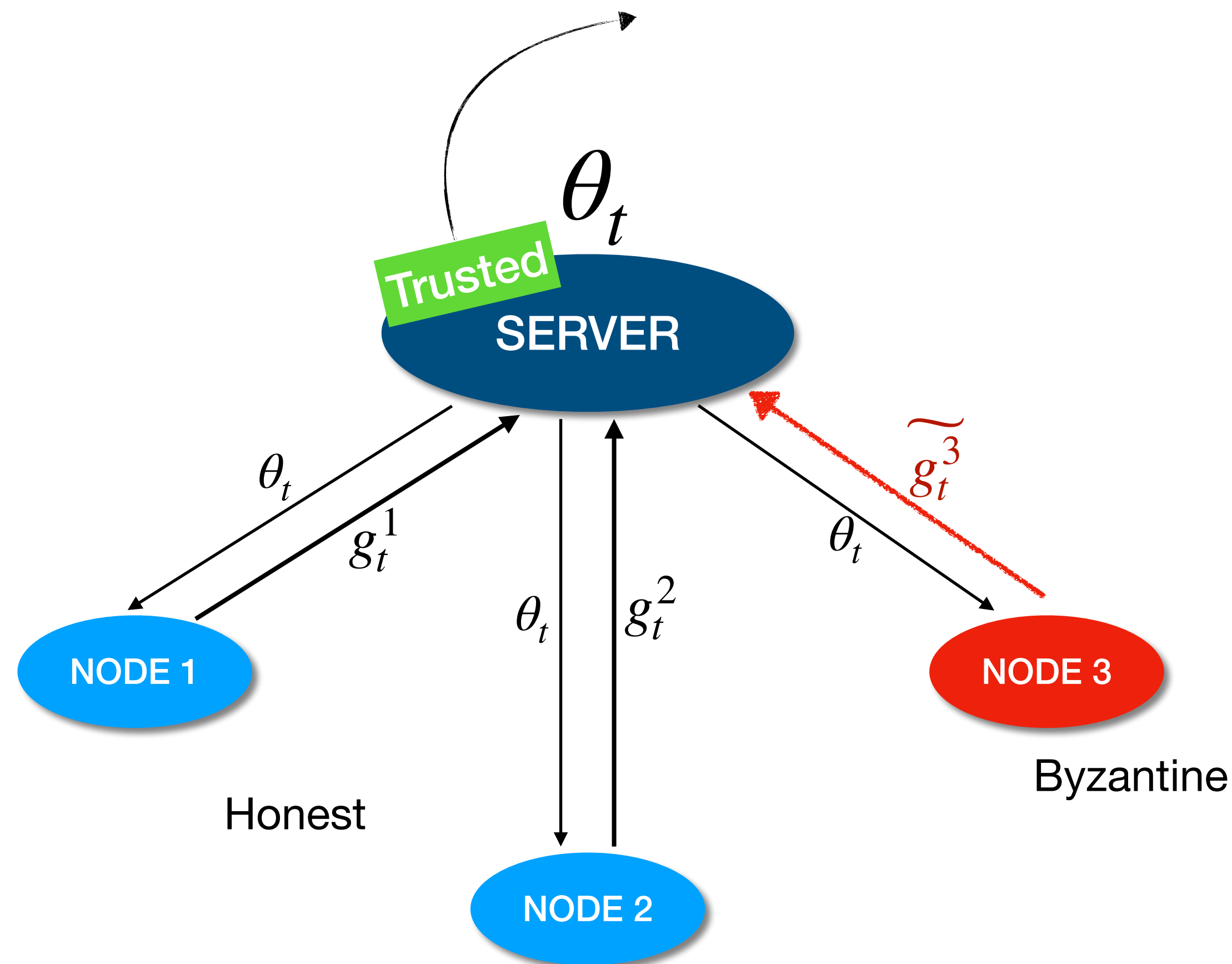
**Feasible** as nodes are assumed
to sample data from the same distribution*

$(f, \epsilon)$**-Resilience**

An algorithm that allows the honest

nodes to compute $\widehat{\theta}$ such that

$$\|\nabla Q\left(\widehat{\theta}\right)\|^2 \le \epsilon$$

$$\theta_t$$

Trusted

SERVER

$\theta_t$

$g_t^1$

$\theta_t$

$g_t^2$

$\widetilde{g_t^3}$

$\theta_t$

NODE 1

NODE 3

NODE 2

Honest

Byzantine

* In general, when data distributions differ, we require $2f$-*redundancy* in data (Gupta and Vaidya, 2021)

# Byzantine Resilience in Distributed Learning

$f$ **out of** $n$ **nodes are Byzantine faulty**

obtain insights that can be applied
to non-trusted server settings

CAN THE SERVER STILL OBTAIN

$$\theta* \in \left( \theta \; ; \; \nabla Q(\theta) = 0 \right)$$
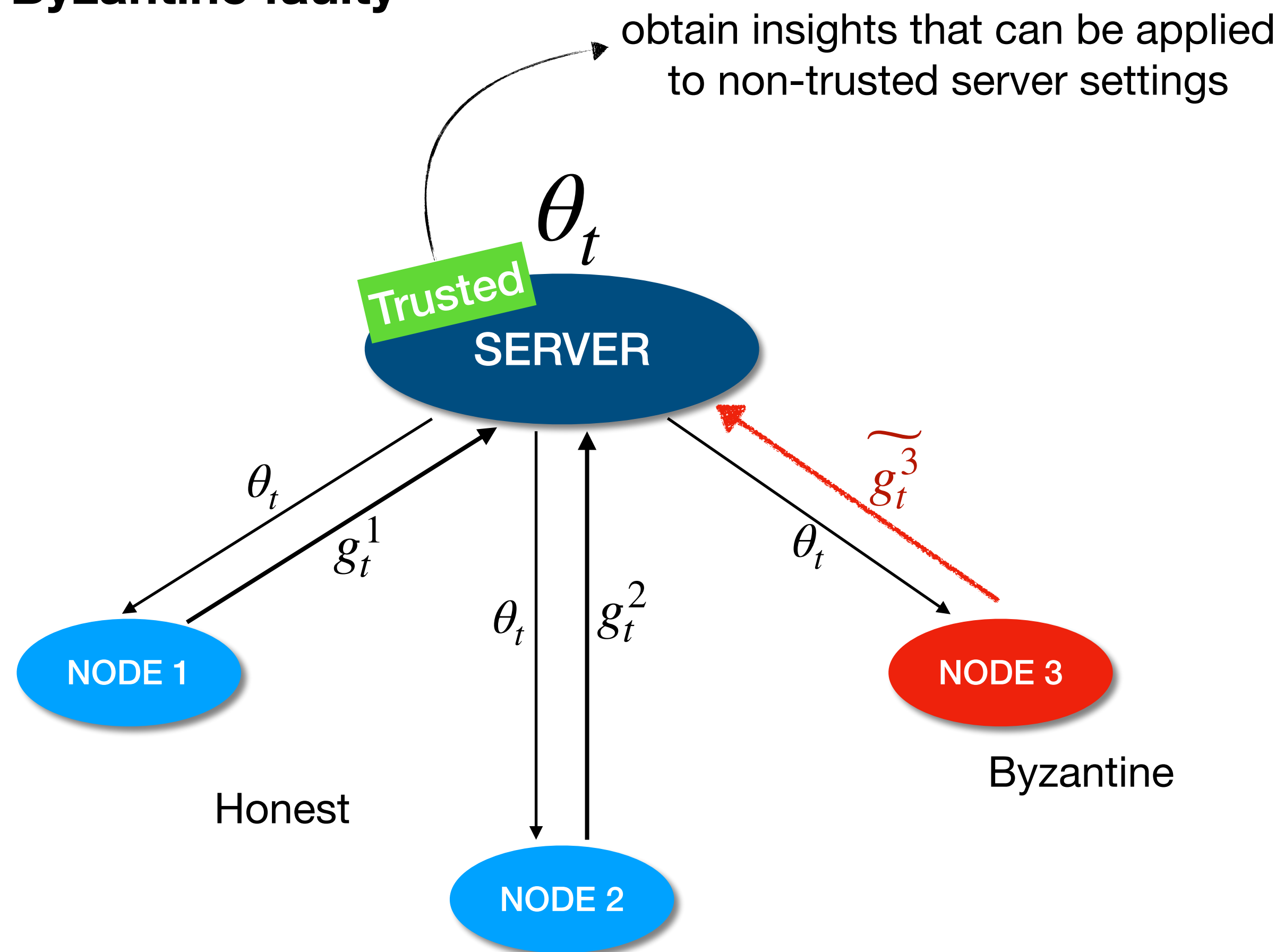
$\theta_t$

Trusted
SERVER

**Feasible** as nodes are assumed
to sample data from the same distribution*

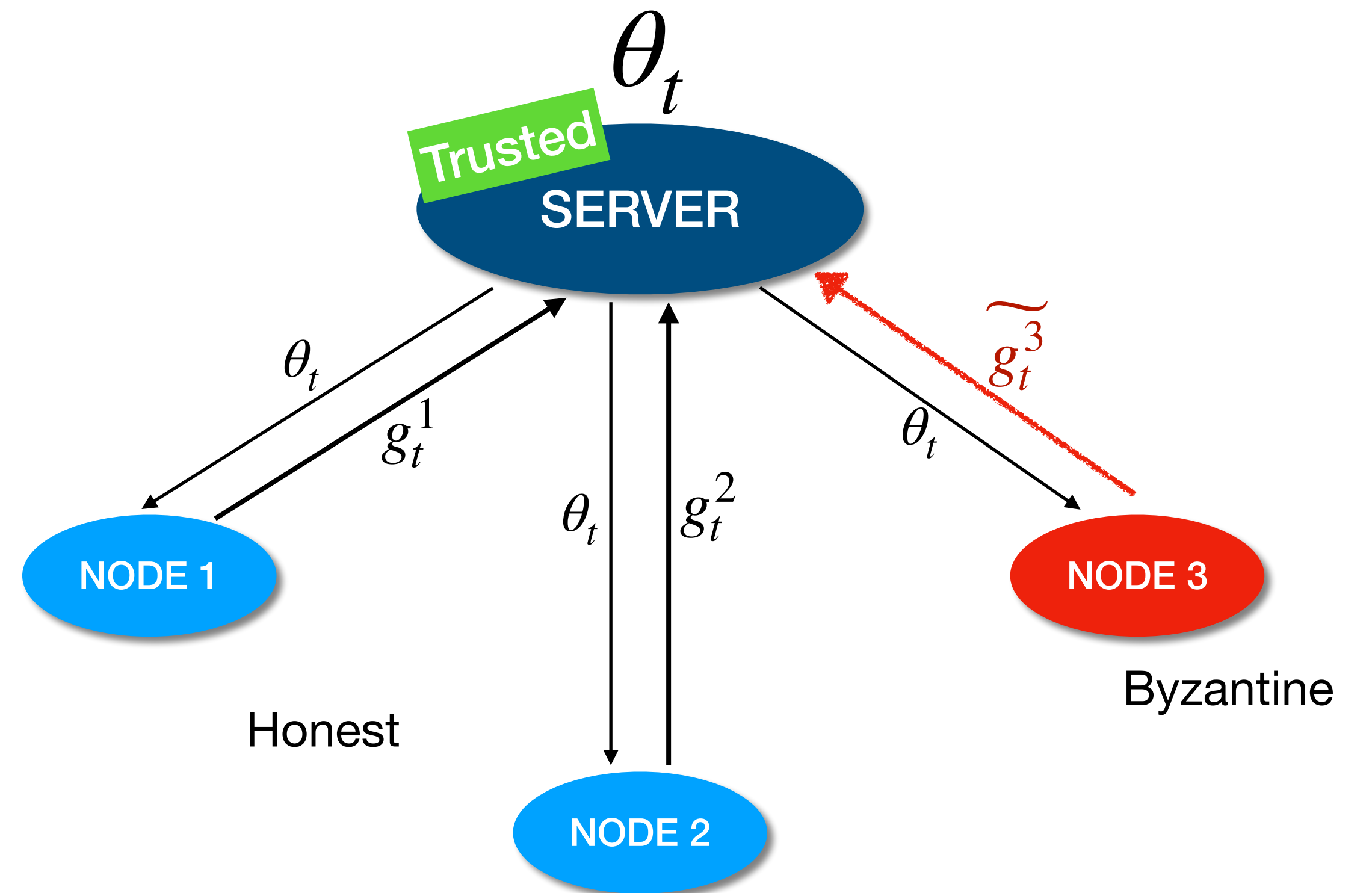$(f, \epsilon)$**-Resilience**

An algorithm that allows the honest

nodes to compute $\widehat{\theta}$ such that

$$\|\nabla Q\left(\widehat{\theta}\right)\|^2 \leq \epsilon$$

$\theta_t$

$g_t^1$

$\widetilde{g_t^3}$

$\theta_t$

$\theta_t$

$g_t^2$
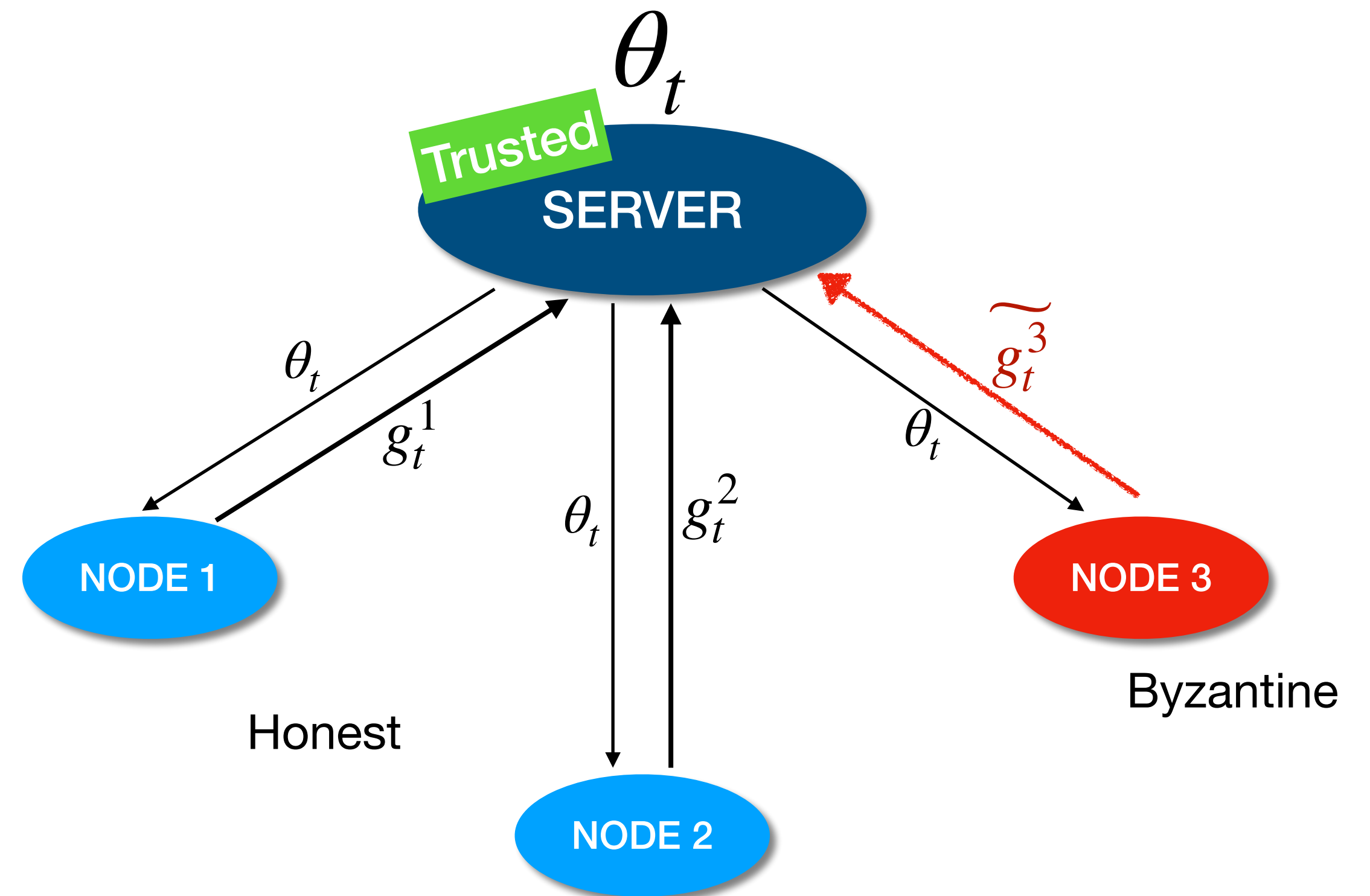
NODE 1

NODE 3

Honest

Byzantine

NODE 2

* In general, when data distributions differ, we require $2f$-*redundancy* in data (Gupta and Vaidya, 2021)

# Making D-SGD Byzantine Resilient

# Making D-SGD Byzantine Resilient

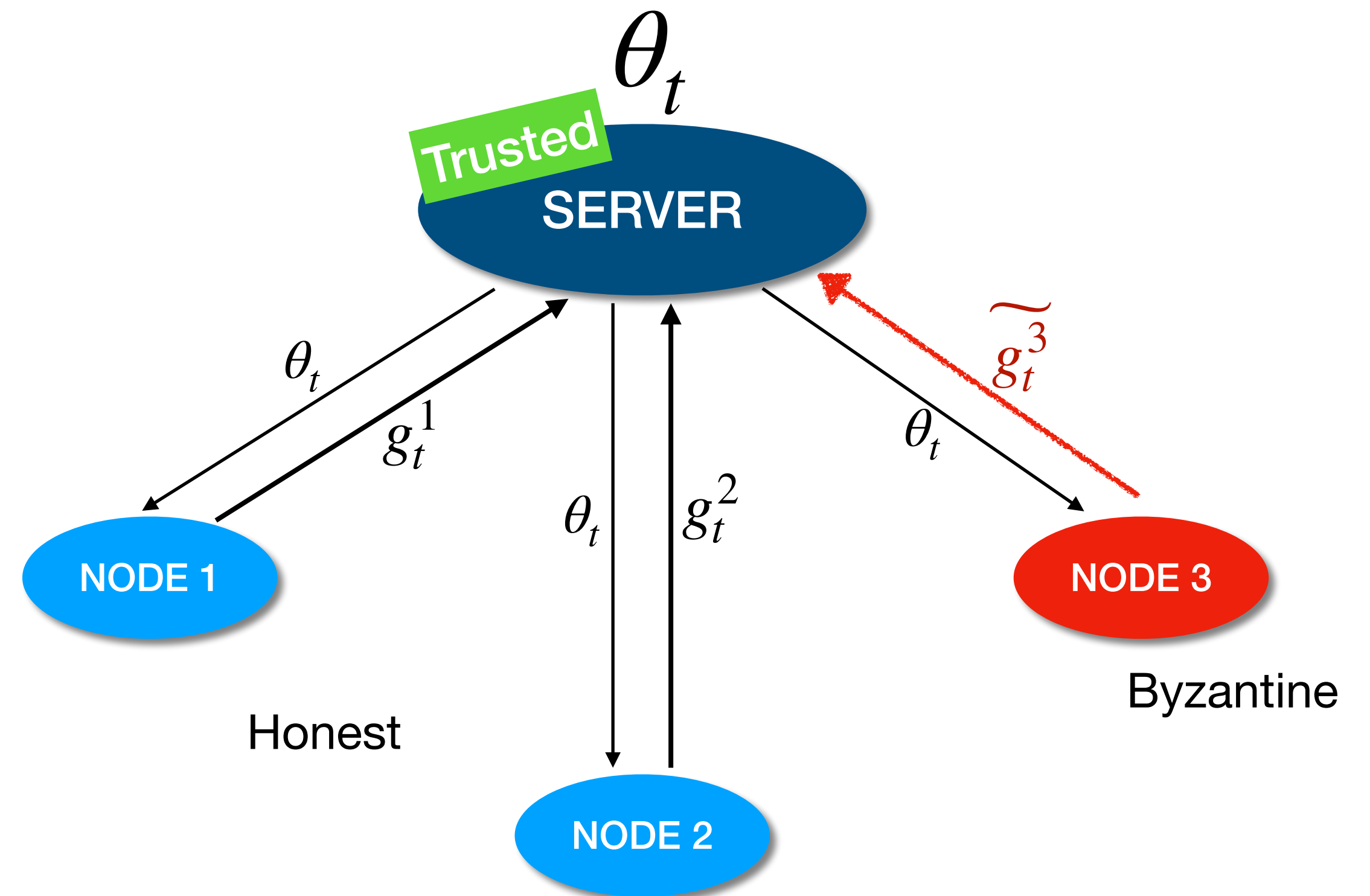**Classical D-SGD -** Server updates $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \; \widehat{g}_t$

# Making D-SGD Byzantine Resilient

**Classical D-SGD -** Server updates $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \, \widehat{g}_t$

We **cannot use**, for updates, $\widehat{g}_t = \dfrac{1}{n} \sum_i g_t^i$
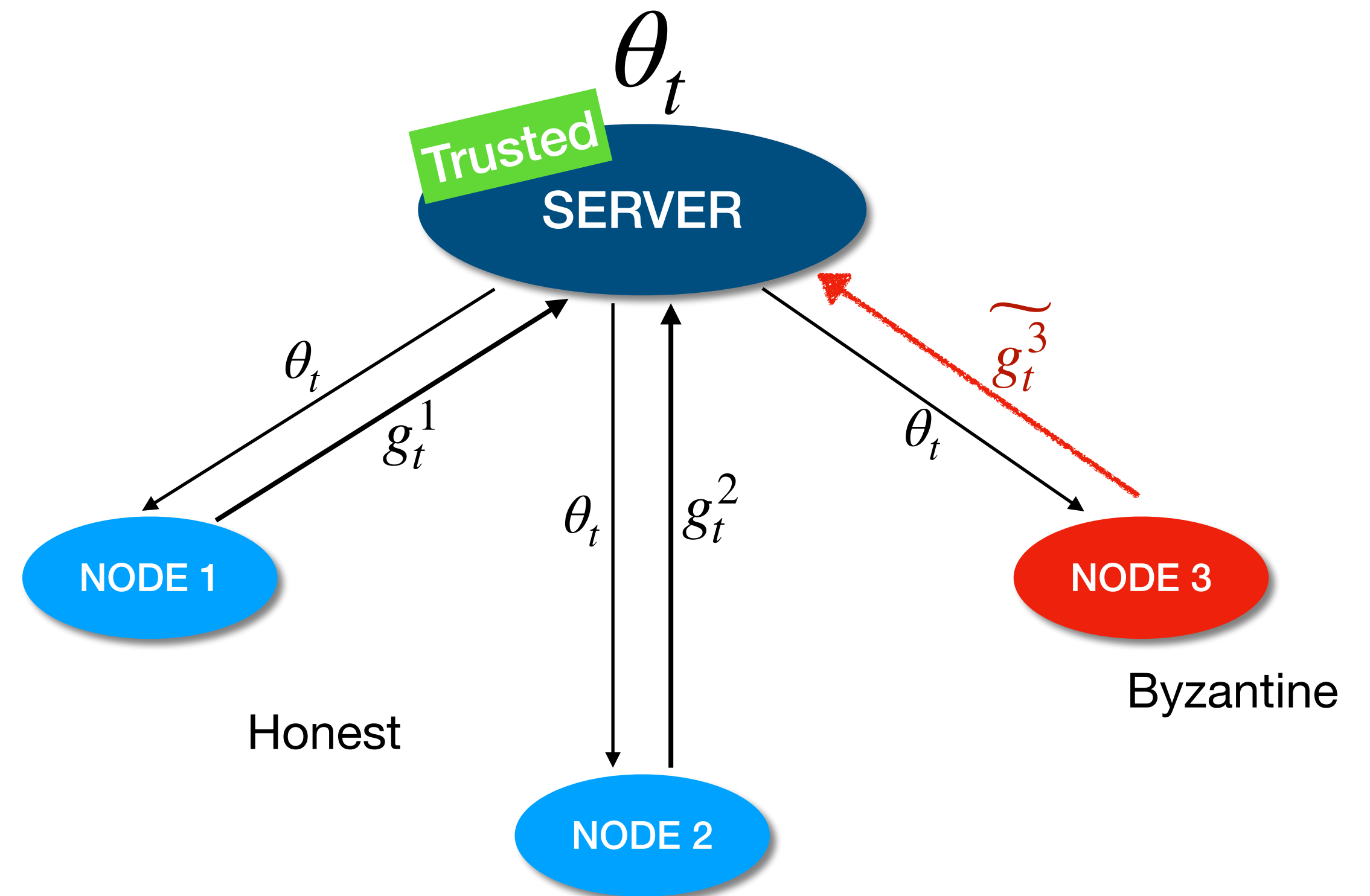
(or any other linear function)

# Making D-SGD Byzantine Resilient

**Classical D-SGD -** Server updates $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \, \widehat{g}_t$

We **cannot use**, for updates, $\widehat{g}_t = \dfrac{1}{n} \sum_i g_t^i$

(or any other linear function)

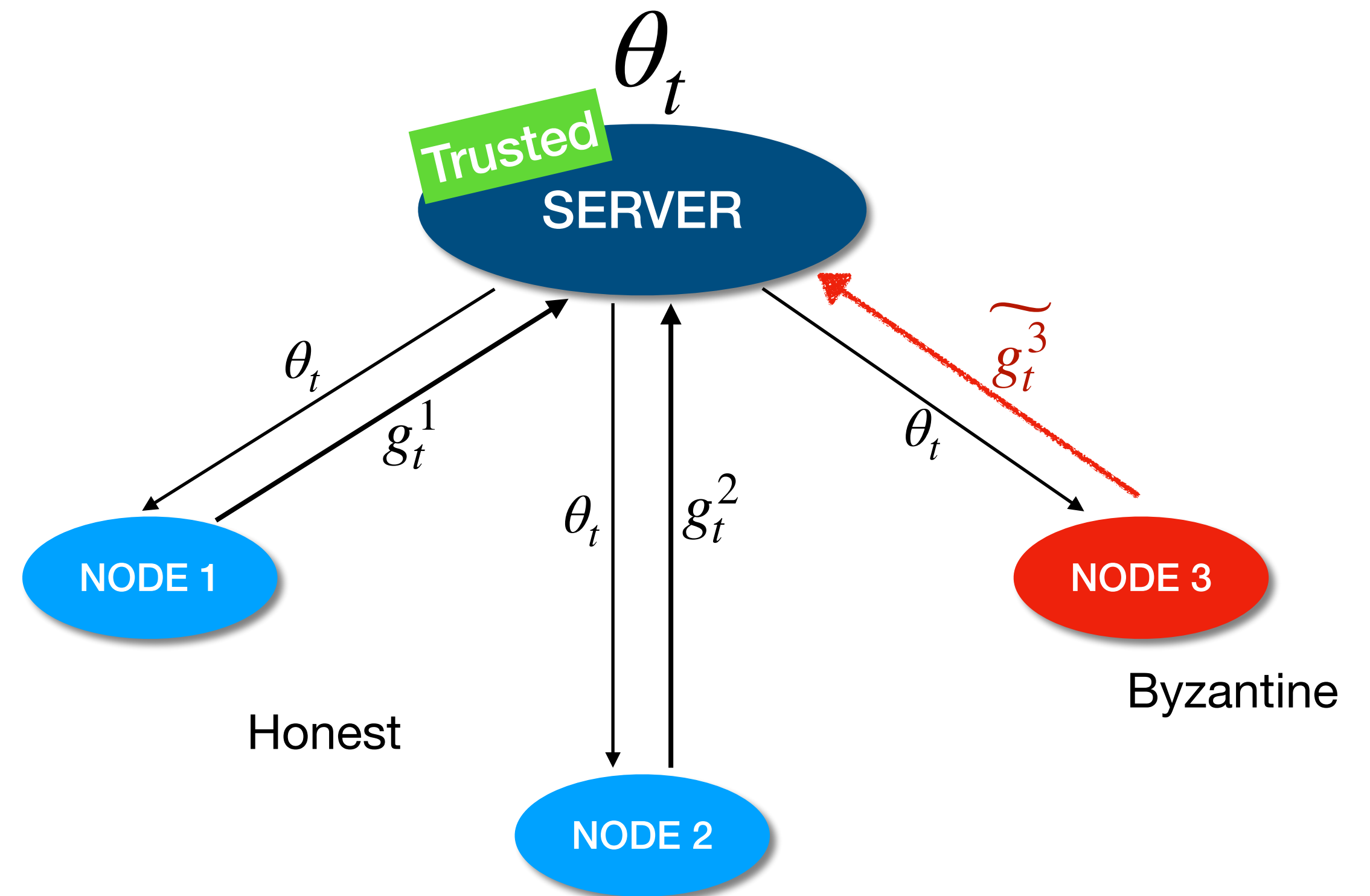Use a **"robust aggregation"** $R_t = F\left(g_t^1, \ldots, g_t^n\right)*$



$\theta_t$

Trusted

SERVER

$\theta_t$

$g_t^1$

$\theta_t$

$g_t^2$

$\widetilde{g_t^3}$

$\theta_t$

NODE 1

NODE 3

Byzantine

Honest

NODE 2

# Making D-SGD Byzantine Resilient
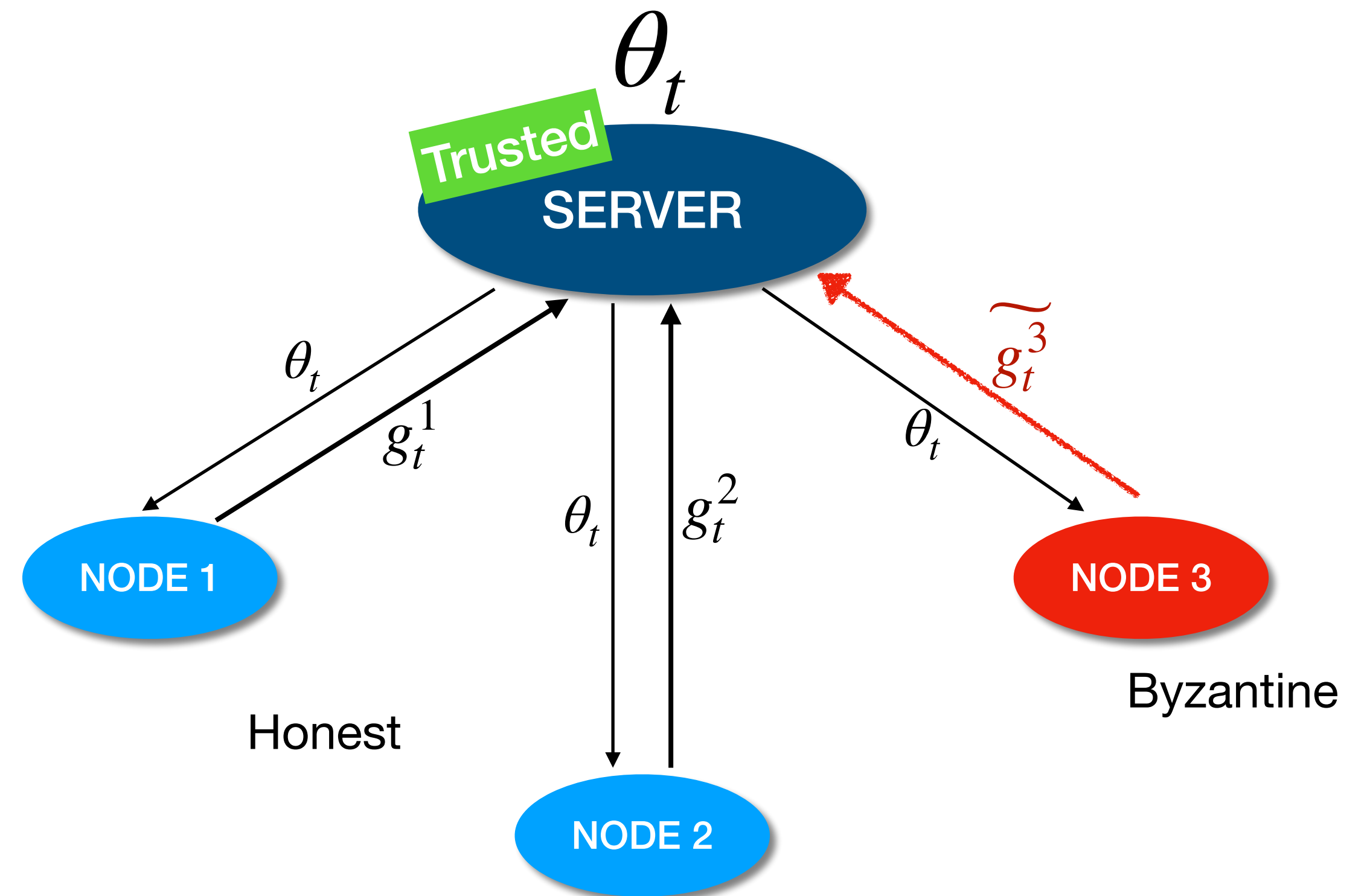
**Classical D-SGD -** Server updates $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \; \widehat{g}_t$

We **cannot use**, for updates, $\;\widehat{g}_t = \dfrac{1}{n} \sum_i g_t^i$

(or any other linear function)

Use a **"robust aggregation"** $\; R_t = F\left(g_t^1, \ldots, g_t^n\right)*$



$\theta_t$

Trusted

SERVER

$\theta_t$

$g_t^1$

NODE 1

$\theta_t$ $\quad g_t^2$

$\theta_t$

$\widetilde{g}_t^3$

NODE 3

Byzantine

Honest

NODE 2

* As the identity of Byzantine nodes is unknown we also denote their gradients by $g$

# Making D-SGD Byzantine Resilient

**Classical D-SGD -** Server updates $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \; \widehat{g}_t$

We **cannot use**, for updates, $\widehat{g}_t = \dfrac{1}{n} \sum_i g_t^i$

(or any other linear function)

Use a **"robust aggregation"** $R_t = F\left(g_t^1, \ldots, g_t^n\right)$*

**Server updates** $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \; R_t$

$$\theta_t$$



* As the identity of Byzantine nodes is unknown we also denote their gradients by $g$

# What is a Robust Aggregation ?

# What is a Robust Aggregation ?

Reduces the impact of Byzantine inputs

# What is a Robust Aggregation ?

Reduces the impact of Byzantine inputs

**Examples:** n = 5 and f = 2

# What is a Robust Aggregation ?

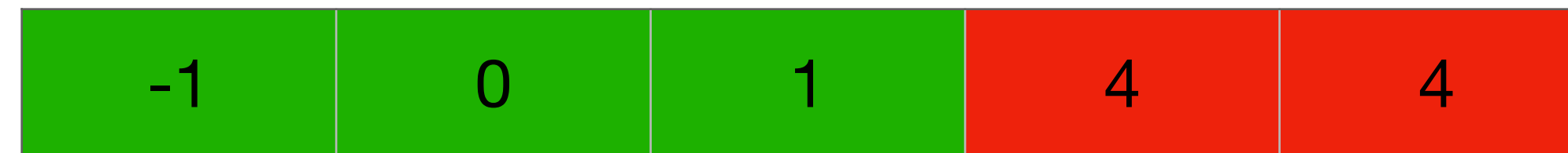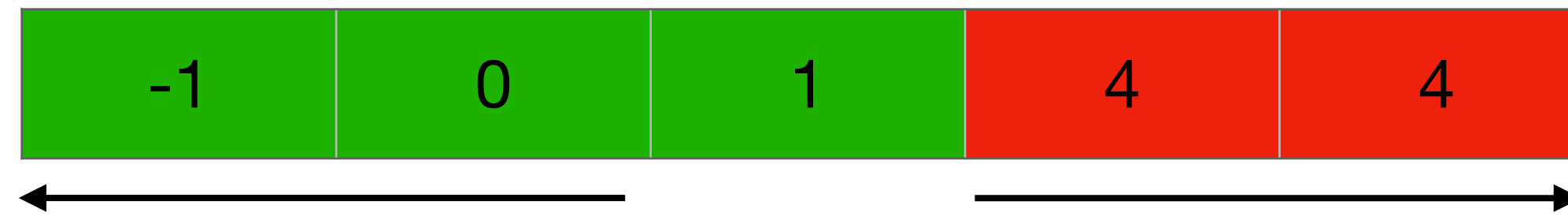Reduces the impact of Byzantine inputs

**Examples:** n = 5 and f = 2

**Trimmed Mean (TM):**

# What is a Robust Aggregation ?

Reduces the impact of Byzantine inputs
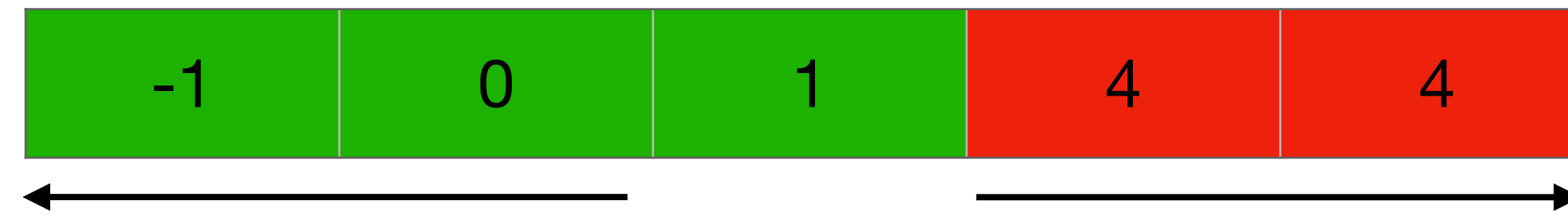
**Examples:** n = 5 and f = 2

