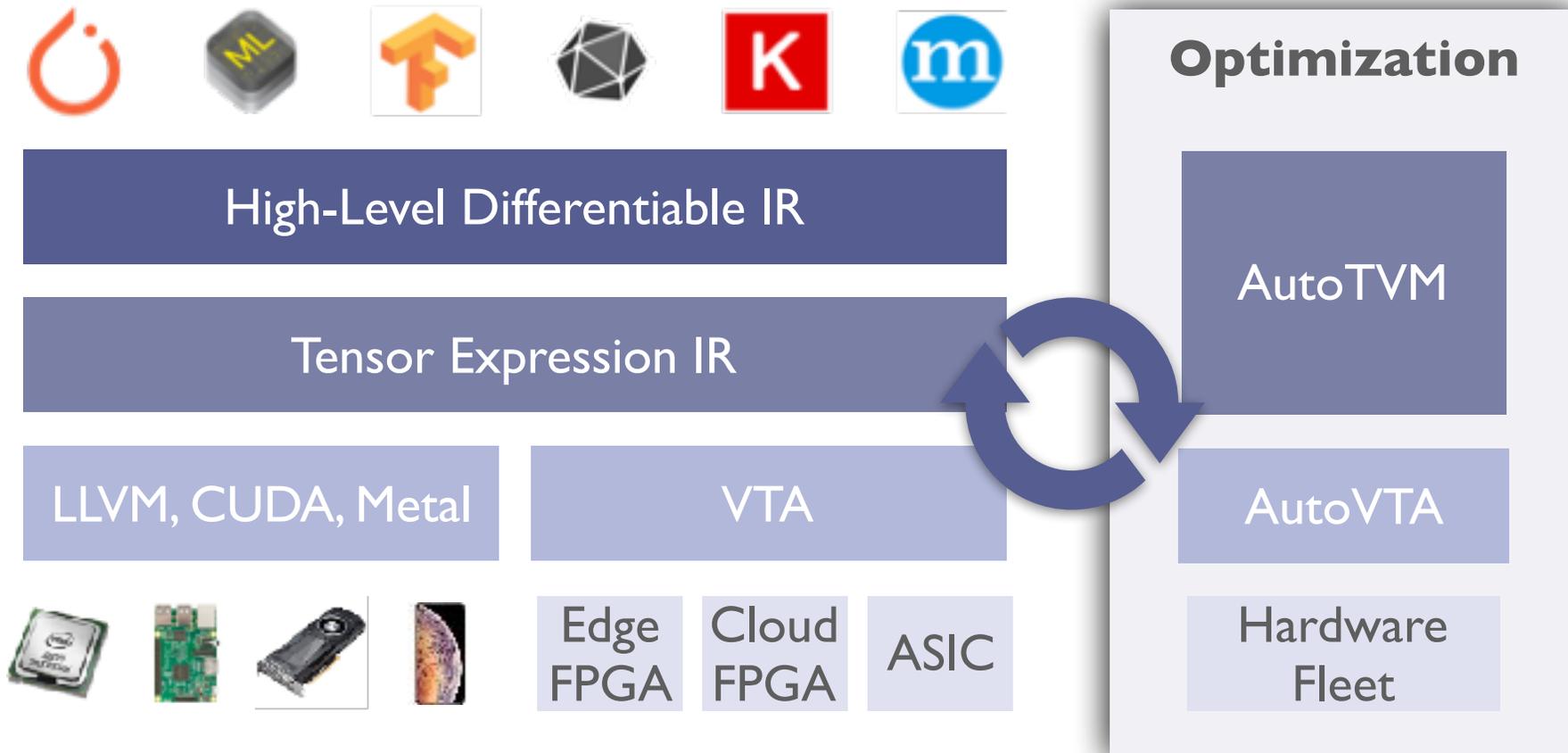
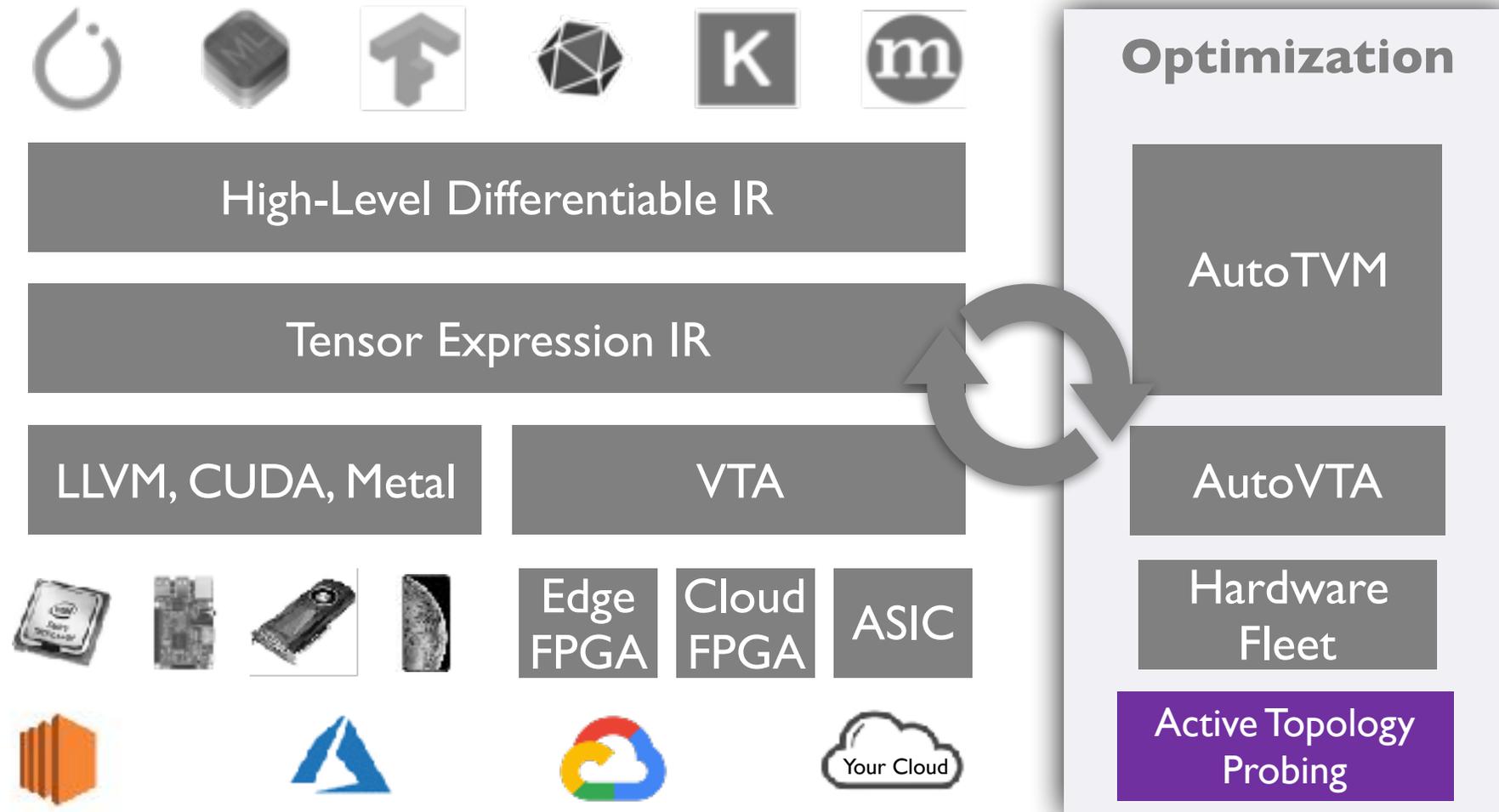


# Scalable Distributed Training with Parameter Hub: a whirlwind tour

# TVM Stack





**Groundwork for bringing TVM to the distributed world for training and inference, on commercial cloud, or in your own cluster.**

# Parameter Hub

**Optimized, topology-aware and dynamic mechanism for  
inter-machine communication**



# Parameter Hub

**Optimized, topology-aware and dynamic mechanism for inter-machine communication**

\* In the cloud-based training context



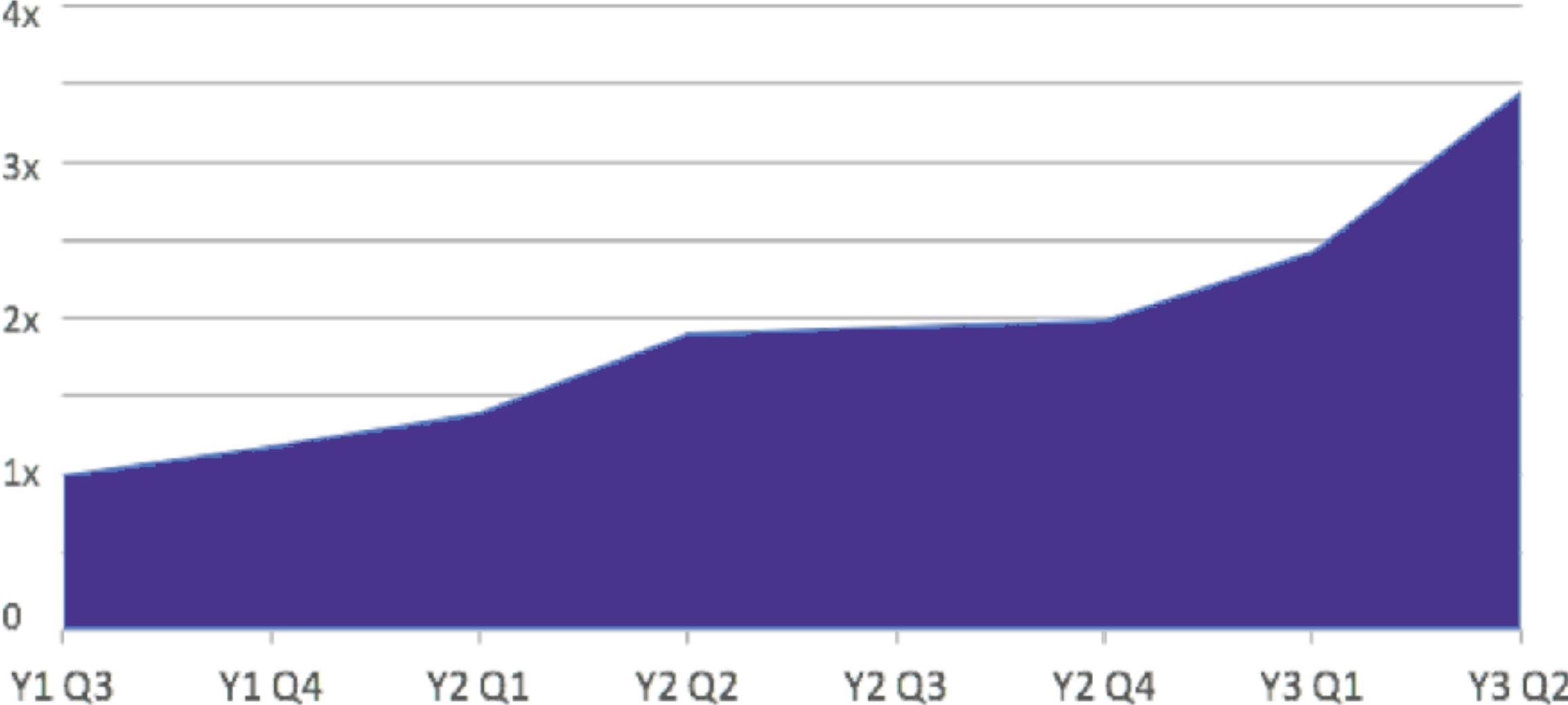
Deep Learning constitutes an important workload in cloud today.

Major cloud providers all have an ecosystem for cloud learning.

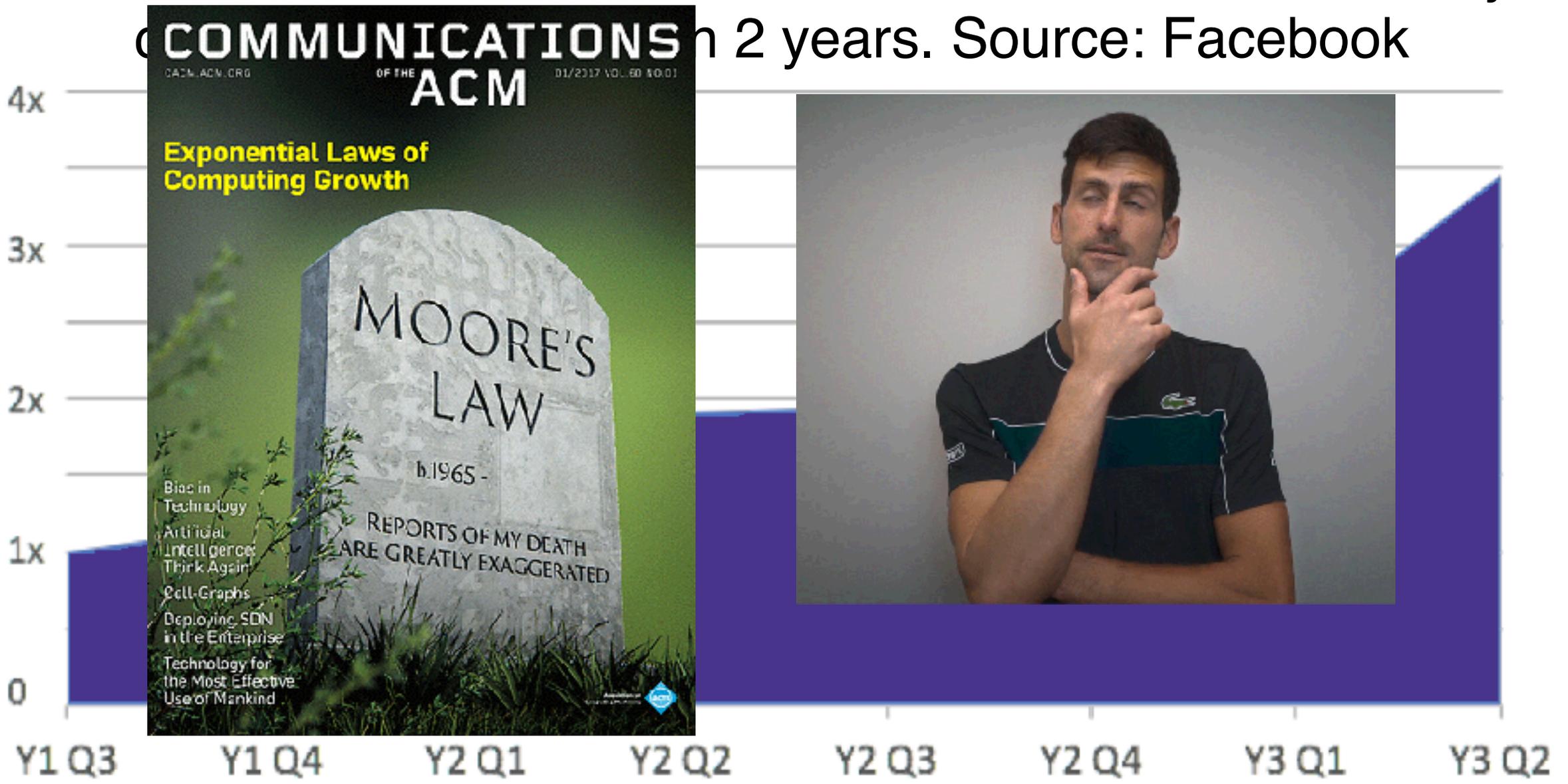
Deep Learning constitutes an important workload in cloud today.

Major cloud providers all have an ecosystem for cloud learning.

# Server demand for DL inference across data centers nearly quadrupled in less than 2 years. Source: Facebook



Server demand for DL inference across data centers nearly 4x in 2 years. Source: Facebook



si-q8389f8m	instance	g3.8xlarge	closed	i-084c42d56d09...	instance terminated no c...	one time	27 minutes ago	\$2.28
si-ffbga5fp	instance	g3.0xlarge	closed	i-Cf040f9bdGa70...	instance-terminated-no-c...	one-time	27 minutes ago	\$2.20
si-7ierb21m	instance	g3.8xlarge	closed	i-0d8444616a2e...	instance-terminated-capa...	one-time	an hour ago	\$2.28
si-pwvia6bq	instance	g3.8xlarge	closed	i-09dd933e7856...	instance terminated no c...	one time	an hour ago	\$2.28
si-rwqradim	instance	g3.0xlarge	closed	i-Cf149c2c1505	instance-terminated-no-c	one-time	an hour ago	\$2.20
si-r1bi81mq	instance	g3.8xlarge	closed	i-087784d7ad8b...	instance-terminated-capa...	one-time	an hour ago	\$2.28
si-33xi8qxm	instance	g3.8xlarge	closed	i-0fb0b38522a0...	instance-terminated-capa...	one-time	an hour ago	\$2.28
si-vfrdradsn	instance	g3.8xlarge	closed	i-009cd0d289b86	instance-terminated-capa	one-time	an hour ago	\$2.28
si-94er9wmn	instance	g3.8xlarge	closed	i-0208ba570400...	Instance-terminated-capa...	one-time	an hour ago	\$2.28
si-bmc8brsq	instance	g3.8xlarge	closed	i-0c16a6cb261b...	instance-terminated-capa...	one-time	an hour ago	\$2.28
si-k3rgamdj	instance	g3.8xlarge	closed	i-0e565c2ccc68	instance-terminated-no-c	one-time	an hour ago	\$2.28
si-fgnl9wan	Instance	g3.8xlarge	closed	i-07090471bd76...	Instance-terminated-no-c...	one-time	an hour ago	\$2.28
si-wzvra82n	instance	g3.8xlarge	closed	i-0c57cac2c5c8...	instance-terminated-capa...	one-time	an hour ago	\$2.28
si-ex789i4p	instance	g3.8xlarge	closed	i-03b1e92f296d	instance-terminated-capa	one-time	an hour ago	\$2.28
si-vvhran2q	Instance	g3.8xlarge	closed	i-087663582f2f2...	Instance-terminated-capa...	one-time	an hour ago	\$2.28
si-irk87nq	instance	g3.8xlarge	closed	i-0d29ad96ca42...	instance-terminated-no-c...	one-time	an hour ago	\$2.28
si-w2eg88pm	instance	g3.8xlarge	closed	i-079a78d77ea0	instance-terminated-no-c	one-time	an hour ago	\$2.28
si-mpghb1rq	Instance	g3.8xlarge	closed	i-0d12485d32f3...	Instance-terminated-capa...	one-time	an hour ago	\$2.28

EC2 reclaims your GPU instances as they run out of capacity

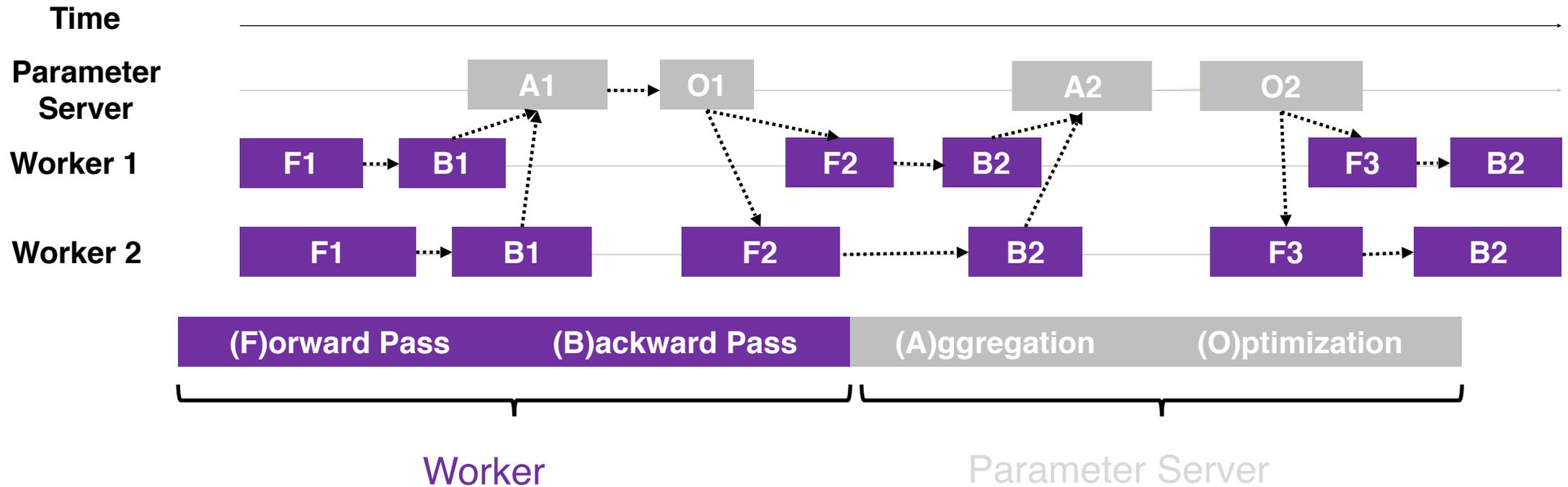
si-q8389r8m	instance	g3.8xlarge	closed	i-084c42d56d09...	instance-terminated-no-c...	one-time	27 minutes ago	\$2.28
si-ffbga5fp	instance	g3.0xlarge	closed	i-Cf040f9bdGa70...	instance-terminated-no-c...	one-time	27 minutes ago	\$2.20
si-7ierb21m	instance	g3.8xlarge	closed	i-Cd8444616a2e...	instance-terminated-capa...	one-time	an hour ago	\$2.28
si-pwvia6bq	instance						an hour ago	\$2.28
si-rwqradim	instance						an hour ago	\$2.20
si-r1bi81mq	instance						an hour ago	\$2.28
si-33x18qxm	instance						an hour ago	\$2.28
si-vfrdradsn	instance						an hour ago	\$2.28
si-94er9wmn	instance						an hour ago	\$2.28
si-bmc8brsq	instance						an hour ago	\$2.28
si-k3rgamdin	instance						an hour ago	\$2.28
si-fgn19wan	instance						an hour ago	\$2.28
si-wzvra82n	instance						an hour ago	\$2.28
si-ex789i4p	instance						an hour ago	\$2.28
si-vvhran2q	instance						an hour ago	\$2.28
si-irkr87nq	instance						an hour ago	\$2.28
si-w2eg88pm	instance	g3.8xlarge	closed	i-079a78d77ea0	instance-terminated-no-c...	one-time	an hour ago	\$2.28
si-mphgb1rq	instance	g3.8xlarge	closed	i-Cd12185d32f3...	Instance-terminated-capa...	one-time	an hour ago	\$2.28



