

# Distributed Machine Learning

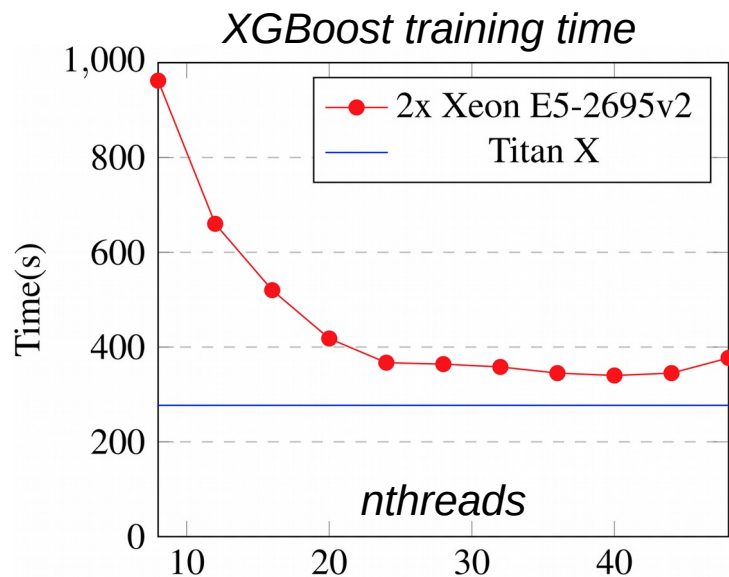
## Episode I, YSDA'21

Yandex  
Research

LAMBDA 



# Зачем это всё?

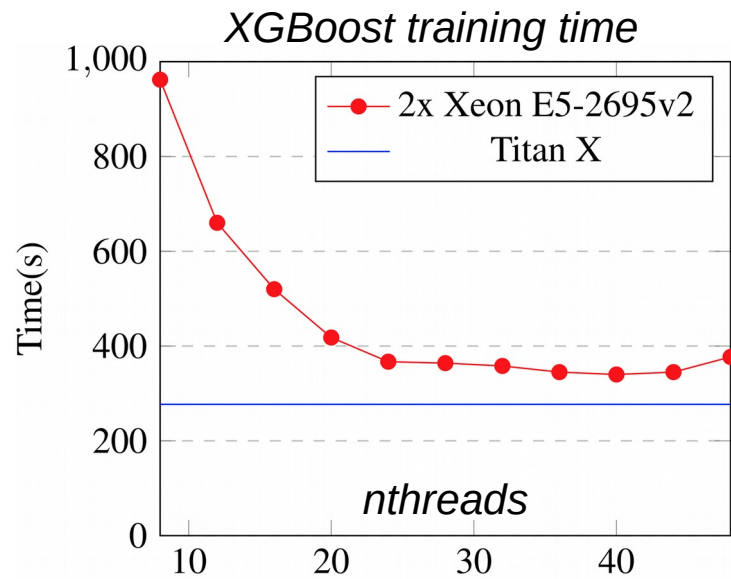


**BERT-Large Training Times on GPUs**

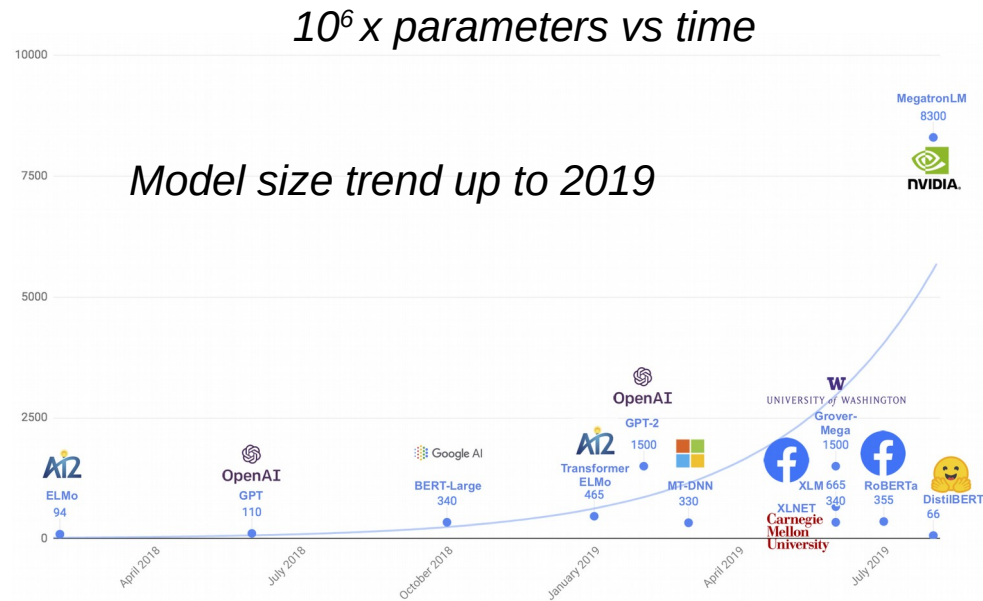
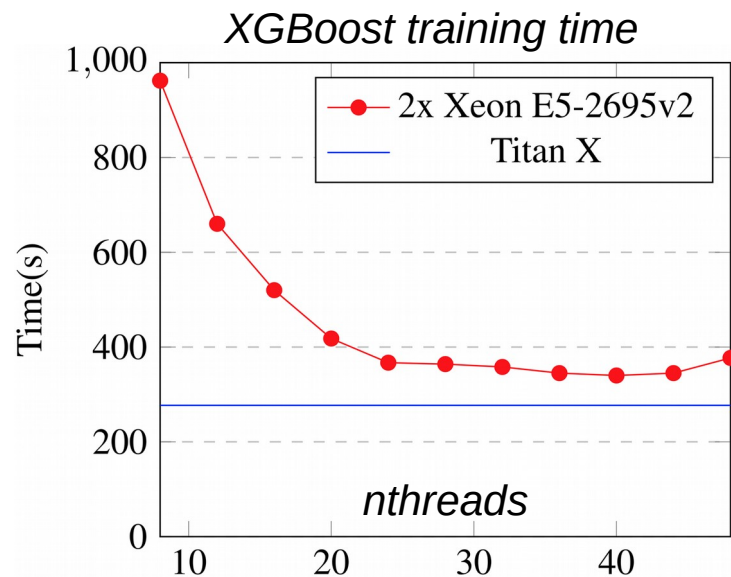
Time	System	Number of Nodes	Number of V100 GPUs
47 min	DGX SuperPOD	92 x DGX-2H	1,472
67 min	DGX SuperPOD	64 x DGX-2H	1,024
236 min	DGX SuperPOD	16 x DGX-2H	256

(single V100 – **over 2 weeks**)

# Зачем это всё?



# Зачем это всё?

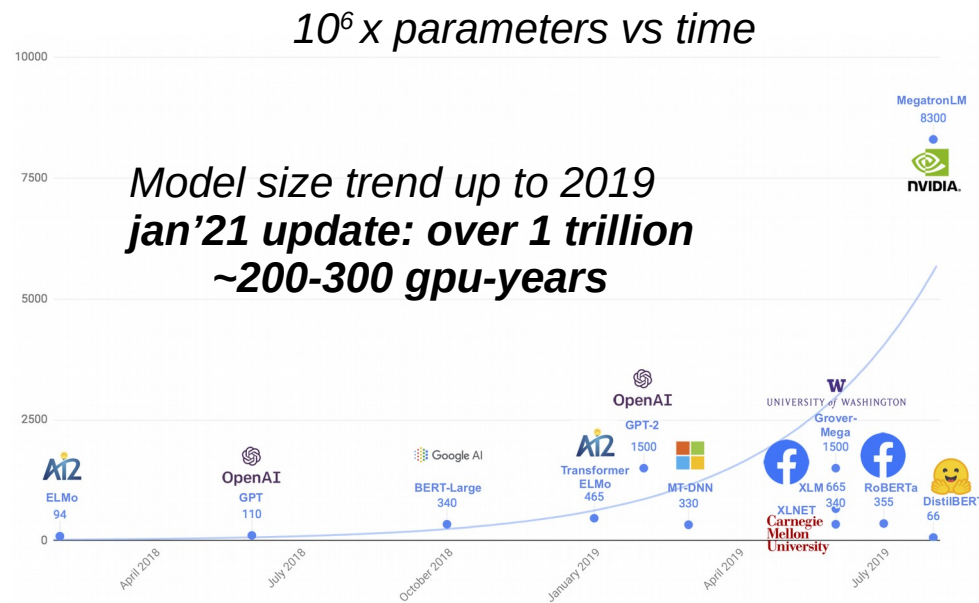
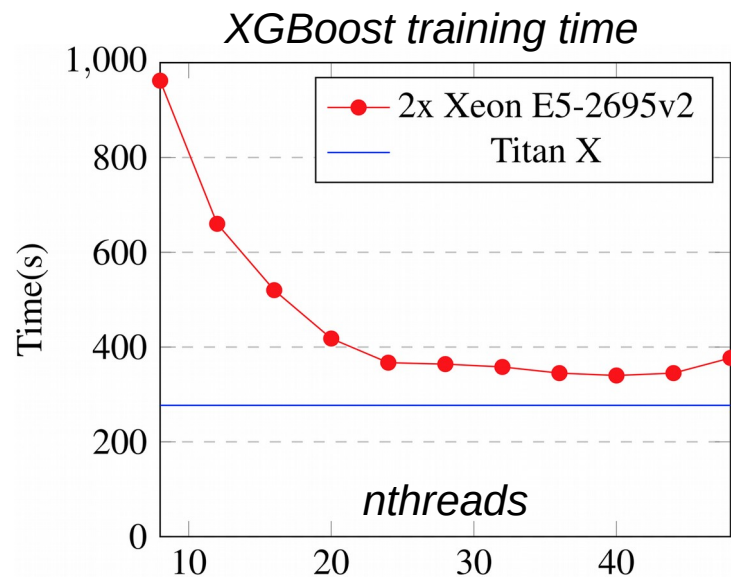