**Trimmed Mean (TM):**

| -1 | 0 | 1 | 4 | 4 |
|----|---|---|---|---|

# What is a Robust Aggregation ?

Reduces the impact of Byzantine inputs

**Examples:** n = 5 and f = 2

**Trimmed Mean (TM):**

# What is a Robust Aggregation ?

Reduces the impact of Byzantine inputs

**Examples:** n = 5 and f = 2

**Trimmed Mean (TM):**
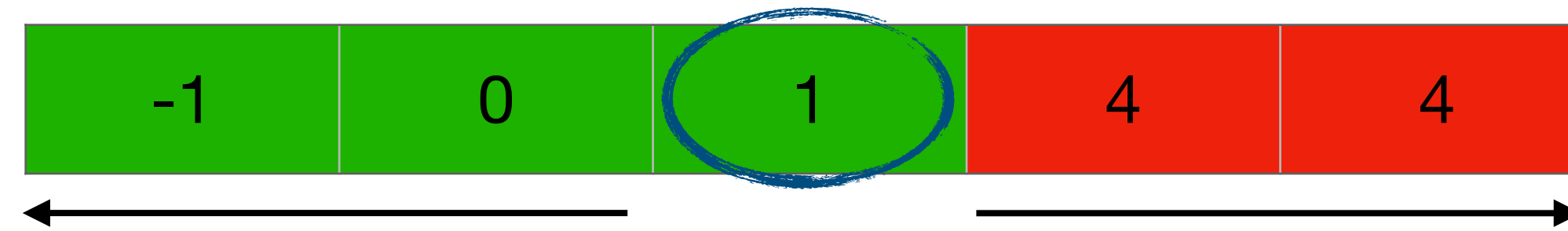
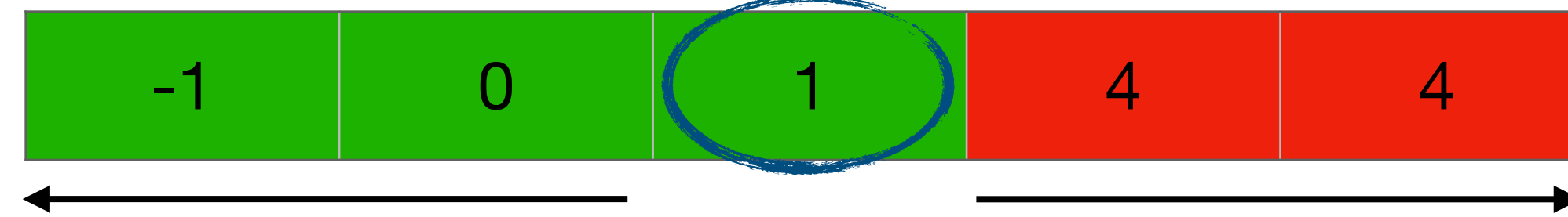| -1 | 0 | 1 | 4 | 4 |
|----|---|---|---|---|

Eliminate f extreme values
Average the remaining

# What is a Robust Aggregation ?

Reduces the impact of Byzantine inputs

**Examples:** n = 5 and f = 2

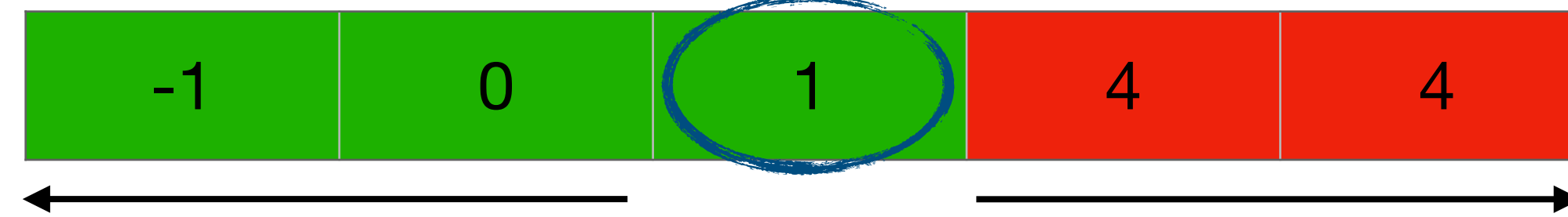**Trimmed Mean (TM):**



Eliminate f extreme values
Average the remaining

# What is a Robust Aggregation ?

Reduces the impact of Byzantine inputs

**Examples:** n = 5 and f = 2

**Trimmed Mean (TM):**

| -1 | 0 | 1 | 4 | 4 |
|---|---|---|---|---|

Eliminate f extreme values
Average the remaining

# What is a Robust Aggregation ?

Reduces the impact of Byzantine inputs

**Examples:** n = 5 and f = 2

**Trimmed Mean (TM):**



| -1 | 0 | 1 | 4 | 4 |

Eliminate f extreme values
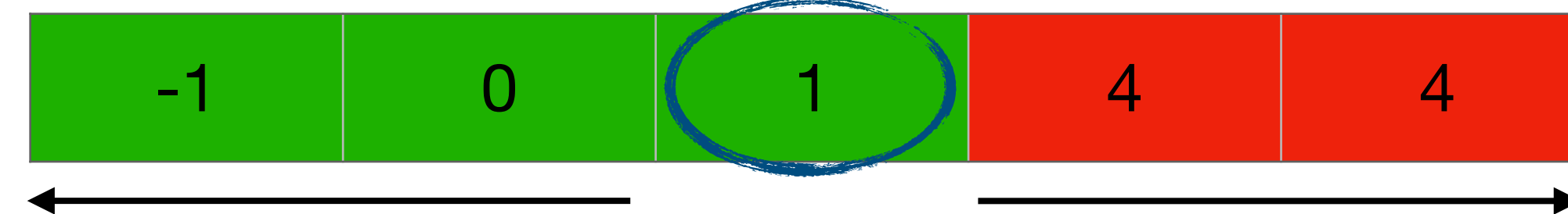Average the remaining

The output lies in the
convex hull of honest inputs

# What is a Robust Aggregation ?

Reduces the impact of Byzantine inputs

**Examples:** n = 5 and f = 2

**Trimmed Mean (TM):**

| -1 | 0 | 1 | 4 | 4 |
|----|---|---|---|---|

Eliminate f extreme values
Average the remaining

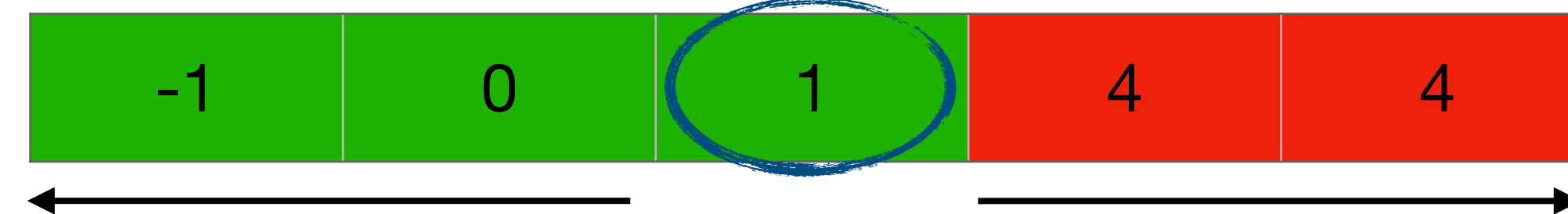The output lies in the
convex hull of honest inputs

**Minimum Diameter Averaging* (MDA):**

# What is a Robust Aggregation ?

## Reduces the impact of Byzantine inputs

**Examples:** $n = 5$ and $f = 2$

**Trimmed Mean (TM):**

| -1 | 0 | 1 | 4 | 4 |
|----|---|---|---|---|

Eliminate f extreme values
Average the remaining

The output lies in the
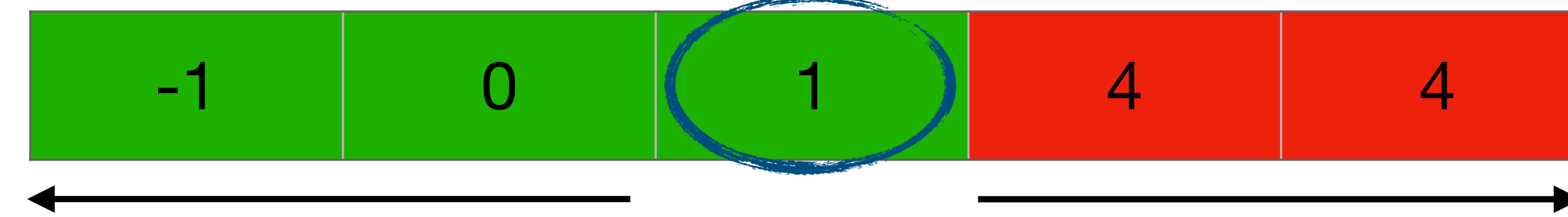convex hull of honest inputs

**Minimum Diameter Averaging* (MDA):**

| -1 | 0 | 1 | 4 | 4 |
|----|---|---|---|---|

# What is a Robust Aggregation ?

Reduces the impact of Byzantine inputs

**Examples:** $n = 5$ and $f = 2$
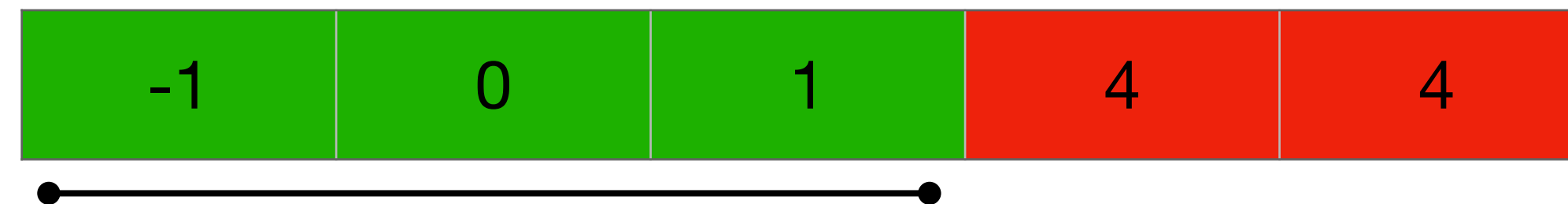
**Trimmed Mean (TM):**

| -1 | 0 | 1 | 4 | 4 |
|----|---|---|---|---|

Eliminate f extreme values
Average the remaining

The output lies in the
convex hull of honest inputs

**Minimum Diameter Averaging* (MDA):**

| -1 | 0 | 1 | 4 | 4 |
|----|---|---|---|---|

# What is a Robust Aggregation ?

Reduces the impact of Byzantine inputs

**Examples:** n = 5 and f = 2

**Trimmed Mean (TM):**

| -1 | 0 | 1 | 4 | 4 |

Eliminate f extreme values
Average the remaining

The output lies in the
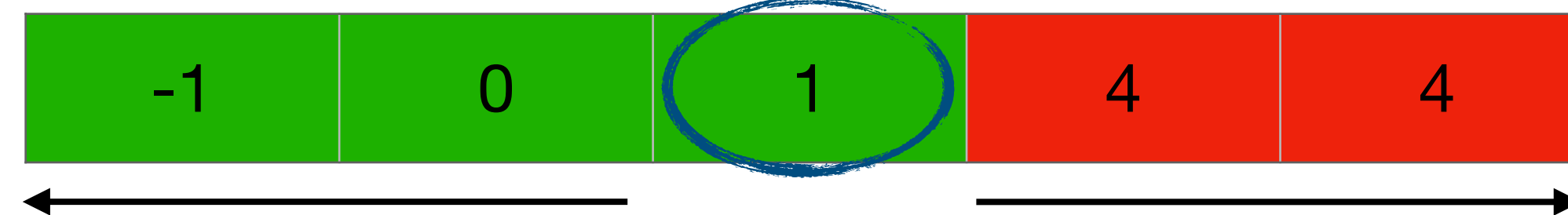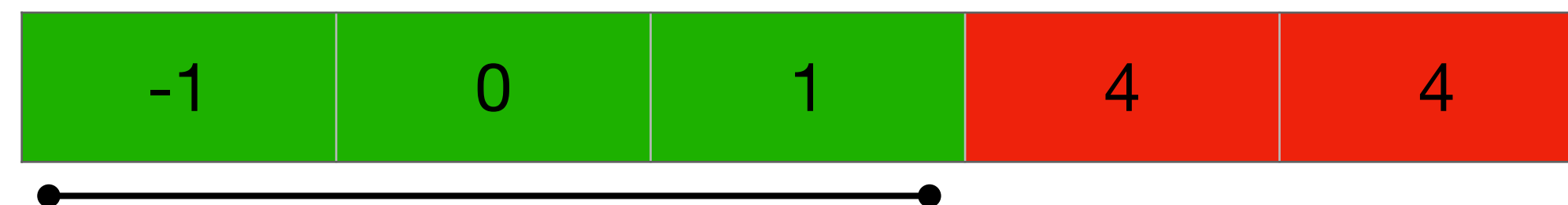convex hull of honest inputs

**Minimum Diameter Averaging* (MDA):**

| -1 | 0 | 1 | 4 | 4 |

Average n-f values
with smallest diameter

# What is a Robust Aggregation ?

## Reduces the impact of Byzantine inputs

**Examples:** $n = 5$ and $f = 2$

**Trimmed Mean (TM):**

| -1 | 0 | 1 | 4 | 4 |
|----|---|---|---|---|

Eliminate f extreme values
Average the remaining

The output lies in the convex hull of honest inputs
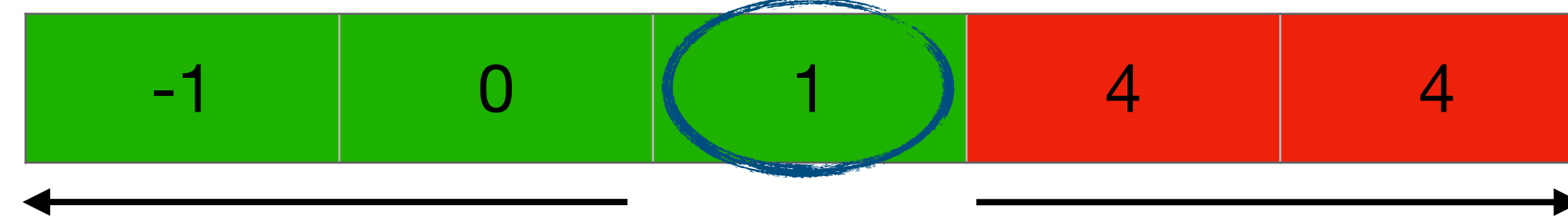
**Minimum Diameter Averaging* (MDA):**

| -1 | 0 | 1 | 4 | 4 |
|----|---|---|---|---|

Average n-f values
with smallest diameter

# What is a Robust Aggregation ?

Reduces the impact of Byzantine inputs

**Examples:** n = 5 and f = 2

**Trimmed Mean (TM):**

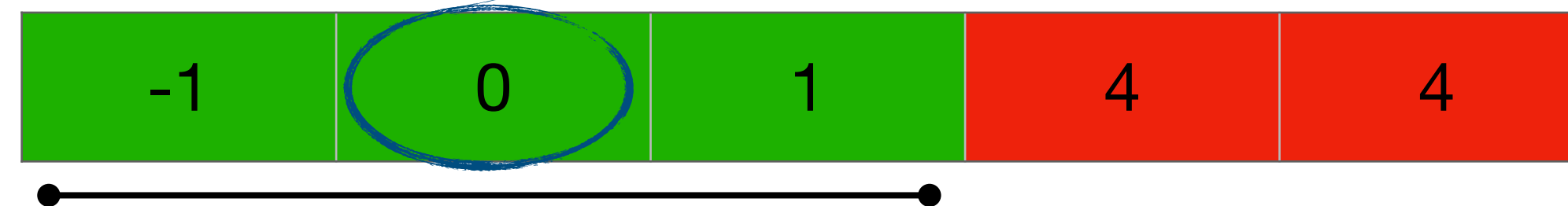| -1 | 0 | 1 | 4 | 4 |
|----|---|---|---|---|

Eliminate f extreme values
Average the remaining

The output lies in the
convex hull of honest inputs

**Minimum Diameter Averaging* (MDA):**
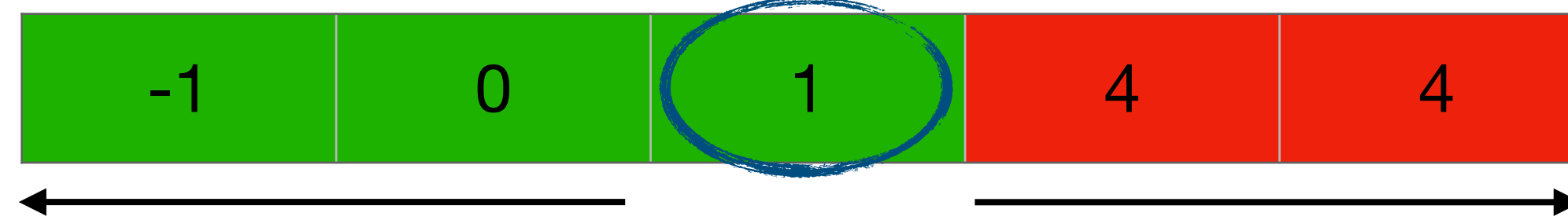
| -1 | 0 | 1 | 4 | 4 |
|----|---|---|---|---|

Average n-f values
with smallest diameter

# What is a Robust Aggregation ?

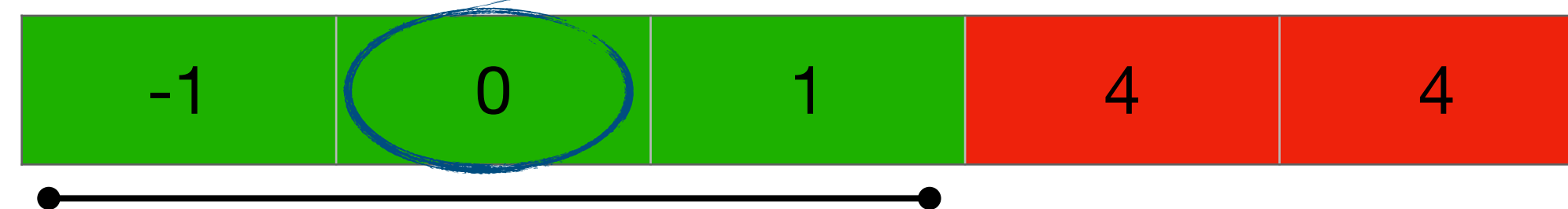Reduces the impact of Byzantine inputs

**Examples:** $n = 5$ and $f = 2$

**Trimmed Mean (TM):**

| -1 | 0 | 1 | 4 | 4 |
|----|---|---|---|---|

Eliminate f extreme values
Average the remaining

The output lies in the convex hull of honest inputs

**Minimum Diameter Averaging* (MDA):**

| -1 | 0 | 1 | 4 | 4 |
|----|---|---|---|---|

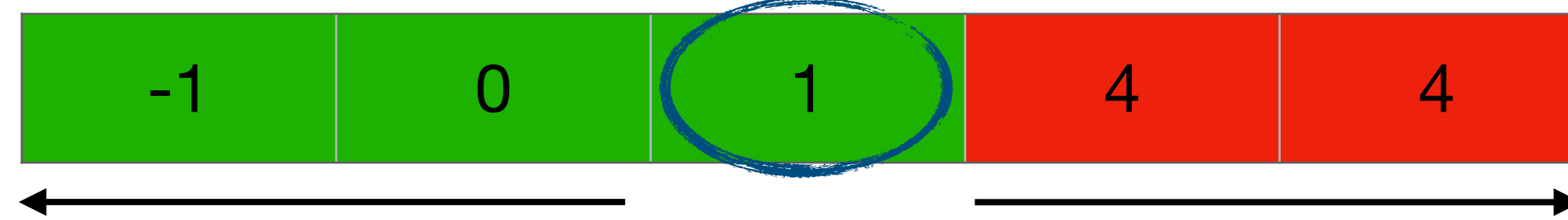Average n-f values
with smallest diameter

The output is the best possible estimate of the honest average

# What is a Robust Aggregation ?

Reduces the impact of Byzantine inputs

**Examples:** n = 5 and f = 2

**Trimmed Mean (TM):**

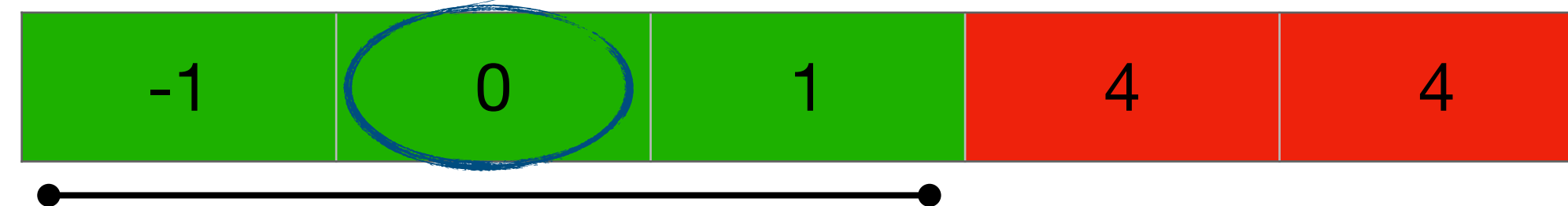| -1 | 0 | 1 | 4 | 4 |
|----|---|---|---|---|

Eliminate f extreme values
Average the remaining

The output lies in the
convex hull of honest inputs

**Minimum Diameter Averaging* (MDA):**

| -1 | 0 | 1 | 4 | 4 |
|----|---|---|---|---|

Average n-f values
with smallest diameter

The output is the best possible
estimate of the honest average

* Very costly in high-dimension

# Does Robust Aggregation Suffice

# Does Robust Aggregation Suffice

It depends

# Does Robust Aggregation Suffice

It depends

| Aggregation Rule/Scheme | Additional Assumption |
|---|---|
| Trimmed Mean<br>*(Yin et al., 2018)* | Sub-exponential stochasticity |
| Geometric Median<br>*(Chen et al., 2017)* | Sub-exponential stochasticity |
| Krum/multi-Krum<br>*(Blanchard et al., 2017)* | Vanishing variance |
| Monitoring temporal averages of gradients<br>*(Alistarh et al., 2018)* | (a.s.) Absolutely boundedness |

# Does Robust Aggregation Suffice

It depends

| Aggregation Rule/Scheme | Additional Assumption |
|---|---|
| Trimmed Mean<br>*(Yin et al., 2018)* | Sub-exponential stochasticity |
| Geometric Median<br>*(Chen et al., 2017)* | Sub-exponential stochasticity |
| Krum/multi-Krum<br>*(Blanchard et al., 2017)* | Vanishing variance |
| Monitoring temporal averages of gradients<br>*(Alistarh et al., 2018)* | (a.s.) Absolutely boundedness |

Most give resilience only against a small fraction of Byzantines << 1/2
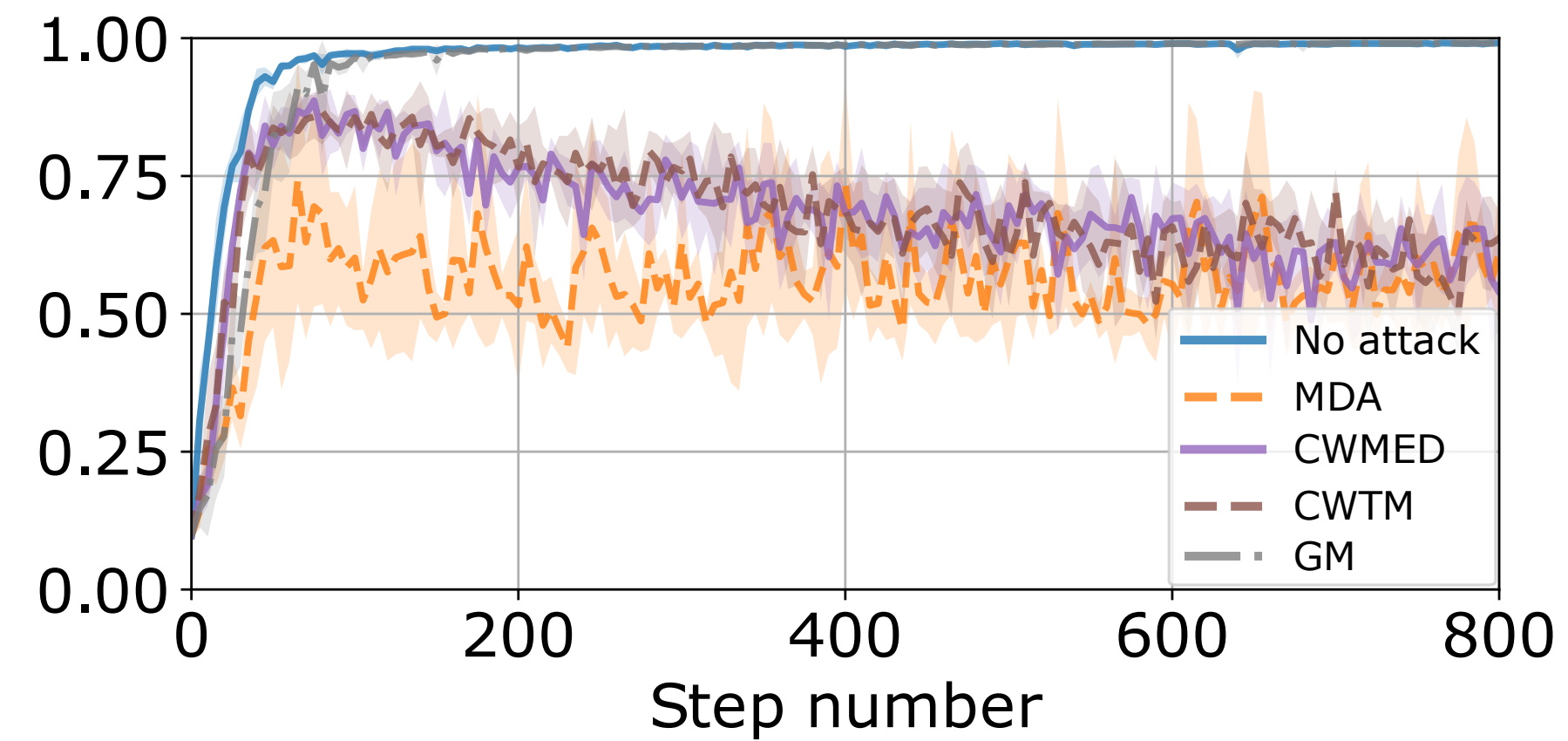
# Vulnerable due to Stringent Assumptions

# Vulnerable due to Stringent Assumptions

\* MNIST classification task

# Vulnerable due to Stringent Assumptions

5 out of 25 nodes are Byzantine faulty

\* MNIST classification task

# Vulnerable due to Stringent Assumptions

5 out of 25 nodes are Byzantine faulty



Legend:
- No attack
- MDA
- CWMED
- CWTM
- GM

X-axis: Step number (0 to 800)
Y-axis: 0.00 to 1.00

* MNIST classification task

# Vulnerable due to Stringent Assumptions

5 out of 25 nodes are Byzantine faulty



**Label-flipping**

* MNIST classification task

# Vulnerable due to Stringent Assumptions

5 out of 25 nodes are Byzantine faulty



**Label-flipping**

* MNIST classification task

# Vulnerable due to Stringent Assumptions

5 out of 25 nodes are Byzantine faulty



**Label-flipping**



**Little is enough** *(Baruch et al., 2019)*

* MNIST classification task
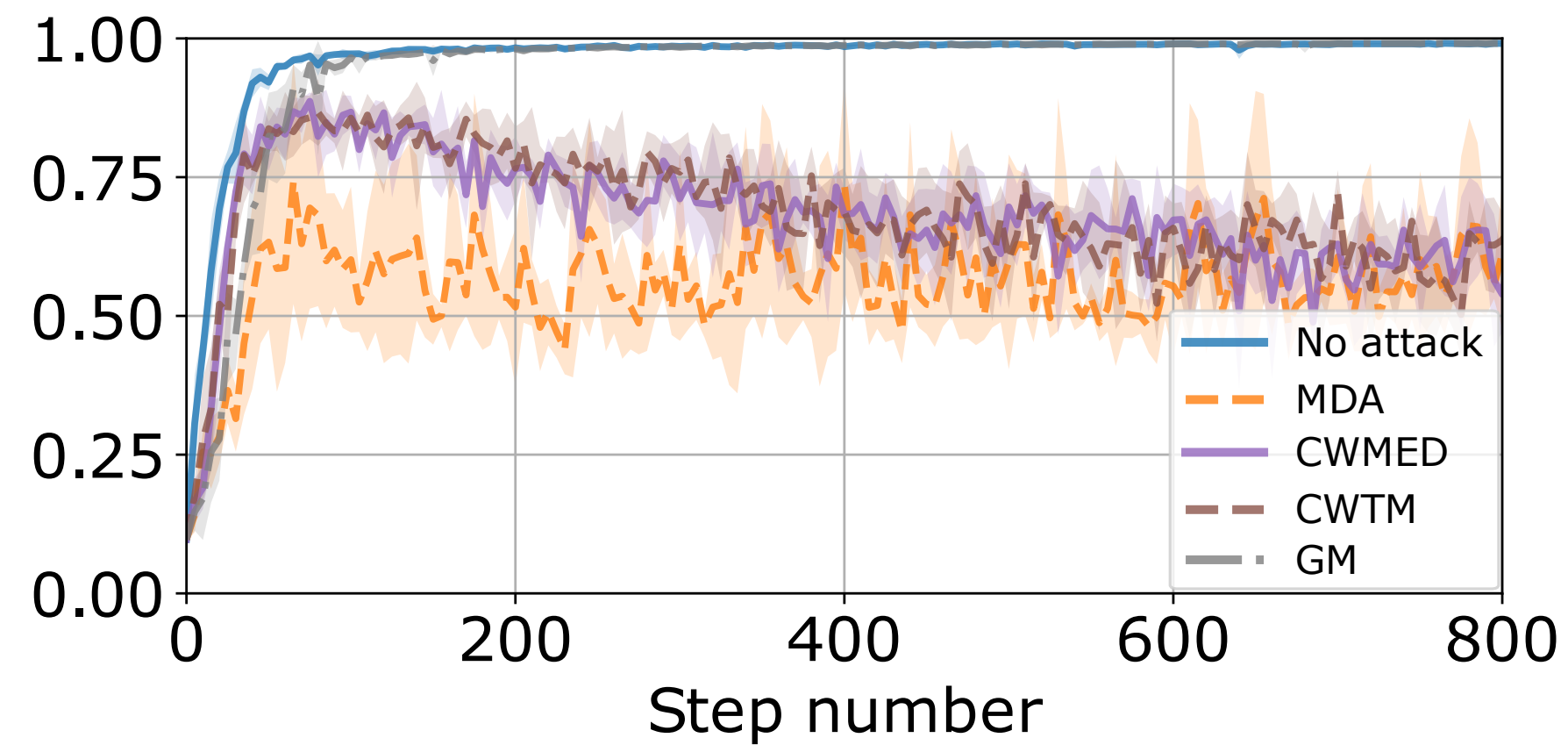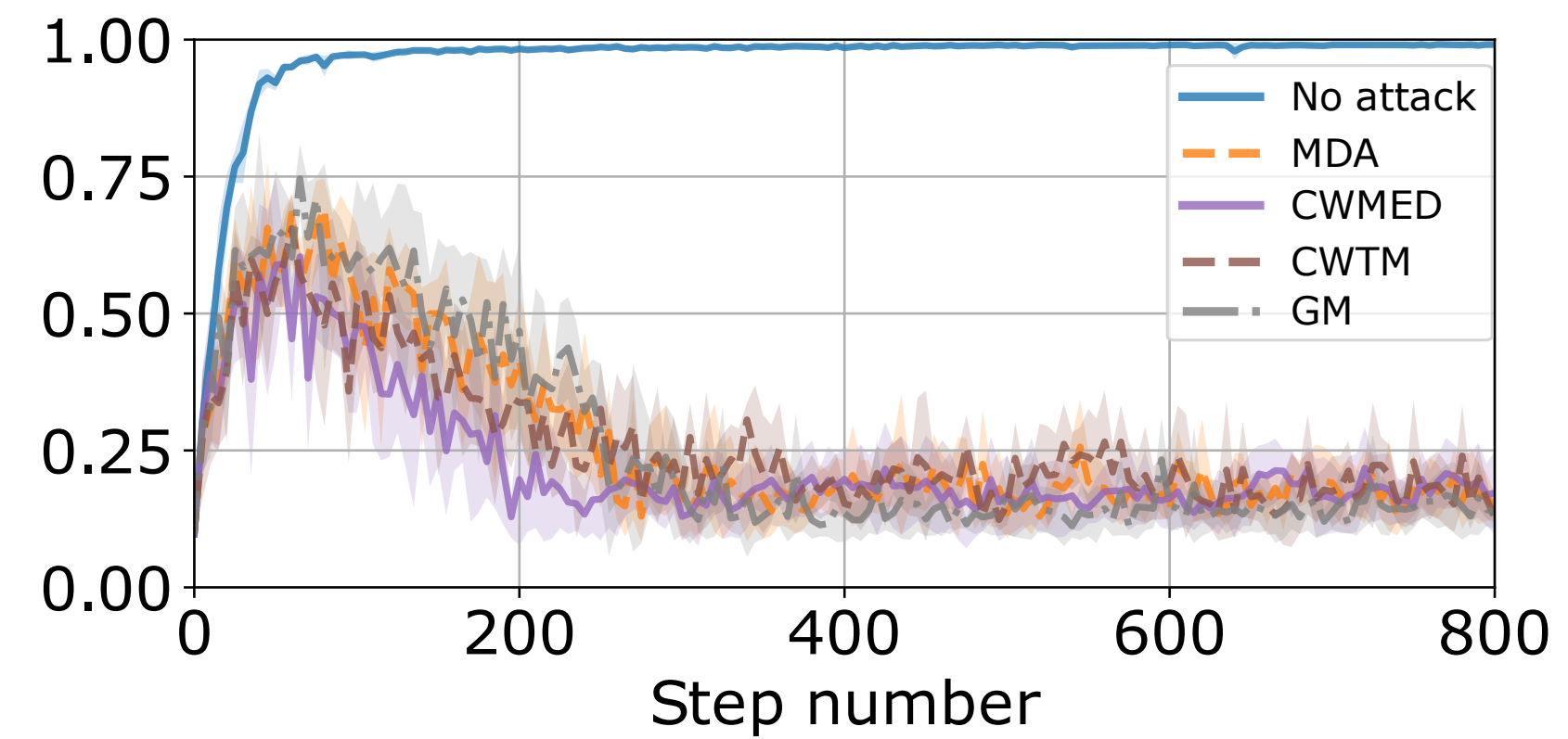
# Vulnerable due to Stringent Assumptions

5 out of 25 nodes are Byzantine faulty



**Label-flipping**



**Little is enough** *(Baruch et al., 2019)*

Memoryless robust aggregation need not be sufficient *(Karimireddy et al., 2021)*

* MNIST classification task

# What is the Problem?

Non-trivial variance of stochastic gradients

# What is the Problem?

Non-trivial variance of stochastic gradients

•

# What is the Problem?

Non-trivial variance of stochastic gradients

$\theta_t$ •

# What is the Problem?

Non-trivial variance of stochastic gradients

$\theta_t$

# What is the Problem?

Non-trivial variance of stochastic gradients

# What is the Problem?

Non-trivial variance of stochastic gradients

True gradient $\nabla Q(\theta_t)$

$\theta_t$

# What is the Problem?

Non-trivial variance of stochastic gradients

True gradient $\nabla Q(\theta_t)$

$\theta_t$

$\theta*$

# What is the Problem?

Non-trivial variance of stochastic gradients

True gradient $\nabla Q(\theta_t)$

$\theta_t$

$\theta*$

# What is the Problem?

Non-trivial variance of stochastic gradients

$g_t^1$

True gradient $\nabla Q(\theta_t)$
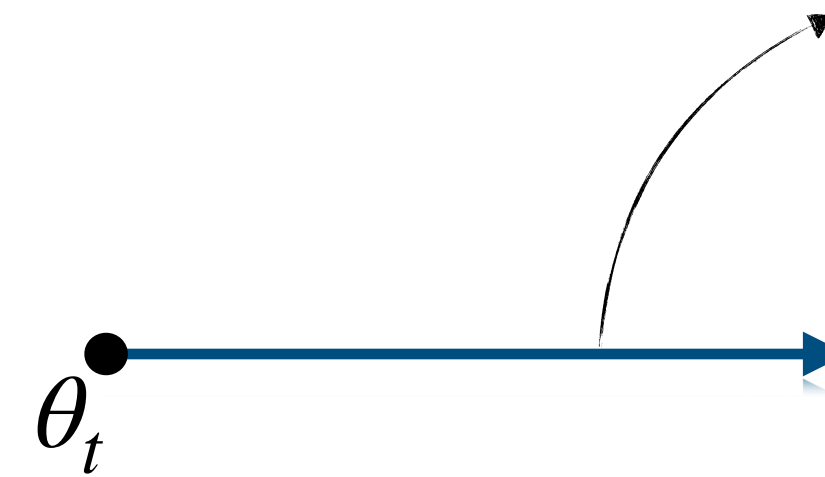
$\theta_t$

$\theta*$

# What is the Problem?

Non-trivial variance of stochastic gradients

# What is the Problem?

Non-trivial variance of stochastic gradients

# What is the Problem?

Non-trivial variance of stochastic gradients

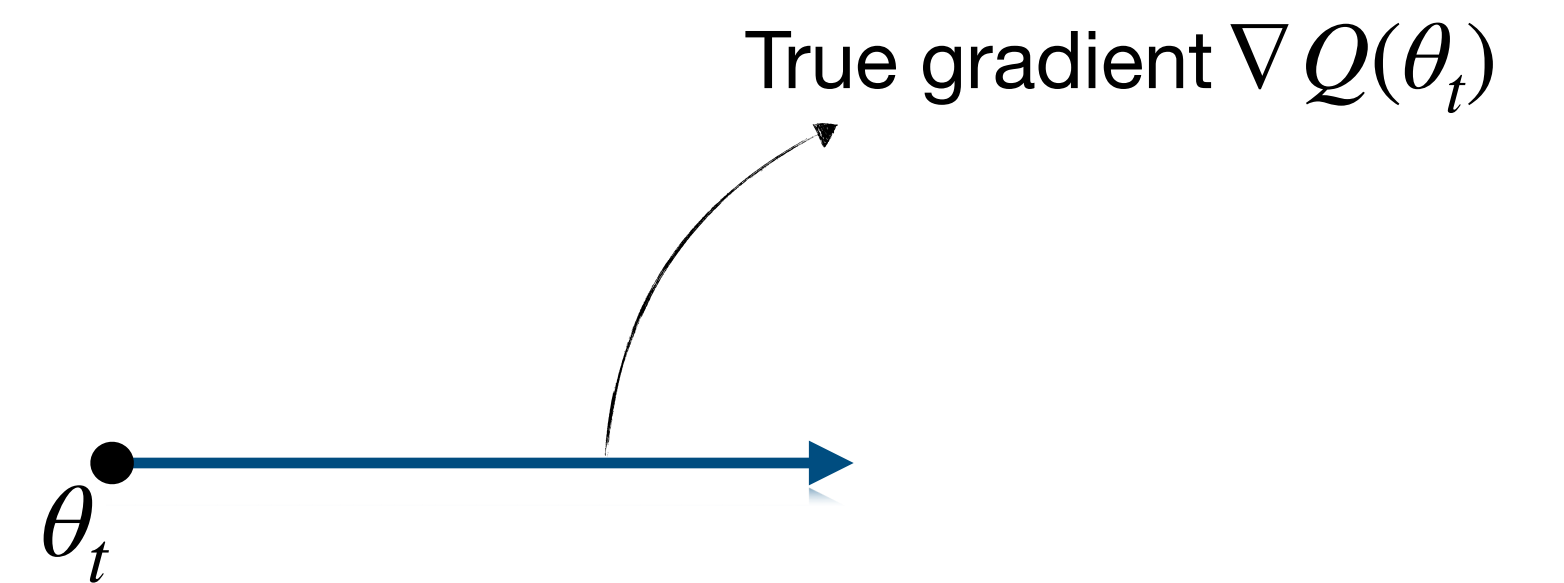# What is the Problem?