EC2 reclaims your GPU instances as they run out of capacity

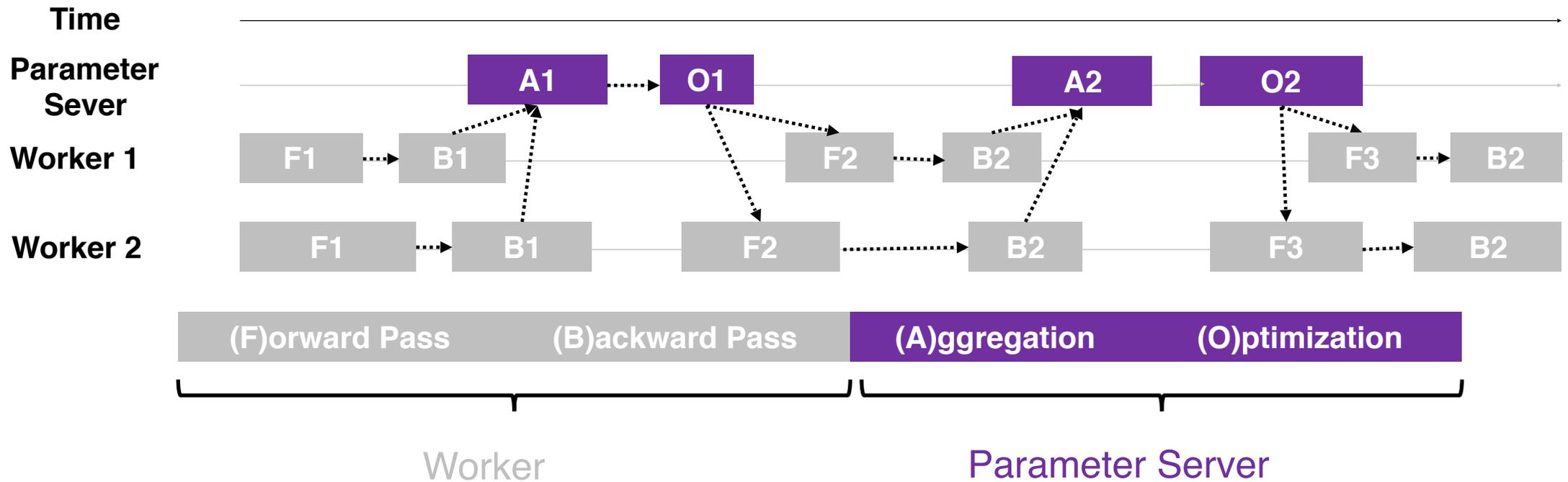
# Distributed Training

INDEPENDENT FORWARD/BACKWARD PASSES +  
COORDINATED PARAMETER EXCHANGE



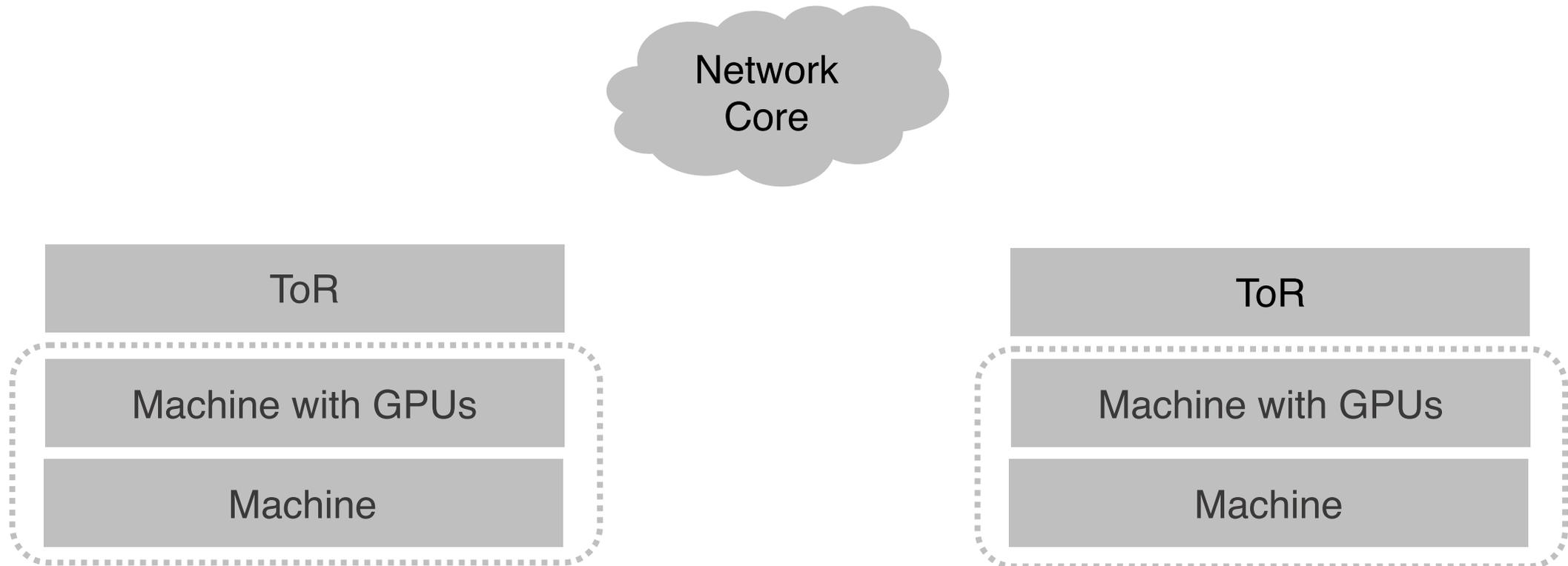
# Distributed Training

INDEPENDENT FORWARD/BACKWARD PASSES +  
COORDINATED PARAMETER EXCHANGE



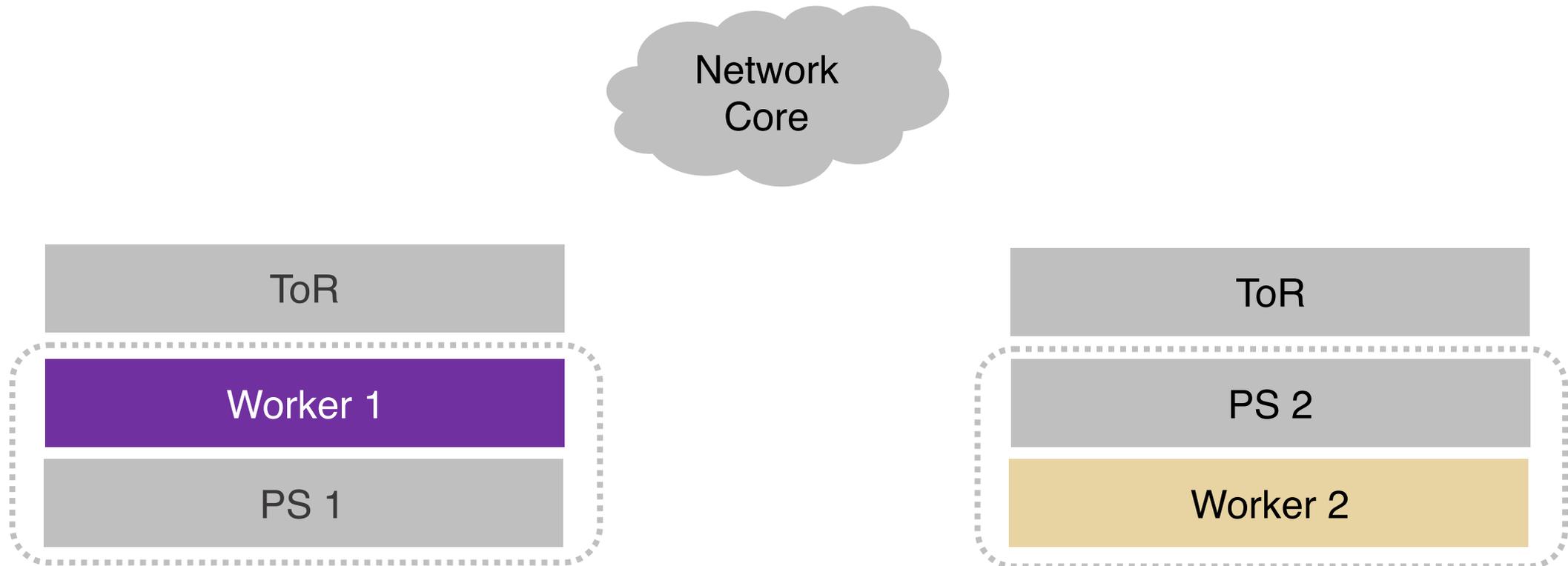
# Distributed Training Today

## IN THE CONTEXT OF THE CLOUD



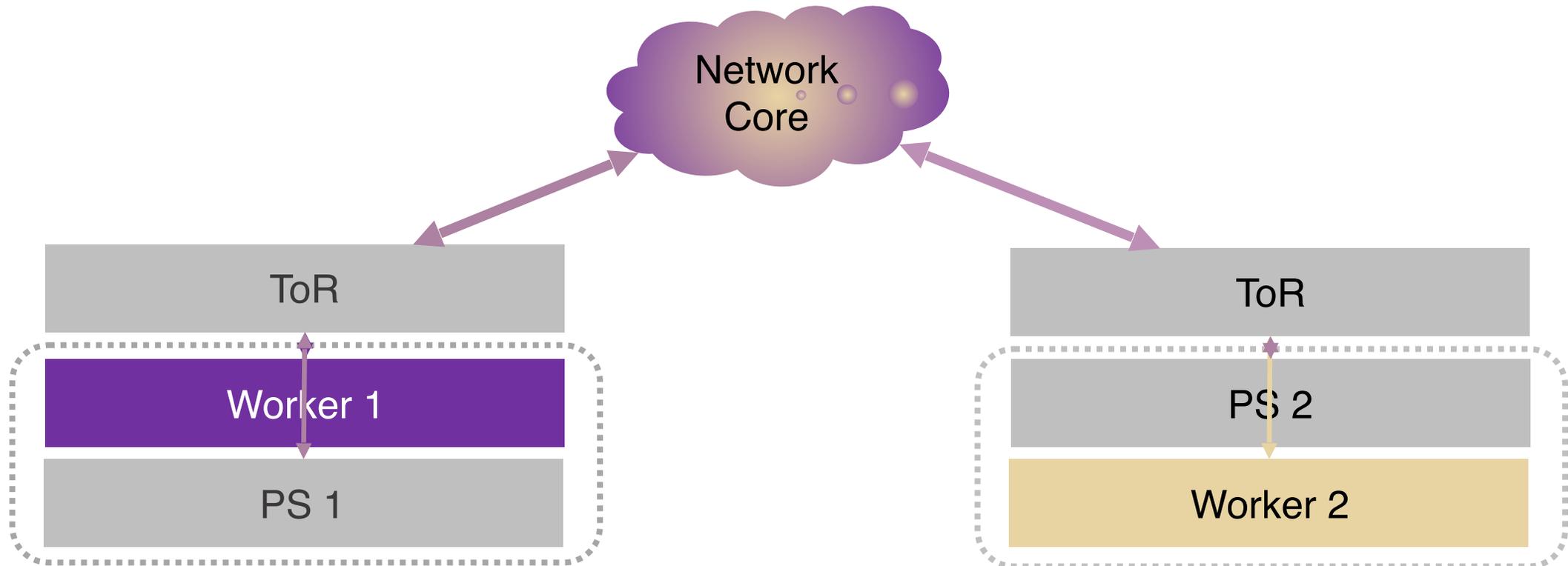
# Distributed Training Today

FORWARD AND BACKWARD PASSES IN WORKER



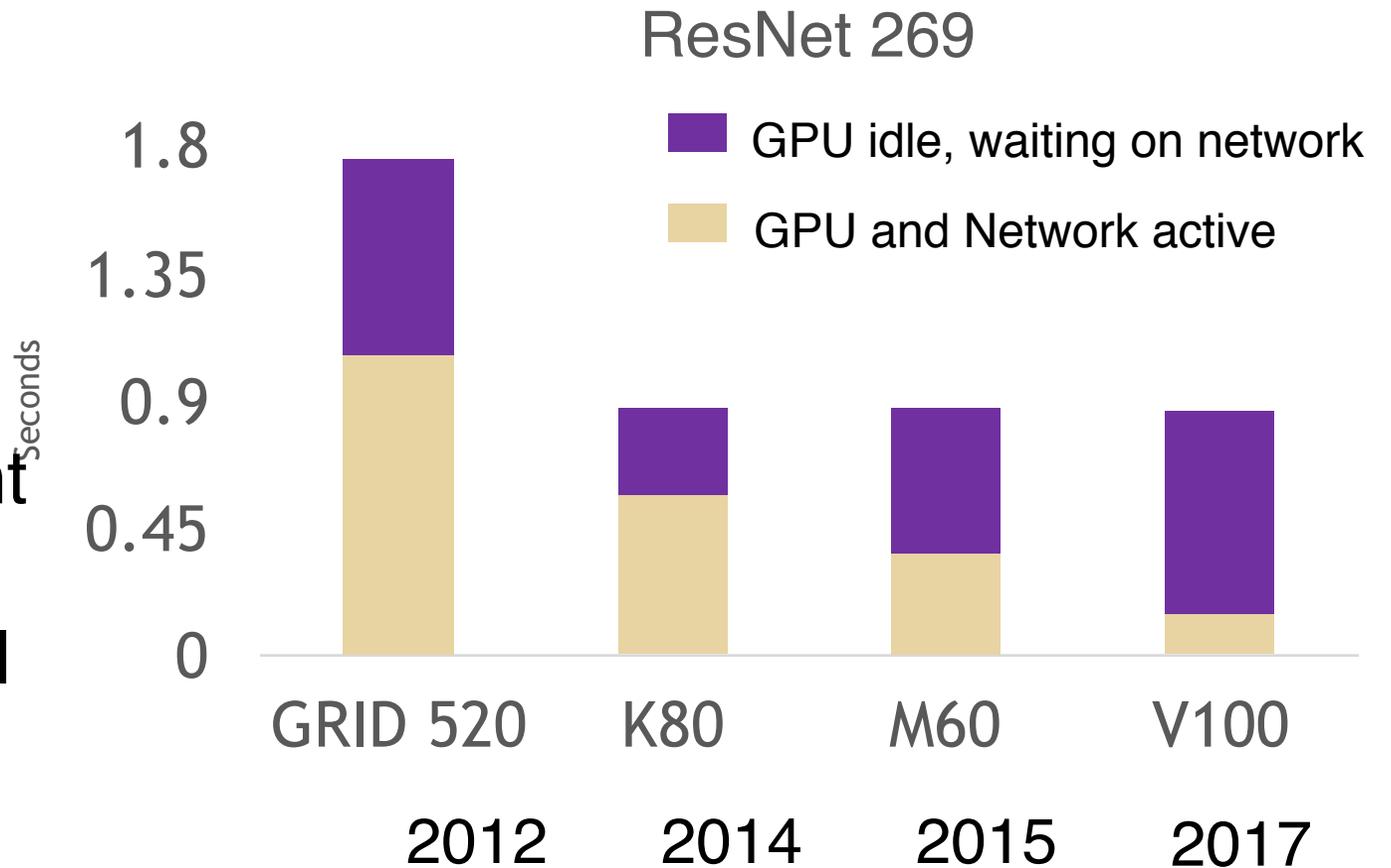
# Distributed Training Today

## AGGREGATION AND OPTIMIZATION IN PS



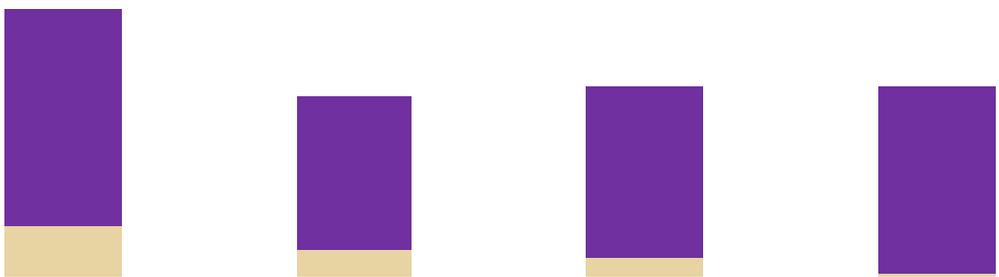
# Distributed training is communication bound

- Problem gets worse over time: shifting bottleneck.
- With modern GPUs most of the time is spent on **communication**.
- Making GPUs faster will **do little to increase throughput**
- Wasting compute resources.

