

Distributed Learning and You

Ruksi Laine
@ruk_si
<https://valohai.com>
@valohaiai

Me?

Hi, I'm Ruksi.

Machine Learning Engineer

<https://ruk.si>

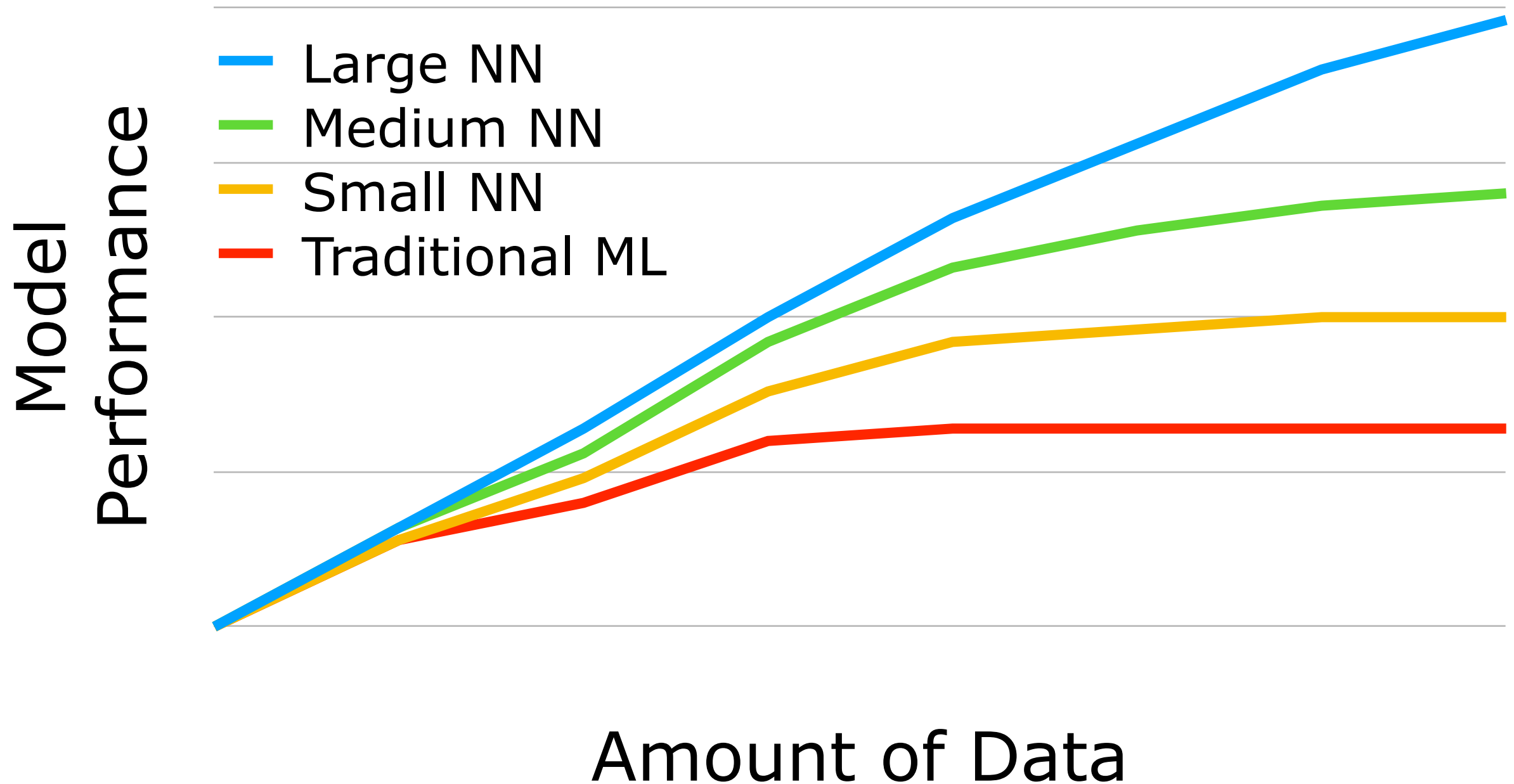
@ruk_si



@valohaiaia

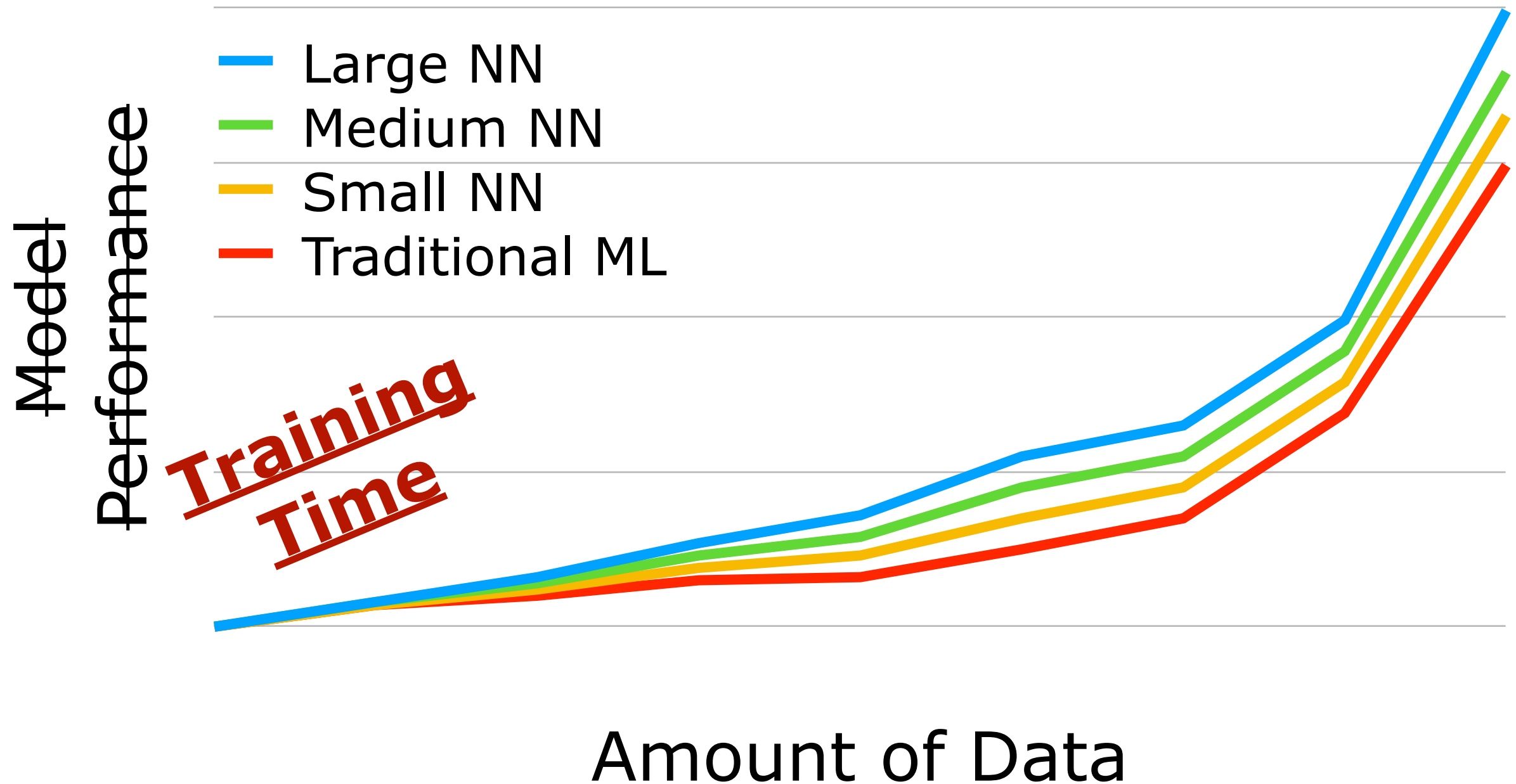


Scale Drives Deep Learning Process



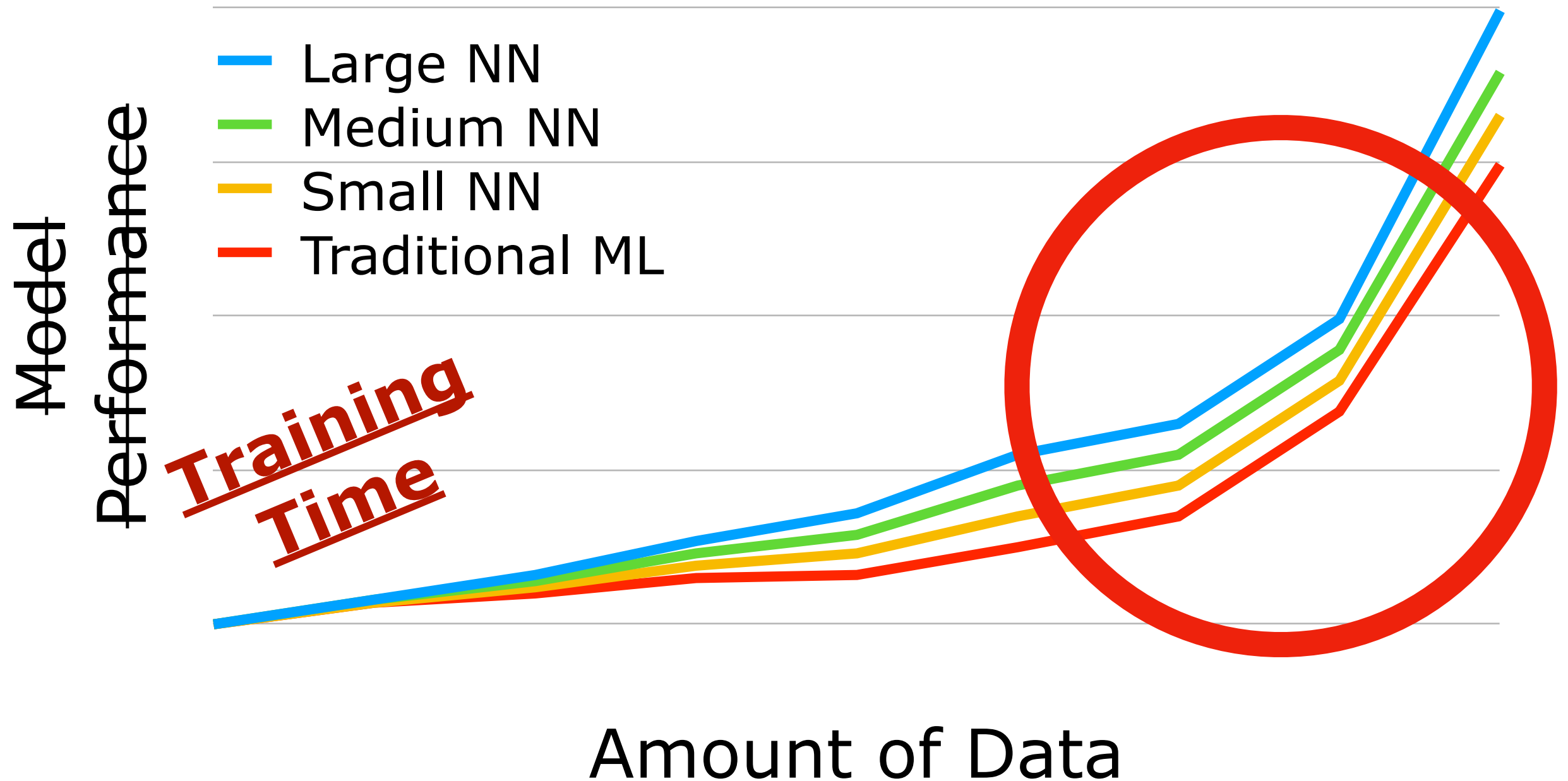
Scale Makes Everything Slow

~~Scale Drives Deep Learning Process~~



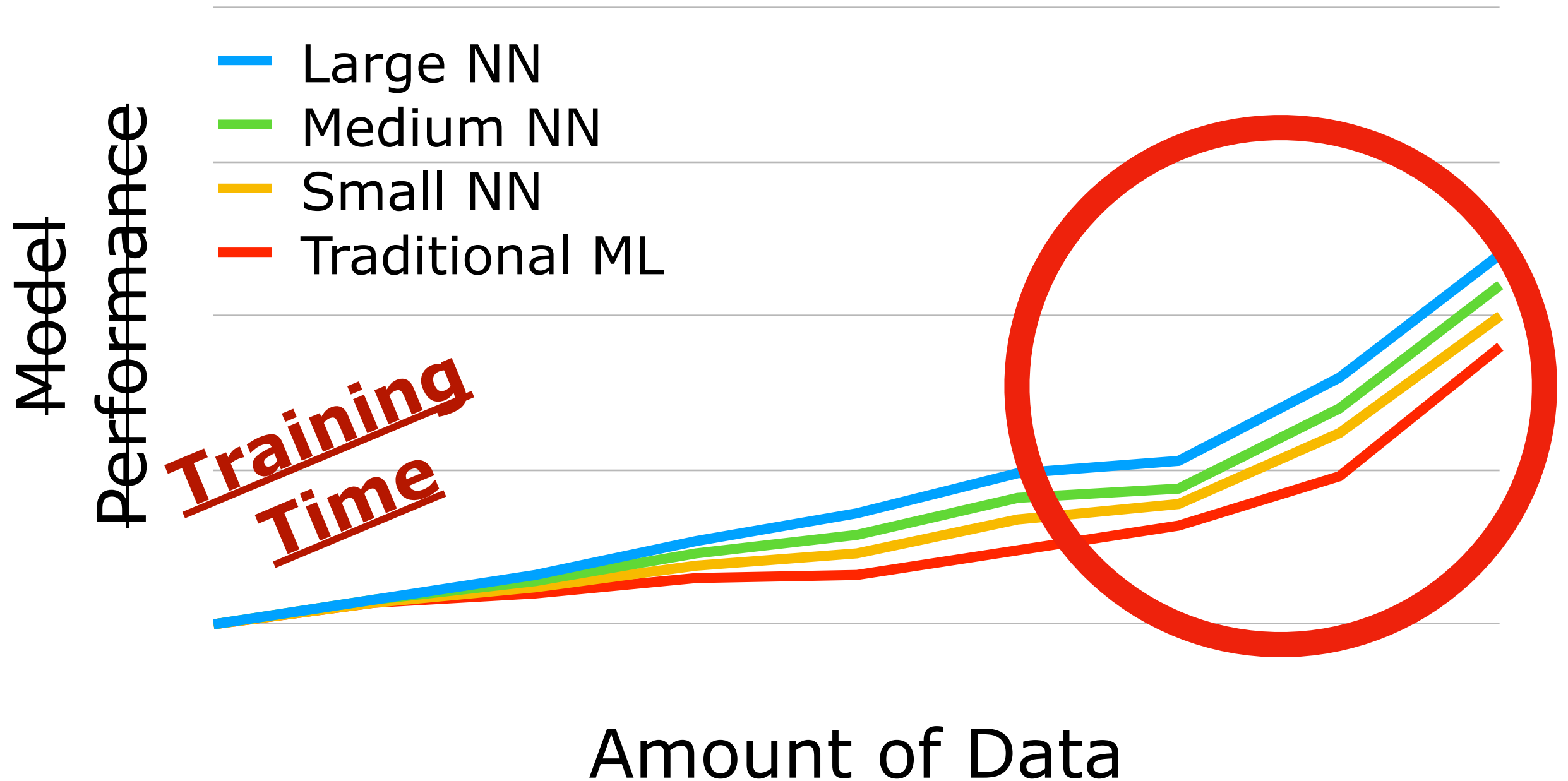
Scale Makes Everything Slow

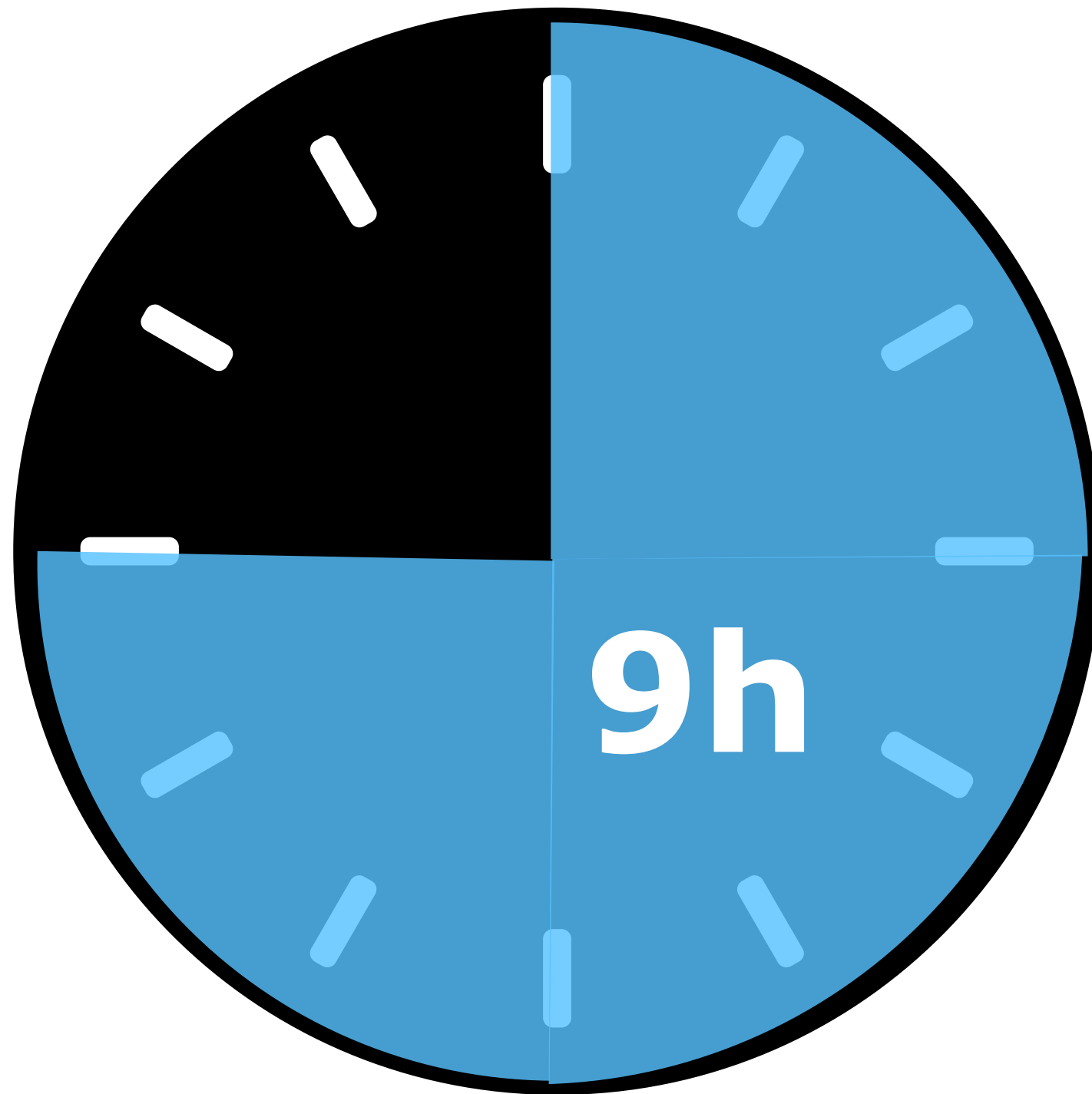
~~Scale Drives Deep Learning Process~~

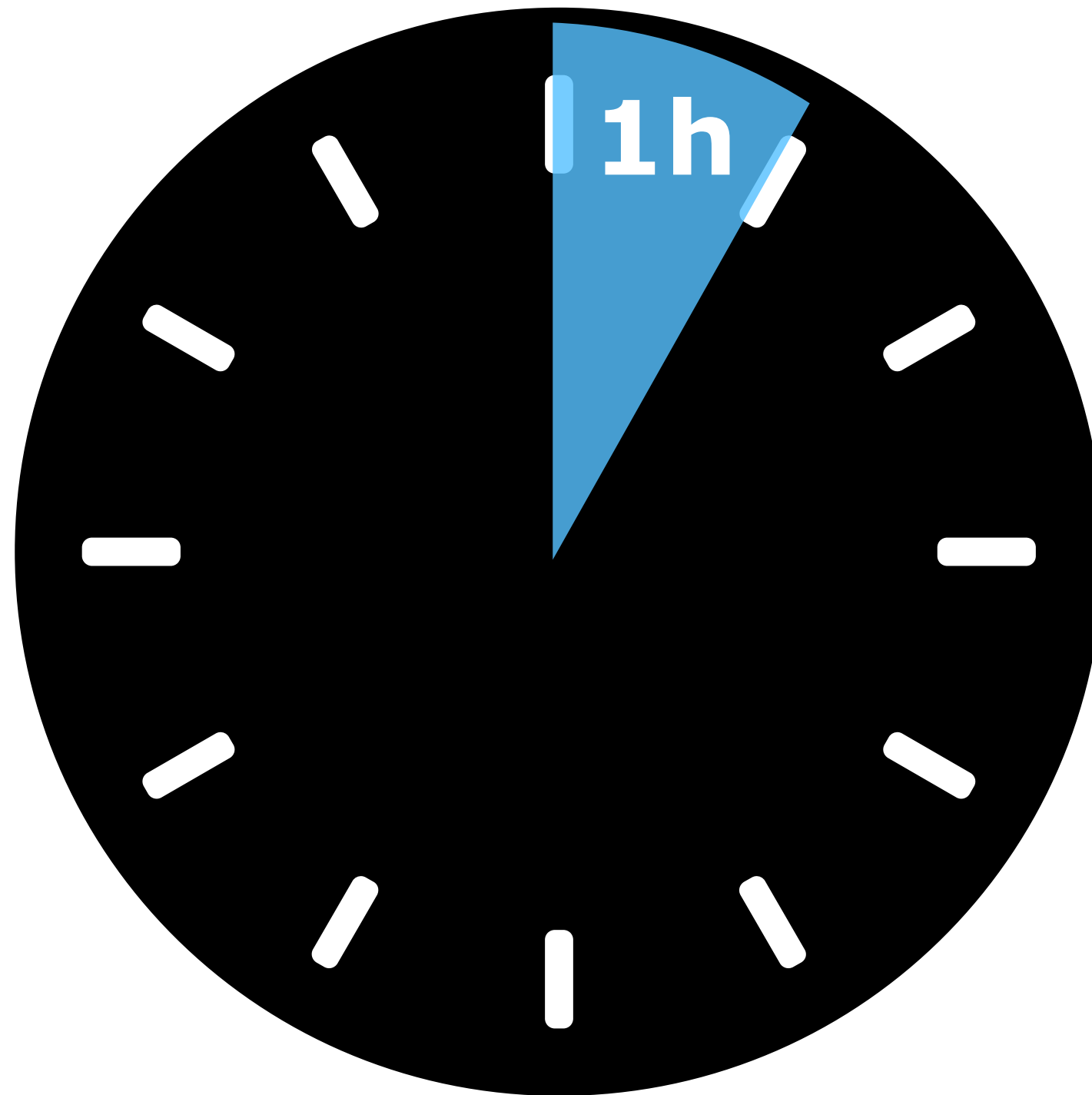


Scale Makes Everything Slow

~~Scale Drives Deep Learning Process~~













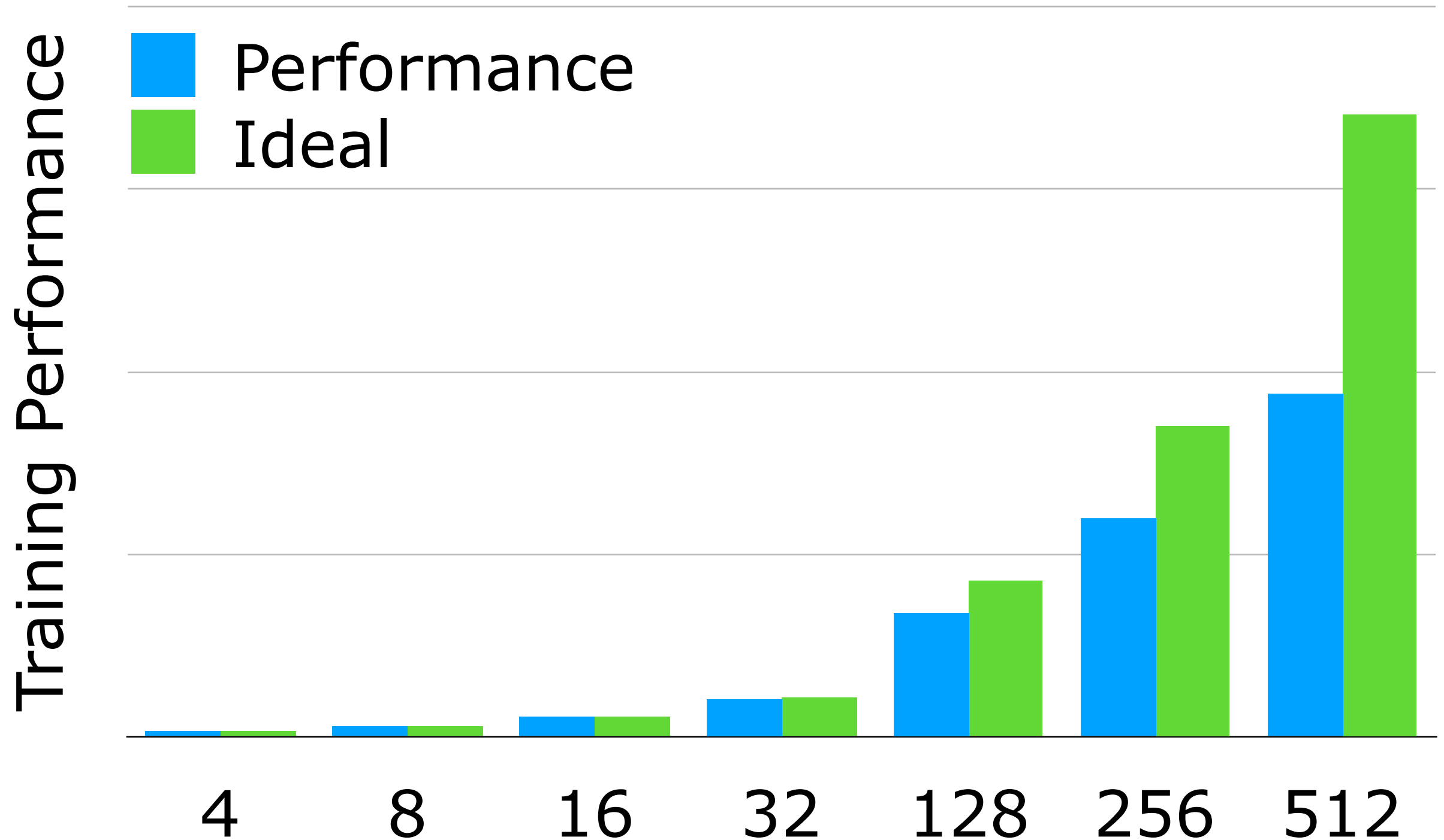






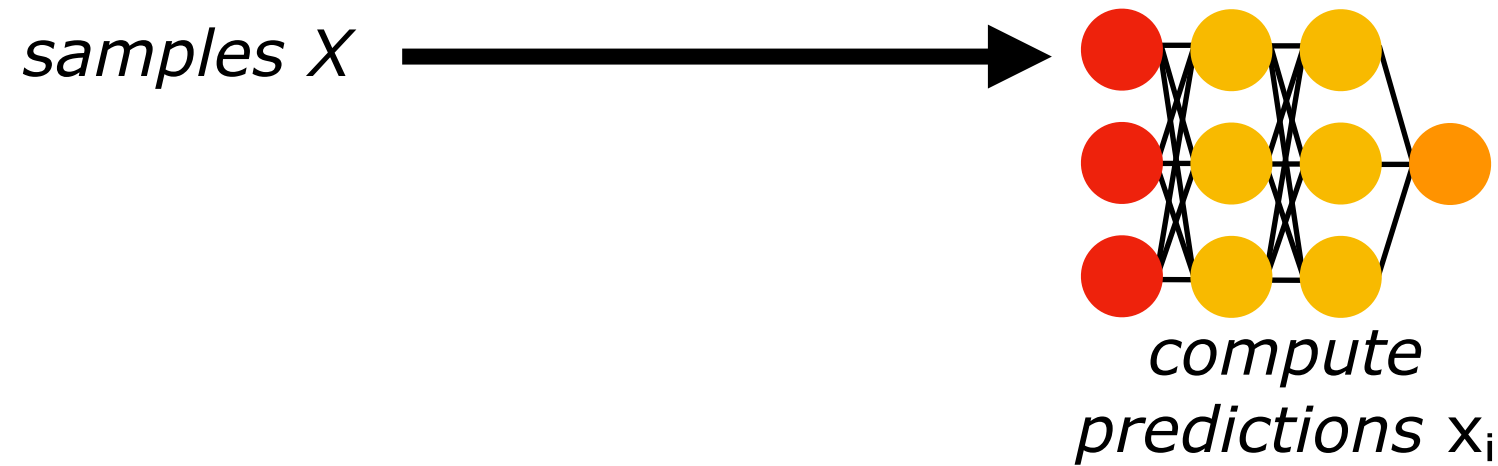


Scaling Efficiency

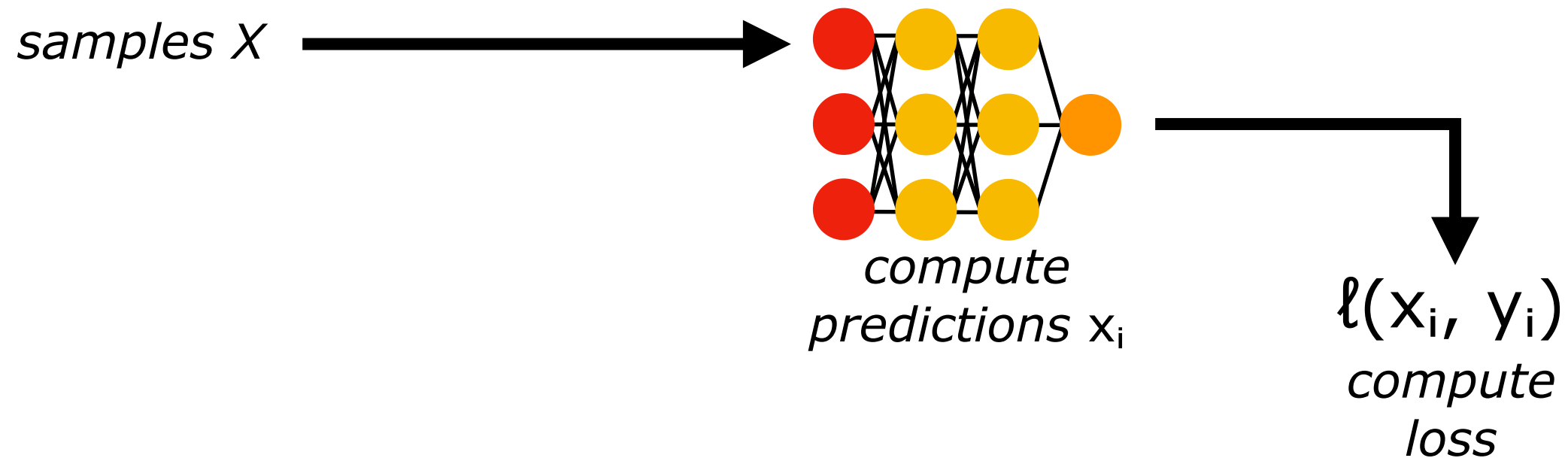


Deep Learning is Innately Sequential

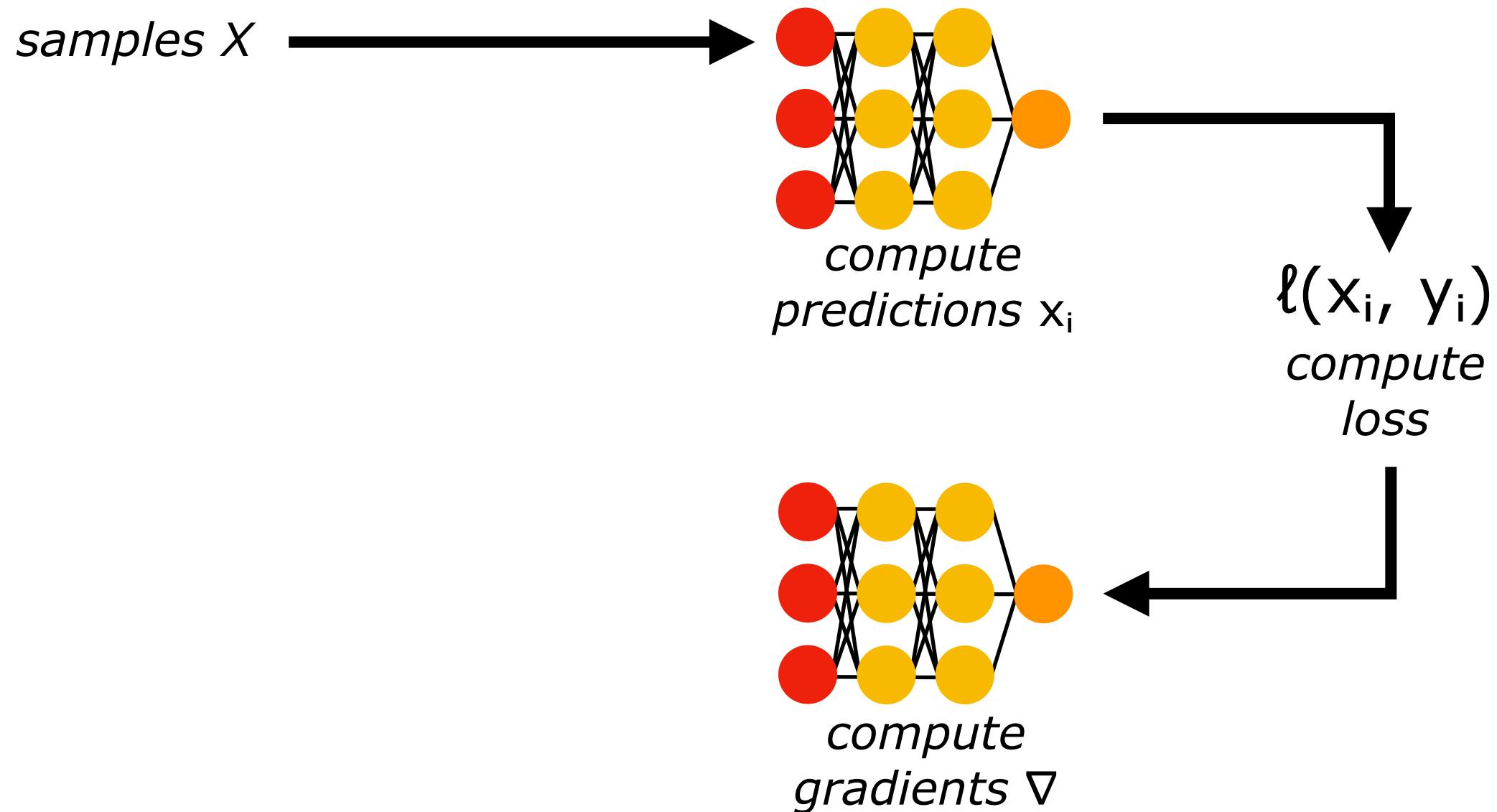
Training Loop



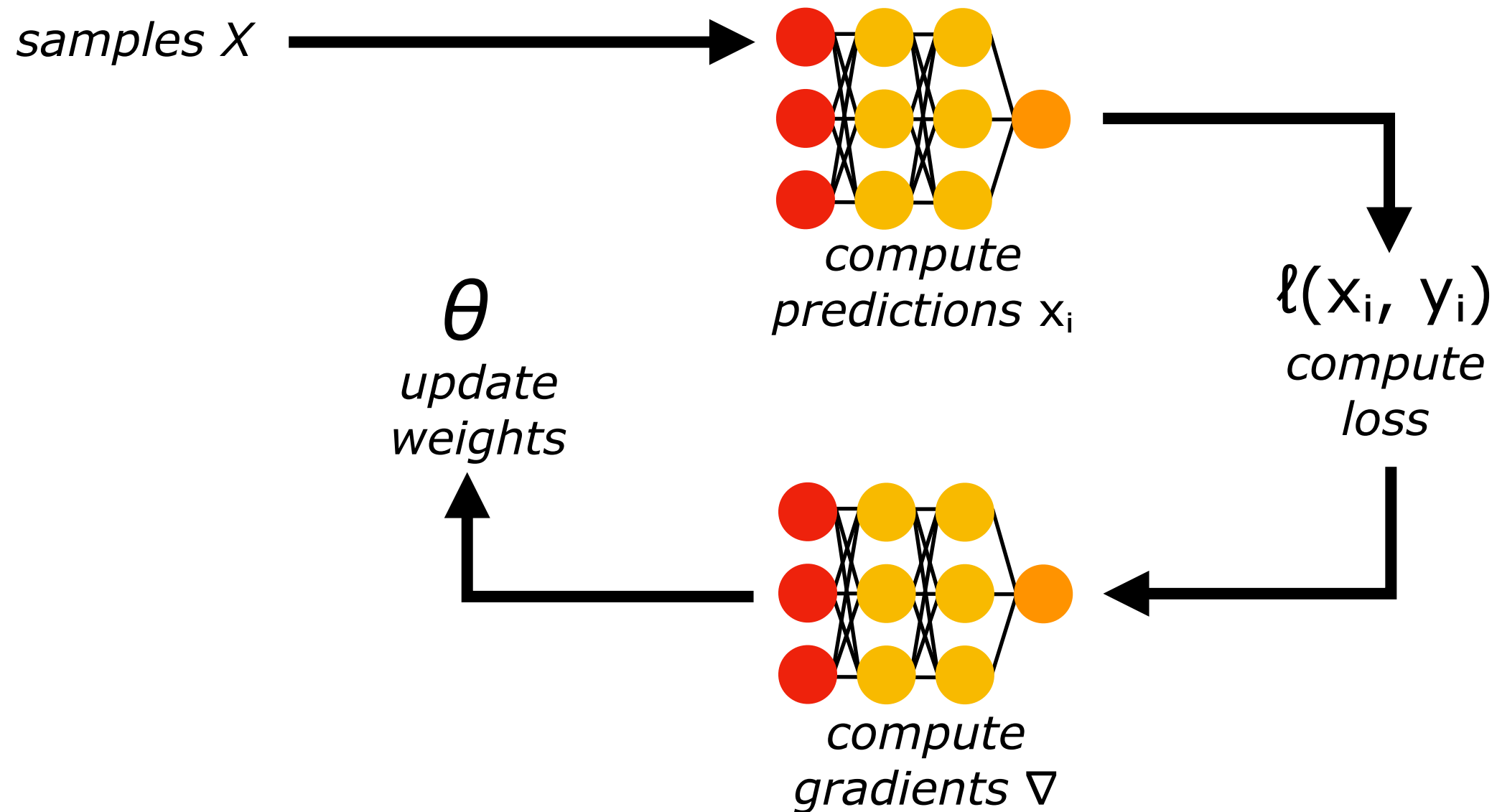
Training Loop