Non-trivial variance of stochastic gradients

# What is the Problem?

Non-trivial variance of stochastic gradients

# What is the Problem?

Non-trivial variance of stochastic gradients

# What is the Problem?

Non-trivial variance of stochastic gradients

# What is the Problem?

Non-trivial variance of stochastic gradients

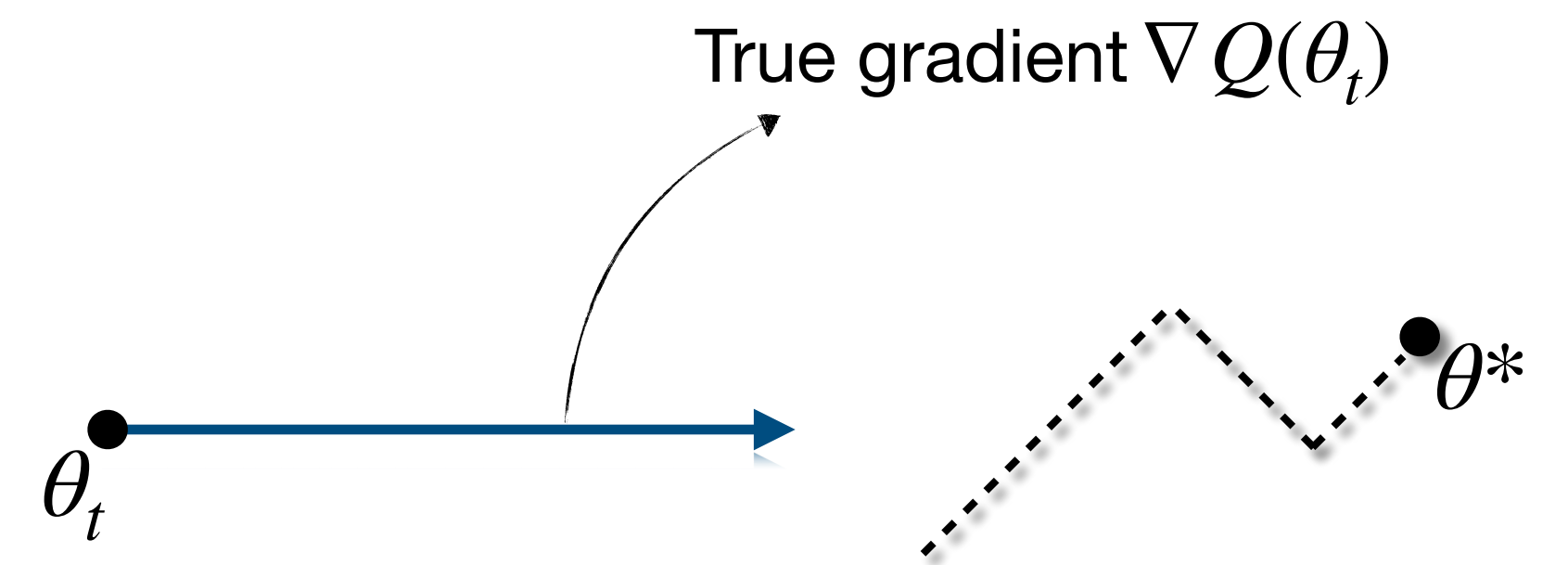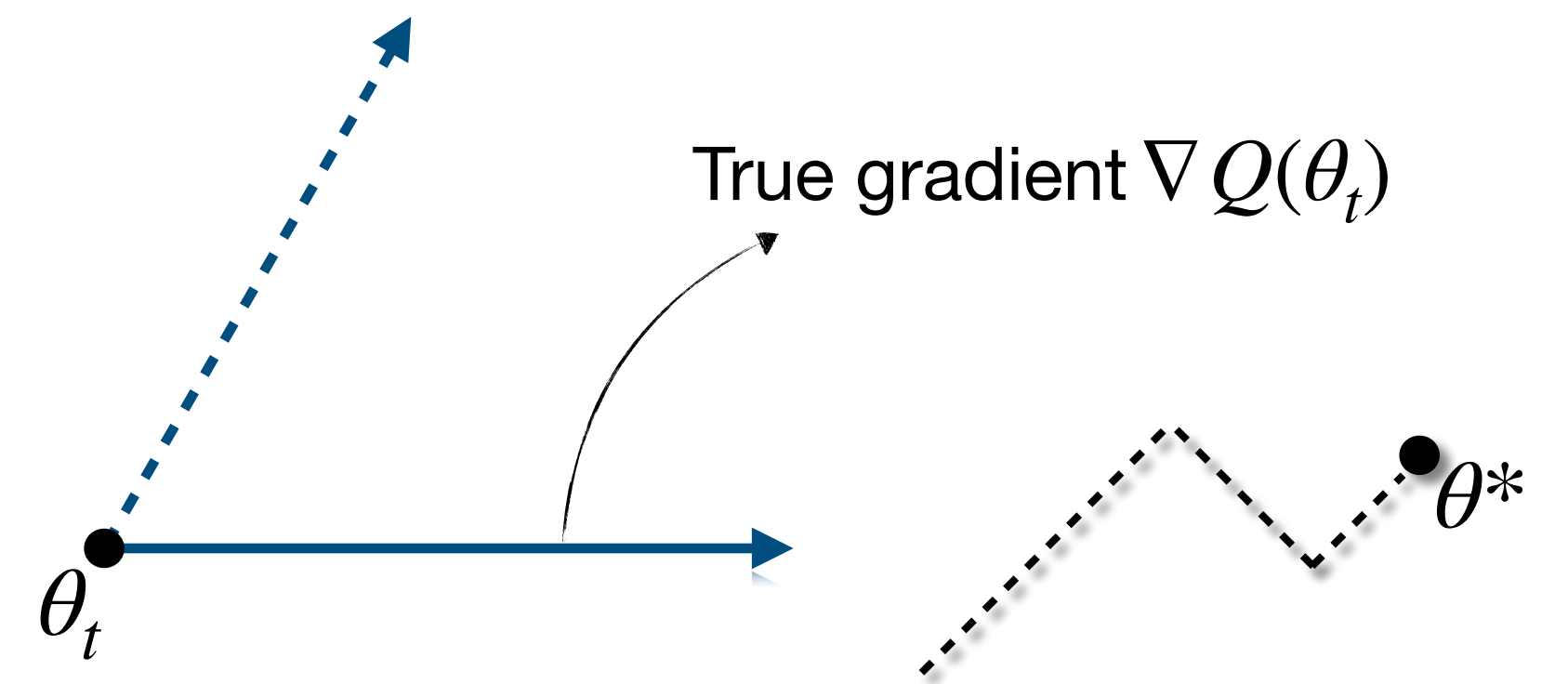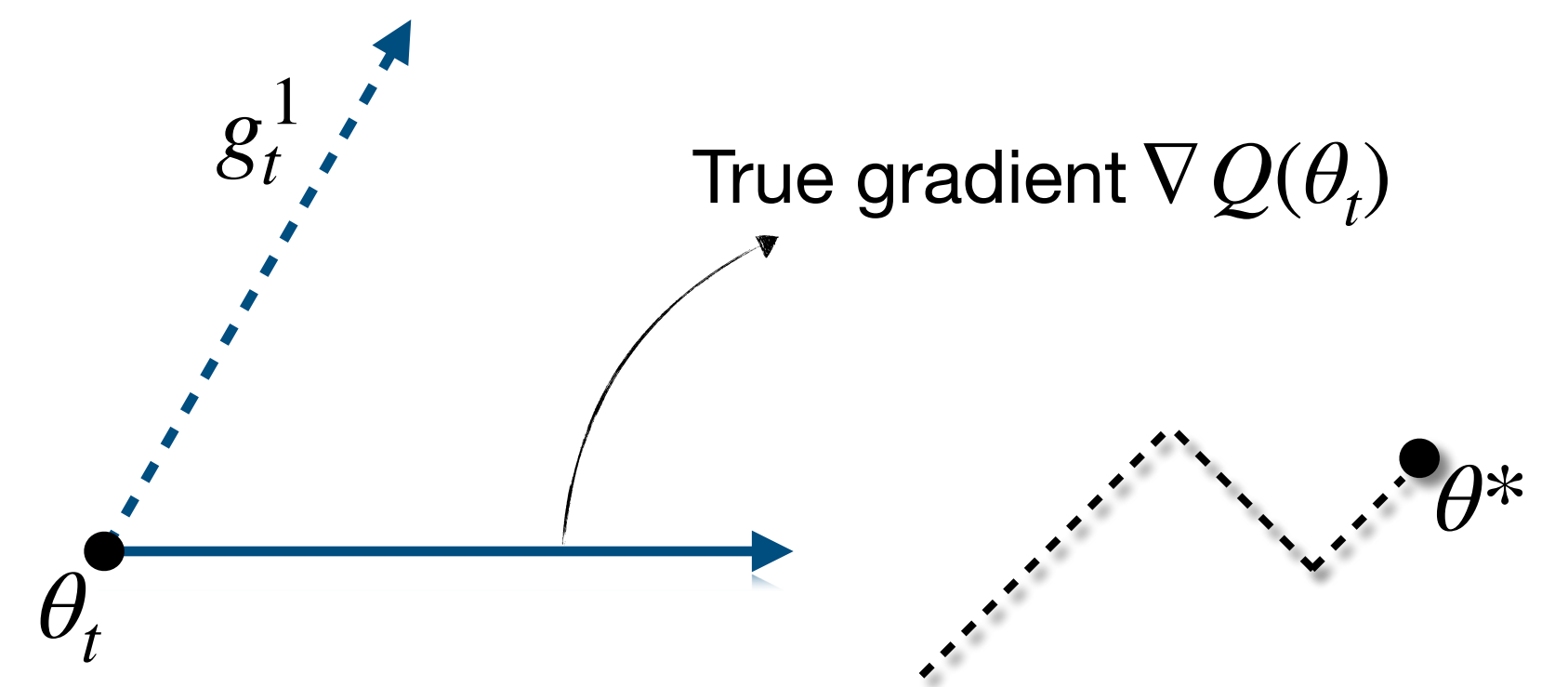# What is the Problem?

Non-trivial variance of stochastic gradients

# What is the Problem?
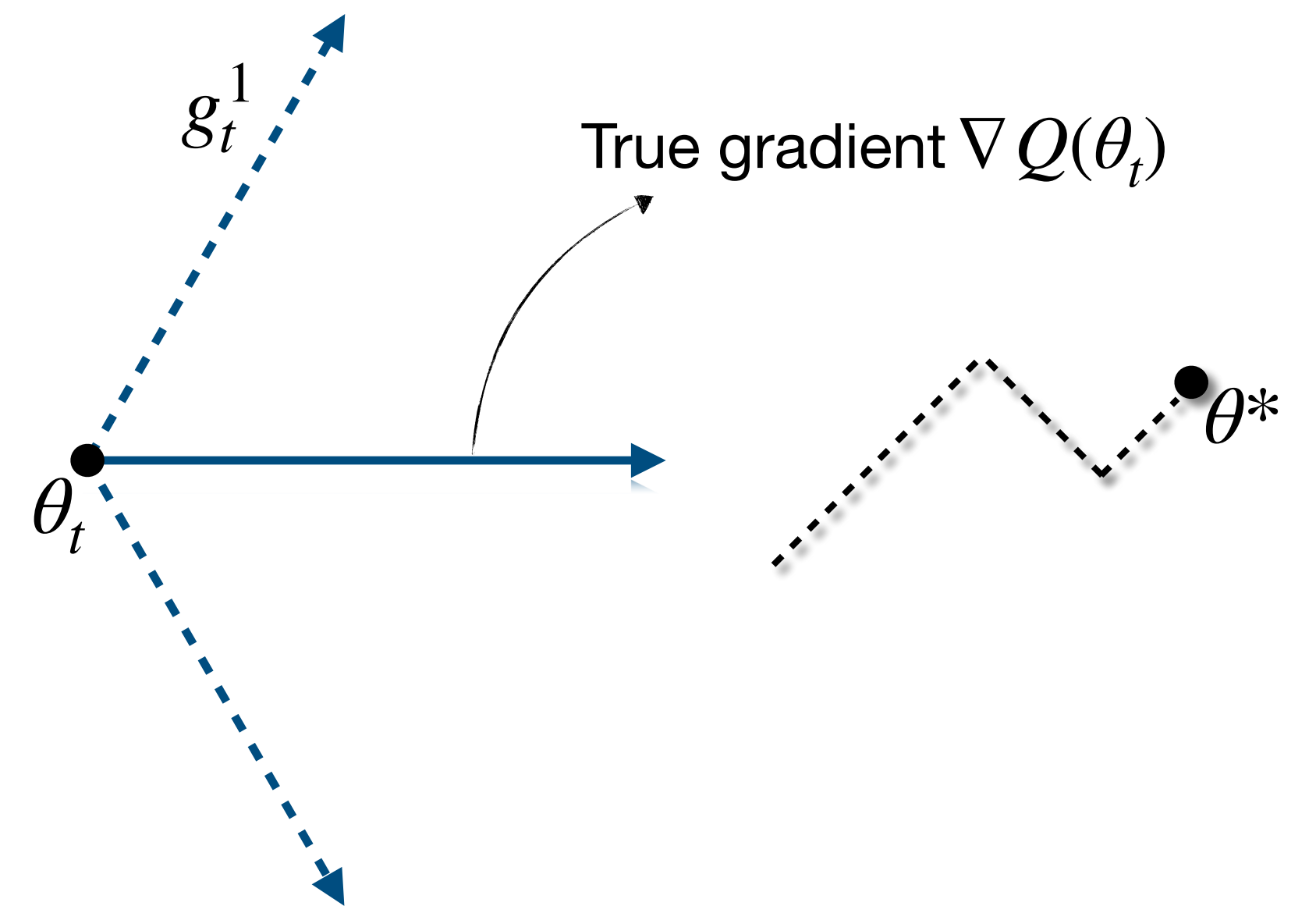
Non-trivial variance of stochastic gradients
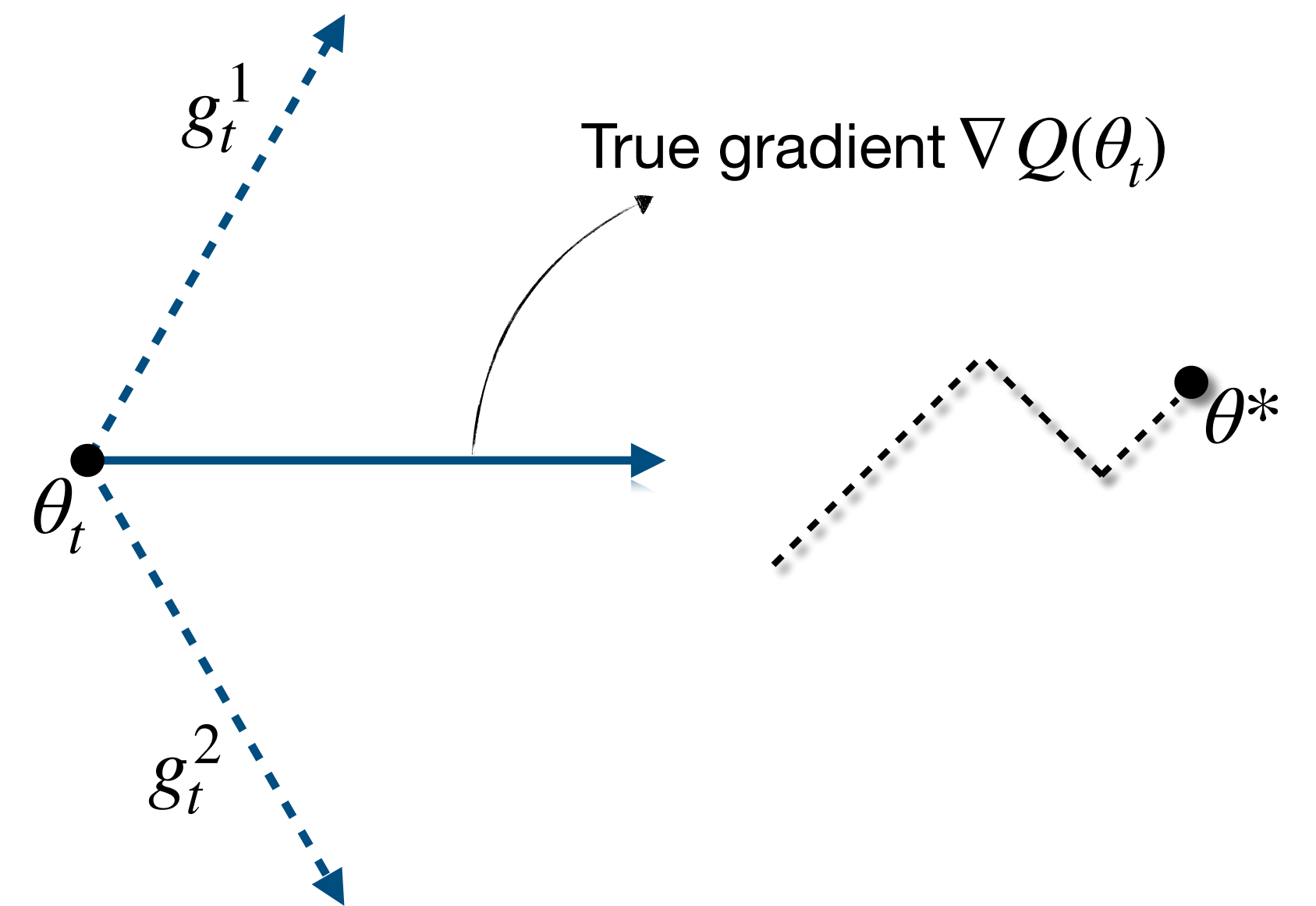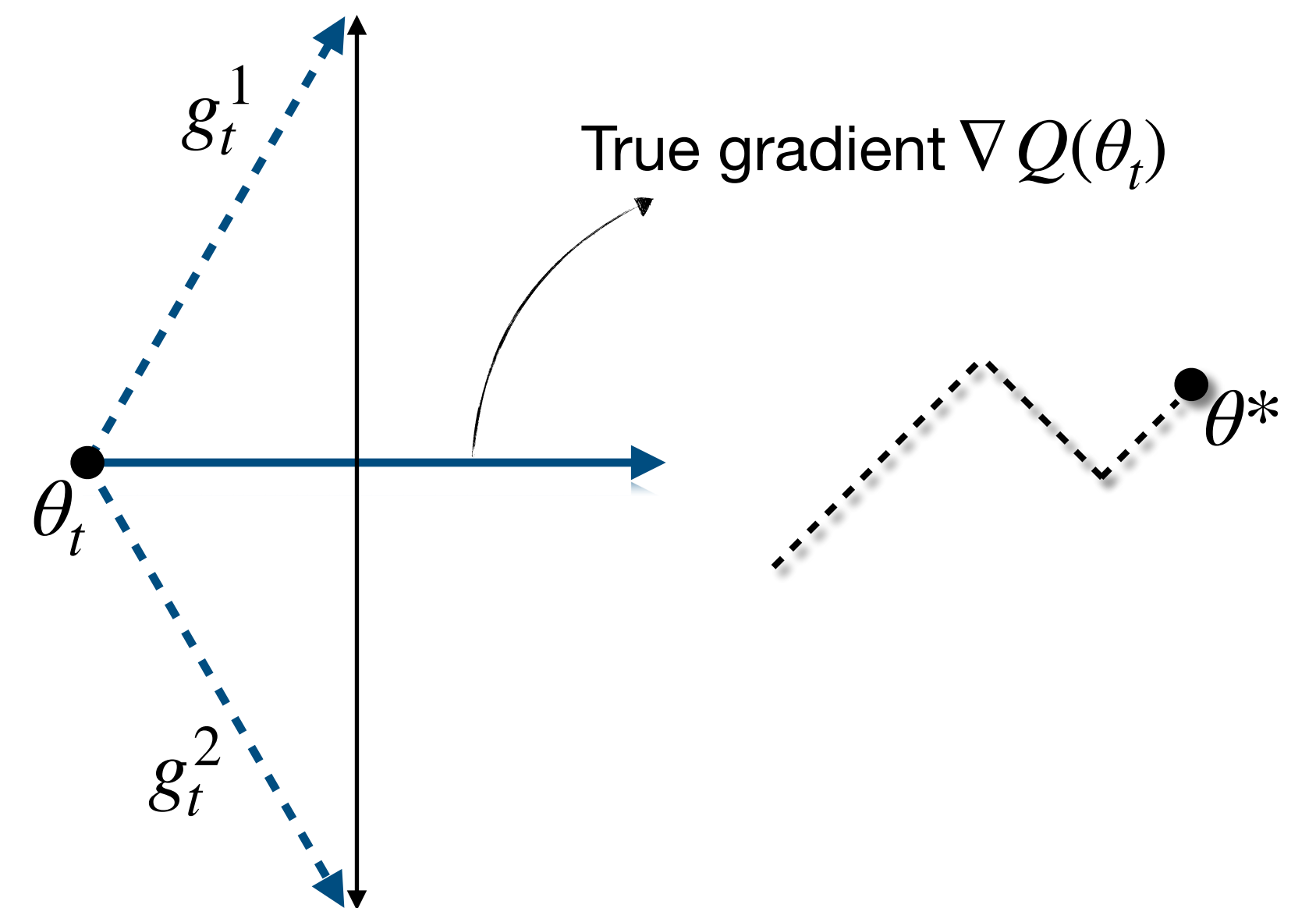
# What is the Problem?

Non-trivial variance of stochastic gradients

# What is the Problem?

Non-trivial variance of stochastic gradients

# What is the Problem?

Non-trivial variance of stochastic gradients

The worst-case guarantees
are quite bad



$g_t^1$

True gradient $\nabla Q(\theta_t)$

$\widetilde{g_t^3}$

$\theta_t$

$\theta*$

Erroneous but
viable outputs

$g_t^2$

variance

# What is the Problem?

Non-trivial variance of stochastic gradients

The worst-case guarantees
are quite bad

Variance reduction schemes -



$g_t^1$

True gradient $\nabla Q(\theta_t)$

$\widetilde{g_t^3}$

$\theta_t$

Erroneous but
viable outputs

$g_t^2$

variance

$\theta*$

# What is the Problem?

Non-trivial variance of stochastic gradients

The worst-case guarantees
are quite bad

Variance reduction schemes -

Computationally **expensive** (querying multiple gradients)

$g_t^1$

True gradient $\nabla Q(\theta_t)$

$\widetilde{g_t^3}$

$\theta_t$

Erroneous but
viable outputs

$g_t^2$

variance

$\theta*$

# What is the Problem?

## Non-trivial variance of stochastic gradients

The worst-case guarantees
are quite bad

Variance reduction schemes -

Computationally **expensive** (querying multiple gradients)

Leads to **bias in gradient estimate**

$g_t^1$

$\widetilde{g_t^3}$

True gradient $\nabla Q(\theta_t)$

$\theta^*$

$\theta_t$

Erroneous but
viable outputs

$g_t^2$

variance

# What is the Problem?

Non-trivial variance of stochastic gradients

The worst-case guarantees
are quite bad

Variance reduction schemes -

Computationally **expensive** (querying multiple gradients)

Leads to **bias in gradient estimate**



$g_t^1$

True gradient $\nabla Q(\theta_t)$

$\widetilde{g_t^3}$

$\theta_t$

Erroneous but
viable outputs

$g_t^2$

variance

$\theta*$

# What is the Problem?

Non-trivial variance of stochastic gradients

The worst-case guarantees
are quite bad

Variance reduction schemes -

Computationally **expensive** (querying multiple gradients)

Leads to **bias in gradient estimate**

Threatens the convergence



$g_t^1$

True gradient $\nabla Q(\theta_t)$

$\widetilde{g_t^3}$

$\theta_t$

Erroneous but
viable outputs

$g_t^2$

variance

$\theta*$

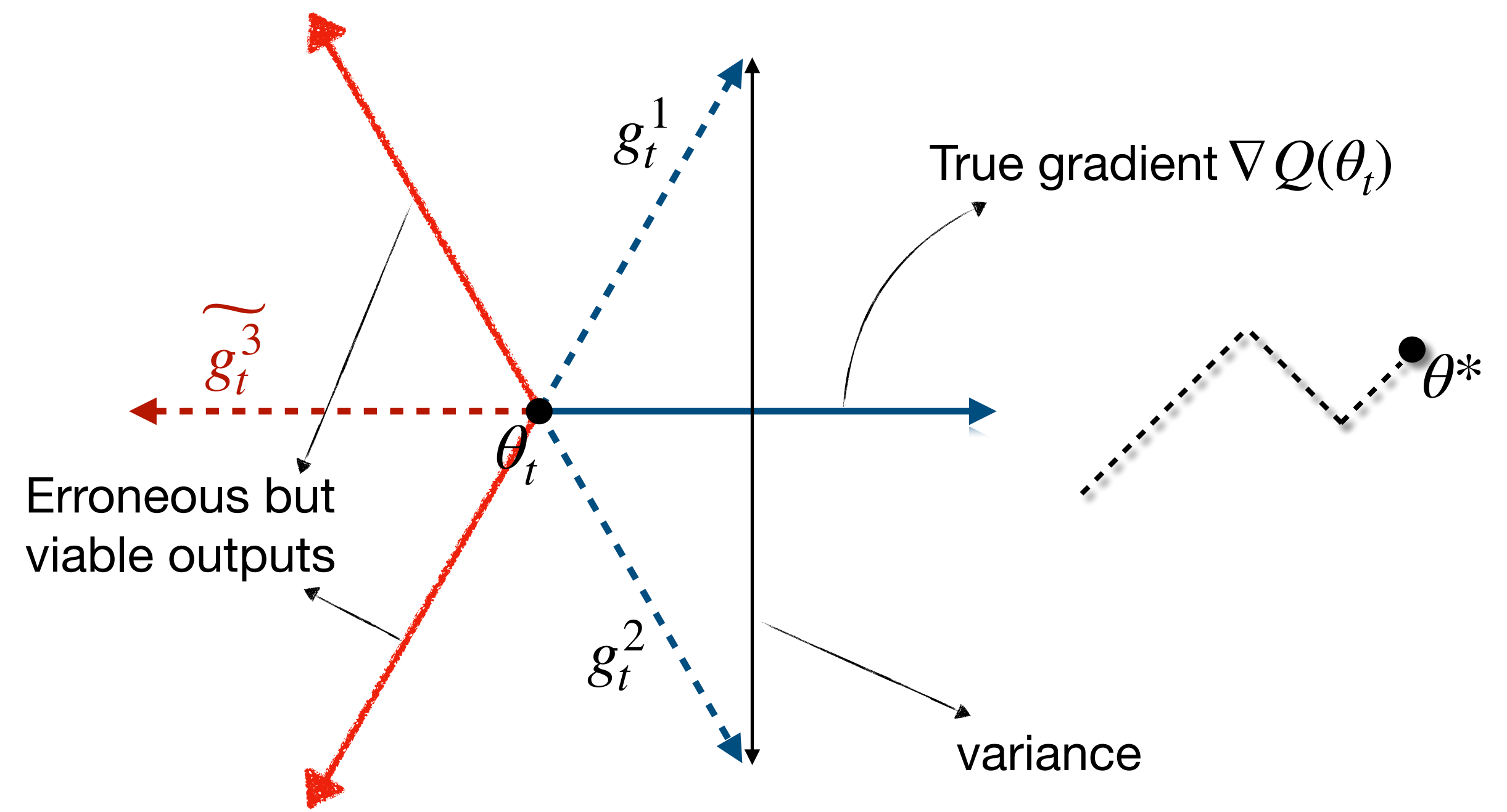# What is the Problem?

Non-trivial variance of stochastic gradients

The worst-case guarantees
are quite bad

Variance reduction schemes -

Computationally **expensive** (querying multiple gradients)

Leads to **bias in gradient estimate**

Threatens the convergence



$g_t^1$

True gradient $\nabla Q(\theta_t)$

$\widetilde{g_t^3}$

$\theta_t$

$\theta*$

Erroneous but
viable outputs

$g_t^2$

variance

# "With every mistake, me must surely be learning …

While my *GPU* gently weeps."

- George Harrison

# Local Gradient Momentum

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \beta \, m_{t-1}^i + (1 - \beta) \, g_t^i$$

$g_t^1$

$\nabla Q(\theta_t)$

$\theta_t$

$g_t^2$

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \beta\, m_{t-1}^i + (1 - \beta)\, g_t^i$$

$\beta \in [0, 1)$ is referred as the momentum coefficient

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \beta \, m_{t-1}^i + (1 - \beta) \, g_t^i$$

$\beta \in [0, 1)$ is referred as the momentum coefficient

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \beta\, m_{t-1}^i + (1 - \beta)\, g_t^i$$

$\beta \in [0, 1)$ is referred as the momentum coefficient

$g_t^1$

$\nabla Q(\theta_t)$

$\theta_t$

$g_t^2$

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \beta\, m_{t-1}^i + (1 - \beta)\, g_t^i$$

Past gradients

$\beta \in [0, 1)$ is referred as the momentum coefficient

$g_t^1$

$\nabla Q(\theta_t)$

$\theta_t$

$g_t^2$

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \beta\, m_{t-1}^i + (1 - \beta)\, g_t^i$$

Past gradients

$\beta \in [0, 1)$ is referred as the momentum coefficient

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \beta\, m_{t-1}^i + (1 - \beta)\, g_t^i$$

Past gradients

$\beta \in [0, 1)$ is referred as the momentum coefficient

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \boxed{\beta\, m_{t-1}^i} + (1 - \beta)\, g_t^i$$

Past gradients

$\beta \in [0, 1)$ is referred as the momentum coefficient

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \boxed{\beta\, m_{t-1}^i} + (1 - \beta)\, g_t^i$$

Past gradients

$\beta \in [0, 1)$ is referred as the momentum coefficient

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \boxed{\beta\, m_{t-1}^i} + (1 - \beta)\, g_t^i \longrightarrow \text{Past gradients}$$

$\beta \in [0, 1)$ is referred as the momentum coefficient

**Server aggregates the momentums** $R_t = F\left(m_t^1, \ldots, m_t^n\right)$

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \beta\, m_{t-1}^i + (1 - \beta)\, g_t^i$$

Past gradients

$\beta \in [0, 1)$ is referred as the momentum coefficient

**Server aggregates the momentums** $R_t = F\left(m_t^1, \ldots, m_t^n\right)$
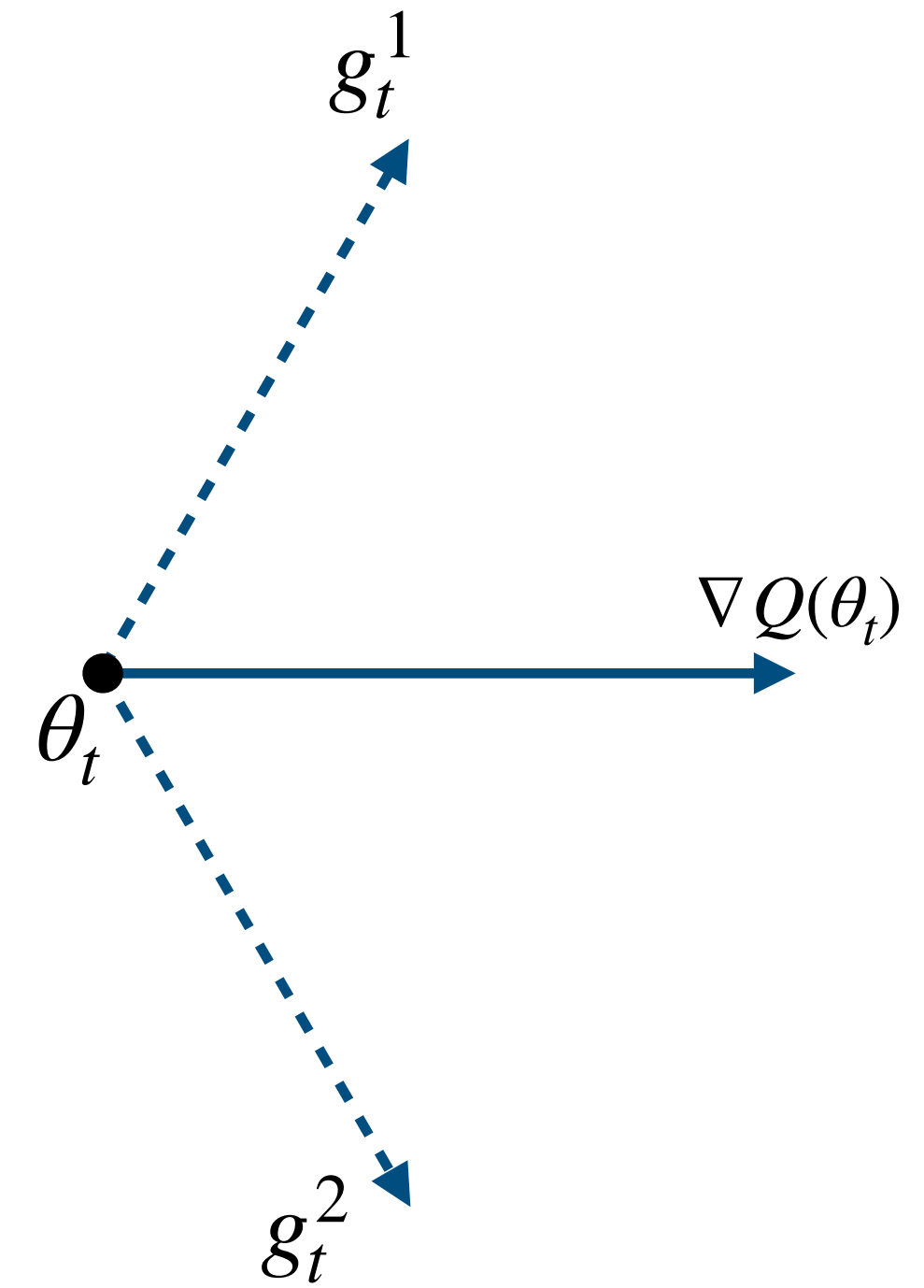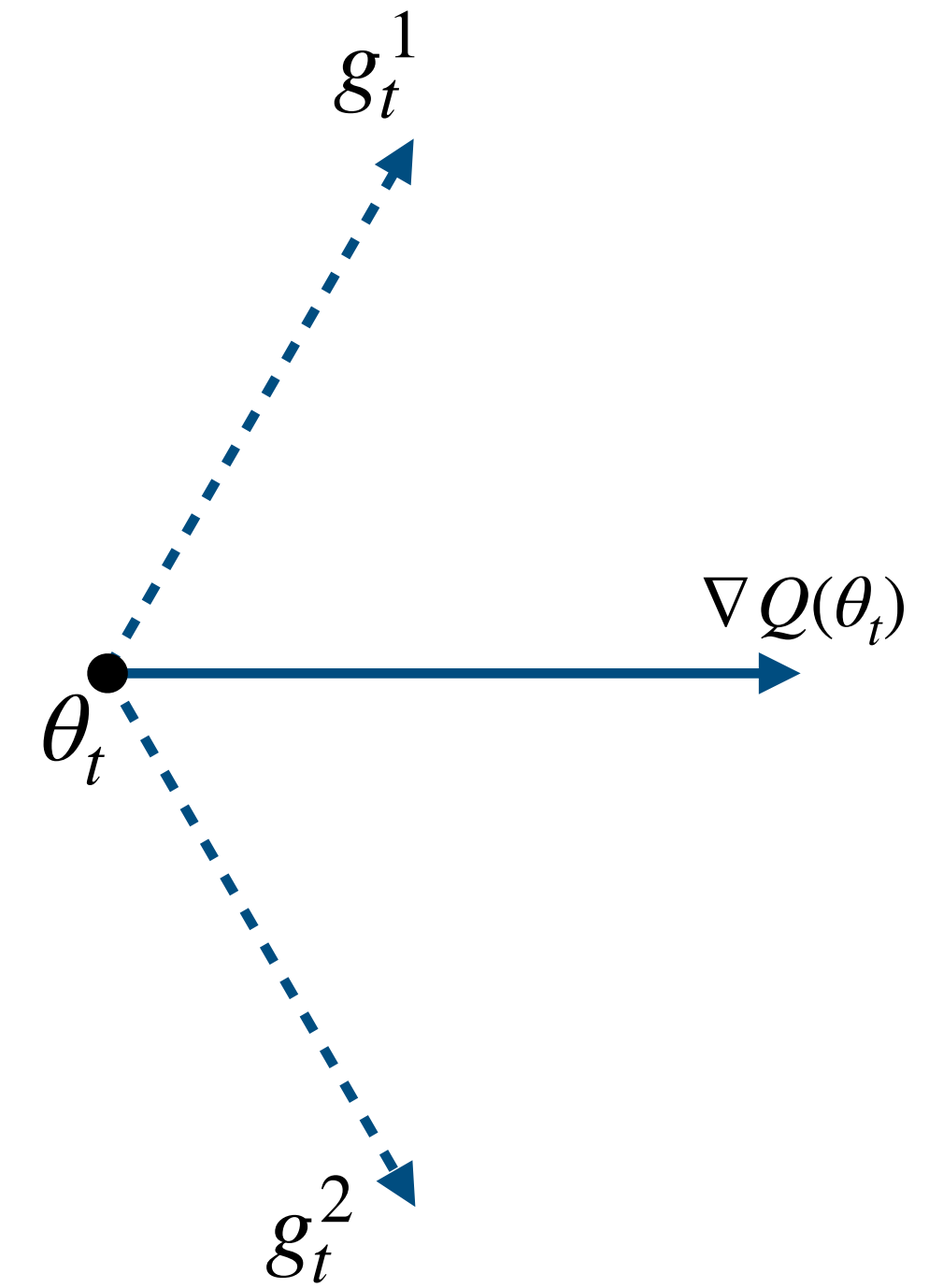
(Rest remains the same)