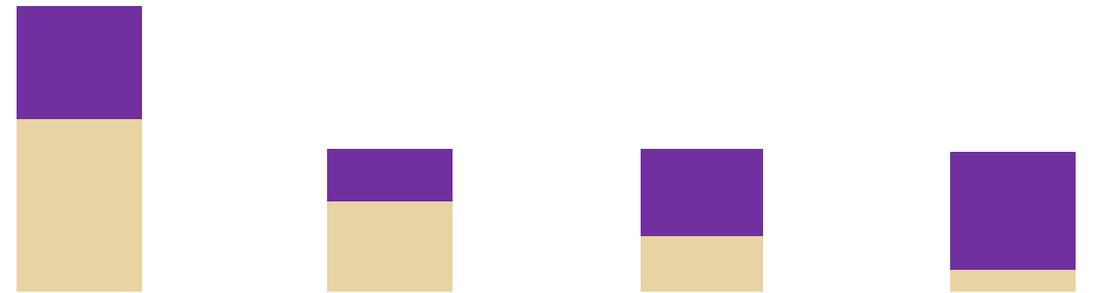


# Distributed training is communication bound

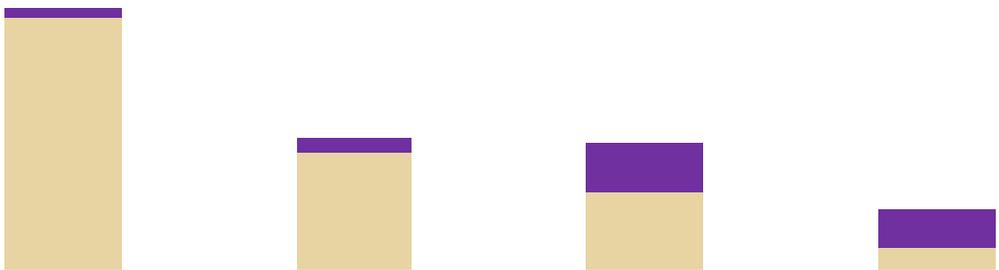
AlexNet



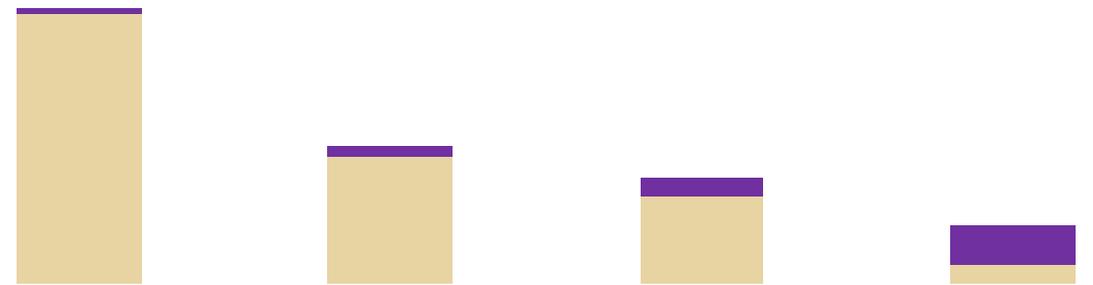
ResNet 269



Inception V3

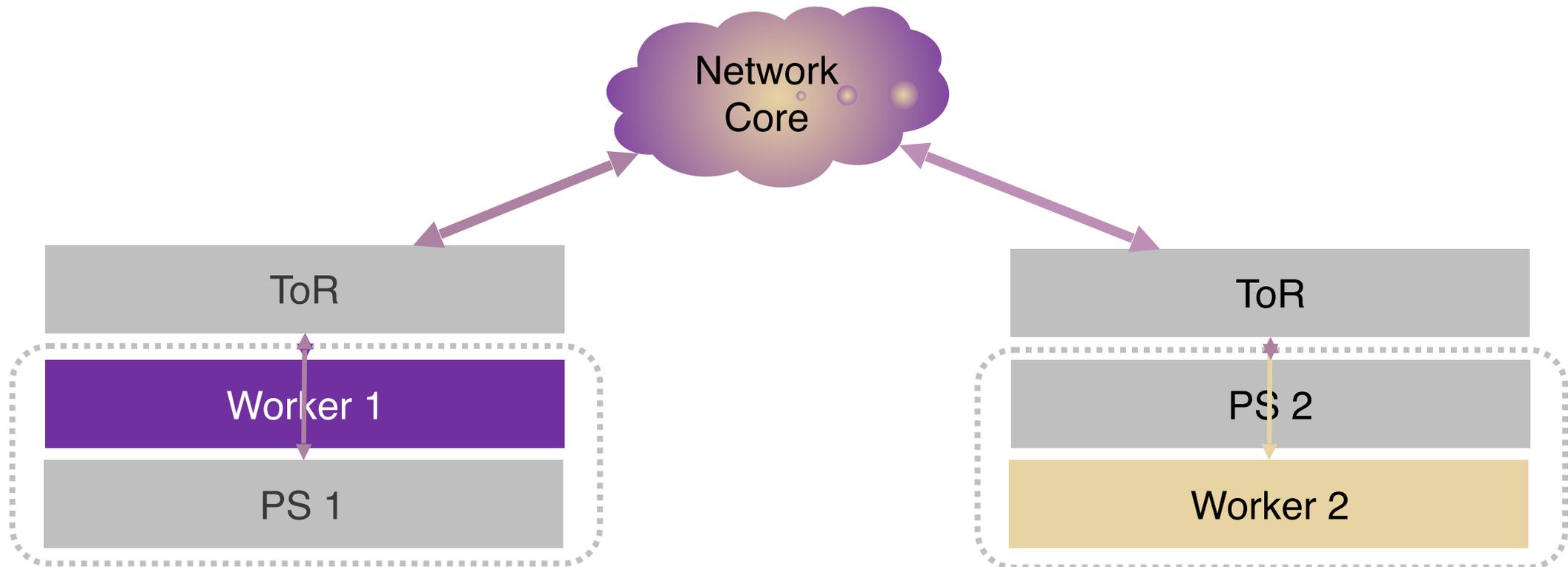


GoogleNet



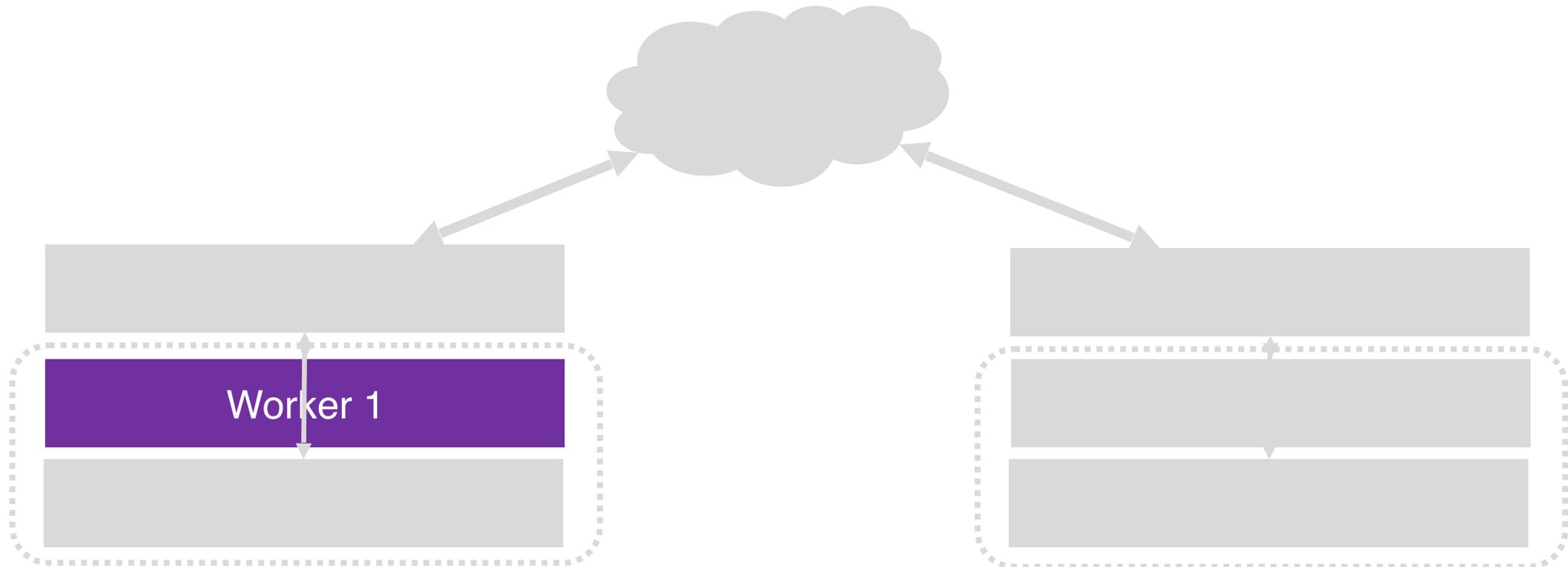
# Bottlenecks in DDNN training

MAPPING OF TRAINING WORKLOAD TO THE CLOUD IS INEFFICIENT.



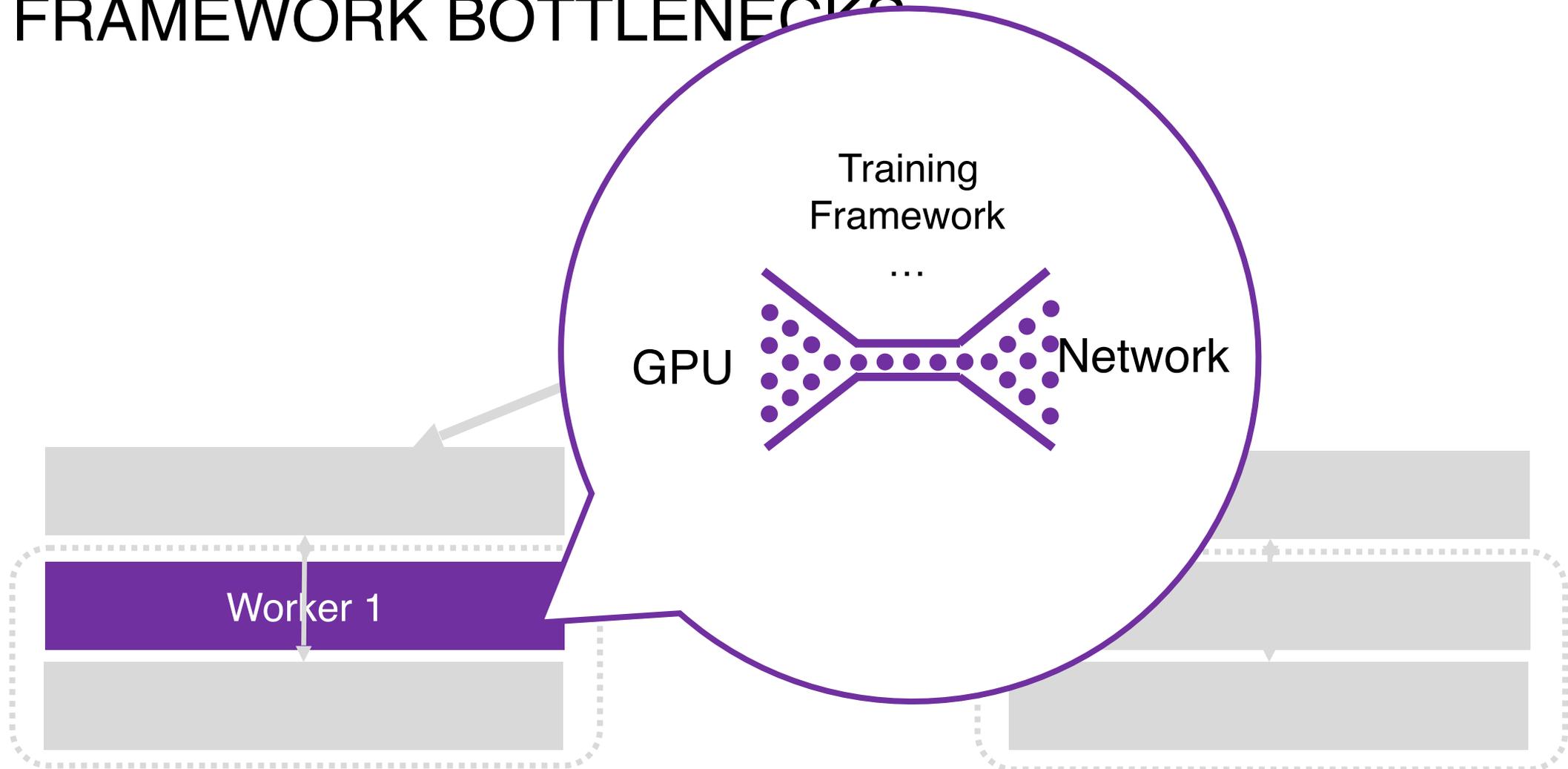
# Bottlenecks in DDNN training

## FRAMEWORK BOTTLENECKS



# Bottlenecks in DDNN training

## FRAMEWORK BOTTLENECKS

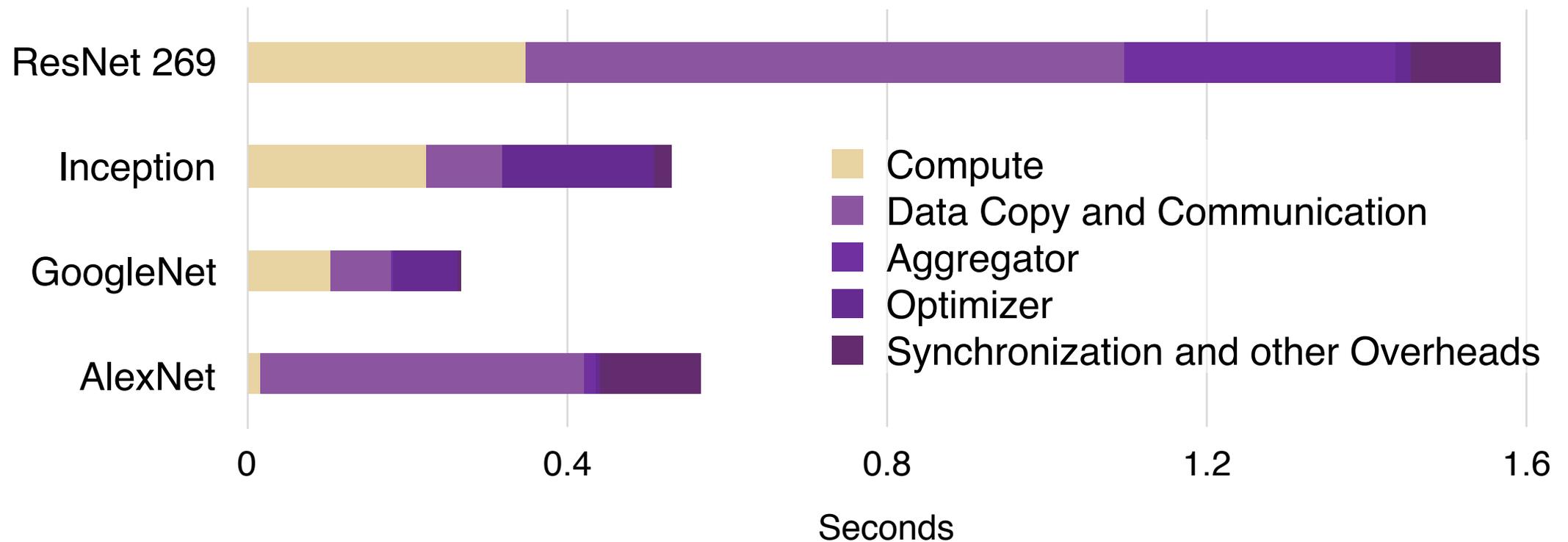


# **Bottlenecks in DDNN training**

## FRAMEWORK BOTTLENECKS

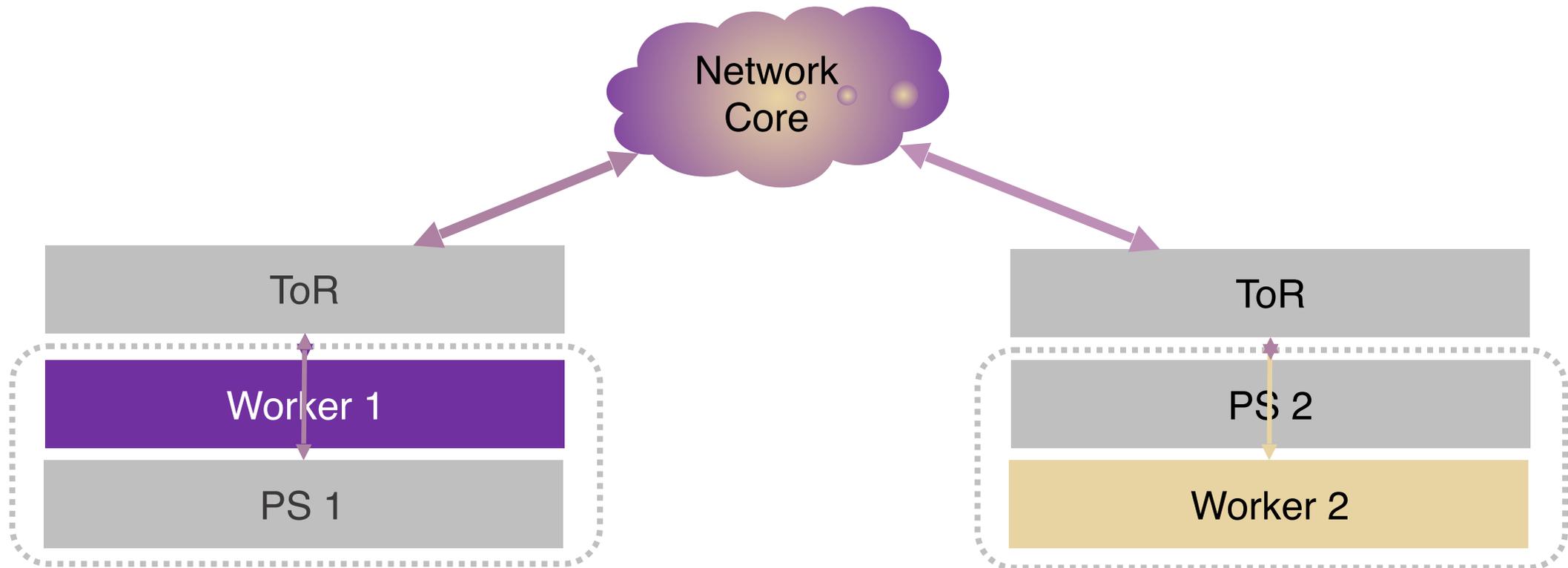
# Bottlenecks in DDNN training

## FRAMEWORK BOTTLENECKS



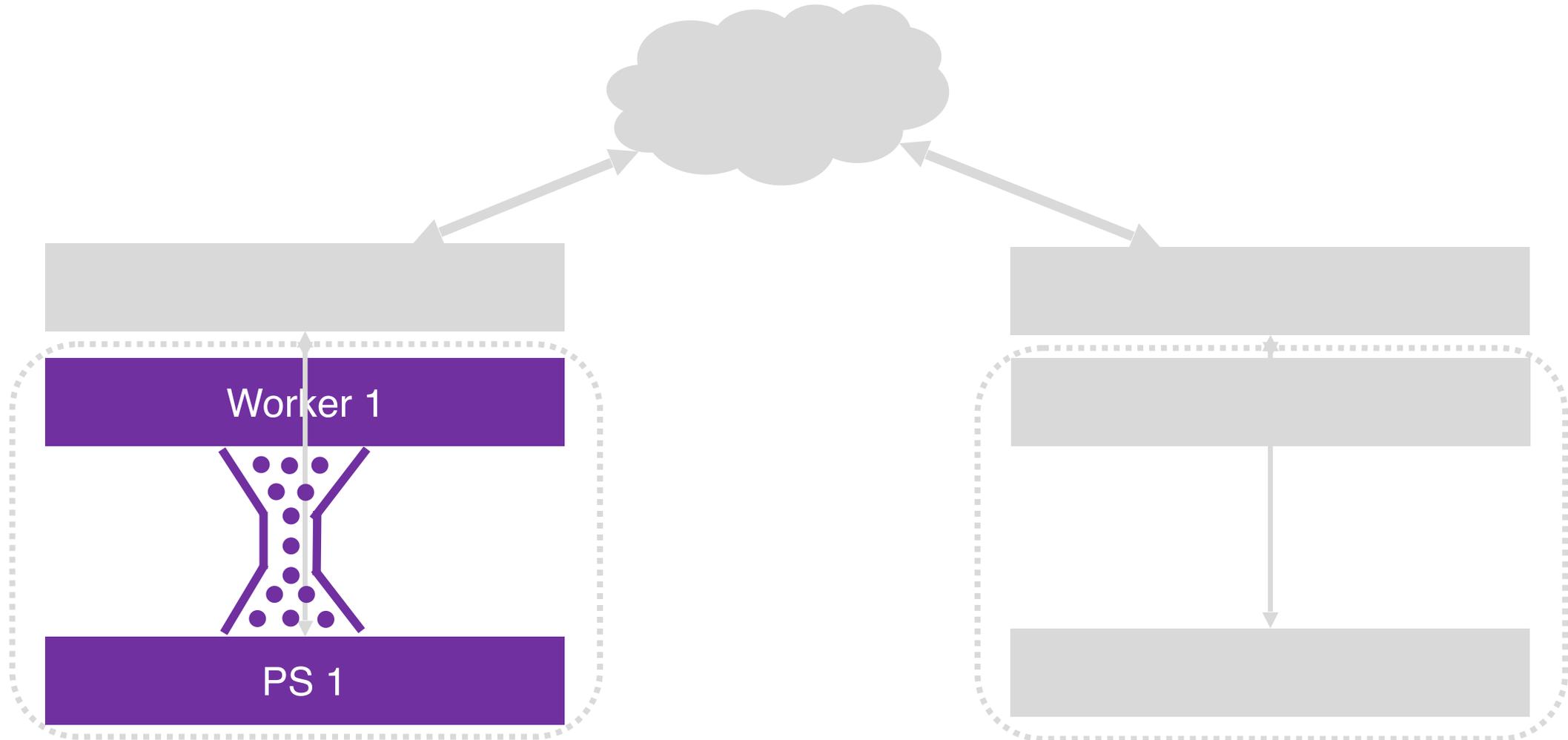
# Bottlenecks in DDNN training

MAPPING OF TRAINING WORKLOAD TO THE CLOUD IS INEFFICIENT.



# Bottlenecks in DDNN training

## BANDWIDTH BOTTLENECK



# Bottlenecks in Cloud-based DDNN training

## INSUFFICIENT BANDWIDTH

Minimum bandwidth required for each of the popular NNs for communication to not bottleneck computation?

8 workers, GTX 1080 Ti,  
central parameter servers.  
MxNet

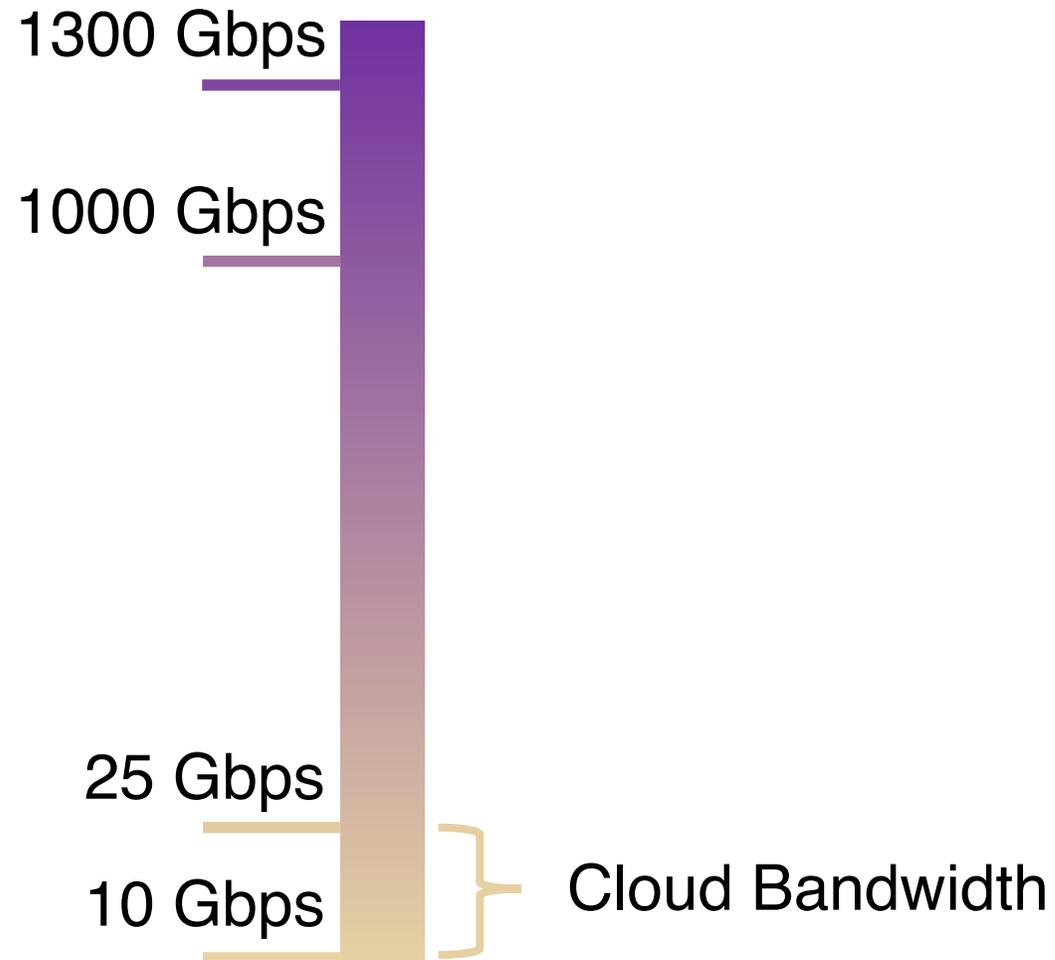


# Bottlenecks in Cloud-based DDNN training

## INSUFFICIENT BANDWIDTH

Minimum bandwidth required for each of the popular NNs for communication to not bottleneck computation?