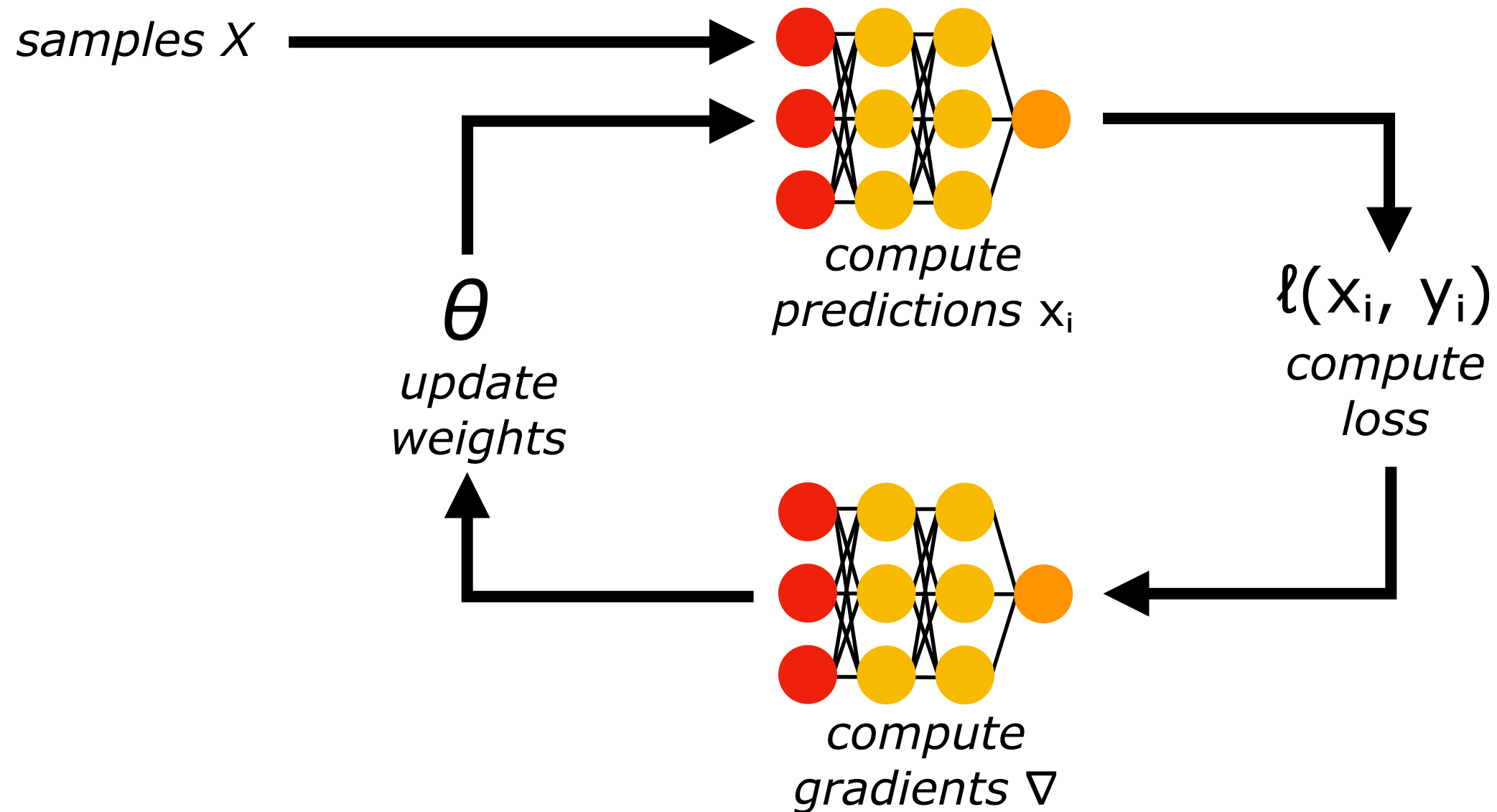
Training Loop



Training Loop



Training Loop

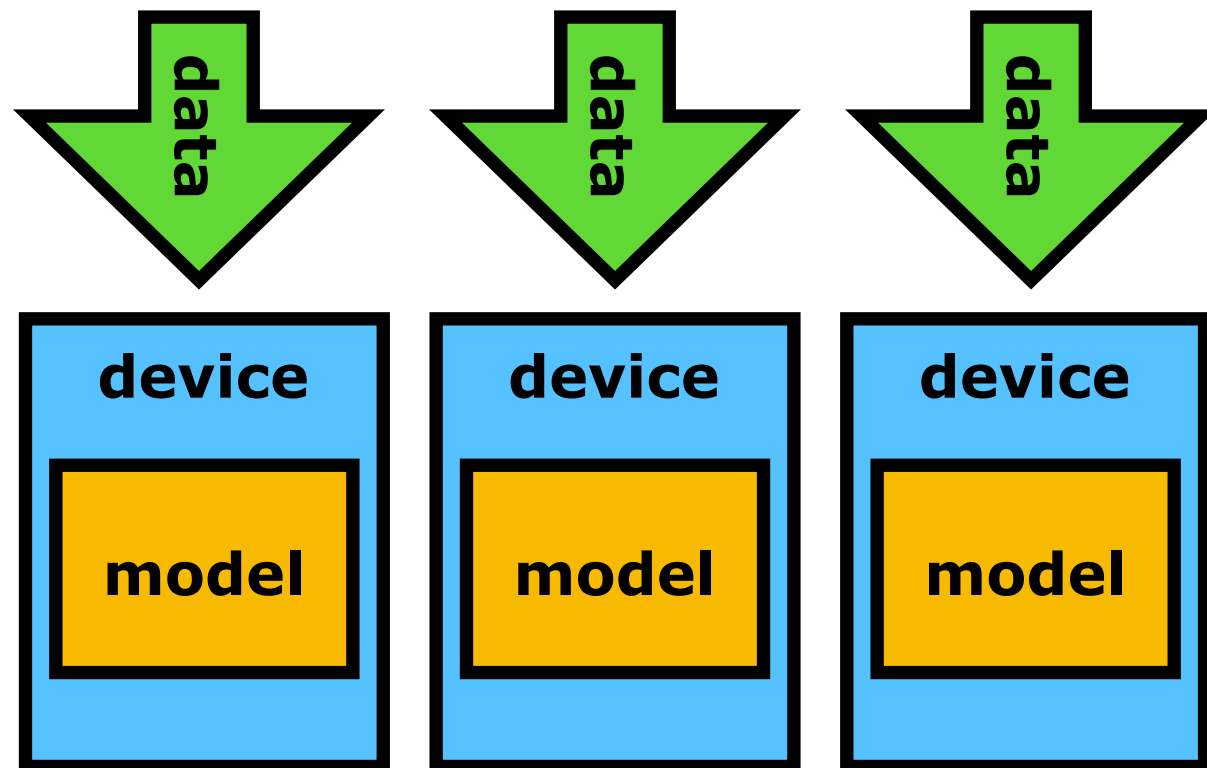


Two Main Branches of Distributed Learning

Data
Parallelization

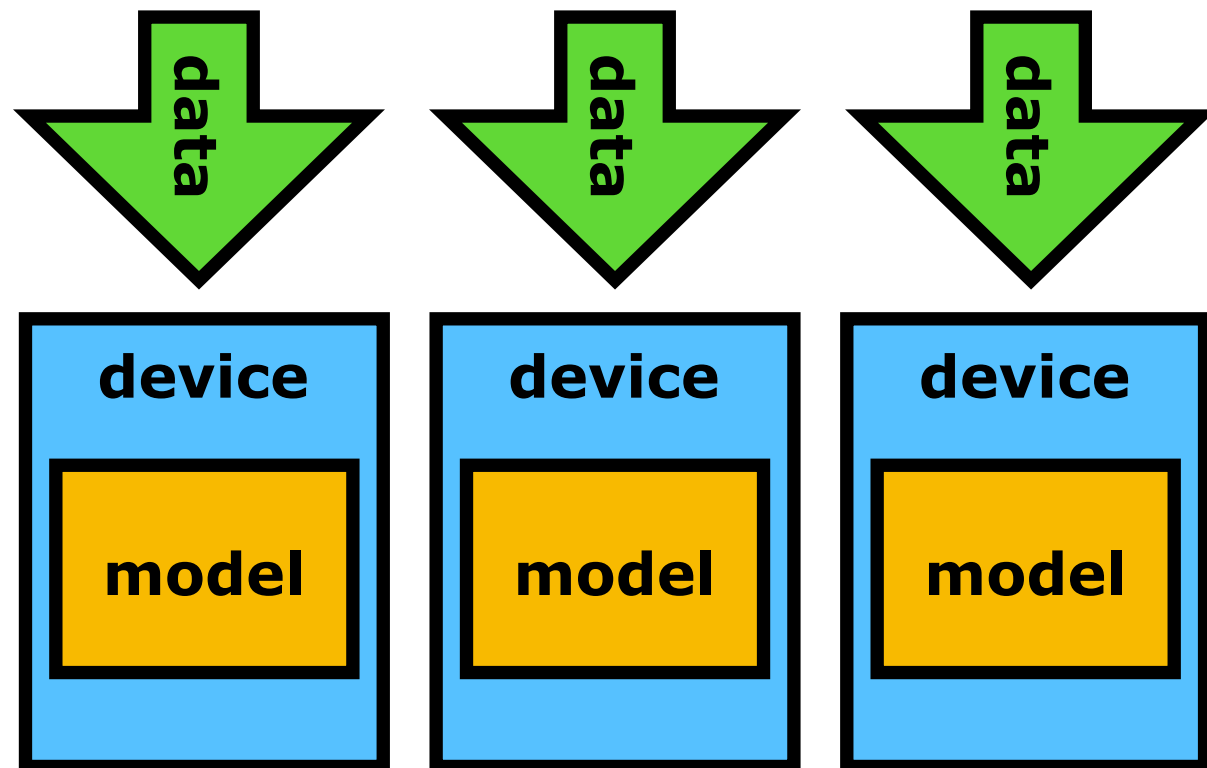
Model
Parallelization

Data Parallelization

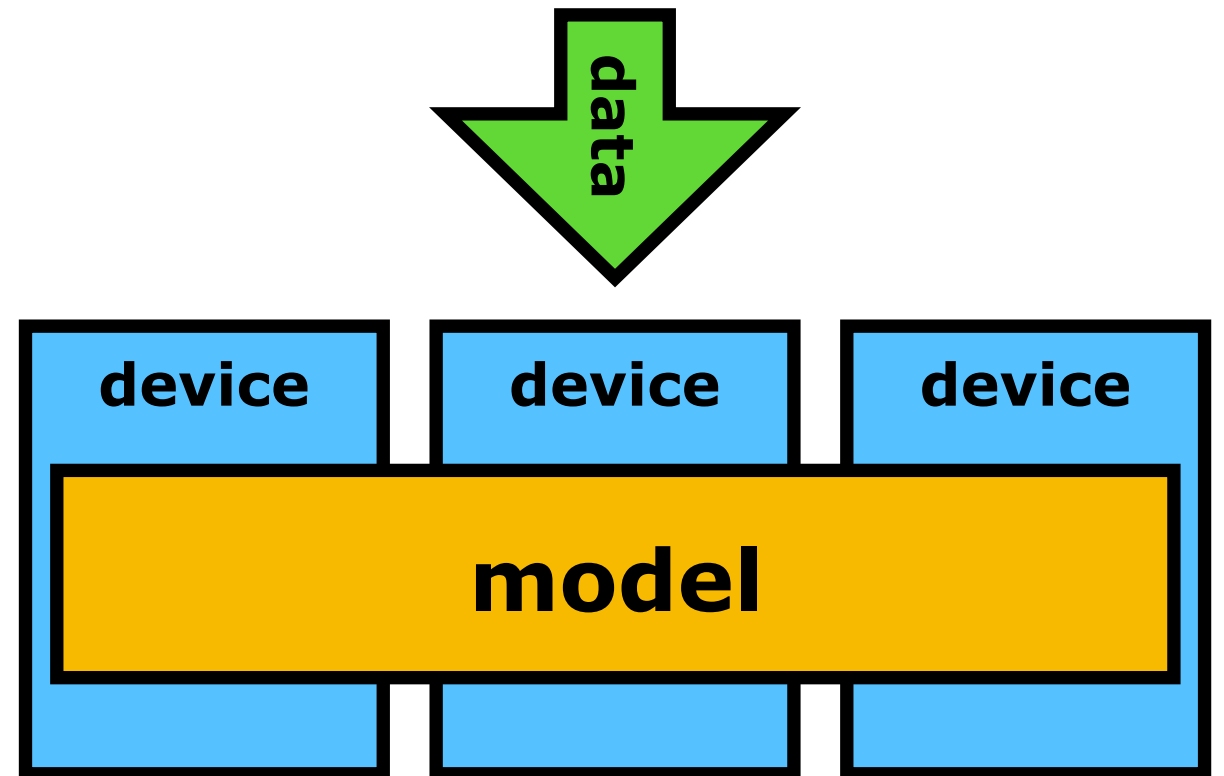


Model Parallelization

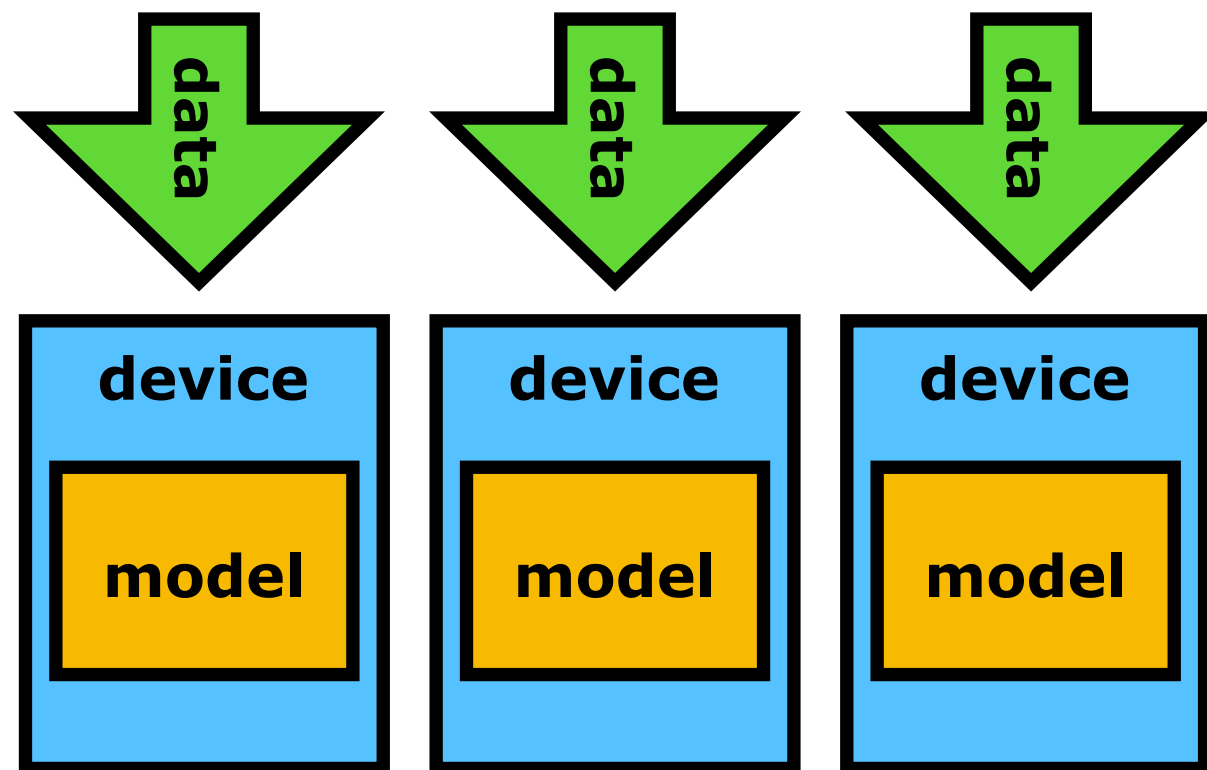
Data Parallelization



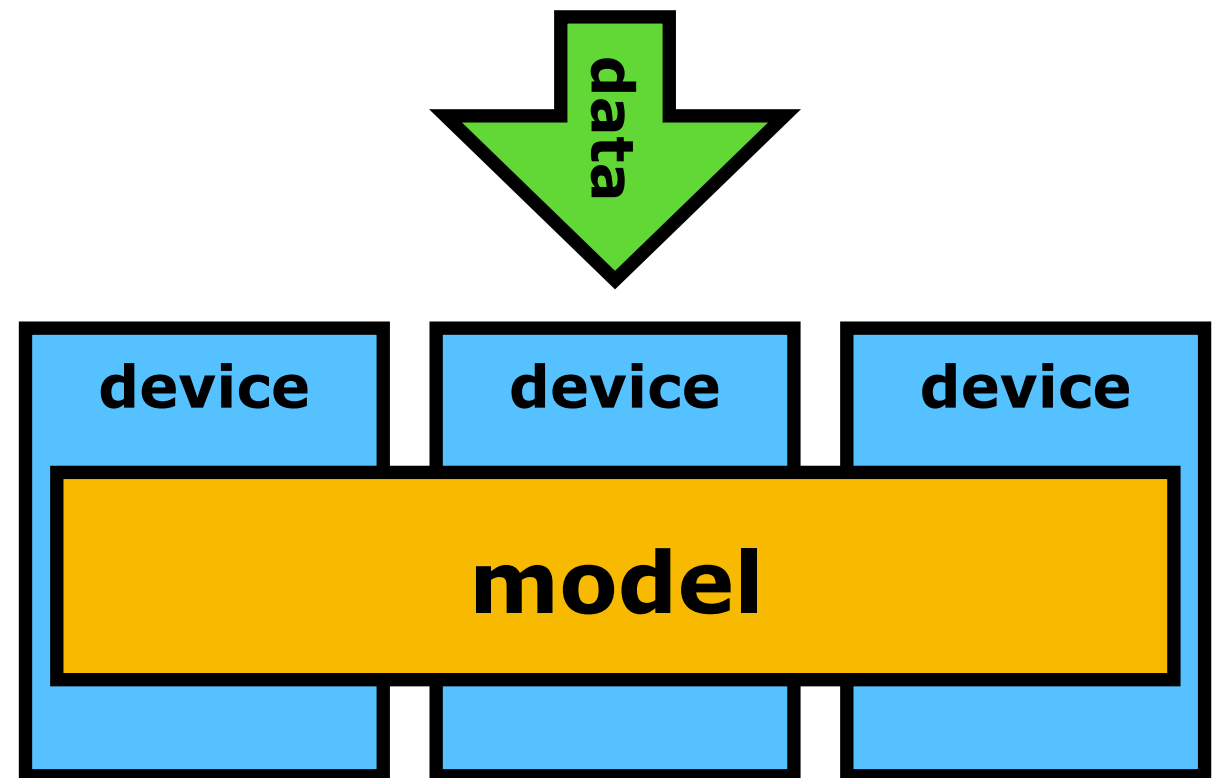
Model Parallelization



Data Parallelization

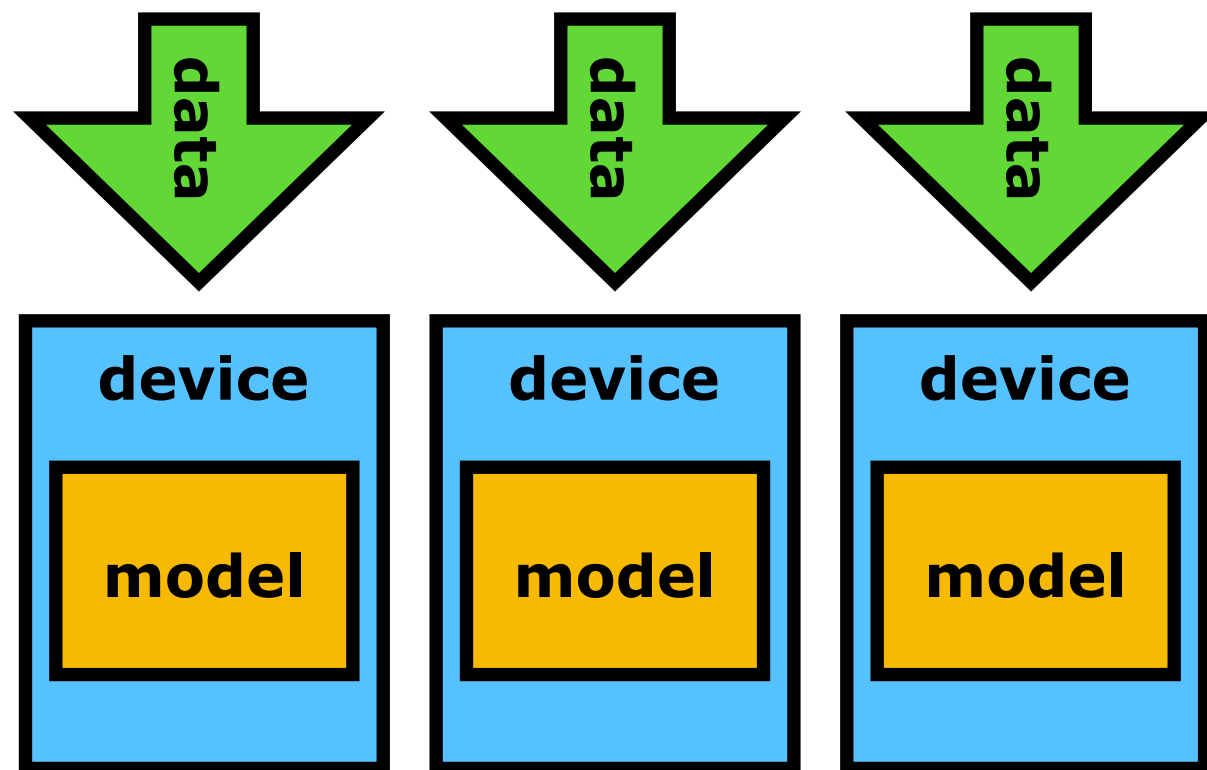


Model Parallelization



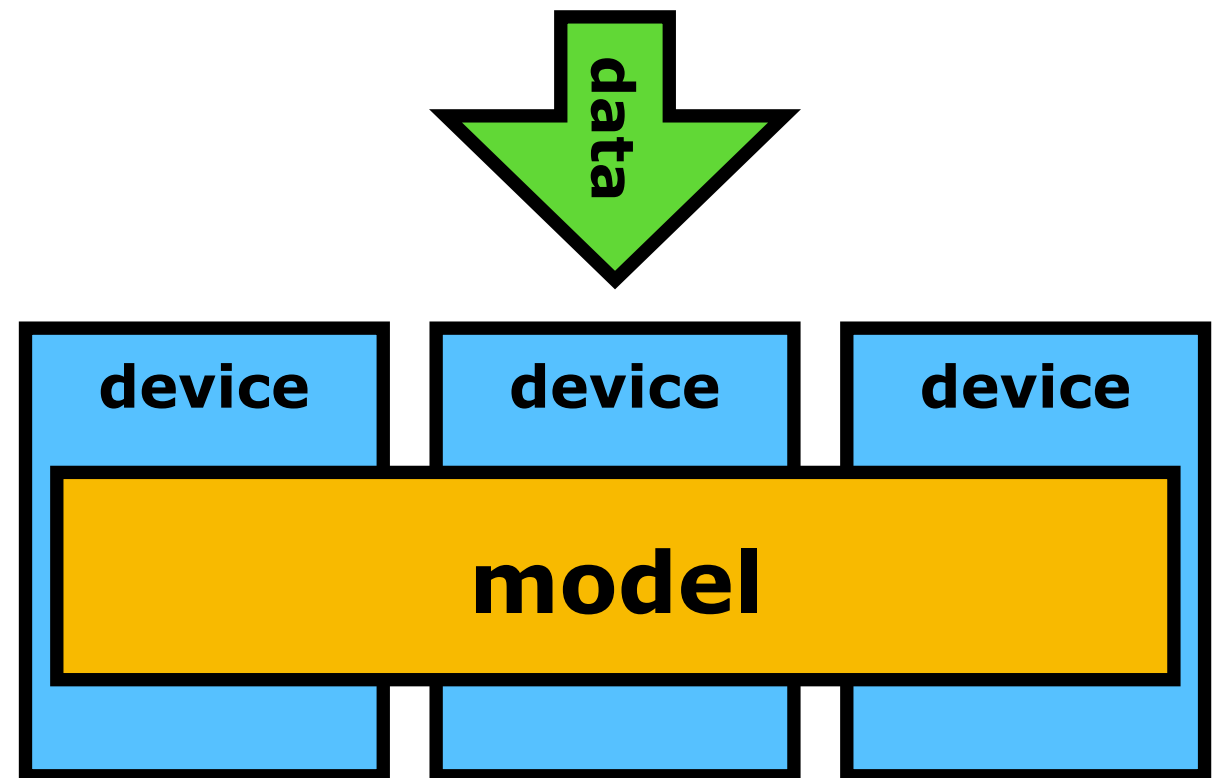
Synchronize weights θ to
train models fast.

Data Parallelization



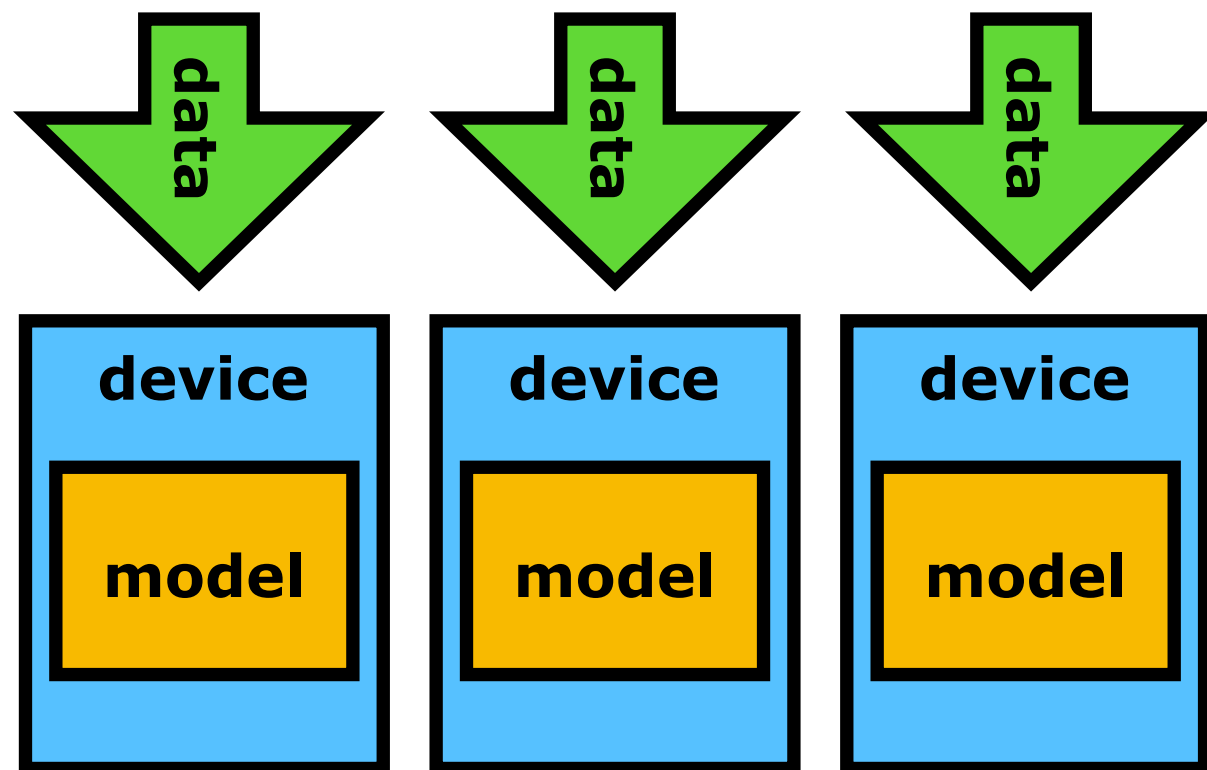
Synchronize weights θ to
train models fast.

Model Parallelization



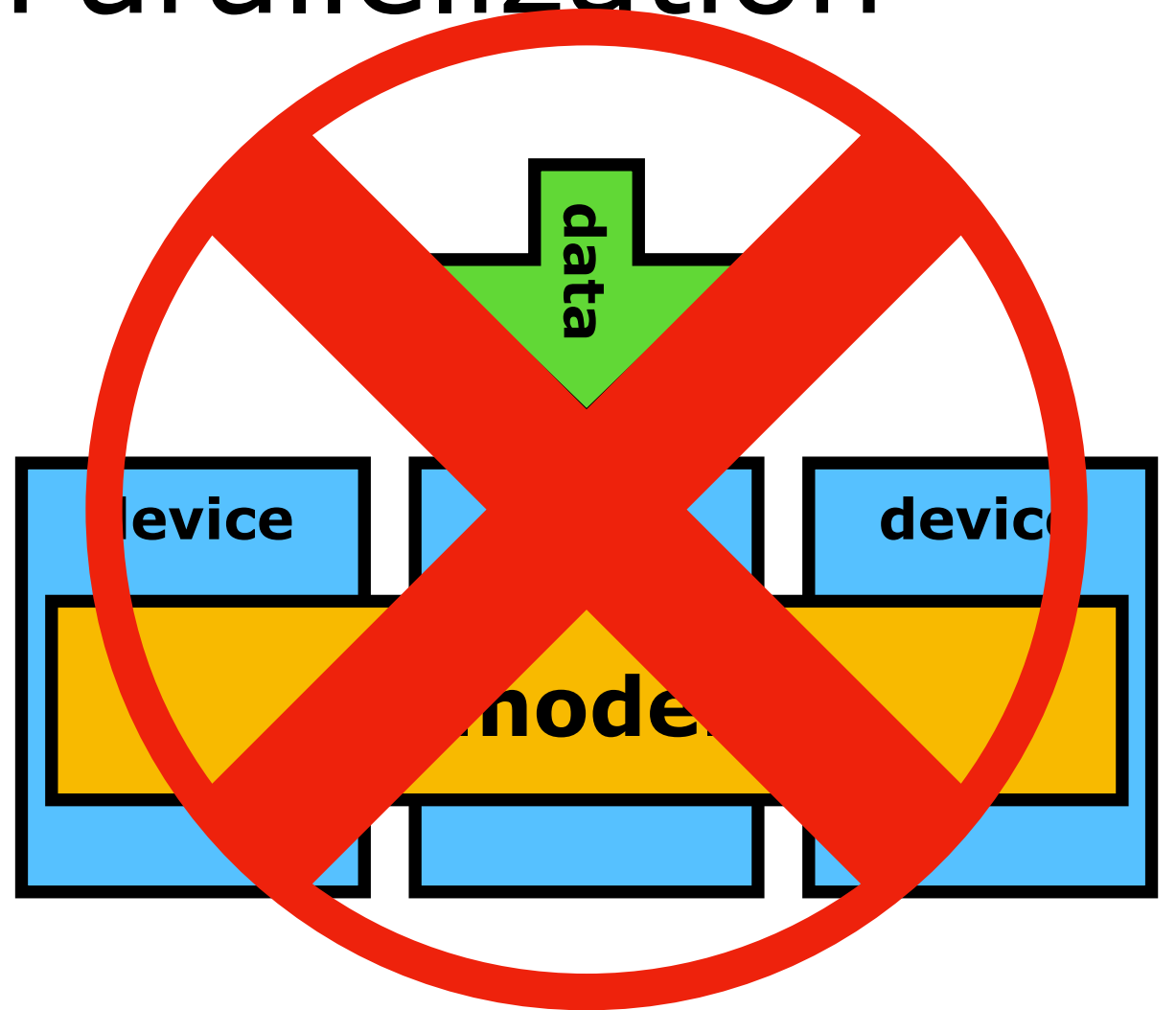
Partition a model to
build huge models.

Data Parallelization



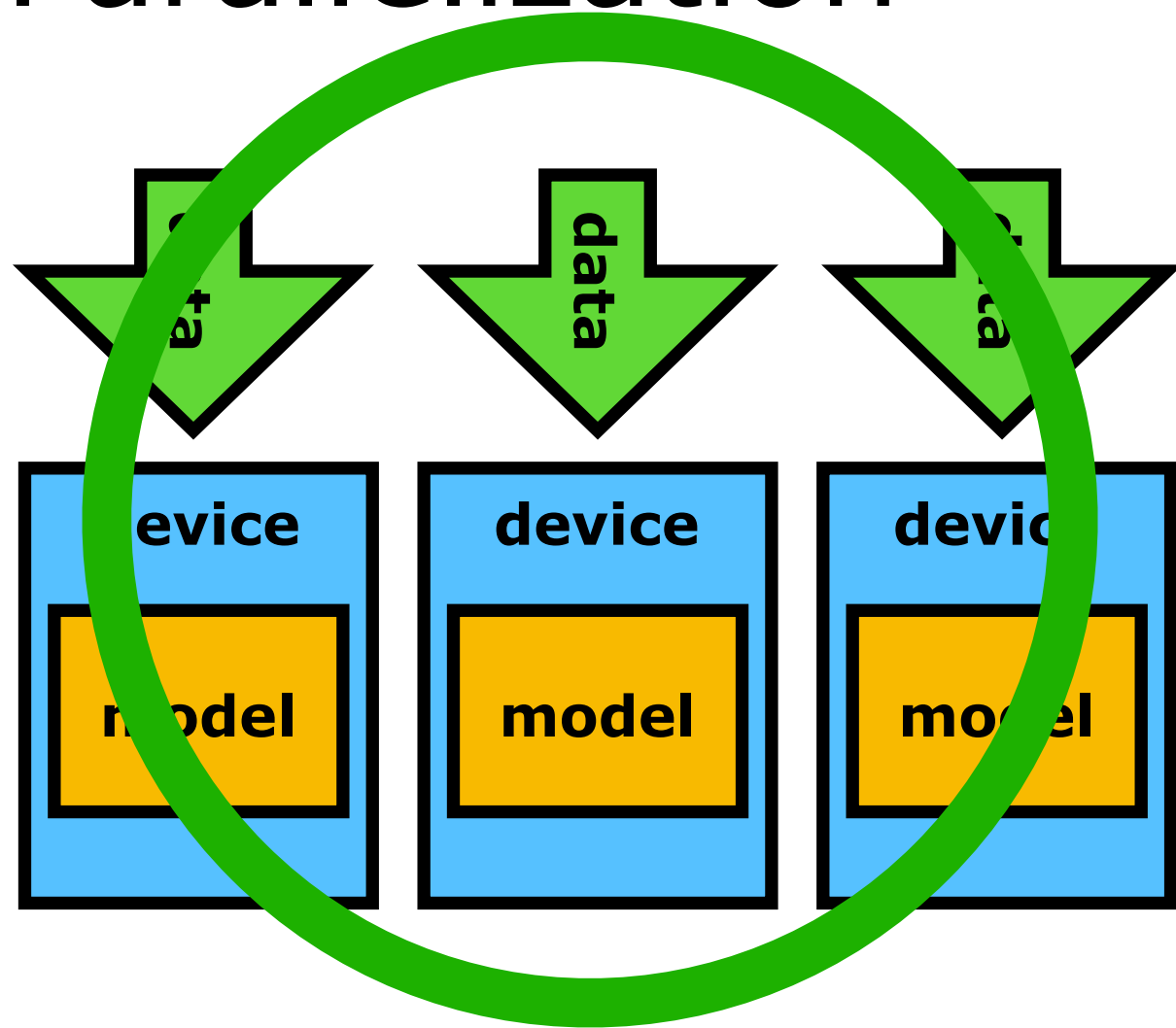
Synchronize weights θ to
train models fast.