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \boxed{\beta\, m_{t-1}^i} + (1 - \beta)\, g_t^i$$

Past gradients

$\beta \in [0, 1)$ is referred as the momentum coefficient

**Server aggregates the momentums** $R_t = F\left(m_t^1, \ldots, m_t^n\right)$
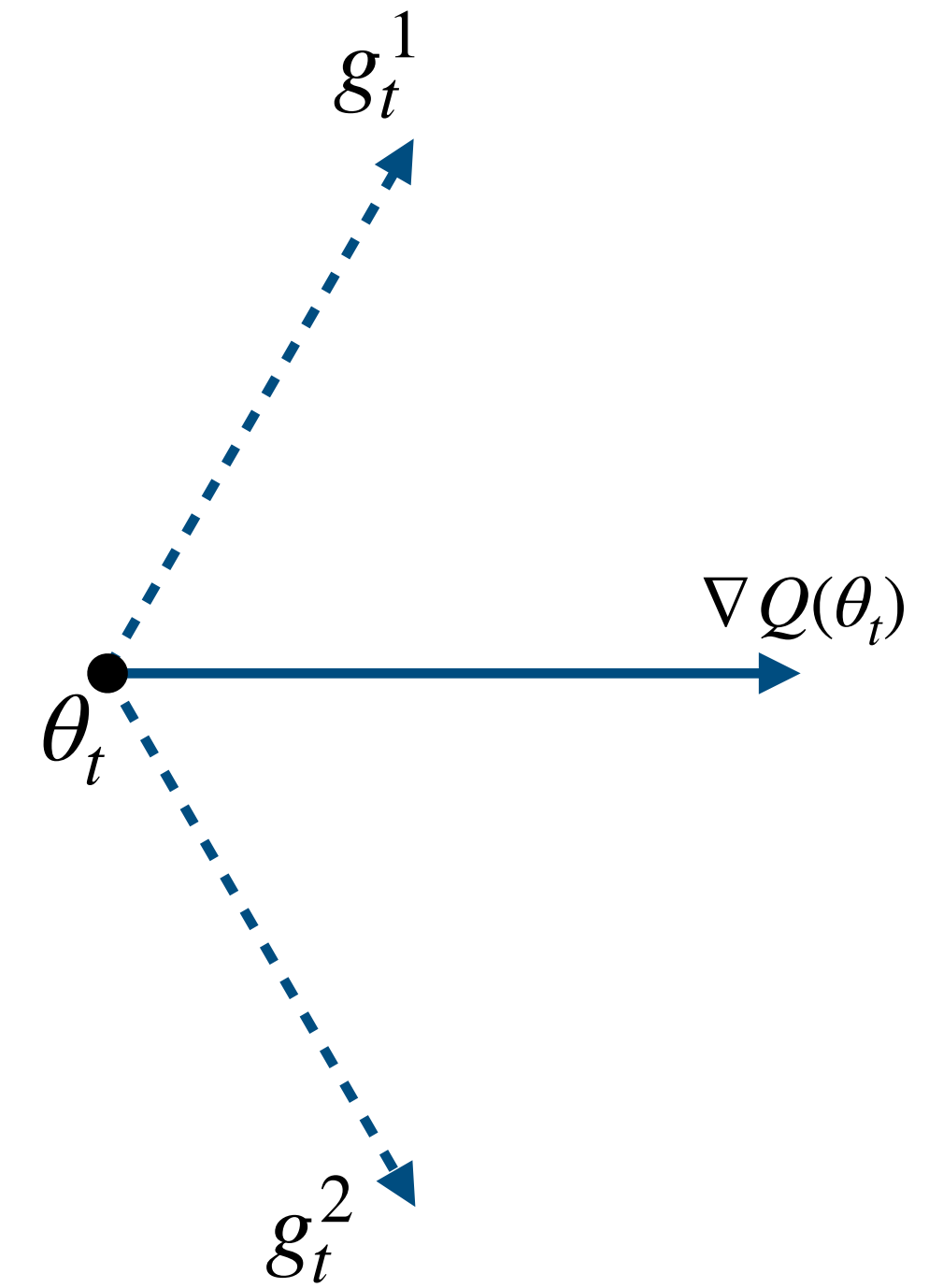
(Rest remains the same)

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \boxed{\beta\, m_{t-1}^i} + (1 - \beta)\, g_t^i$$

Past gradients

$\beta \in [0, 1)$ is referred as the momentum coefficient

**Server aggregates the momentums** $R_t = F\left(m_t^1, \ldots, m_t^n\right)$

(Rest remains the same)

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \boxed{\beta\, m_{t-1}^i} + (1-\beta)\, g_t^i$$  ← Past gradients

$\beta \in [0,1)$ is referred as the momentum coefficient

**Server aggregates the momentums** $R_t = F\left(m_t^1, \ldots, m_t^n\right)$

(Rest remains the same)

$g_t^1$

$m_t^1$

$\nabla Q(\theta_t)$

$\theta_t$

$m_t^2$

$g_t^2$

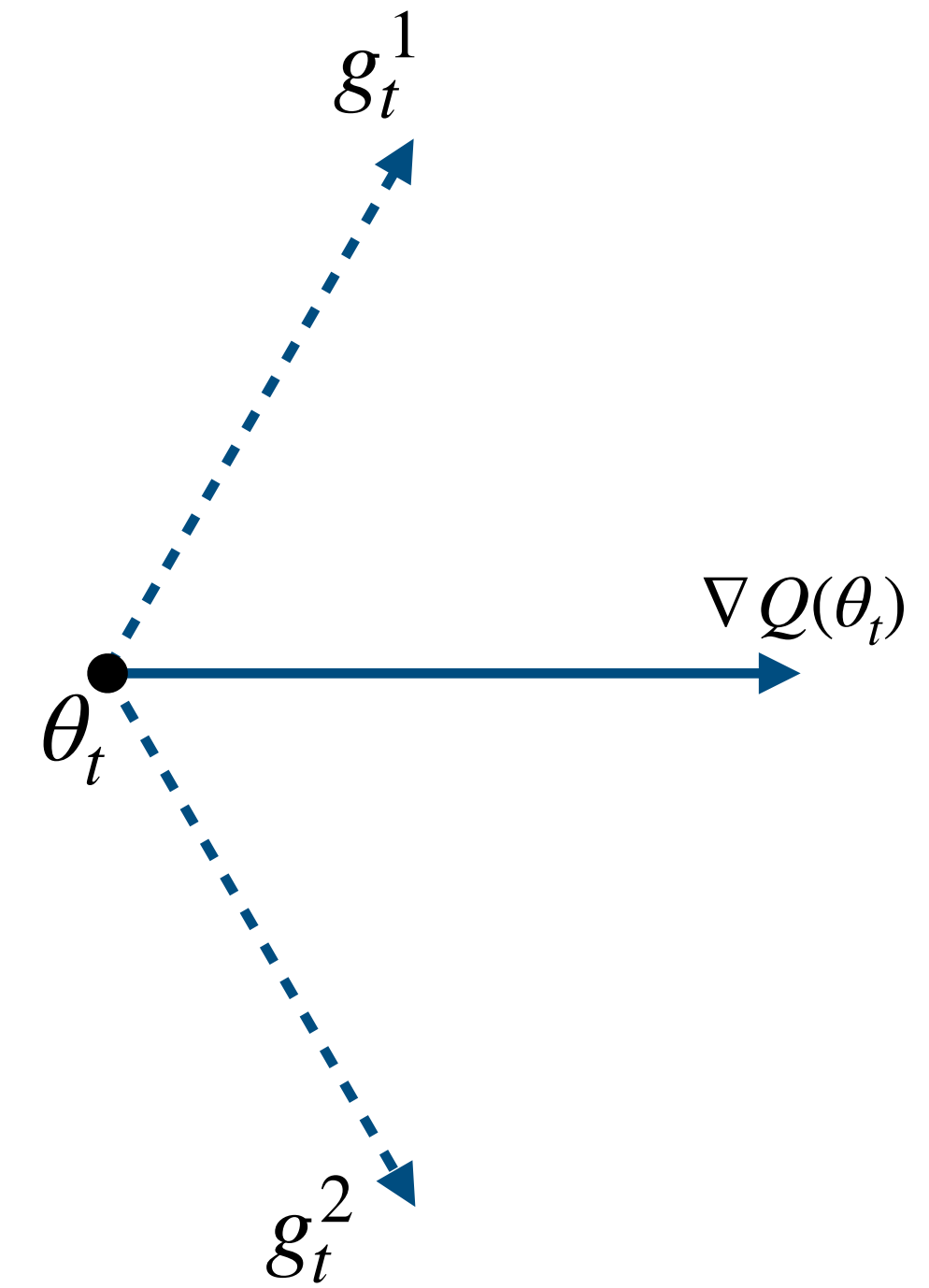The variance of the momentums can be reduced -

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \beta\, m_{t-1}^i + (1-\beta)\, g_t^i$$

Past gradients

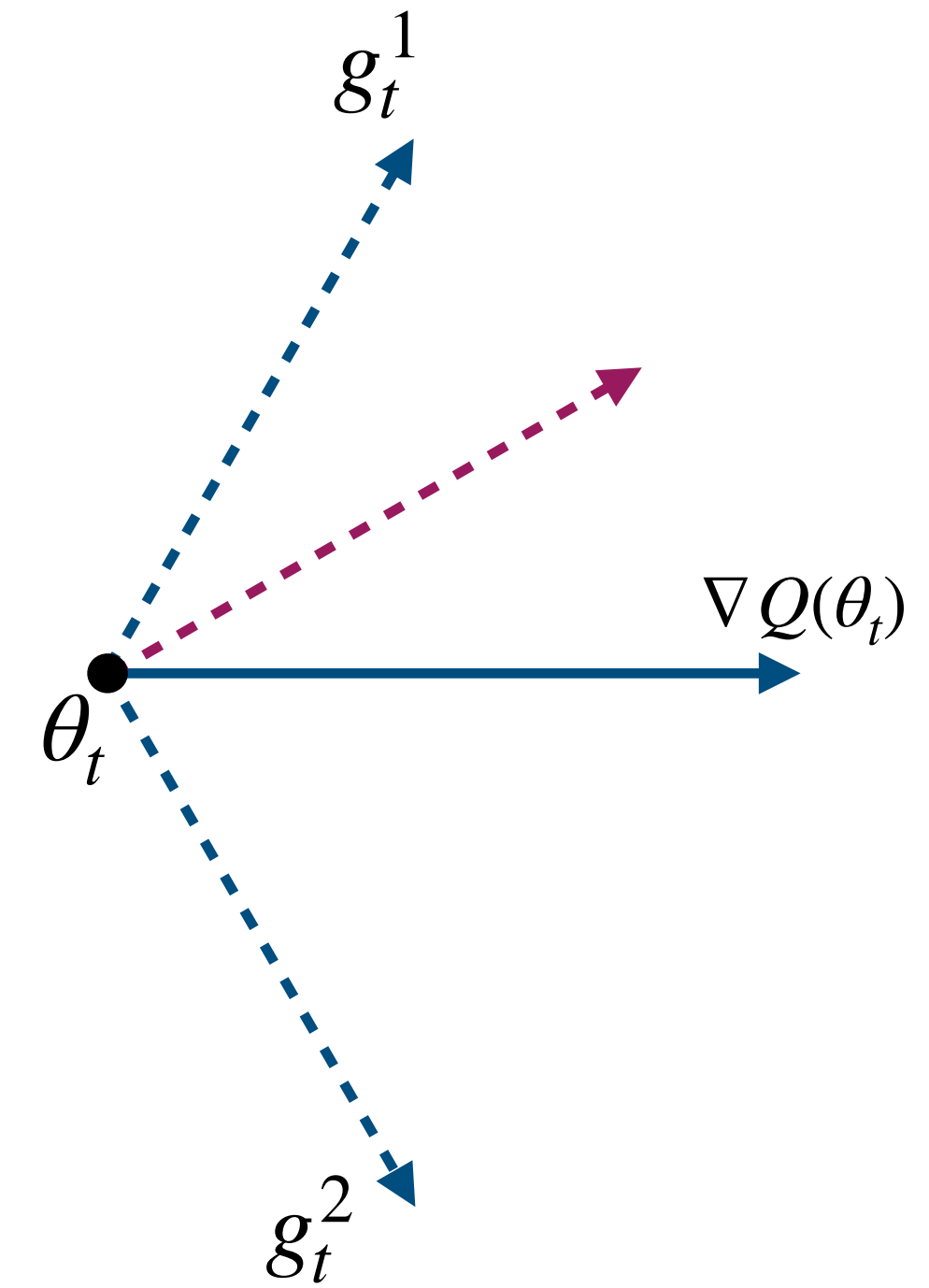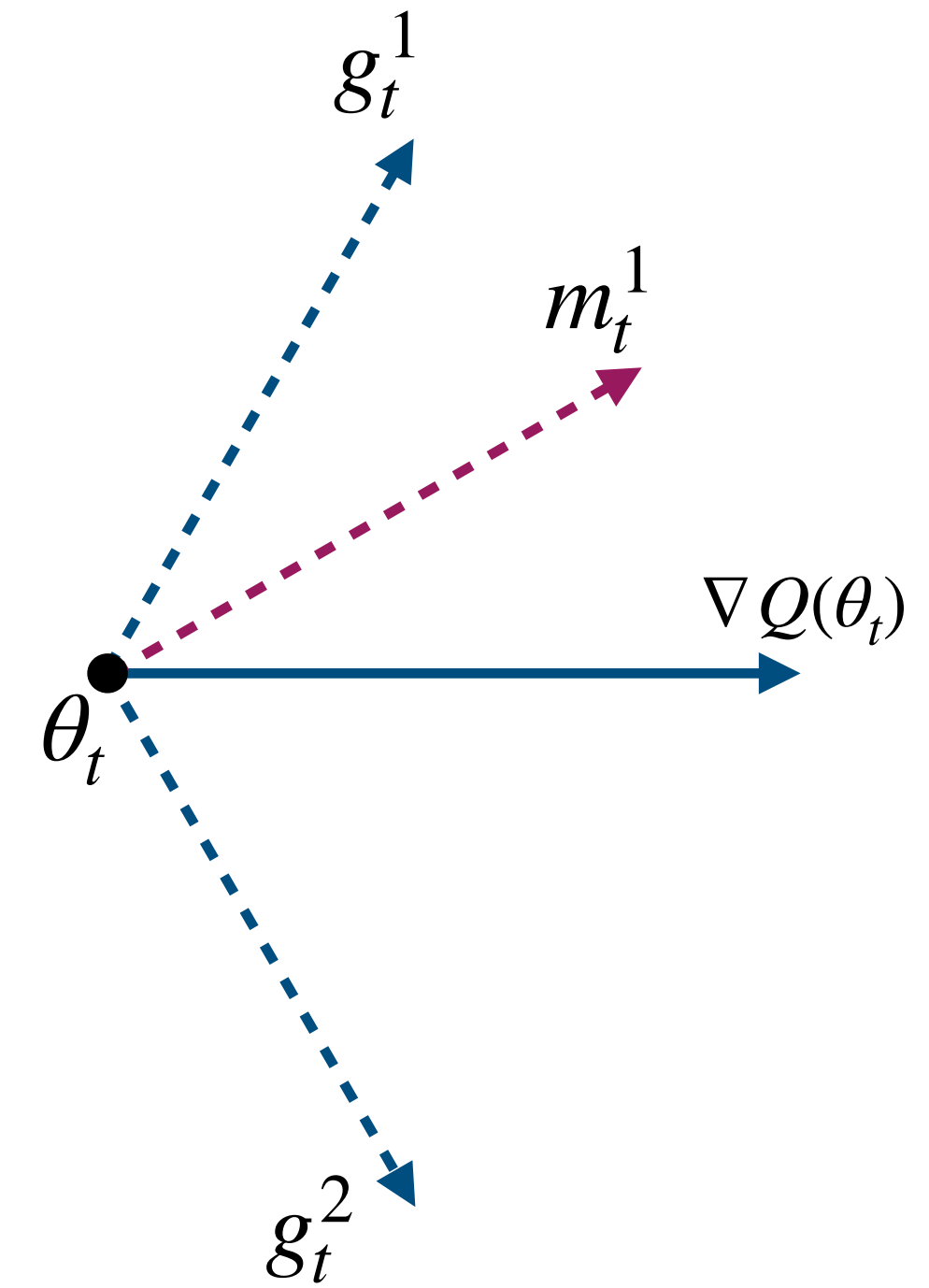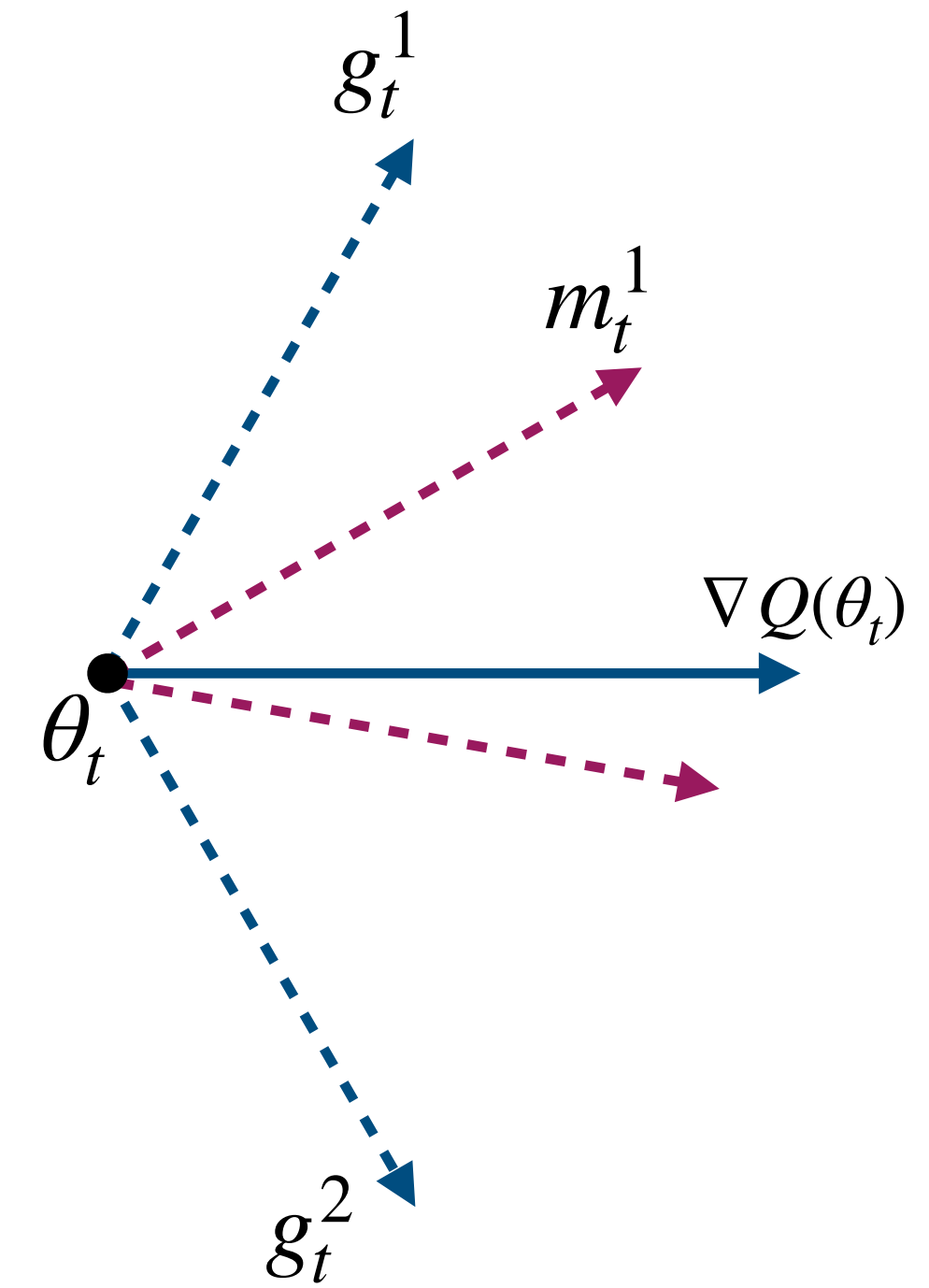$\beta \in [0, 1)$ is referred as the momentum coefficient

**Server aggregates the momentums** $R_t = F\left(m_t^1, \ldots, m_t^n\right)$

(Rest remains the same)

The variance of the momentums can be reduced -

$$\mathbb{E}\left[\|m_t^i - m_t^j\|^2\right] \leq (1-\beta)\, c\sigma^2 \,, \quad \forall \ \textbf{honest} \ \ i, j$$

$g_t^1$

$m_t^1$

$\nabla Q(\theta_t)$

$\theta_t$

$m_t^2$

$g_t^2$

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \boxed{\beta\, m_{t-1}^i} + (1 - \beta)\, g_t^i$$

Past gradients

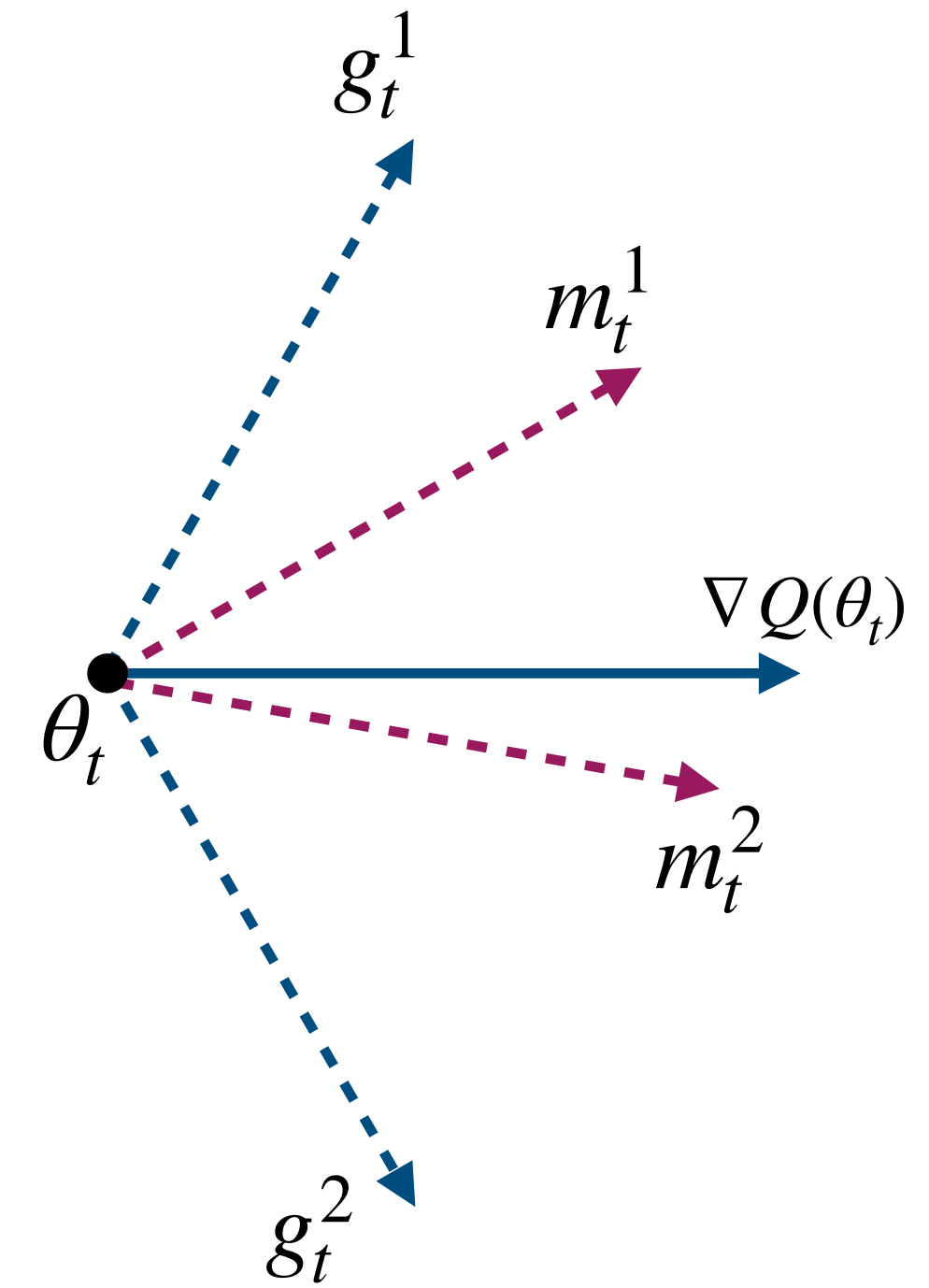$\beta \in [0, 1)$ is referred as the momentum coefficient

**Server aggregates the momentums** $R_t = F\left(m_t^1, \ldots, m_t^n\right)$

(Rest remains the same)

$g_t^1$

$m_t^1$

$\nabla Q(\theta_t)$

$\theta_t$

$m_t^2$

$g_t^2$

The variance of the momentums can be reduced -

$$\mathbb{E}\left[\|m_t^i - m_t^j\|^2\right] \leq (1 - \beta)\, c\sigma^2\,, \quad \forall \ \textbf{honest} \ i, j$$
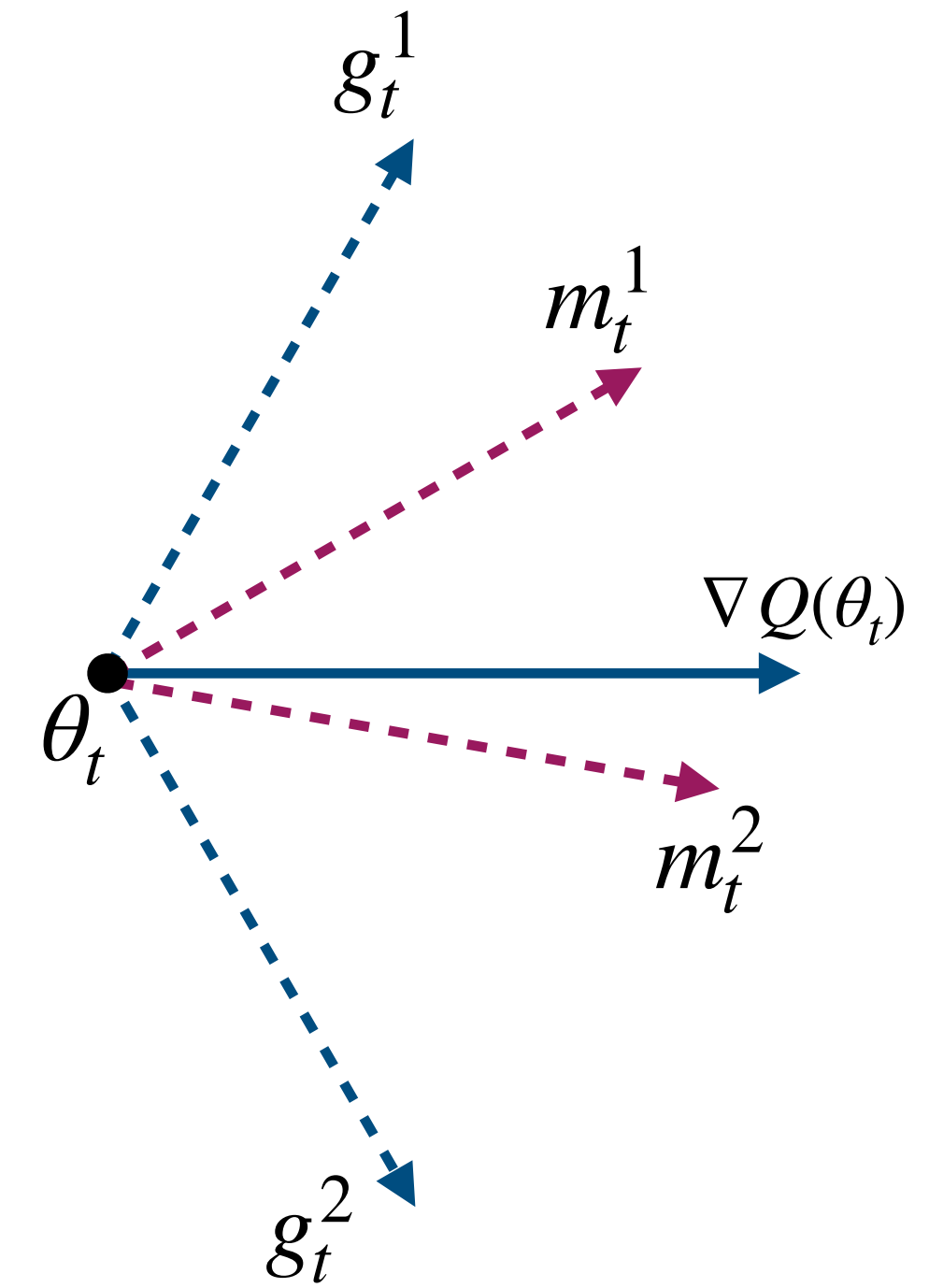
# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,

$$m_t^i = \boxed{\beta\, m_{t-1}^i} + (1 - \beta)\, g_t^i$$

Past gradients

$\beta \in [0, 1)$ is referred as the momentum coefficient

**Server aggregates the momentums** $R_t = F\left(m_t^1, \ldots, m_t^n\right)$

(Rest remains the same)

The variance of the momentums can be reduced -

$$\mathbb{E}\left[\|m_t^i - m_t^j\|^2\right] \leq (1 - \beta)\, c\sigma^2\,, \quad \forall\ \text{honest}\ i, j$$

# Local Gradient Momentum

**Nodes** compute and **send momentums** of their gradients,
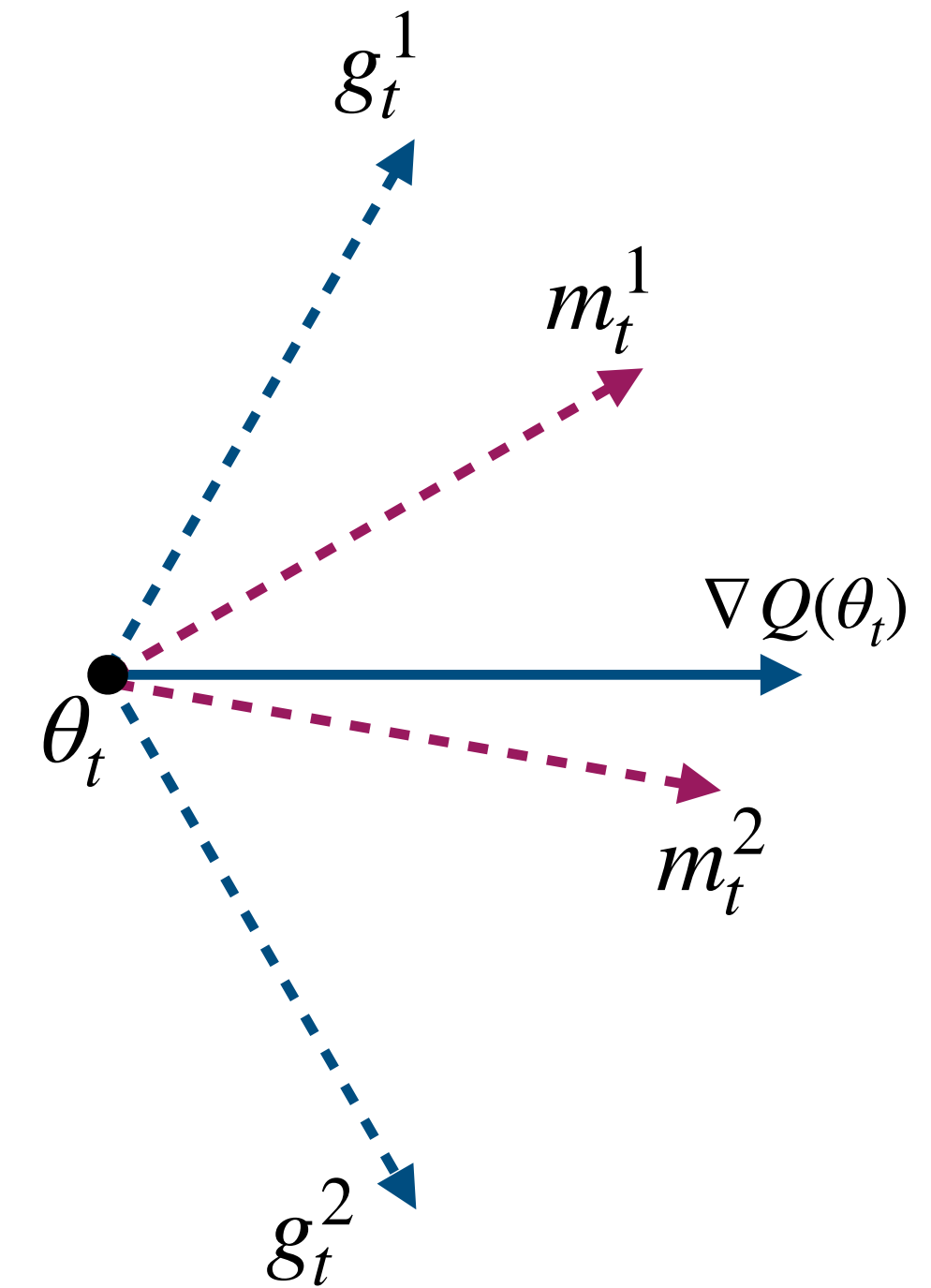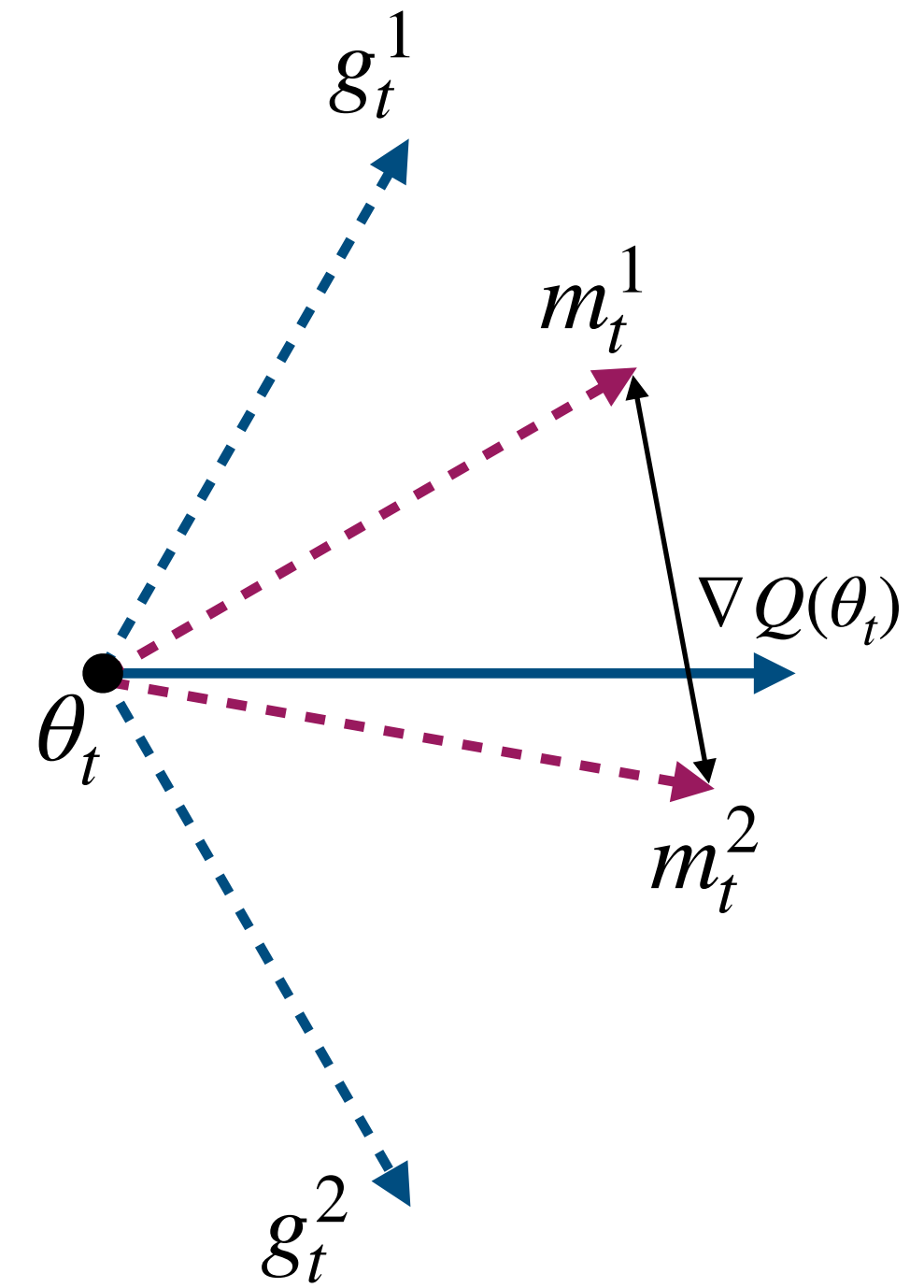
$$m_t^i = \boxed{\beta \, m_{t-1}^i} + (1 - \beta) \, g_t^i$$