**BERT-Large Training Times on GPUs**

Time	System	Number of Nodes	Number of V100 GPUs
47 min	DGX SuperPOD	92 x DGX-2H	1,472
67 min	DGX SuperPOD	64 x DGX-2H	1,024
236 min	DGX SuperPOD	16 x DGX-2H	256

(single V100 – over 2 weeks)

# Зачем это всё?



**BERT-Large Training Times on GPUs**

Time	System	Number of Nodes	Number of V100 GPUs
47 min	DGX SuperPOD	92 x DGX-2H	1,472
67 min	DGX SuperPOD	64 x DGX-2H	1,024
236 min	DGX SuperPOD	16 x DGX-2H	256

(single V100 – over 2 weeks)

# Зачем мы тут?

Заставить много железяк вместе учить одну модель



# Зачем мы тут?

Заставить много железяк вместе учить одну модель



понять общие подходы  
закодировать своими руками  
на python / pytorch



# TL;DR this course

*with sample problems*

## 1) Distributed machine learning

*Embeddings or log.regression with tons of training data*



# TL;DR this course

*with sample problems*

## 1) Distributed machine learning

*Embeddings or log.regression with tons of training data*

## 2) Data-parallel deep learning

*Train BERT-base on wikipedia in 20 minutes or less*

# TL;DR this course

*with sample problems*

## 1) Distributed machine learning

*Embeddings or log.regression with tons of training data*

## 2) Data-parallel deep learning

*Train BERT-base on wikipedia in 20 minutes or less*

## 3) Model-parallel deep learning

*Train BERT-xxxxxxxxxxxxlarge in less than a lifetime*

# TL;DR this course

*with sample problems*

## 1) Distributed machine learning

*Embeddings or log.regression with tons of training data*

## 2) Data-parallel deep learning

*Train BERT-base on wikipedia in 20 minutes or less*

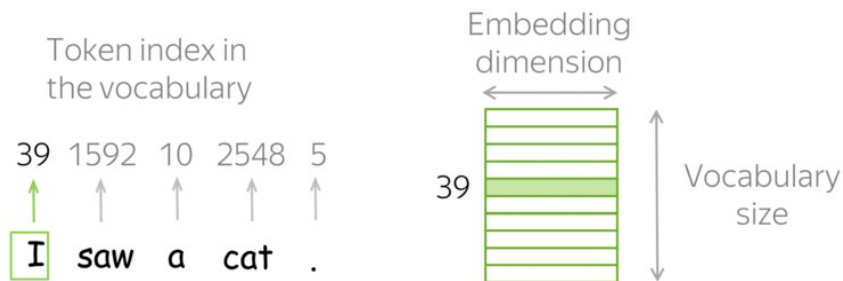
## 3) Model-parallel deep learning

*Train BERT-xxxxxxxxxxxxlarge in less than a lifetime*

## 4) Decentralized deep learning

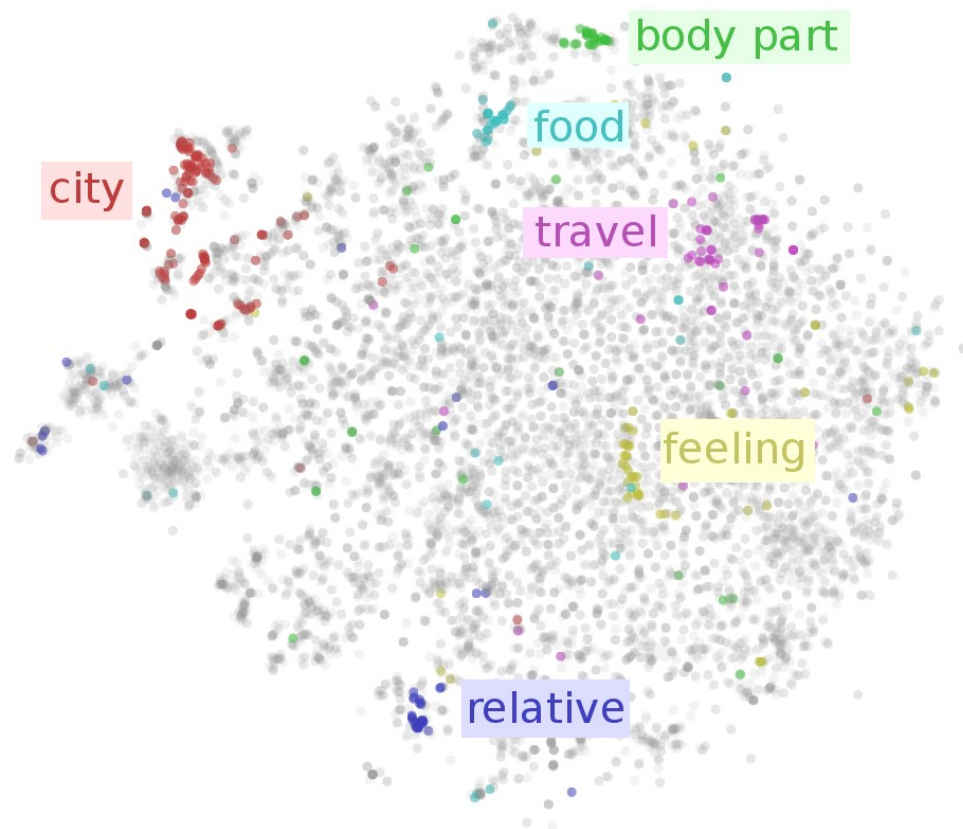
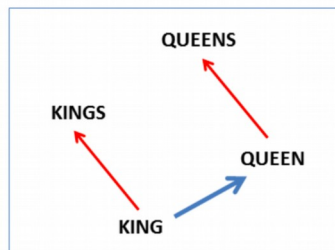
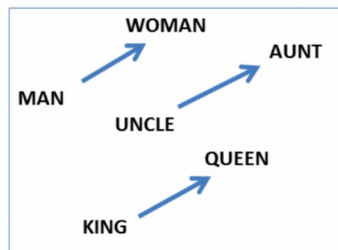
*Train ~something~ with a million smart teapots*

# Problem of the day: word embeddings

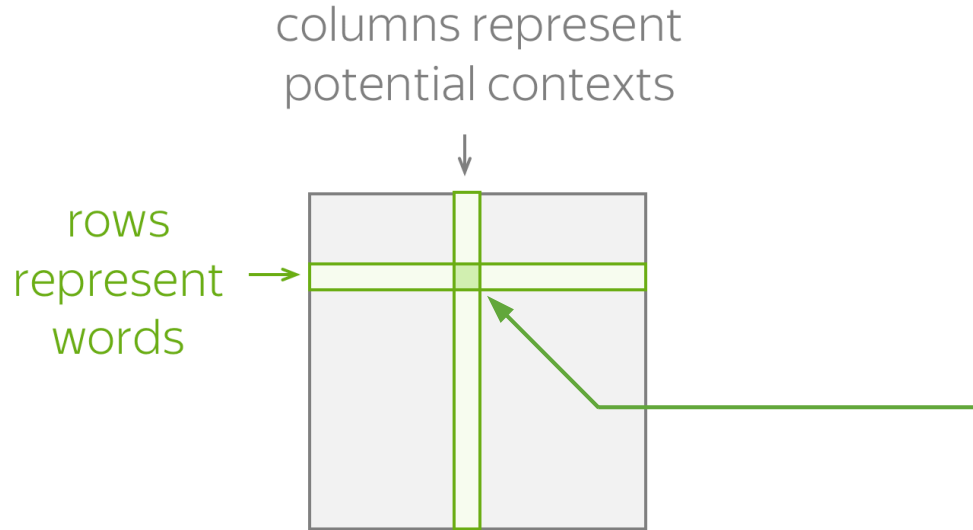


semantic:  $v(\text{king}) - v(\text{man}) + v(\text{woman}) \approx v(\text{queen})$

syntactic:  $v(\text{kings}) - v(\text{king}) + v(\text{queen}) \approx v(\text{queens})$