Model Parallelization



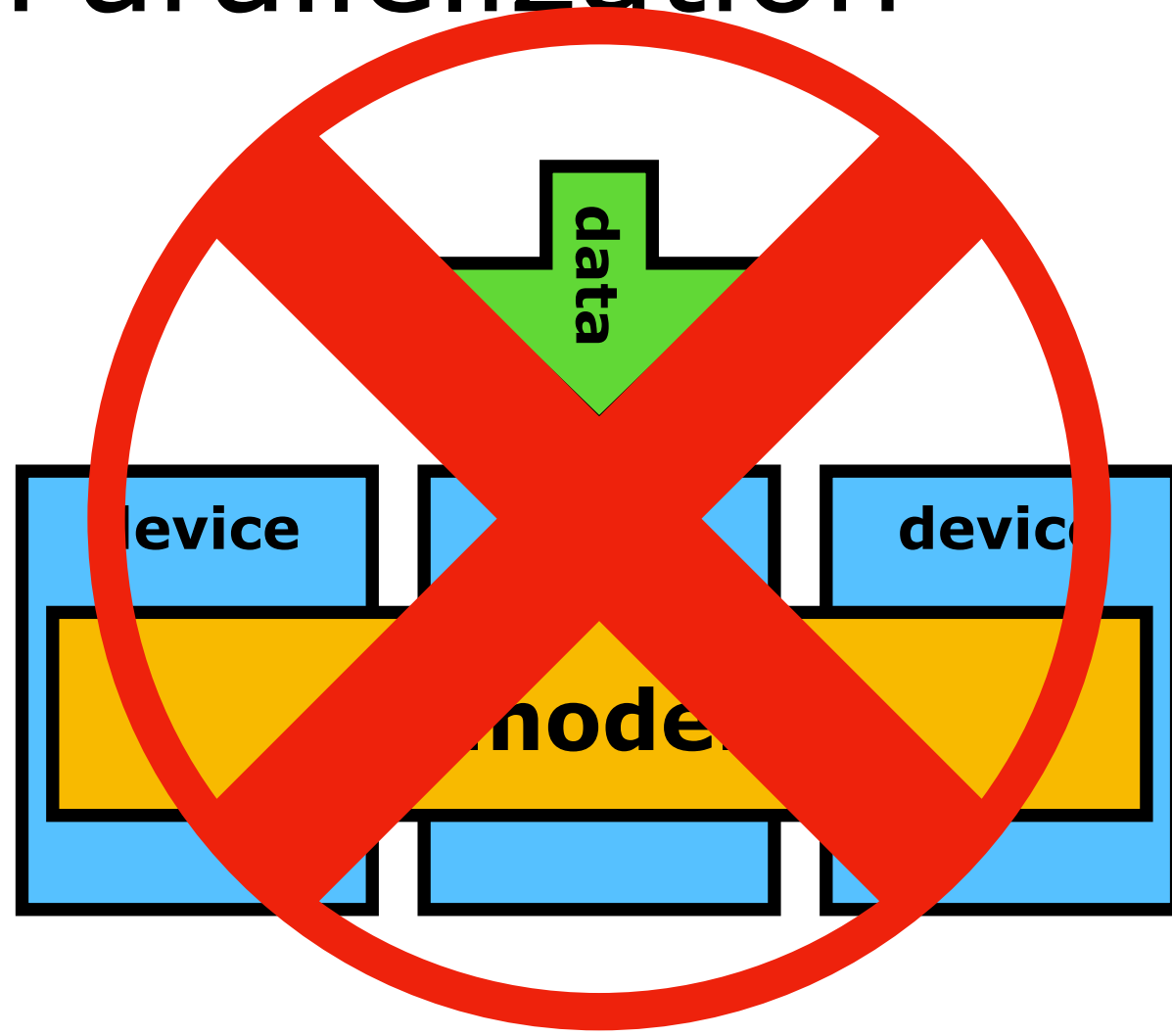
Partition a model to
build huge models.

Data Parallelization



Synchronize weights θ to
train models fast.

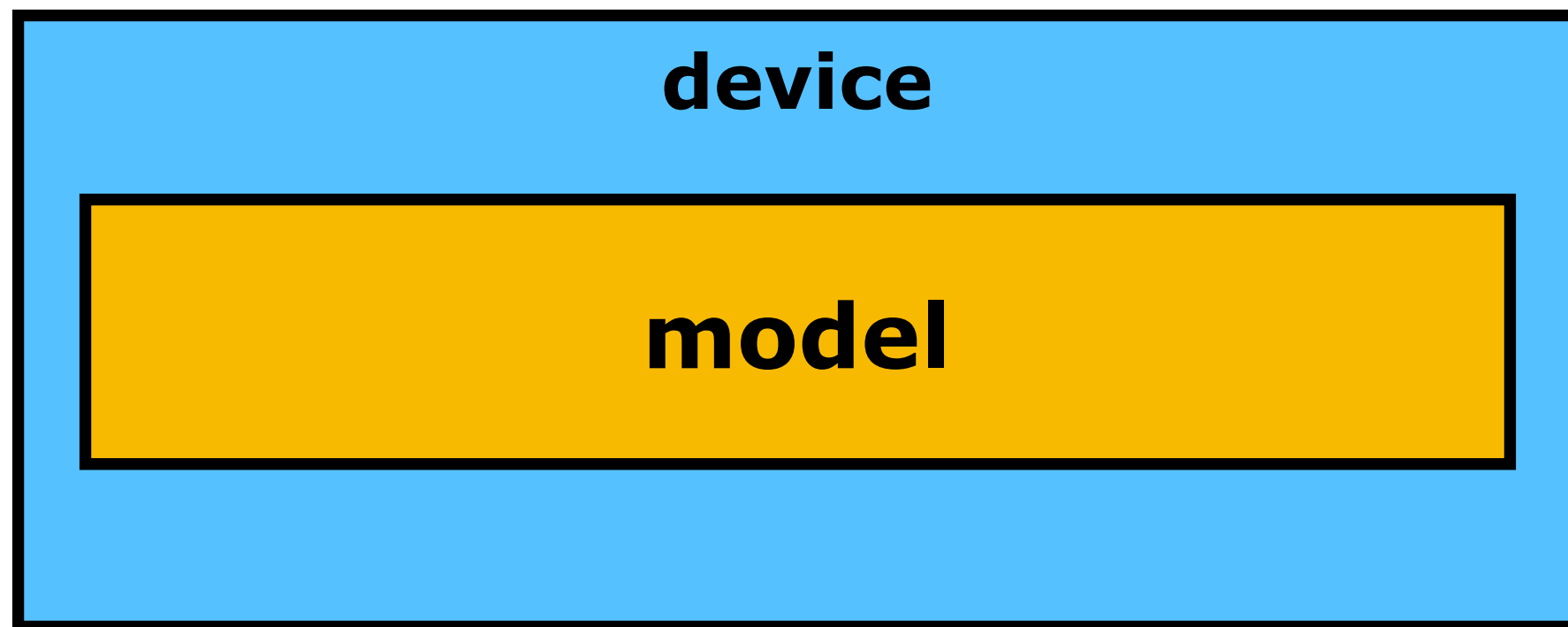
Model Parallelization



Partition a model to
build huge models.

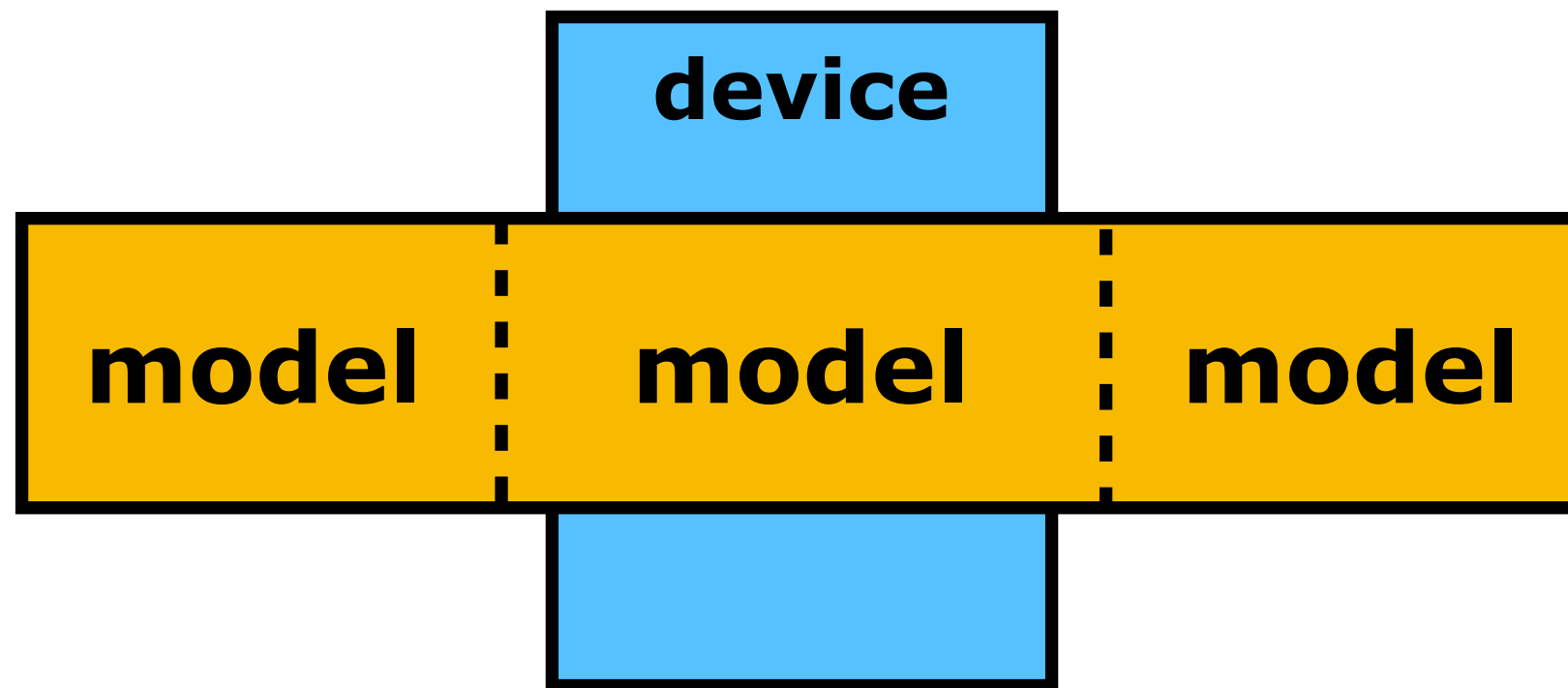
device

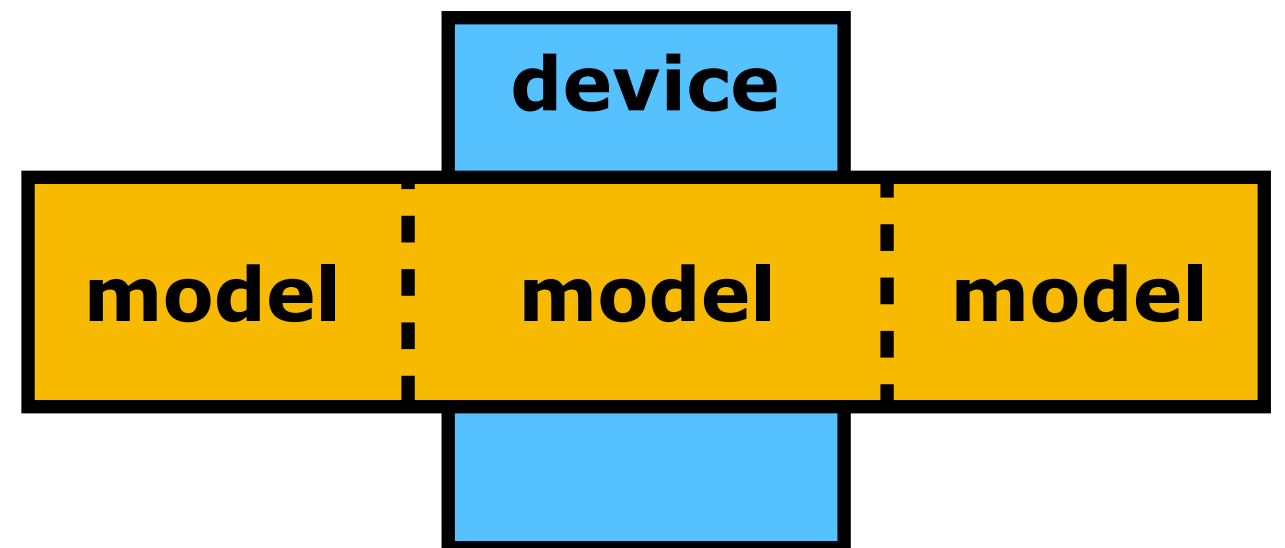
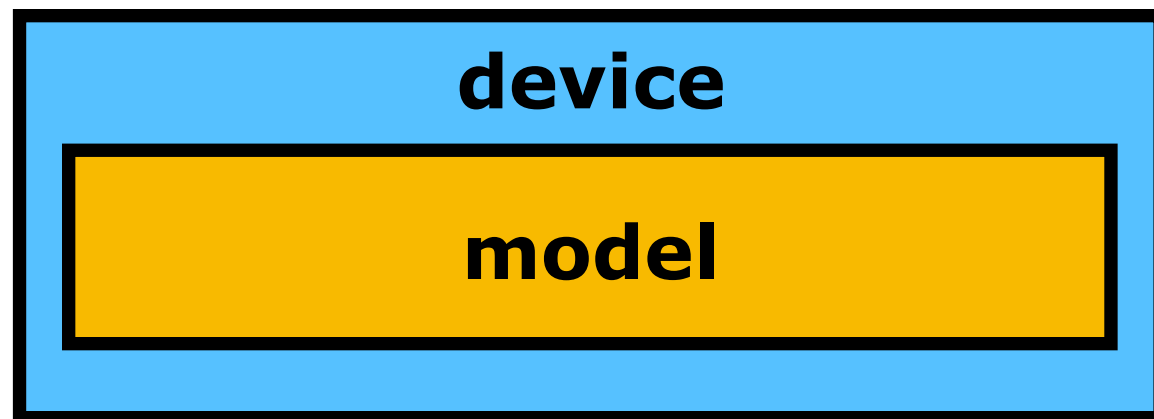
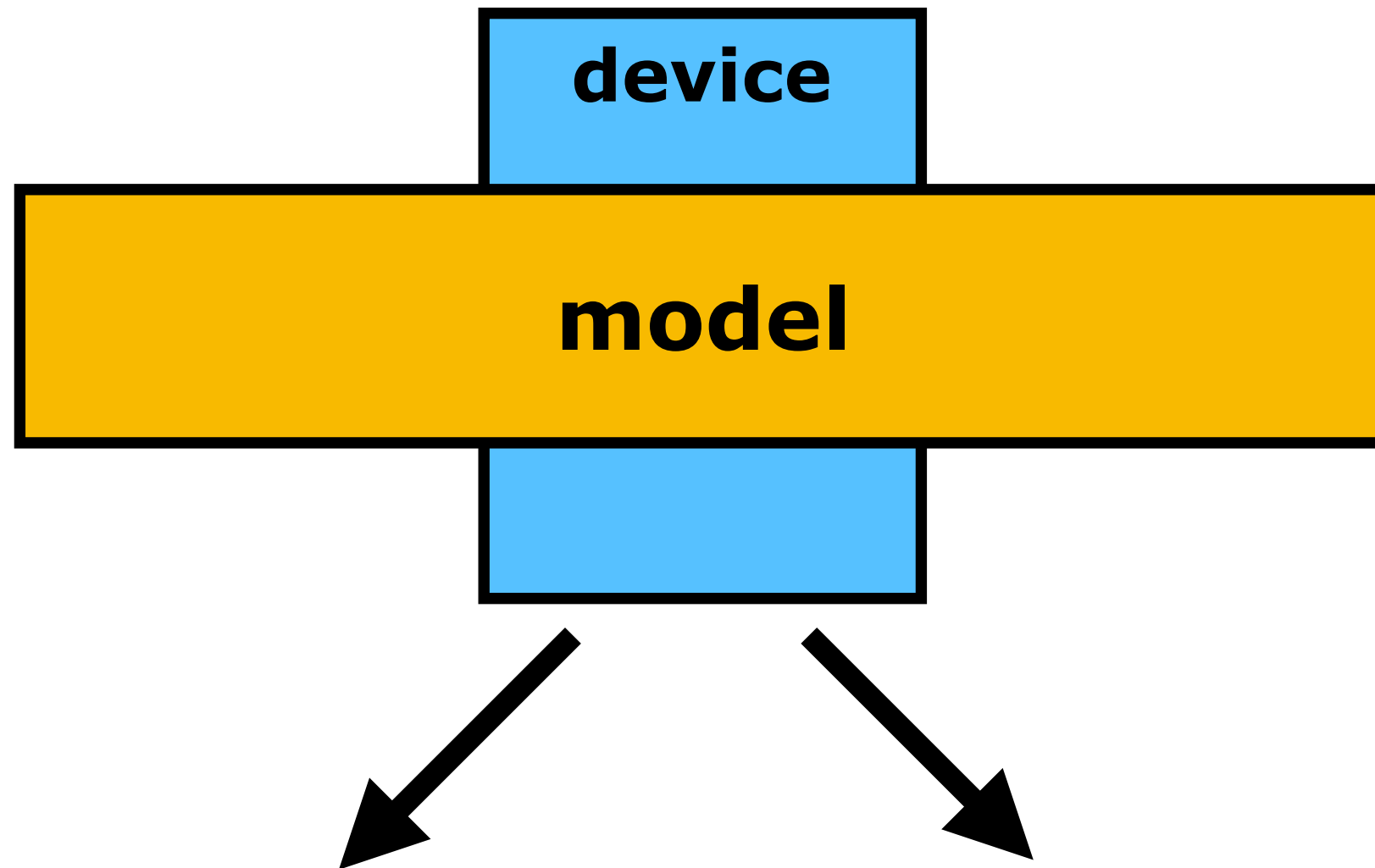
model



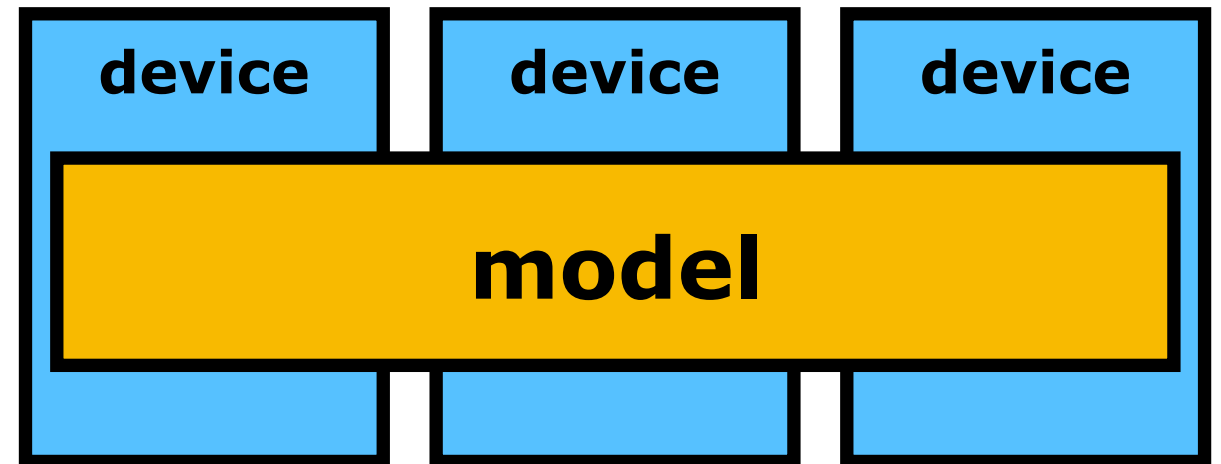
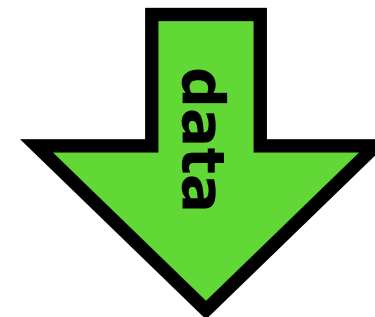
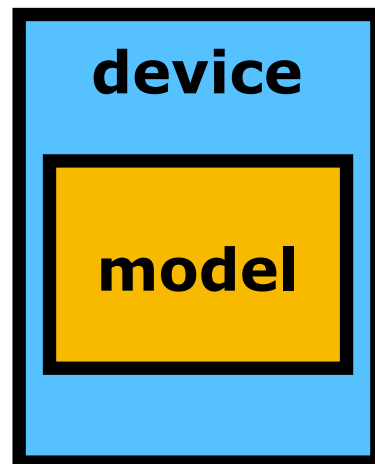
device

model

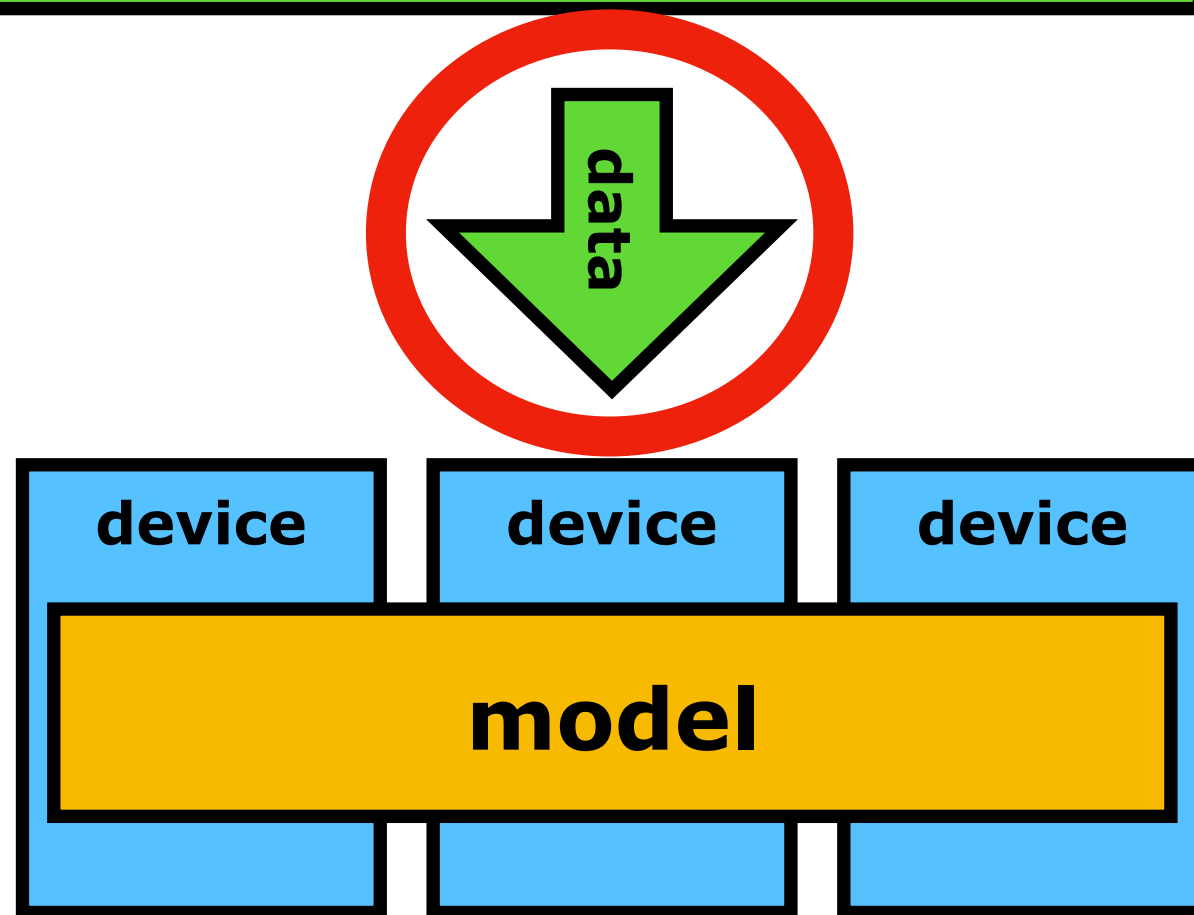
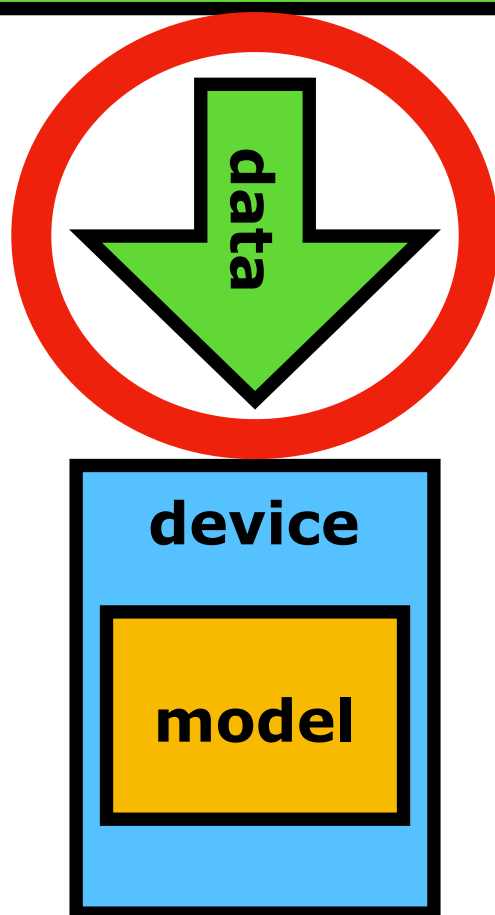




