# HetPipe: Enabling Large DNN Training on (Whimpy) Heterogeneous GPU Clusters through Integration of Pipelined Model Parallelism and Data Parallelism

USENIX ATC 2020

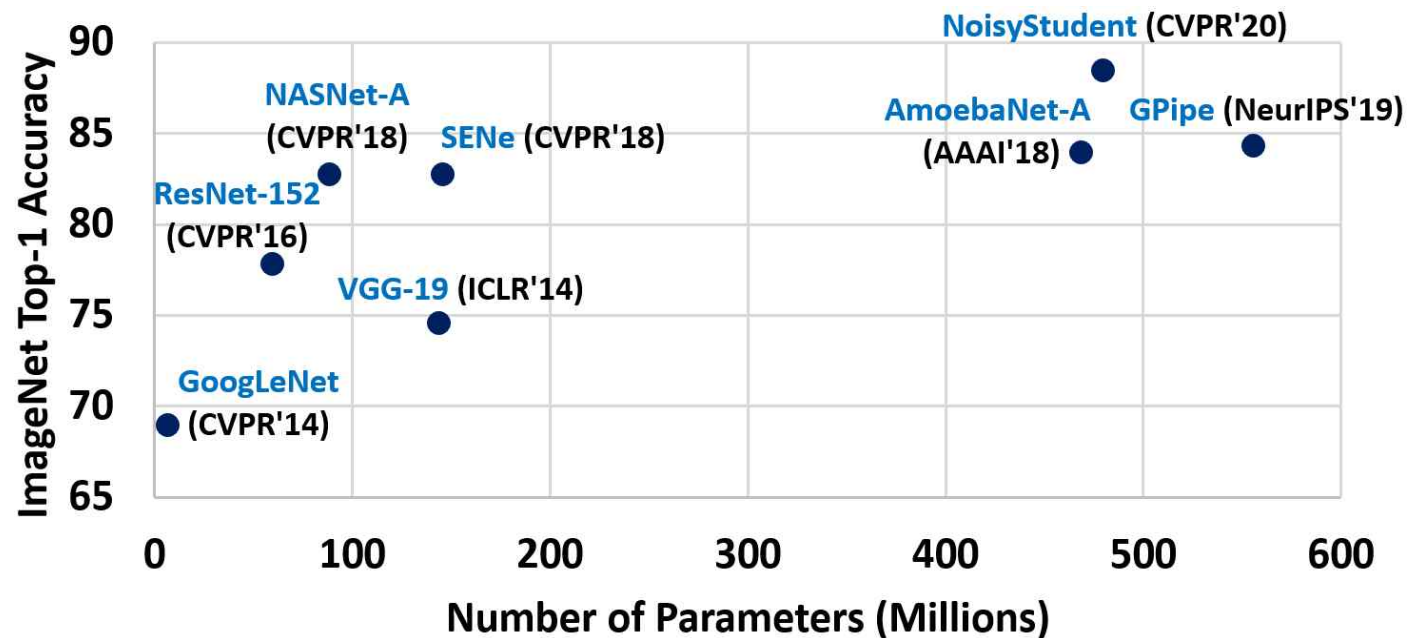Jay H. Park, Gyeongchan Yun, Chang M. Yi, Nguyen T. Nguyen, Seungmin Lee, Jaesik Choi[†], Sam H. Noh, and Young-ri Choi