8 workers, GTX 1080 Ti,  
central parameter servers.  
MxNet

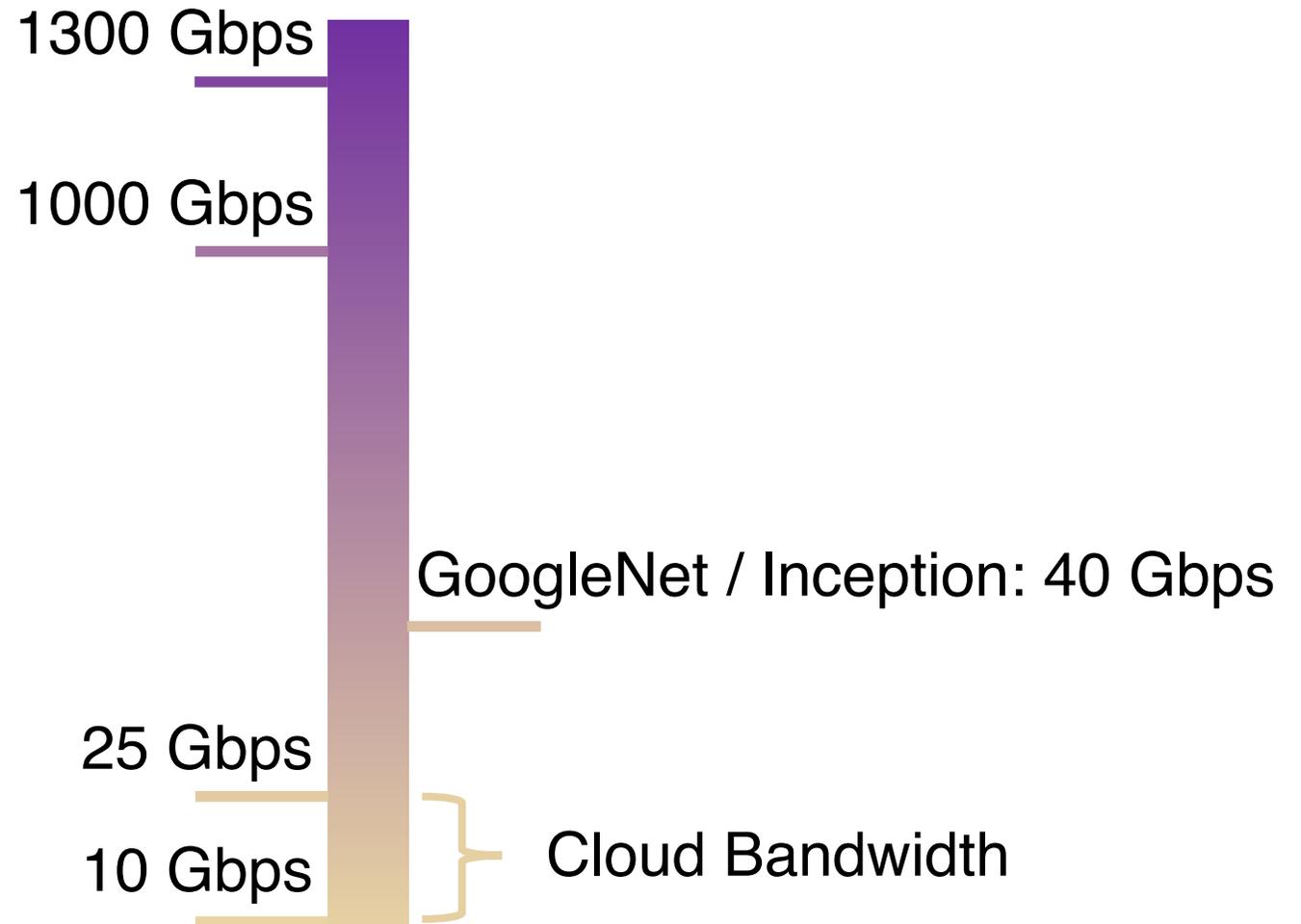


# Bottlenecks in Cloud-based DDNN training

## INSUFFICIENT BANDWIDTH

Minimum bandwidth required for each of the popular NNs for communication to not bottleneck computation?

8 workers, GTX 1080 Ti,  
central parameter servers.  
MxNet

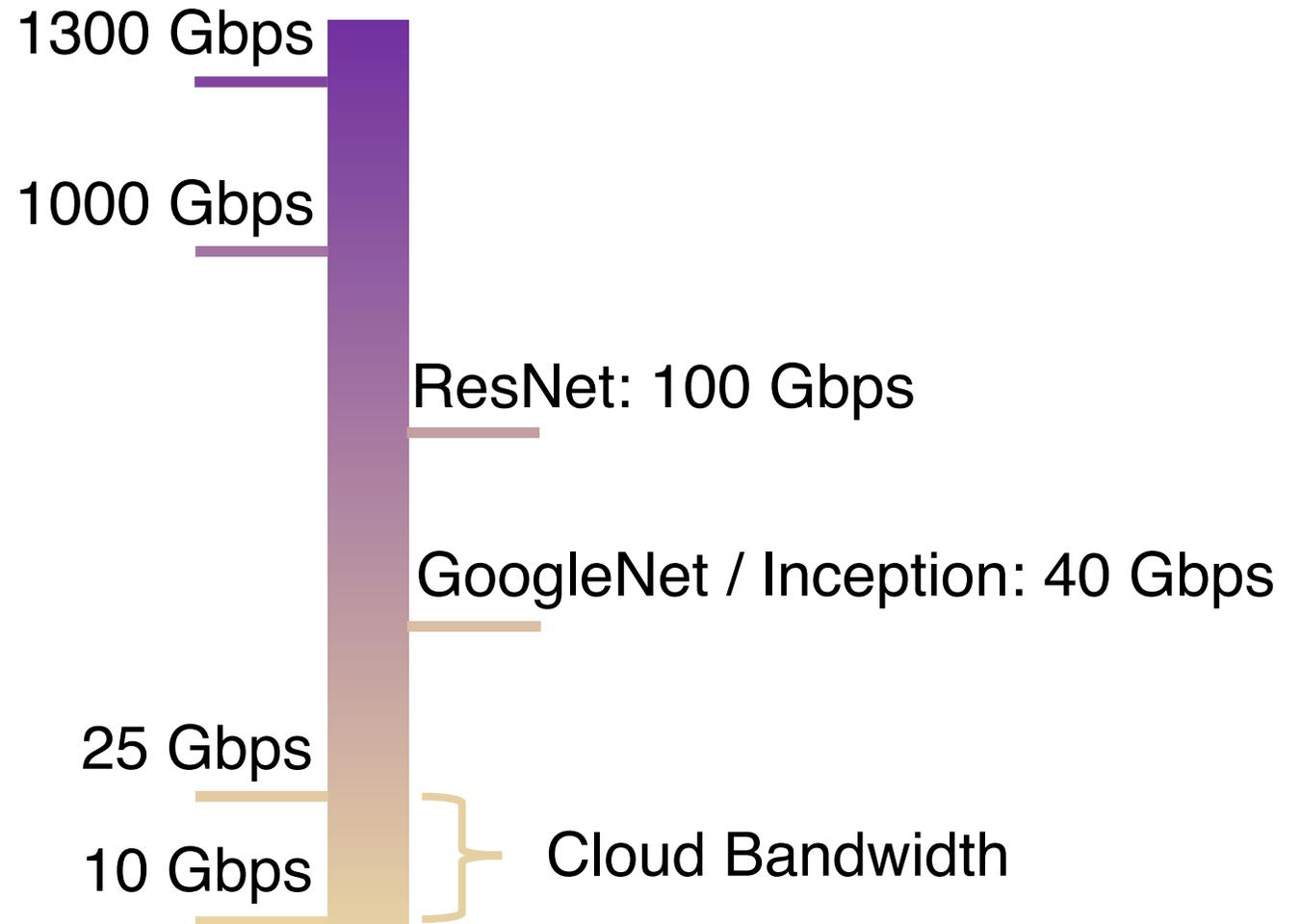


# Bottlenecks in Cloud-based DDNN training

## INSUFFICIENT BANDWIDTH

Minimum bandwidth required for each of the popular NNs for communication to not bottleneck computation?

8 workers, GTX 1080 Ti,  
central parameter servers.  
MxNet

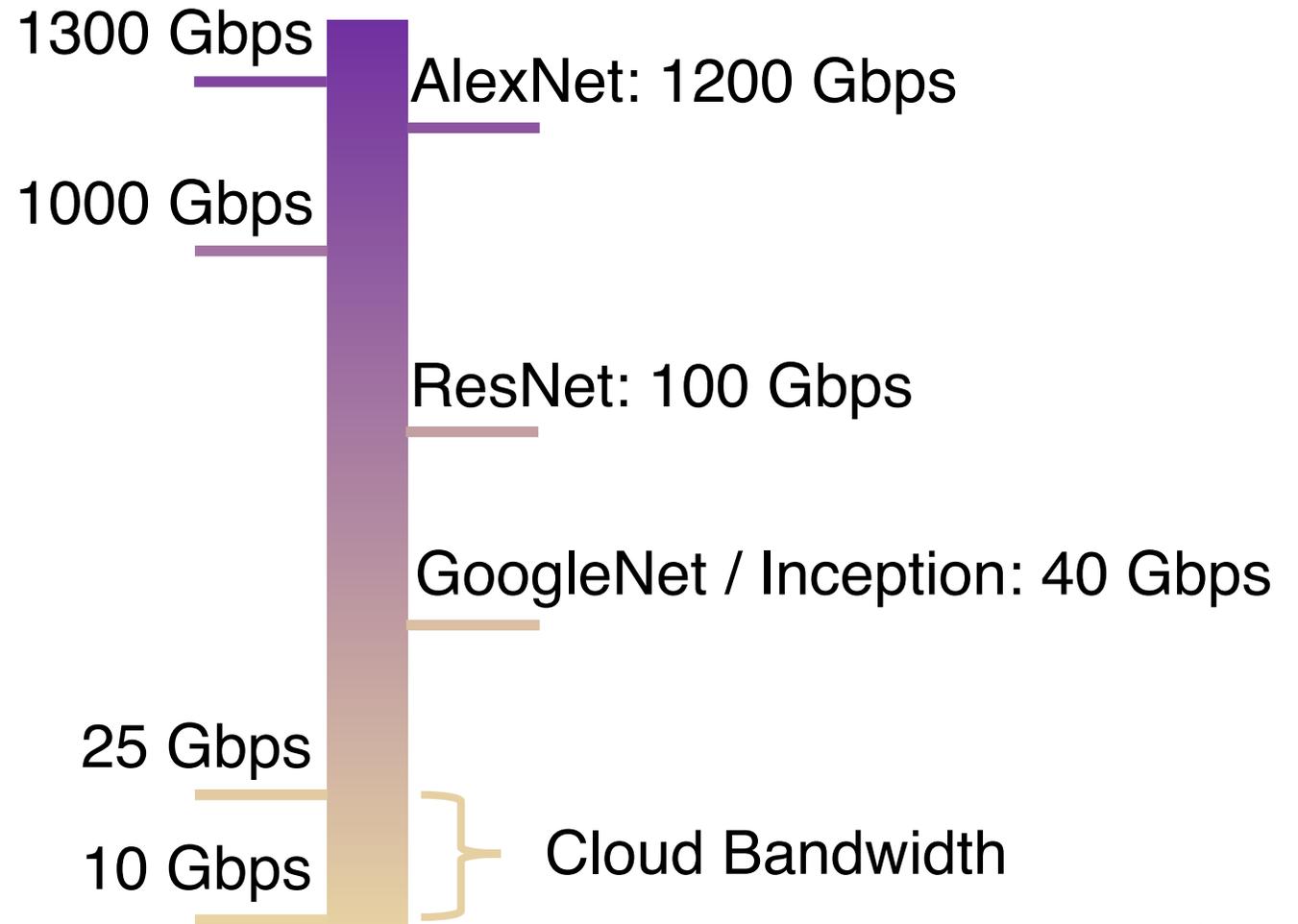


# Bottlenecks in Cloud-based DDNN training

## INSUFFICIENT BANDWIDTH

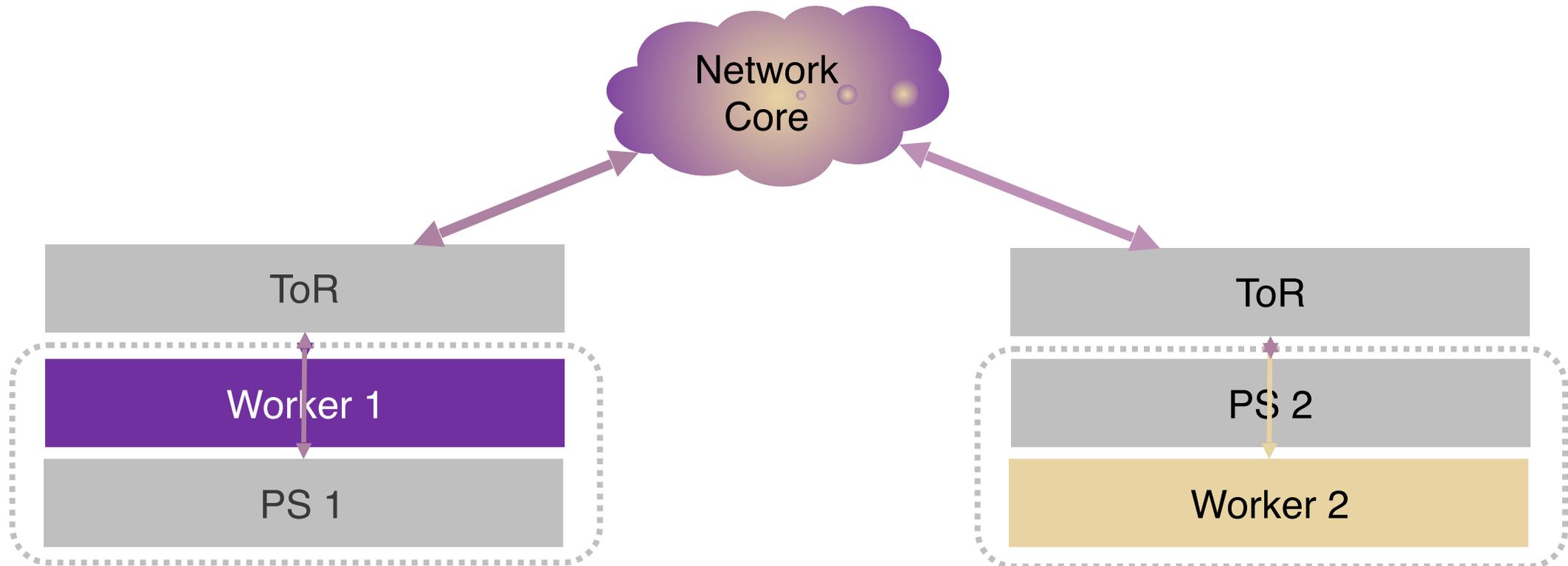
Minimum bandwidth required for each of the popular NNs for communication to not bottleneck computation?

8 workers, GTX 1080 Ti,  
central parameter servers.  
MxNet



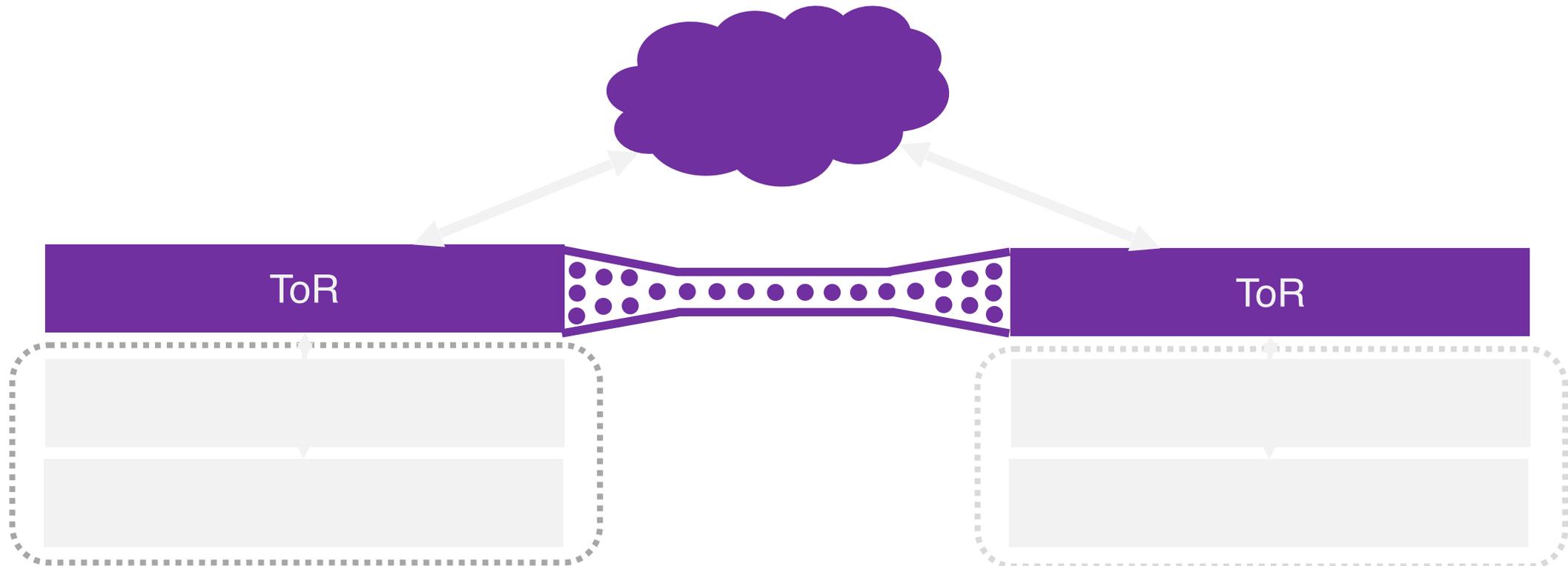
# Bottlenecks in Cloud-based DDNN training

MAPPING OF TRAINING WORKLOAD TO THE CLOUD IS INEFFICIENT.



# Bottlenecks in Cloud-based DDNN training

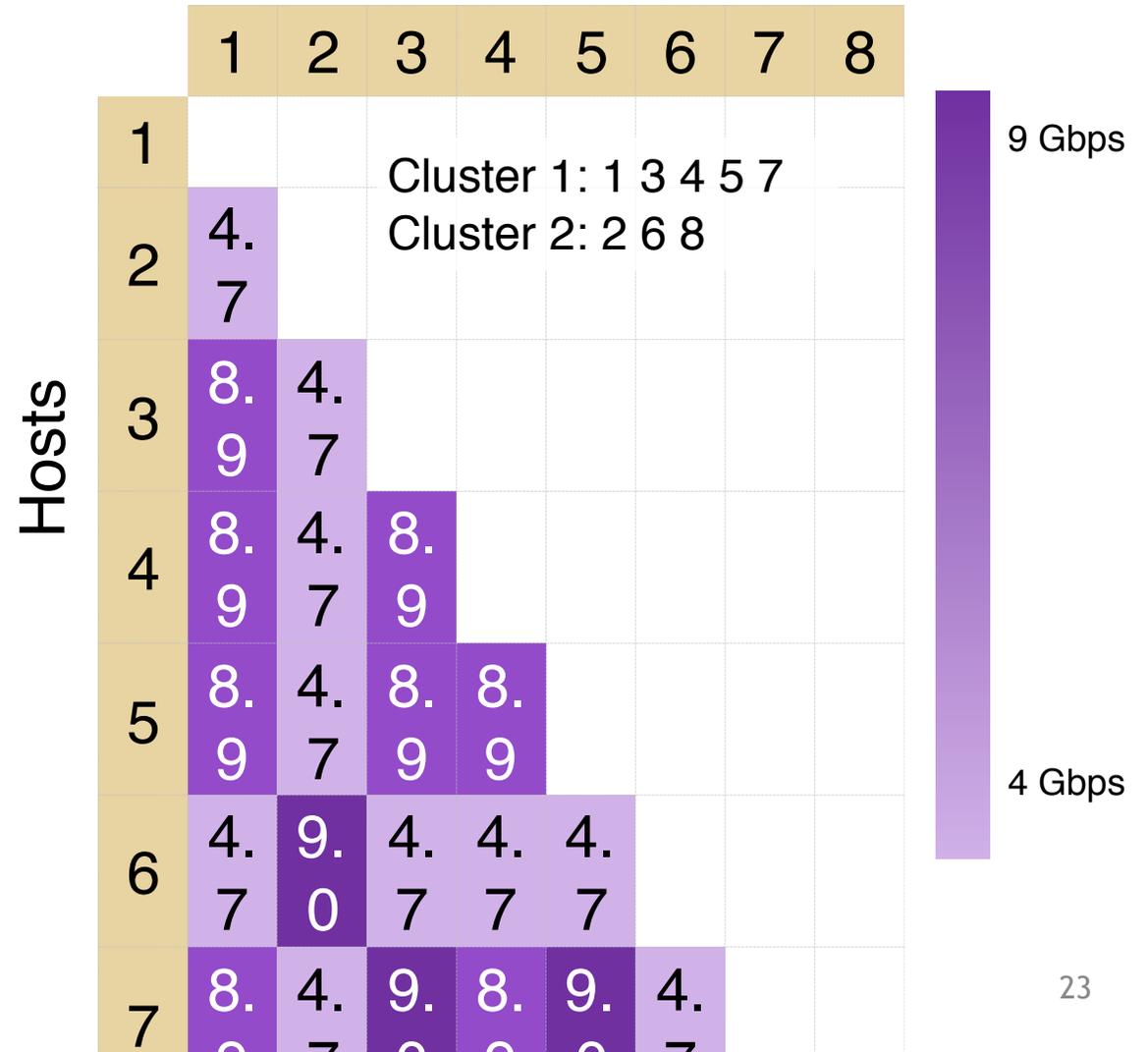
## DEPLOYMENT-RELATED OVERHEAD



# Bottlenecks in Cloud-based DDNN training

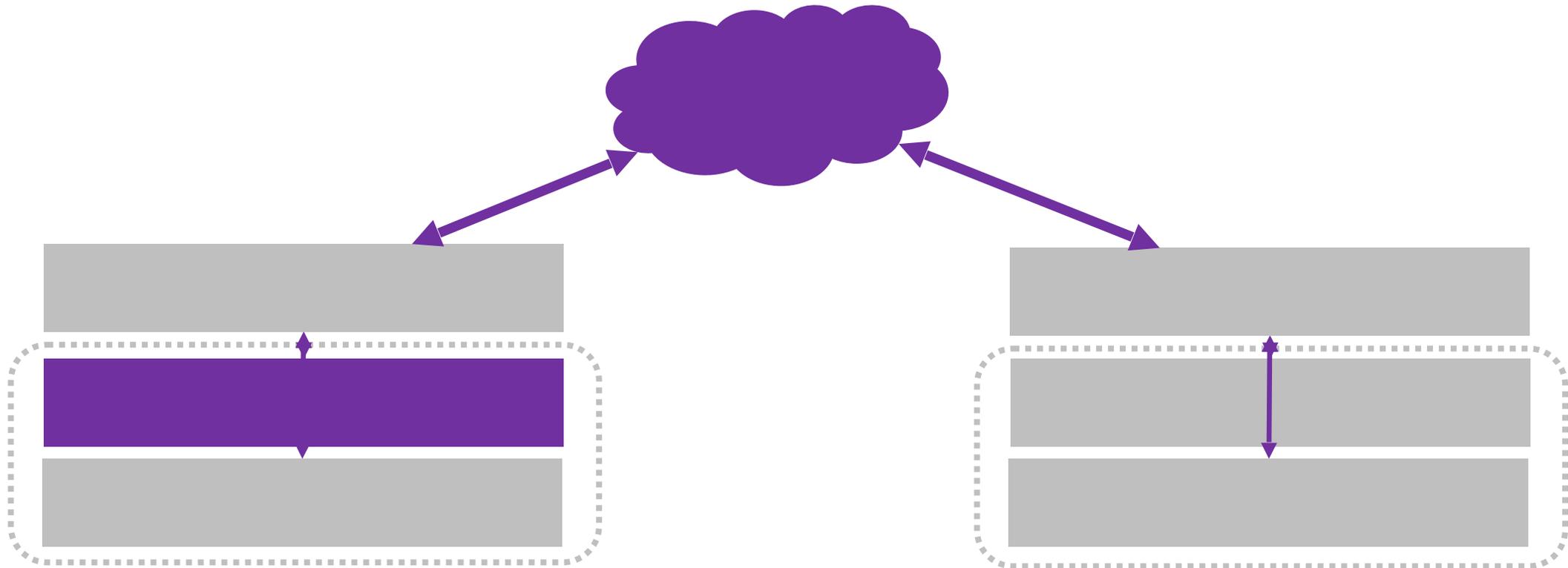
## DEPLOYMENT-RELATED OVERHEAD

- Transient congestion, or oversubscription by design
- Cross-rack communication cost is higher than Intra-rack communication.