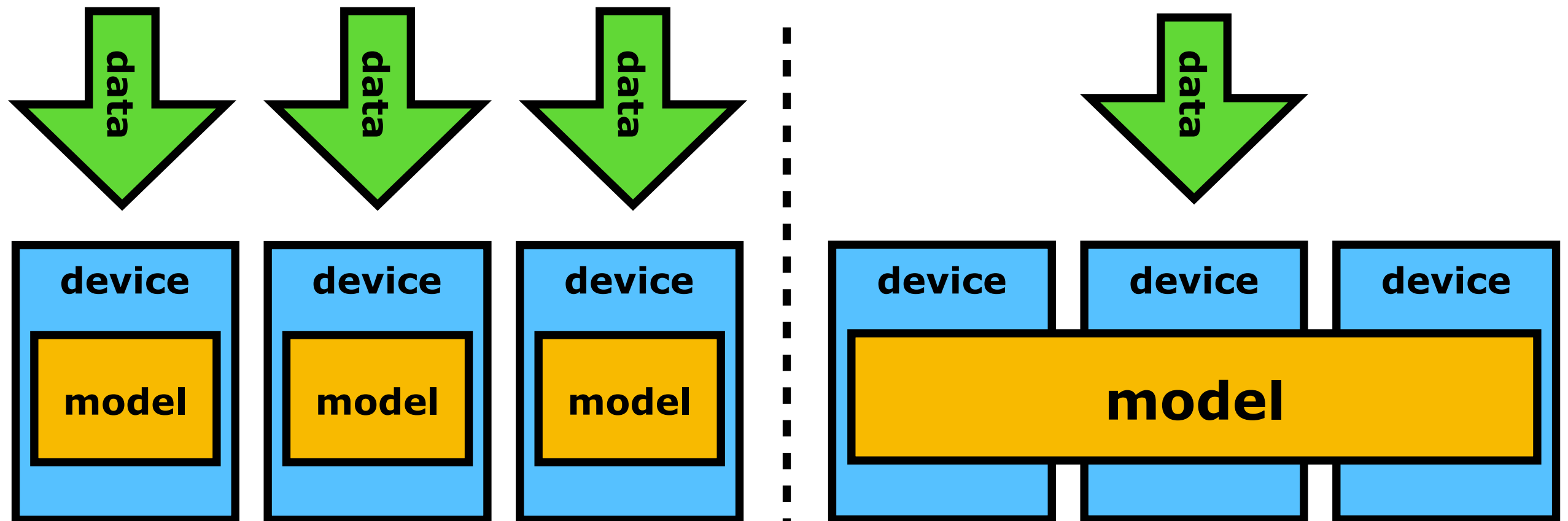
some training data



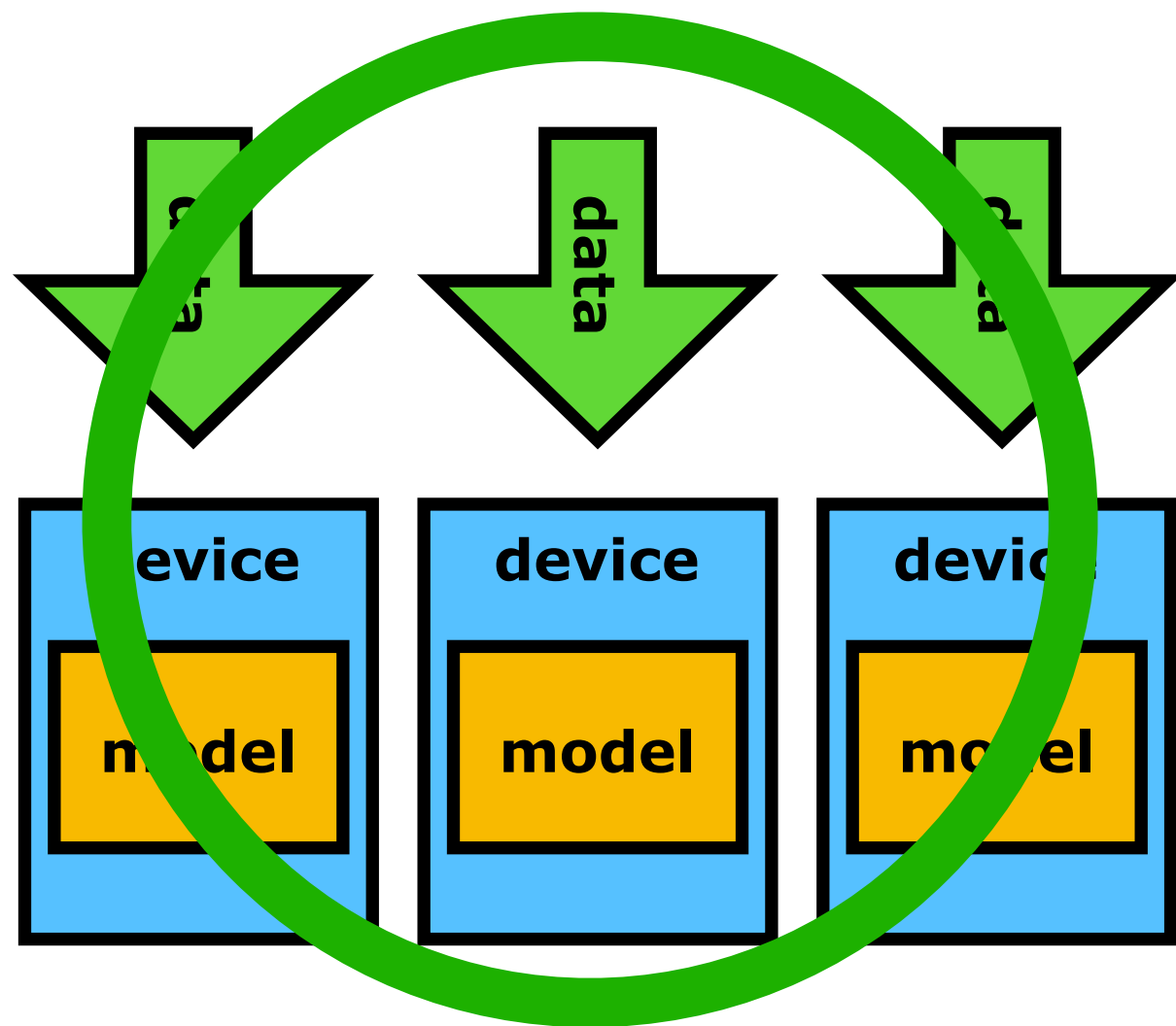
a lot of training data



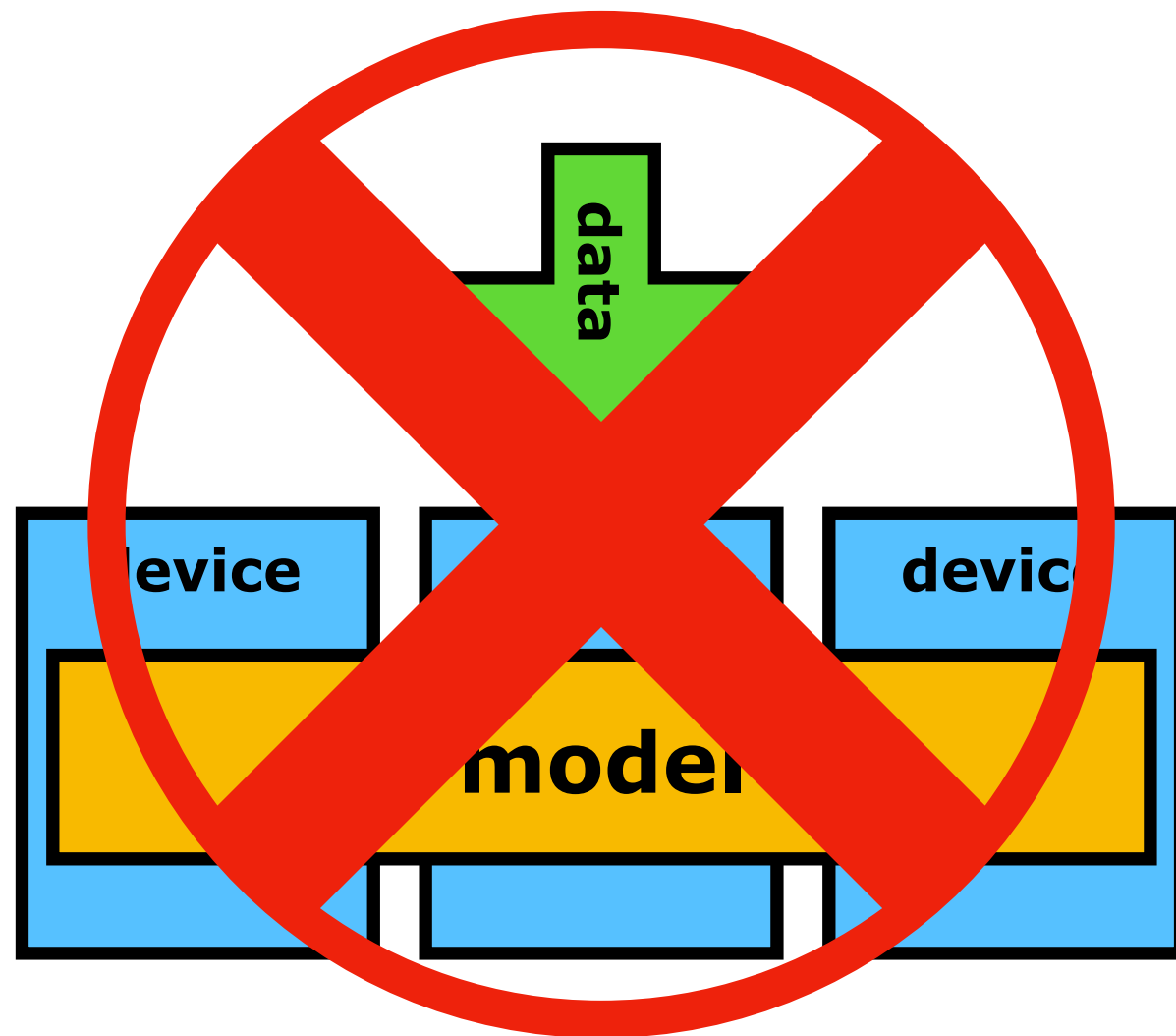
a lot of training data



Data Parallelization

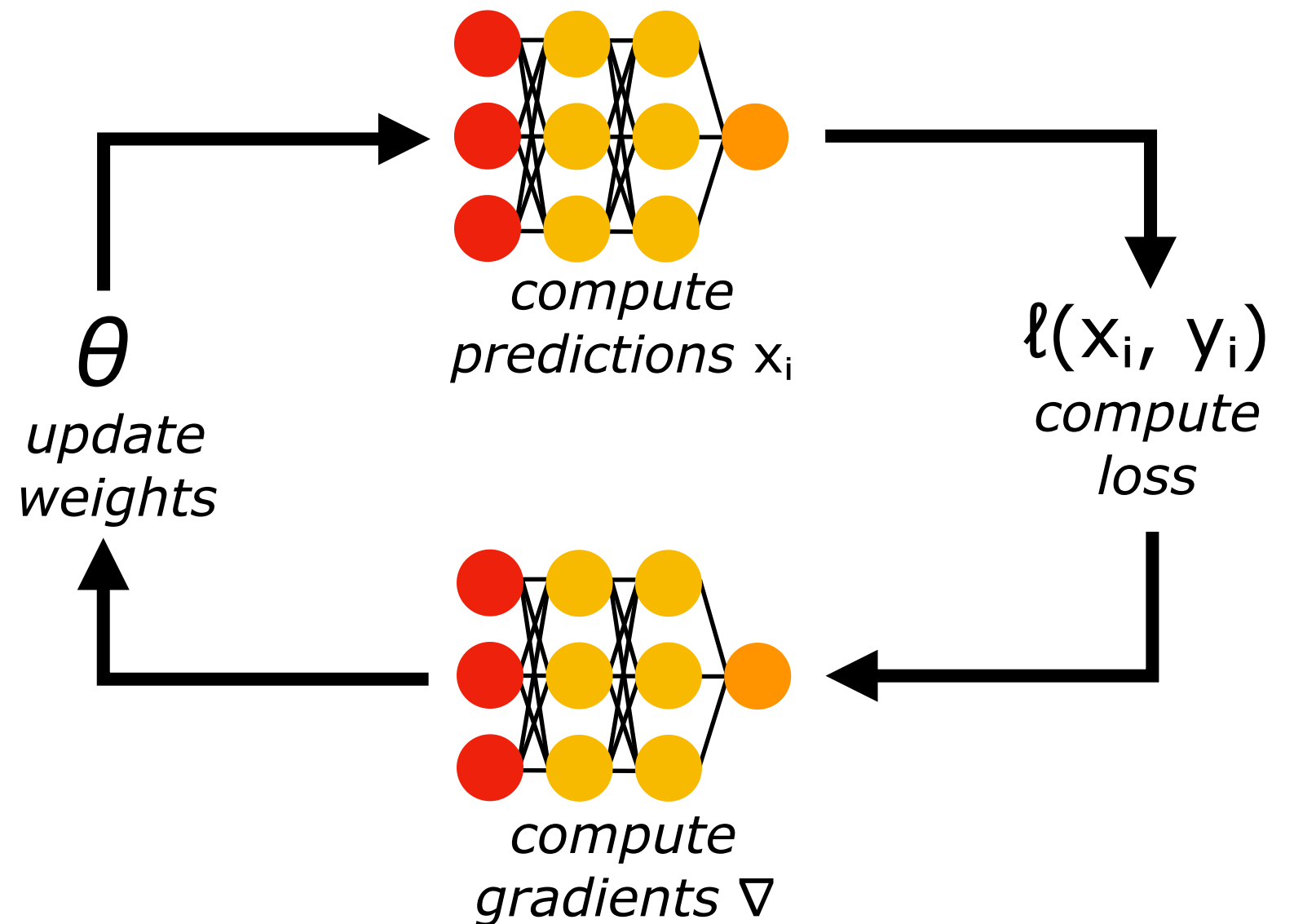


Model Parallelization

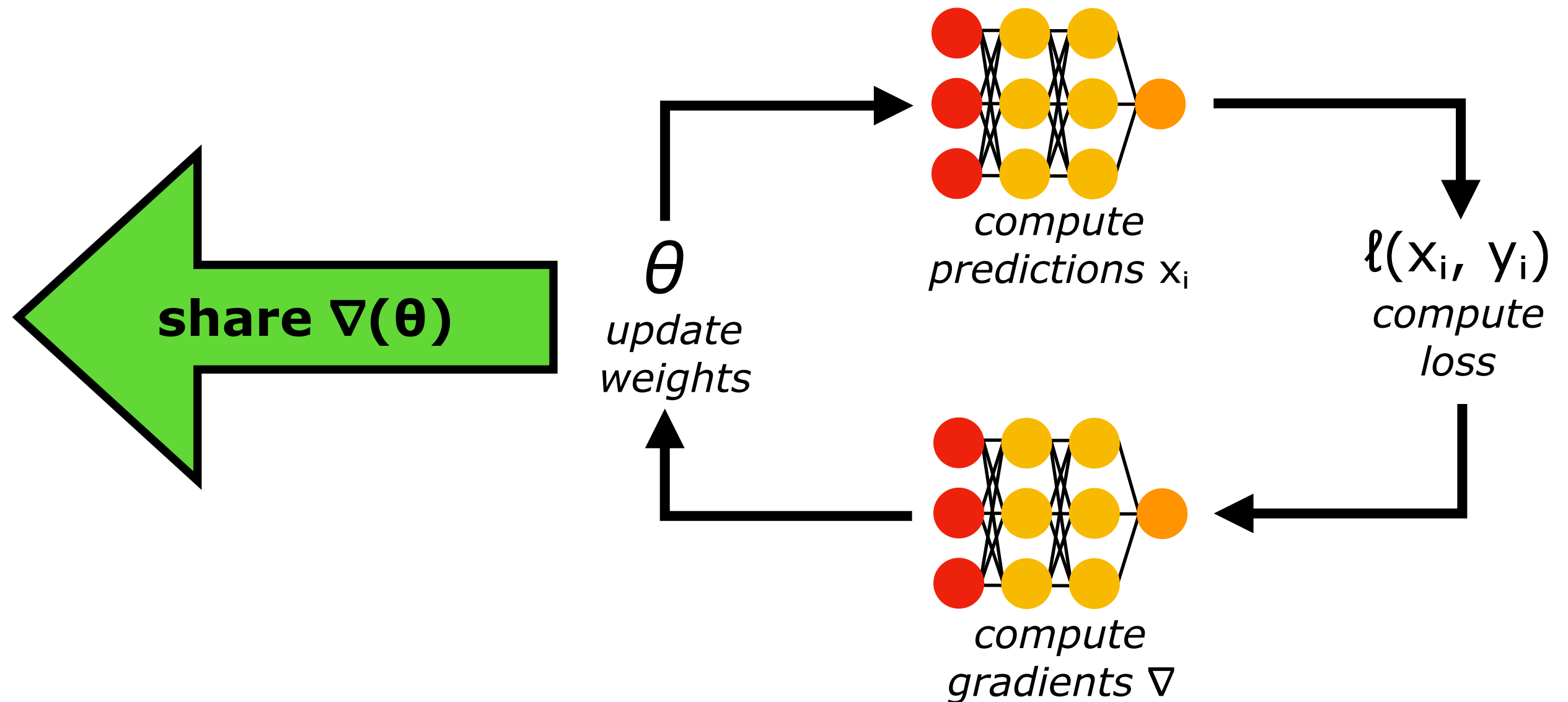


The Brief History of Distributed Learning

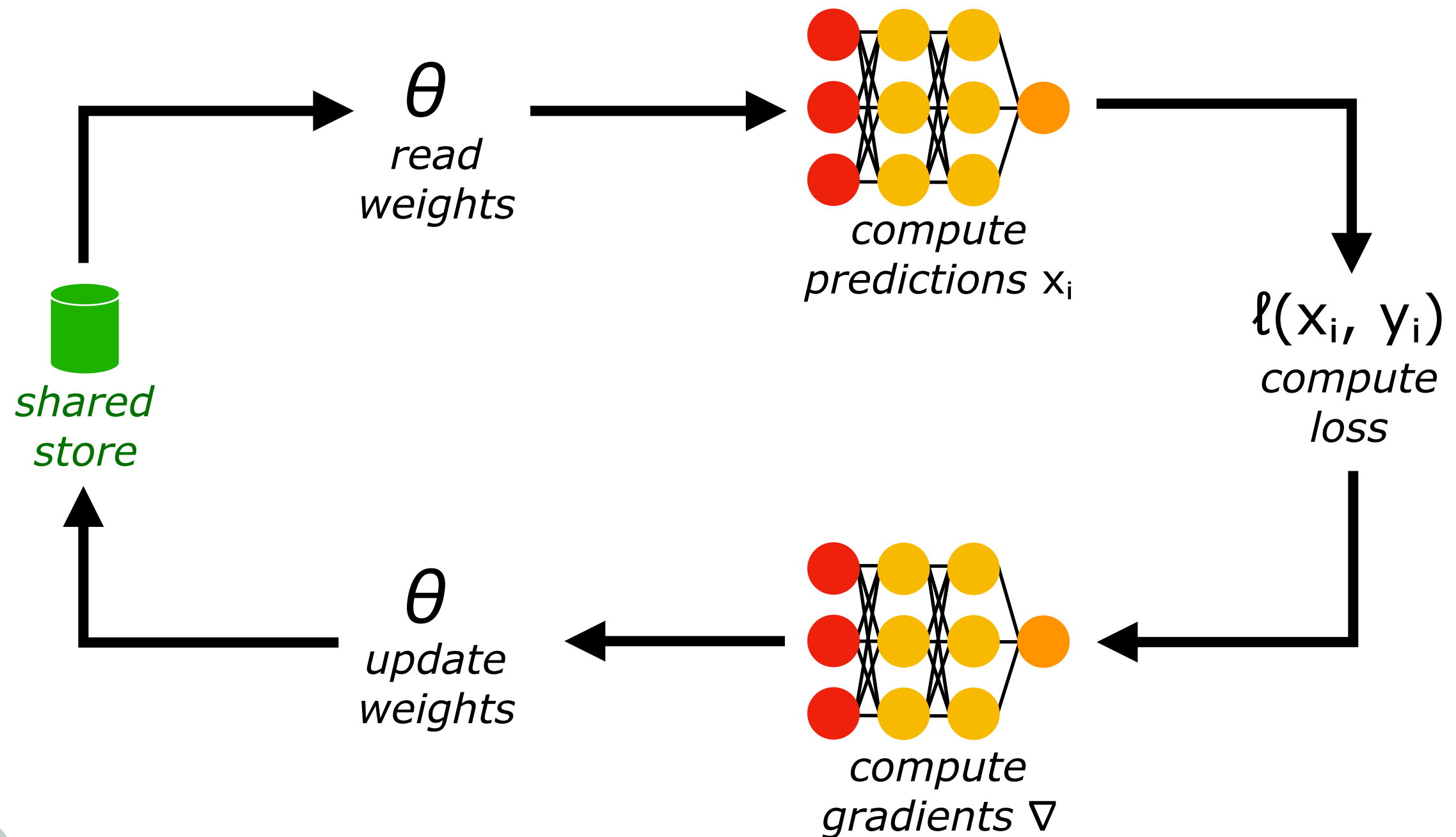
Shared Training Loop



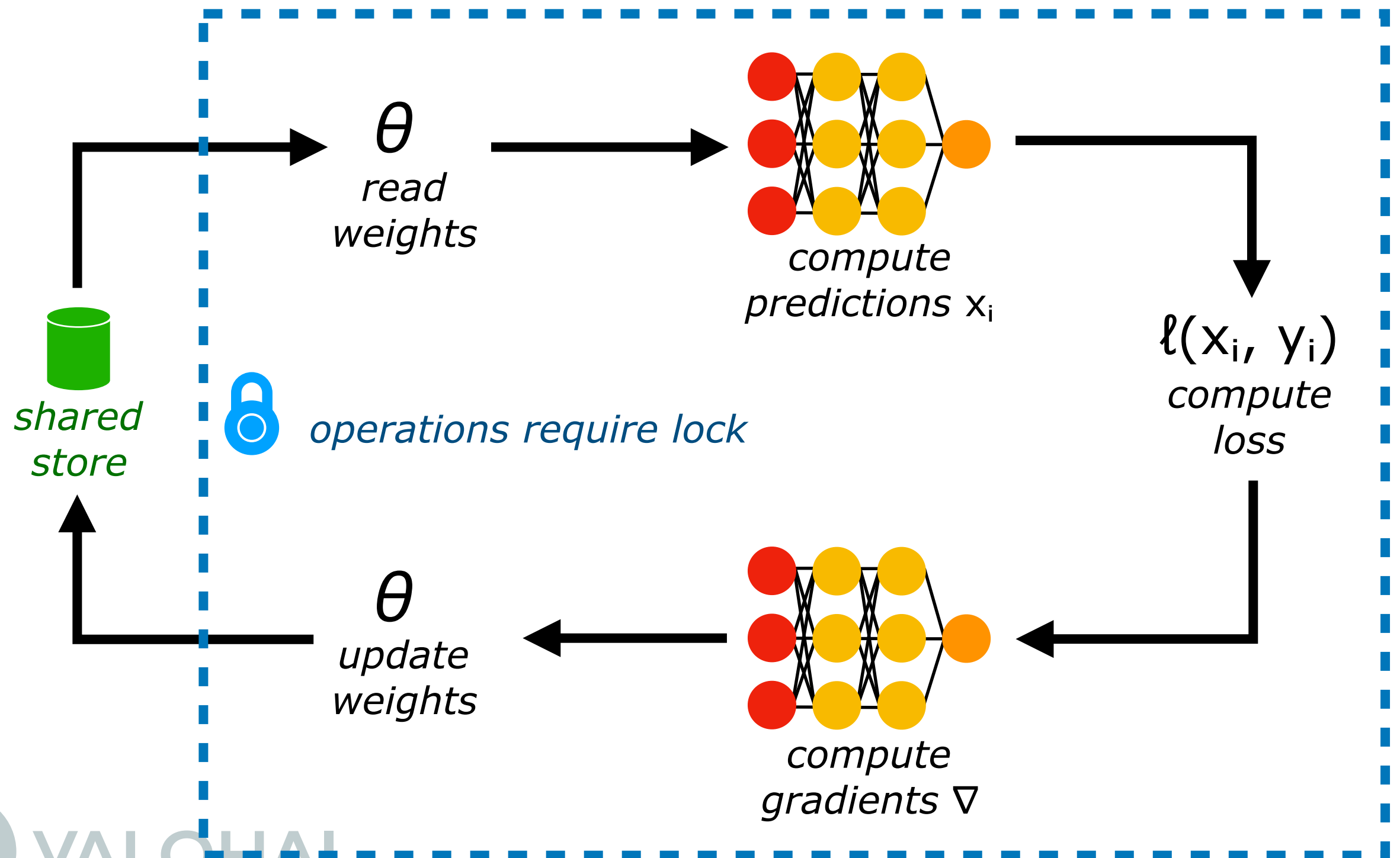
Shared Training Loop



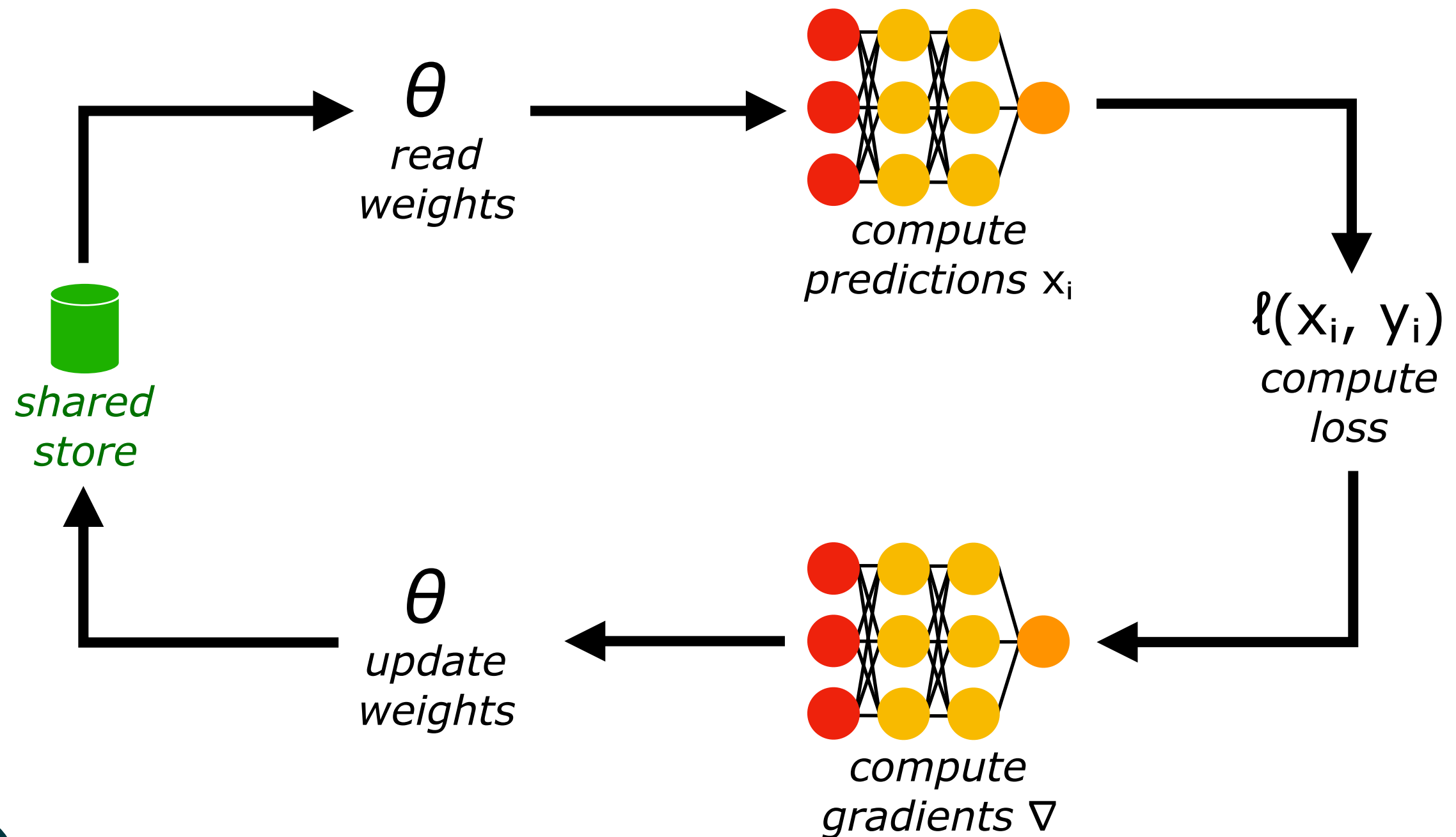
Shared Training Loop



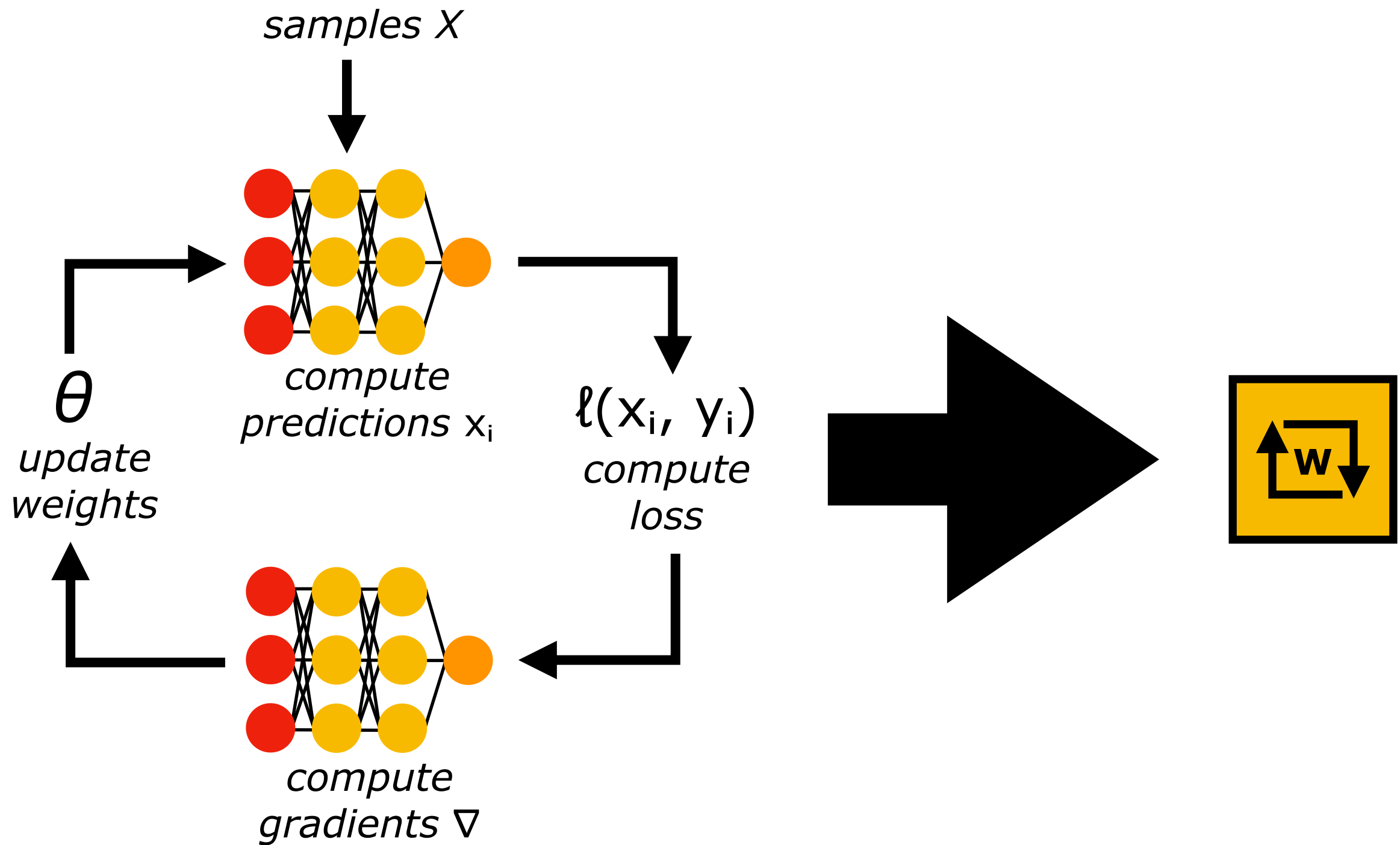
Locking Updates



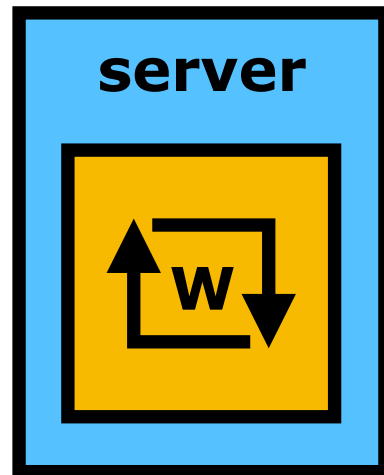
Lock-free Updates (The Hogwild)



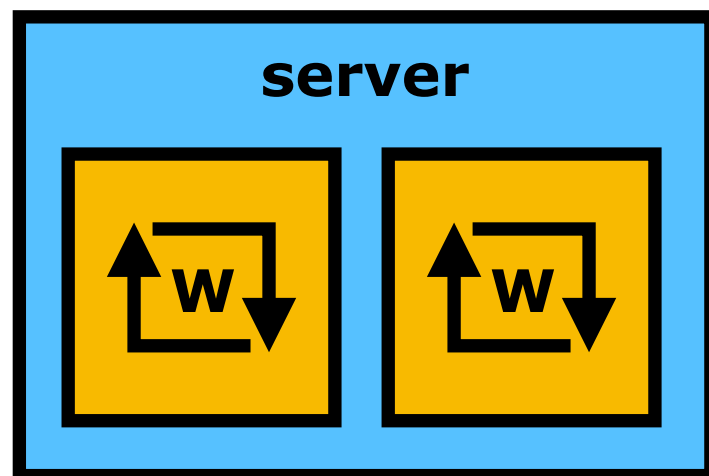
Notation



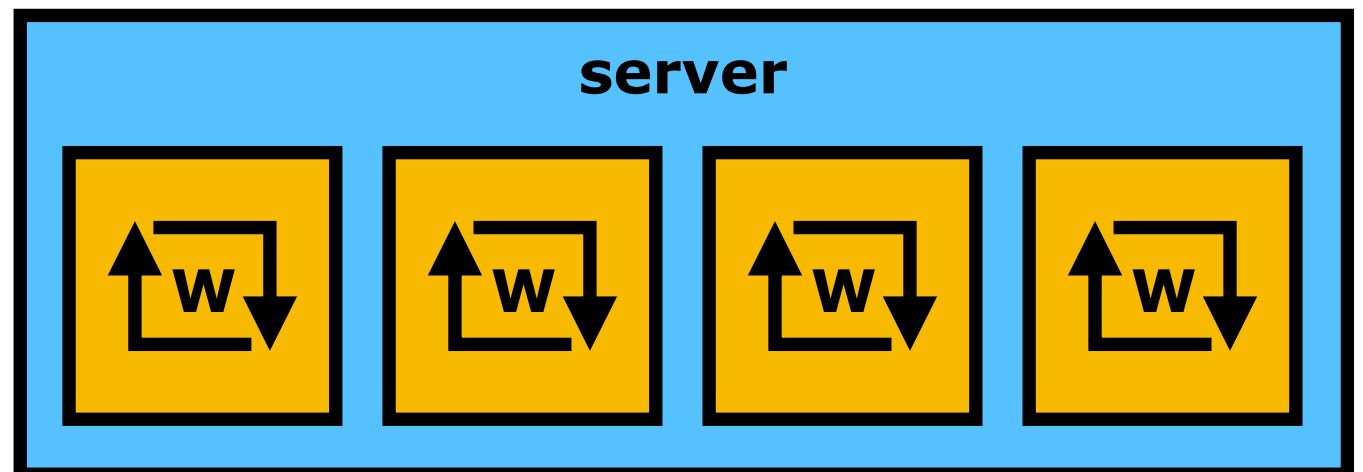
Notation

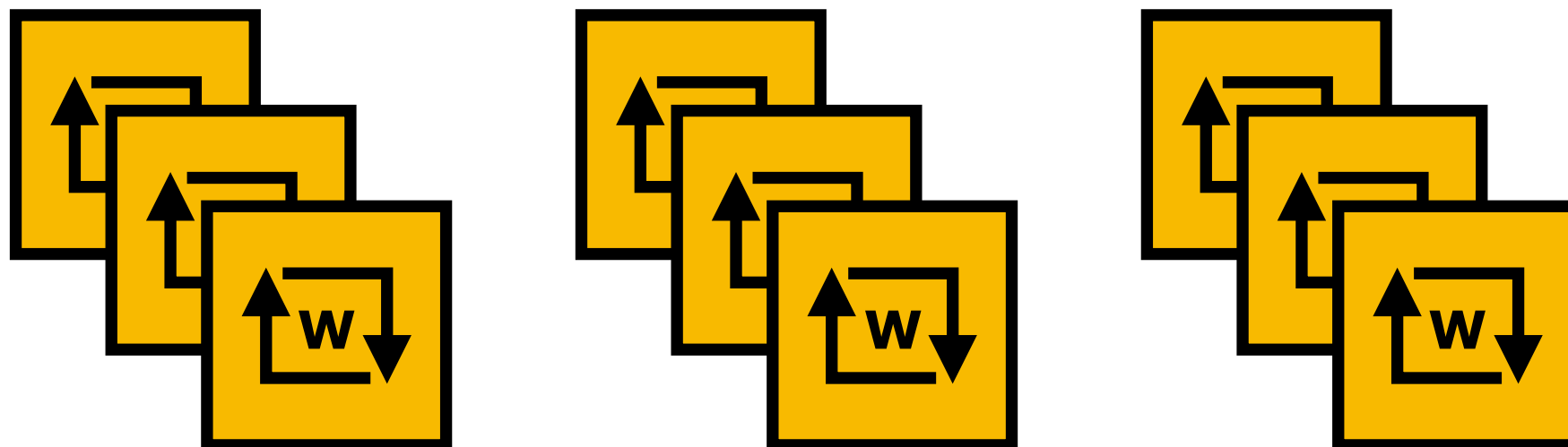
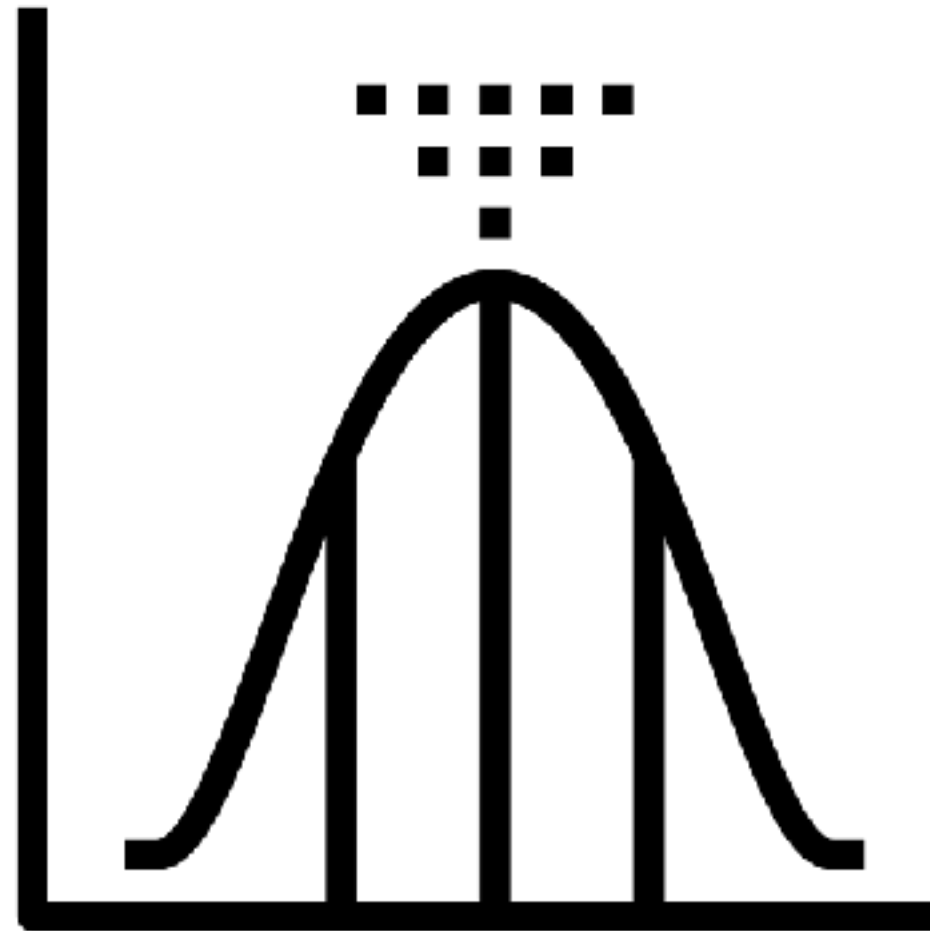


Usually
worker count == number of GPUs

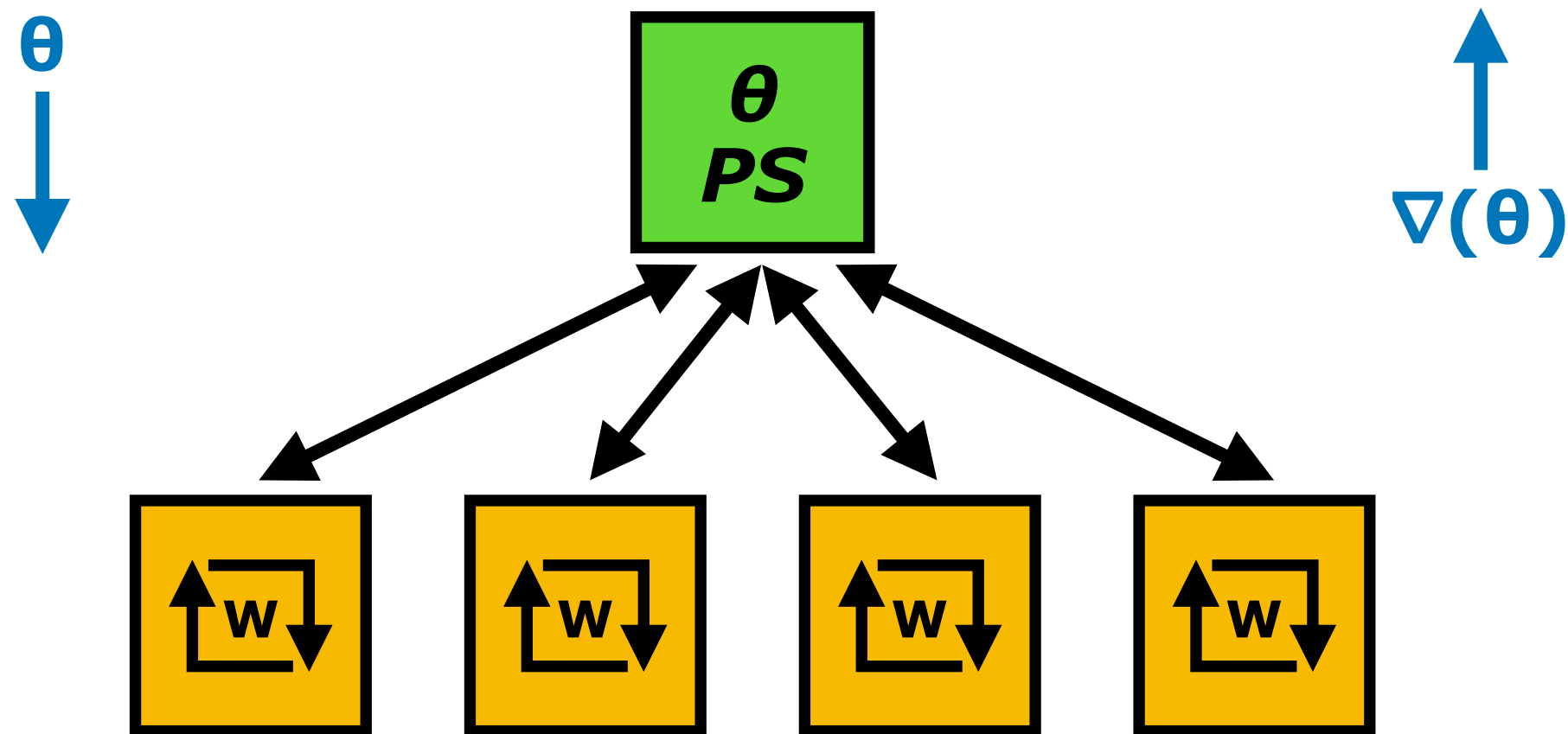


Each server has
limited network
bandwidth.