→ Past gradients

$\beta \in [0, 1)$ is referred as the momentum coefficient

**Server aggregates the momentums** $R_t = F\left(m_t^1, \ldots, m_t^n\right)$

(Rest remains the same)

The variance of the momentums can be reduced -

$$\mathbb{E}\left[\|m_t^i - m_t^j\|^2\right] \leq (1 - \beta) c \sigma^2 \, , \quad \forall \; \textbf{honest} \;\; i, j$$

→ Ideally large

$g_t^1$

$m_t^1$

$\nabla Q(\theta_t)$

$\theta_t$

$m_t^2$

$g_t^2$

# Caveat of Momentum

# Caveat of Momentum

# Caveat of Momentum

# Caveat of Momentum

# Caveat of Momentum

# Caveat of Momentum

Local momentums are no longer
unbiased gradient estimators

$g_t^1$

$m_t^1$

$\overline{m}_t$

$\theta_t$

$\nabla Q(\theta_t)$

$m_t^2$

$g_t^2$

# Caveat of Momentum

Local momentums are no longer unbiased gradient estimators

$$\mathbb{E}\left[\|\overline{m}_t - \nabla Q\left(\theta_t\right)\|^2\right] \leq \beta^2 \, \mathbb{E}\left[\|\overline{m}_{t-1} - \nabla Q\left(\theta_{t-1}\right)\|^2\right] + (1-\beta)^2 \frac{\sigma^2}{n-f} + \epsilon(t)$$

# Caveat of Momentum

Local momentums are no longer
unbiased gradient estimators

$$\mathbb{E}\left[\|\overline{m}_t - \nabla Q\left(\theta_t\right)\|^2\right] \leq \beta^2 \,\mathbb{E}\left[\|\overline{m}_{t-1} - \nabla Q\left(\theta_{t-1}\right)\|^2\right] + (1-\beta)^2 \frac{\sigma^2}{n-f} + \epsilon(t)$$

$g_t^1$

$m_t^1$

$\overline{m}_t$

$\theta_t$

$m_t^2$

$\nabla Q(\theta_t)$

$g_t^2$

# Caveat of Momentum

Local momentums are no longer
unbiased gradient estimators

$$\mathbb{E}\left[\|\overline{m}_t - \nabla Q\left(\theta_t\right)\|^2\right] \leq \beta^2\,\mathbb{E}\left[\|\overline{m}_{t-1} - \nabla Q\left(\theta_{t-1}\right)\|^2\right] + (1-\beta)^2\frac{\sigma^2}{n-f} + \epsilon(t)$$

$g_t^1$

$m_t^1$

$\overline{m}_t$

$\theta_t$

$m_t^2$

$\nabla Q(\theta_t)$

$g_t^2$

# Caveat of Momentum

Local momentums are no longer unbiased gradient estimators

$$\mathbb{E}\left[\|\overline{m}_t - \nabla Q\left(\theta_t\right)\|^2\right] \le \boxed{\beta^2} \mathbb{E}\left[\|\overline{m}_{t-1} - \nabla Q\left(\theta_{t-1}\right)\|^2\right] + (1-\beta)^2 \frac{\sigma^2}{n-f} + \epsilon(t)$$

$g_t^1$

$m_t^1$

$\overline{m}_t$

$\theta_t$

$m_t^2$

$\nabla Q(\theta_t)$

$g_t^2$

# Caveat of Momentum

Local momentums are no longer
unbiased gradient estimators

$$\mathbb{E}\left[\|\overline{m}_t - \nabla Q\left(\theta_t\right)\|^2\right] \leq \beta^2 \mathbb{E}\left[\|\overline{m}_{t-1} - \nabla Q(\theta_{t-1})\|^2\right] + (1-\beta)^2\frac{\sigma^2}{n-f} + \epsilon(t)$$

$g_t^1$

$m_t^1$

$\overline{m}_t$

$\theta_t$

$\nabla Q(\theta_t)$

$m_t^2$

$g_t^2$

# Caveat of Momentum

Local momentums are no longer
unbiased gradient estimators

$$\mathbb{E}\left[\|\overline{m}_t - \nabla Q\left(\theta_t\right)\|^2\right] \leq \beta^2 \mathbb{E}\left[\|\overline{m}_{t-1} - \nabla Q\left(\theta_{t-1}\right)\|^2\right] + (1-\beta)^2\frac{\sigma^2}{n-f} + \epsilon(t)$$

Ideally small

$g_t^1$

$m_t^1$

$\overline{m}_t$

$\theta_t$

$m_t^2$

$\nabla Q(\theta_t)$

$g_t^2$

# Caveat of Momentum

Local momentums are no longer unbiased gradient estimators

$$\mathbb{E}\left[\|\overline{m}_t - \nabla Q\left(\theta_t\right)\|^2\right] \leq \beta^2 \mathbb{E}\left[\|\overline{m}_{t-1} - \nabla Q\left(\theta_{t-1}\right)\|^2\right] + (1-\beta)^2 \frac{\sigma^2}{n-f} + \epsilon(t)$$

Ideally small

!! Good News !!

$g_t^1$

$m_t^1$

$\overline{m}_t$

$\theta_t$

$m_t^2$

$\nabla Q(\theta_t)$

$g_t^2$

# Caveat of Momentum

Local momentums are no longer unbiased gradient estimators

$$\mathbb{E}\left[\|\overline{m}_t - \nabla Q\left(\theta_t\right)\|^2\right] \leq \beta^2 \mathbb{E}\left[\|\overline{m}_{t-1} - \nabla Q\left(\theta_{t-1}\right)\|^2\right] + (1-\beta)^2 \frac{\sigma^2}{n-f} + \epsilon(t)$$

Ideally small

!! Good News !!

**IF** the **momentum** coefficient is **chosen wisely**

$g_t^1$

$m_t^1$

$\overline{m}_t$

$\theta_t$

$m_t^2$ $\nabla Q(\theta_t)$

$g_t^2$

# Caveat of Momentum

Local momentums are no longer
unbiased gradient estimators

$$\mathbb{E}\left[\|\overline{m}_t - \nabla Q\left(\theta_t\right)\|^2\right] \leq \beta^2 \mathbb{E}\left[\|\overline{m}_{t-1} - \nabla Q\left(\theta_{t-1}\right)\|^2\right] + (1-\beta)^2 \frac{\sigma^2}{n-f} + \epsilon(t)$$

Ideally small

!! Good News !!

**IF** the **momentum** coefficient is **chosen wisely**

**THEN** the **good** (reduced variance)
**Outweighs the bad** (biased gradient estimation)



$g_t^1$

$m_t^1$

$\overline{m}_t$

$\theta_t$

$m_t^2$ $\nabla Q(\theta_t)$

$g_t^2$

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

$(f, \lambda)$-*resilient averaging (RESA)* -

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

$(f, \lambda)$-***resilient averaging (RESA)*** -

$\exists\ \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall\ S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i, j \in S} \|x_i - x_j\|$$

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

$(f, \lambda)$-**resilient averaging (RESA)** -

$\exists\ \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall\ S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\| F(x_1, \ldots, x_n) - \bar{x}_S \| \leq \lambda \max_{i,j \in S} \| x_i - x_j \|$$

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

$(f, \lambda)$-**resilient averaging (RESA)** -

$\exists \; \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall \; S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i,j \in S} \|x_i - x_j\|$$

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

$(f, \lambda)$-*resilient averaging (RESA)* -

$\exists \, \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall \, S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i,j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

$(f, \lambda)$-**resilient averaging (RESA)** -

$\exists \, \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall \, S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i,j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

$(f, \lambda)$-***resilient averaging (RESA)*** -

$\exists \, \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall \, S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i, j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

$(f, \lambda)$-*resilient averaging (RESA)* -

$\exists \, \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall \, S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i,j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

$(f, \lambda)$-***resilient averaging (RESA)*** -

$\exists\ \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall\ S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i, j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

$(f, \lambda)$-***resilient averaging (RESA)*** -

$\exists \; \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall \; S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i,j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

$(f, \lambda)$-*resilient averaging (RESA)* -

$\exists \ \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall \ S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i,j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient

Byzantine

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

$(f, \lambda)$-*resilient averaging (RESA)* -

$\exists \; \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall \; S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\| F(x_1, \ldots, x_n) - \bar{x}_S \| \leq \lambda \max_{i, j \in S} \| x_i - x_j \|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient

Byzantine

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

$(f, \lambda)$-**resilient averaging (RESA)** -

$\exists\ \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall\ S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i,j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient

Honest

Byzantine

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

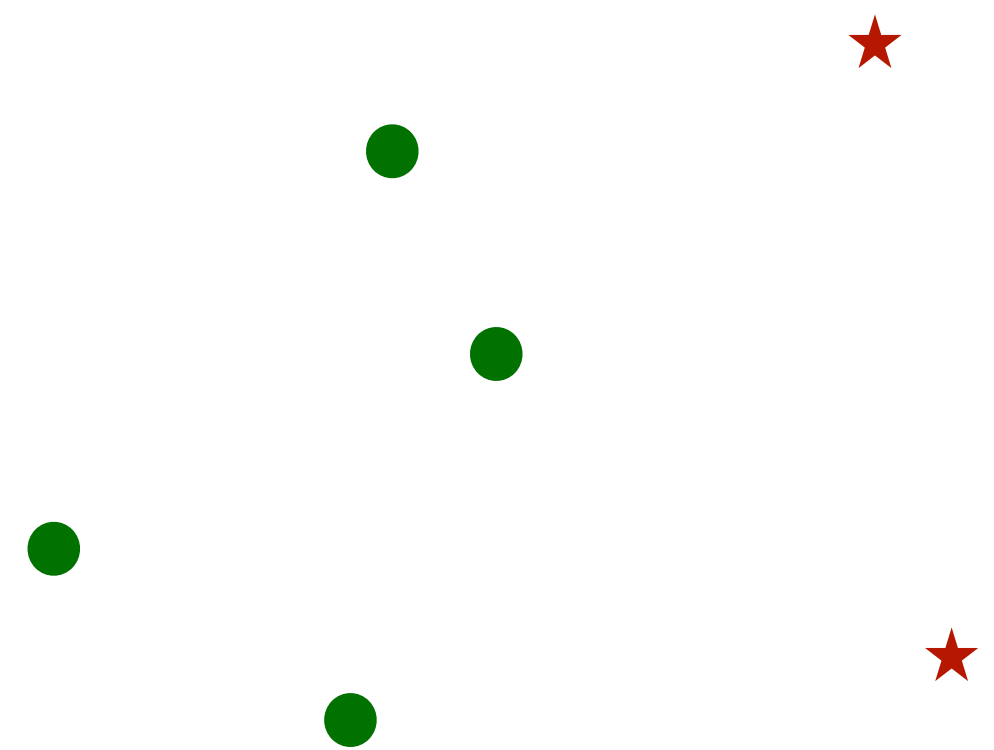$(f, \lambda)$-*resilient averaging (RESA)* -

$\exists \; \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall \; S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i, j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient

Honest

Byzantine

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

$(f, \lambda)$-**resilient averaging (RESA)** -
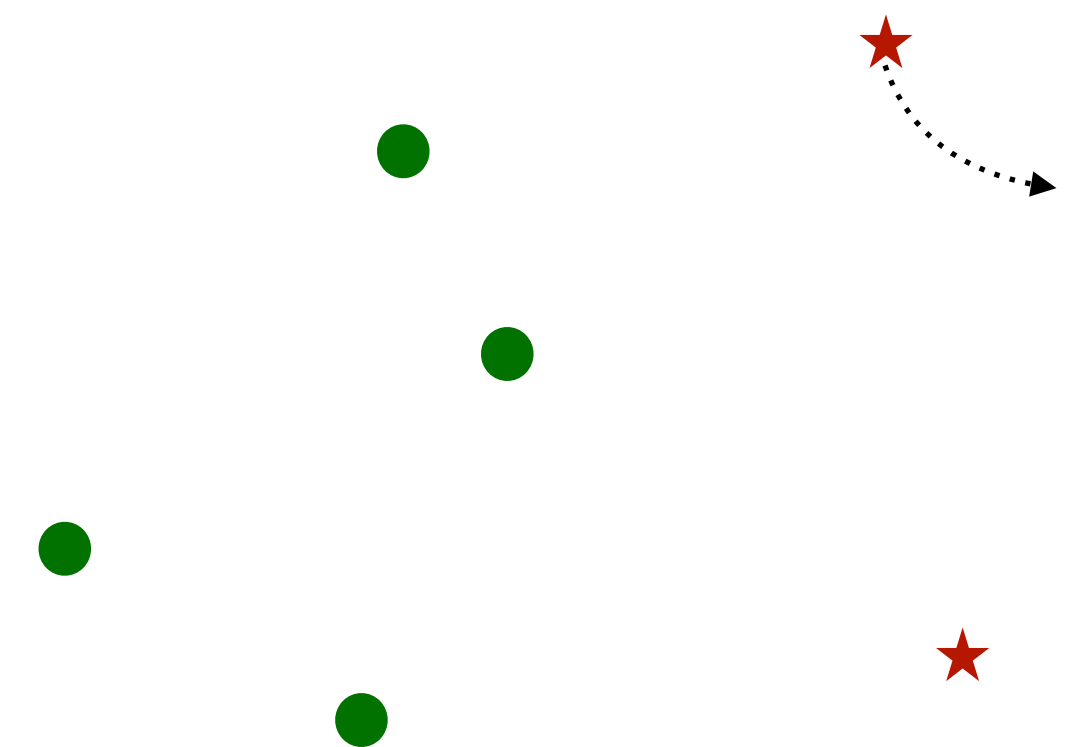
$\exists\ \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall\ S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i,j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient

Honest

Byzantine

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

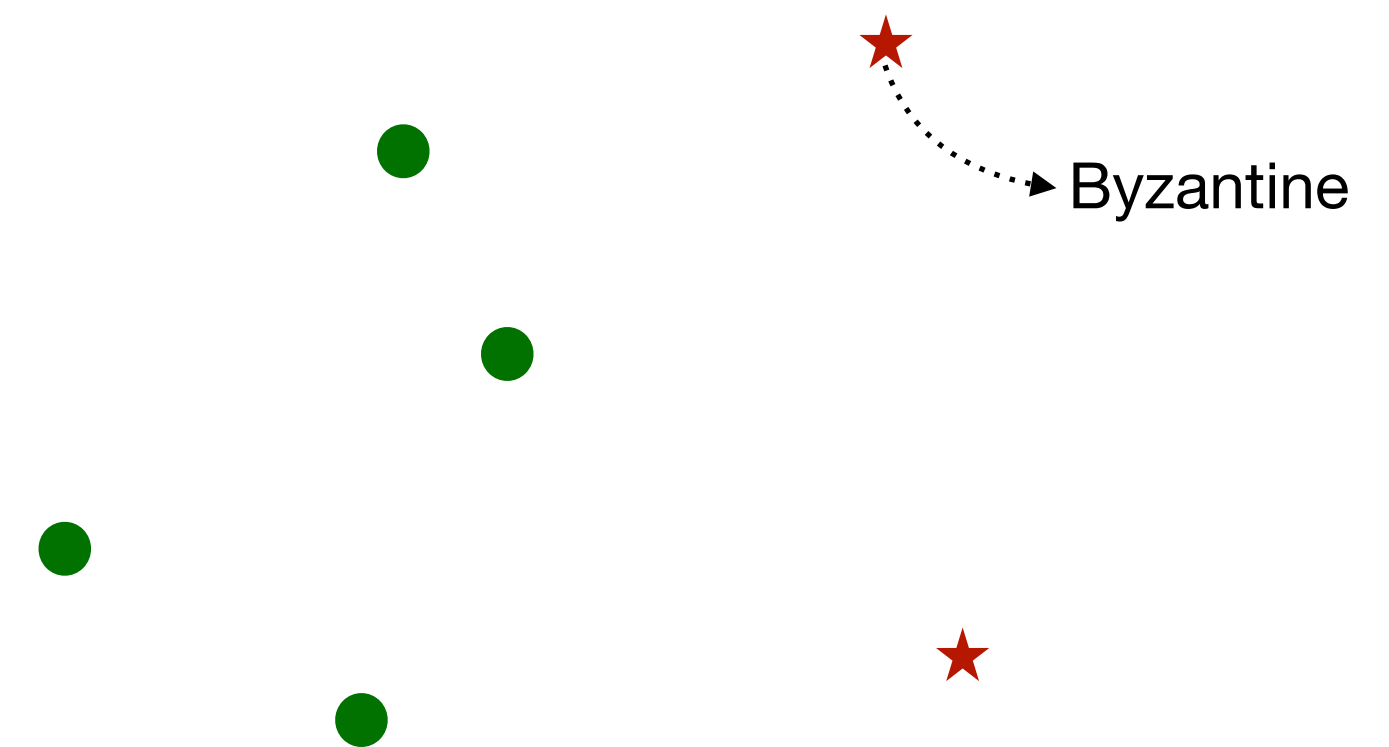$(f, \lambda)$-***resilient averaging (RESA)*** -

$\exists\ \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall\ S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i,j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient

Honest

Byzantine

Honest diameter

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

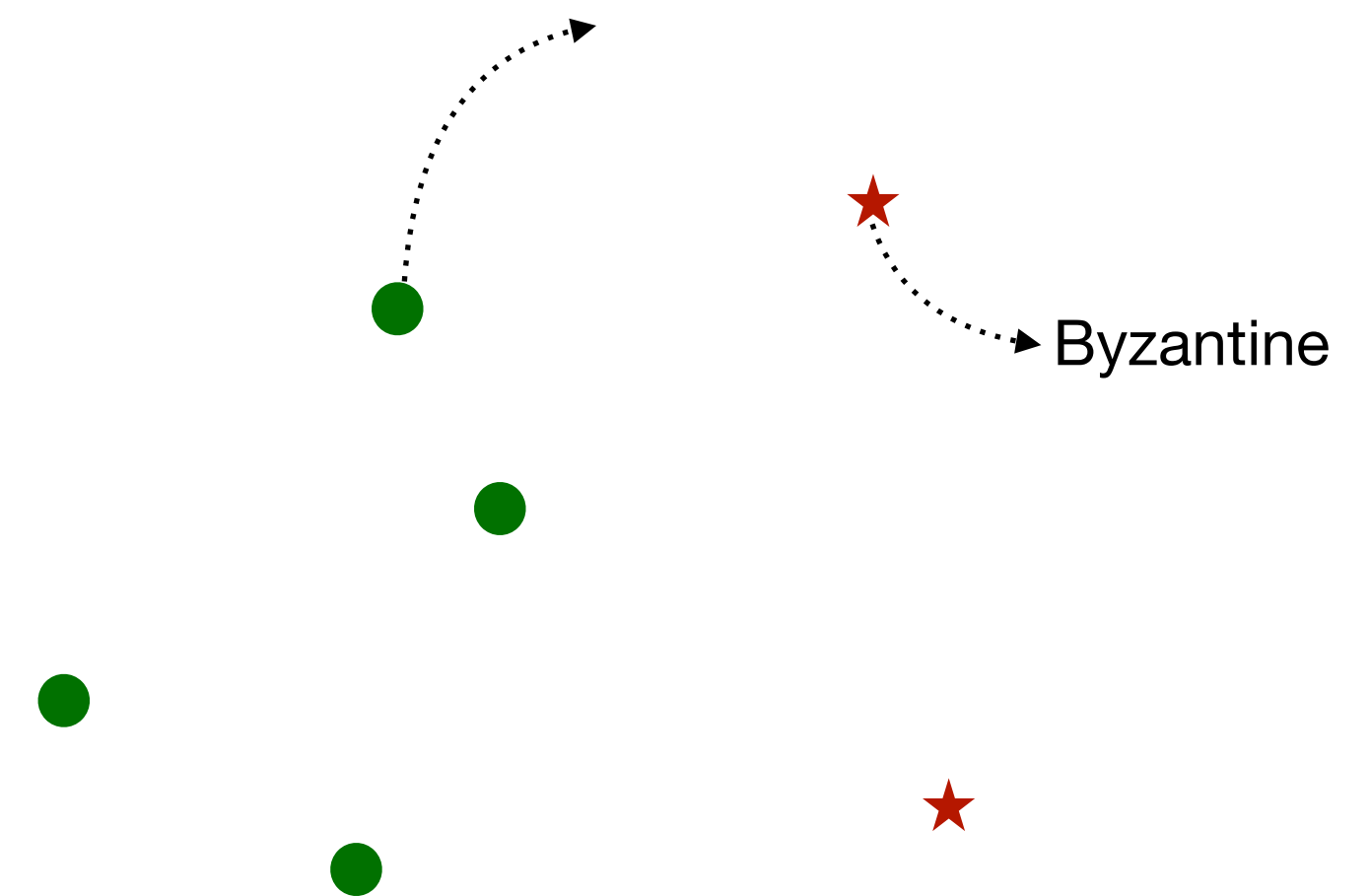$(f, \lambda)$-***resilient averaging (RESA)*** -

$\exists\ \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall\ S \subset \{1,\ldots,n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i,j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient



Honest

Byzantine

Honest diameter

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

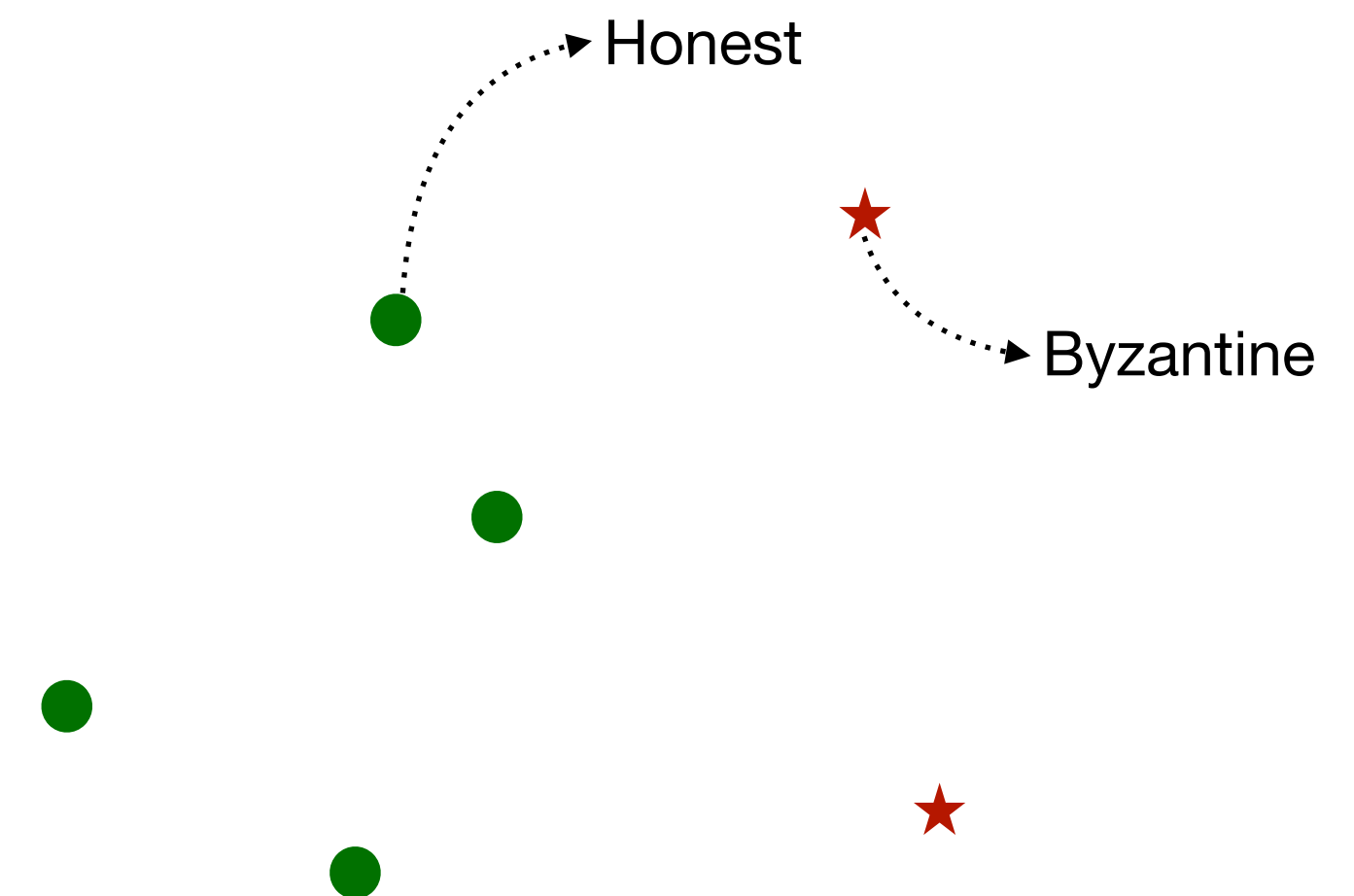$(f, \lambda)$-*resilient averaging (RESA)* -

$\exists \, \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall \, S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i, j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient



Honest

Byzantine

Honest diameter

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

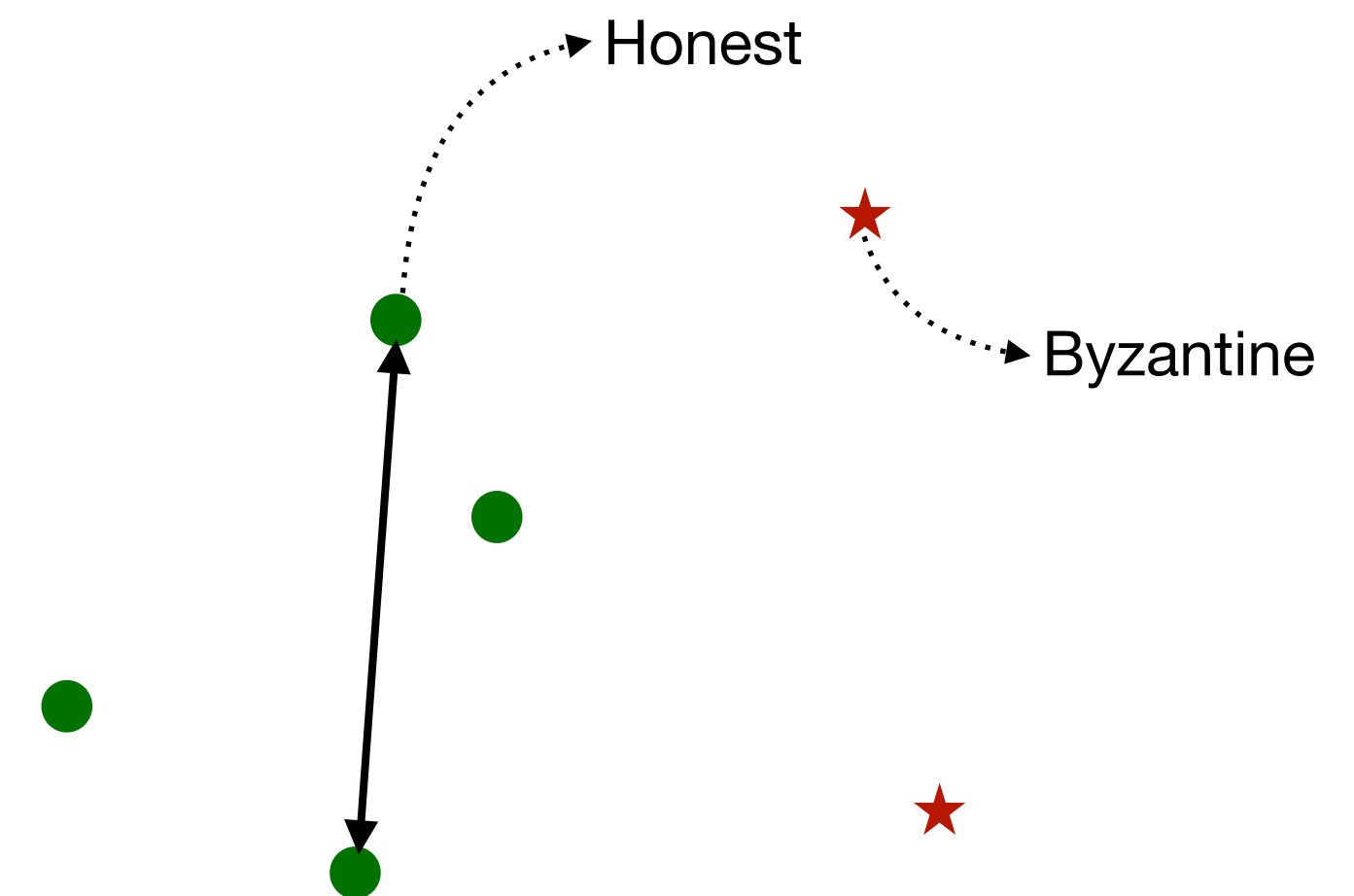**$(f, \lambda)$-*resilient averaging (RESA)* -**

$\exists\ \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall\ S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i,j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient

Honest

Byzantine

Honest average

Honest diameter

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

$(f, \lambda)$-*resilient averaging (RESA)* -
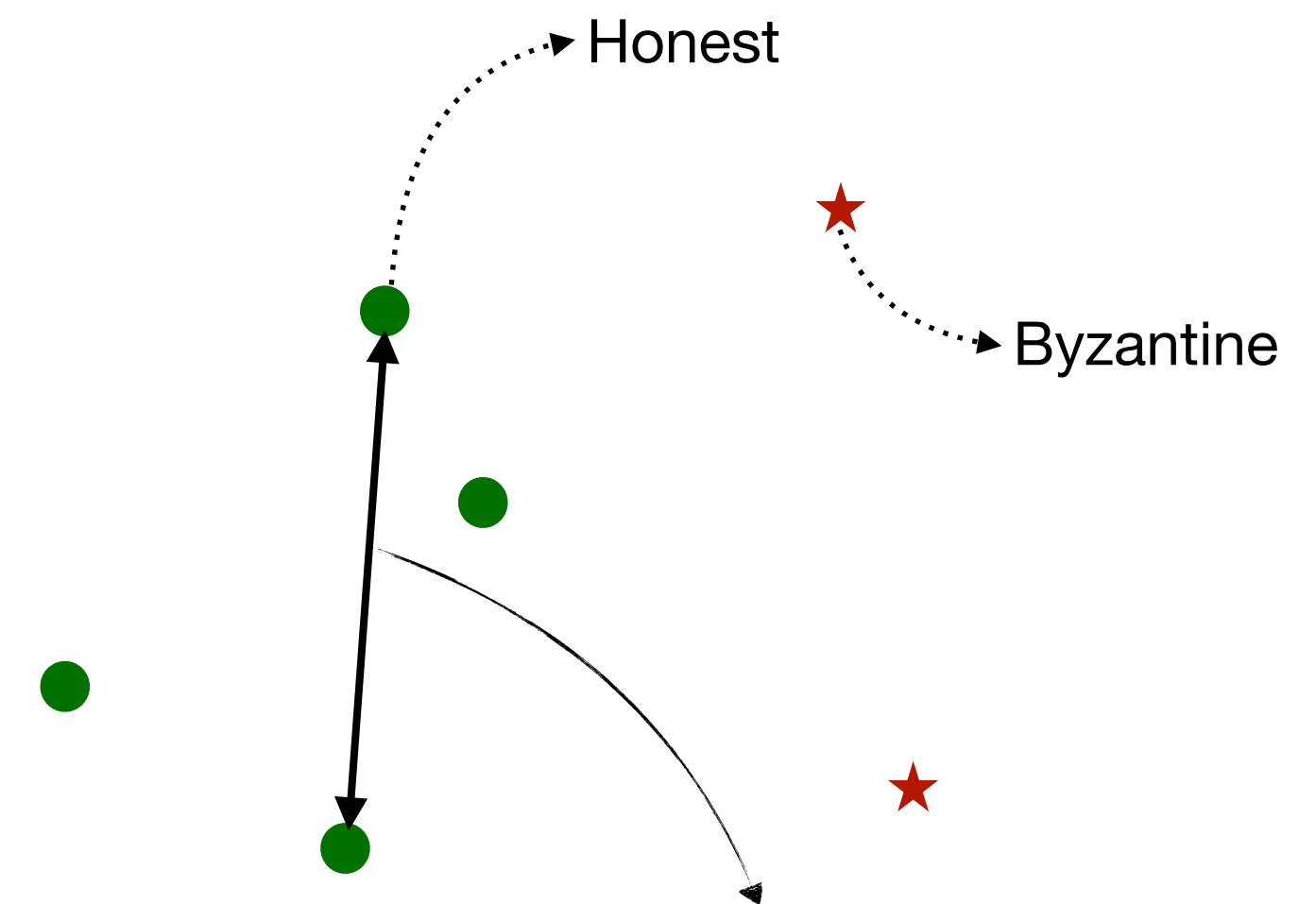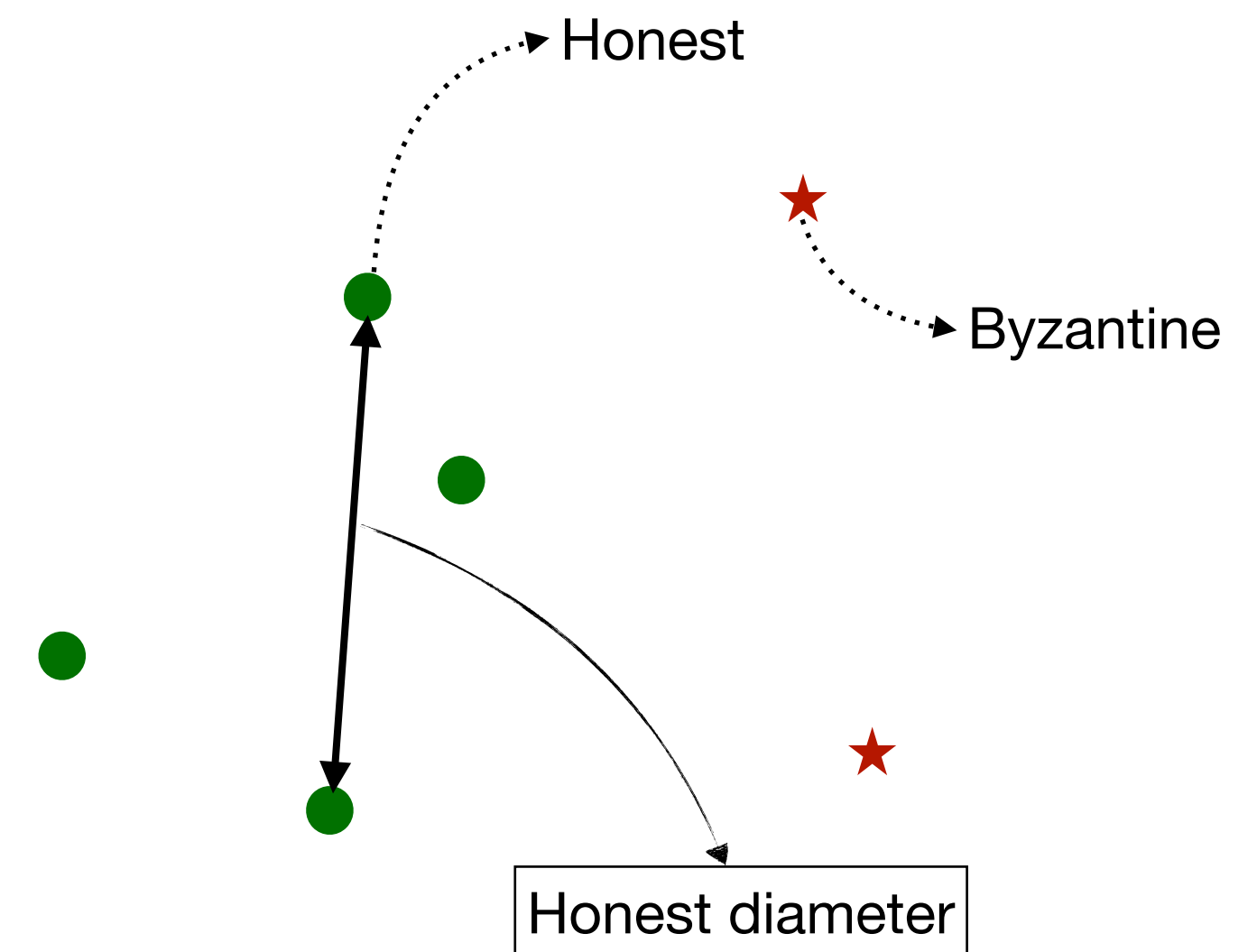
$\exists\ \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall\ S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i, j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient