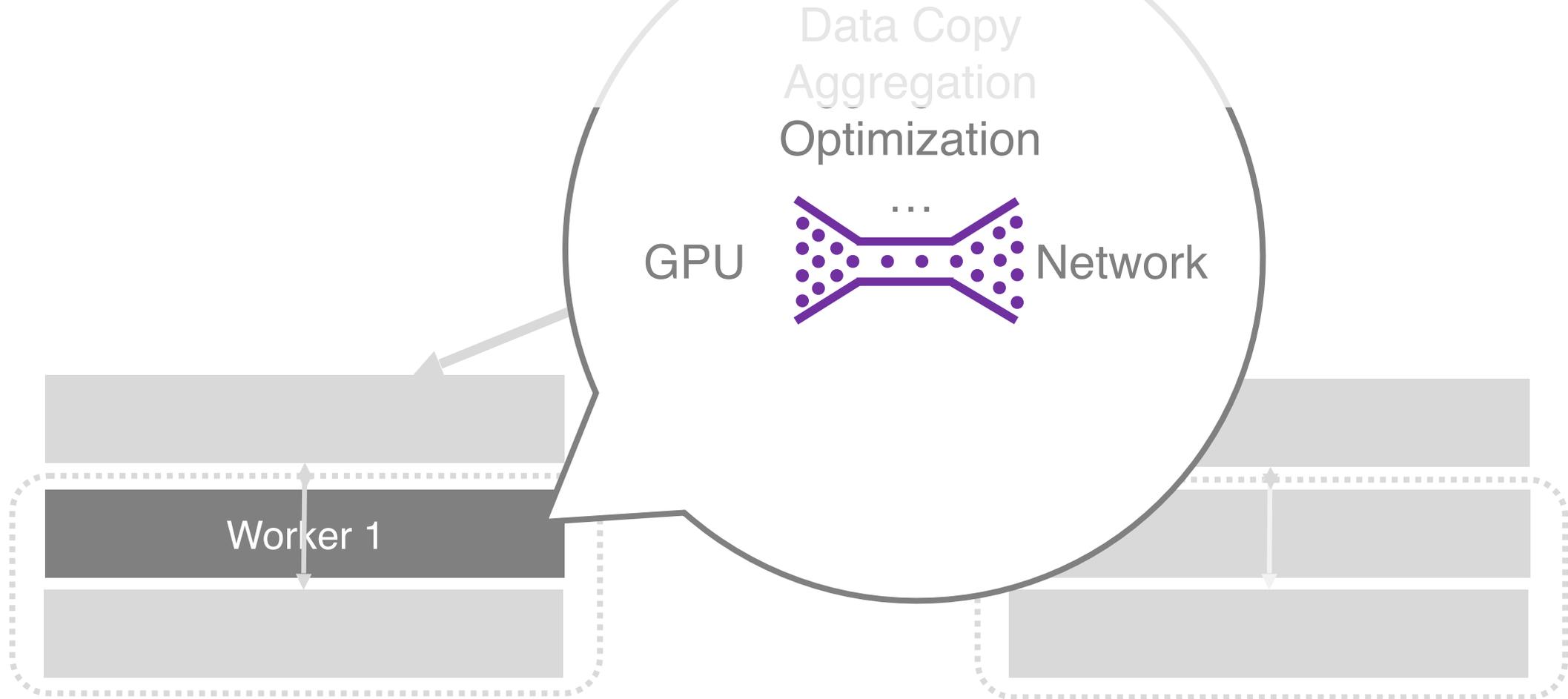
# Parameter Hub Optimizations

CODESIGNING SOFTWARE, HARDWARE WITH  
CLUSTER CONFIGURATION FOR EFFICIENT CLOUD-  
BASED DDNN TRAINING



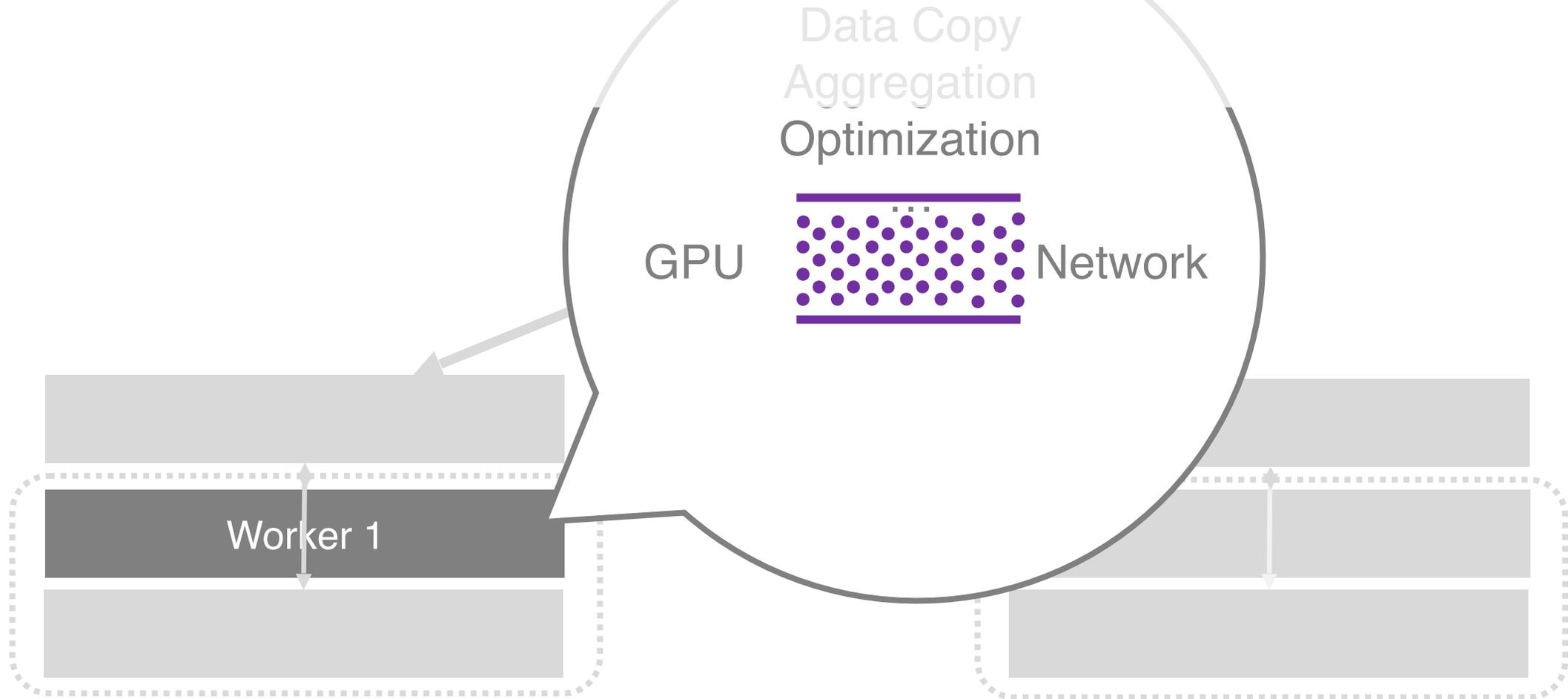
# Eliminating framework bottlenecks:

## PHub Optimizations: streamlining DDNN training pipeline

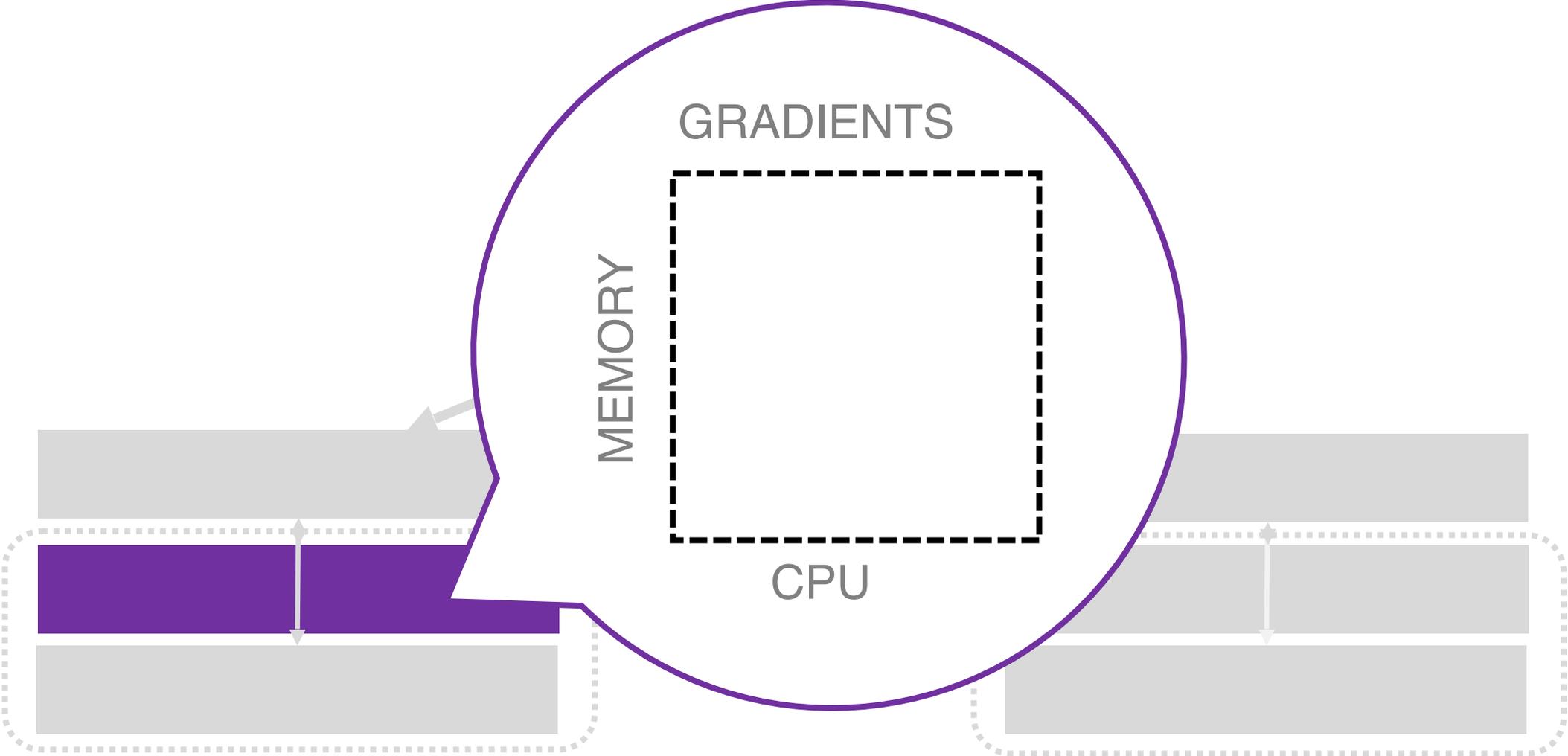


# Eliminating framework bottlenecks:

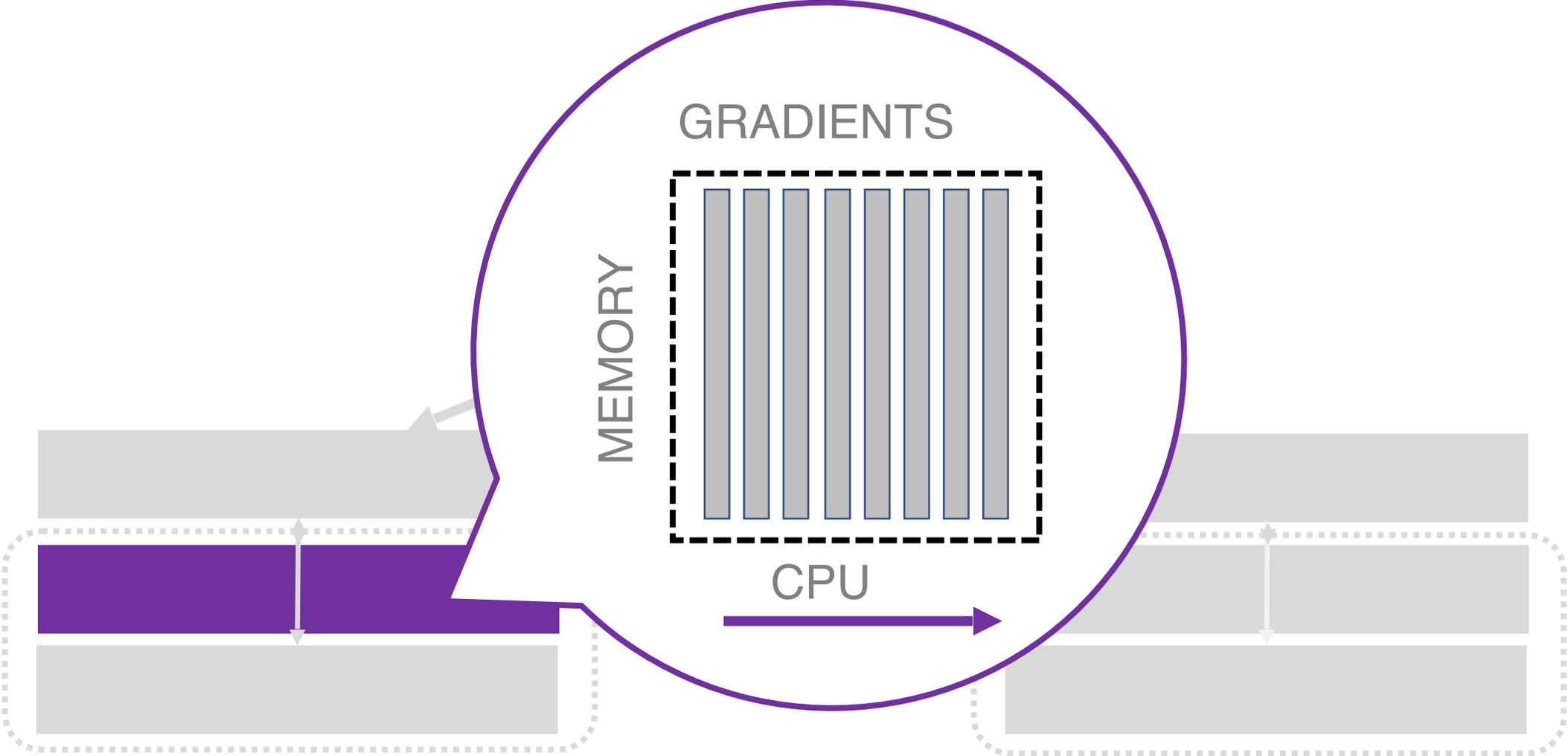
## PHub Optimizations: streamlining DDNN training pipeline



# Software Optimizations



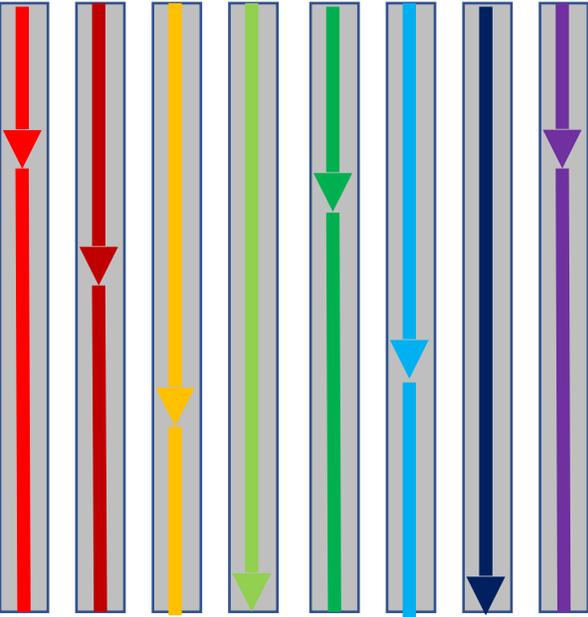
# Software Optimizations



# Software Optimizations

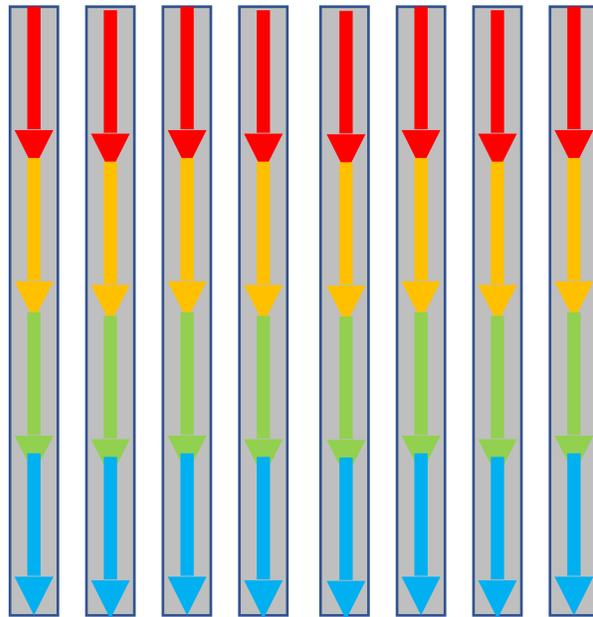
## GRADIENT AGGREGATION AND OPTIMIZATION

Requires synchronization.



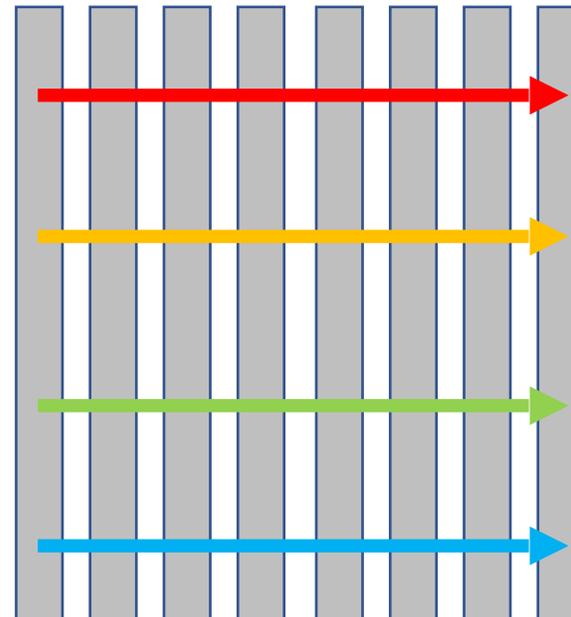
Each core reads the input Q from different workers and writes to different locations to the output queue

Great locality. No synchronization



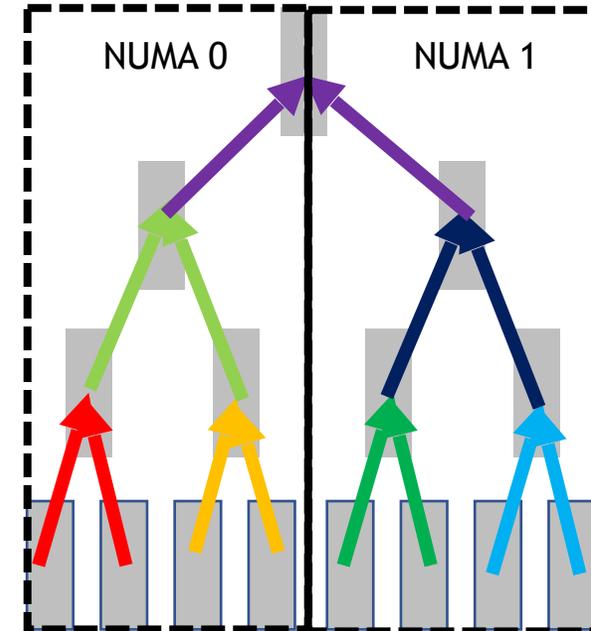
For each input Q, launch a series of threads for aggregation. This is used in MxNet. (Wide Aggregation)

Great locality. No synchronization



Sequentially aggregates the same portion of gradients within each queue. (Tall Aggregation)

Too much coherence and synchronization



Organize processors into hierarchy. Perform NUMA aware tree reduction.

# Software Optimizations

## GRADIENT AGGREGATION AND OPTIMIZATION

Requires synchronization.