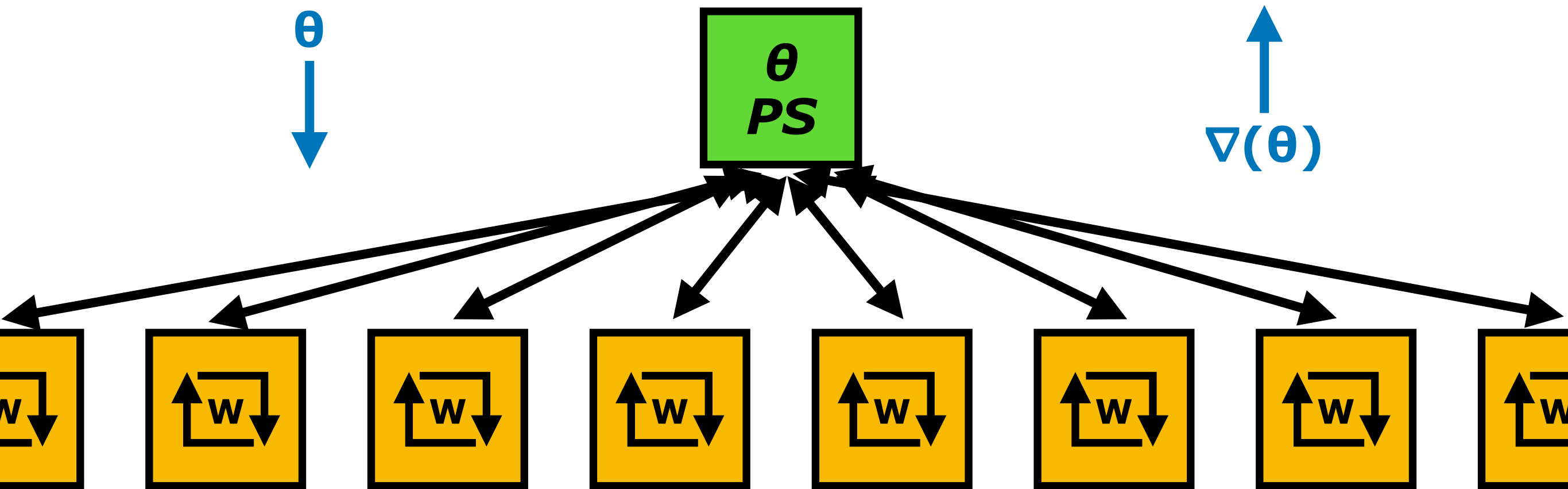




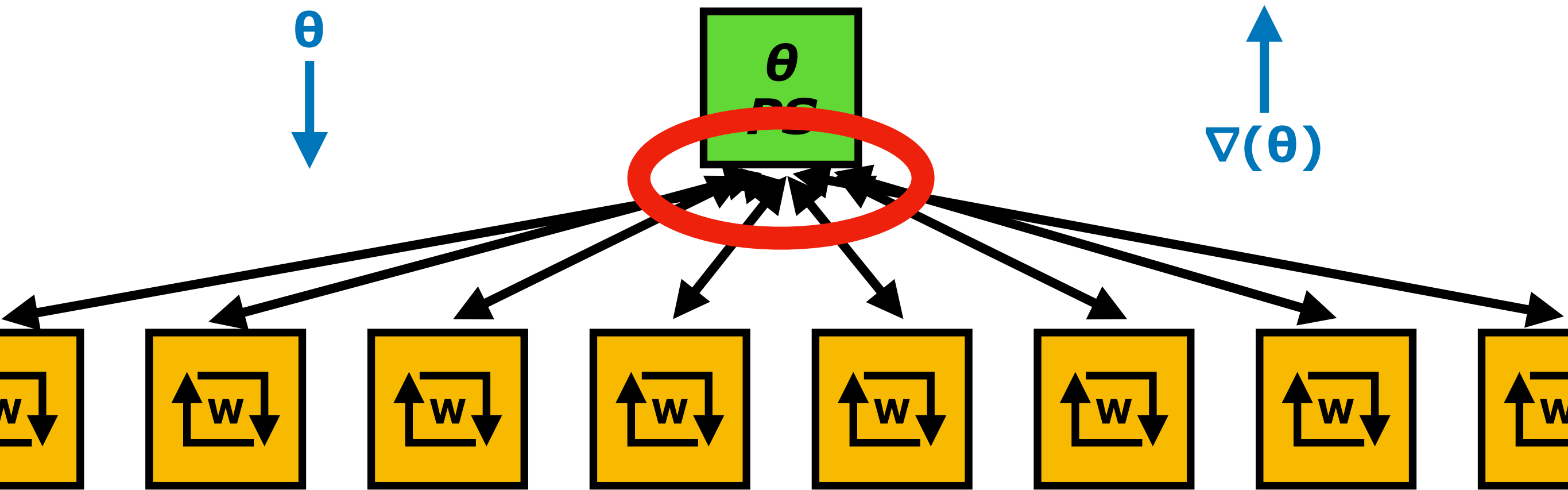
Asynchronous Stochastic Gradient Descent (Parameter Server)



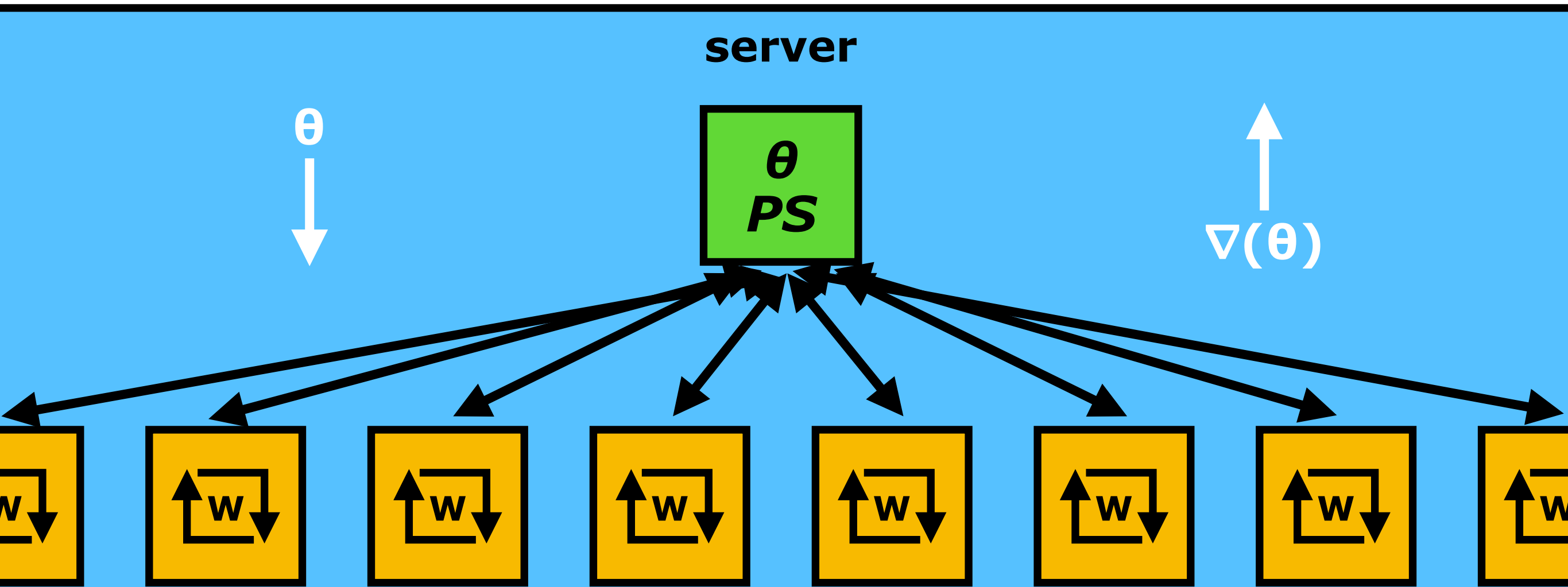
Asynchronous Stochastic Gradient Descent (Parameter Server)



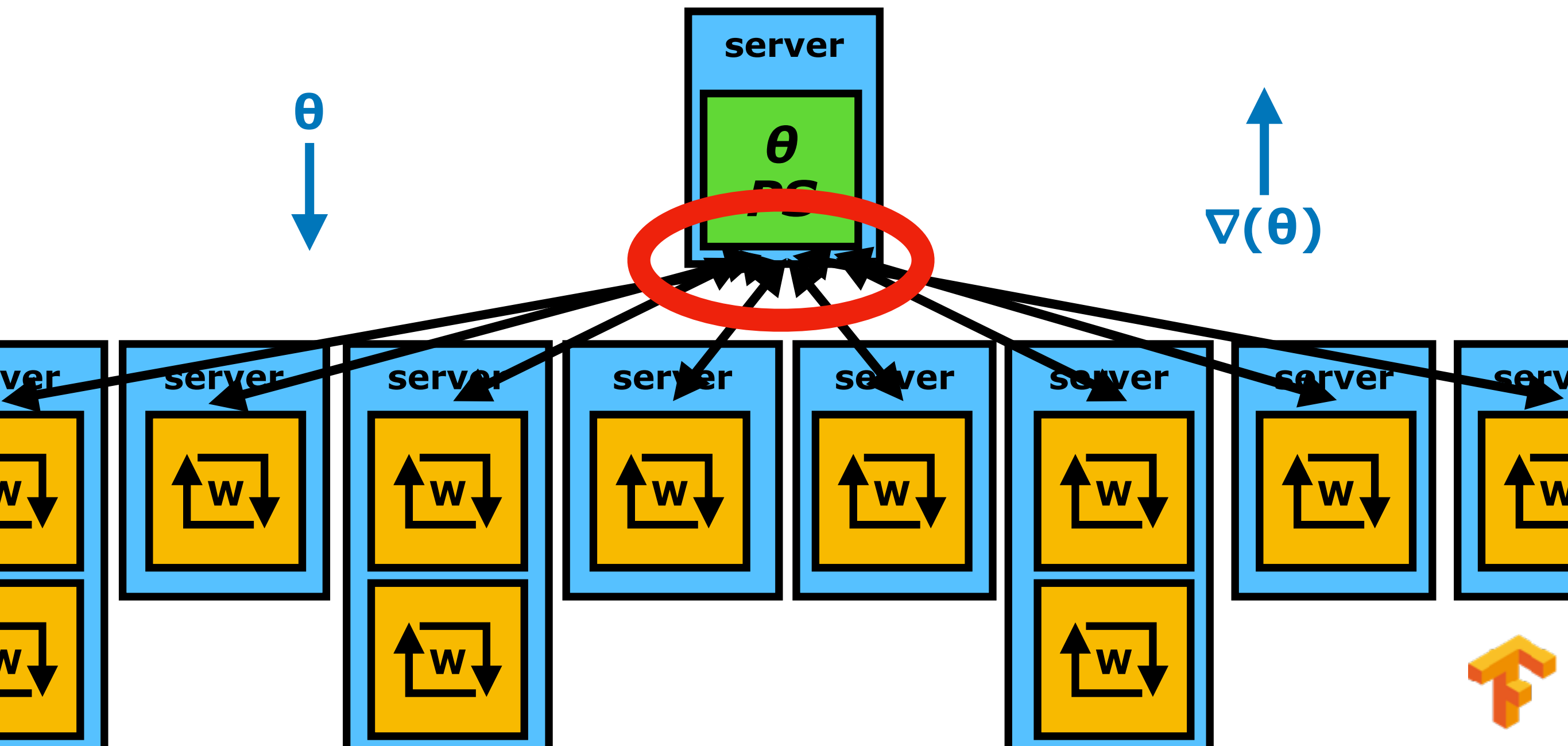
Asynchronous Stochastic Gradient Descent (Parameter Server)



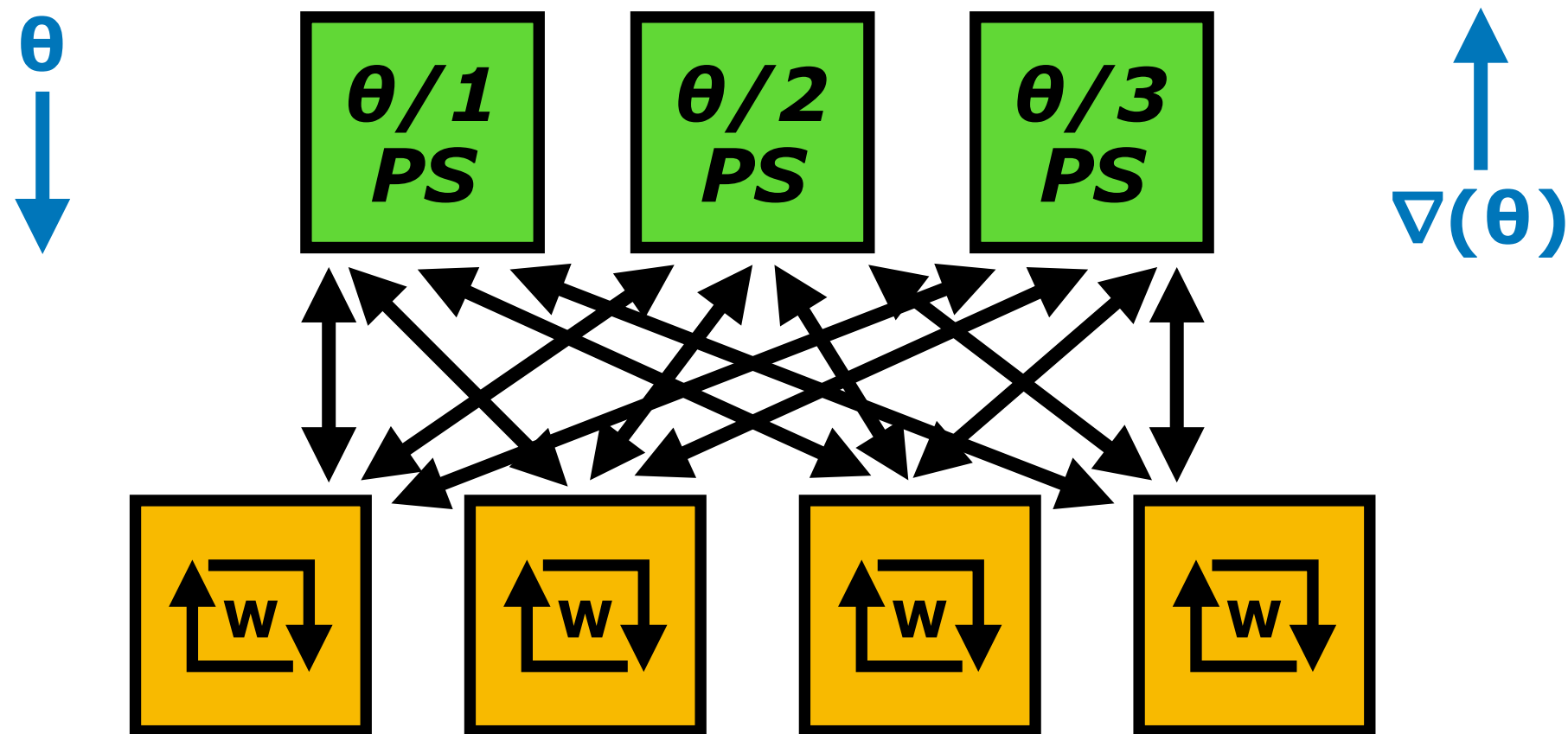
Asynchronous Stochastic Gradient Descent (Parameter Server)



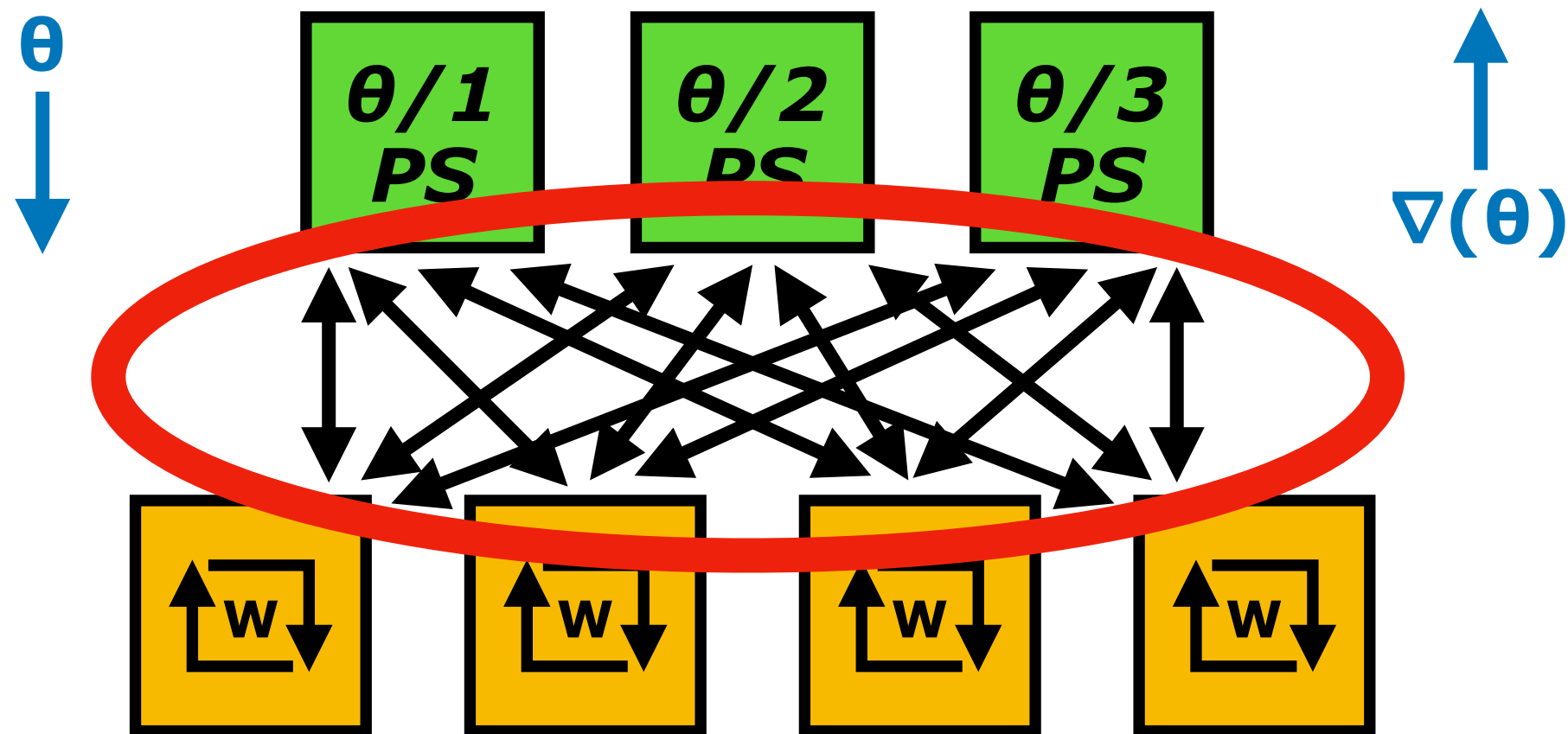
Asynchronous Stochastic Gradient Descent (Parameter Server)



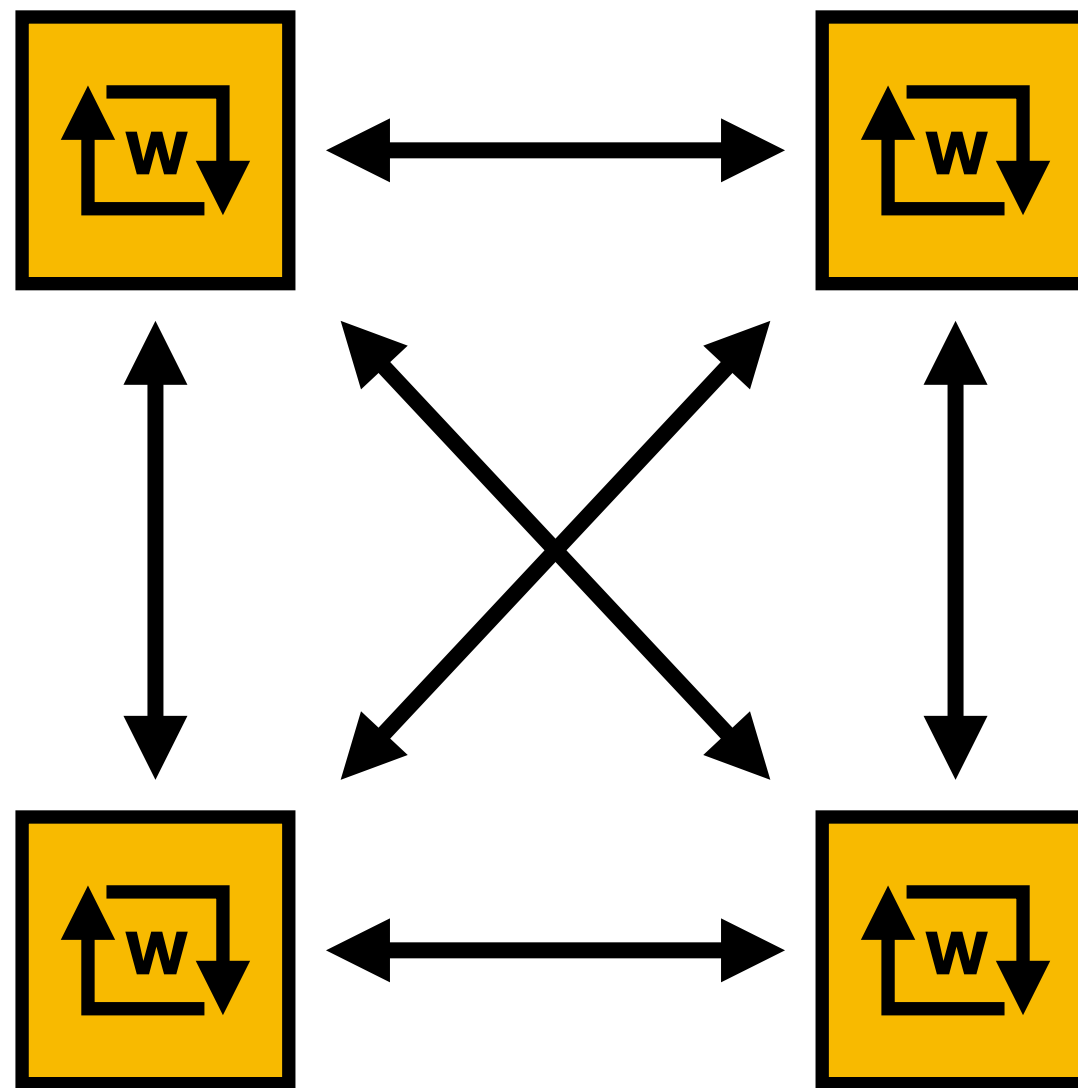
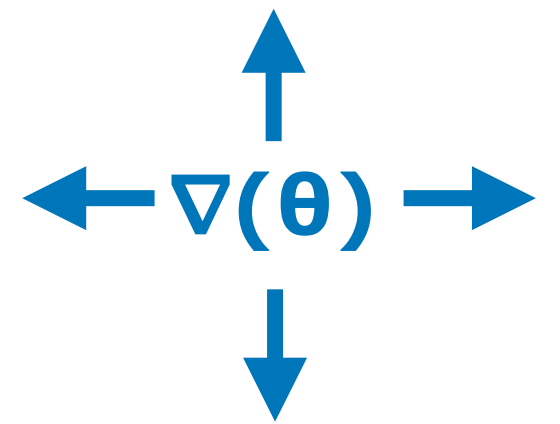
Asynchronous Stochastic Gradient Descent (Parameter Server Cluster)



Asynchronous Stochastic Gradient Descent (Parameter Server Cluster)