Honest

Byzantine

Honest average

Honest diameter

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

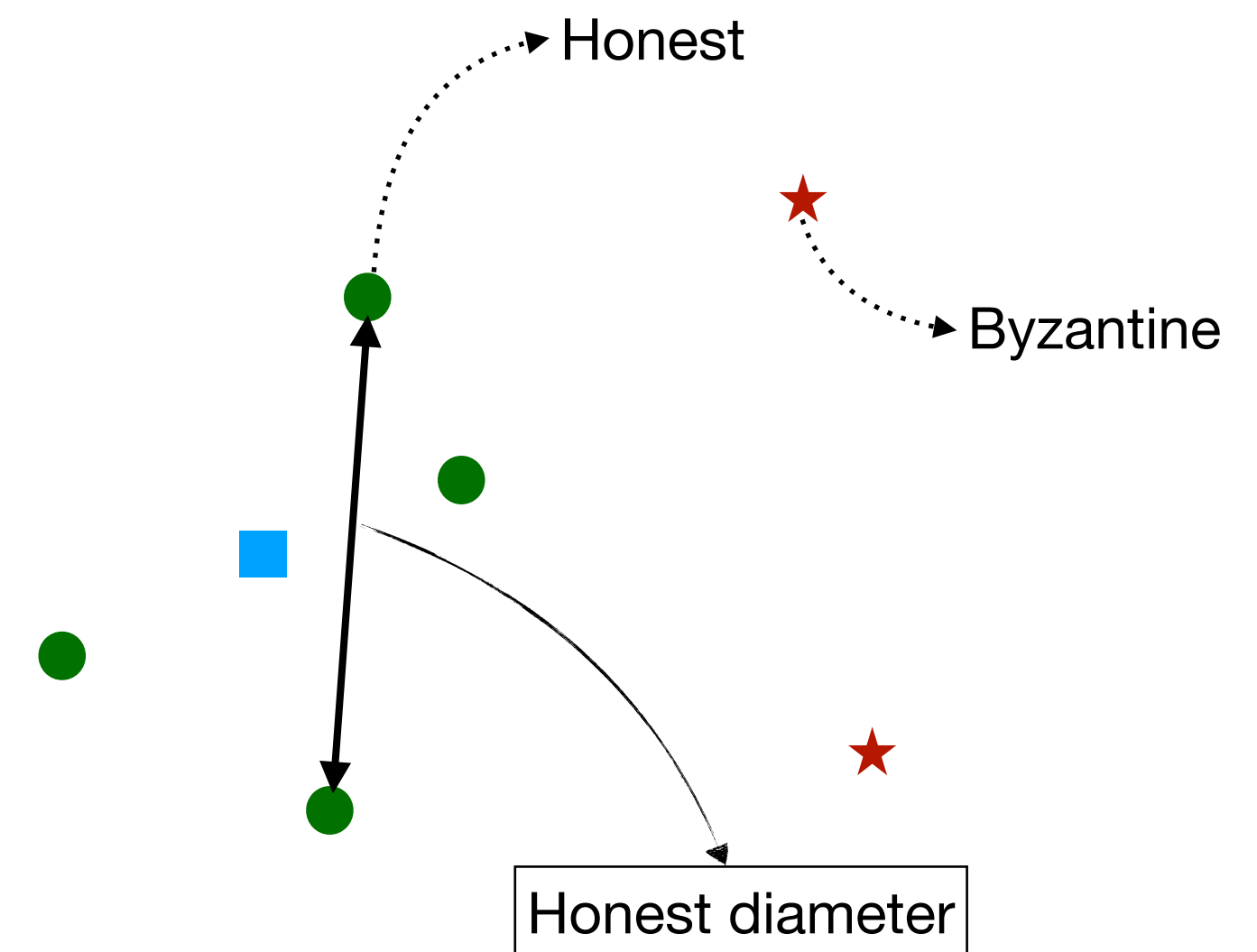$(f, \lambda)$-*resilient averaging (RESA)* -

$\exists\ \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall\ S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i, j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient



Honest

Byzantine

Honest average

Honest diameter

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

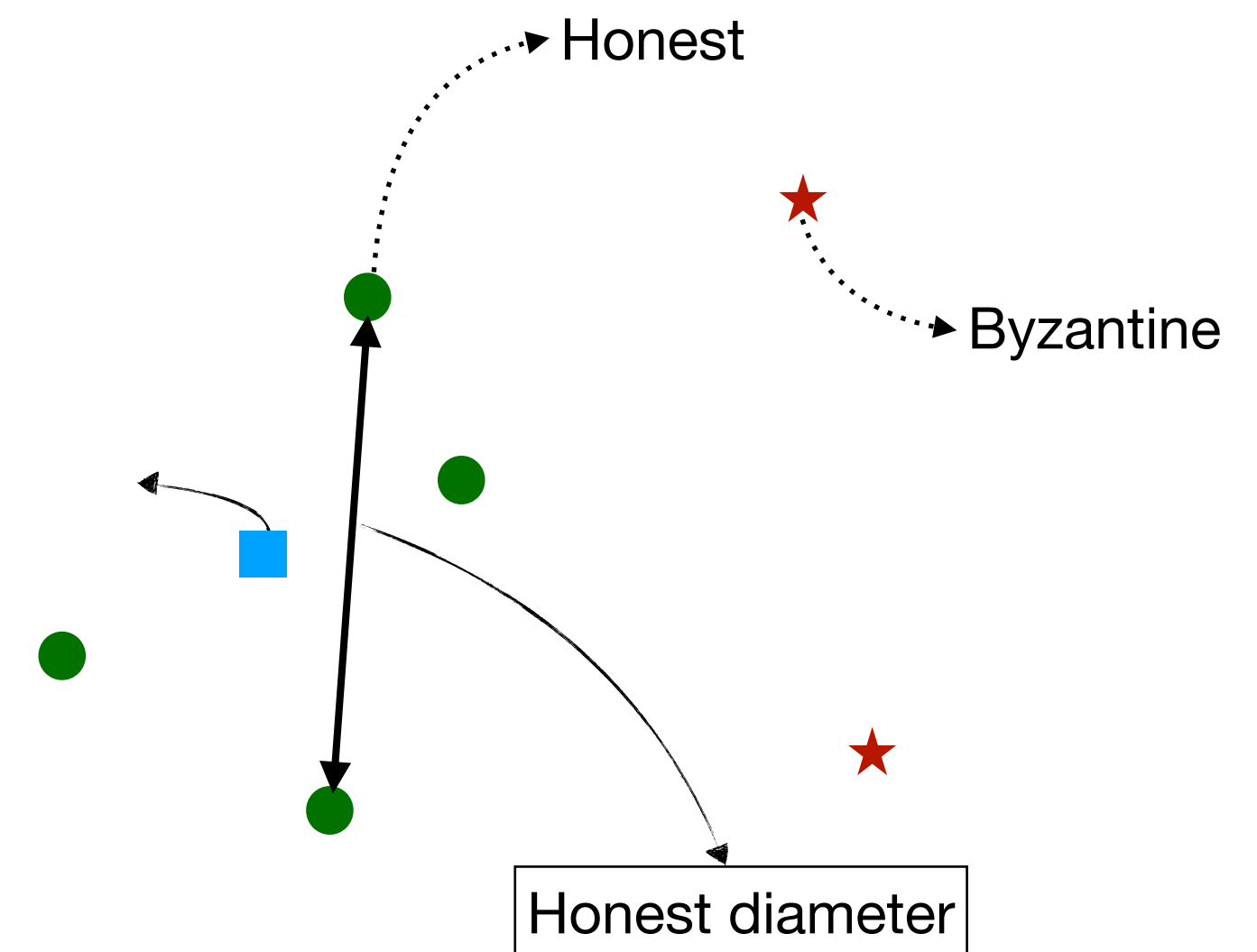$(f, \lambda)$-**resilient averaging (RESA)** -

$\exists\ \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall\ S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i,j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient

Region of output

Honest

Byzantine

Honest average

Honest diameter

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

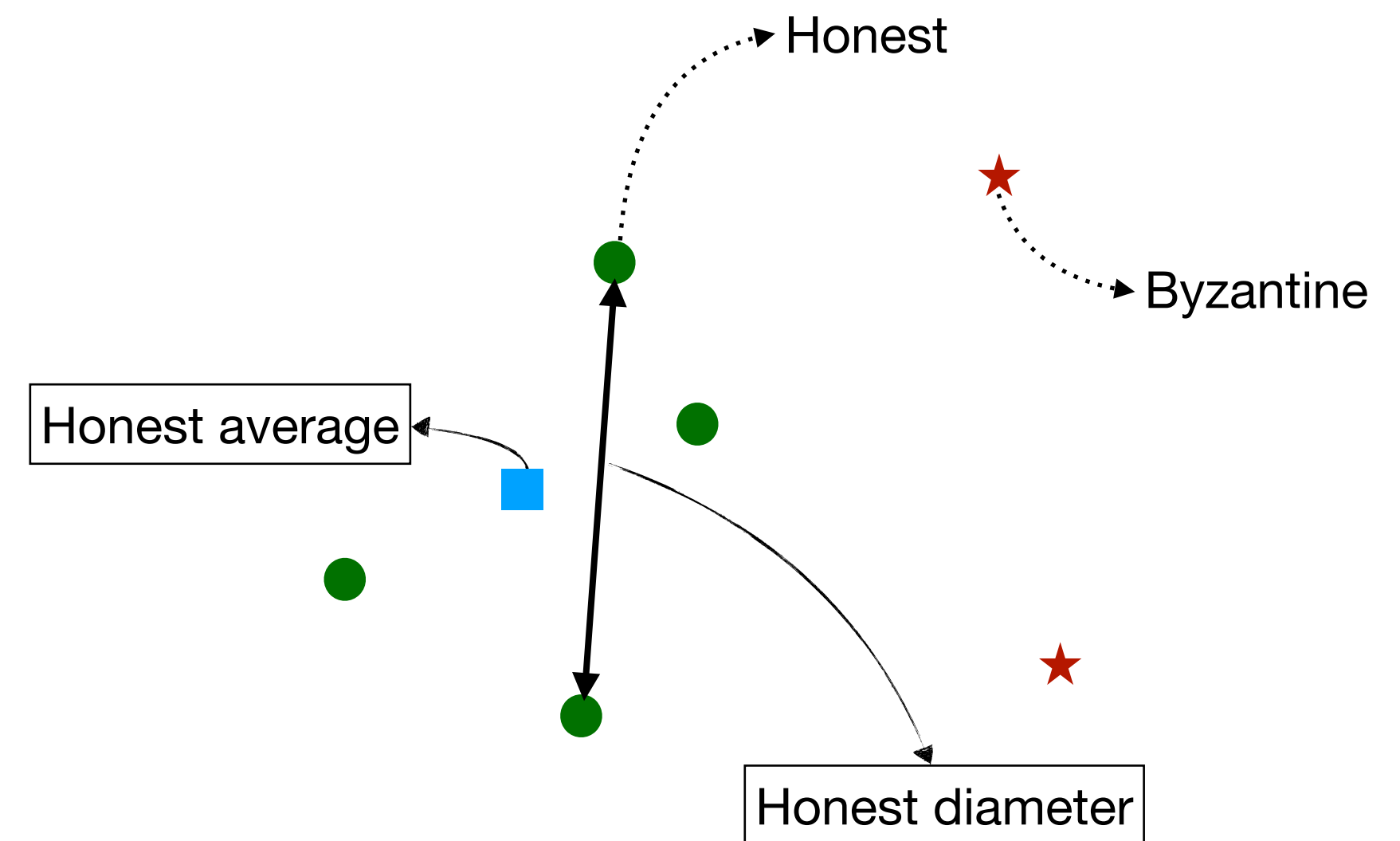$(f, \lambda)$-***resilient averaging (RESA)*** -

$\exists\ \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall\ S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i,j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient



Honest

Region of output

Byzantine

Honest average

Honest diameter

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

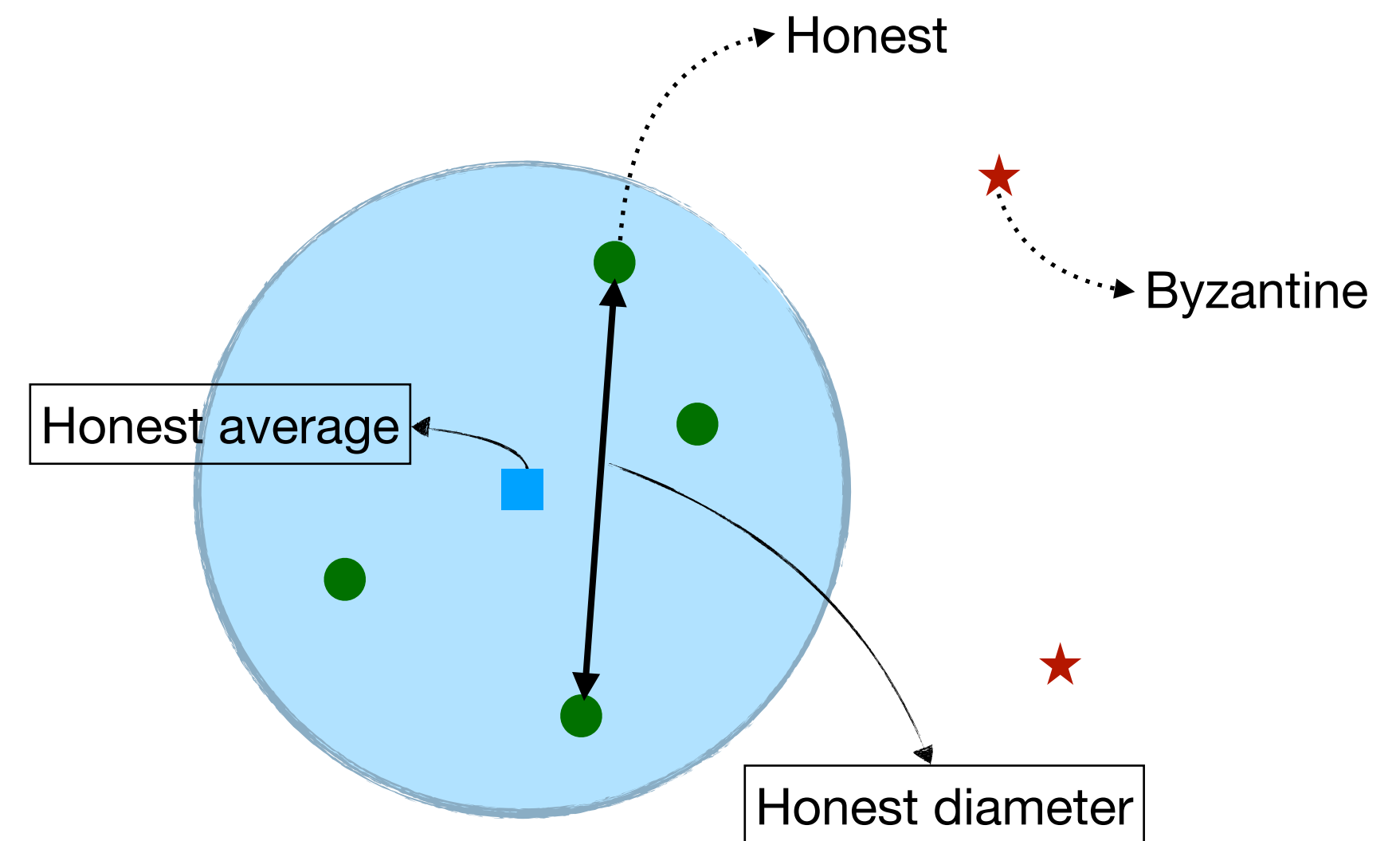$(f, \lambda)$-*resilient averaging (RESA)* -

$\exists \ \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall \ S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i, j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient

Honest

Region of output

Byzantine

Honest average

Honest diameter

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

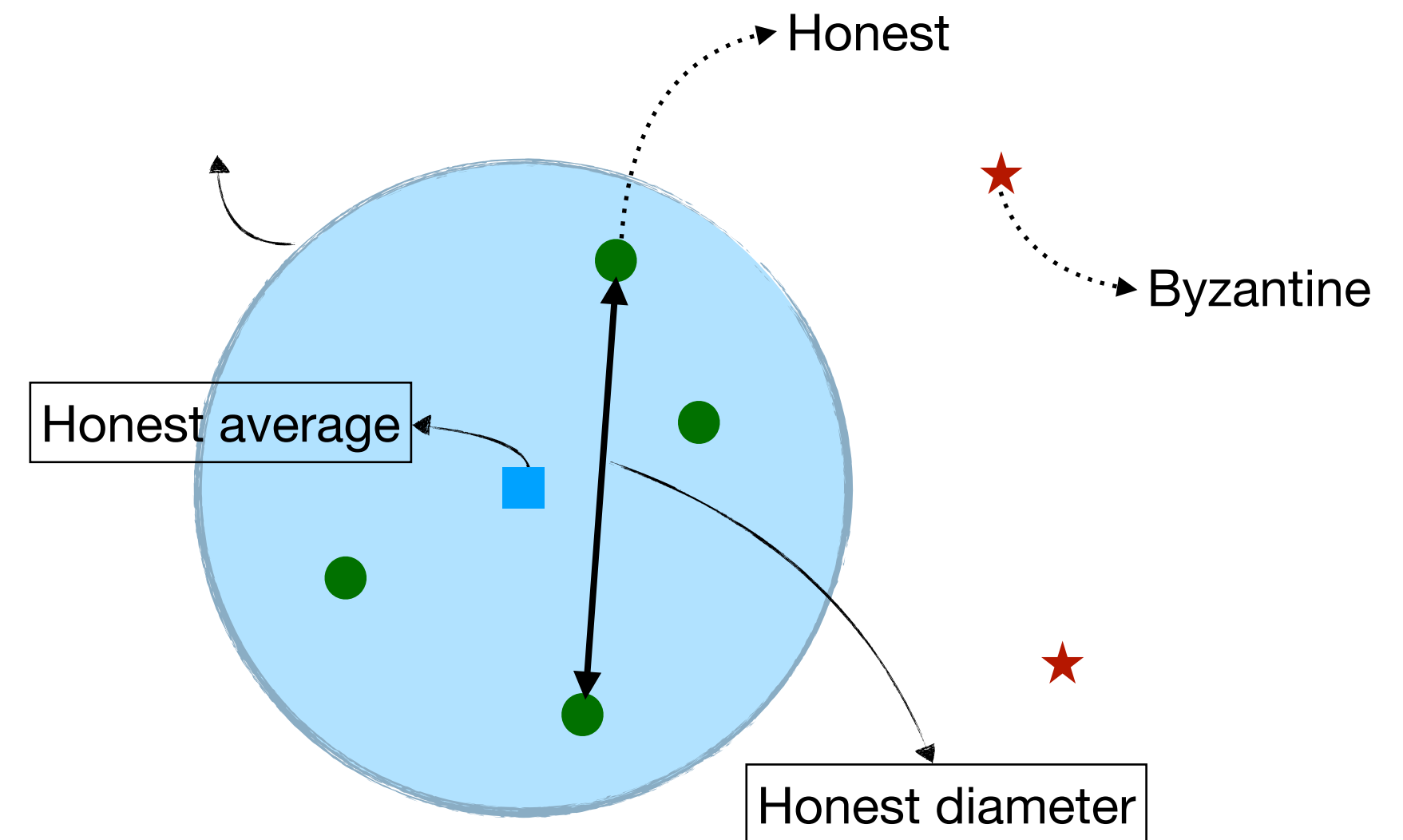$(f, \lambda)$-***resilient averaging (RESA)*** -

$\exists \ \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall \ S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i,j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient



Honest

Region of output

Byzantine

Honest average

Honest diameter

$\lambda \times$ honest diameter

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

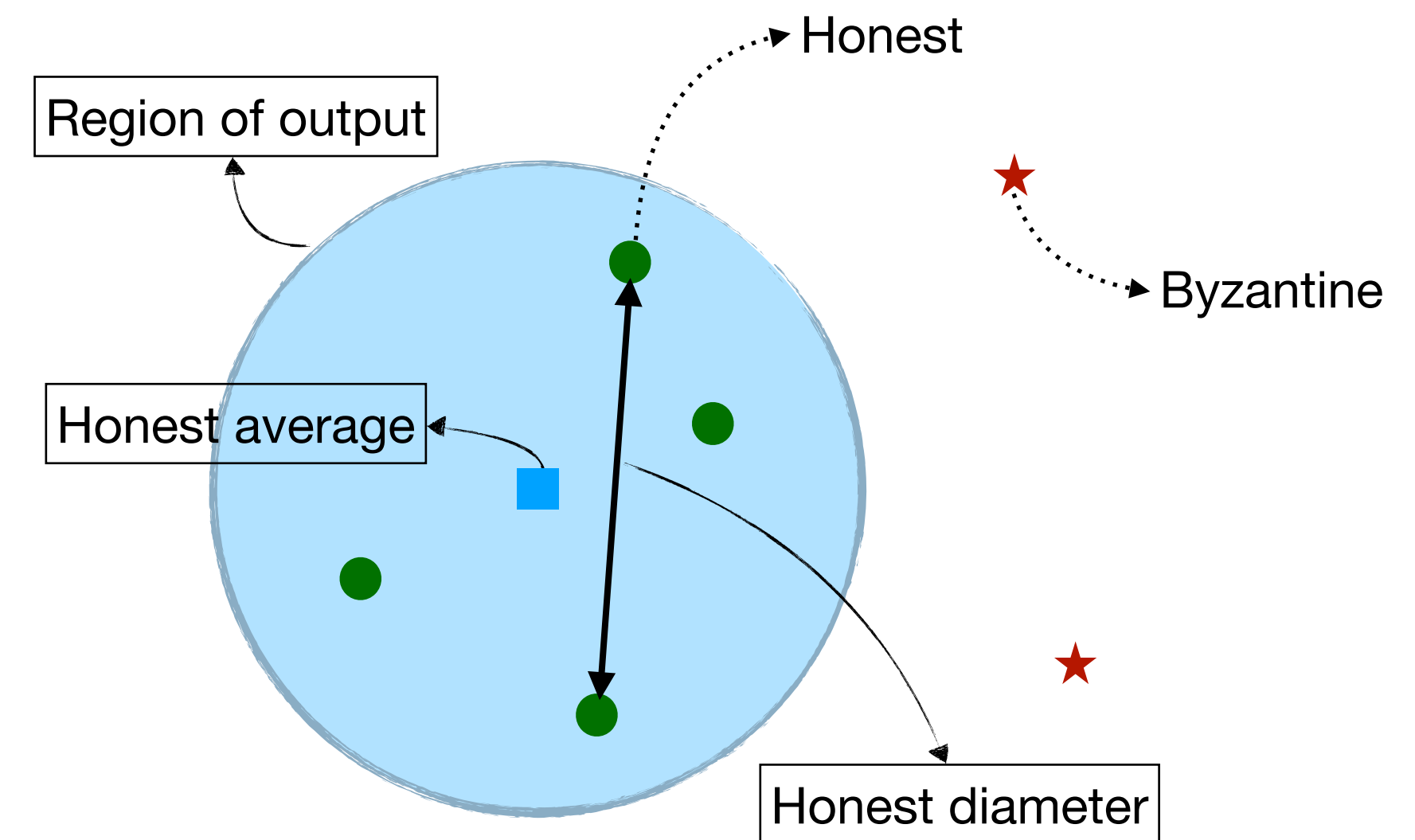**$(f, \lambda)$-*resilient averaging (RESA)* -**

$\exists\ \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall\ S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i,j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient

Region of output

Honest

Byzantine

Honest average

Honest diameter

$\lambda \times$ honest diameter

If all honest values are equal then RESA outputs that value.

# Resilient Averaging (RESA)

Robustness property of $F$ - **key** to optimally utilize Momentum

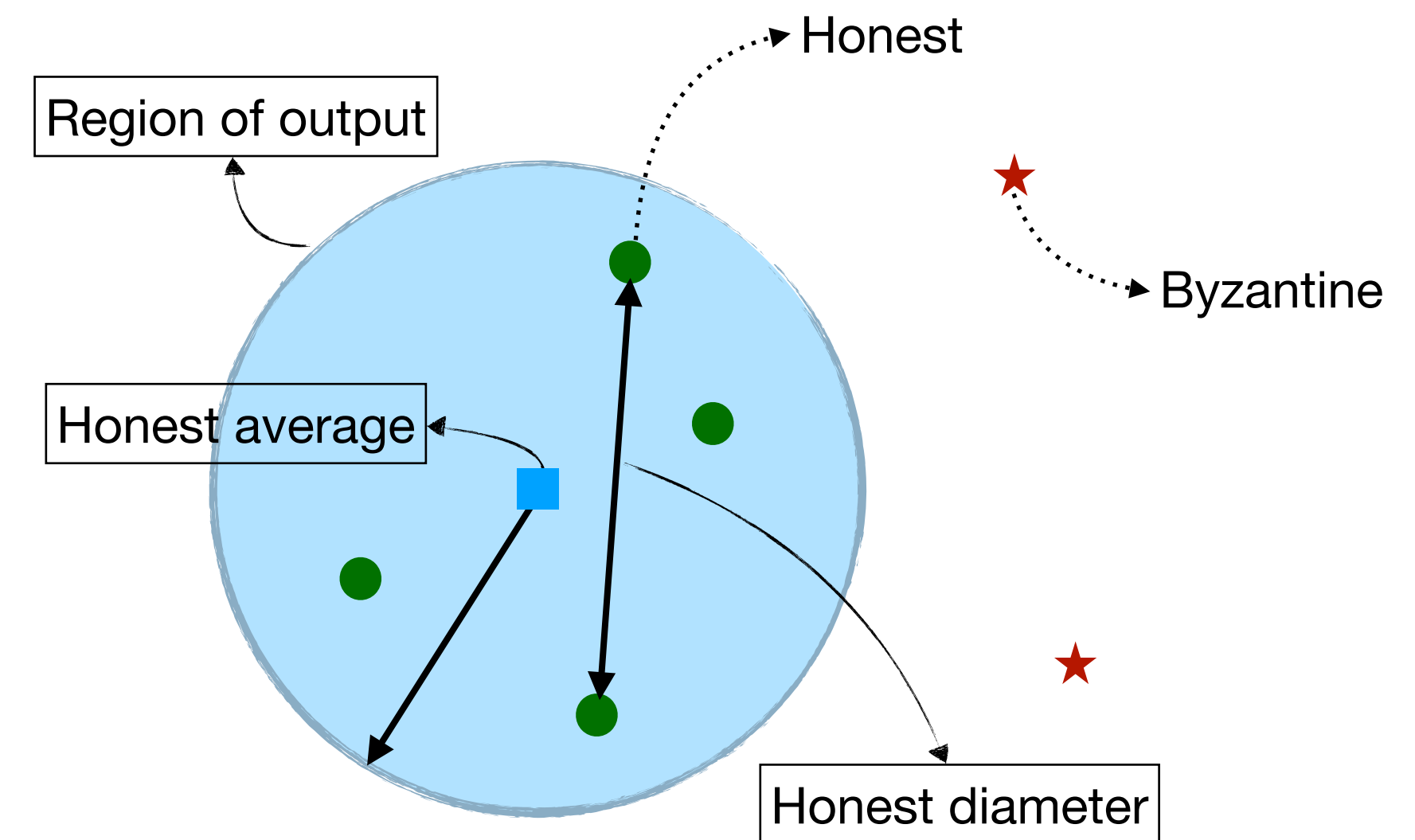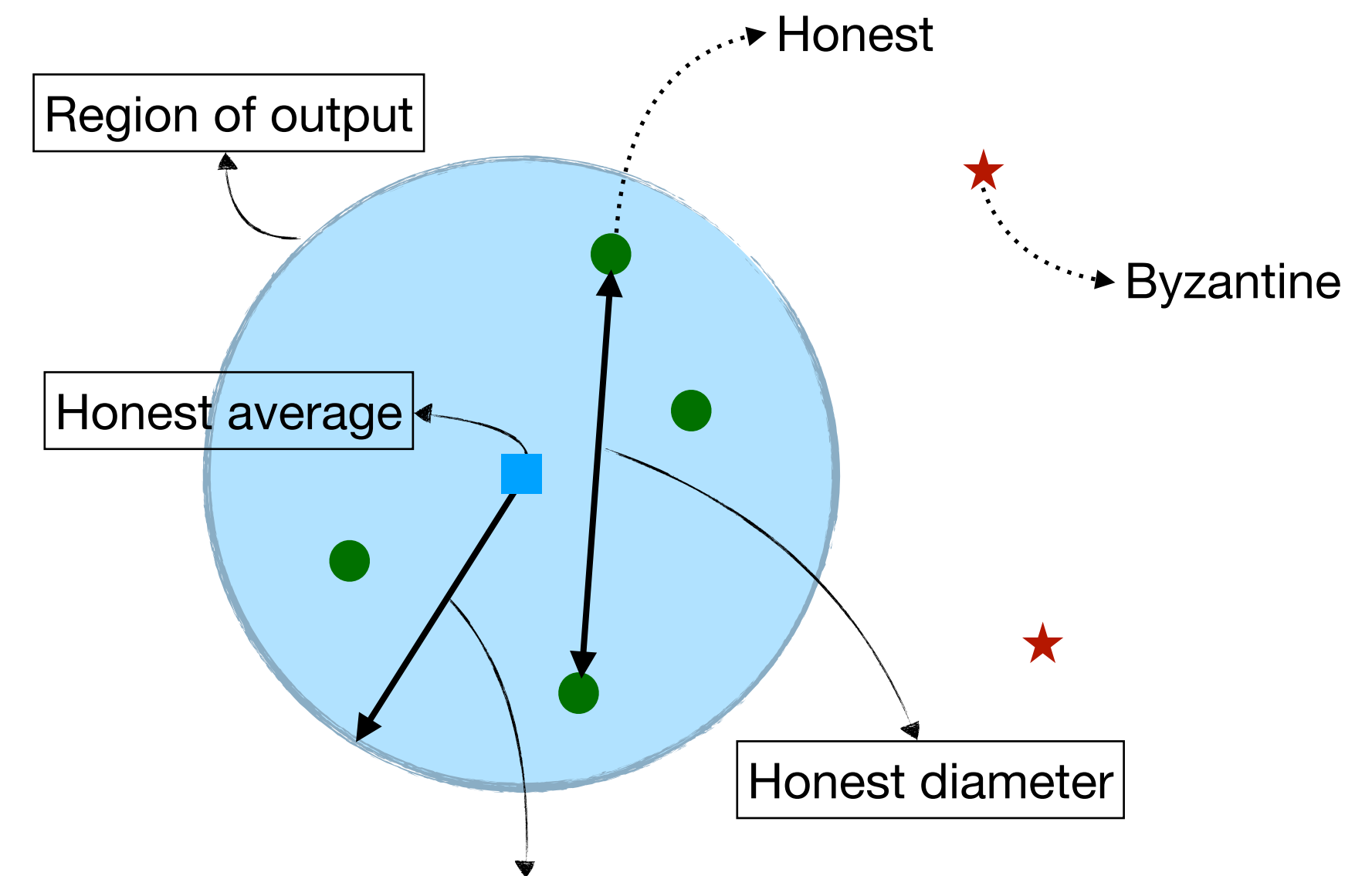$(f, \lambda)$-**resilient averaging (RESA)** -

$\exists \; \lambda \geq 0$ s.t. $\forall$ set of $n$ vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

and $\forall \; S \subset \{1, \ldots, n\}$ of size $n - f$,

$$\|F(x_1, \ldots, x_n) - \bar{x}_S\| \leq \lambda \max_{i, j \in S} \|x_i - x_j\|$$

Avg. of $\{x_i\}_{i \in S}$

Resilience coefficient

Honest

Region of output

Byzantine

Honest average

Honest diameter

$\lambda \times$ honest diameter

**Sanity check**

If all honest values are equal then RESA outputs that value.