# Contents

- **Motivation & Background**

- **HetPipe in a Nutshell**

- **Our System: HetPipe**

- **Evaluation**

- **Conclusion**

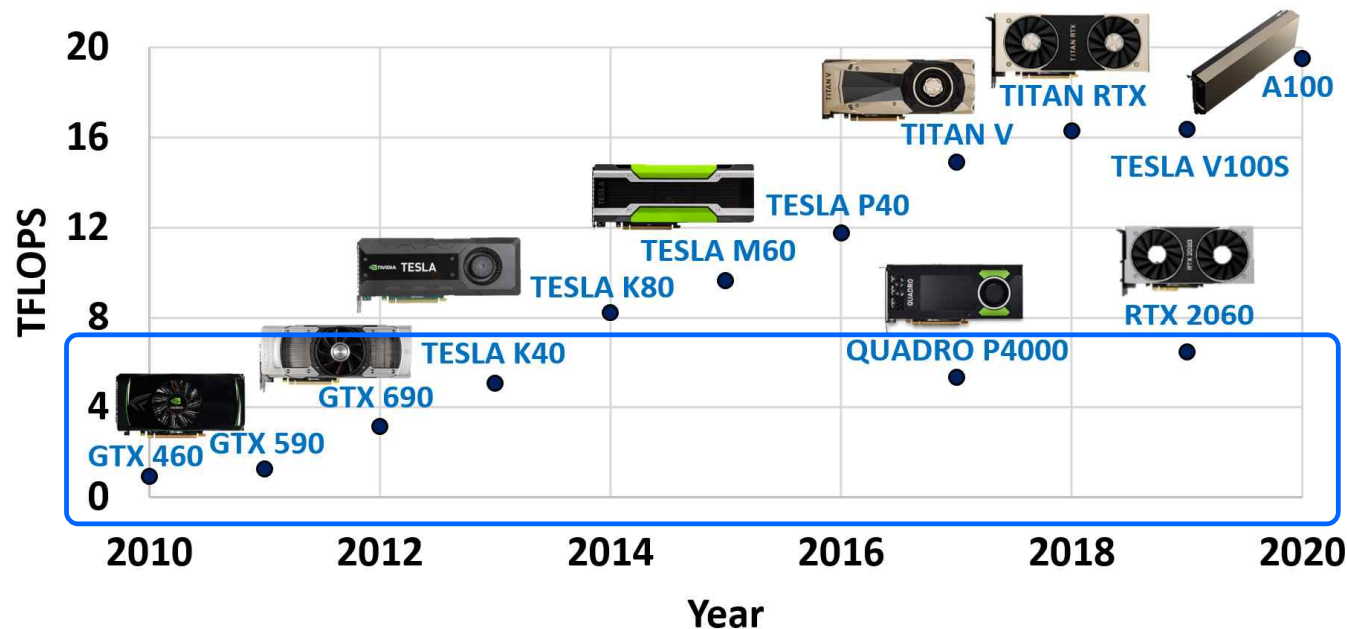- **DNN (Deep Neural Network) models continue to grow**



- Need more powerful GPUs for training!

3
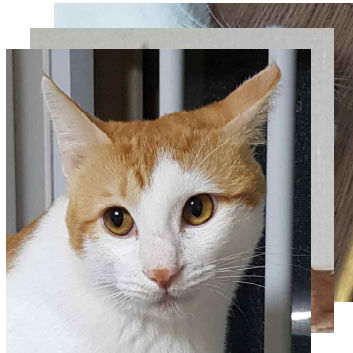
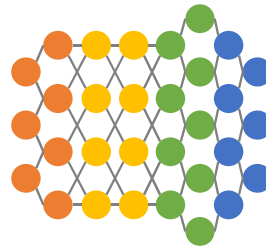- **Short release cycle of new GPU architectures**



- • Use of heterogeneous GPUs is inevitable!
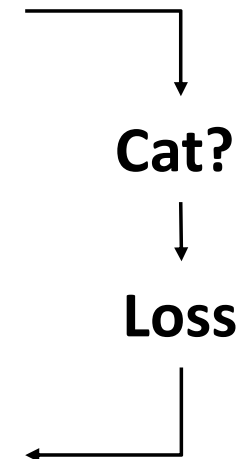- • What to do with *whimpy* GPUs?

# DNN Training



**Minibatch $i$**
**(Training Data)**

Forward Pass $i$

Weight Parameter $w$

Backward Pass $i$

Cat?