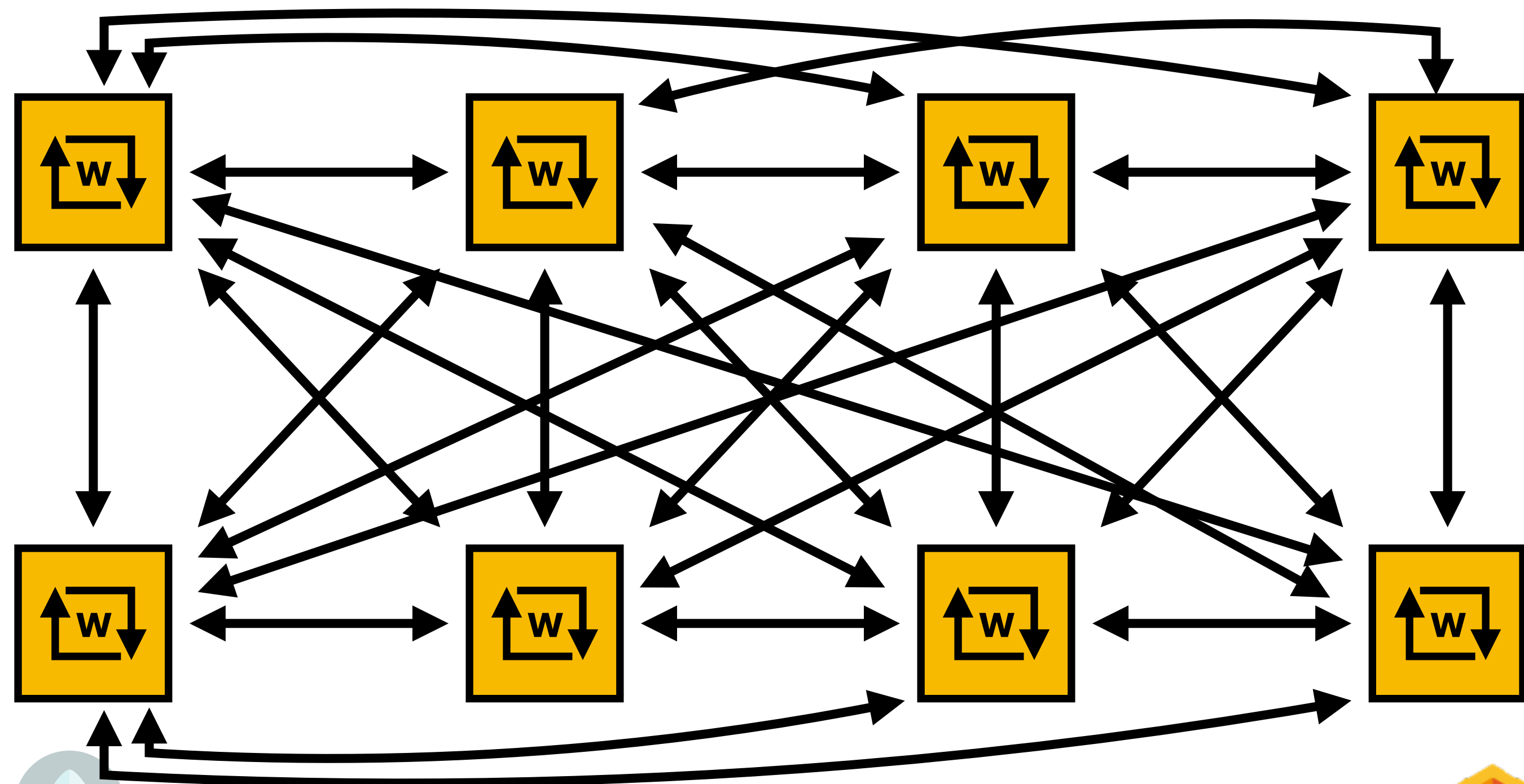


Synchronous Allreduce



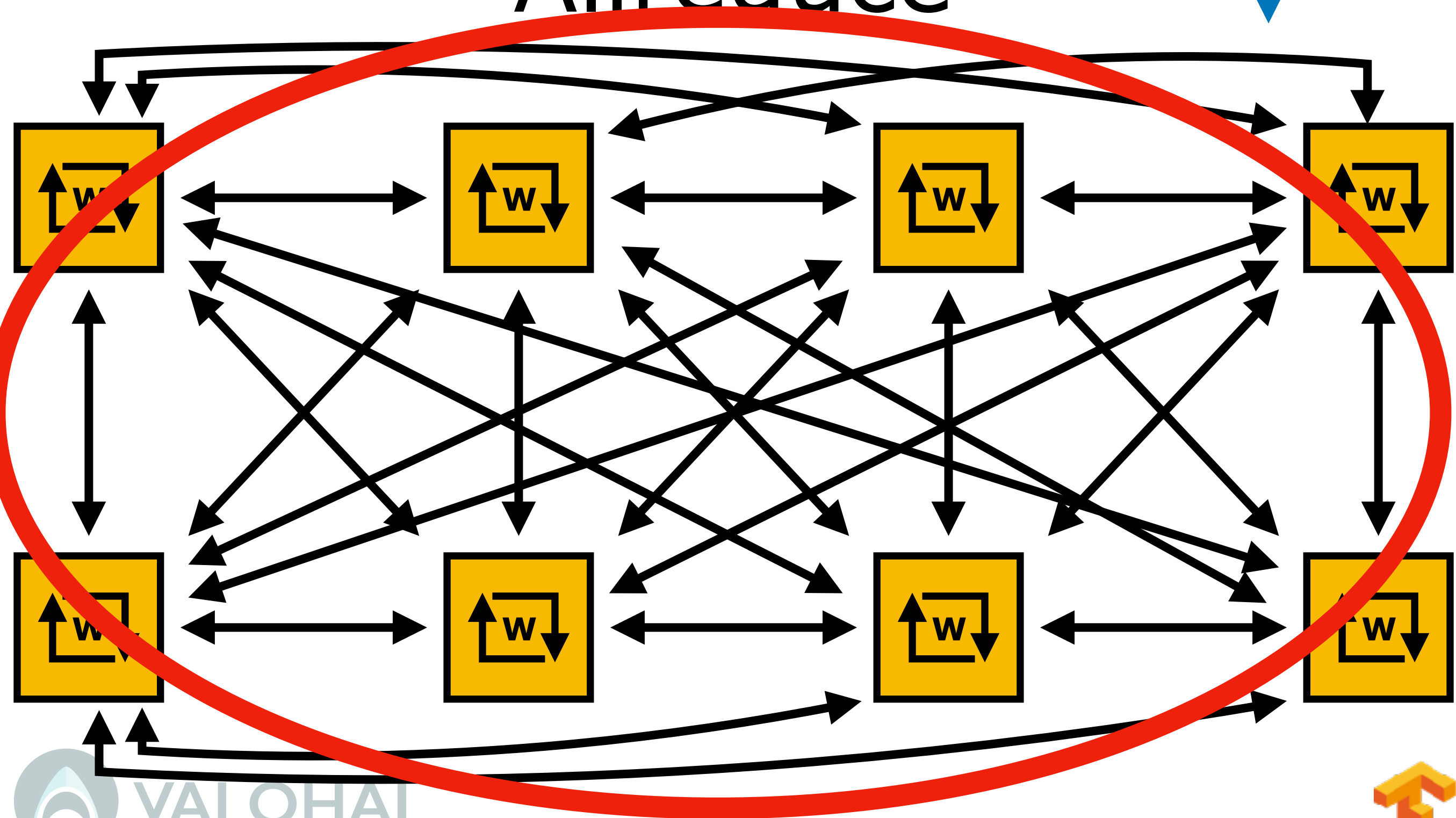
Synchronous Allreduce

$$\nabla(\theta)$$



Synchronous Allreduce

$$\nabla(\theta)$$





Butterfly
Allreduce

1

Worker 1

0

Worker 2

1

Worker 3

2

Worker 4

3

Worker 5

4

Worker 6

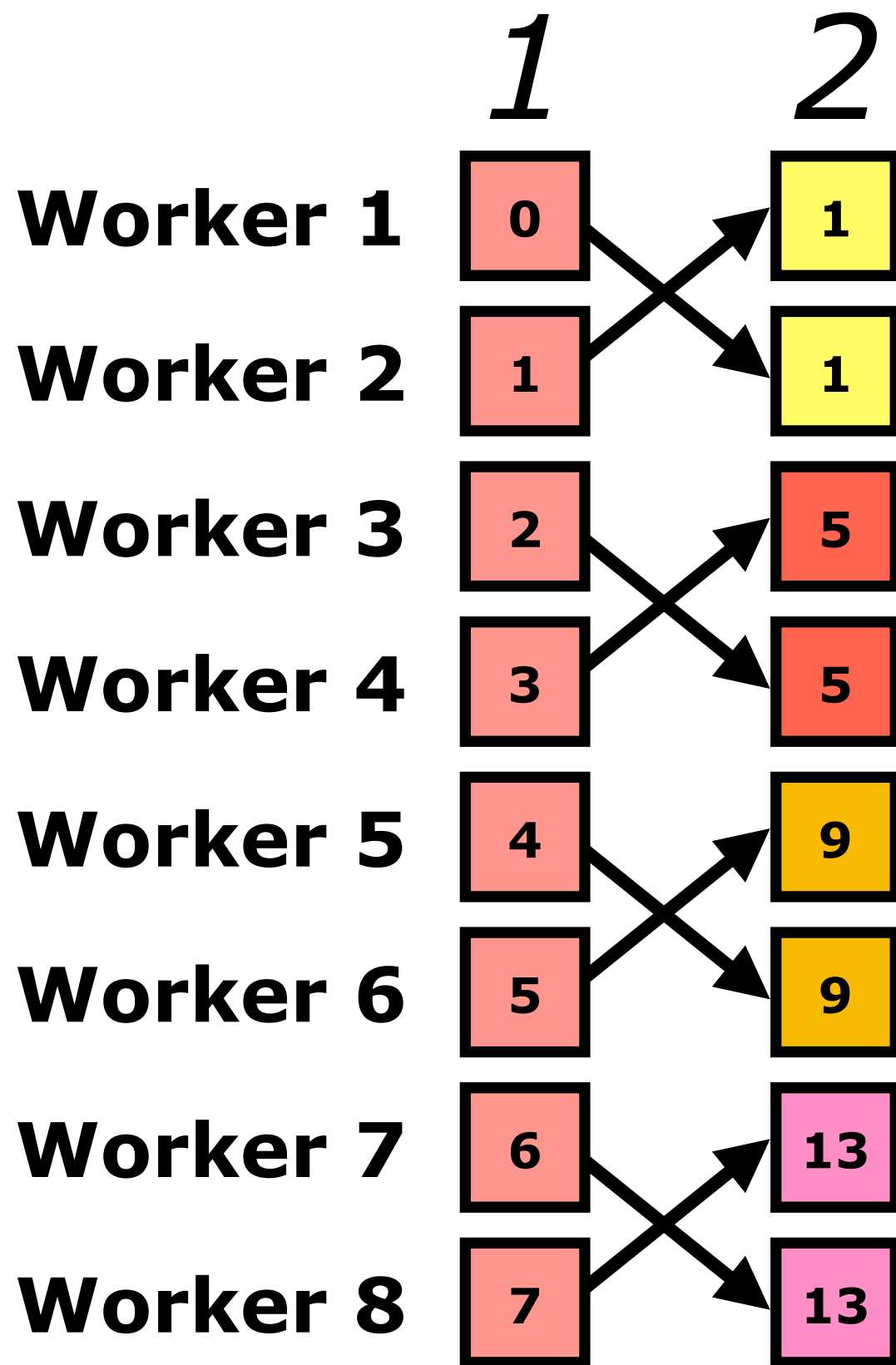
5

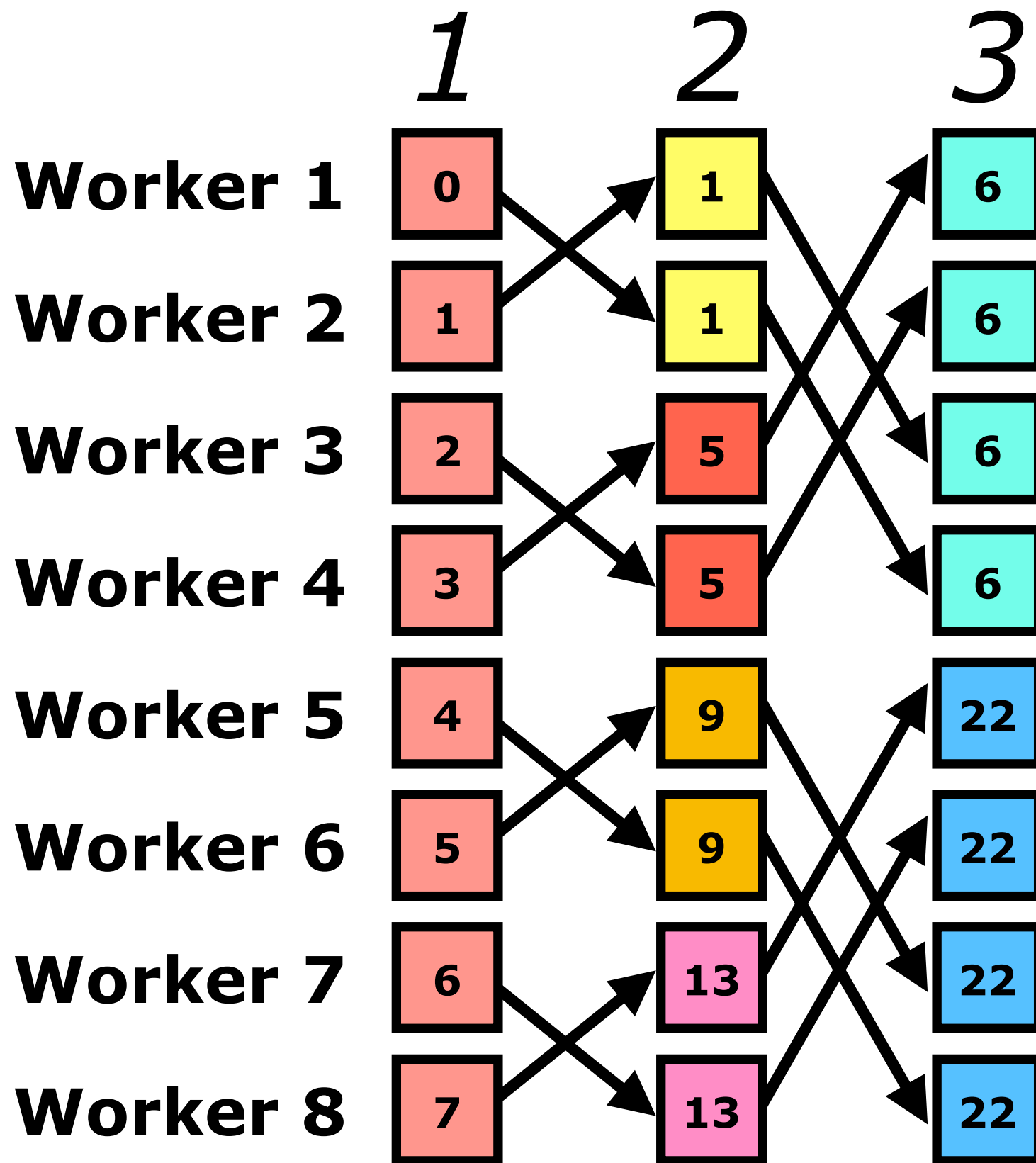
Worker 7

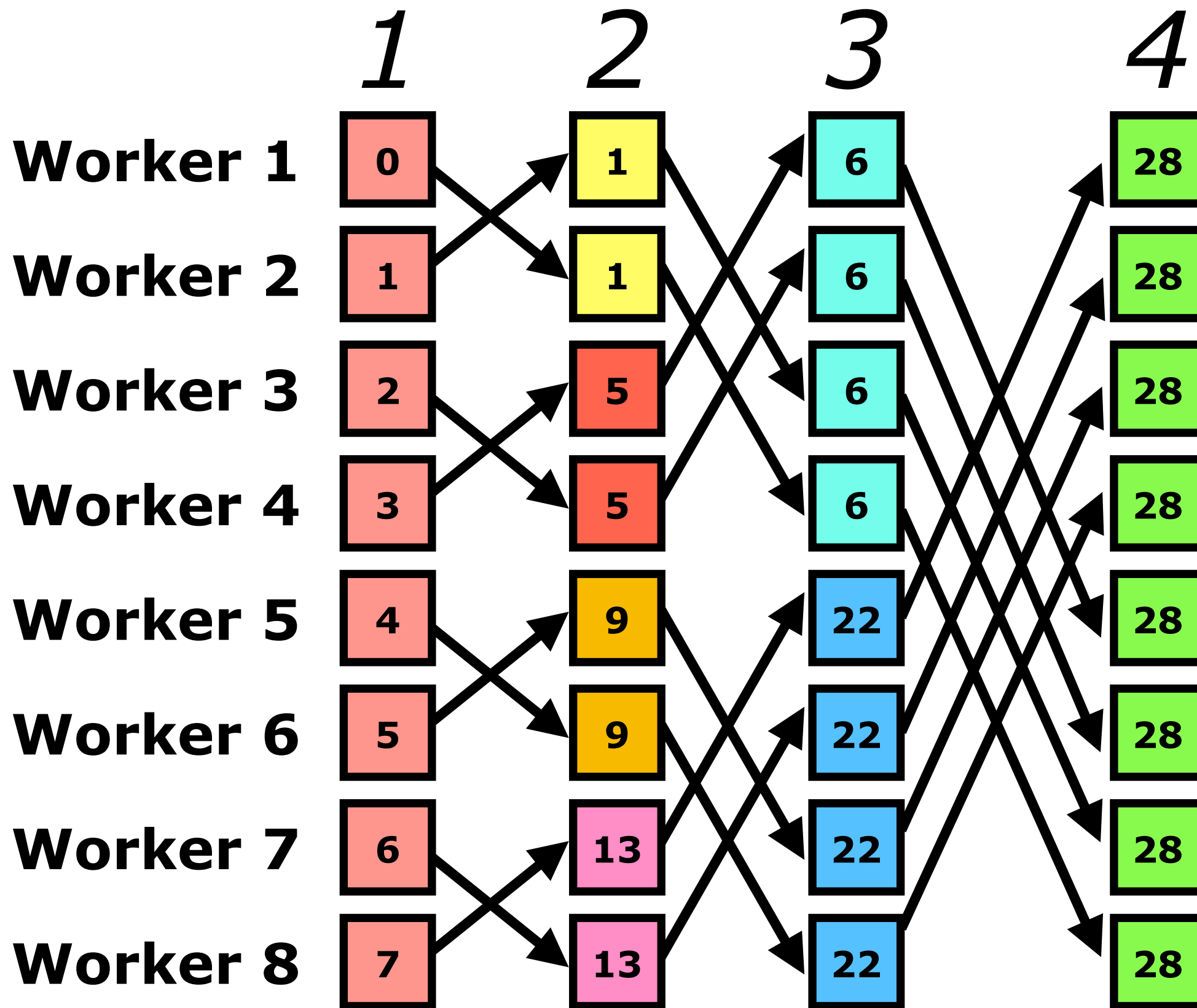
6

Worker 8

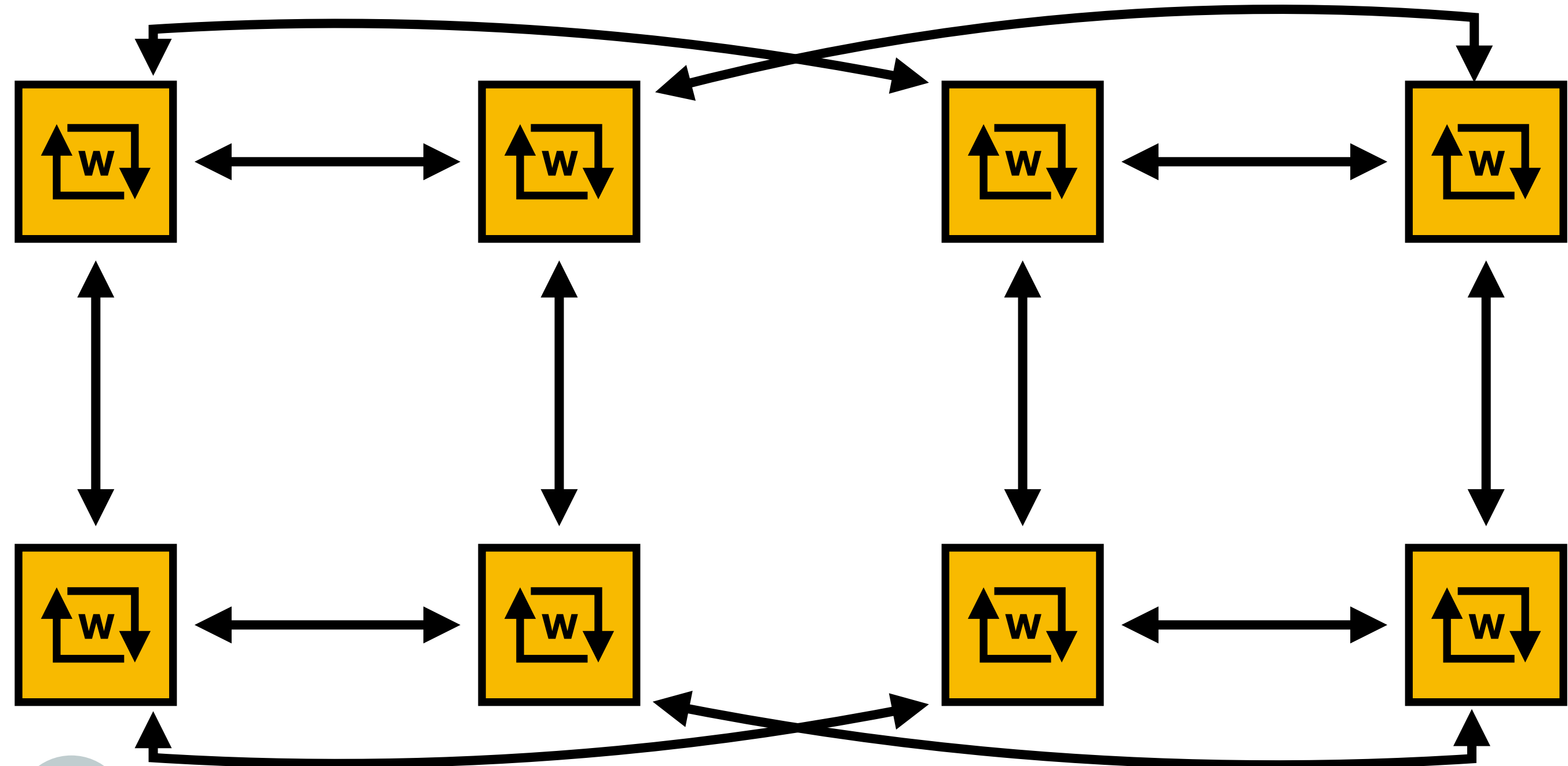
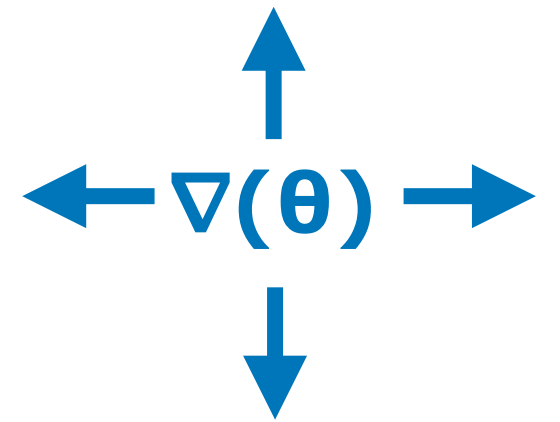
7







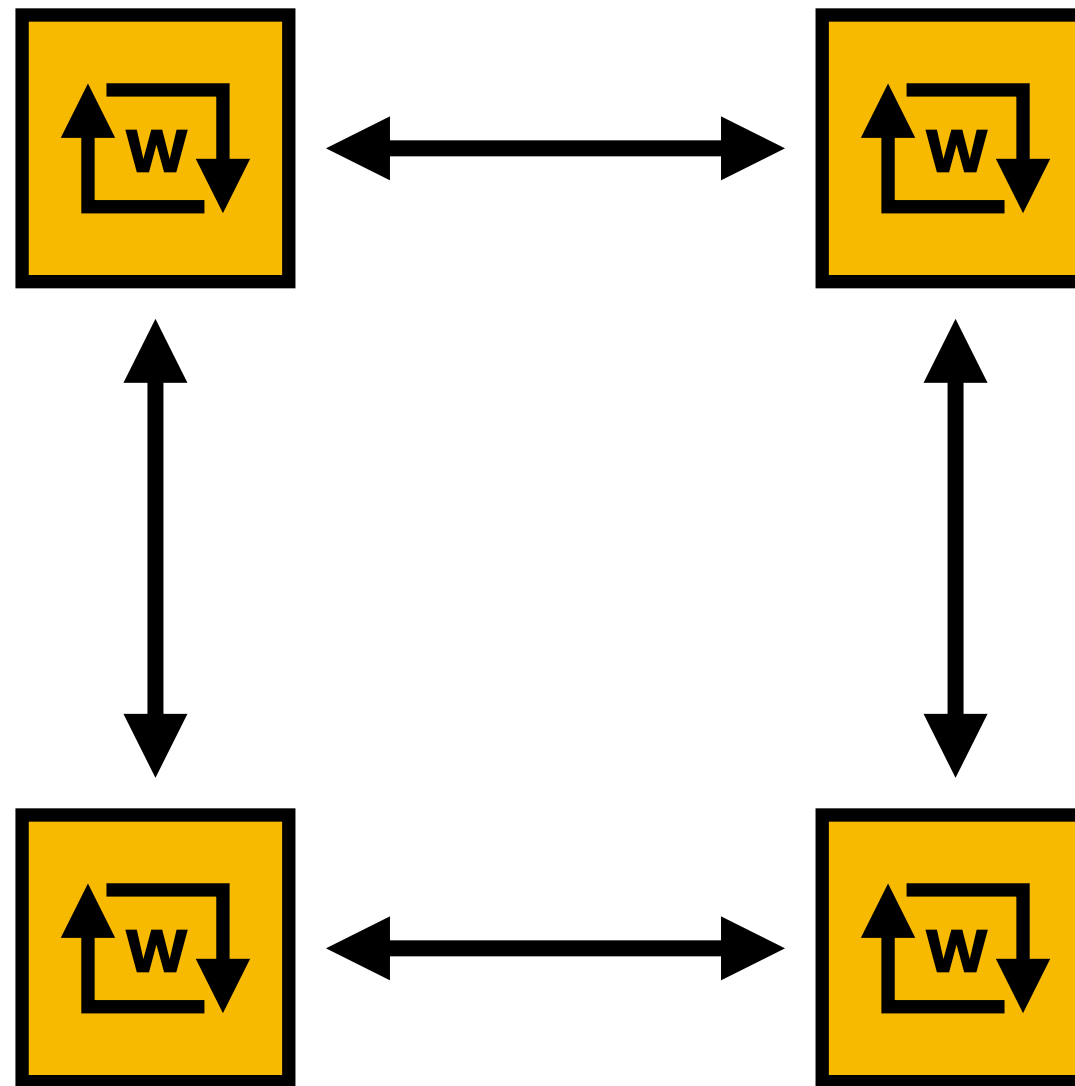
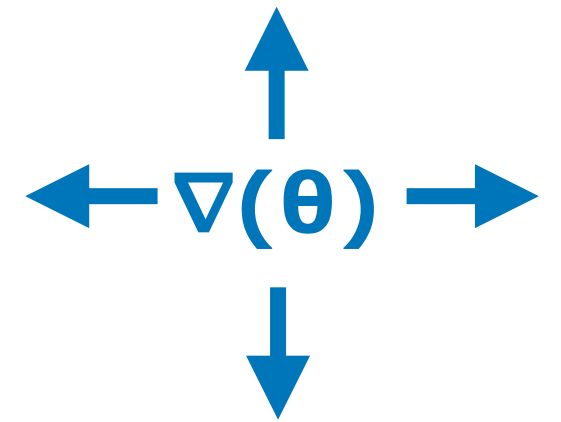
Butterfly Allreduce



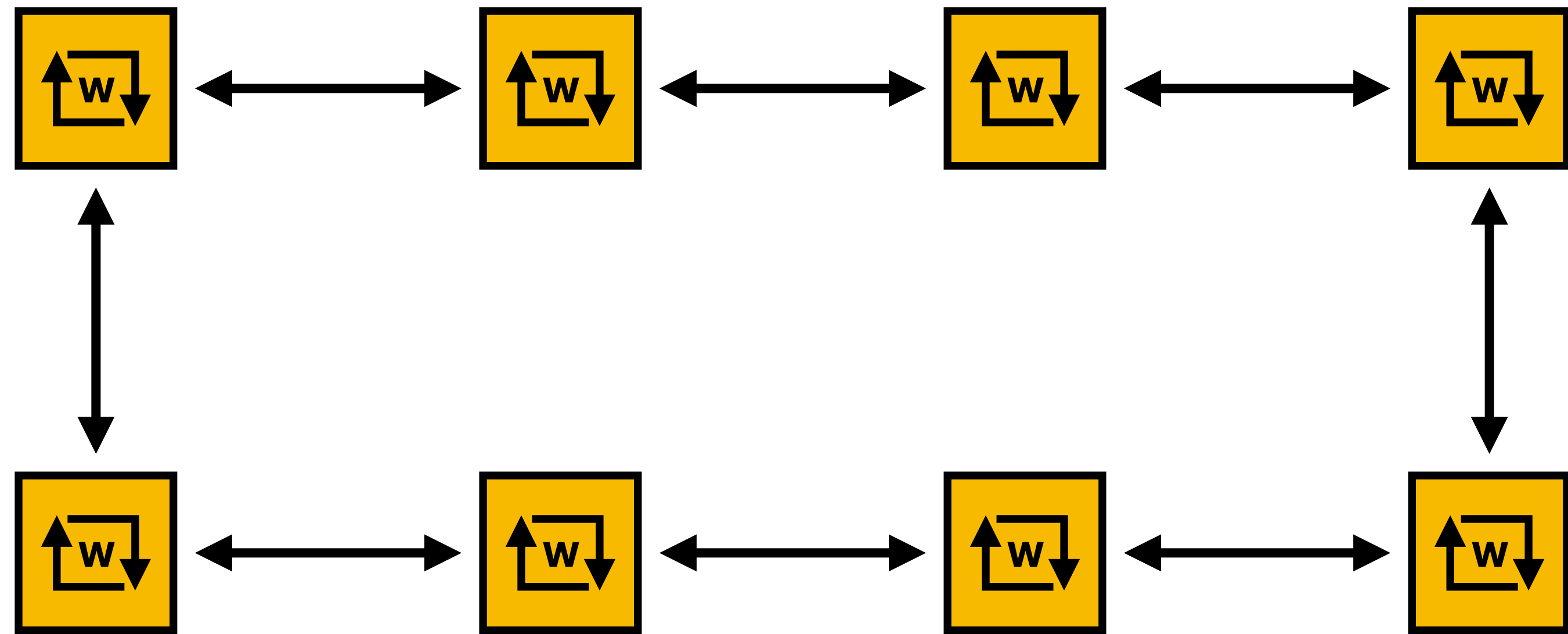
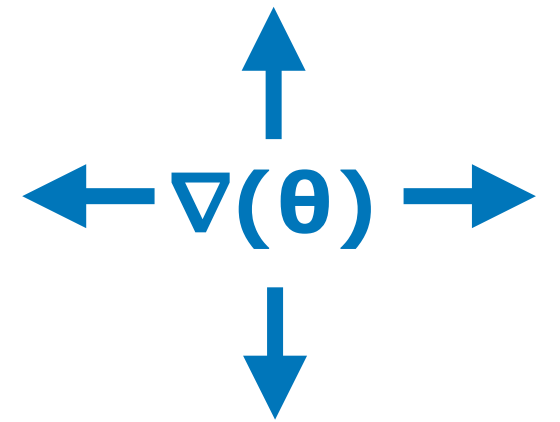


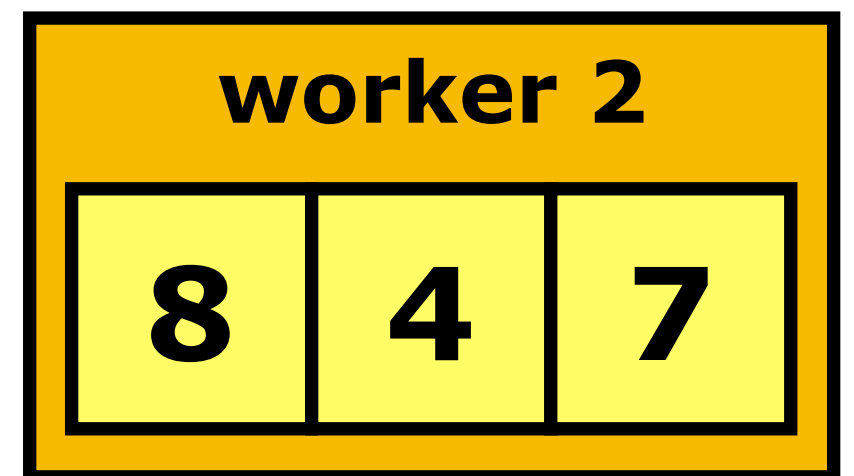
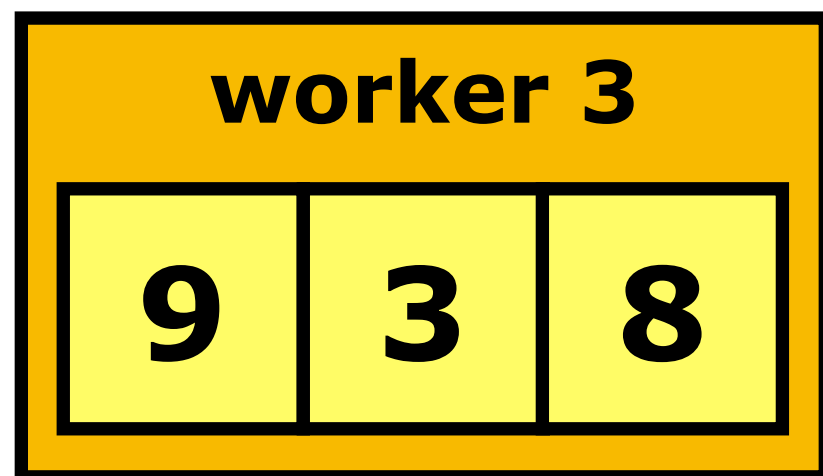
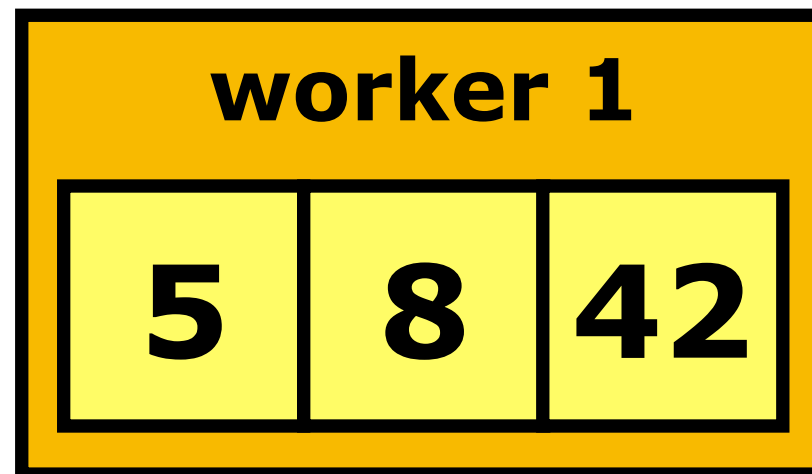
Ring
Allreduce

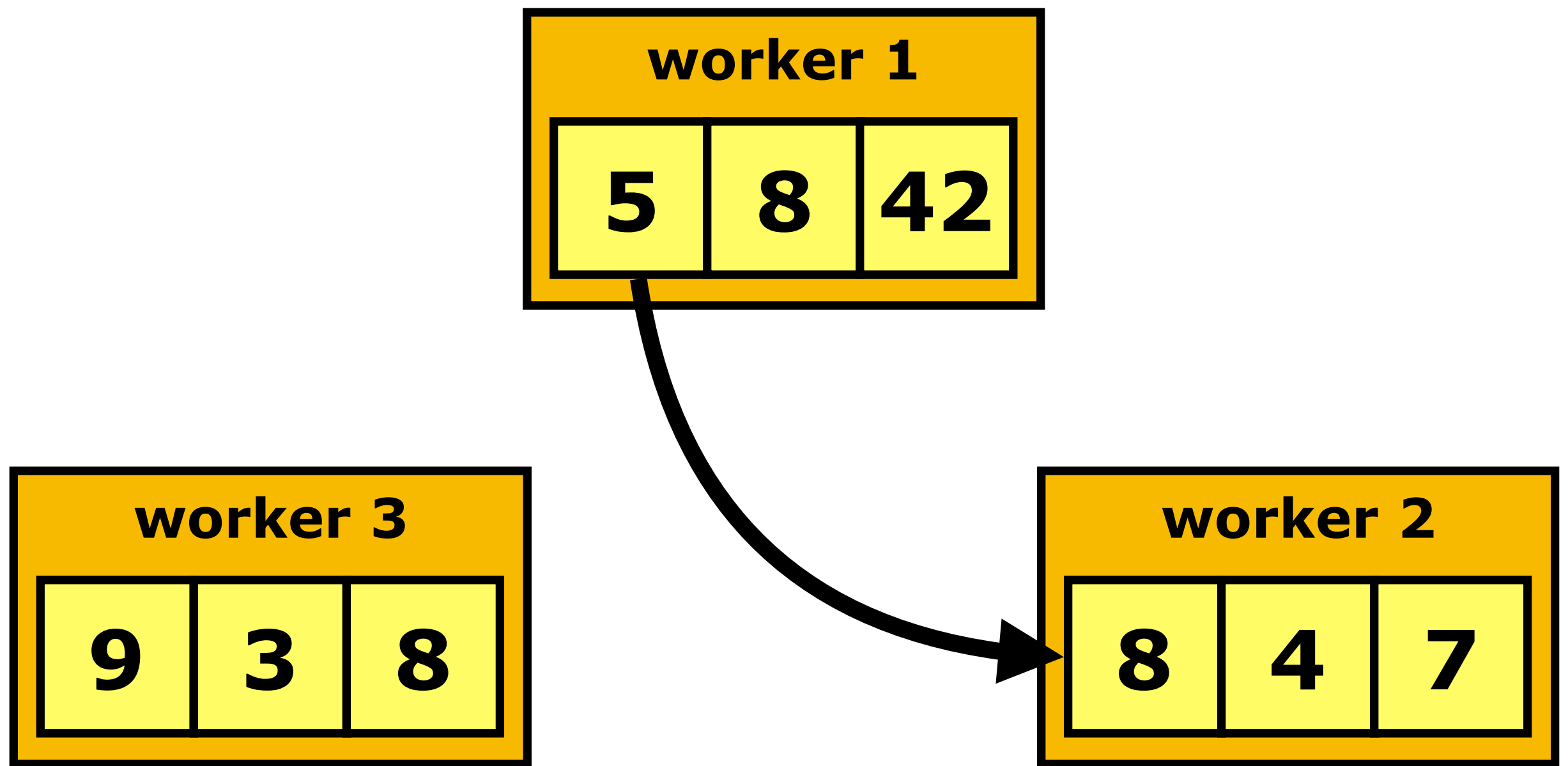
Ring Allreduce

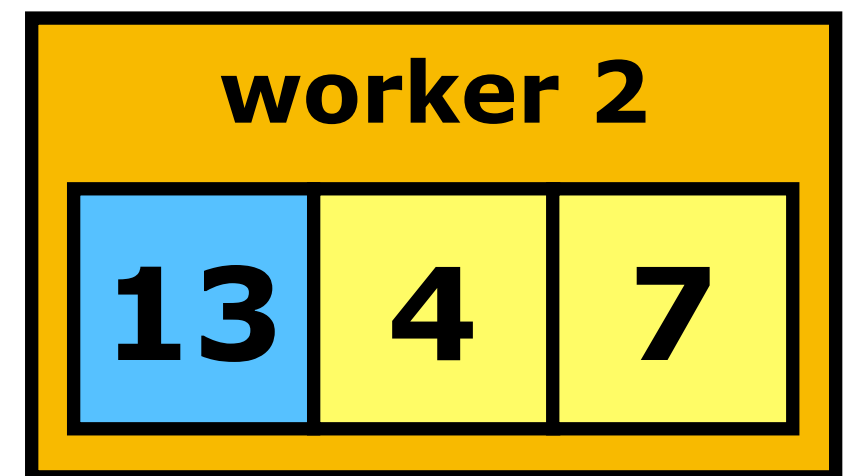
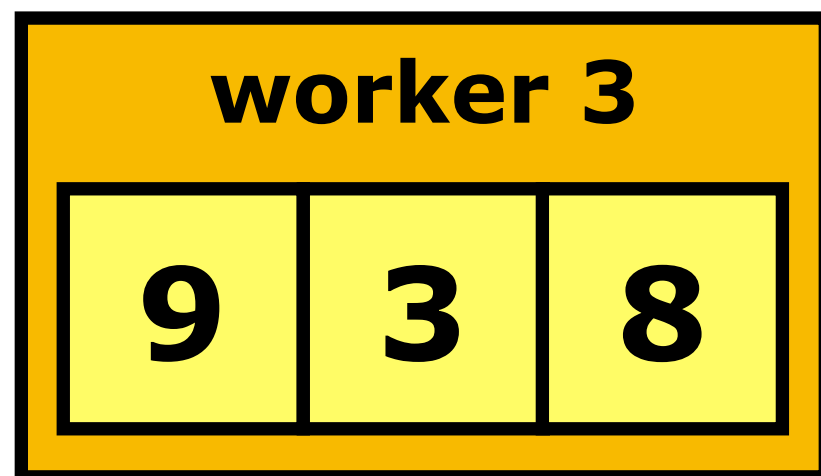
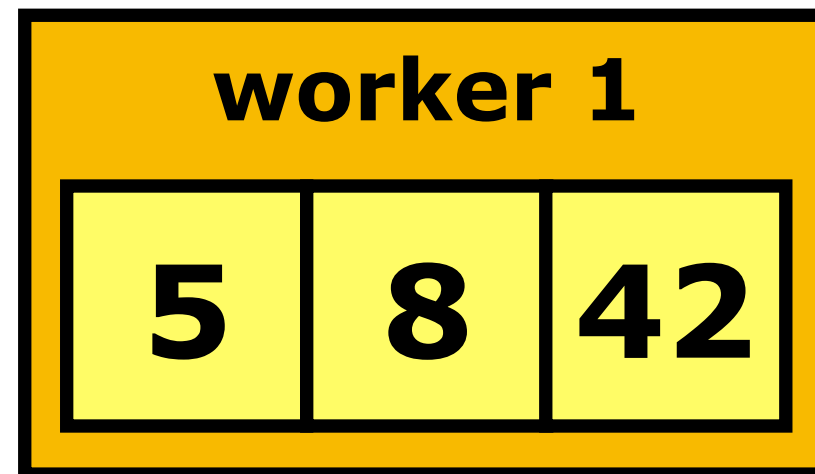


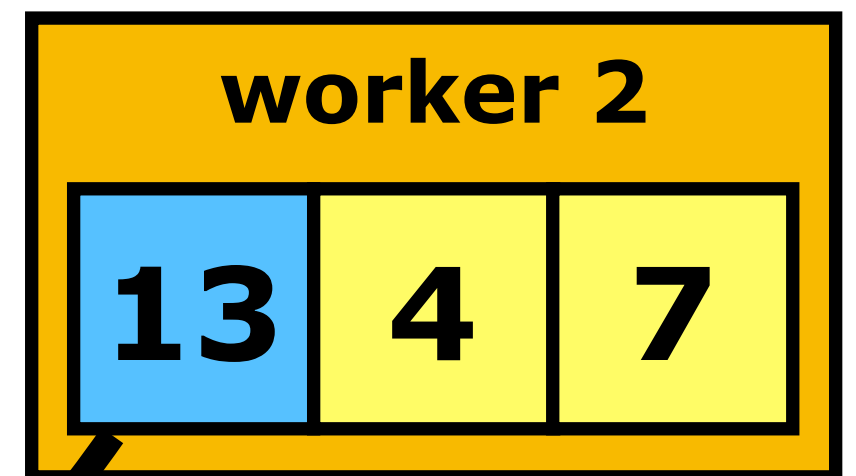
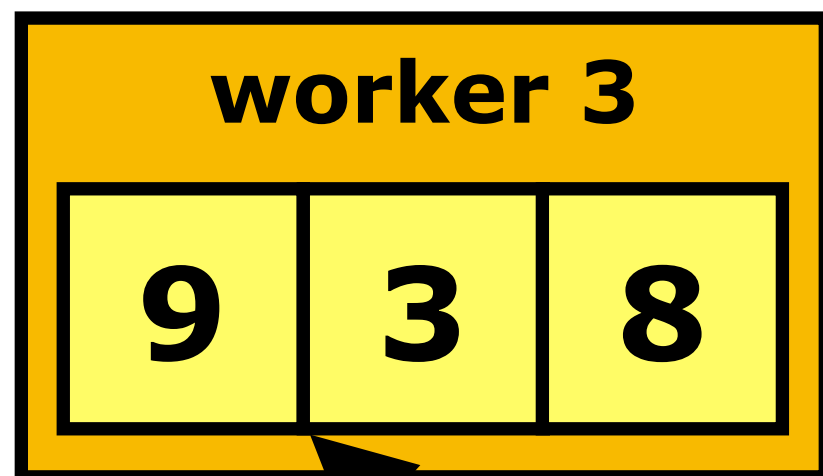
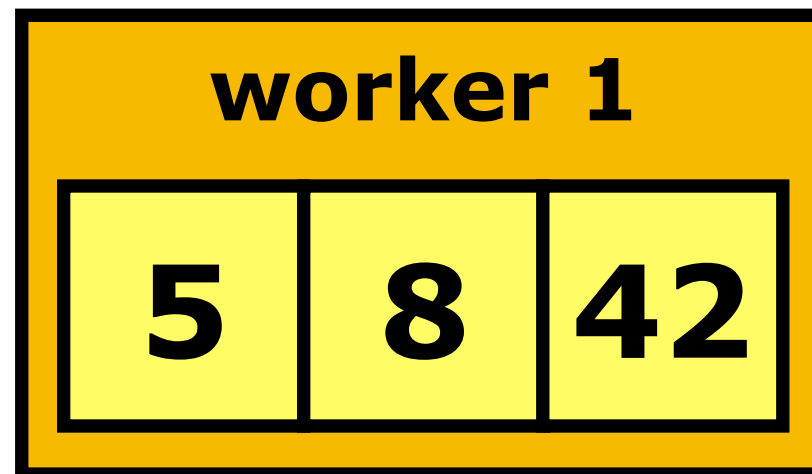
Ring Allreduce

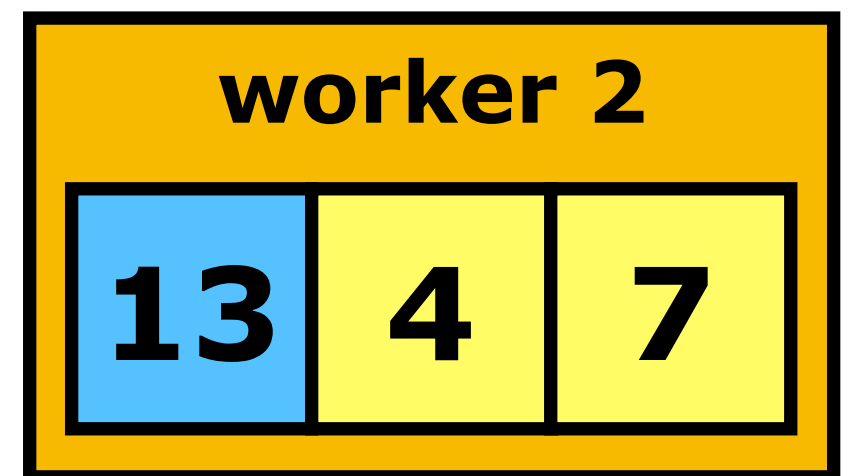
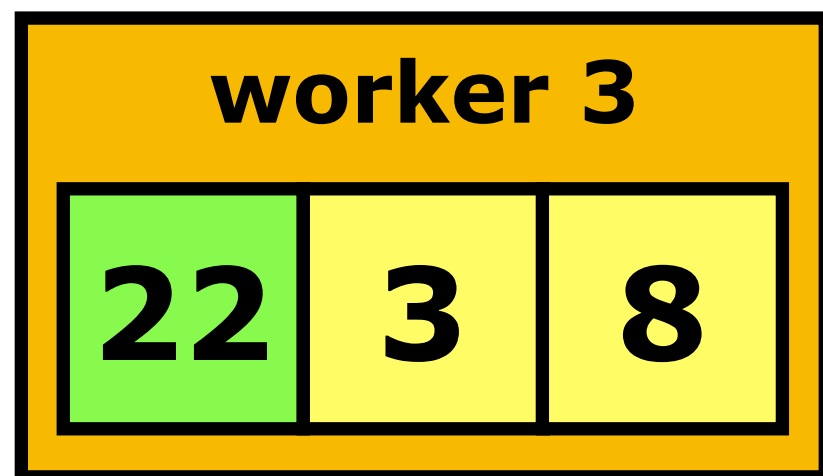
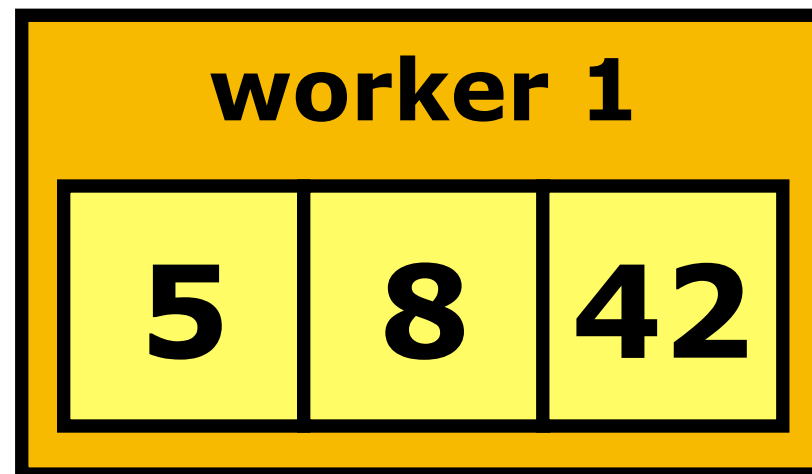