# Resilience Coefficients of some Aggregation Rules

# Resilience Coefficients of some Aggregation Rules

Smaller the $\lambda$ better the Byzantine resilience

# Resilience Coefficients of some Aggregation Rules

Smaller the $\lambda$ better the Byzantine resilience

**Lower bound:** $\lambda \geq \dfrac{f}{n-f}$

# Resilience Coefficients of some Aggregation Rules

Smaller the $\lambda$ better the Byzantine resilience

**Lower bound:** $\lambda \geq \dfrac{f}{n-f}$

| Aggregation: | MDA | TM | Geometric Median | Krum | Median |
|---|---|---|---|---|---|
|  |  |  |  |  |  |

# Resilience Coefficients of some Aggregation Rules

Smaller the $\lambda$ better the Byzantine resilience

**Lower bound:** $\lambda \geq \dfrac{f}{n-f}$

| Aggregation: | MDA | TM | Geometric Median | Krum | Median |
|---|---|---|---|---|---|
| $\lambda$ | | | | | |

# Resilience Coefficients of some Aggregation Rules

Smaller the $\lambda$ better the Byzantine resilience

**Lower bound:** $\lambda \geq \dfrac{f}{n-f}$

| Aggregation: | MDA | TM | Geometric Median | Krum | Median |
|---|---|---|---|---|---|
| $\lambda$ | $\dfrac{2f}{n-f}$ | | | | |

# Resilience Coefficients of some Aggregation Rules

Smaller the $\lambda$ better the Byzantine resilience

**Lower bound:** $\lambda \geq \dfrac{f}{n-f}$

| Aggregation: | MDA | TM | Geometric Median | Krum | Median |
|:---:|:---:|:---:|:---:|:---:|:---:|
| $\lambda$ | $\dfrac{2f}{n-f}$ | $\dfrac{f}{2\sqrt{n-f}}$ | | | |

# Resilience Coefficients of some Aggregation Rules

Smaller the $\lambda$ better the Byzantine resilience

**Lower bound:** $\lambda \geq \dfrac{f}{n-f}$

| Aggregation: | MDA | TM | Geometric Median | Krum | Median |
|---|---|---|---|---|---|
| $\lambda$ | $\dfrac{2f}{n-f}$ | $\dfrac{f}{2\sqrt{n-f}}$ | $1 + \dfrac{n-f}{2\sqrt{n(n-2f)}}$ | | |

# Resilience Coefficients of some Aggregation Rules

Smaller the $\lambda$ better the Byzantine resilience

**Lower bound:** $\lambda \geq \dfrac{f}{n-f}$

| Aggregation: | MDA | TM | Geometric Median | Krum | Median |
|---|---|---|---|---|---|
| $\lambda$ | $\dfrac{2f}{n-f}$ | $\dfrac{f}{2\sqrt{n-f}}$ | $1 + \dfrac{n-f}{2\sqrt{n(n-2f)}}$ | $1 + \sqrt{\dfrac{n-f}{n-2f}}$ | |

# Resilience Coefficients of some Aggregation Rules

Smaller the $\lambda$ better the Byzantine resilience

**Lower bound:** $\lambda \geq \dfrac{f}{n-f}$

| Aggregation: | MDA | TM | Geometric Median | Krum | Median |
|---|---|---|---|---|---|
| $\lambda$ | $\dfrac{2f}{n-f}$ | $\dfrac{f}{2\sqrt{n-f}}$ | $1 + \dfrac{n-f}{2\sqrt{n(n-2f)}}$ | $1 + \sqrt{\dfrac{n-f}{n-2f}}$ | $\dfrac{n}{2\sqrt{n-f}}$ |

# Resilience Coefficients of some Aggregation Rules

Smaller the $\lambda$ better the Byzantine resilience

**Lower bound:** $\lambda \geq \dfrac{f}{n-f}$

| Aggregation: | MDA | TM | Geometric Median | Krum | Median |
|---|---|---|---|---|---|
| $\lambda$ | $\dfrac{2f}{n-f}$ | $\dfrac{f}{2\sqrt{n-f}}$ | $1 + \dfrac{n-f}{2\sqrt{n(n-2f)}}$ | $1 + \sqrt{\dfrac{n-f}{n-2f}}$ | $\dfrac{n}{2\sqrt{n-f}}$ |

# Resilience Coefficients of some Aggregation Rules

Smaller the $\lambda$ better the Byzantine resilience

**Lower bound:** $\lambda \geq \dfrac{f}{n-f}$

| Aggregation: | MDA | TM | Geometric Median | Krum | Median |
|:---:|:---:|:---:|:---:|:---:|:---:|
| $\lambda$ | $\dfrac{2f}{n-f}$ | $\dfrac{f}{2\sqrt{n-f}}$ | $1 + \dfrac{n-f}{2\sqrt{n(n-2f)}}$ | $1 + \sqrt{\dfrac{n-f}{n-2f}}$ | $\dfrac{n}{2\sqrt{n-f}}$ |

# Resilience Coefficients of some Aggregation Rules

Smaller the $\lambda$ better the Byzantine resilience

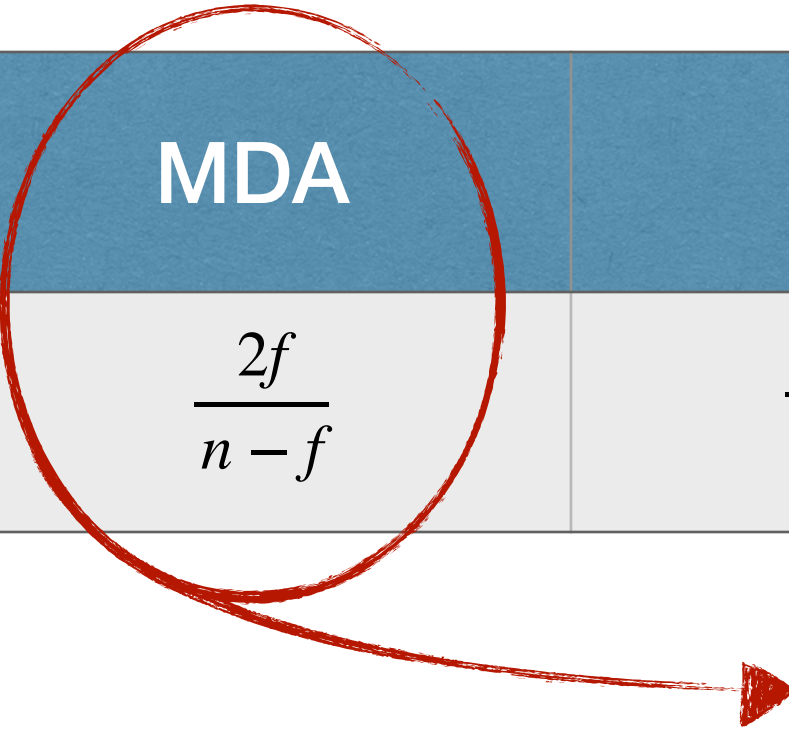**Lower bound:** $\lambda \geq \dfrac{f}{n-f}$

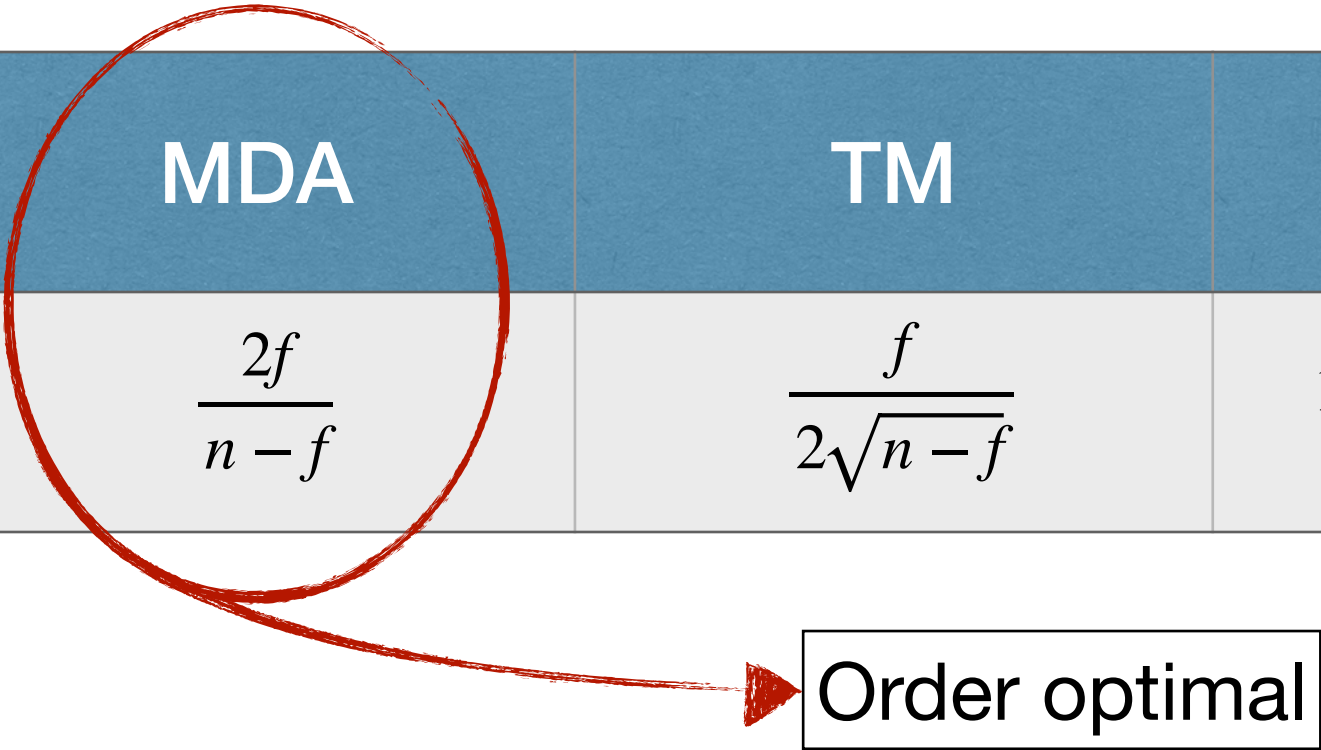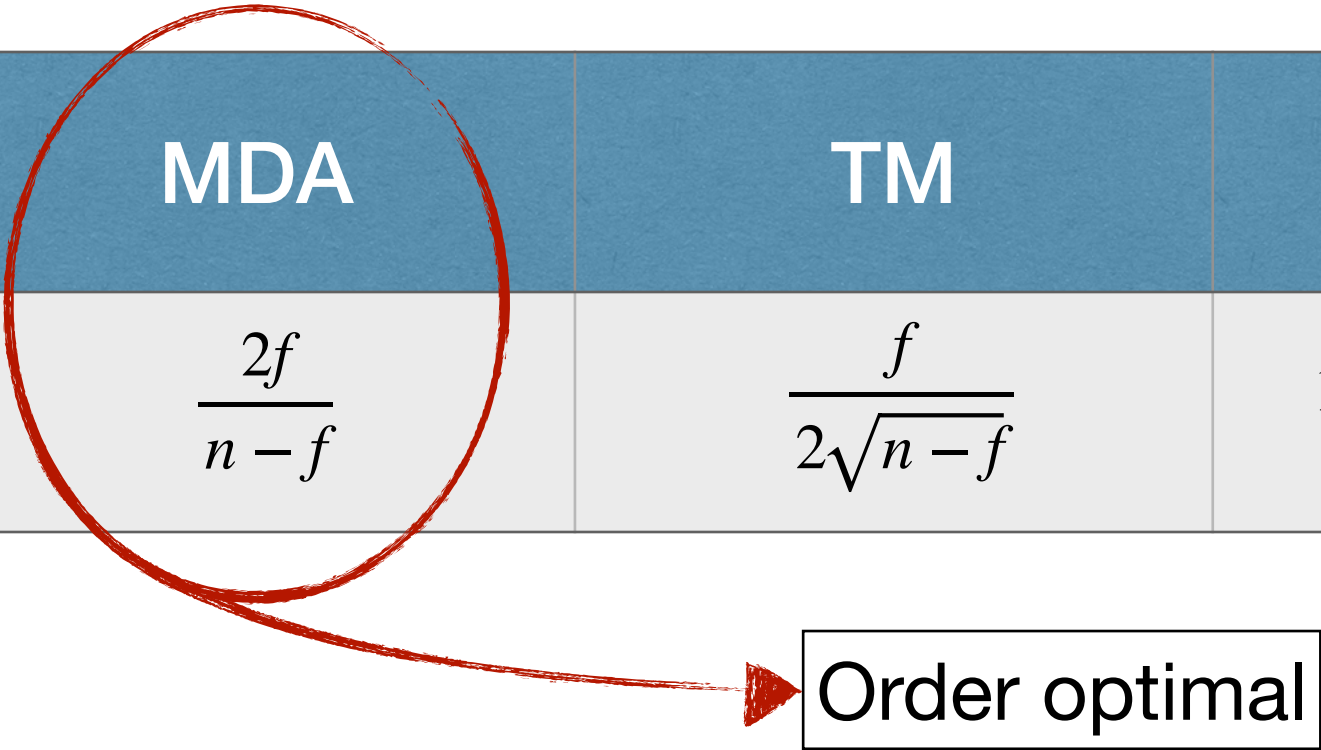| Aggregation: | MDA | TM | Geometric Median | Krum | Median |
|:---:|:---:|:---:|:---:|:---:|:---:|
| $\lambda$ | $\dfrac{2f}{n-f}$ | $\dfrac{f}{2\sqrt{n-f}}$ | $1 + \dfrac{n-f}{2\sqrt{n(n-2f)}}$ | $1 + \sqrt{\dfrac{n-f}{n-2f}}$ | $\dfrac{n}{2\sqrt{n-f}}$ |

Order optimal

# Resilience Coefficients of some Aggregation Rules

Smaller the $\lambda$ better the Byzantine resilience

**Lower bound:** $\lambda \geq \dfrac{f}{n-f}$

| Aggregation: | MDA | TM | Geometric Median | Krum | Median |
|---|---|---|---|---|---|
| $\lambda$ | $\dfrac{2f}{n-f}$ | $\dfrac{f}{2\sqrt{n-f}}$ | $1 + \dfrac{n-f}{2\sqrt{n(n-2f)}}$ | $1 + \sqrt{\dfrac{n-f}{n-2f}}$ | $\dfrac{n}{2\sqrt{n-f}}$ |

Order optimal

The values hold for any $f < n/2$

# Byzantine Resilience Guarantee

# Byzantine Resilience Guarantee

**Server** updates: $\quad \theta_{t+1} \longleftarrow \theta_t - \gamma_t \, R_t$

where $\quad R_t = F\left(m_t^1, \ldots, m_t^n\right)$ and $\quad m_t^i = \beta \, m_{t-1}^i + (1 - \beta) \, g_t^i \quad$ for each honest node $i$

# Byzantine Resilience Guarantee

**Server** updates: $\quad \theta_{t+1} \longleftarrow \theta_t - \gamma_t \, R_t$

where $\quad R_t = F\left(m_t^1, \ldots, m_t^n\right) \quad$ and $\quad m_t^i = \beta \, m_{t-1}^i + (1-\beta) \, g_t^i \quad$ for each honest node $i$

Upon $T$ iterations of **RESAM**

# Byzantine Resilience Guarantee

**Server** updates:  $\theta_{t+1} \longleftarrow \theta_t - \gamma_t \, R_t$

where  $R_t = F\left(m_t^1, \ldots, m_t^n\right)$  and  $m_t^i = \beta \, m_{t-1}^i + (1 - \beta) \, g_t^i$  for each honest node $i$

Upon $T$ iterations of **RESAM**

$(f, \lambda)$**-RESA** aggregation rule $F$ + **M**omentum with $\beta = \sqrt{1 - c/T}$

Let $\widehat{\theta}$ be chosen uniformly from $\theta_1, \ldots, \theta_T$ then

# Byzantine Resilience Guarantee

**Server** updates:   $\theta_{t+1} \longleftarrow \theta_t - \gamma_t\, R_t$

where   $R_t = F\left(m_t^1, \ldots, m_t^n\right)$   and   $m_t^i = \beta\, m_{t-1}^i + (1-\beta)\, g_t^i$   for each honest node $i$

Upon $T$ iterations of **RESAM**

$(f, \lambda)$-**RESA** aggregation rule $F$ + **M**omentum with $\beta = \sqrt{1 - c/T}$

Let  $\widehat{\theta}$  be chosen uniformly from  $\theta_1, \ldots, \theta_T$  then

$$\mathbb{E}\left[\|\nabla Q\left(\widehat{\theta}\right)\|^2\right] \le \epsilon \in \mathcal{O}\left(\sqrt{\left(\frac{1}{n-f} + \lambda^2(n-f)\right)\frac{\sigma^2}{T}}\right)$$

# Byzantine Resilience Guarantee

**Server** updates: $\quad \theta_{t+1} \longleftarrow \theta_t - \gamma_t \, R_t$

where $\quad R_t = F\left(m_t^1, \ldots, m_t^n\right) \quad$ and $\quad m_t^i = \beta \, m_{t-1}^i + (1-\beta) \, g_t^i \quad$ for each honest node $i$

Upon $T$ iterations of **RESAM**

$(f, \lambda)$-**RESA** aggregation rule $F$ + **M**omentum with $\beta = \sqrt{1 - c/T}$

Let $\widehat{\theta}$ be chosen uniformly from $\theta_1, \ldots, \theta_T$ then

$(f, \epsilon)-$Resilience

$$\mathbb{E}\left[\|\nabla Q\left(\widehat{\theta}\right)\|^2\right] \leq \epsilon \in \mathcal{O}\left(\sqrt{\left(\frac{1}{n-f} + \lambda^2(n-f)\right)\frac{\sigma^2}{T}}\right)$$

# Byzantine Resilience Guarantee

**Server** updates: $\quad \theta_{t+1} \longleftarrow \theta_t - \gamma_t\, R_t$

where $\quad R_t = F\left(m_t^1, \ldots, m_t^n\right) \quad$ and $\quad m_t^i = \beta\, m_{t-1}^i + (1 - \beta)\, g_t^i \quad$ for each honest node $i$

Upon $T$ iterations of **RESAM**

$(f, \lambda)$-**RESA** aggregation rule $F$ + **M**omentum with $\beta = \sqrt{1 - c/T}$

Let $\widehat{\theta}$ be chosen uniformly from $\theta_1, \ldots, \theta_T$ then

$$\mathbb{E}\left[\|\nabla Q\left(\widehat{\theta}\right)\|^2\right] \leq \epsilon \in \mathcal{O}\left(\sqrt{\left(\frac{1}{n-f} + \lambda^2(n-f)\right)\frac{\sigma^2}{T}}\right)$$

$(f, \epsilon) -$ Resilience

# Byzantine Resilience Guarantee

**Server** updates: $\quad \theta_{t+1} \longleftarrow \theta_t - \gamma_t \, R_t$

where $\quad R_t = F\left(m_t^1, \ldots, m_t^n\right) \quad$ and $\quad m_t^i = \beta \, m_{t-1}^i + (1-\beta) \, g_t^i \quad$ for each honest node $i$

Upon $T$ iterations of **RESAM**

$(f, \lambda)$-**RESA** aggregation rule $F$ + **M**omentum with $\beta = \sqrt{1 - c/T}$

Let $\widehat{\theta}$ be chosen uniformly from $\theta_1, \ldots, \theta_T$ then

$(f, \epsilon) -$ Resilience

$$\mathbb{E}\left[\|\nabla Q\left(\widehat{\theta}\right)\|^2\right] \leq \epsilon \in \mathcal{O}\left(\sqrt{\left(\frac{1}{n-f} + \lambda^2(n-f)\right)\frac{\sigma^2}{T}}\right)$$

# Byzantine Resilience Guarantee

**Server** updates: $\quad \theta_{t+1} \longleftarrow \theta_t - \gamma_t\, R_t$

where $\quad R_t = F\left(m_t^1, \ldots, m_t^n\right) \quad$ and $\quad m_t^i = \beta\, m_{t-1}^i + (1 - \beta)\, g_t^i \quad$ for each honest node $i$

Upon $T$ iterations of **RESAM**

$(f, \lambda)$-**RESA** aggregation rule $F$ + **M**omentum with $\beta = \sqrt{1 - c/T}$

Let $\hat{\theta}$ be chosen uniformly from $\theta_1, \ldots, \theta_T$ then

$$\mathbb{E}\left[\|\nabla Q\left(\hat{\theta}\right)\|^2\right] \leq \epsilon \in \mathcal{O}\left(\sqrt{\left(\frac{1}{n-f} + \lambda^2(n-f)\right)\frac{\sigma^2}{T}}\right)$$

$(f, \epsilon)-$Resilience

Compares overall efficiency for
different aggregation rules

# Empirical Evidence

# Empirical Evidence

* MNIST classification task

# Empirical Evidence

**Without Momentum**

* MNIST classification task

# Empirical Evidence



Without
Momentum

* MNIST classification task

# Empirical Evidence

**Label-flipping**

* MNIST classification task

# Empirical Evidence

**Label-flipping**

* MNIST classification task

# Empirical Evidence

**Without Momentum**





**Label-flipping**

**Little is enough** *(Baruch et al., 2019)*

* MNIST classification task

# Empirical Evidence



**Without Momentum**

**With Momentum**

**Label-flipping**

**Little is enough** *(Baruch et al., 2019)*

* MNIST classification task

# Empirical Evidence

**Without Momentum**





**With Momentum**



**Label-flipping**

**Little is enough** *(Baruch et al., 2019)*

* MNIST classification task

# Empirical Evidence



**Without Momentum**

**With Momentum**

**Label-flipping**

**Little is enough** *(Baruch et al., 2019)*

* MNIST classification task

# *Efficiency* + Byzantine Resilience ?

# *Efficiency* + Byzantine Resilience ?
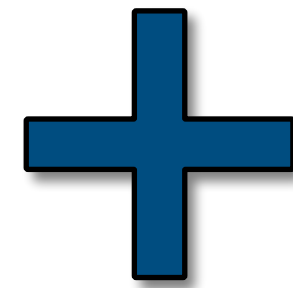
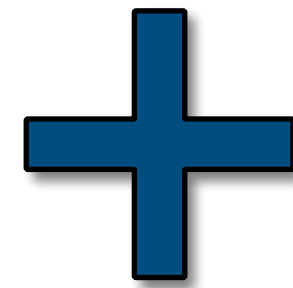**CAN THE SERVER STILL OBTAIN**

$$\theta* \in \left( \theta \; ; \; \nabla Q(\theta) = 0 \right)$$

# *Efficiency* + Byzantine Resilience ?

CAN THE SERVER STILL OBTAIN

$$\theta* \in \left( \theta \ ; \ \nabla Q(\theta) = 0 \right)$$

+

# *Efficiency* + Byzantine Resilience ?

CAN THE SERVER STILL OBTAIN

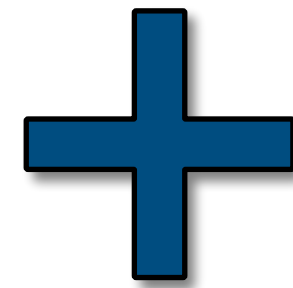$$\theta* \in \big(\theta \; ; \; \nabla Q(\theta) = 0\big)$$

**+**

FRACTIONAL WORKLOAD PER NODE

# *Efficiency* + Byzantine Resilience ?

**CAN THE SERVER STILL OBTAIN**

$$\theta* \in \left(\theta \; ; \; \nabla Q(\theta) = 0\right)$$
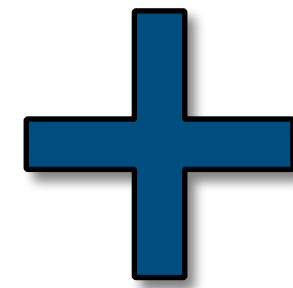
**+**

**FRACTIONAL WORKLOAD PER NODE**

$(f, \epsilon)-$Resilience

# *Efficiency* + Byzantine Resilience ?

CAN THE SERVER STILL OBTAIN

$$\theta^* \in \left(\theta \; ; \; \nabla Q(\theta) = 0\right)$$

**+**

FRACTIONAL WORKLOAD PER NODE

$(f, \epsilon)-$Resilience

$$\epsilon \in \mathcal{O}\left(\sqrt{\left(\frac{1}{n-f} + \lambda^2(n-f)\right)\frac{\sigma^2}{T}}\right)$$

# *Efficiency* + Byzantine Resilience ?

CAN THE SERVER STILL OBTAIN
$$\theta* \in \big(\theta \; ; \; \nabla Q(\theta) = 0\big)$$

**+**

FRACTIONAL WORKLOAD PER NODE

$(f, \epsilon)-$Resilience

$$\epsilon \in \mathcal{O}\left(\sqrt{\left(\frac{1}{n-f} + \lambda^2(n-f)\right)\frac{\sigma^2}{T}}\right)$$

# *Efficiency* + Byzantine Resilience ?

**CAN THE SERVER STILL OBTAIN**

$$\theta* \in \left( \theta \; ; \; \nabla Q(\theta) = 0 \right)$$
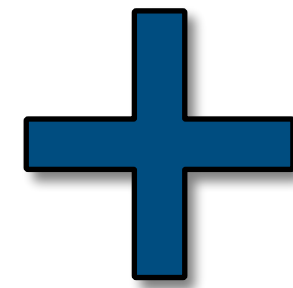
**+**

**FRACTIONAL WORKLOAD PER NODE**
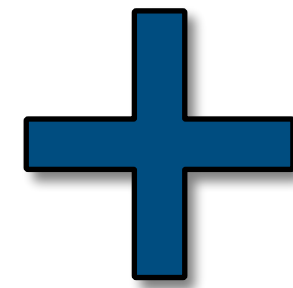
$(f, \epsilon)-$Resilience

$$\epsilon \in \mathcal{O}\left( \sqrt{\left( \frac{1}{n-f} + \lambda^2(n-f) \right) \frac{\sigma^2}{T}} \right)$$

# *Efficiency* + Byzantine Resilience ?

CAN THE SERVER STILL OBTAIN

$$\theta* \in \left( \theta \; ; \; \nabla Q(\theta) = 0 \right)$$

**+**

FRACTIONAL WORKLOAD PER NODE

$(f, \epsilon)-$Resilience

$$\epsilon \in \mathcal{O}\left( \sqrt{\left( \frac{1}{n-f} + \lambda^2(n-f) \right) \frac{\sigma^2}{T}} \right)$$

Additional workload

# *Efficiency* + Byzantine Resilience ?

**CAN THE SERVER STILL OBTAIN**

$$\theta* \in \left( \theta \; ; \; \nabla Q(\theta) = 0 \right)$$

**+**

**FRACTIONAL WORKLOAD PER NODE**

$(f, \epsilon)-$Resilience
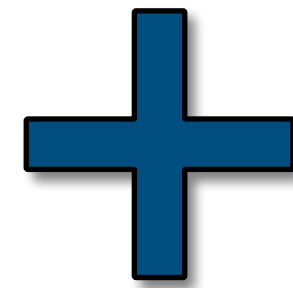
$$\epsilon \in \mathcal{O}\left( \sqrt{\left( \frac{1}{n-f} + \lambda^2(n-f) \right) \frac{\sigma^2}{T}} \right)$$

Additional workload

| Aggregation: | MDA | TM | Geometric Median | Krum* | Median |
|---|---|---|---|---|---|
| $\lambda^2(n-f)$ | $\dfrac{f^2}{n-f}$ | $f^2$ | $\dfrac{(n-f)^3}{(n-2f)n}$ | $\dfrac{(n-f)^2}{n-2f}$ | $n^2$ |

# *Efficiency* + Byzantine Resilience ?

CAN THE SERVER STILL OBTAIN

$$\theta* \in \left( \theta \; ; \; \nabla Q(\theta) = 0 \right)$$

**+**

FRACTIONAL WORKLOAD PER NODE

$(f, \epsilon)-$Resilience
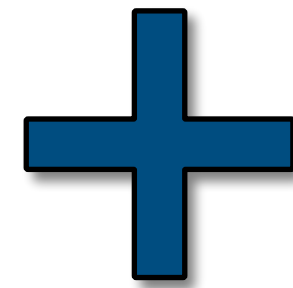
$$\epsilon \in \mathcal{O}\left( \sqrt{\left( \frac{1}{n-f} + \lambda^2(n-f) \right) \frac{\sigma^2}{T}} \right)$$

Additional workload

| Aggregation: | MDA ✅ | TM | Geometric Median | Krum* | Median |
|---|---|---|---|---|---|
| $\lambda^2(n-f)$ | $\dfrac{f^2}{n-f}$ | $f^2$ | $\dfrac{(n-f)^3}{(n-2f)n}$ | $\dfrac{(n-f)^2}{n-2f}$ | $n^2$ |

# *Efficiency* + Byzantine Resilience ?

**CAN THE SERVER STILL OBTAIN**

$$\theta* \in \left(\theta \; ; \; \nabla Q(\theta) = 0\right)$$

**+**

**FRACTIONAL WORKLOAD PER NODE**

$(f, \epsilon)-$Resilience

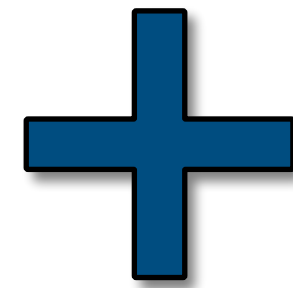$$\epsilon \in \mathcal{O}\left(\sqrt{\left(\frac{1}{n-f} + \lambda^2(n-f)\right)\frac{\sigma^2}{T}}\right)$$

Additional workload

| Aggregation: | MDA ✅ | TM ❌ | Geometric Median ❌ | Krum* ❌ | Median ❌ |
|---|---|---|---|---|---|
| $\lambda^2(n-f)$ | $\dfrac{f^2}{n-f}$ | $f^2$ | $\dfrac{(n-f)^3}{(n-2f)n}$ | $\dfrac{(n-f)^2}{n-2f}$ | $n^2$ |

# *Efficiency* + Byzantine Resilience ?

**CAN THE SERVER STILL OBTAIN**

$$\theta* \in \left( \theta \; ; \; \nabla Q(\theta) = 0 \right)$$

**+**

**FRACTIONAL WORKLOAD PER NODE**

$(f, \epsilon)-$Resilience

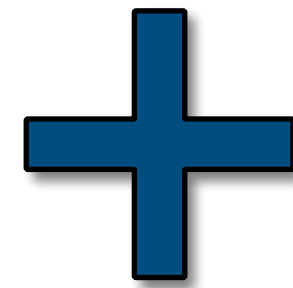$$\epsilon \in \mathcal{O}\left( \sqrt{\left( \frac{1}{n-f} + \lambda^2(n-f) \right) \frac{\sigma^2}{T}} \right)$$

Additional workload

| Aggregation: | MDA ✅ | TM ❌ | Geometric Median ❌ | Krum* ❌ | Median ❌ |
|---|---|---|---|---|---|
| $\lambda^2(n-f)$ | $\dfrac{f^2}{n-f}$ | $f^2$ | $\dfrac{(n-f)^3}{(n-2f)n}$ | $\dfrac{(n-f)^2}{n-2f}$ | $n^2$ |

# *Efficiency* + Byzantine Resilience ?

**CAN THE SERVER STILL OBTAIN**

$$\theta* \in \left(\theta \; ; \; \nabla Q(\theta) = 0\right)$$

**+**

**FRACTIONAL WORKLOAD PER NODE**

$(f, \epsilon)-$Resilience

$$\epsilon \in \mathcal{O}\left(\sqrt{\left(\frac{1}{n-f} + \lambda^2(n-f)\right)\frac{\sigma^2}{T}}\right)$$
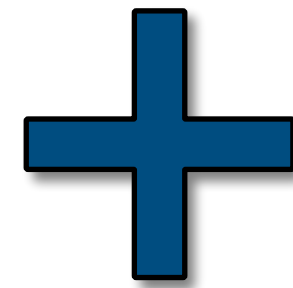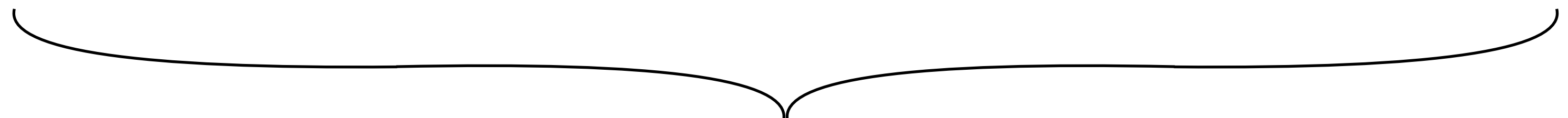
Additional workload

| Aggregation: | MDA ✅ | TM ❌ | Geometric Median ❌ | Krum* ❌ | Median ❌ |
|---|---|---|---|---|---|
| $\lambda^2(n-f)$ | $\dfrac{f^2}{n-f}$ | $f^2$ | $\dfrac{(n-f)^3}{(n-2f)n}$ | $\dfrac{(n-f)^2}{n-2f}$ | $n^2$ |

Median based rules

# *Efficiency* + Byzantine Resilience ?

CAN THE SERVER STILL OBTAIN

$$\theta* \in \left( \theta \; ; \; \nabla Q(\theta) = 0 \right)$$

**+**

FRACTIONAL WORKLOAD PER NODE

$(f, \epsilon)-$Resilience

$$\epsilon \in \mathcal{O}\left( \sqrt{\left( \frac{1}{n-f} + \lambda^2(n-f) \right) \frac{\sigma^2}{T}} \right)$$
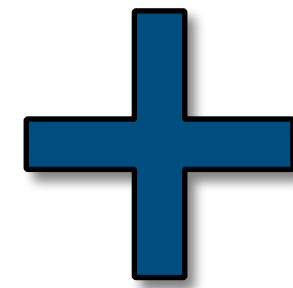
Additional workload

| Aggregation: | MDA ✅ | TM ❌ | Geometric Median ❌ | Krum* ❌ | Median ❌ |
|---|---|---|---|---|---|
| $\lambda^2(n-f)$ | $\dfrac{f^2}{n-f}$ | $f^2$ | $\dfrac{(n-f)^3}{(n-2f)n}$ | $\dfrac{(n-f)^2}{n-2f}$ | $n^2$ |

Averaging Component

Median based rules

# *Efficiency* + Byzantine Resilience ?

**CAN THE SERVER STILL OBTAIN**

$$\theta* \in \left( \theta \ ; \ \nabla Q(\theta) = 0 \right)$$

**+**

**FRACTIONAL WORKLOAD PER NODE**

$(f, \epsilon)-$Resilience

$$\epsilon \in \mathcal{O}\left( \sqrt{\left( \frac{1}{n-f} + \lambda^2(n-f) \right) \frac{\sigma^2}{T}} \right)$$
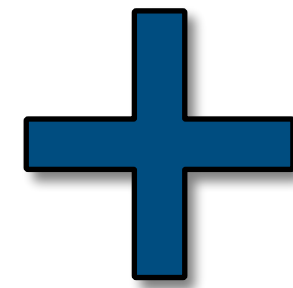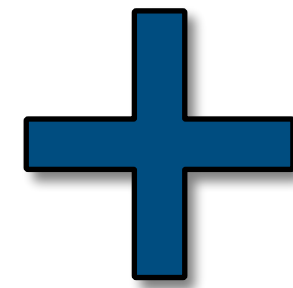
Additional workload

| Aggregation: | MDA ✅ | TM ❌ | Geometric Median ❌ | Krum* ❌ | Median ❌ |
|---|---|---|---|---|---|
| $\lambda^2(n-f)$ | $\dfrac{f^2}{n-f}$ | $f^2$ | $\dfrac{(n-f)^3}{(n-2f)n}$ | $\dfrac{(n-f)^2}{n-2f}$ | $n^2$ |

Averaging Component

Median based rules

* An extension of Krum, called multi-Krum, uses an averaging component to obtain guarantees comparable to MDA, but is computationally cheaper.

# !!Caution!! Momentum may NOT Help Always

# !!Caution!! Momentum may NOT Help Always

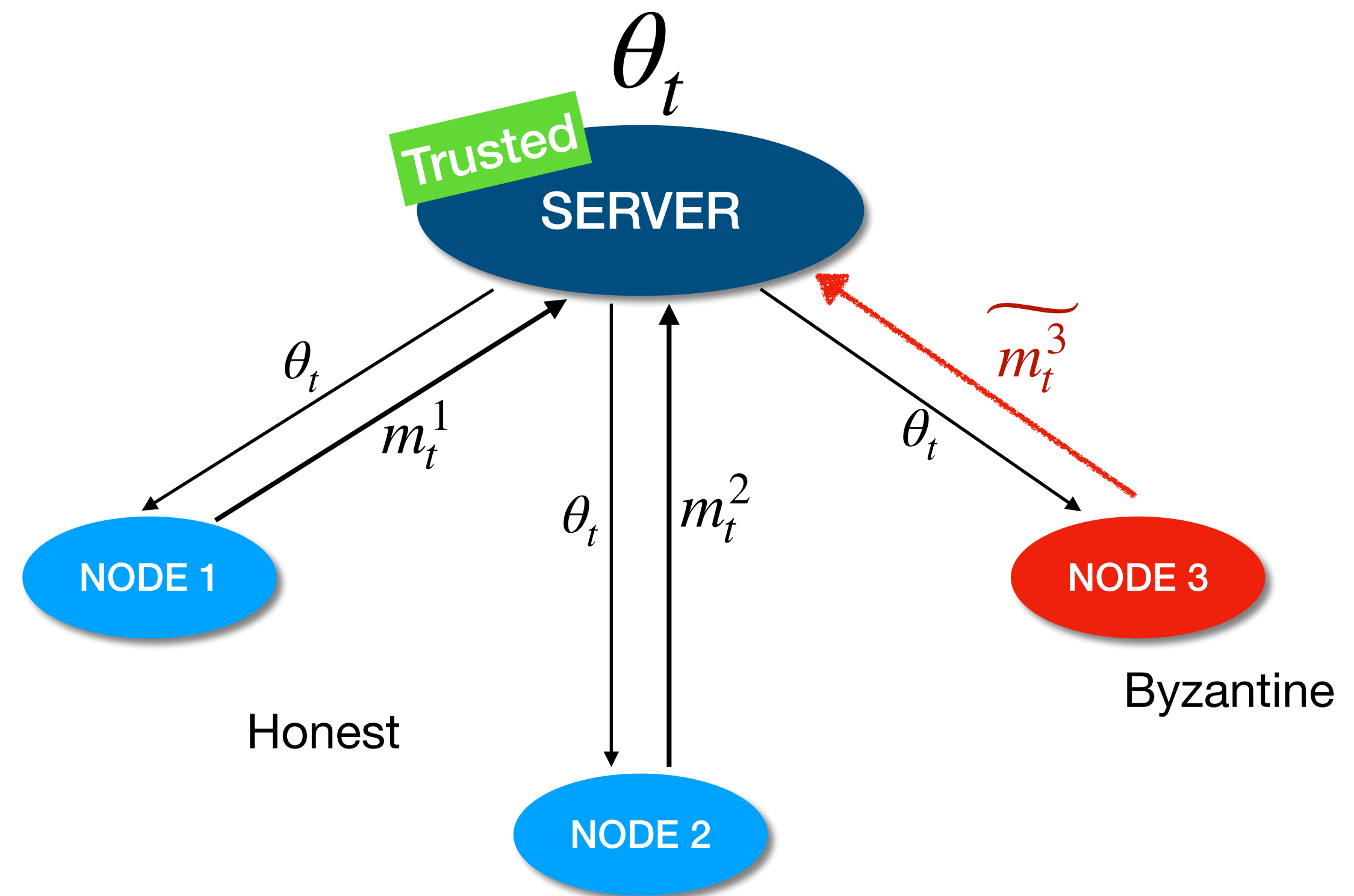When **data distributions** across nodes **are "very" different**

# !!Caution!! Momentum may NOT Help Always

When **data distributions** across nodes **are "very" different**

**For each honest node :**

# !!Caution!! Momentum may NOT Help Always

When **data distributions** across nodes **are "very" different**

**For each honest node :**

$$Q_i(\theta) := \mathbb{E}_{x \sim \mathscr{D}_i} q(\theta, x) \ ;$$

# !!Caution!! Momentum may NOT Help Always

When **data distributions** across nodes **are "very" different**

**For each honest node :**

$$Q_i(\theta) := \mathop{\mathbb{E}}_{x \sim \mathscr{D}_i} q\,(\theta, x) \;\; ; \;\; g_t^i = \nabla Q_i(\theta_t) + u_t^i$$

# !!Caution!! Momentum may NOT Help Always

When **data distributions** across nodes **are "very" different**

**For each honest node :**

$$Q_i(\theta) := \mathop{\mathbb{E}}_{x \sim \mathscr{D}_i} q(\theta, x) \;\; ; \;\; g_t^i = \nabla Q_i(\theta_t) + u_t^i$$