Loss

$$w_{i+1} = w_i - \eta \cdot u_i$$

# Parallelizing DNN Training

- **Data parallelism (DP)**



Parameter Server (PS)

Worker 1 ... Worker $n$

Weights synchronized through **PS** or **AllReduce**

- GPU memory limitation

- **Model parallelism (MP)**



P1    P2    P3    P4

Forward pass
Backward pass

- Low GPU utilization

- **Attempts to improve MP utilization**
  - Pipelined model parallelism (PMP)

- **PipeDream** [SOSP'19]
- **GPipe** [NIPS'19]



- Designed for homogeneous GPUs
- Designed for a single PMP worker

# Challenges in integration PMP+DP in Heterogeneous GPUs

- What weight version should be used by each VW to synchronize with other VWs?

- How do we reduce virtual worker stragglers when we consider DP?

⋮

**Many more in the paper**

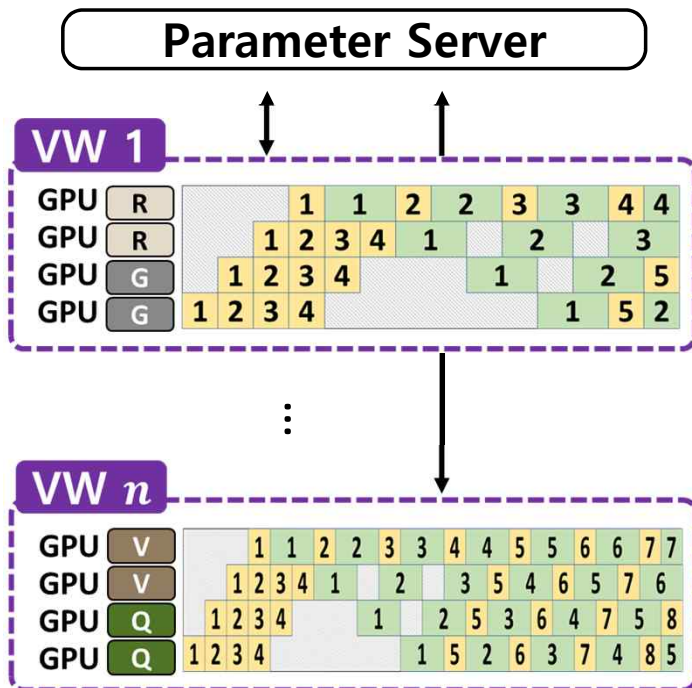**Enable Large DNN Training on Heterogeneous GPUs**

Aggregate heterogeneous resources

Reduce the straggler problem

**Integrates PMP + DP**

Novel parameter synchronization model

WSP (Wave Synchronous Parallel)

**Proof of WSP Convergence**

**Cluster Configuration**

**Resource Allocator**

Assign $k$ GPUs to each virtual worker

**DNN Model**

**Model Partitioner**

Divide model into $k$ partitions

VW 1
P1   P2   P3   P4   ...   VW $n$
P1'  P2'  P3'  P4'

VW 1
P4 → R
P3 → R
P2 → G
P1 → G
Time   ...

PS

VW $n$
P4' → V
P3' → V
P2' → Q
P1' → Q
...

11

# Outline

- Motivation & Background

- HetPipe in a Nutshell

- **Our System: HetPipe**
  - **Pipelined Model Parallelism Within a VW**
  - Data Parallelism with Multiple VWs