# Co-occurrence matrix

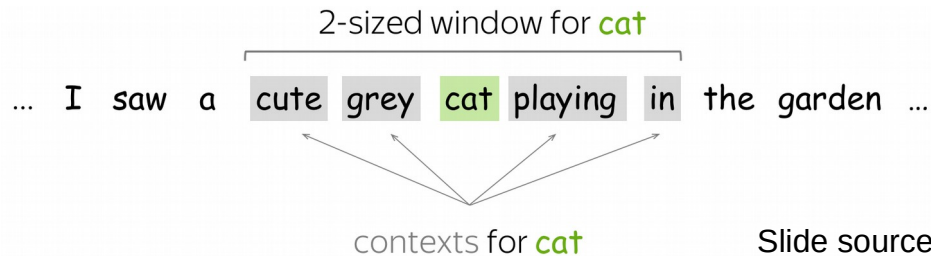


Context:

- surrounding words in a L-sized window

Matrix element:


- $N(w, c)$  – number of times word  $w$  appears in context  $c$



**Note:** in our case,  $N$  is symmetric!

# GloVe

context vector      word vector      bias terms (also learned)


$$L = \sum_{i \neq j} w(N(i, j)) \cdot (\langle \vec{v}_i, \vec{v}_j \rangle + b_i + b_j - \log N(i, j))^2$$

# GloVe

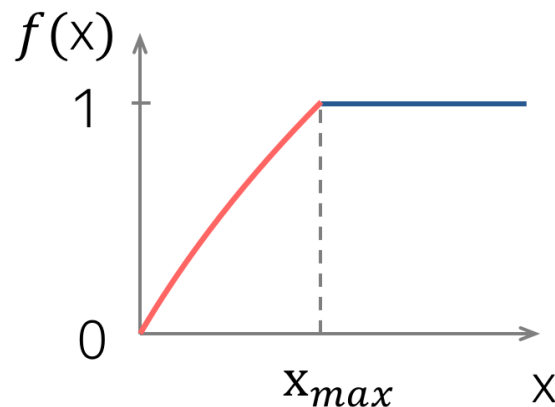
context vector      word vector      bias terms (also learned)

$$L = \sum_{i \neq j} w(N(i, j)) \cdot (\langle \vec{v}_i, \vec{v}_j \rangle + b_i + b_j - \log N(i, j))^2$$

↓

Weighting function to:


- penalize rare events
- not to over-weight frequent events



$$\begin{cases} (x/x_{max})^\alpha & \text{if } x < x_{max}, \\ 1 & \text{otherwise.} \end{cases}$$

$$\alpha = 0.75, x_{max} = 100$$

# GloVe

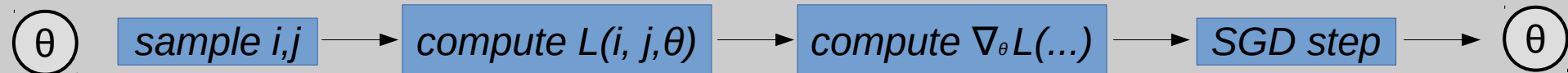

$$L = \sum_{i \neq j} w(N(i, j)) \cdot (\langle \vec{v}_i, \vec{v}_j \rangle + b_i + b_j - \log N(i, j))^2$$

**Learn more:** [lena-voita.github.io/nlp\\_course/word\\_embeddings.html](https://lena-voita.github.io/nlp_course/word_embeddings.html)

*So how do we train 'em?*



# Training Step

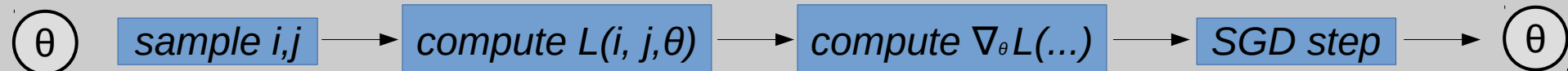


current  
params  
 $\theta: \{v, b\}$

$$L = \sum_{i \neq j} w(N(i, j)) \cdot (\langle \vec{v}_i, \vec{v}_j \rangle + b_i + b_j - \log N(i, j))^2$$

updated  
params

# Training Step



current  
params  
 $\theta: \{v, b\}$

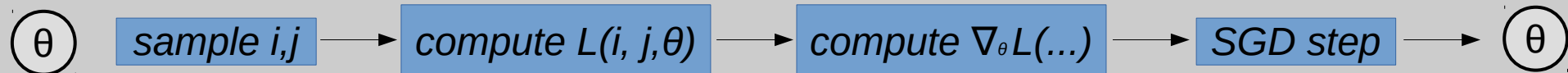
$$L = \sum_{i \neq j} w(N(i, j)) \cdot (\langle \vec{v}_i, \vec{v}_j \rangle + b_i + b_j - \log N(i, j))^2$$

updated  
params

**Trainable parameters:**

	V matrix	b
a		
cat		
saw		

# Training Step

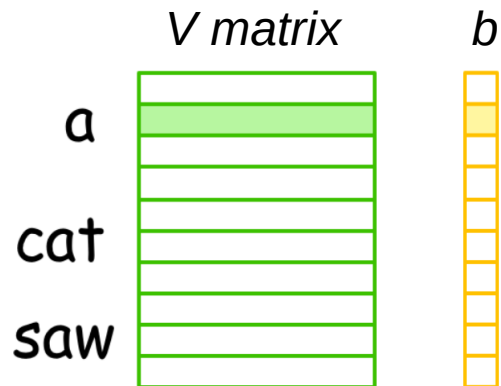


current  
params  
 $\theta: \{v, b\}$

$$L = \sum_{i \neq j} w(N(i, j)) \cdot (\langle \vec{v}_i, \vec{v}_j \rangle + b_i + b_j - \log N(i, j))^2$$

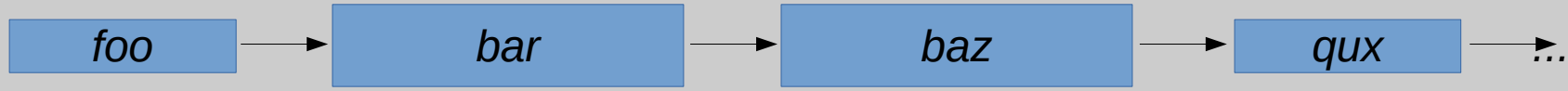
updated  
params

**Trainable parameters:**



**How do we go faster  
with 8 CPU cores?**

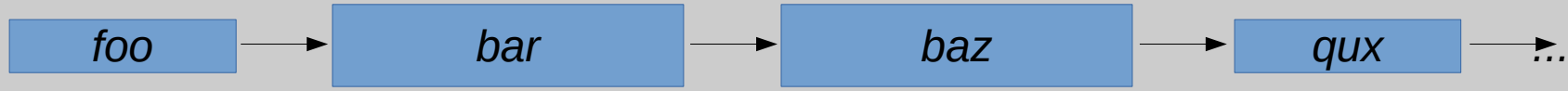
# Rules: Process



## Process:

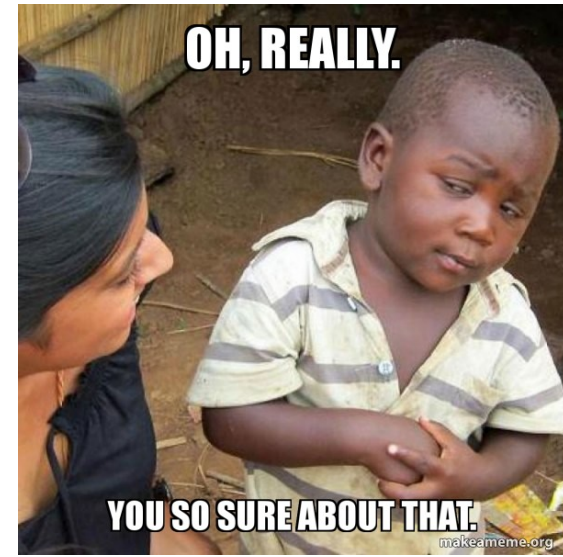
- Runs some code
- Has some memory
- No one else can access your memory