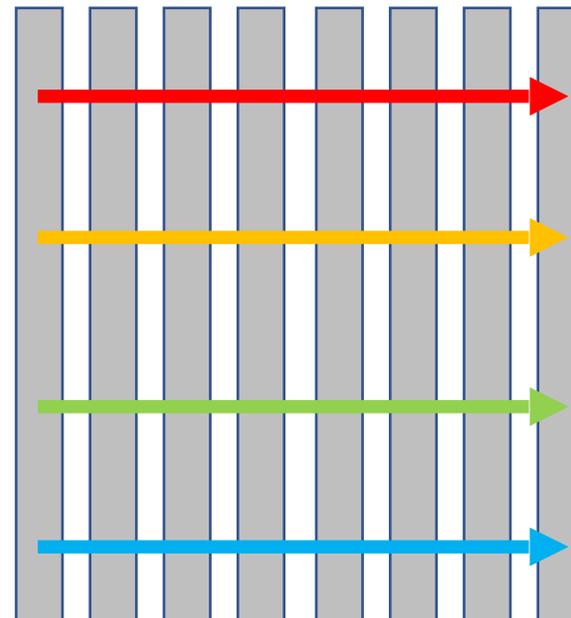
Each core reads the input  $Q$  from different workers and writes to different locations to the output queue

Great locality. No synchronization



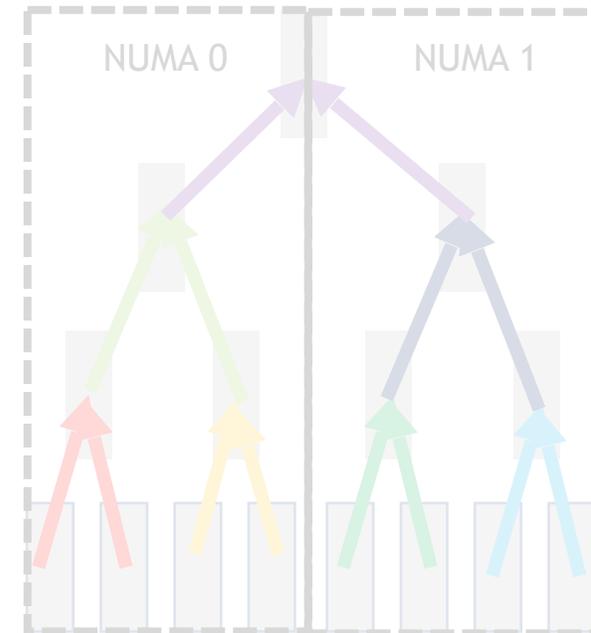
For each input  $Q$ , launch a series of threads for aggregation. This is used in MxNet. (Wide Aggregation)

Great locality. No synchronization



Sequentially aggregates the same portion of gradients within each queue. (Tall Aggregation)

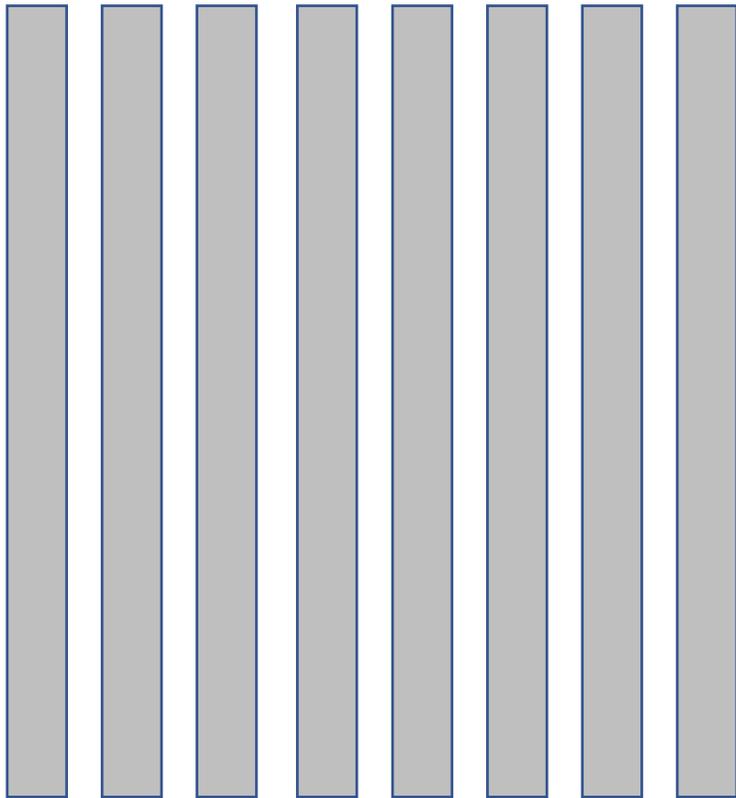
Too much coherence and synchronization



Organize processors into hierarchy. Perform NUMA aware tree reduction.

# Software Optimizations

## TALL AGGREGATION AND OPTIMIZATION

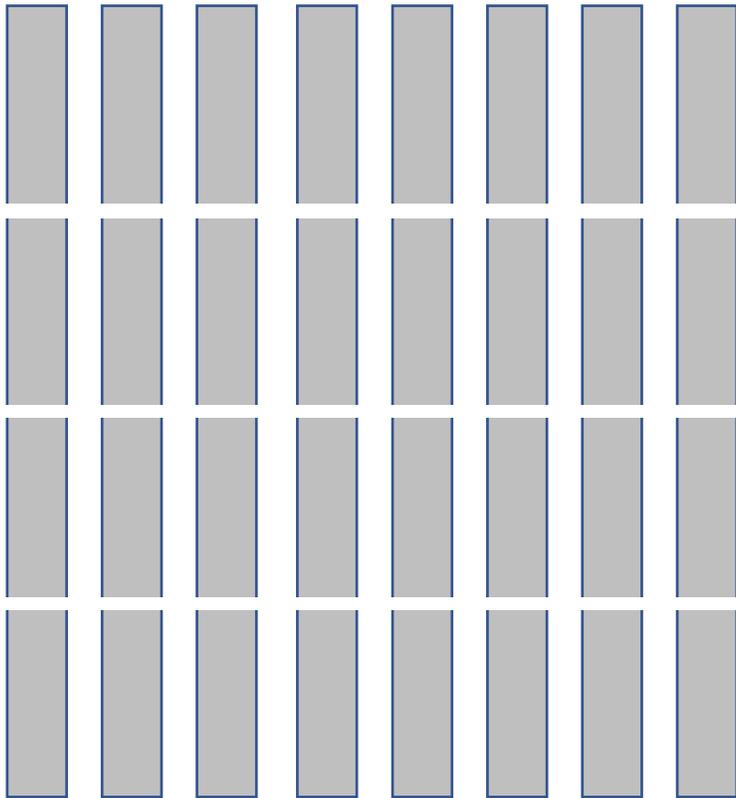


Gradient Array for Key 0 from 8 workers

- Chunk a gradient into a series of virtual gradients deterministically.
- A virtual gradient is mapped to a particular core on the server.

# Software Optimizations

## TALL AGGREGATION AND OPTIMIZATION

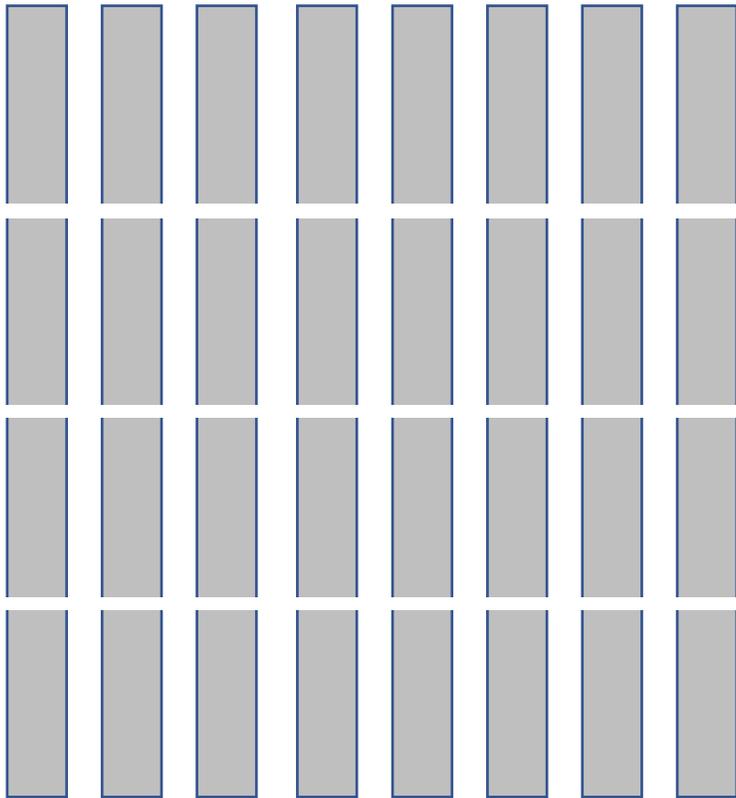


- Chunk a gradient into a series of virtual gradients deterministically.
- A virtual gradient is mapped to a particular core on the server.

Gradient Array for Key 0 from 8 workers

# Software Optimizations

## TALL AGGREGATION AND OPTIMIZATION

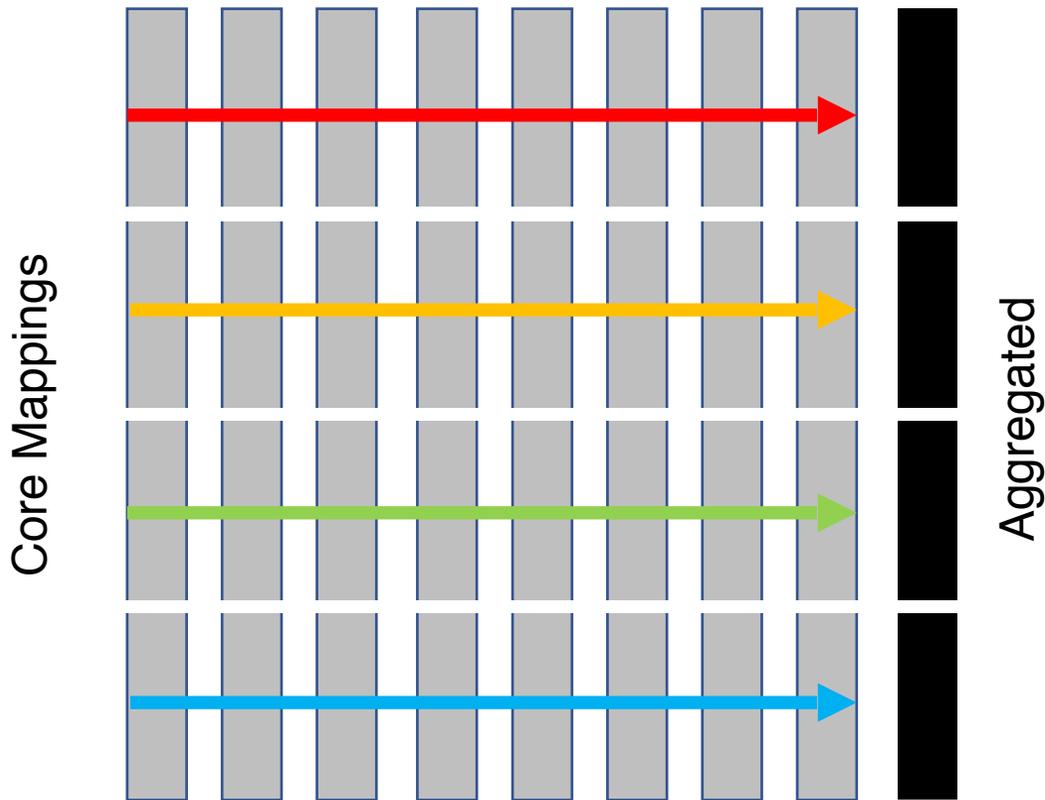


Gradient Array for Key 0 from 8 workers

- Chunk a gradient into a series of **virtual gradients deterministically**.
- A virtual gradient is mapped to a particular core on the server.
- Virtual gradients are transferred **independently**.

# Software Optimizations

## TALL AGGREGATION AND OPTIMIZATION

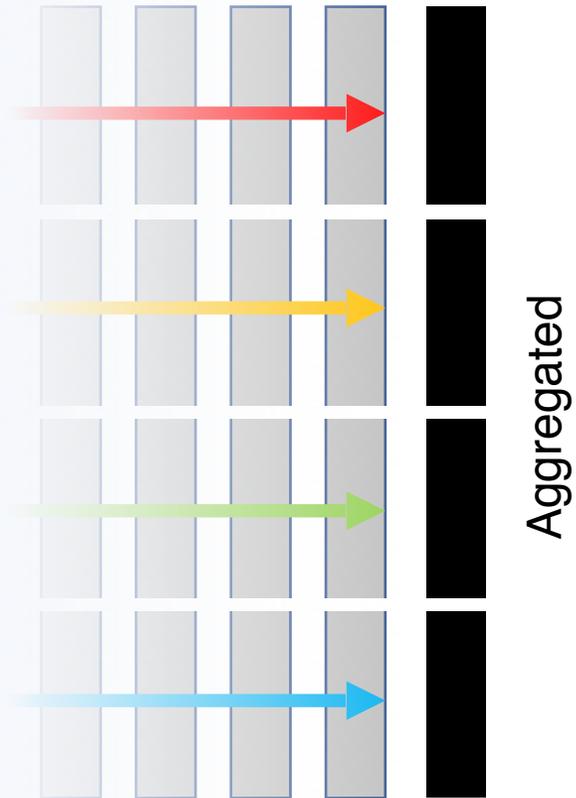


Gradient Array for Key 0 from 8 workers

- Chunk a gradient into a series of **virtual gradients deterministically**.
- A virtual gradient is mapped to a particular core on the server.
- Virtual gradients are transferred **independently**.
- A chunk is only processed by a single core : maintaining maximum locality.

# Software Optimizations

## TALL AGGREGATION AND OPTIMIZATION



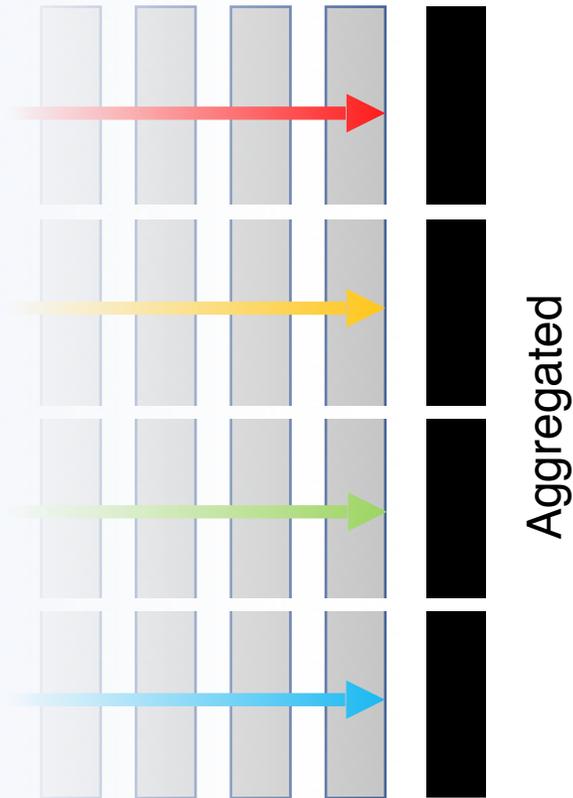
Key 0 from 8 workers

When Aggregation is done, PHub:

- PHub optimizes a chunk with the **same core** that aggregates that chunk.

# Software Optimizations

## TALL AGGREGATION AND OPTIMIZATION



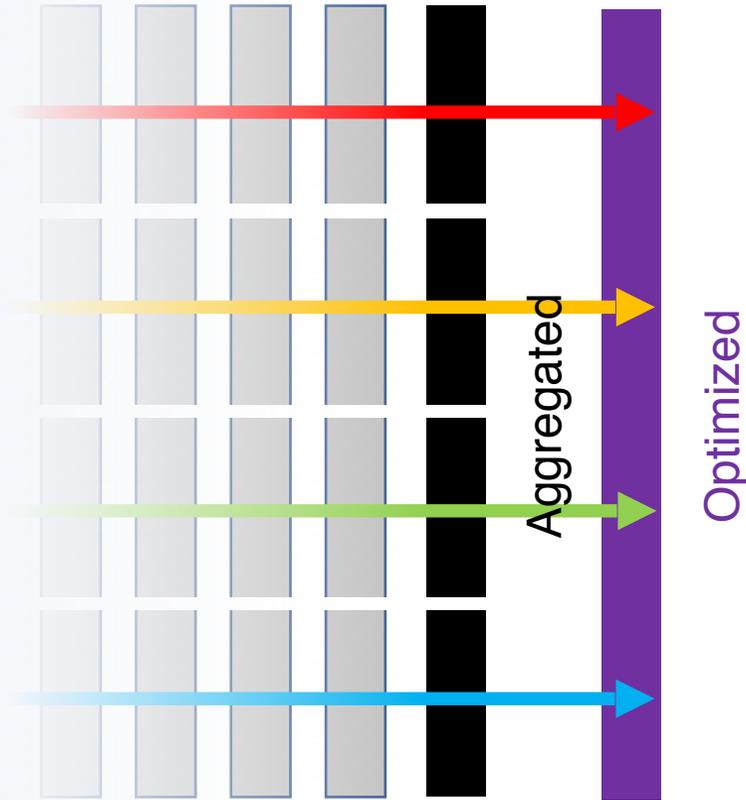
Key 0 from 8 workers

When Aggregation is done, PHub:

- PHub optimizes a chunk with the **same core** that aggregates that chunk.
- **FP32-level streaming** aggregation and optimization to hide communication latency.

# Software Optimizations

## TALL AGGREGATION AND OPTIMIZATION



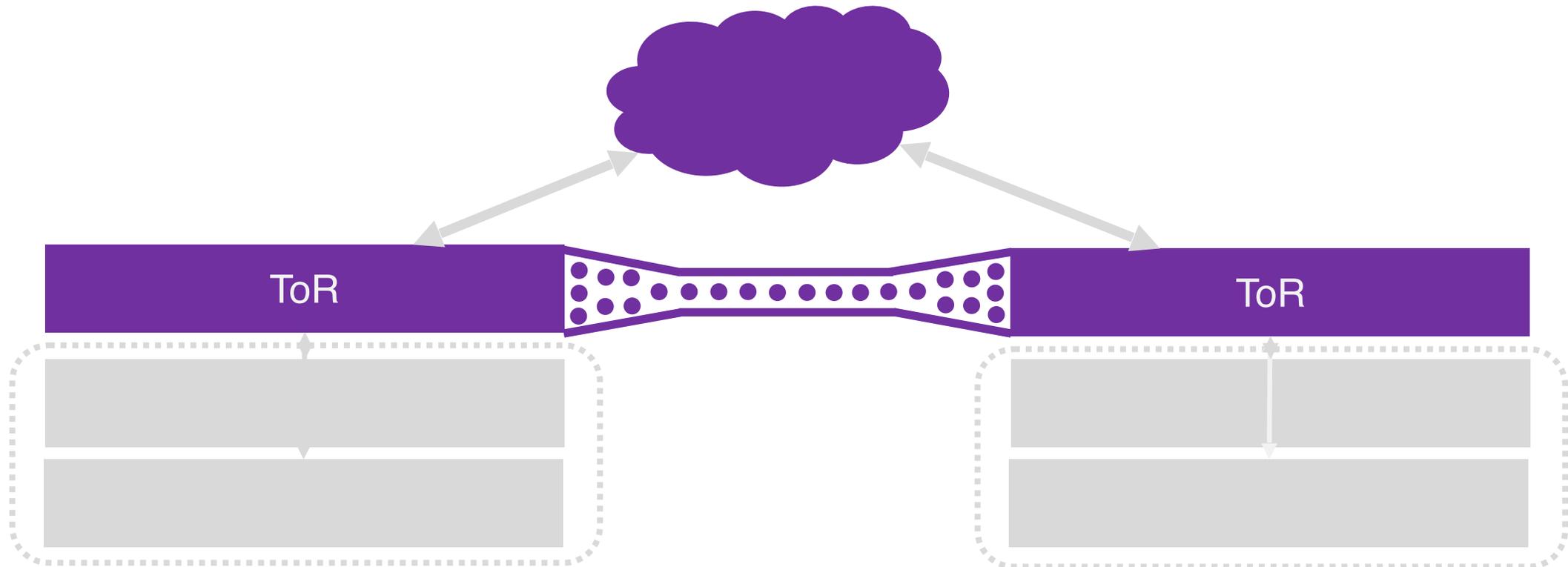
Key 0 from 8 workers

When Aggregation is done, PHub:

- PHub optimizes a chunk with the **same core** that aggregates that chunk.
- **FP32-level streaming** aggregation and optimization to hide communication latency.