- Evaluation

- Conclusion

- **Execution of a virtual worker**

Time



$N_m$ minibatches processed concurrently in pipeline manner

$$W_{bcal}$$

$W_{bcal}$ is a consistent version of weights within a VW

- **Weight management procedure**

Time

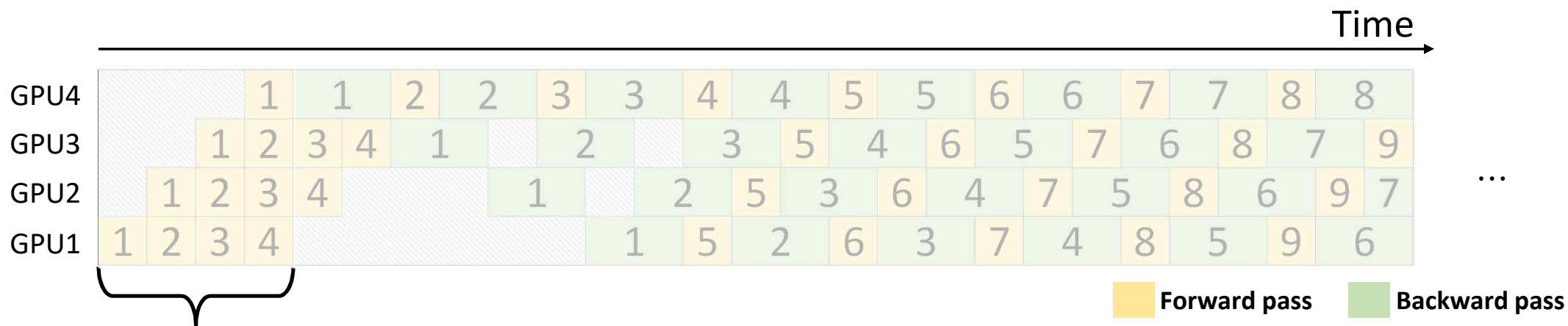| | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GPU4 | | | 1 | 1 | 2 | 2 | 3 | 3 | 4 | 4 | 5 | 5 | 6 | 6 | 7 | 7 | 8 | 8 |
| GPU3 | | 1 | 2 | 3 | 4 | 1 | | 2 | | 3 | 5 | 4 | 6 | 5 | 7 | 6 | 8 | 7 | 9 |
| GPU2 | 1 | 2 | 3 | 4 | | 1 | | 2 | 5 | 3 | 6 | 4 | 7 | 5 | 8 | 6 | 9 | 7 |
| GPU1 | 1 | 2 | 3 | 4 | | | | 1 | 5 | 2 | 6 | 3 | 7 | 4 | 8 | 5 | 9 | 6 |

$\uparrow$ $\uparrow$ $\uparrow$ $\uparrow$
$w_1$ $w_2$ $w_3$ $w_4$

**Update $u_1$**

$w_5 \leftarrow w_{bcal}$

Forward pass     Backward pass

**Initial weight version ($w_0$)**

$w_{bcal} = w_0 = w_1 = w_2 = w_3 = w_4$

$$w_{bcal}$$

$$w_{bcal} \leftarrow w_{bcal} + u_1$$

- **Local staleness ($S_{bcal}$)**: maximum missing updates

Time

| | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GPU4 | | | | 1 | 1 | 2 | 2 | 3 | 3 | 4 | 4 | 5 | 5 | 6 | 6 | 7 | 7 | 8 | 8 |
| GPU3 | | | 1 | 2 | 3 | 4 | 1 | | 2 | | 3 | 5 | 4 | 6 | 5 | 7 | 6 | 8 | 7 | 9 |
| GPU2 | | 1 | 2 | 3 | 4 | | | 1 | | 2 | 5 | 3 | 6 | 4 | 7 | 5 | 8 | 6 | 9 | 7 |
| GPU1 | 1 | 2 | 3 | 4 | | | | | 1 | 5 | 2 | 6 | 3 | 7 | 4 | 8 | 5 | 9 | 6 |