# Rules: Process

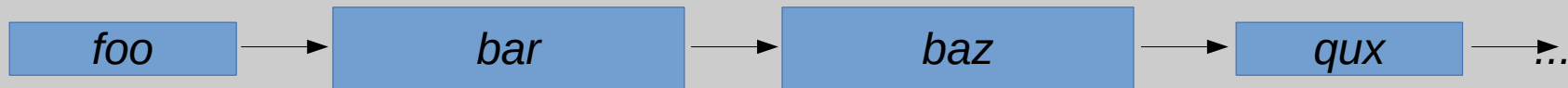


## Process:

- Runs some code
- Has some memory
- No one else can access your memory



# Rules: Process



## Process:

- Runs some code
- Has some memory
- No one else can access your memory\*

\* – not if you use shared memory

# Rules: Process



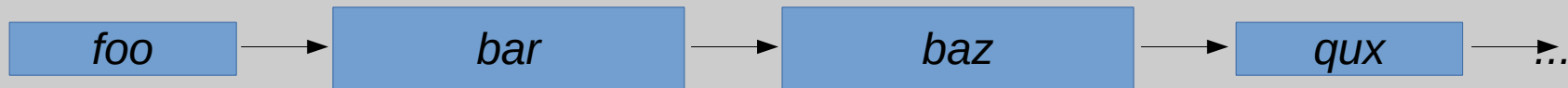
## Process:

- Runs some code
- Has some memory
- No one else can access your memory<sup>\*†</sup>

<sup>\*</sup> – not if you use shared memory

<sup>†</sup> – superuser can still do that (os-dependent)

# Rules: Process



## Process:

- Runs some code
- Has some memory
- No one else can access your memory<sup>\*†‡</sup>

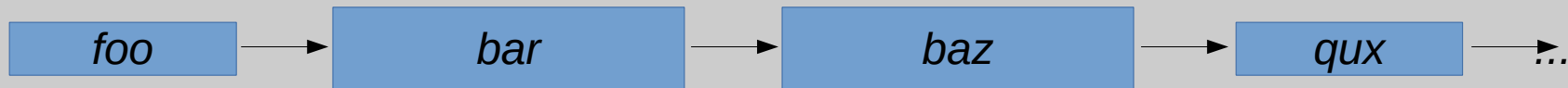
\* – not if you use shared memory

† – superuser can still do that (os-dependent)

‡ – attacker can do that through spectre/meltdown/etc



# Rules: Process



## Process:

- Runs some code
- Has some memory
- No one else **should** access your memory<sup>\*†‡</sup>

<sup>\*†‡</sup> – not relevant for this course

# Rules: Process



## Process:

- Runs some code
- Has some memory
- No one else **should** access your memory<sup>\*†‡</sup>

<sup>\*†‡</sup> – not relevant for this course

**Q:** How do we make processes work together?

# Rules: Channel / Pipe

Process A:



Process B:

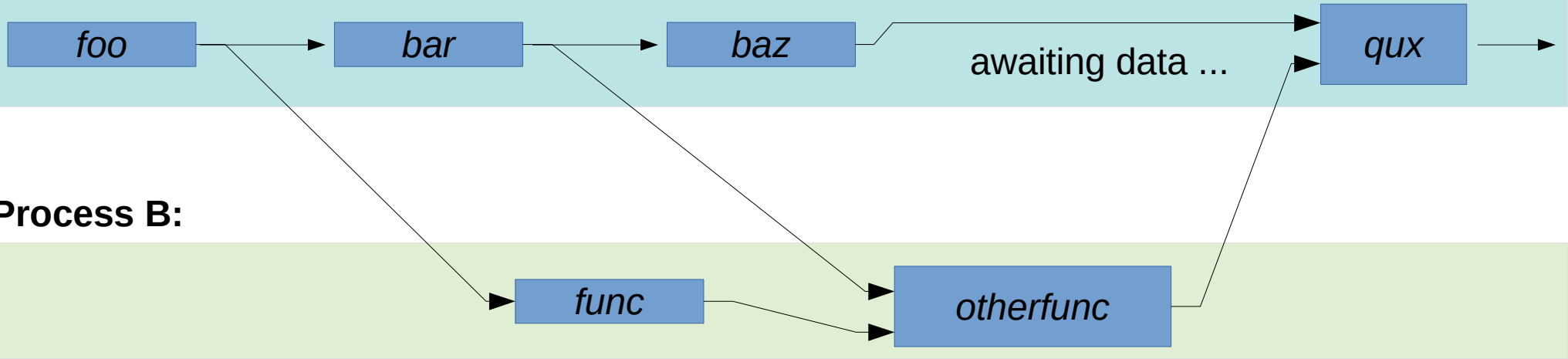


## Channel (pipe):

- Communication in  $O(\text{message size})$
- Asynchronous read/write

# MP Rules

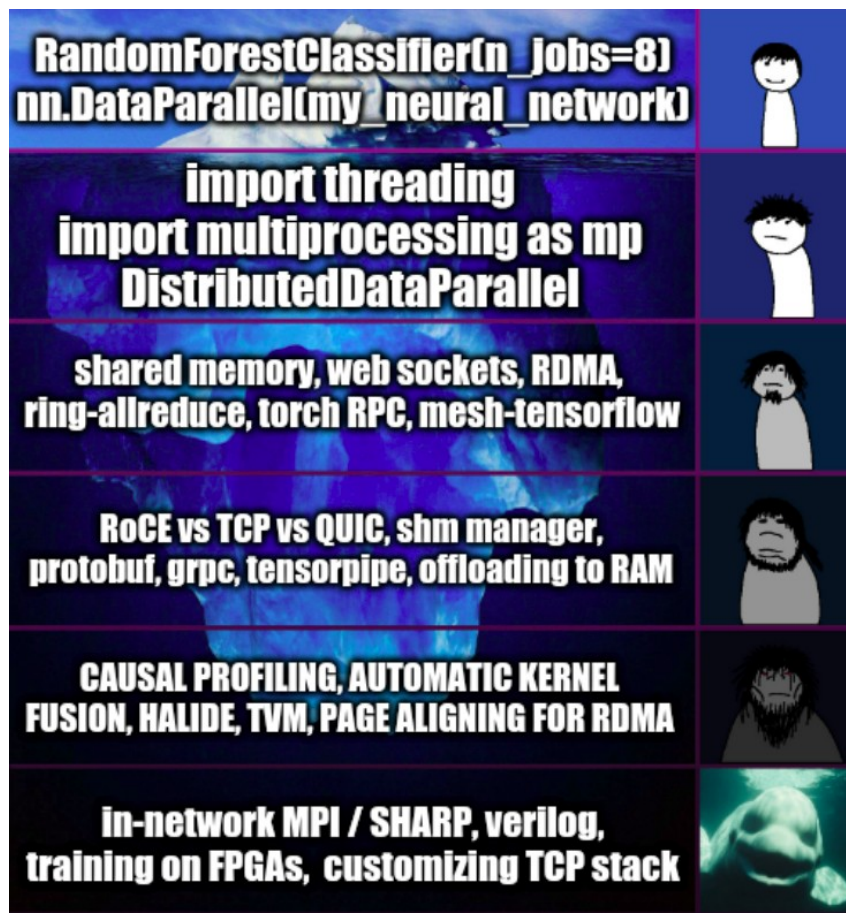
**Process A:** not waiting on write



## Channel (pipe):

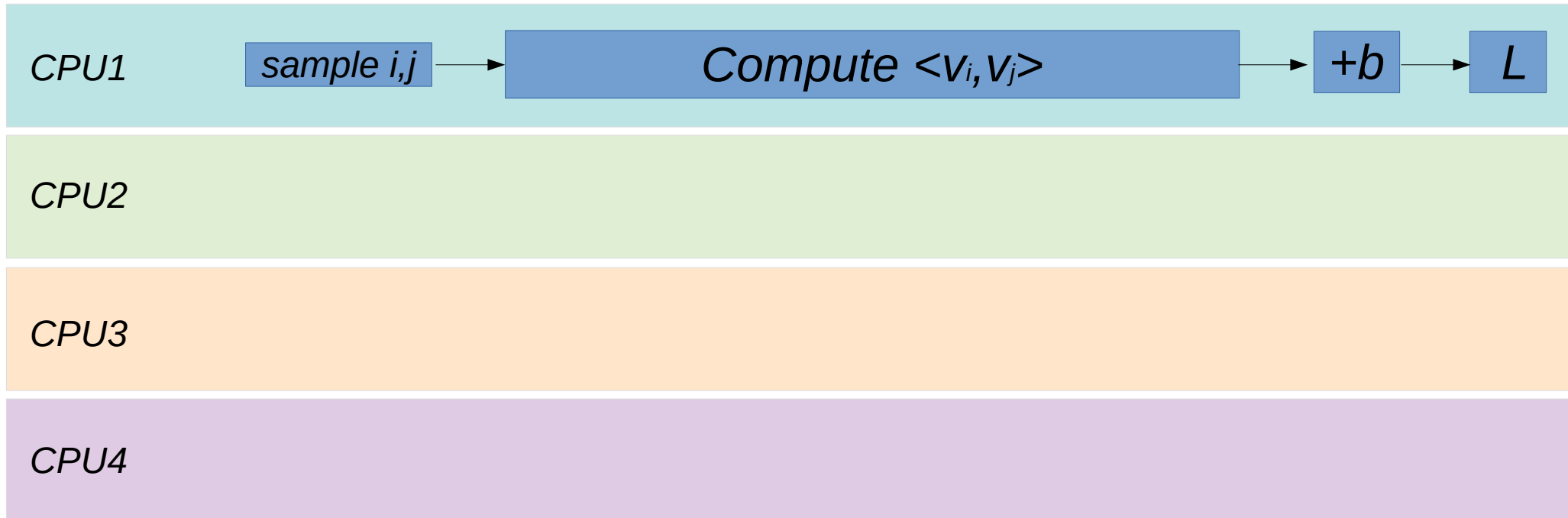
- Communication in  $O(\text{message size})$
- **Asynchronous** read/write

# Details are (not) important



# Operation parallelism

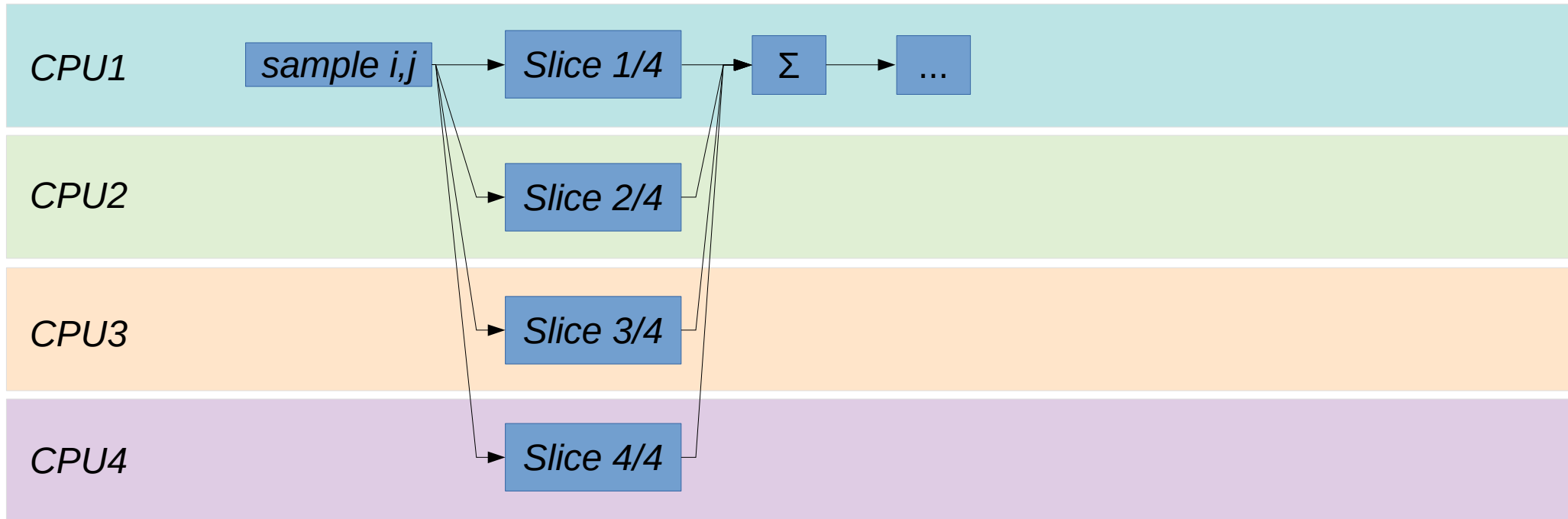
*run algorithm in parallel without changing the math*



# Operation parallelism

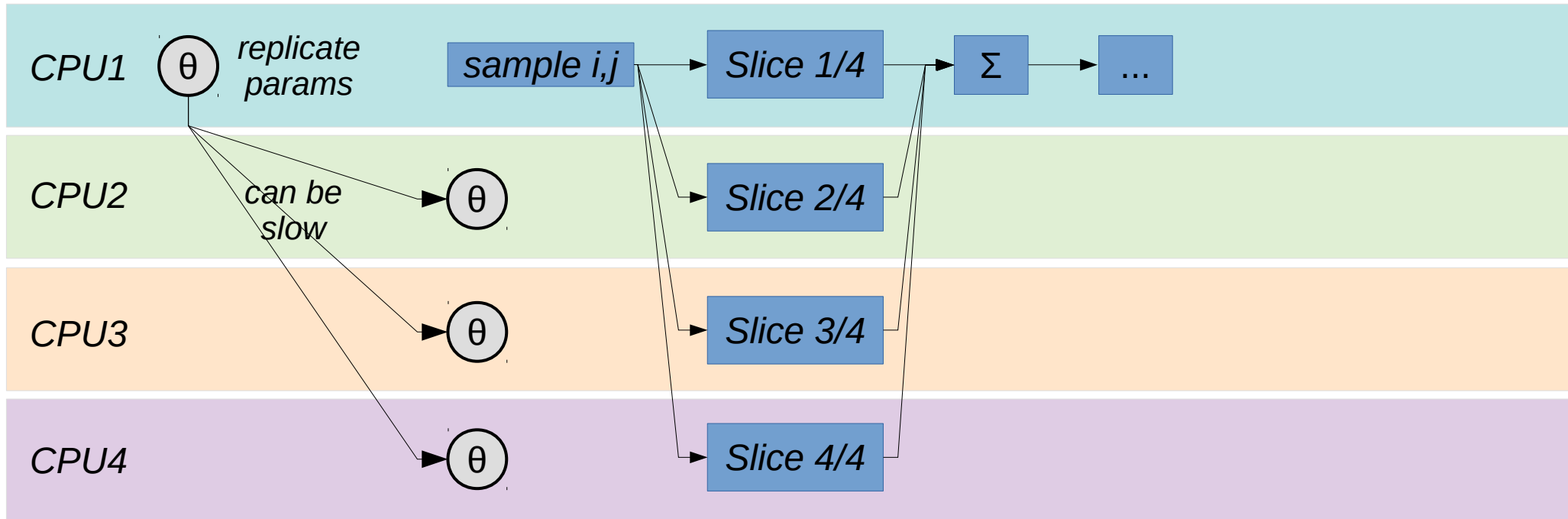
*run algorithm in parallel without changing the math*

$\langle V_i, V_j \rangle$



# Operation parallelism

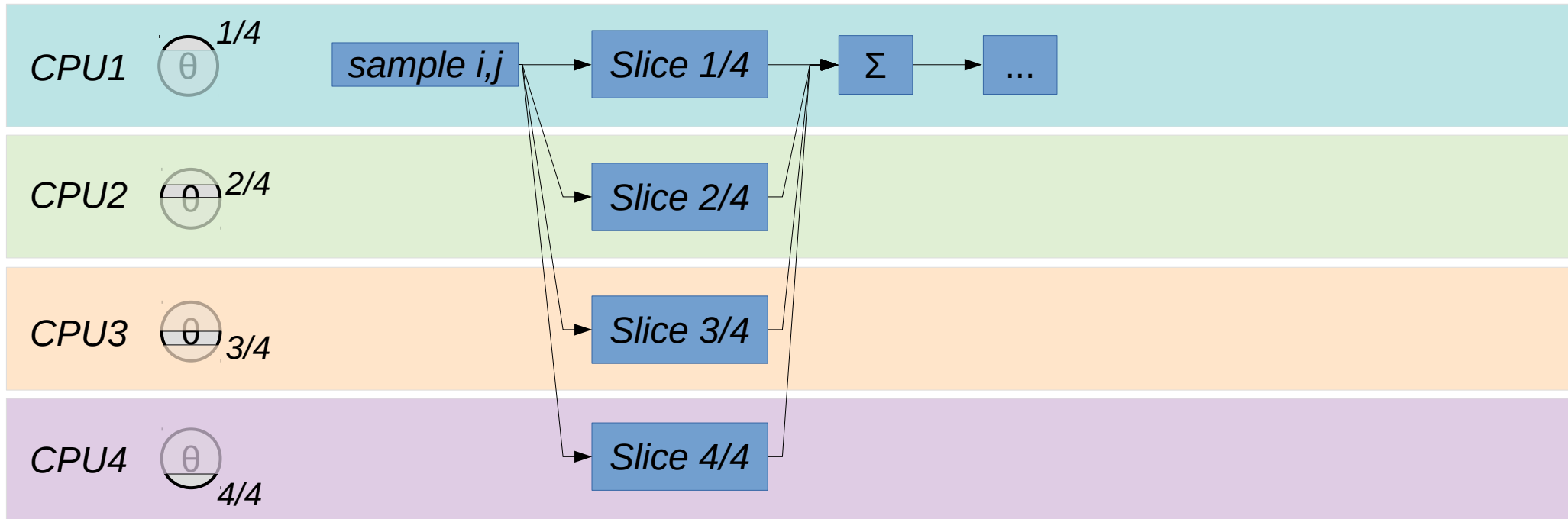
*How to organize parameters across processes?*