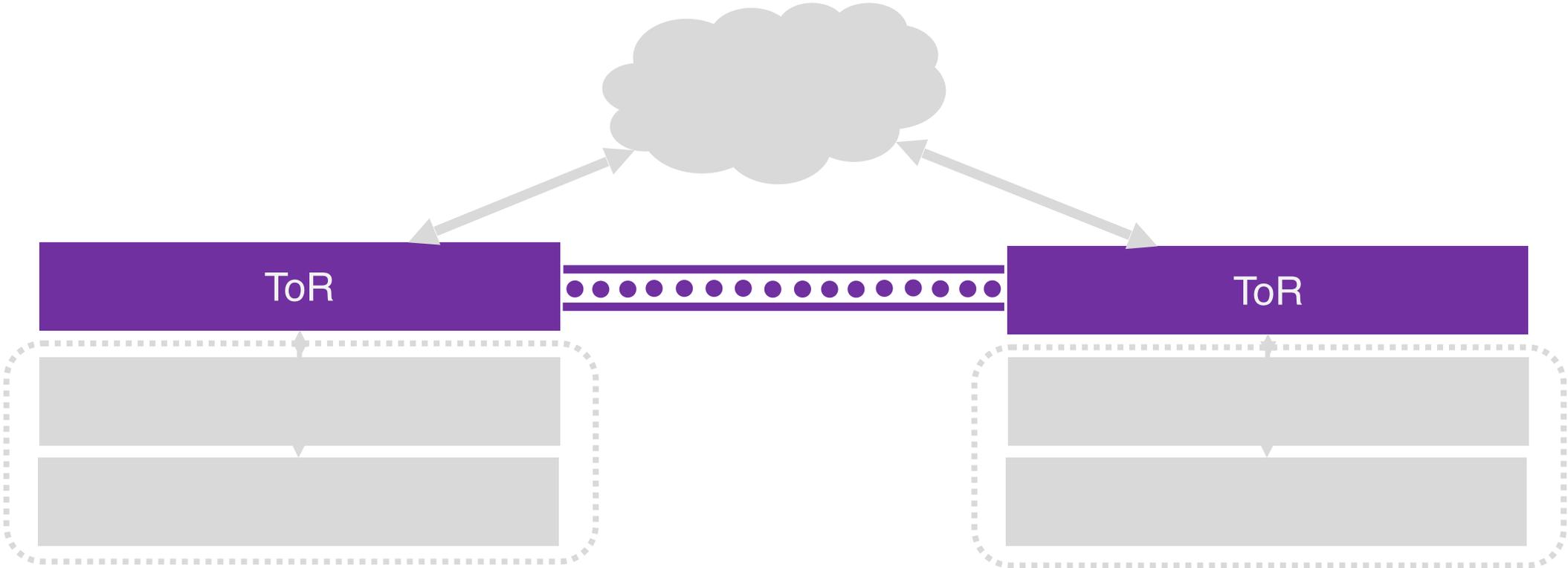
# Eliminating deployment bottlenecks:

PHub hierarchical reduction: reducing cross rack traffic



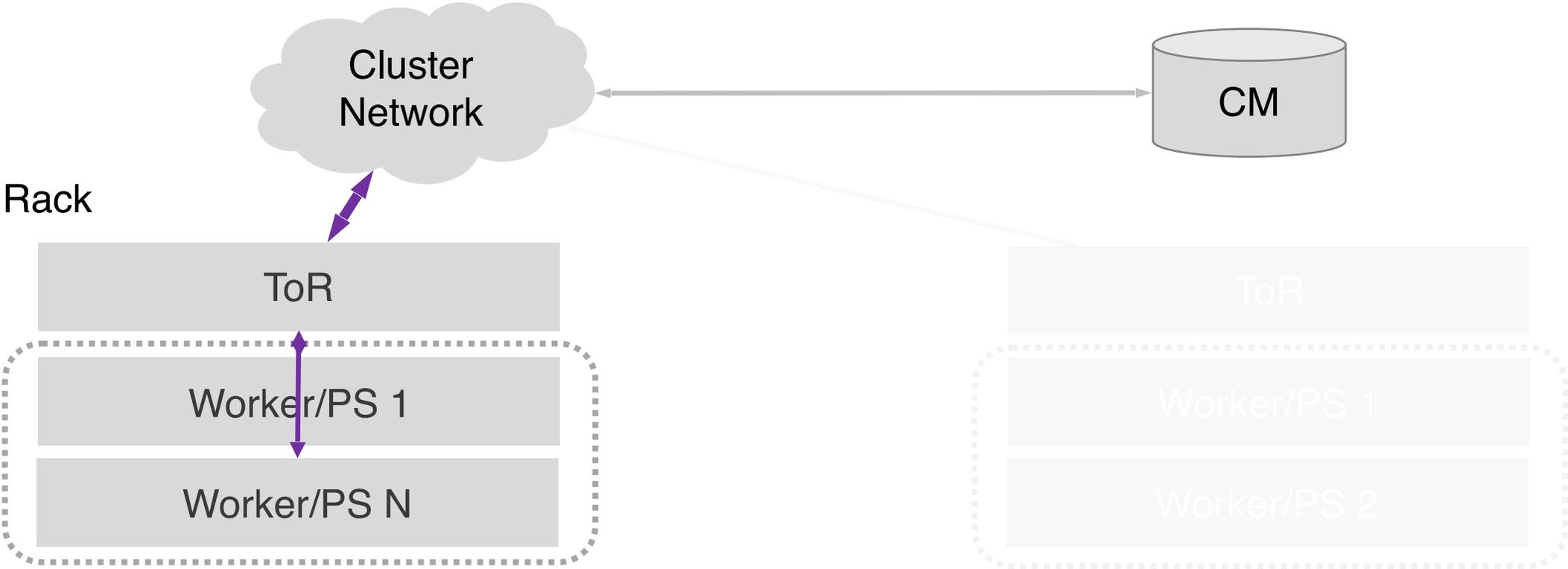
# Eliminating deployment bottlenecks:

## PHub hierarchical reduction: reducing cross rack traffic



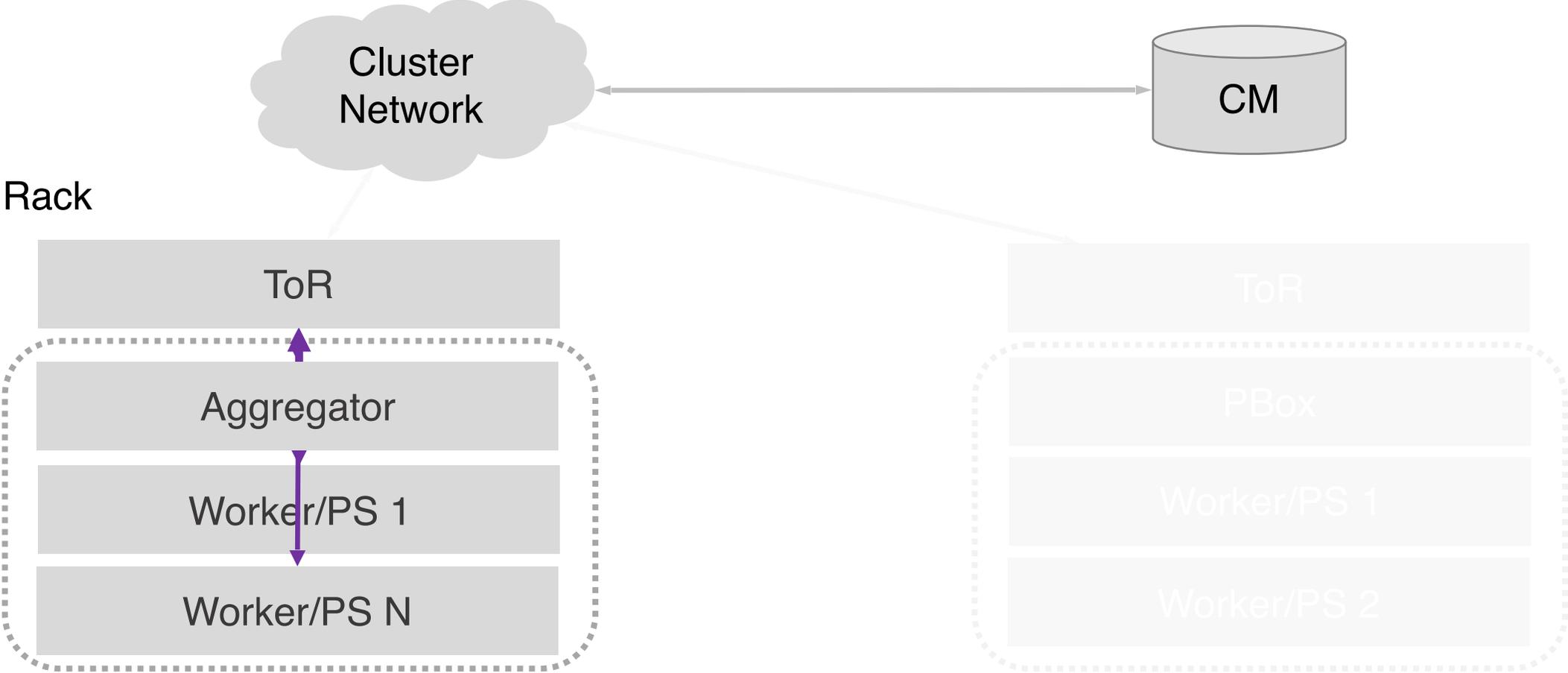
# Two-Phase Hierarchical Aggregation

## RACK SCALE PARAMETER SERVICE



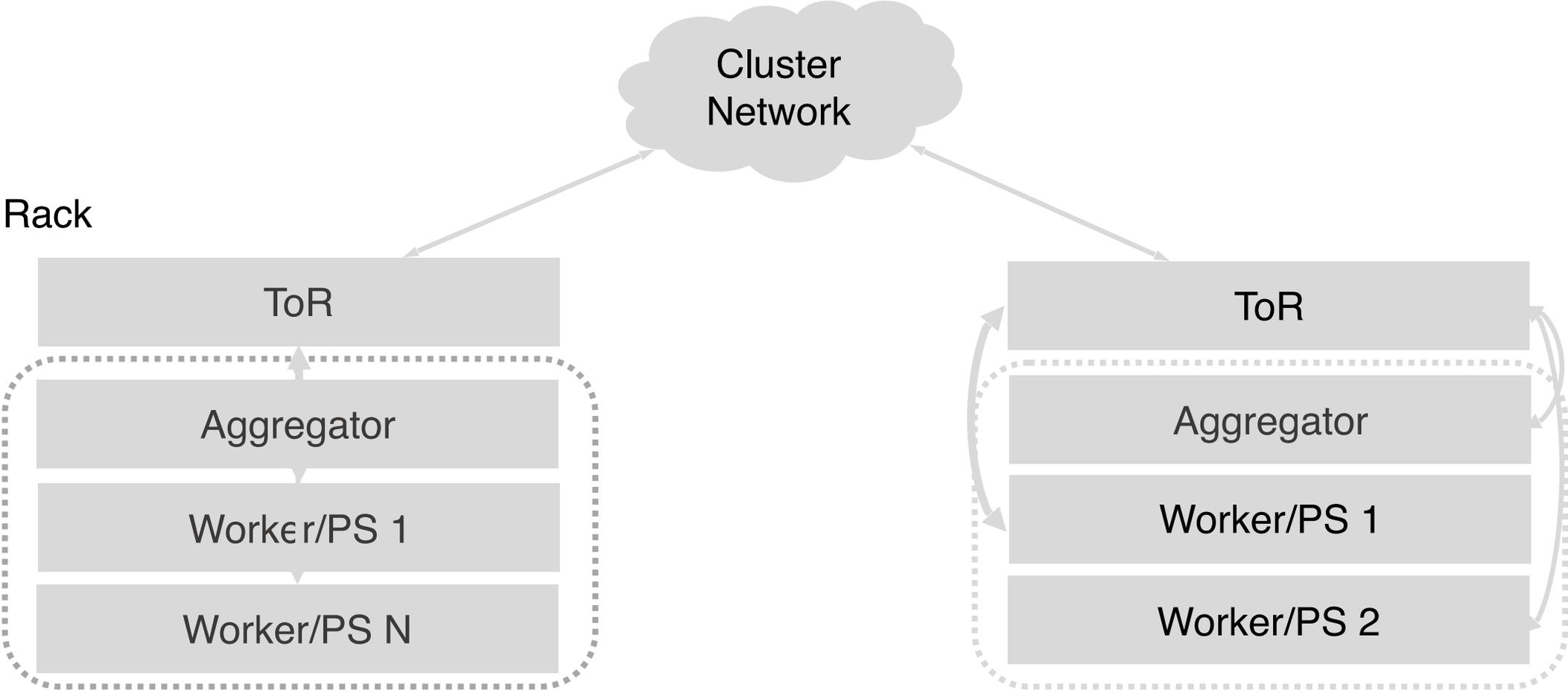
# Two-Phase Hierarchical Aggregation

## RACK SCALE PARAMETER SERVICE



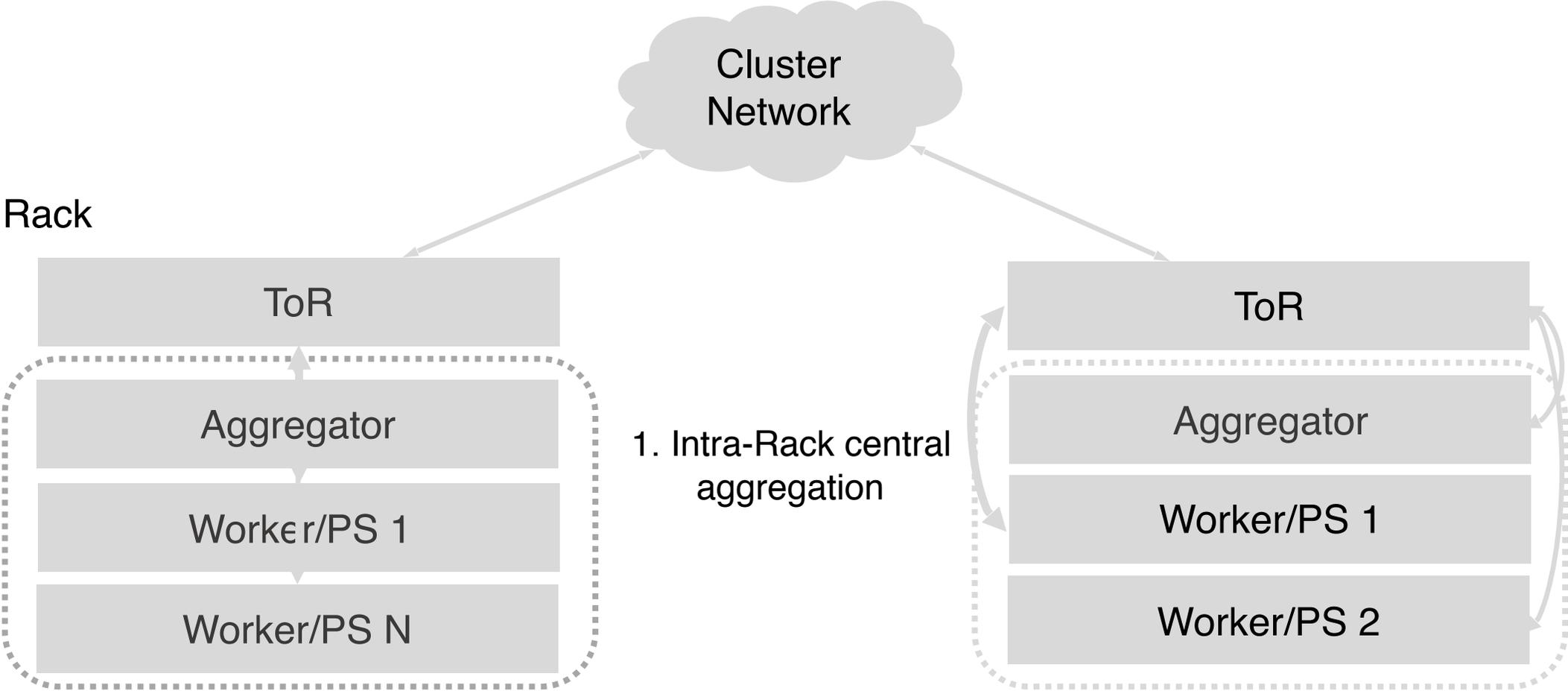
# Two-Phase Hierarchical Aggregation

## ADAPTING TO THE DATACENTER NETWORK TOPOLOGY



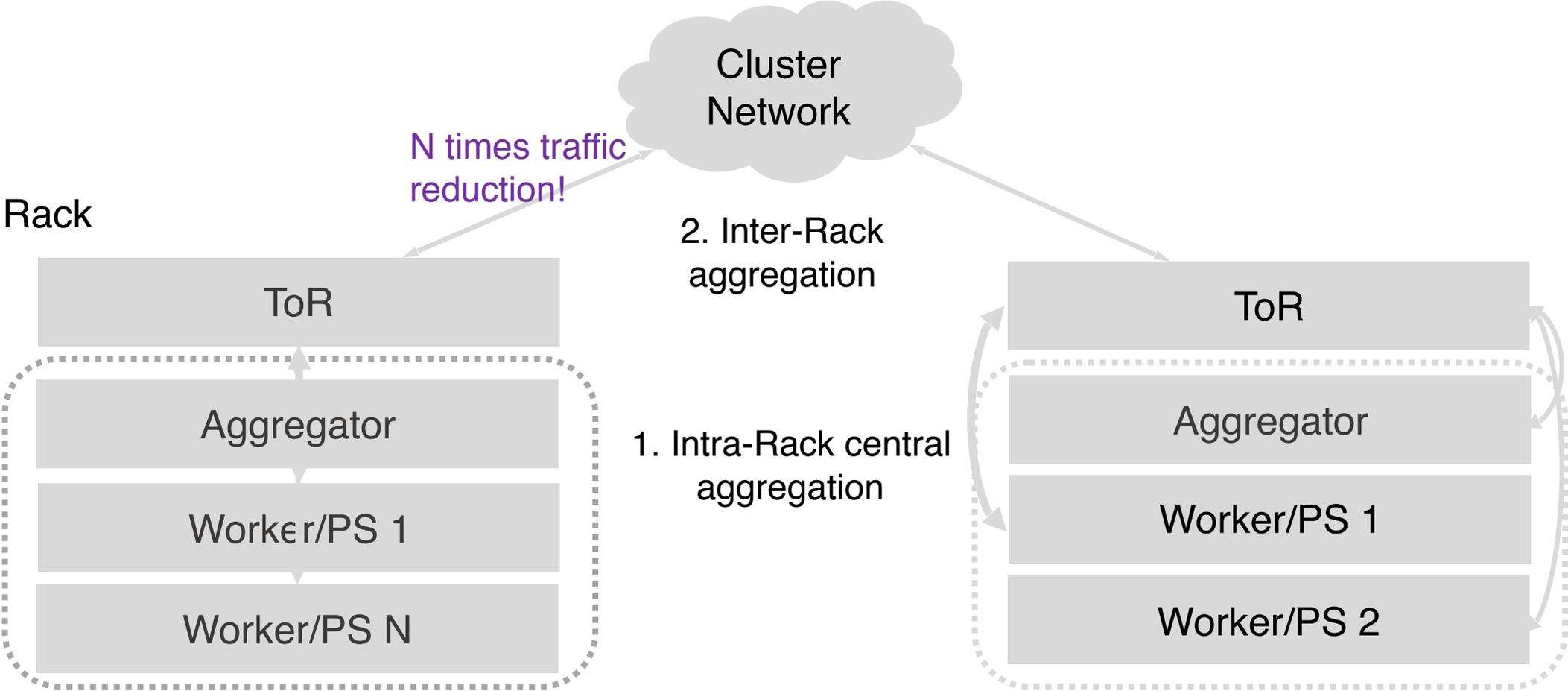
# Two-Phase Hierarchical Aggregation

## ADAPTING TO THE DATACENTER NETWORK TOPOLOGY



# Two-Phase Hierarchical Aggregation

## ADAPTING TO THE DATACENTER NETWORK TOPOLOGY



# Efficient DDNN Training in Commercial Cloud

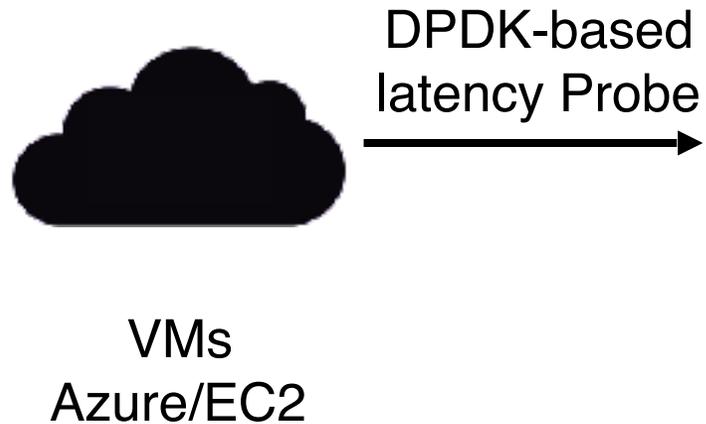
## ACTIVE TOPOLOGY PROBING



VMs  
Azure/EC2

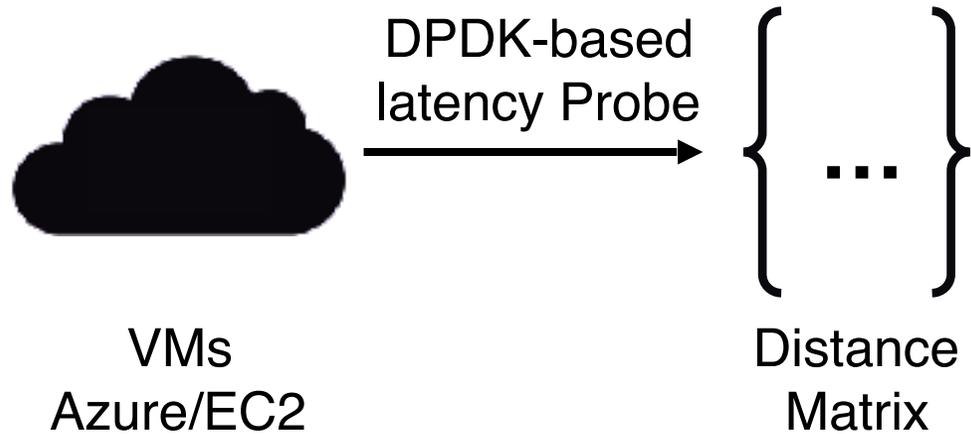
# Efficient DDNN Training in Commercial Cloud

## ACTIVE TOPOLOGY PROBING



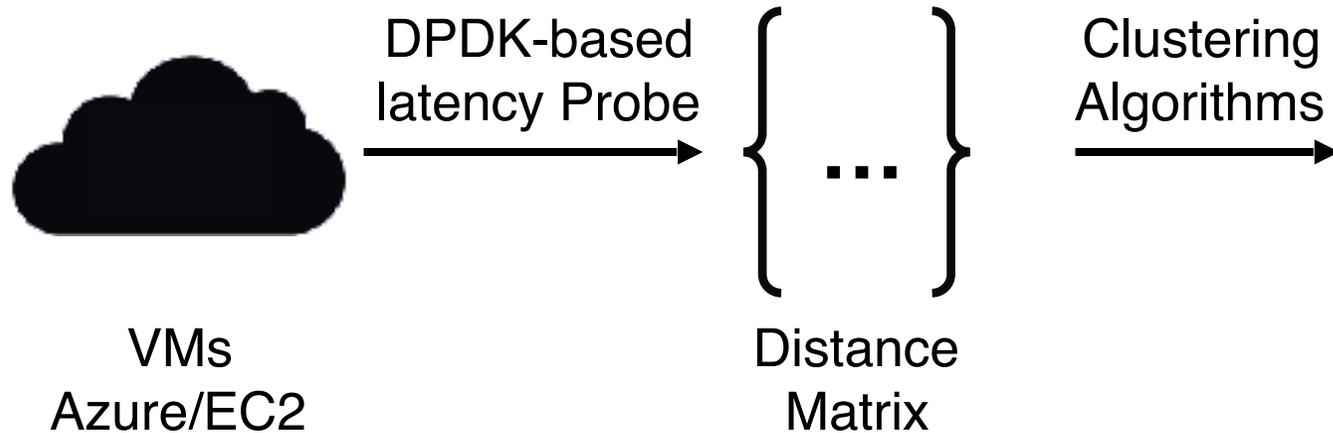
# Efficient DDNN Training in Commercial Cloud