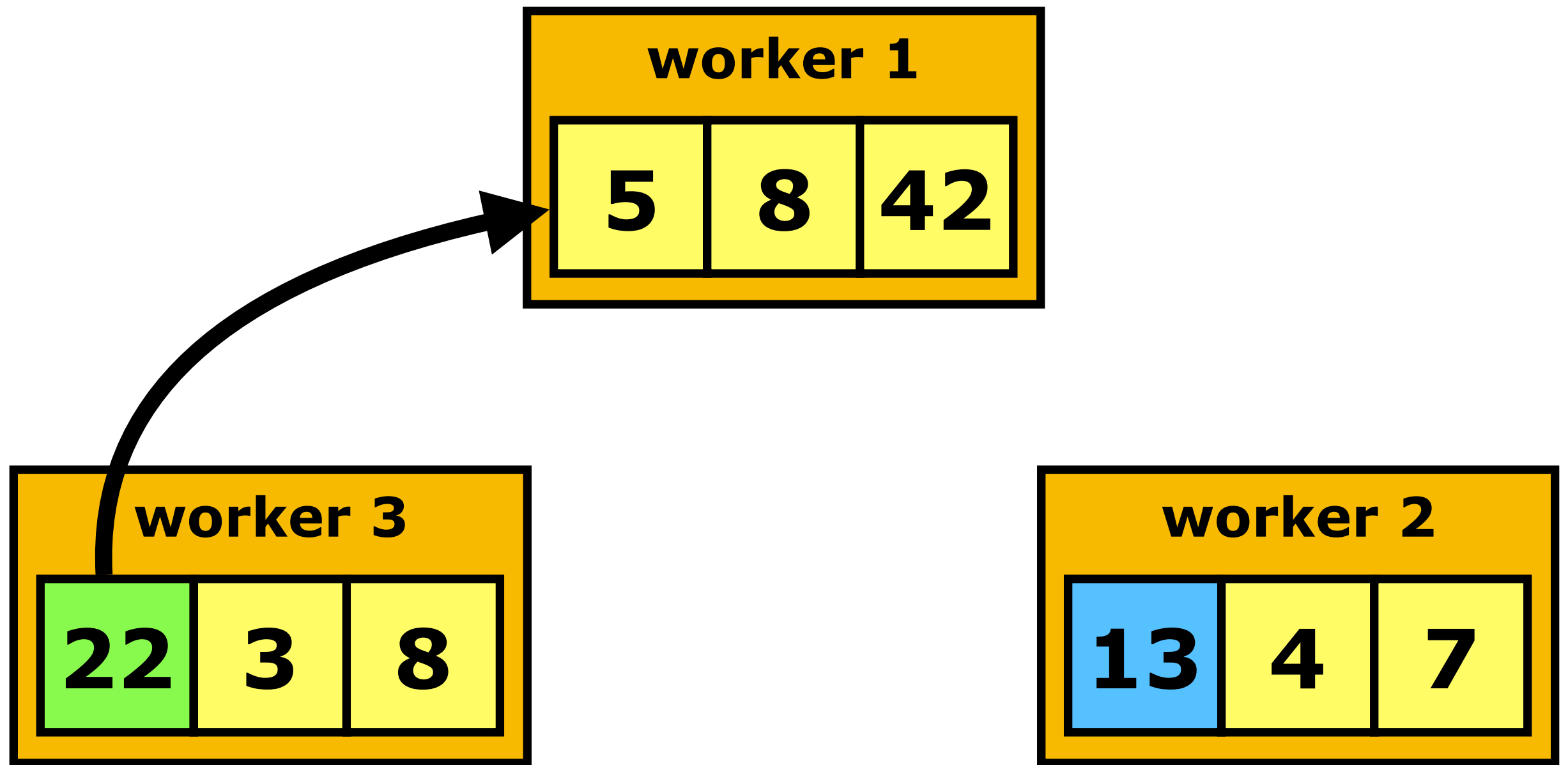


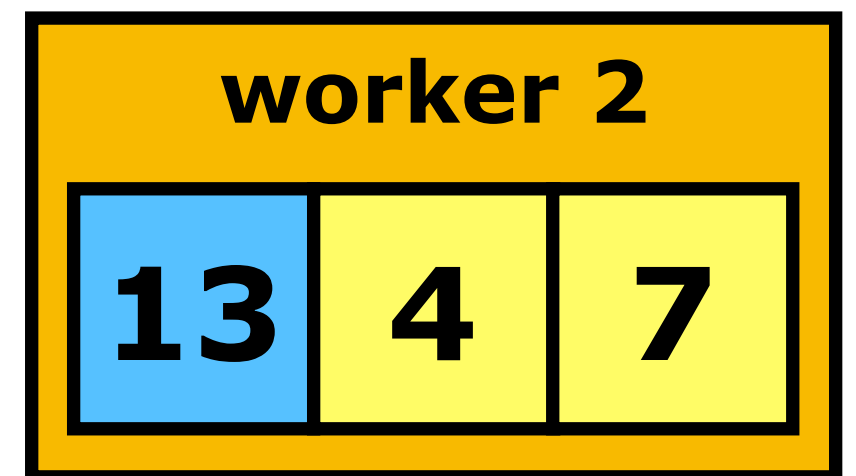
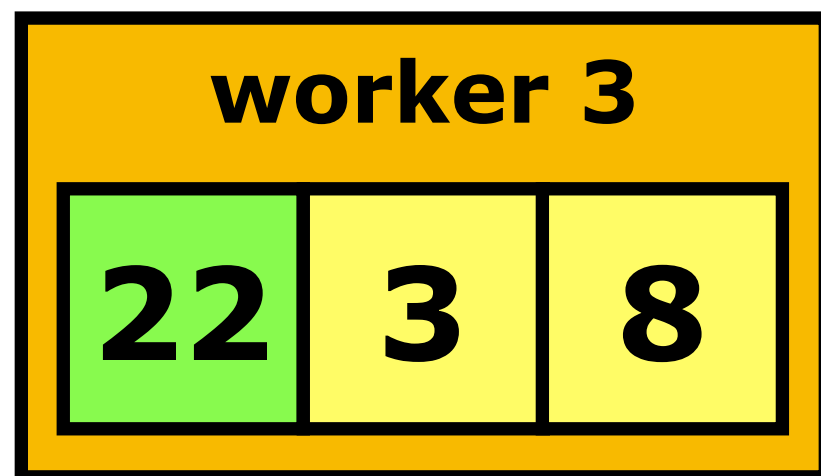
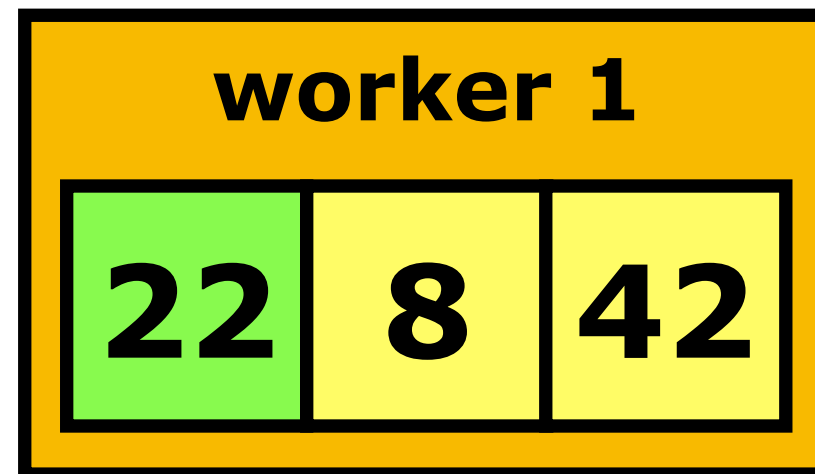


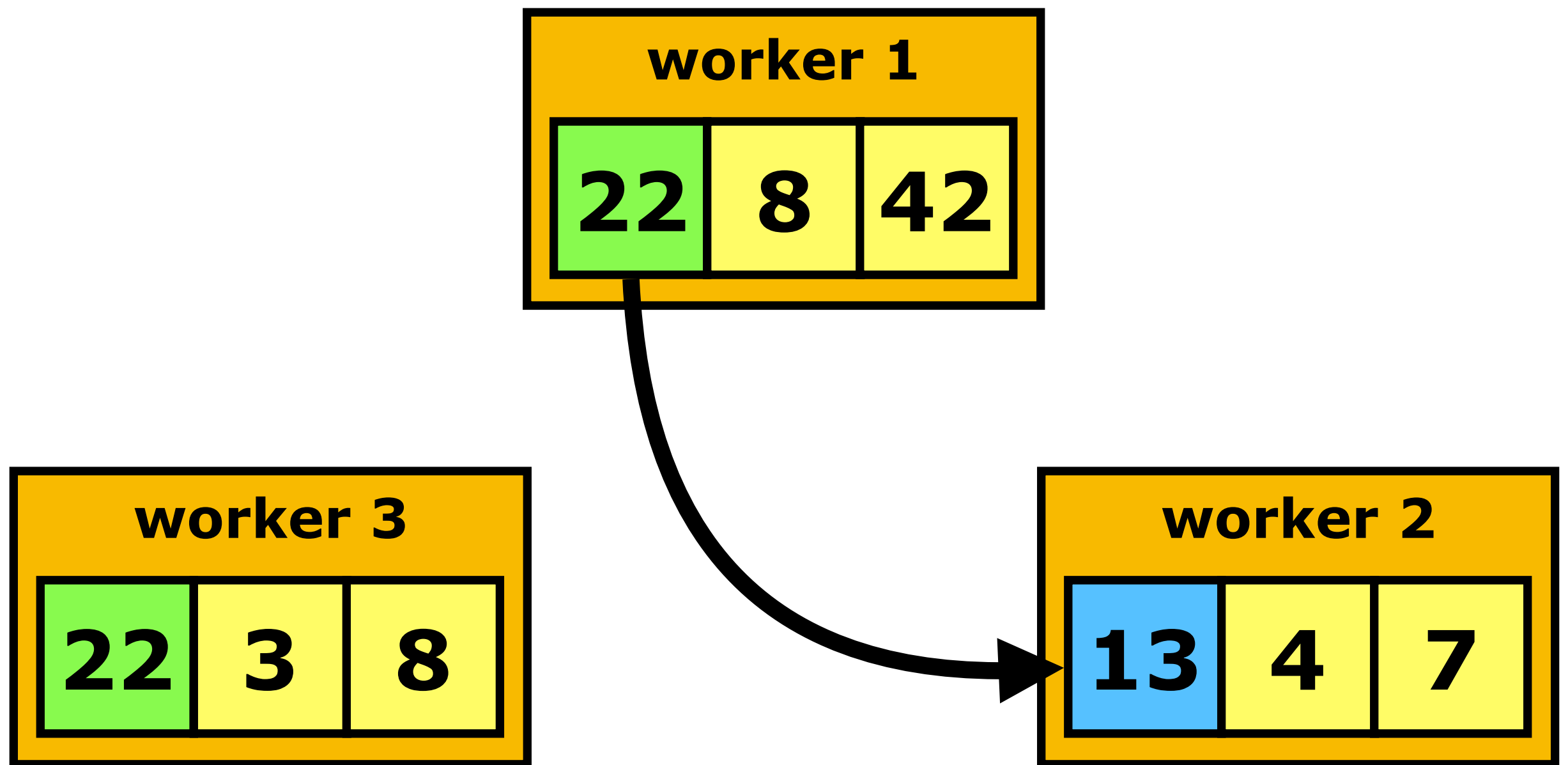


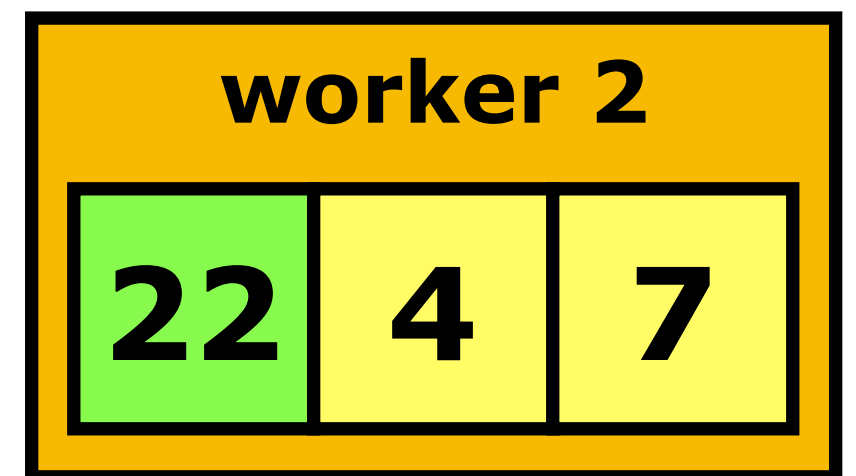
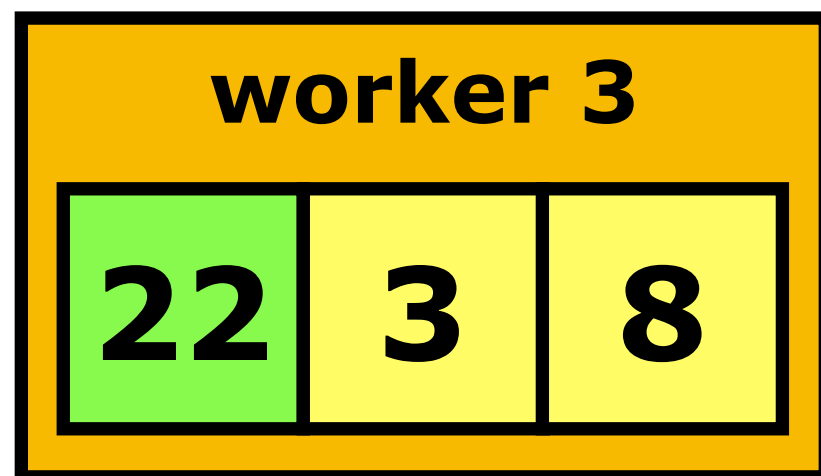
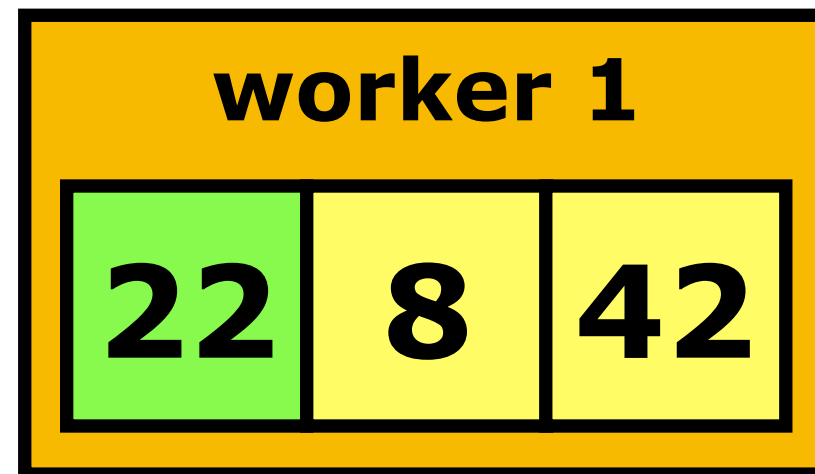


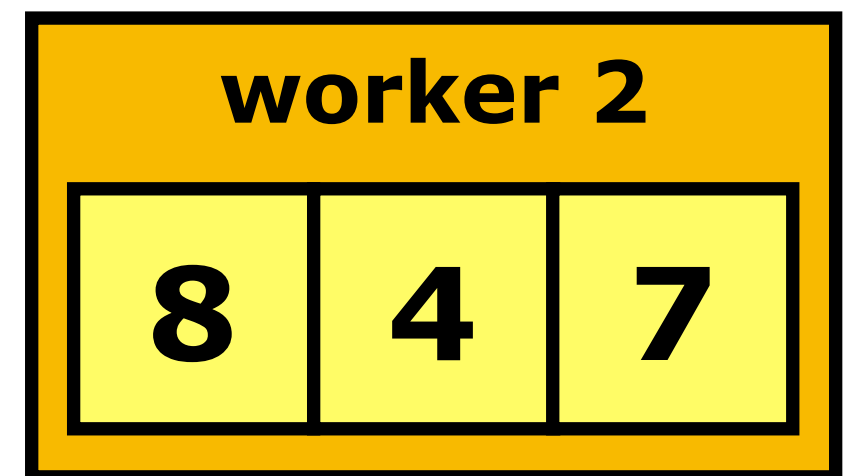
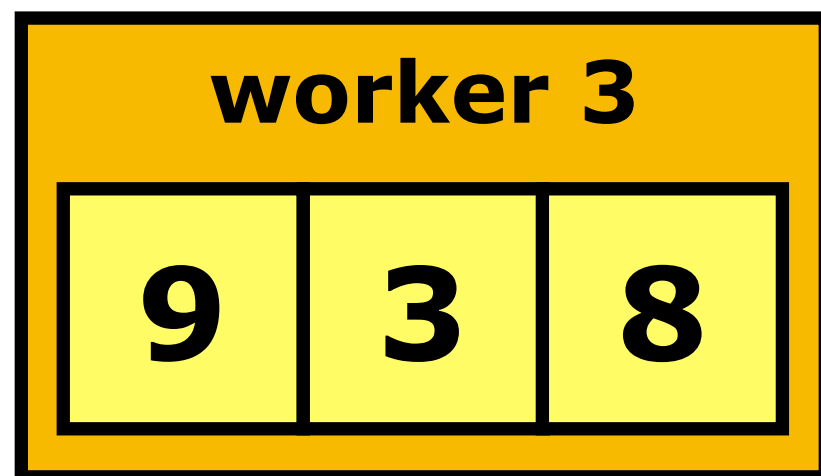
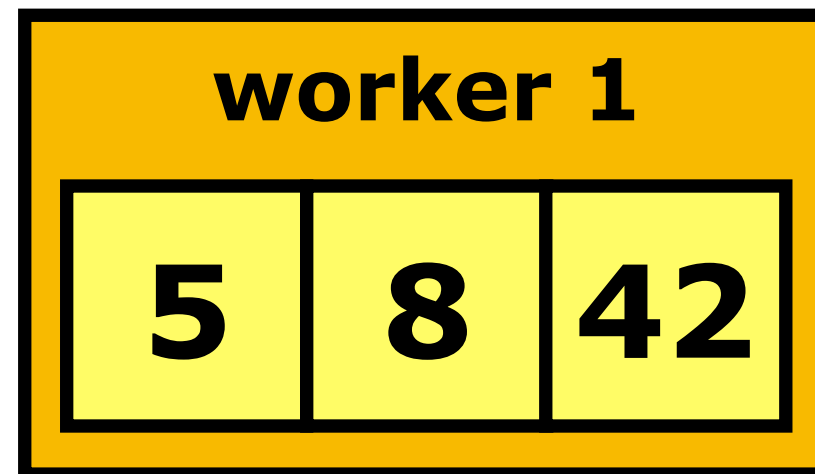


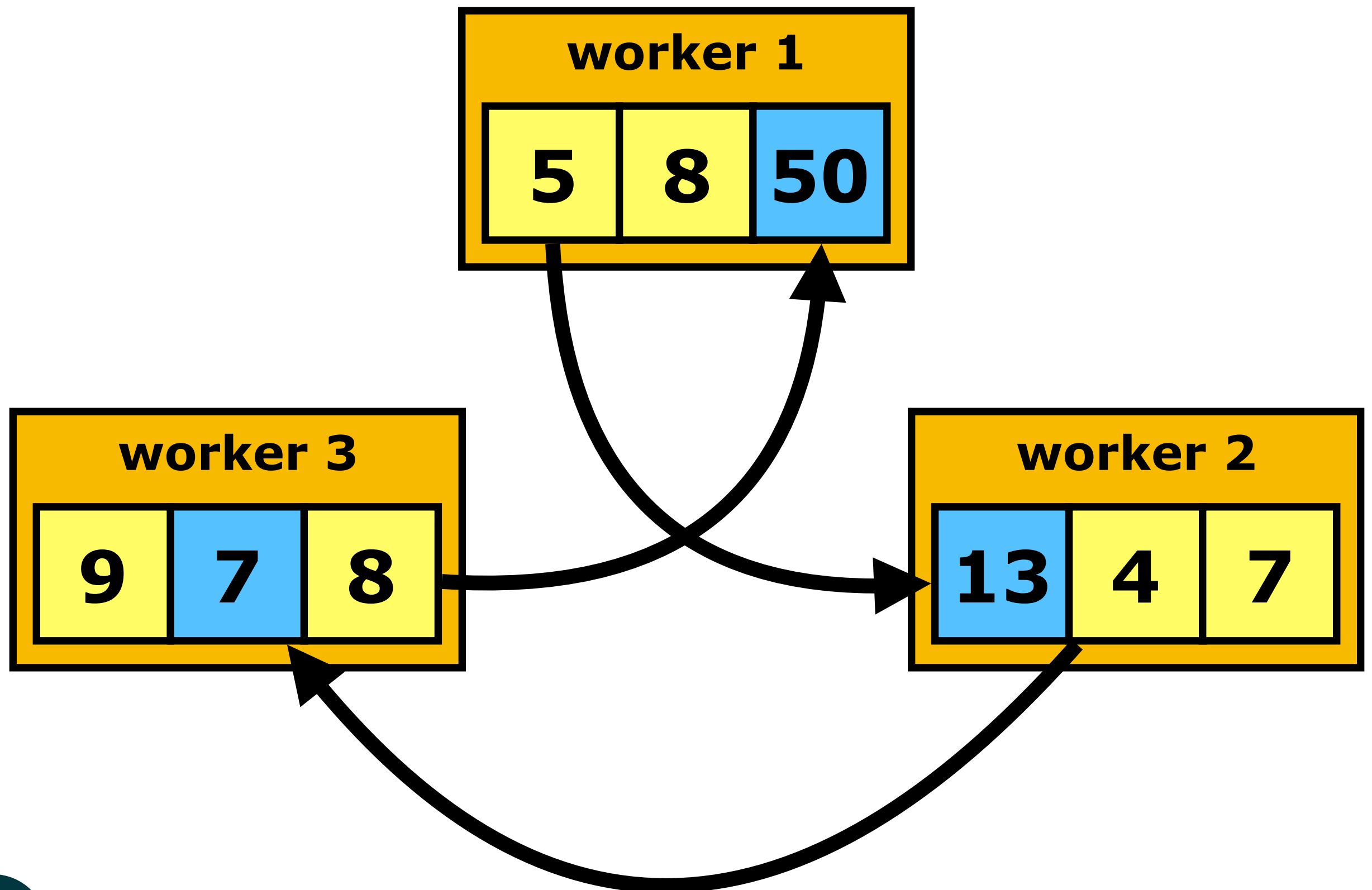


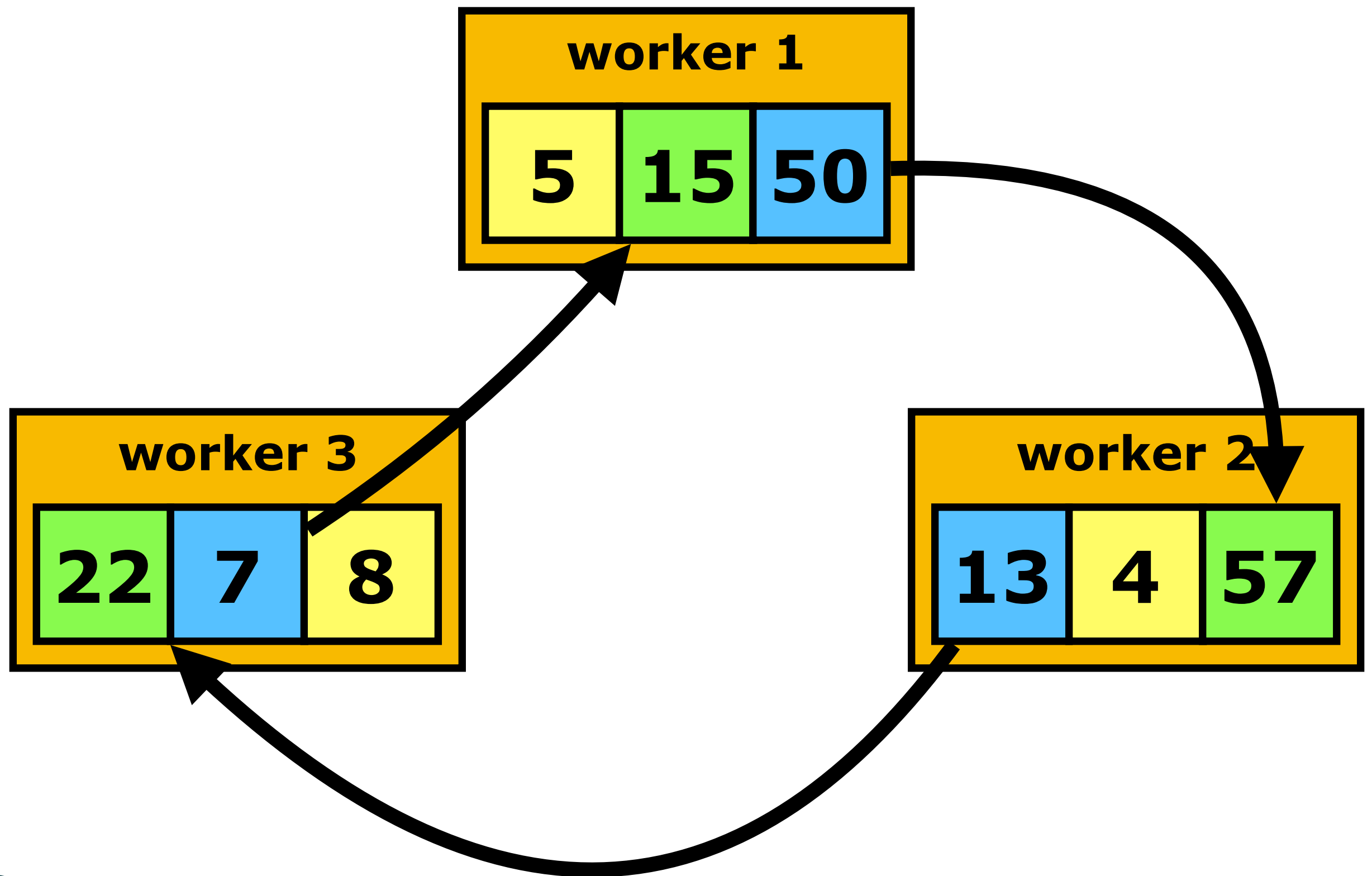


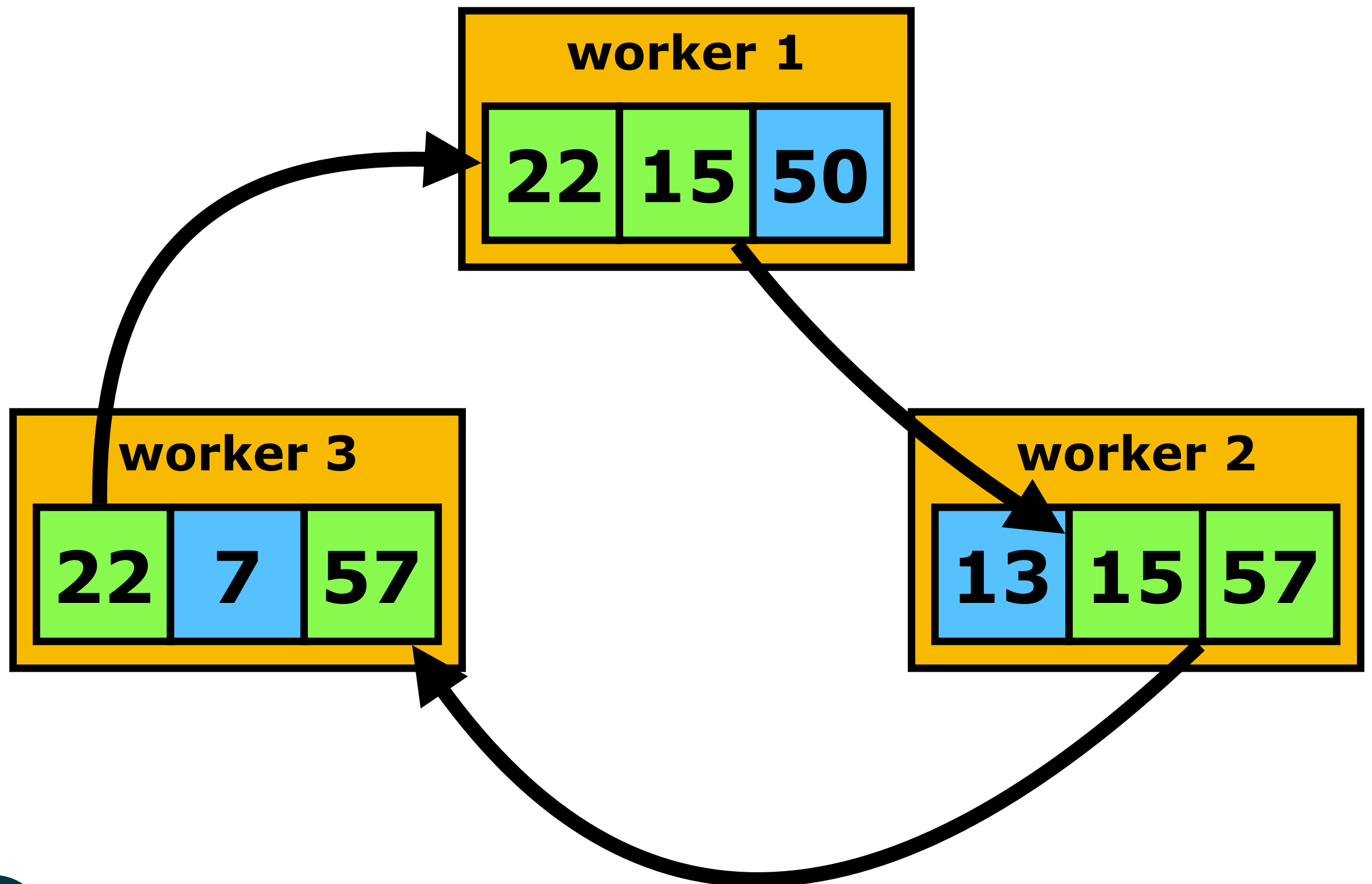


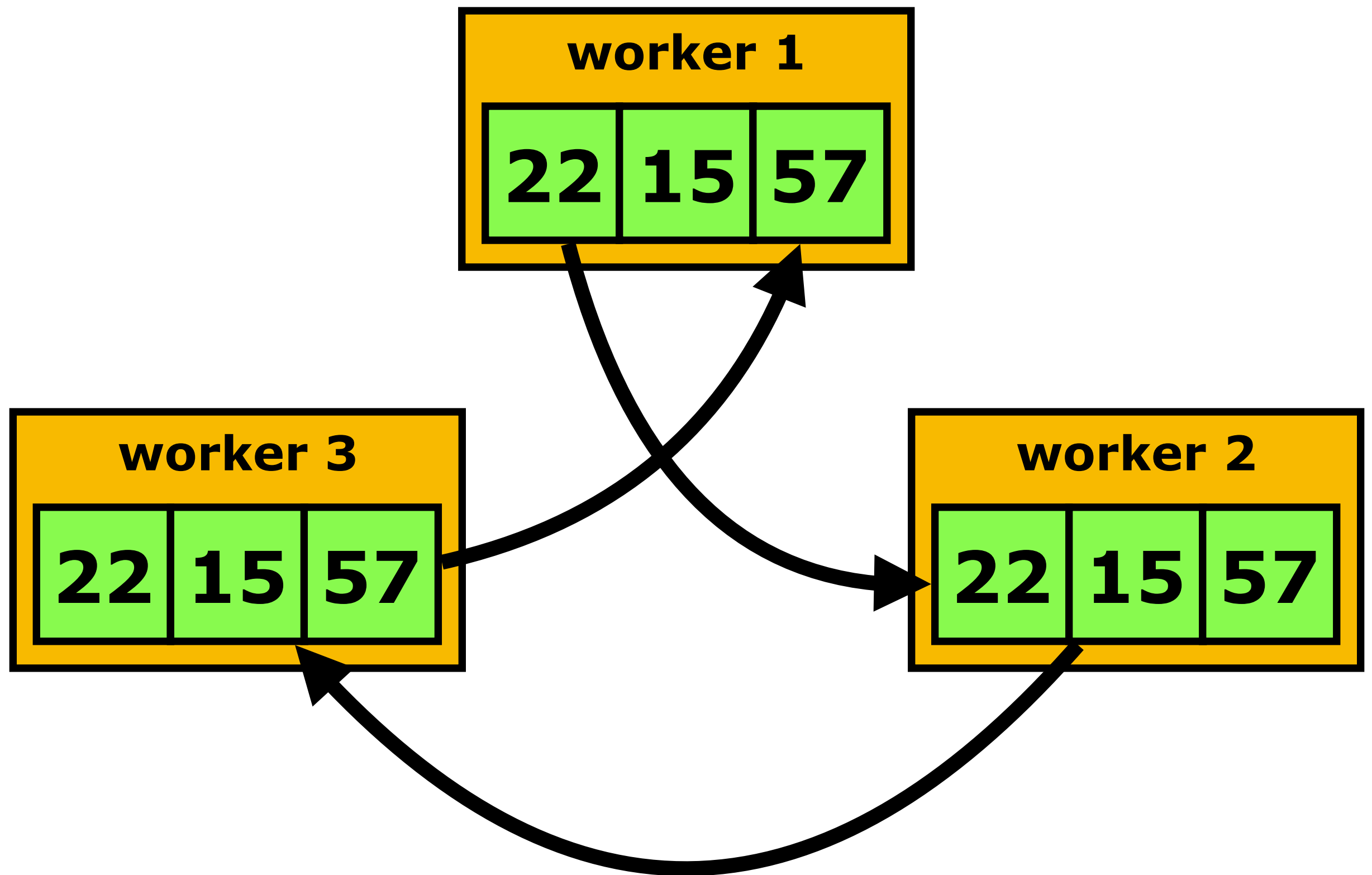


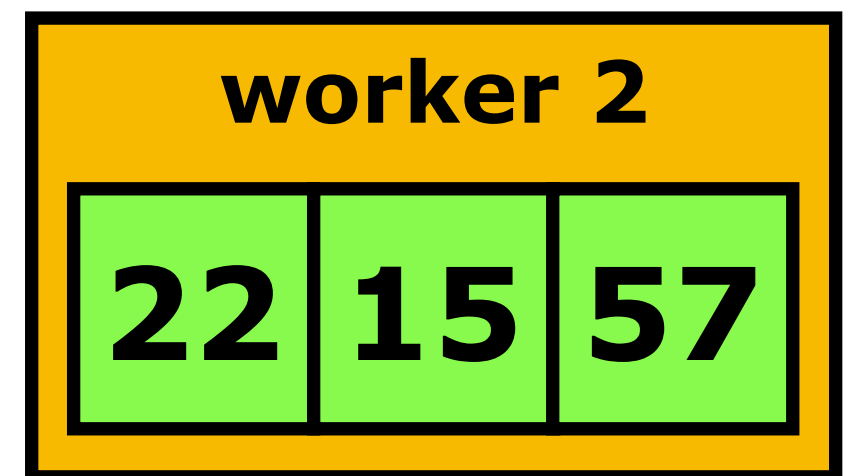
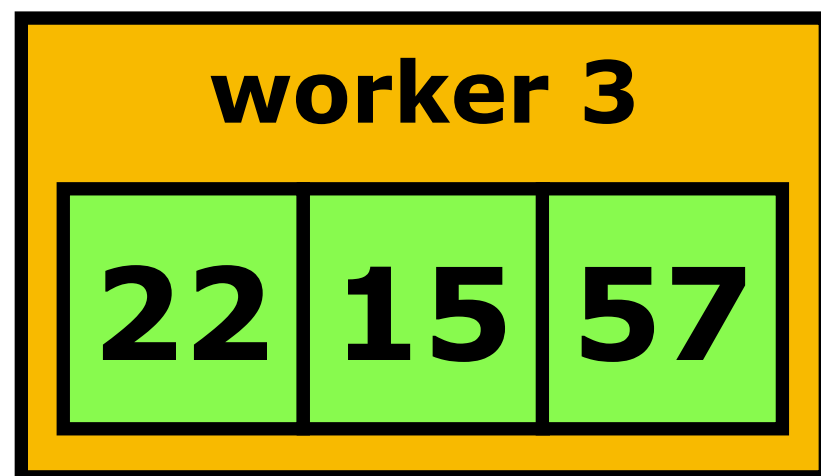
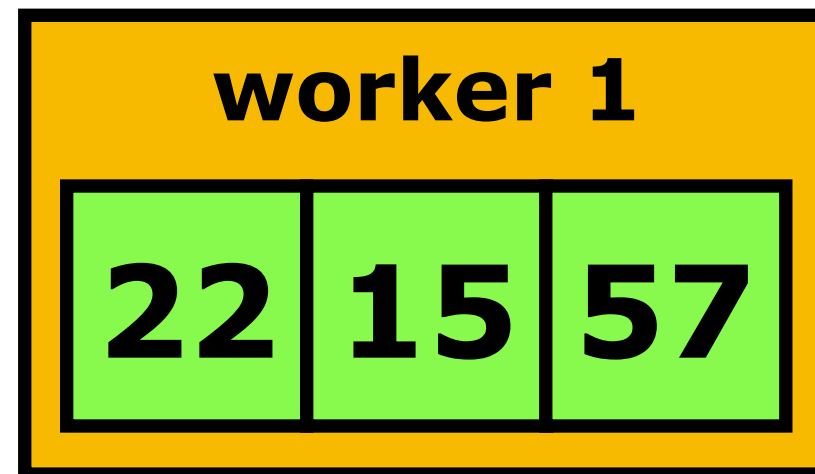