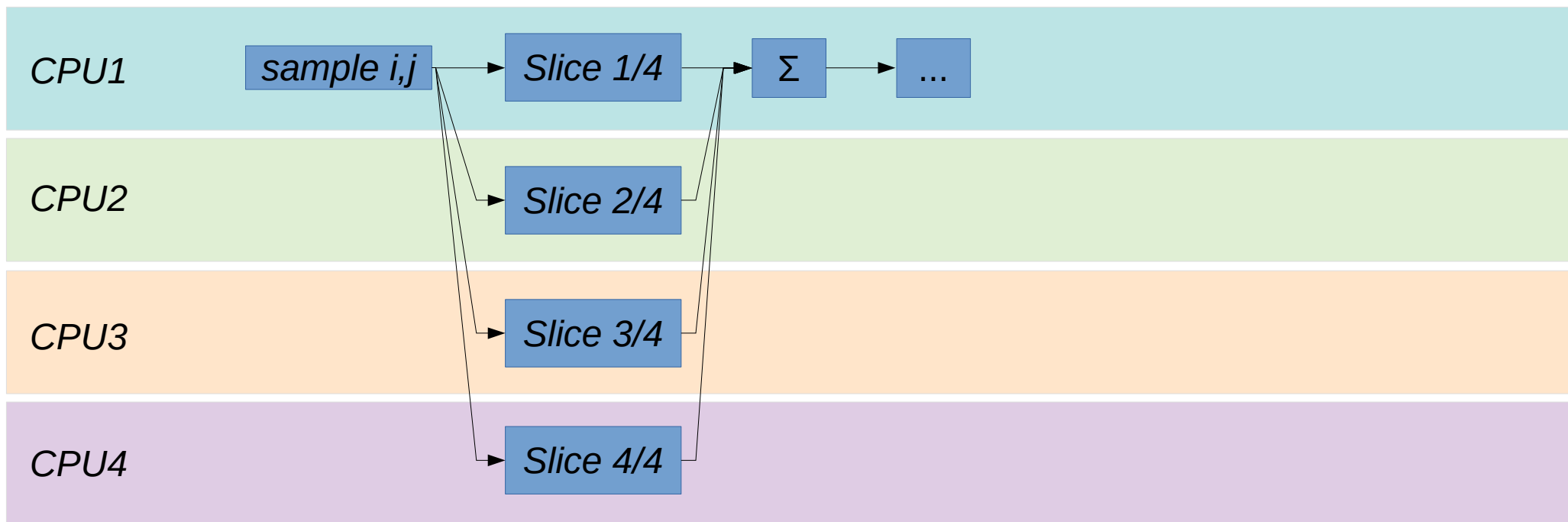
# Operation parallelism

*Each process holds one shard of parameters  
(no transfer required)*



# Operation parallelism

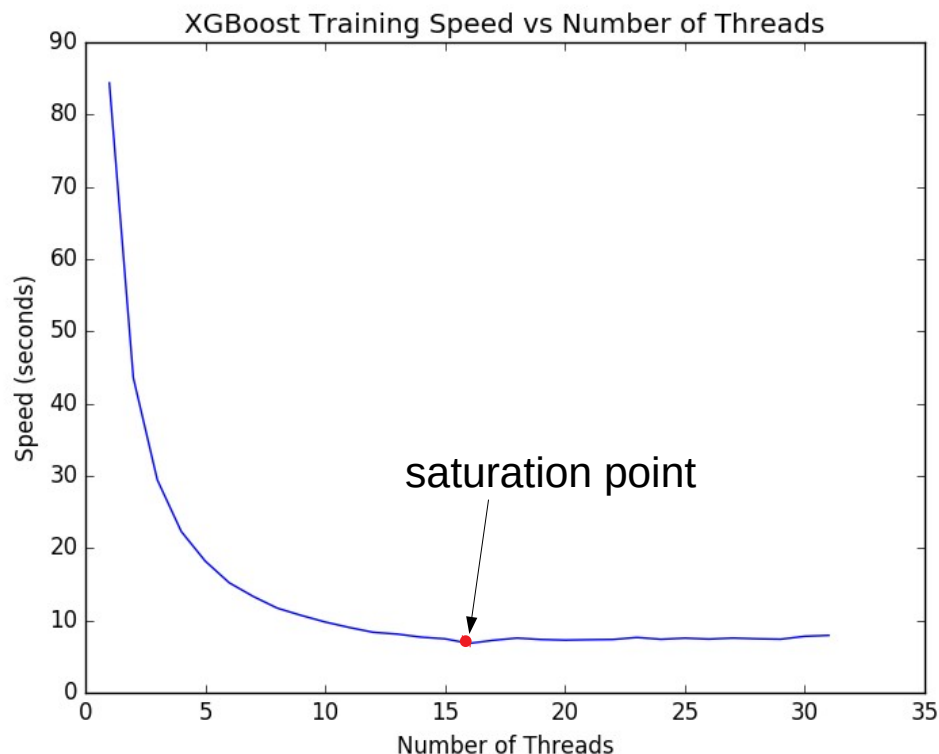
*What can we improve for  $10^6$ -dim vectors and 256 cores?*



# Operation parallelism

*Если мы идём быстро, тут можно поговорить про бустинг.*

# Operation parallelism

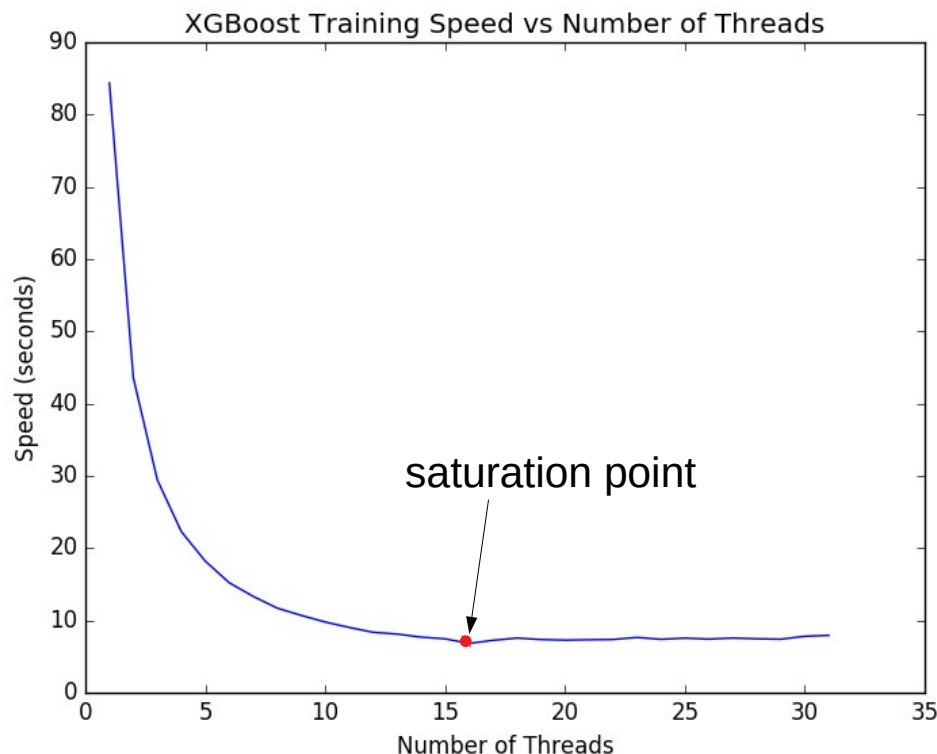


*More processes = more overhead*

- *waiting for each other*
- *sending data over the network*
- *performance fluctuations*