## ACTIVE TOPOLOGY PROBING



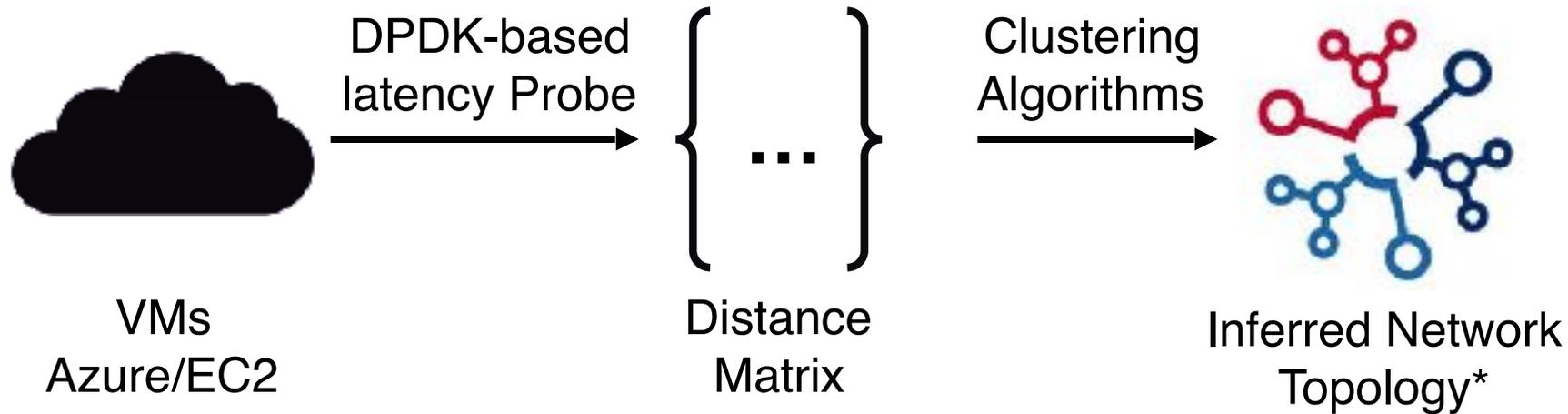
# Efficient DDNN Training in Commercial Cloud

## ACTIVE TOPOLOGY PROBING



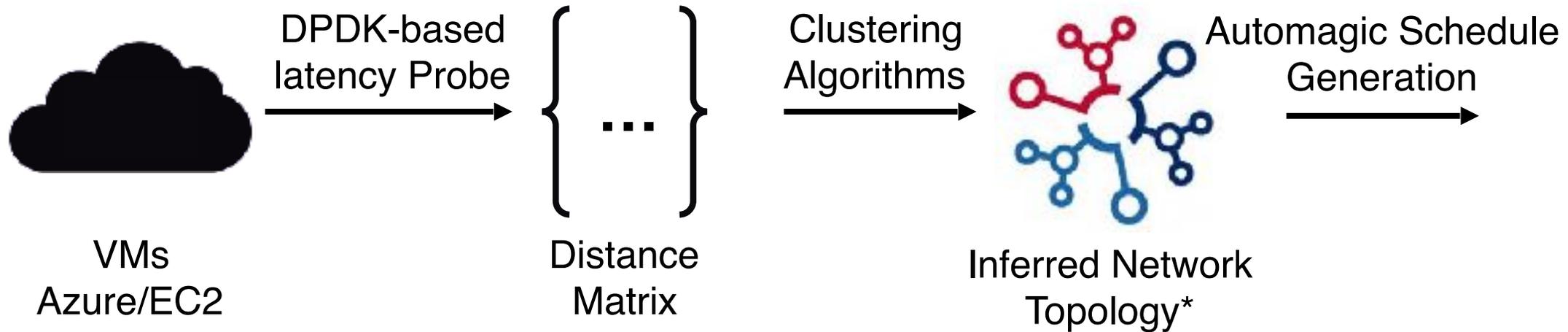
# Efficient DDNN Training in Commercial Cloud

## ACTIVE TOPOLOGY PROBING



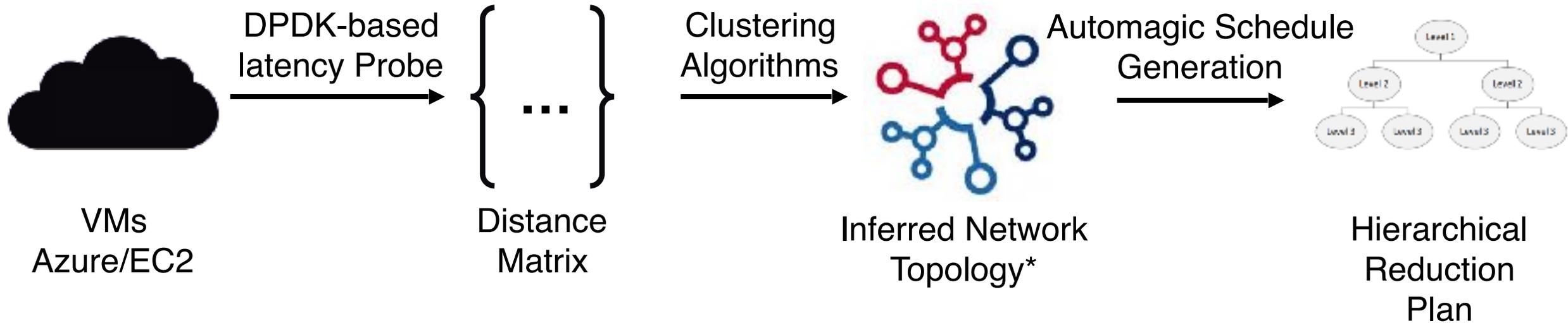
# Efficient DDNN Training in Commercial Cloud

## ACTIVE TOPOLOGY PROBING



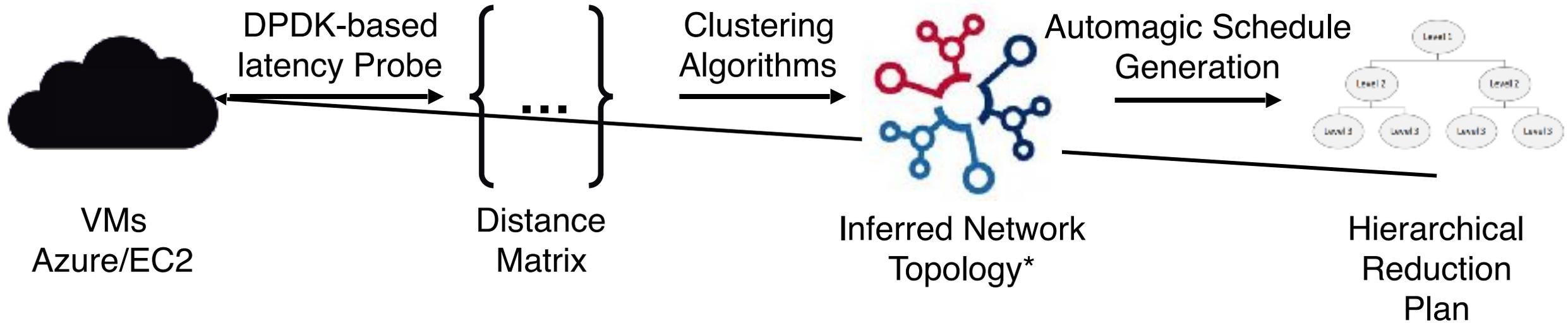
# Efficient DDNN Training in Commercial Cloud

## ACTIVE TOPOLOGY PROBING

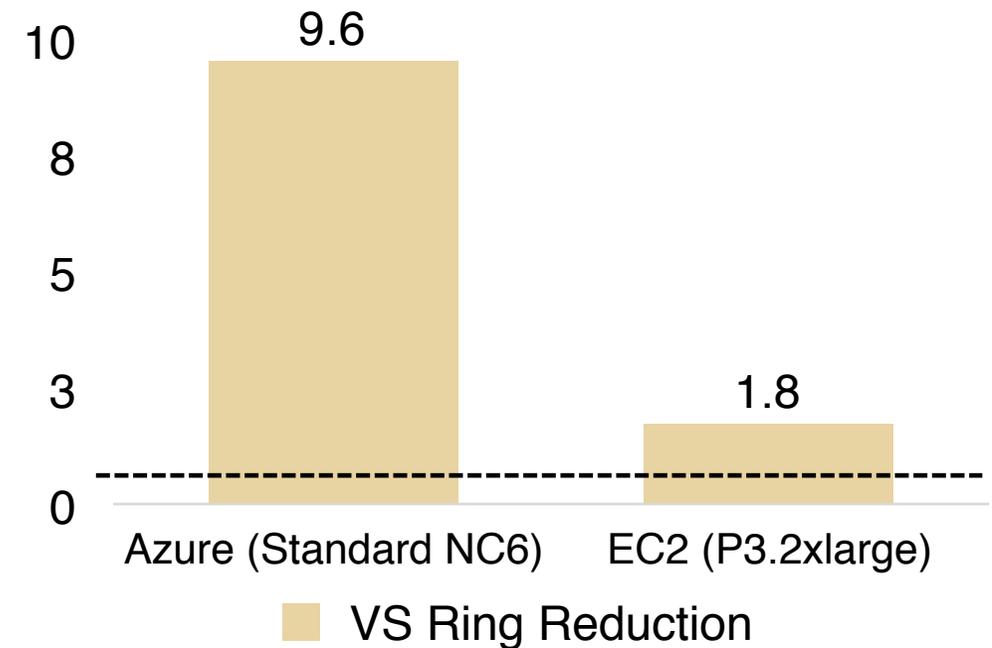
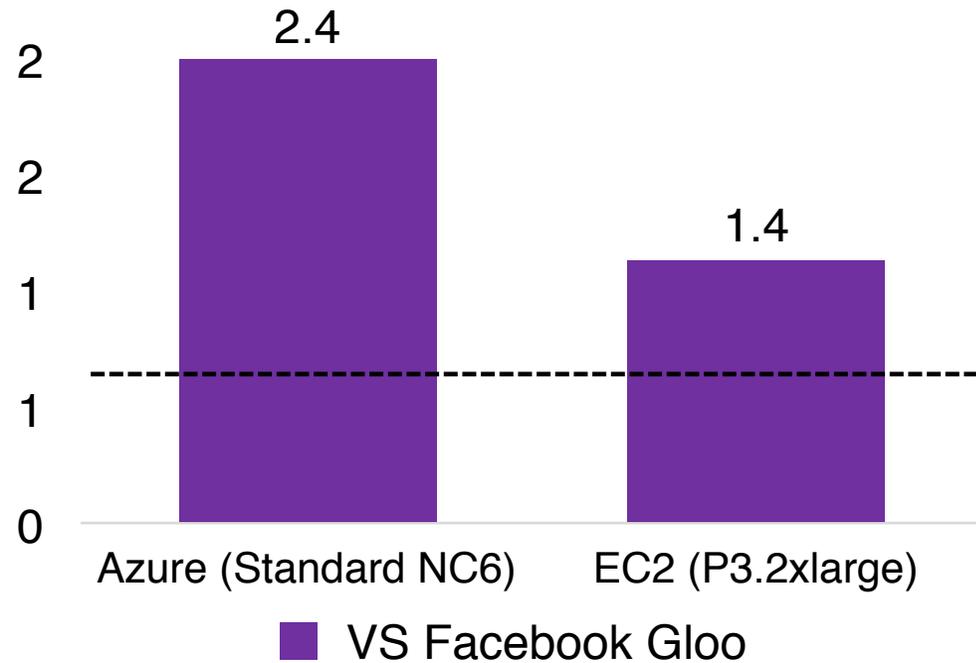


# Efficient DDNN Training in Commercial Cloud

## ACTIVE TOPOLOGY PROBING



# Performance in commercial cloud with PHub

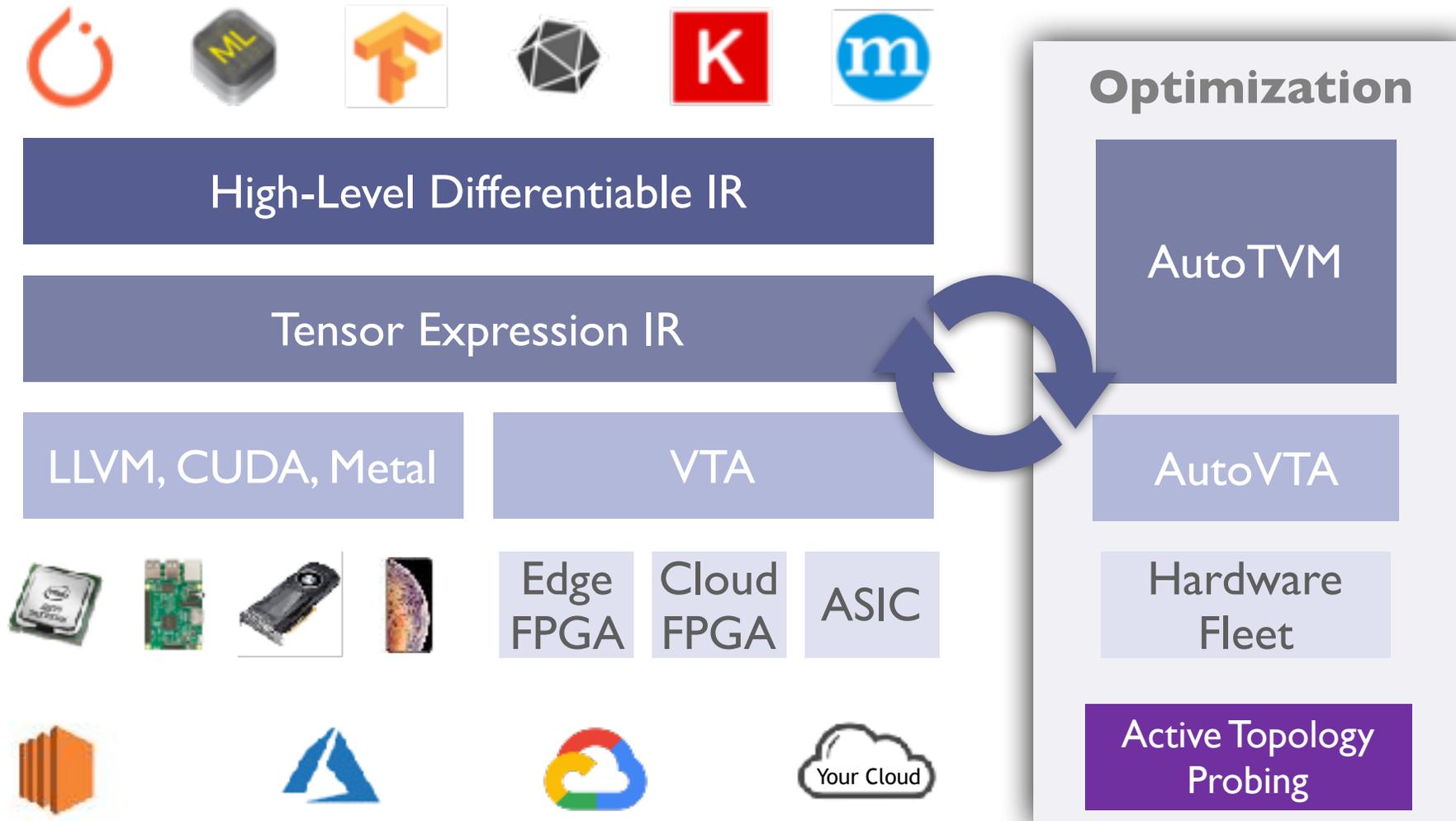


Windows Azure and Amazon EC2. 32 instances. Up to 10 Gbps. Standard\_NC6: Nvidia K80. Batch Size = 512. P3.2xLarge: Nvidia V100. Batch Size = 512. Facebook Caffe2/Pytorch. ResNet 50.

# Framework Integration

Support for Mxnet/Pytorch/Caffe2.

```
var pHub = std::make_shared<PHub>(cfg.redisIp, nMap, keySize, appAddrs, cntr,  
sizeof(float), cfg.rank, plp);  
pHub->ToggleUseSchedule(pSchedule);  
pHub->Reduce();
```



**Groundwork for bringing TVM to the distributed world for training and inference, on commercial cloud, or in your own cluster.**



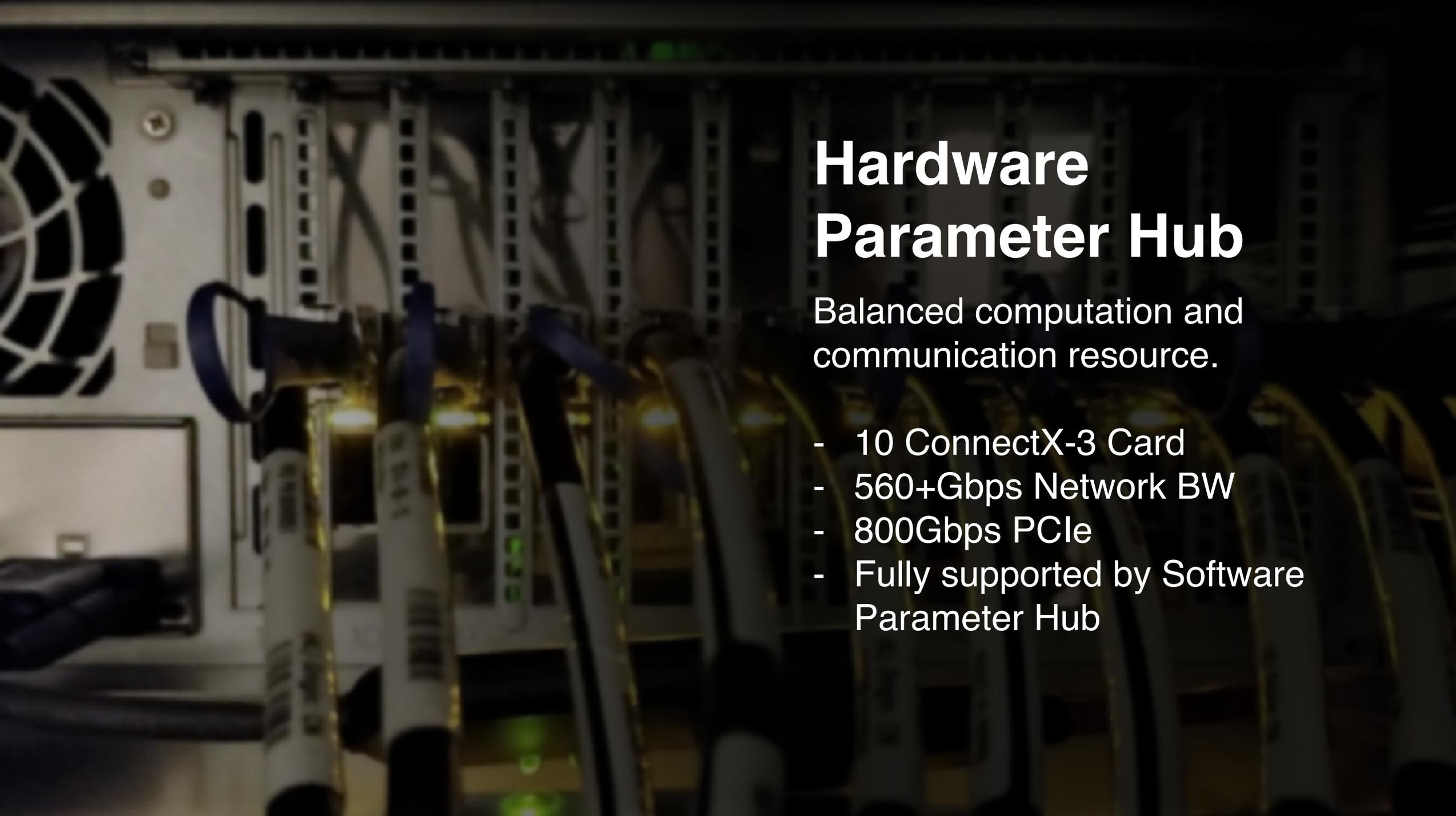




# Hardware Parameter Hub

# Hardware Parameter Hub



A server rack with network cables and a fan. The background is dark, and the text is overlaid on the right side. The server rack is filled with various components, including network cables and a fan. The text is overlaid on the right side of the image.

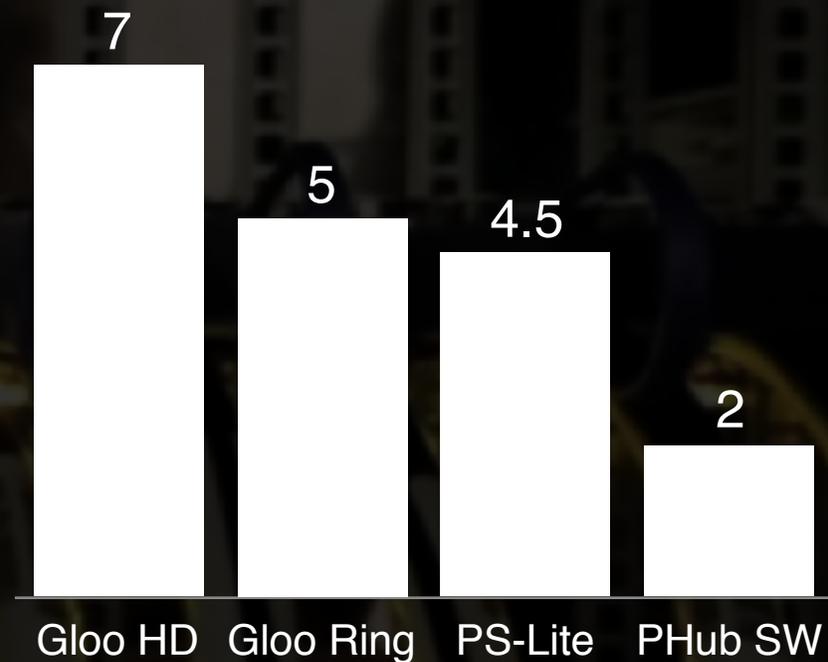
# Hardware Parameter Hub

Balanced computation and communication resource.

- 10 ConnectX-3 Card
- 560+Gbps Network BW
- 800Gbps PCIe
- Fully supported by Software Parameter Hub

# Hardware Parameter Hub

35GB/s aggregation throughput.  
Supports 100+ ResNet-50  
training nodes with a single  
machine.



# Hardware Parameter Hub

ResNet-50.

See paper for detailed estimates.

Better training throughput/\$.

# Hardware Parameter Hub

ResNet-50.

See paper for detailed estimates.

25%

Better training throughput/\$.