**The goal is to compute**



$\theta_t$

Trusted

SERVER

$\theta_t$

$m_t^1$

$\theta_t$

$m_t^2$

$\theta_t$

$\widetilde{m_t^3}$

NODE 1

NODE 2
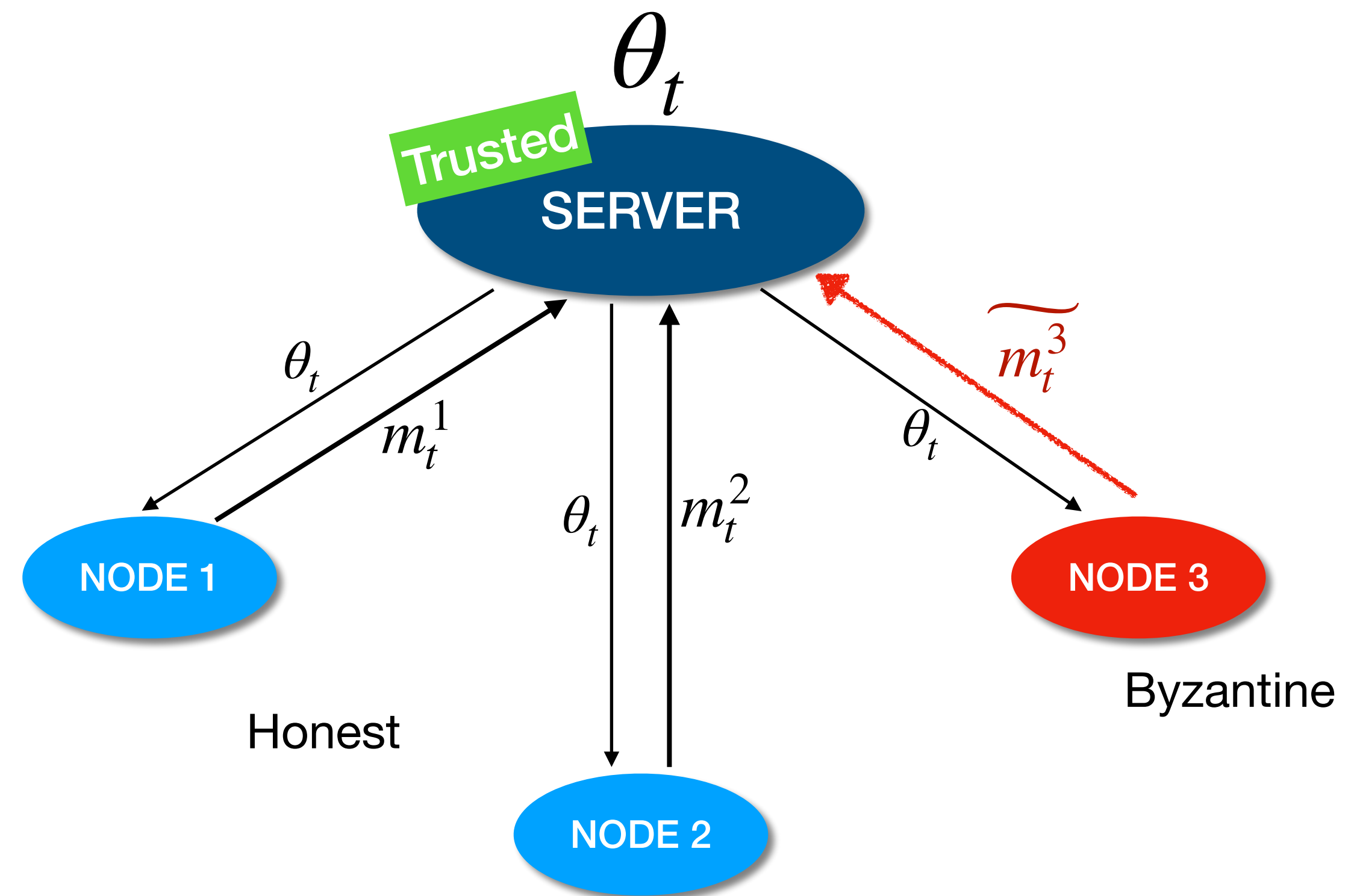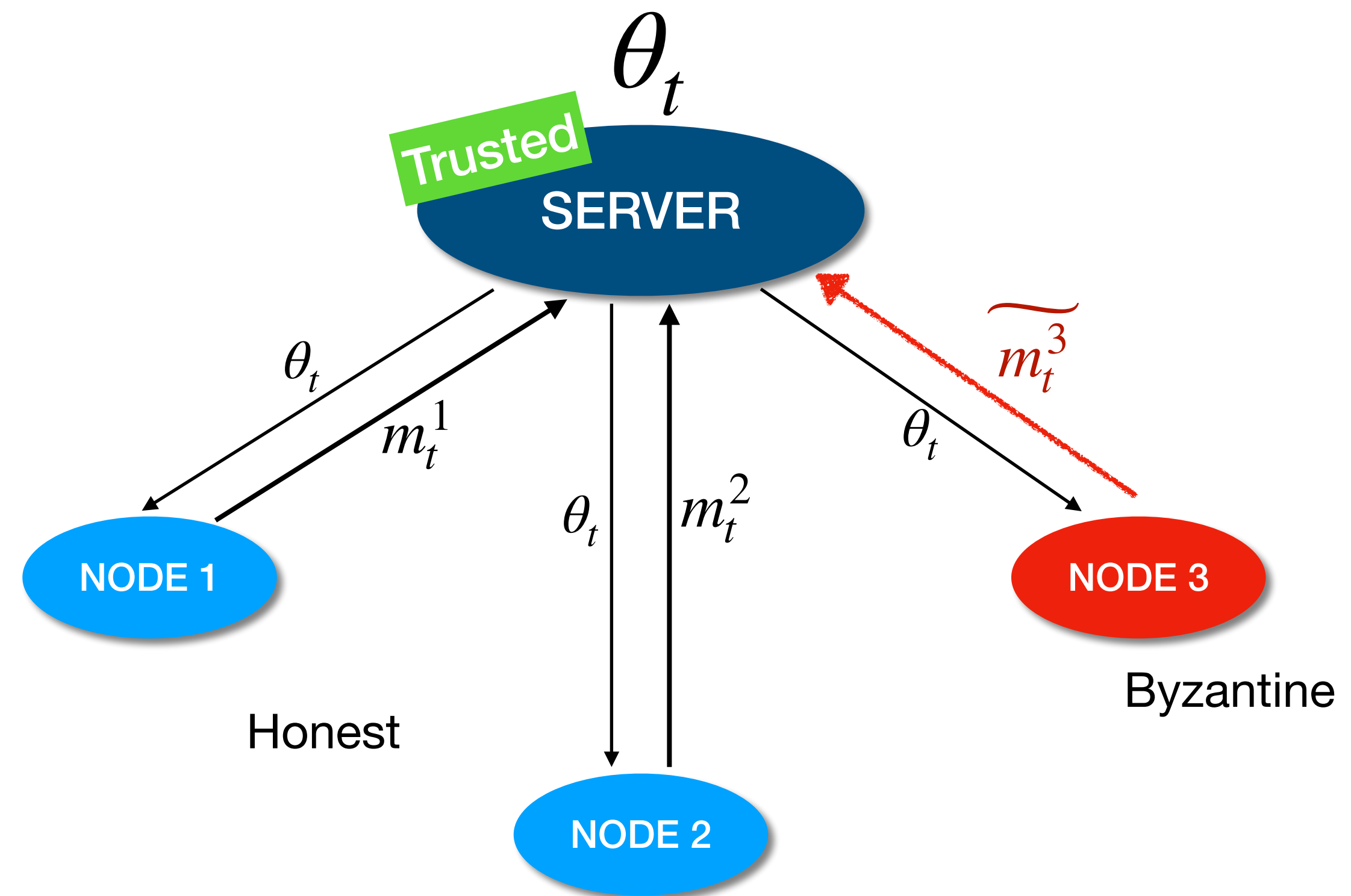
NODE 3

Honest

Byzantine

# !!Caution!! Momentum may NOT Help Always

When **data distributions** across nodes **are "very" different**

**For each honest node :**

$$Q_i(\theta) := \mathop{\mathbb{E}}_{x \sim \mathscr{D}_i} q\,(\theta, x) \;\; ; \;\; g_t^i = \nabla Q_i(\theta_t) + u_t^i$$

**The goal is to compute** $\;\; \theta* \in \left( \theta \; ; \; \dfrac{1}{H} \sum_{i \in H} \nabla Q_i(\theta) = 0 \right)$
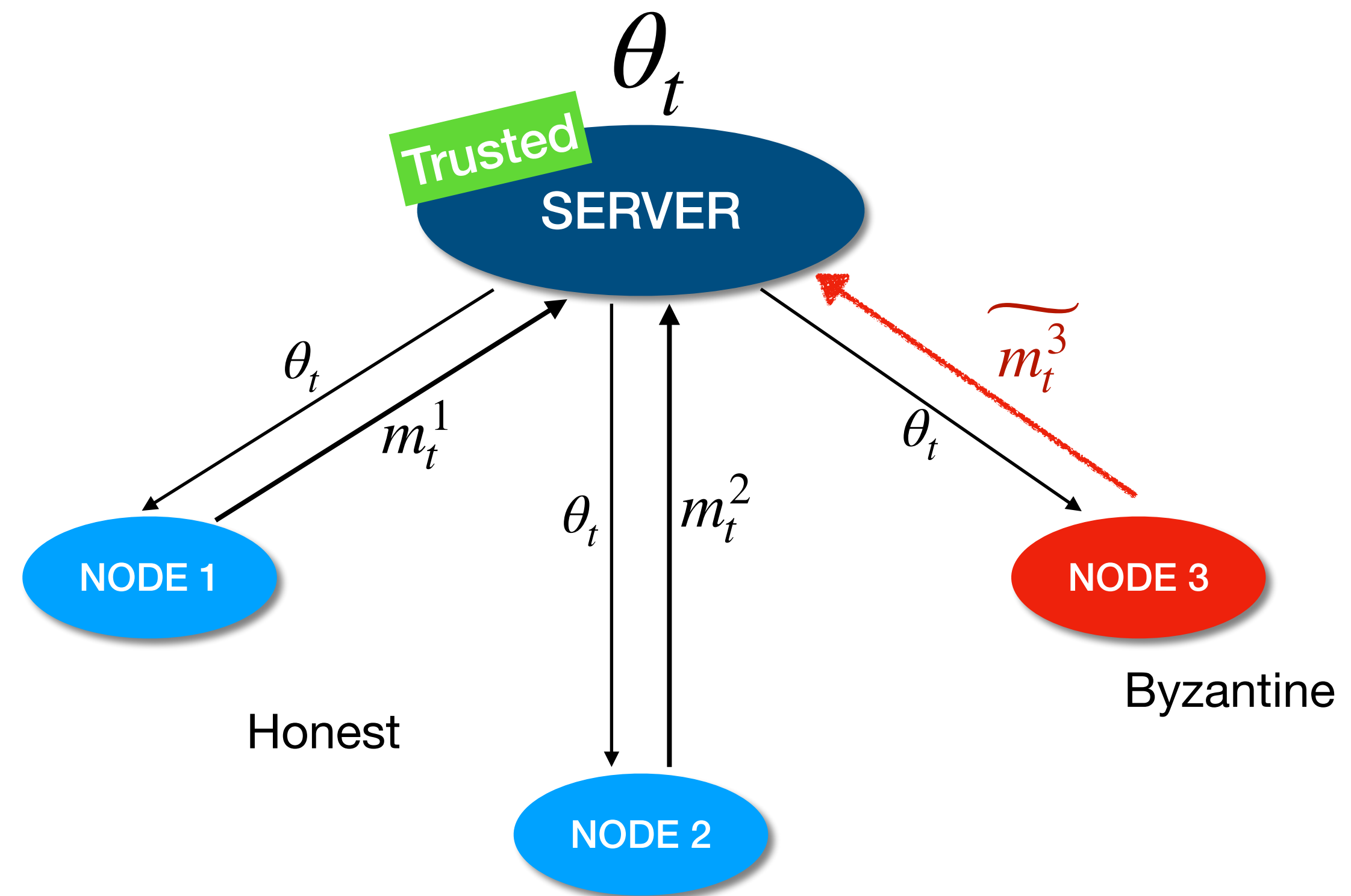
# !!Caution!! Momentum may NOT Help Always

When **data distributions** across nodes **are "very" different**

**For each honest node :**

$$Q_i(\theta) := \mathop{\mathbb{E}}_{x \sim \mathscr{D}_i} q\,(\theta, x) \;\; ; \;\; g_t^i = \nabla Q_i(\theta_t) + u_t^i$$

**The goal is to compute** $\theta* \in \left( \theta \; ; \; \dfrac{1}{H} \sum_{i \in H} \nabla Q_i(\theta) = 0 \right)$
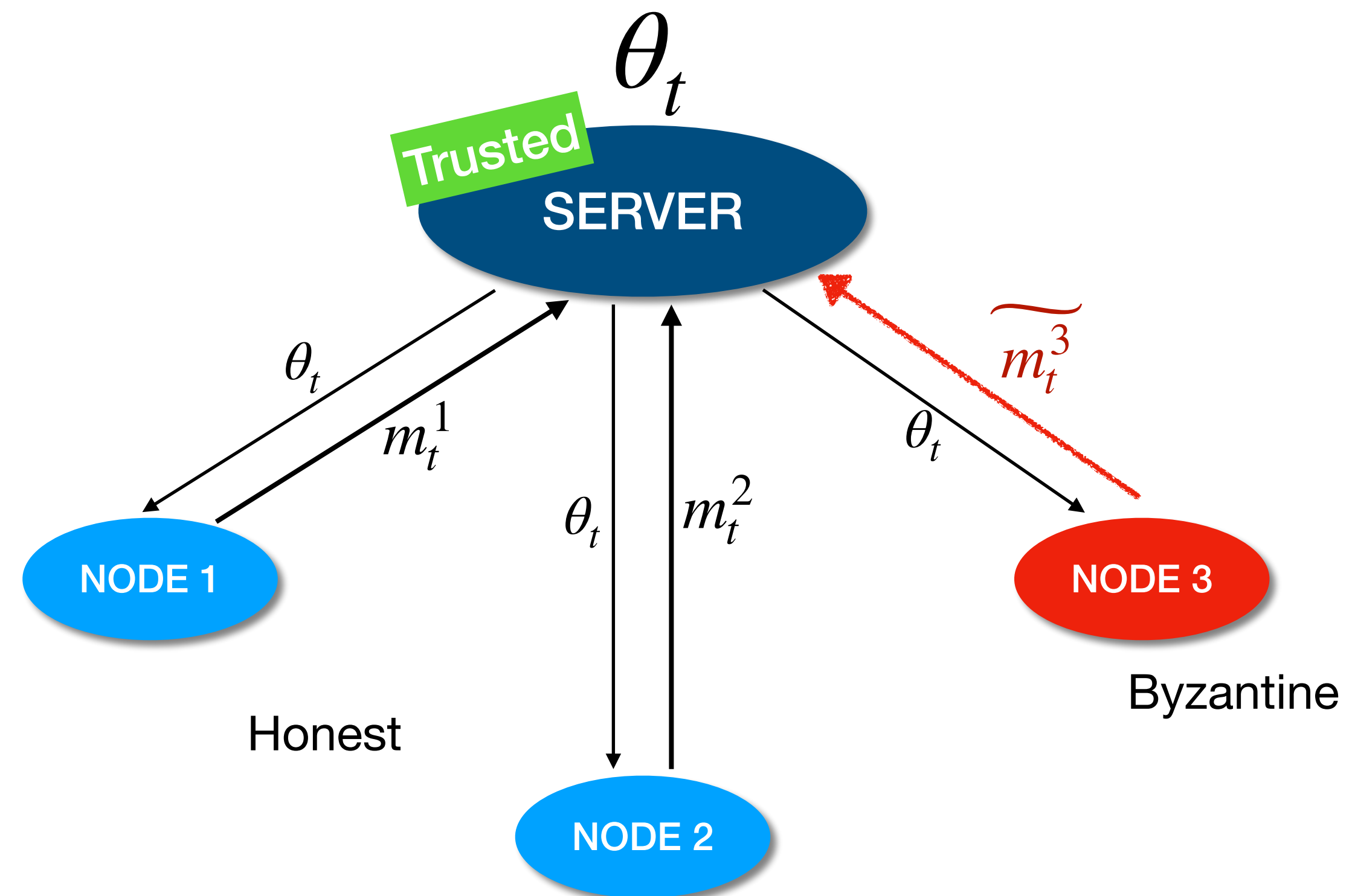
Momentum accumulates the heterogeneity

# !!Caution!! Momentum may NOT Help Always

When **data distributions** across nodes **are "very" different**

**For each honest node :**

$$Q_i(\theta) := \mathop{\mathbb{E}}_{x \sim \mathscr{D}_i} q(\theta, x) \ ; \quad g_t^i = \nabla Q_i(\theta_t) + u_t^i$$

**The goal is to compute** $\quad \theta^* \in \left( \theta \ ; \ \dfrac{1}{H} \sum_{i \in H} \nabla Q_i(\theta) = 0 \right)$

Momentum accumulates the heterogeneity

$$\mathbb{E}\left[ \|m_t^i - m_t^j\|^2 \right] \le c_1(1 - \beta)\sigma^2 + c_2 \sum_{k=1}^{t} \|\nabla Q_i(\theta_k) - \nabla Q_j(\theta_k)\|^2$$

$\theta_t$

Trusted

SERVER

$\theta_t$

$m_t^1$

$\widetilde{m_t^3}$

$\theta_t$

NODE 1

$\theta_t$

$m_t^2$

NODE 3

Byzantine

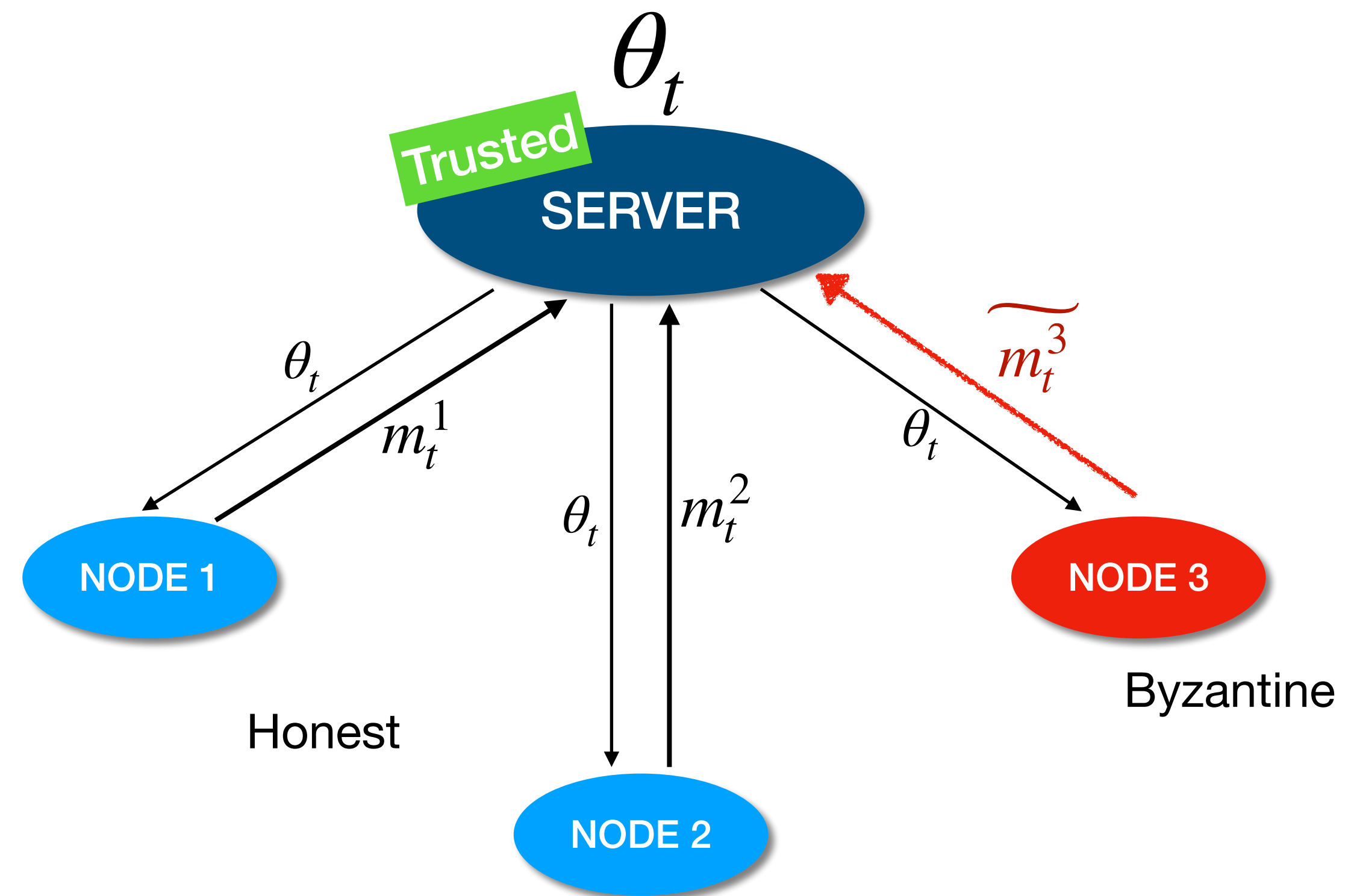Honest

NODE 2

# !!Caution!! Momentum may NOT Help Always

When **data distributions** across nodes **are "very" different**

**For each honest node :**

$$Q_i(\theta) := \mathop{\mathbb{E}}_{x \sim \mathscr{D}_i} q\left(\theta, x\right) \; ; \quad g_t^i = \nabla Q_i(\theta_t) + u_t^i$$

**The goal is to compute** $\theta^* \in \left( \theta \; ; \; \dfrac{1}{H} \displaystyle\sum_{i \in H} \nabla Q_i(\theta) = 0 \right)$

Momentum accumulates the heterogeneity

$$\mathbb{E}\left[ \|m_t^i - m_t^j\|^2 \right] \le c_1(1 - \beta)\sigma^2 + c_2 \sum_{k=1}^{t} \|\nabla Q_i(\theta_k) - \nabla Q_j(\theta_k)\|^2$$

$\theta_t$

Trusted

SERVER

$\theta_t$

$m_t^1$

$\widetilde{m_t^3}$

$\theta_t$ | $m_t^2$

$\theta_t$

NODE 1

NODE 3

Honest

Byzantine

NODE 2

# !!Caution!! Momentum may NOT Help Always

When **data distributions** across nodes **are "very" different**
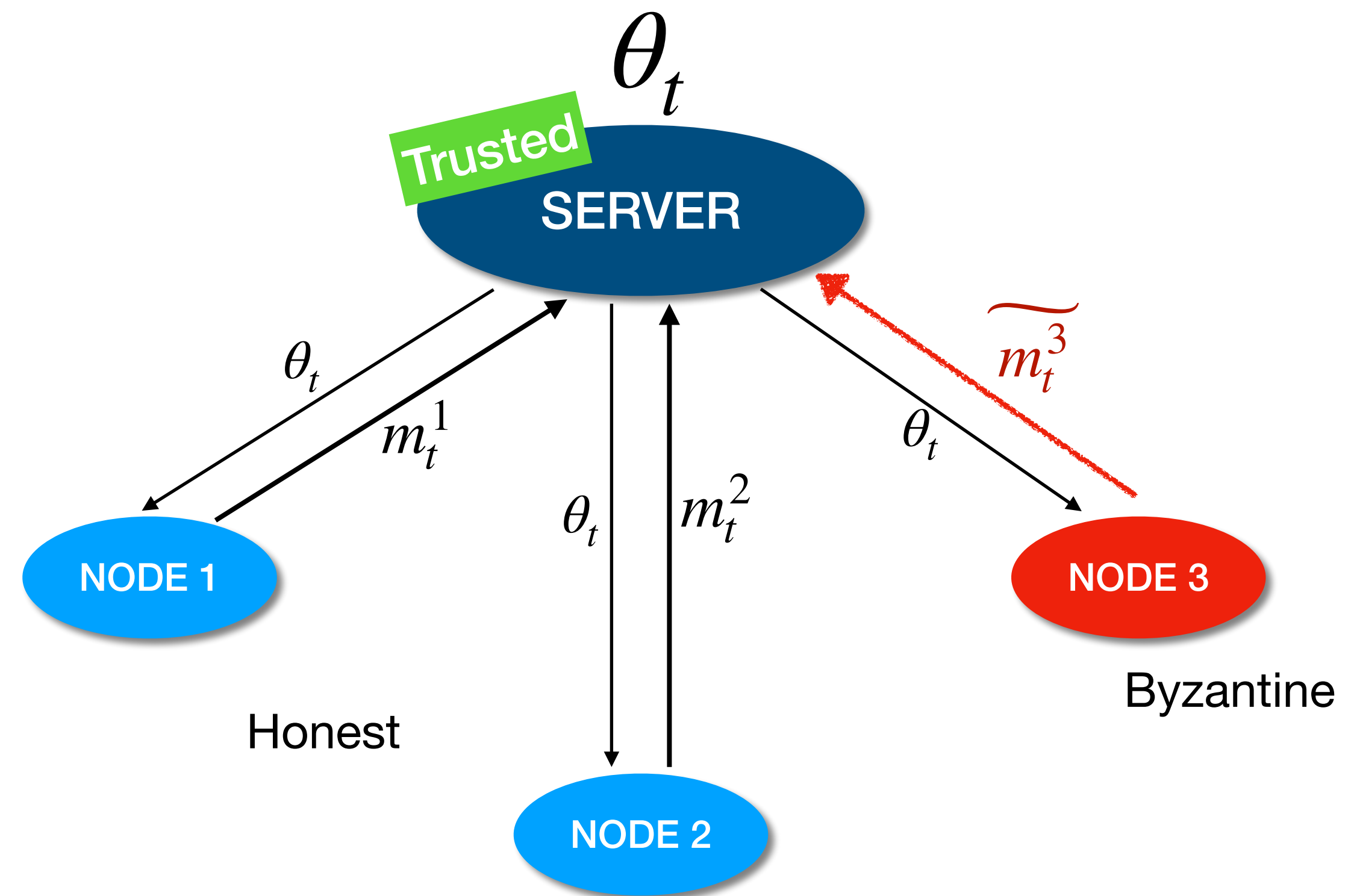
**For each honest node :**

$$Q_i(\theta) := \mathop{\mathbb{E}}_{x \sim \mathcal{D}_i} q\,(\theta, x) \;\; ; \;\; g_t^i = \nabla Q_i(\theta_t) + u_t^i$$

**The goal is to compute** $\theta* \in \left( \theta \; ; \; \dfrac{1}{H} \sum_{i \in H} \nabla Q_i(\theta) = 0 \right)$

Momentum accumulates the heterogeneity

$$\mathbb{E}\left[ \|m_t^i - m_t^j\|^2 \right] \le c_1(1 - \beta)\sigma^2 + c_2 \sum_{k=1}^{t} \|\nabla Q_i(\theta_k) - \nabla Q_j(\theta_k)\|^2$$

On the other hand,

# !!Caution!! Momentum may NOT Help Always

When **data distributions** across nodes **are "very" different**
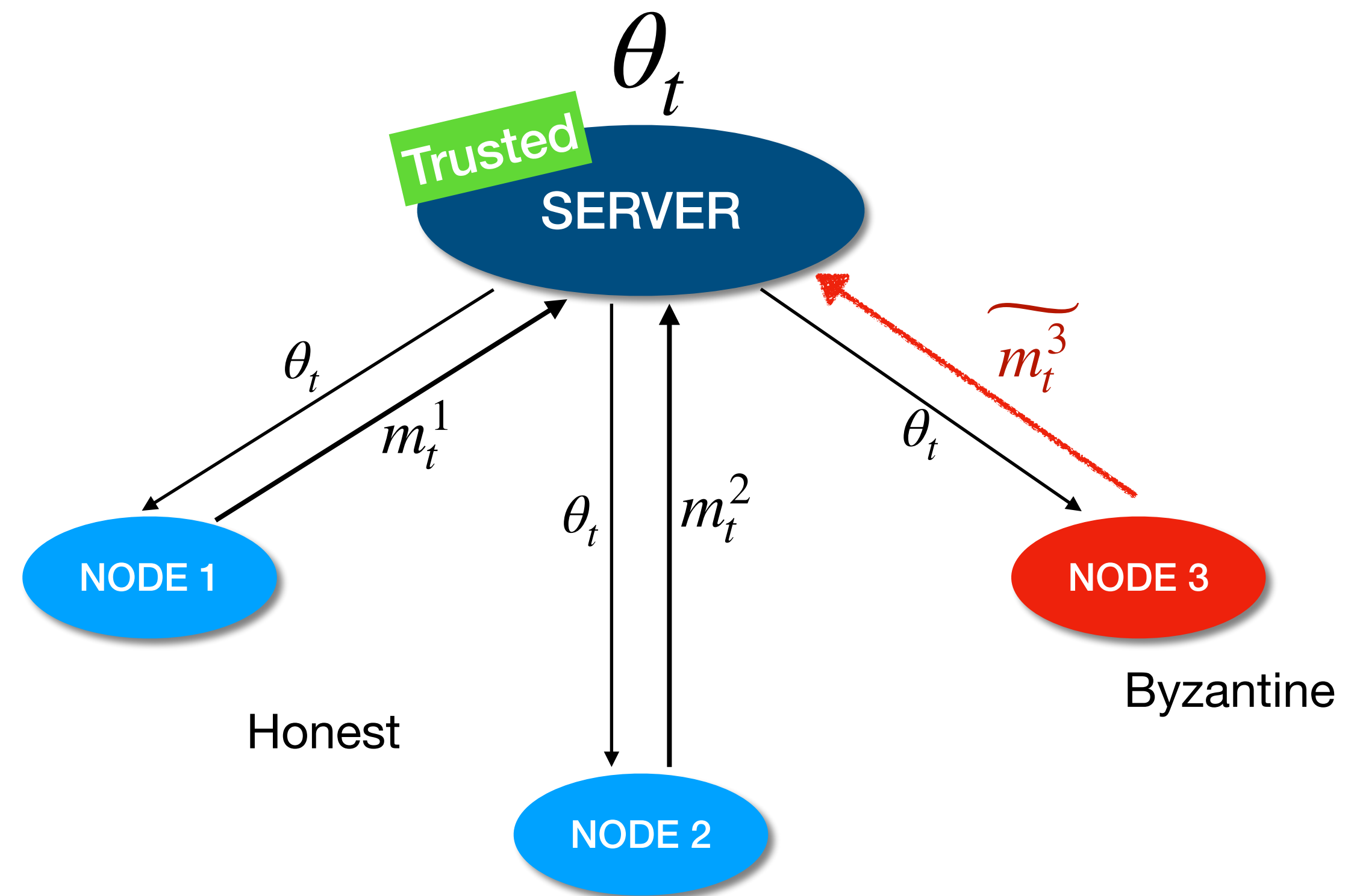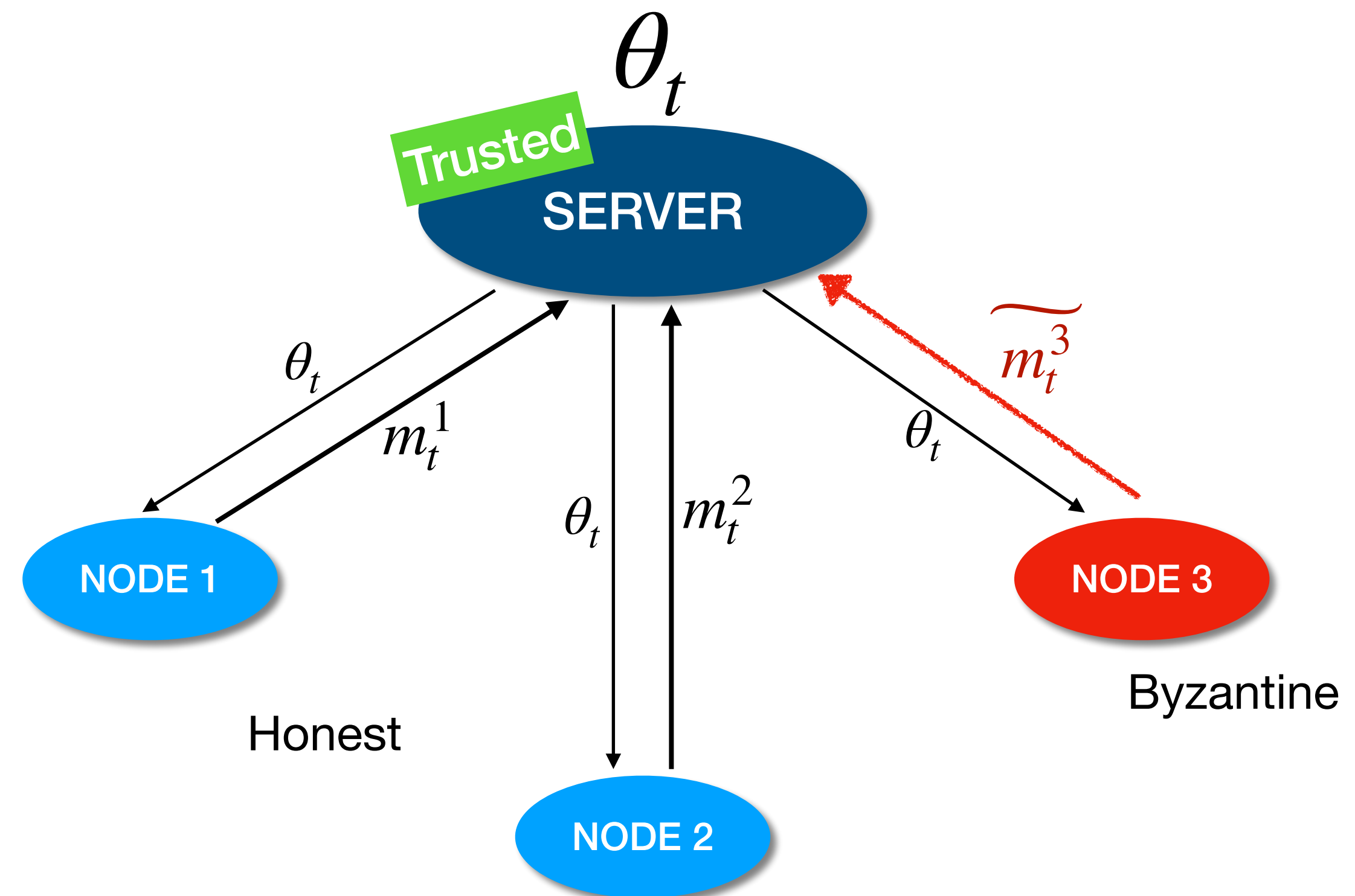
**For each honest node :**

$$Q_i(\theta) := \mathop{\mathbb{E}}_{x \sim \mathscr{D}_i} q(\theta, x) \; ; \quad g_t^i = \nabla Q_i(\theta_t) + u_t^i$$

**The goal is to compute** $\theta* \in \left( \theta \; ; \; \dfrac{1}{H} \sum_{i \in H} \nabla Q_i(\theta) = 0 \right)$

Momentum accumulates the heterogeneity

$$\mathbb{E}\left[ \|m_t^i - m_t^j\|^2 \right] \leq c_1(1-\beta)\sigma^2 + c_2 \sum_{k=1}^{t} \|\nabla Q_i(\theta_k) - \nabla Q_j(\theta_k)\|^2$$

On the other hand, $\mathbb{E}\left[ \|g_t^i - g_t^j\|^2 \right] \leq c_1 \sigma^2 + c_2 \|\nabla Q_i(\theta_k) - \nabla Q_j(\theta_k)\|^2$

# What's Next?!

# What's Next?!

Study the impact of momentum with **heterogeneity**.

# What's Next?!

Study the impact of momentum with **heterogeneity**.

Explore **other variance reduction strategies**, e.g., MVR*.

# What's Next?!

Study the impact of momentum with **heterogeneity**.

Explore **other variance reduction strategies**, e.g., MVR*.

* Momentum-based Variance Reduction (MVR) has optimal convergence rate for a first-order stochastic algorithm *(Cutkosky and Orabona, 2019)*

# What's Next?!

Study the impact of momentum with **heterogeneity**.

Explore **other variance reduction strategies**, e.g., MVR*.

**Improve the robustness condition** to guarantee *efficiency*

* Momentum-based Variance Reduction (MVR) has optimal convergence rate for a first-order stochastic algorithm *(Cutkosky and Orabona, 2019)*

# What's Next?!

Study the impact of momentum with **heterogeneity**.

Explore **other variance reduction strategies**, e.g., MVR*.

**Improve the robustness condition** to guarantee *efficiency*

Does use of local momentum improve **privacy**?

* Momentum-based Variance Reduction (MVR) has optimal convergence rate for a first-order stochastic algorithm *(Cutkosky and Orabona, 2019)*
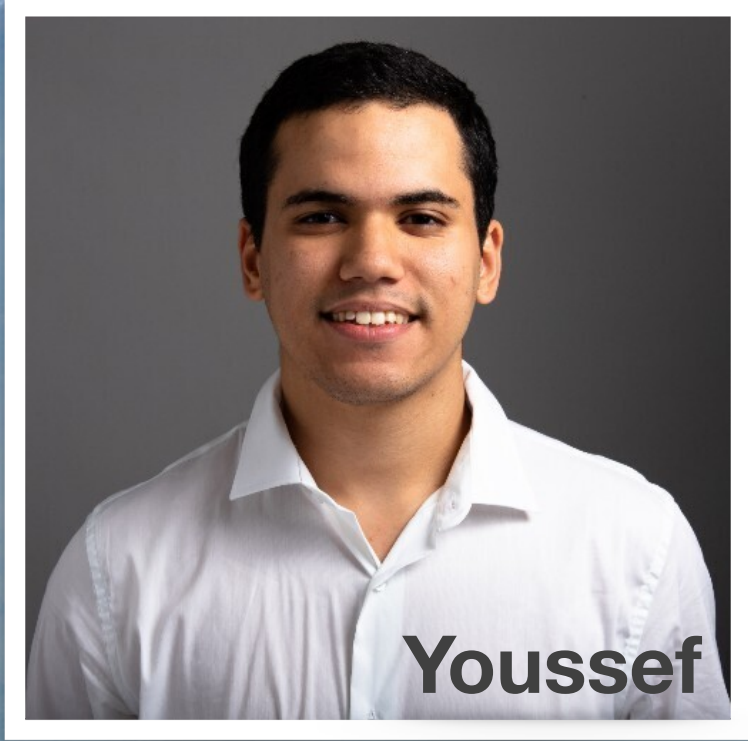
Thanks to ...

John

Sadegh

Youssef

Rafael

Rachid

Lê Nguyên

# Readings

Farhadkhani, Sadegh; Rachid Guerraoui; Nirupam Gupta; Rafael Pinot and John Stephan.
**Byzantine Machine Learning Made Easy by Resilient Averaging of Momentums.** *ICML 2022.*

El-Mhamdi El-Mahdi; Rachid Guerraoui and Sébastien Rouault.
**Distributed momentum for byzantine-resilient learning.** *ICLR 2021.*

Karimireddy Sai Praneeth; Lie He and Martin Jaggi.
**Learning from history for byzantine robust optimization.** *ICML 2021.*

**Momentum for
Byzantine resilience**

Alistarh Dan; Zeyuan Allen-Zhu and Jerry Li.
**Byzantine stochastic gradient descent.** *NeurIPS 2018.*

El-Mhamdi El-Mahdi; Guerraoui Rachid, and Sébastien Rouault.
**The hidden vulnerability of distributed learning in Byzantium.** *ICML 2018.*

Chen Yudong; Lili Su and Jiaming Xu.
**Distributed statistical machine learning in adversarial settings: Byzantine gradient descent.**
*Proceedings of the ACM on Measurement and Analysis of Computing Systems 2017.*

**Notable prior work
On Byzantine resilience**

# Thank You