*Eventually adding more threads will no longer boost performance*

# Operation parallelism



*More processes = more overhead*

- *waiting for each other*
- *sending data over the network*
- *performance fluctuations*

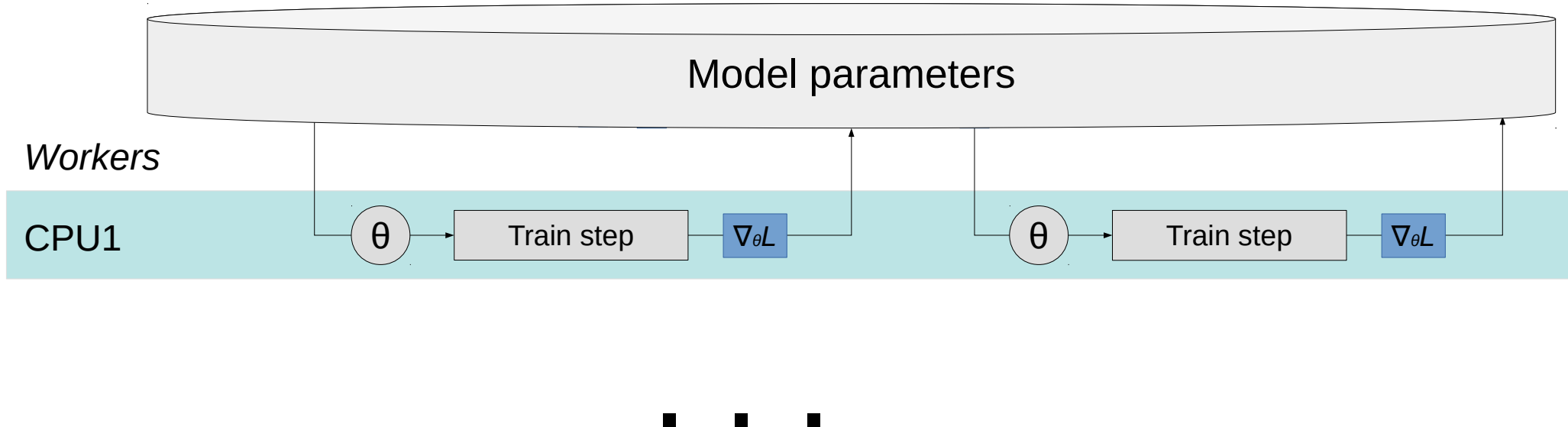
*Eventually adding more threads will no longer boost performance*

*How do we push this point further?*

# Parameter Server

Paper: [Smola et al. \(2010\)](#)

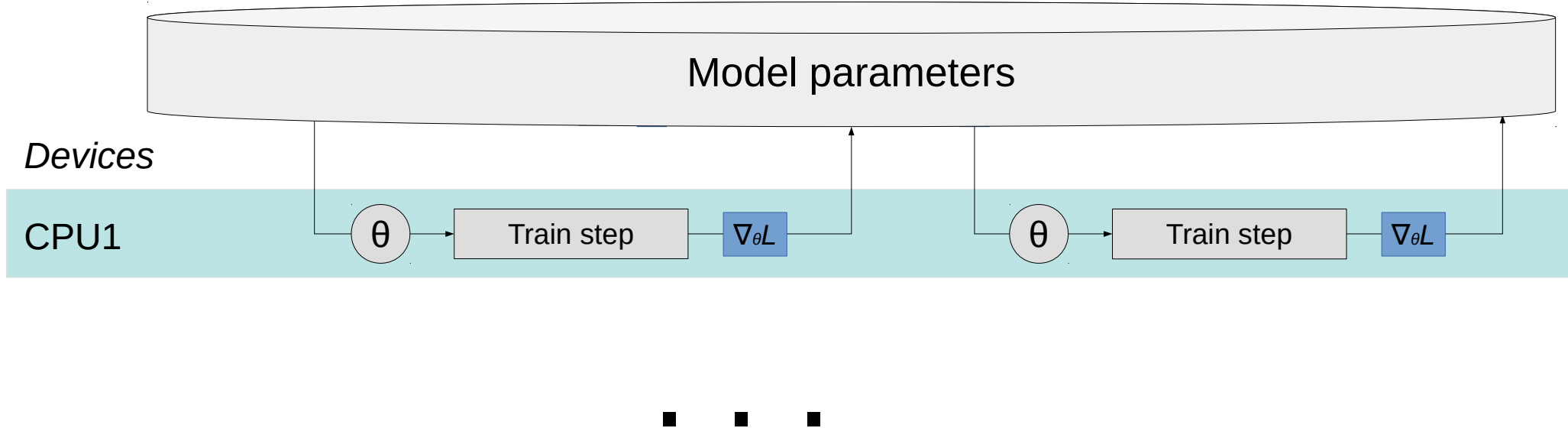
Make a dedicated process for parameters & optimizer



# Asynchronous training

HOGWILD! [arxiv.org/abs/1106.5730](https://arxiv.org/abs/1106.5730)

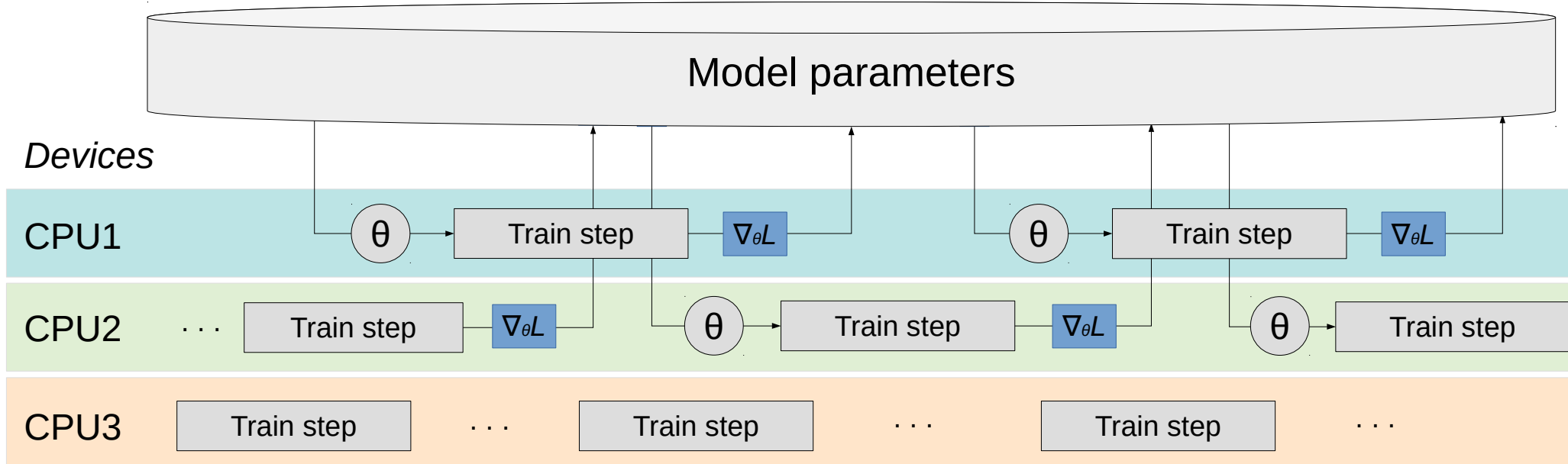
**Idea:** remove synchronization step altogether, use parameter server



# Asynchronous training

HOGWILD! [arxiv.org/abs/1106.5730](https://arxiv.org/abs/1106.5730)