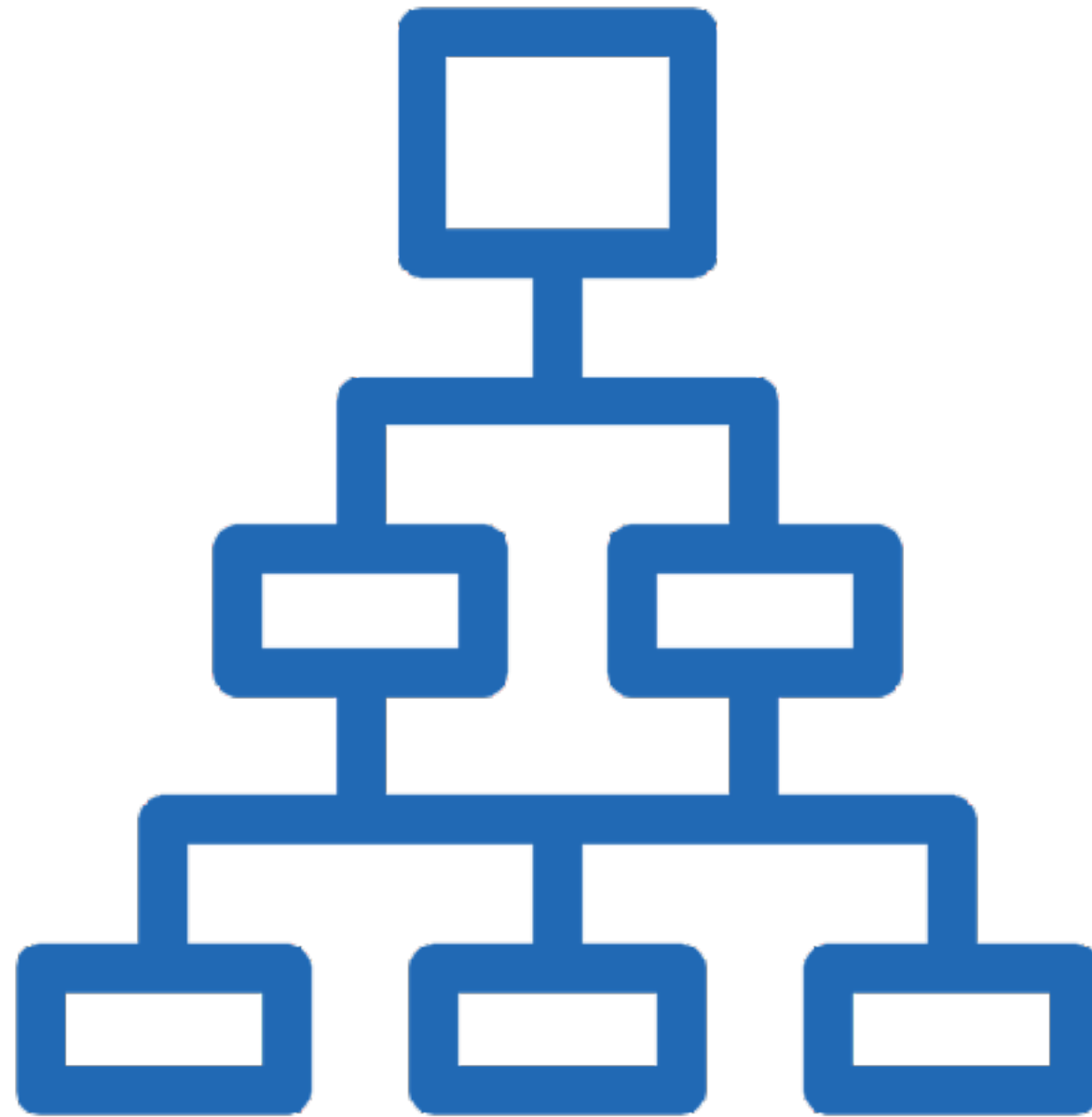






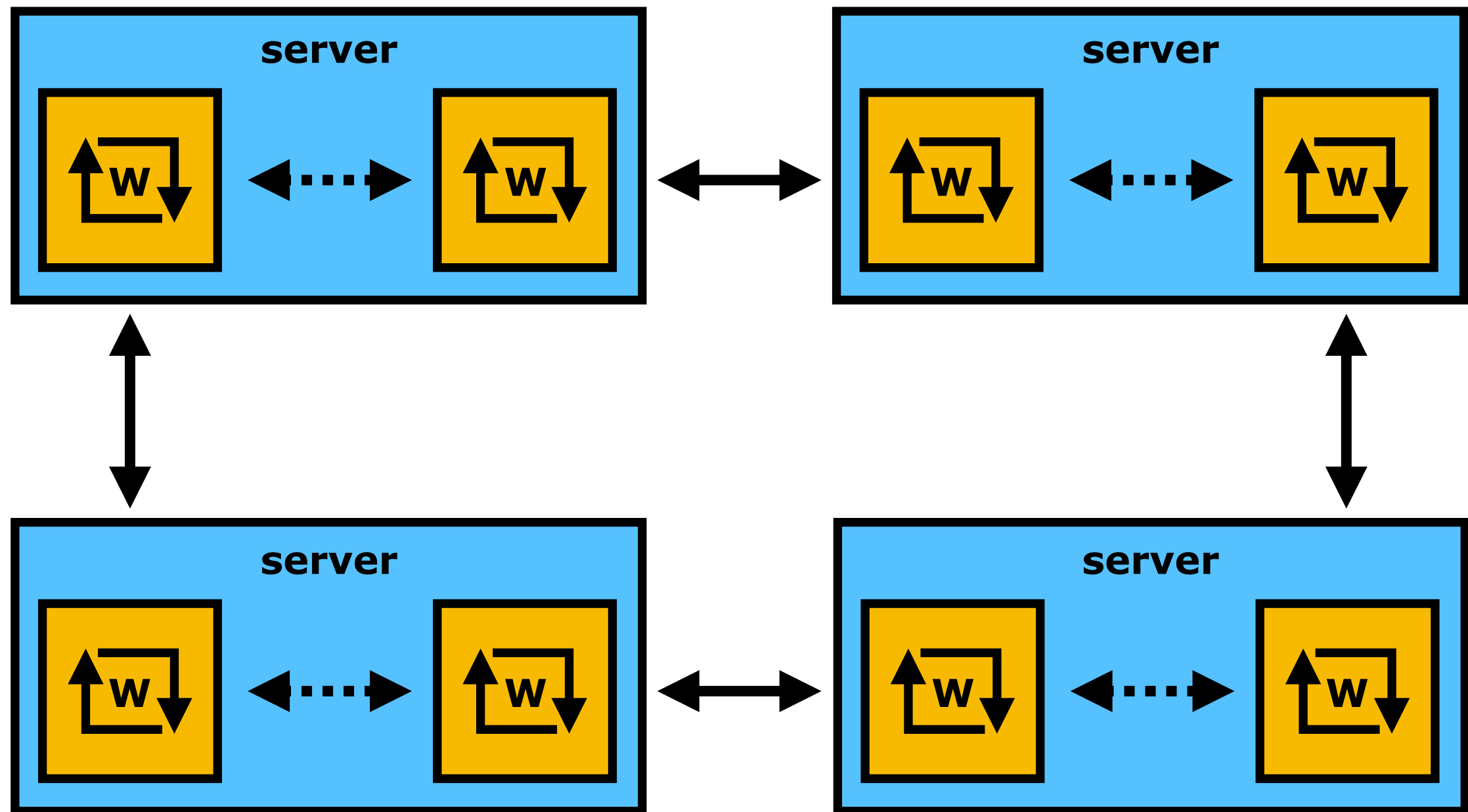
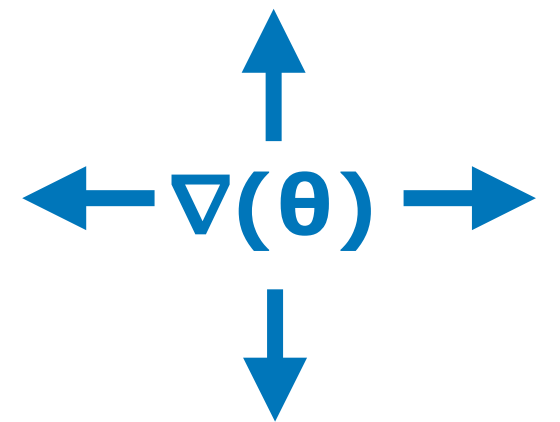


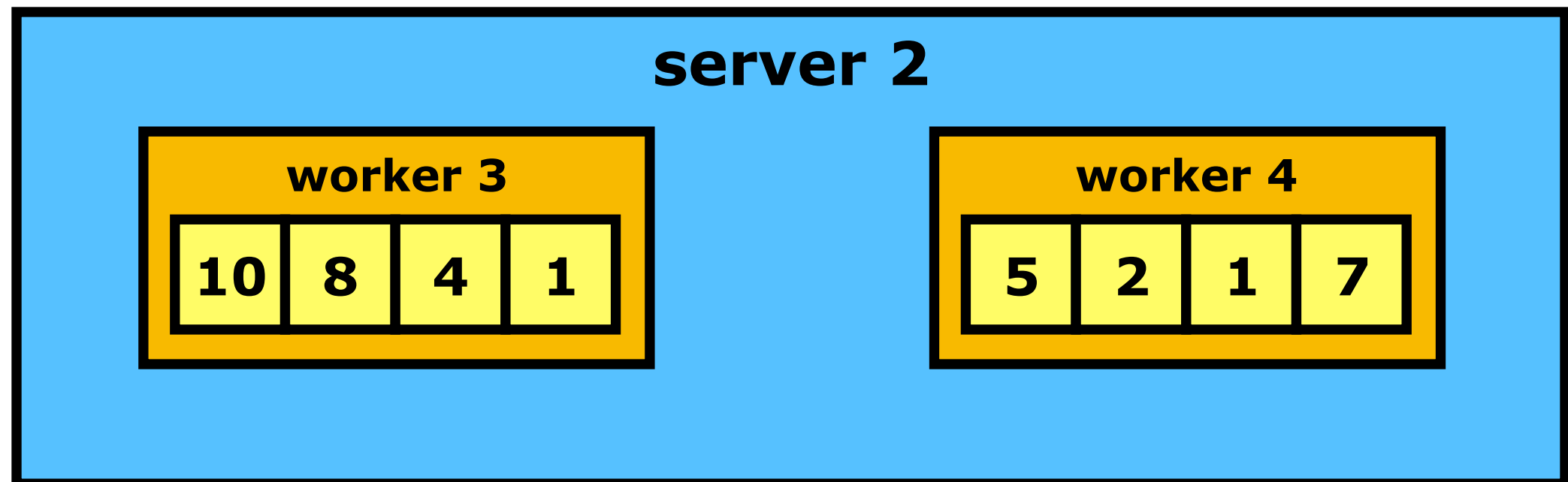
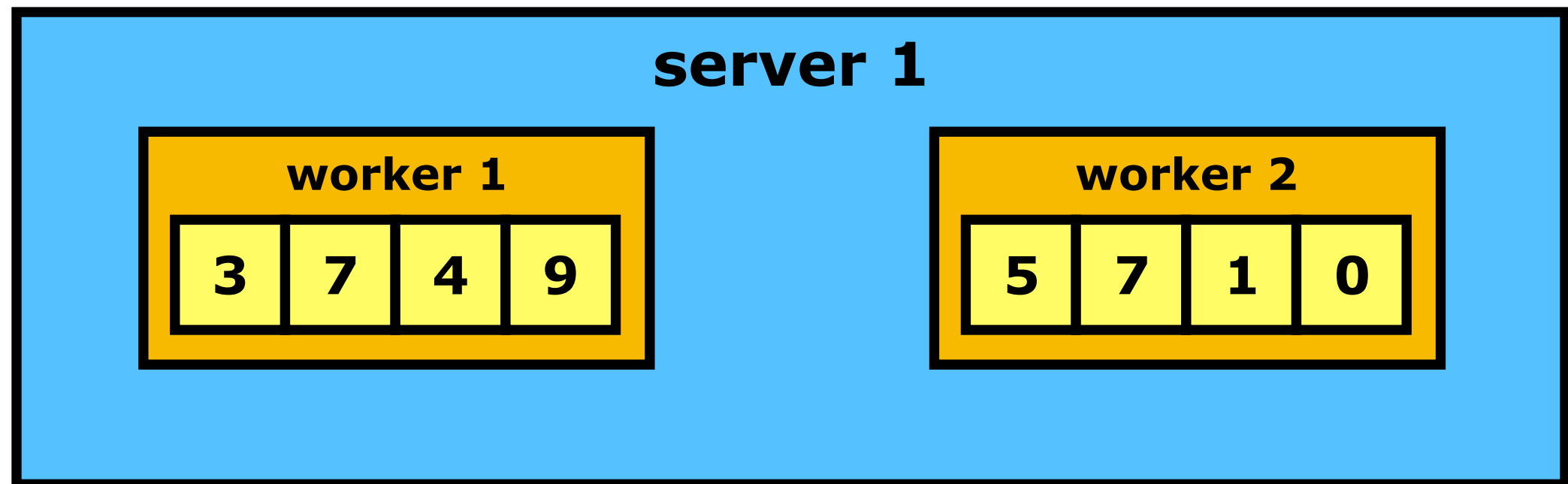


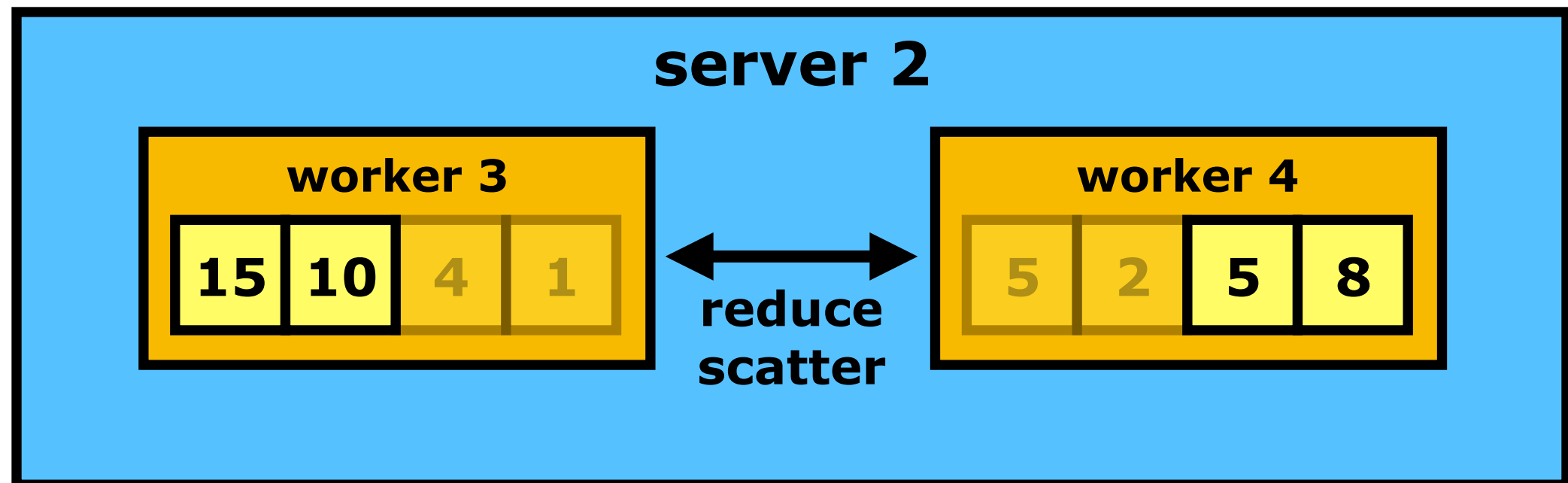
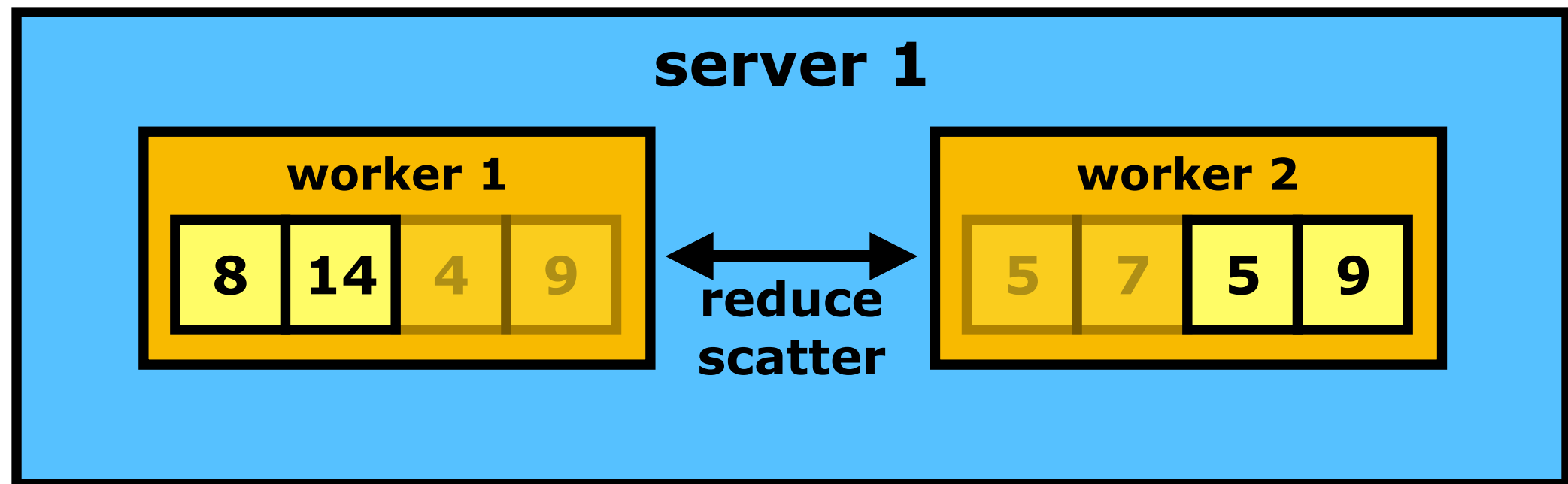


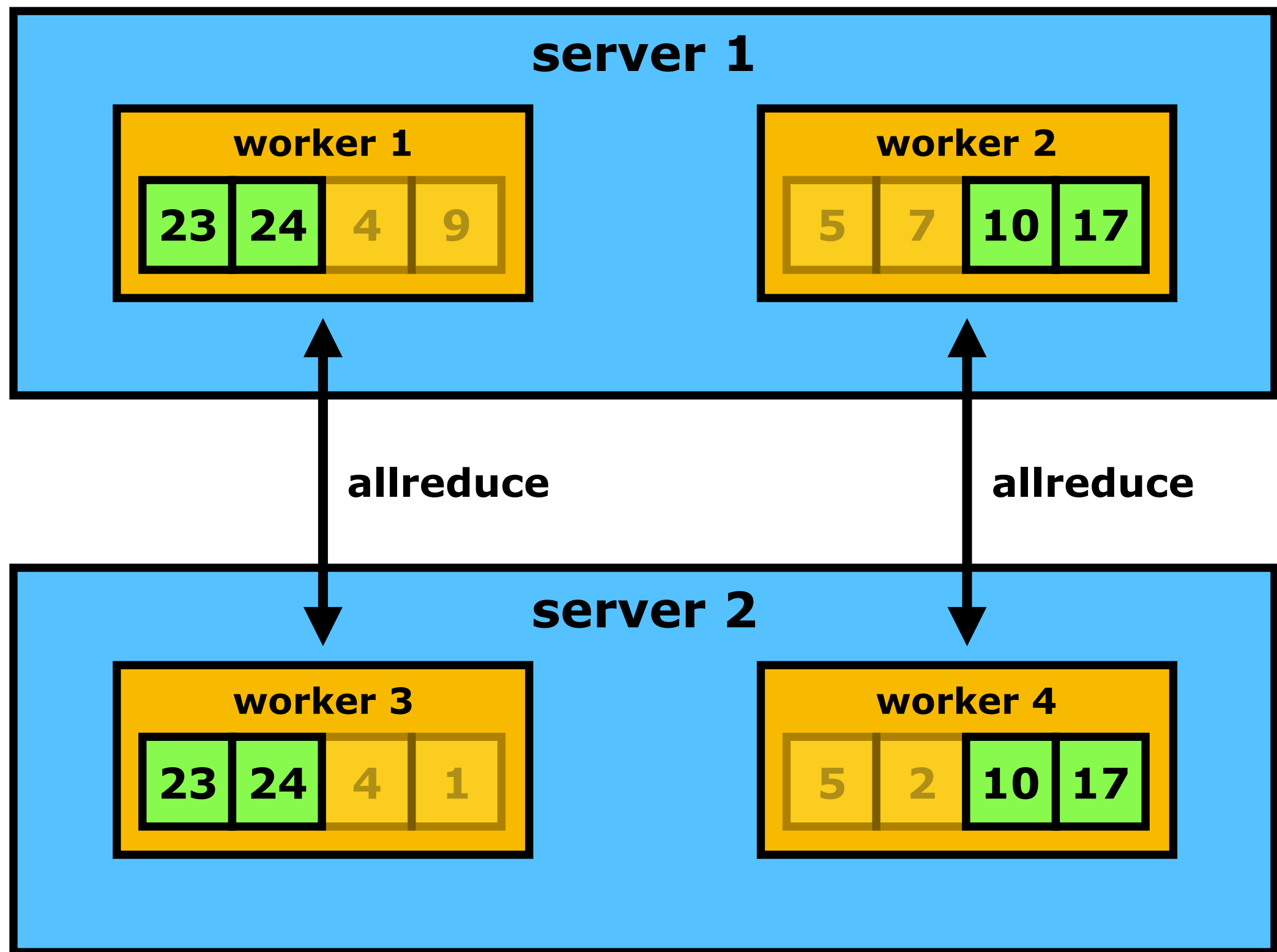
Hierarchical
Allreduce

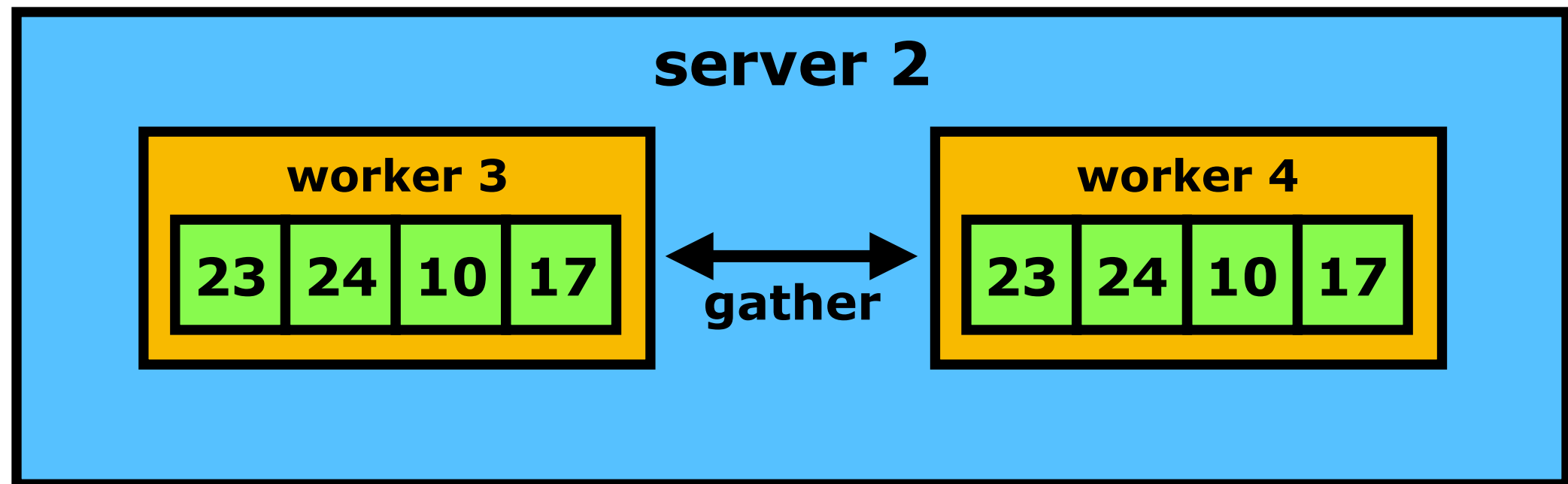
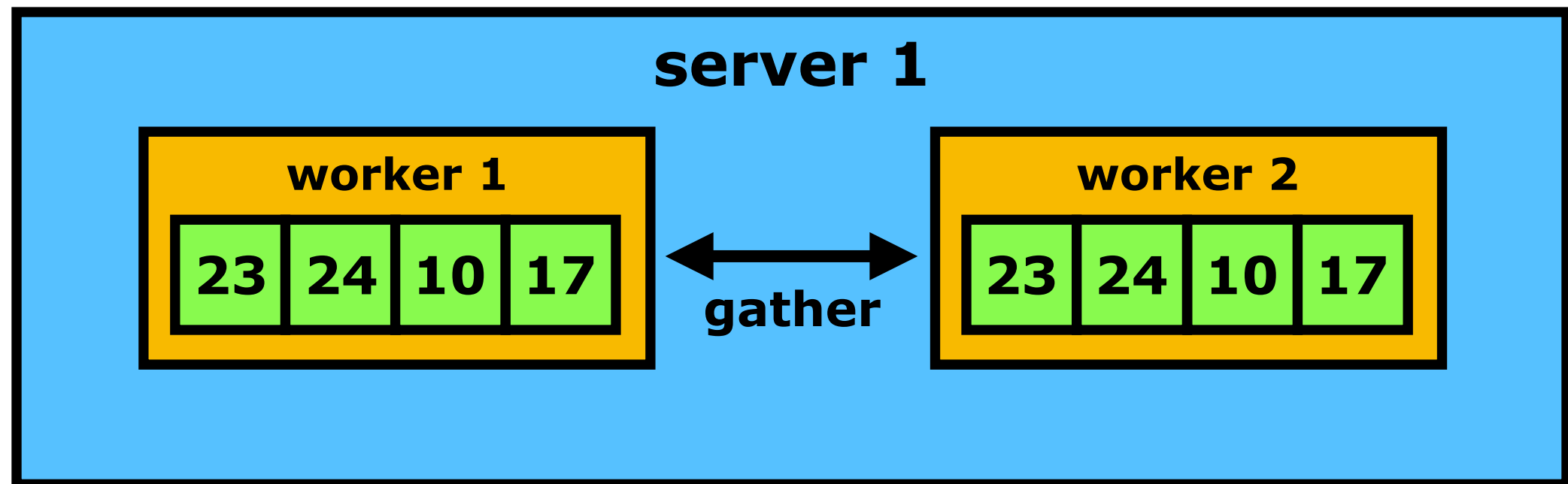
Hierarchical Allreduce



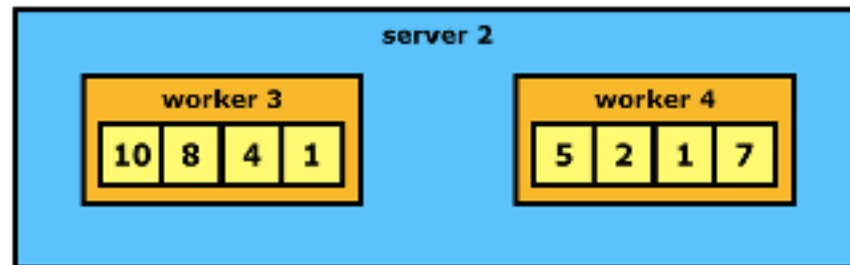
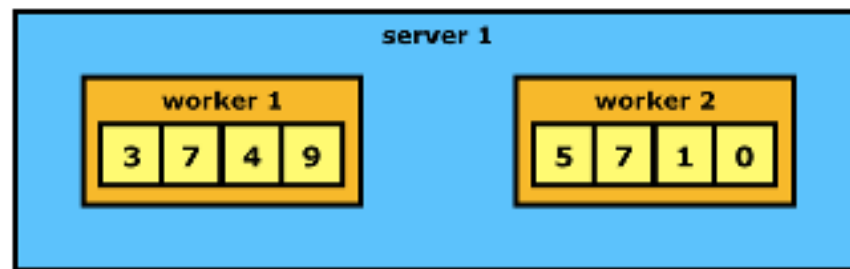




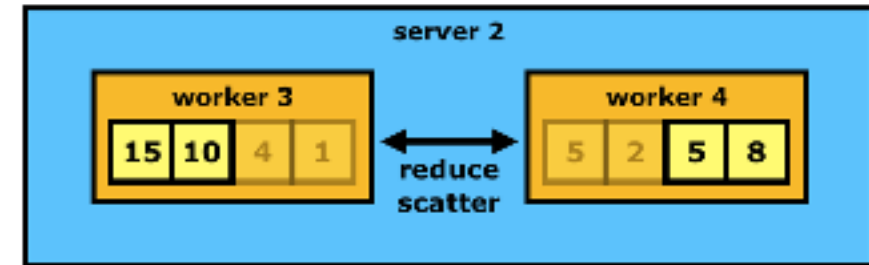
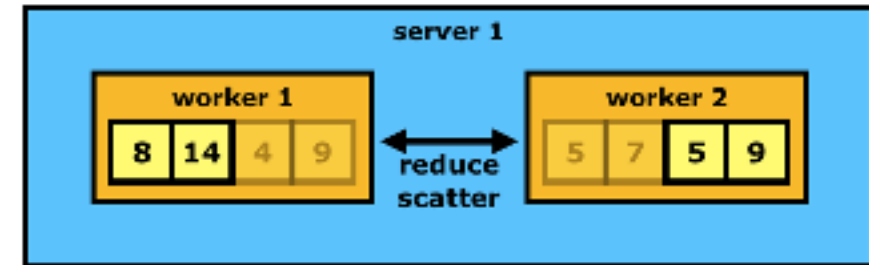




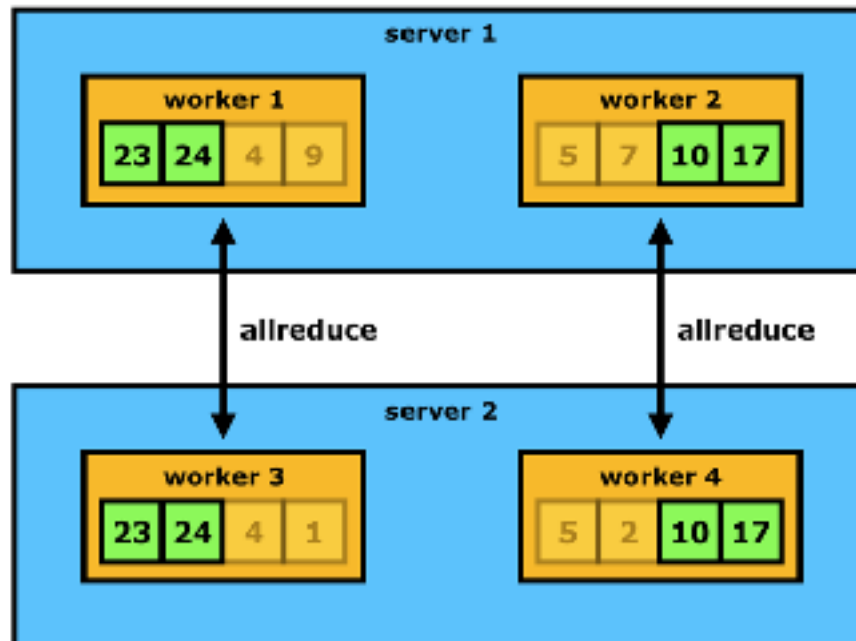
1. Setup



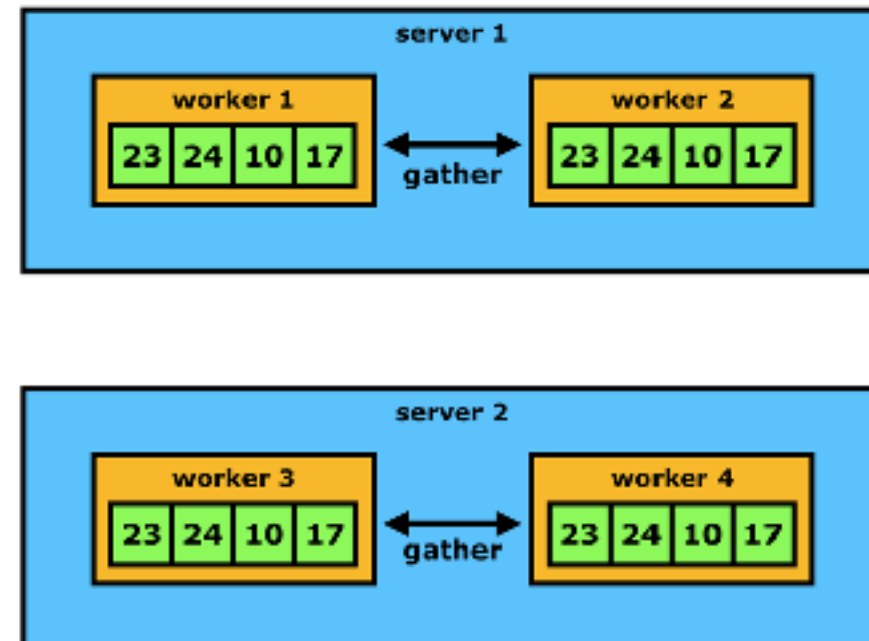
2. Local ReduceScatter



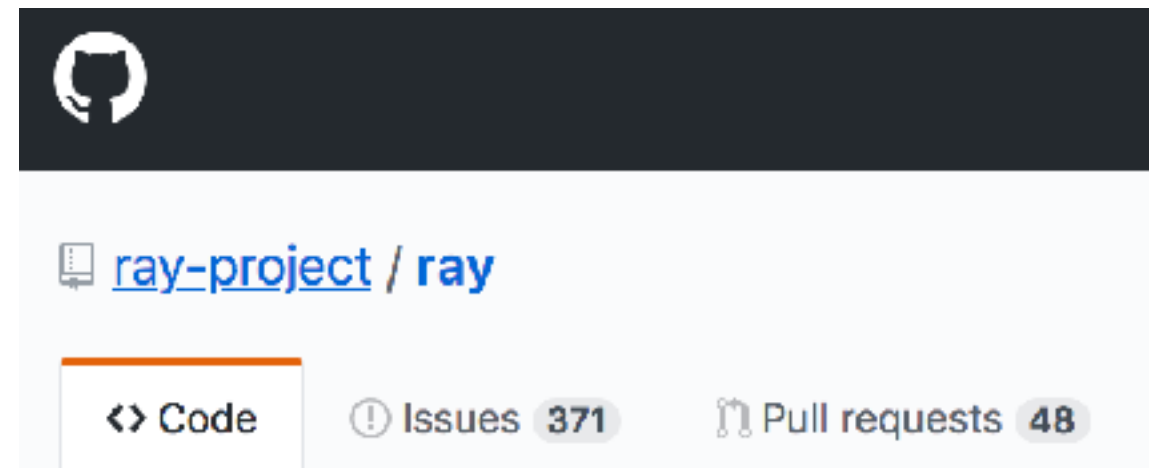
3. Global Allreduce



4. Local Gather



Tools





VALOHAI

Deep Learning Management Platform



Sources and Further Reading

HOGWILD!: A Lock-Free Approach to Parallelizing Stochastic Gradient Descent
Feng Niu, Benjamin Recht, Christopher Ré, Stephen J. Wright

Large Scale Distributed Deep Networks
Jeffrey Dean, Greg S. Corrado, Rajat Monga, Kai Chen, et al.

Bandwidth Optimal All-reduce Algorithms for Clusters of Workstations
Pitch Patarasuk, Xin Yuan

Slim Fly: A Cost Effective Low-Diameter Network Topology
Maciej Besta, Torsten Hoefler

Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour
Priya Goyal, Piotr Dollar, , Ross Girshick, et al.

Horovod: Distributed Deep Learning in 5 Lines of Python
Uber Open Summit 2018

Distributed TensorFlow Documentation
TensorFlow

Thank You!

and Q&A

Ruksi Laine

@ruk_si

Machine Learning Engineer

@valohaiaia

<https://valohai.com/>