**Idea:** remove synchronization step altogether, use parameter server

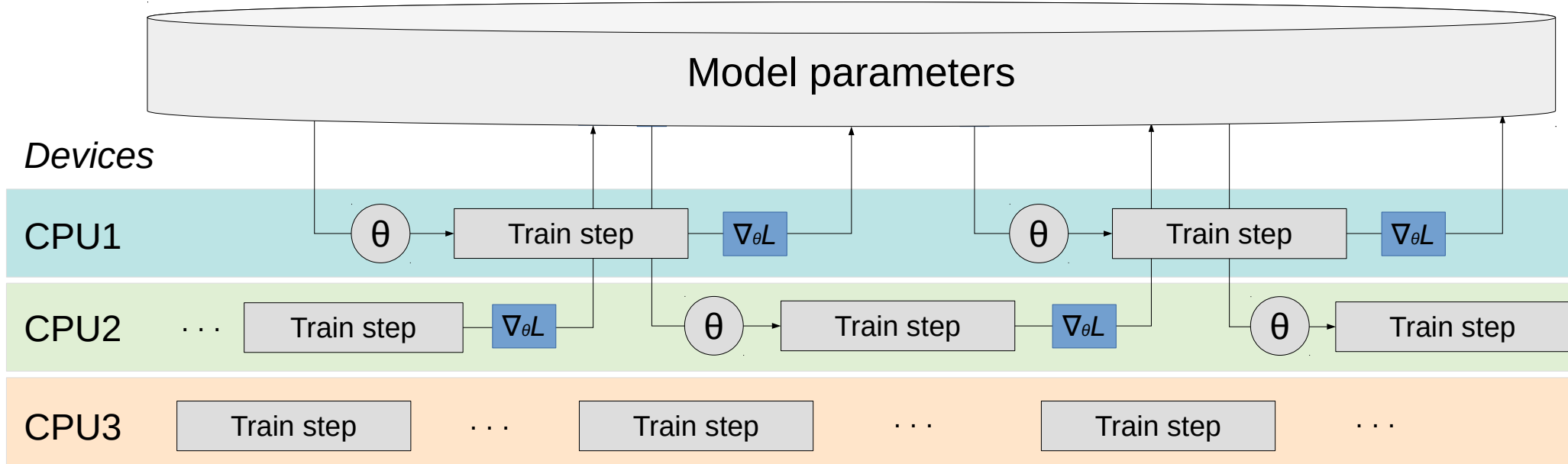




# Asynchronous training

HOGWILD! [arxiv.org/abs/1106.5730](https://arxiv.org/abs/1106.5730)

**Idea:** remove synchronization step altogether, use parameter server

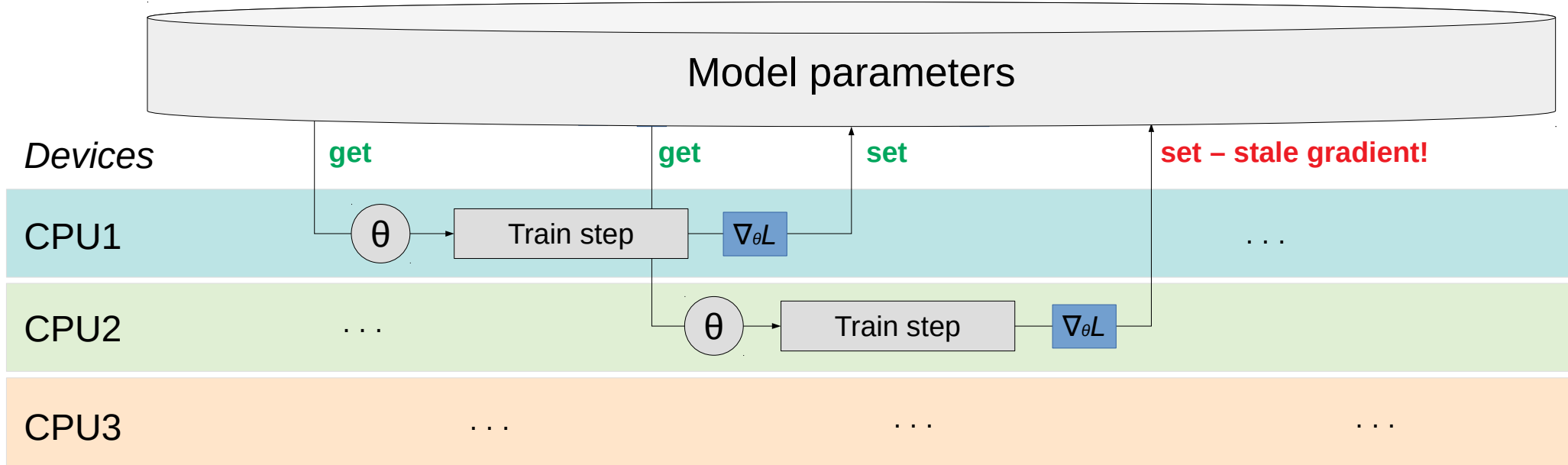


**Q:** have we lost anything by going asynchronous?

# Asynchronous training

HOGWILD! [arxiv.org/abs/1106.5730](https://arxiv.org/abs/1106.5730)

**Idea:** remove synchronization step altogether, use parameter server



# Staleness-aware SGD

Paper: [arxiv.org/abs/1511.05950](https://arxiv.org/abs/1511.05950) & others

**Updates accumulated:**  $c = \lfloor (\lambda/n) \rfloor$

**Average gradient:**  $g_i = \frac{1}{c} \sum_{l=1}^c \alpha(\tau_{i,l}) \Delta\theta_l, \quad l \in \{1, 2, \dots, \lambda\}$

**New parameters:**  $\theta_{i+1} = \theta_i - g_i,$

# Staleness-aware SGD

Paper: [arxiv.org/abs/1511.05950](https://arxiv.org/abs/1511.05950) & others

**Updates accumulated:**  $c = \lfloor (\lambda/n) \rfloor$   **$n$  = total workers**  
 **$\lambda$  = “accumulation factor”**

**Average gradient:** 
$$g_i = \frac{1}{c} \sum_{l=1}^c \alpha(\tau_{i,l}) \Delta \theta_l, \quad l \in \{1, 2, \dots, \lambda\}$$

**New parameters:** 
$$\theta_{i+1} = \theta_i - g_i,$$

# Staleness-aware SGD

Paper: [arxiv.org/abs/1511.05950](https://arxiv.org/abs/1511.05950) & others

Updates accumulated:  $c = \lfloor (\lambda/n) \rfloor$

Average gradient:  $g_i = \frac{1}{c} \sum_{l=1}^c \alpha(\tau_{i,l}) \Delta\theta_l, \quad l \in \{1, 2, \dots, \lambda\}$

**staleness-dependent**

New parameters:  $\theta_{i+1} = \theta_i - g_i,$  **“learning rate”**

# Staleness-aware SGD

Paper: [arxiv.org/abs/1511.05950](https://arxiv.org/abs/1511.05950) & others

**Updates accumulated:**  $c = \lfloor (\lambda/n) \rfloor$

**Average gradient:**  $g_i = \frac{1}{c} \sum_{l=1}^c \alpha(\tau_{i,l}) \Delta\theta_l, \quad l \in \{1, 2, \dots, \lambda\}$

**New parameters:**  $\theta_{i+1} = \theta_i - g_i,$

$$\alpha_{i,l} = \frac{\alpha_0}{\tau_{i,l}}$$

**base learning rate**

**staleness ( $\geq 1$ )**

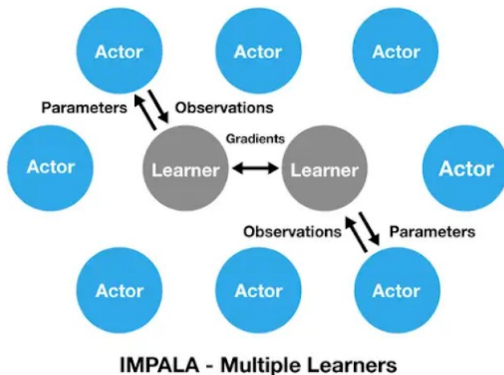
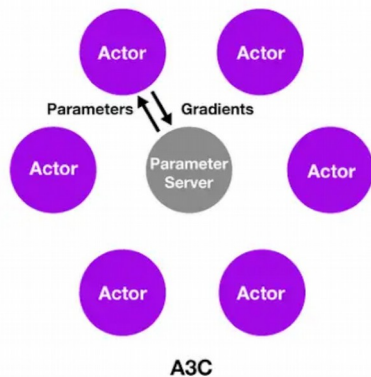
*aka number of skipped updates*

# Parameter Server Applications

**Conventional ML:** *e.g. (Logistic Regression, CNN classifiers)*

Paper (sharded PS): <https://www.cs.cmu.edu/~muli/file/ps.pdf>  
Another paper (optimization tricks): [parameter\\_server\\_nips14.pdf](#)  
[PyTorch tutorial \(hogwild\)](#), [TF tutorial \(parameter server\)](#)

## Reinforcement learning:



Async. RL: [arxiv.org/abs/1602.01783](https://arxiv.org/abs/1602.01783)

IMPALA: [arxiv.org/abs/1802.01561](https://arxiv.org/abs/1802.01561)

SEED RL: [arxiv.org/abs/1910.06591](https://arxiv.org/abs/1910.06591)

## More:

(english) <https://youtu.be/kOy49NqZeqI>

(russian) <https://youtu.be/wswbMkT55ml>

Привет, Ёж. Это – последний слайд.  
Если вы уже здесь, а время ещё осталось,  
разбери-ка вот эту статью:

<https://www.usenix.org/system/files/osdi20-jiang.pdf>