Forward pass   Backward pass

$w_5 \leftarrow w_{bcal}$

$S_{bcal} = 3 \leftarrow$ $w_5$ missing updates of minibatches 2 to 4

$w_{bcal}$

$w_{bcal} \leftarrow w_{bcal} + u_1$

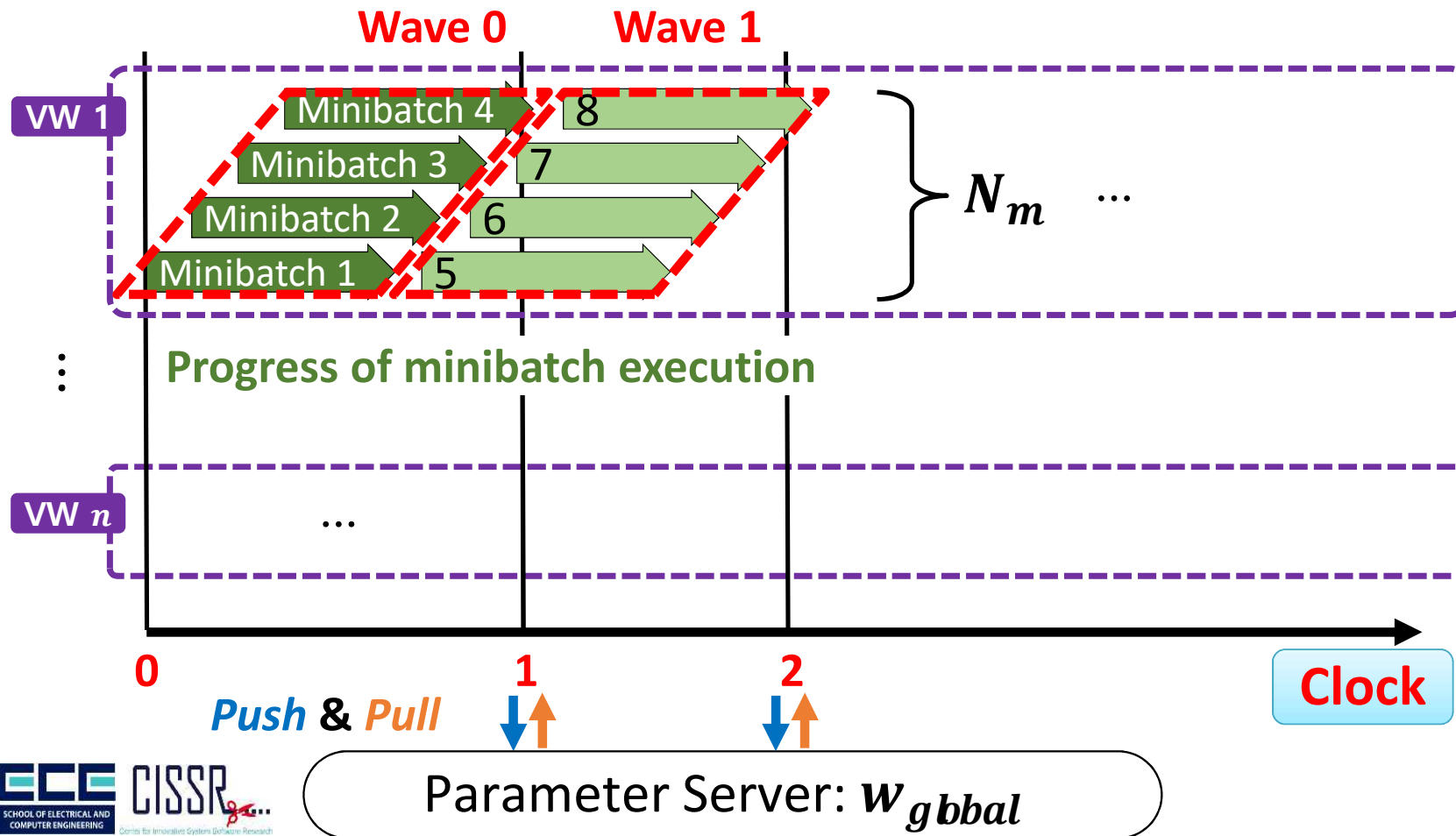# Outline

- <span style="color:gray">**Motivation & Background**</span>
- <span style="color:gray">**HetPipe in a Nutshell**</span>

## Our System: HetPipe

- Pipelined Model Parallelism Within a VW
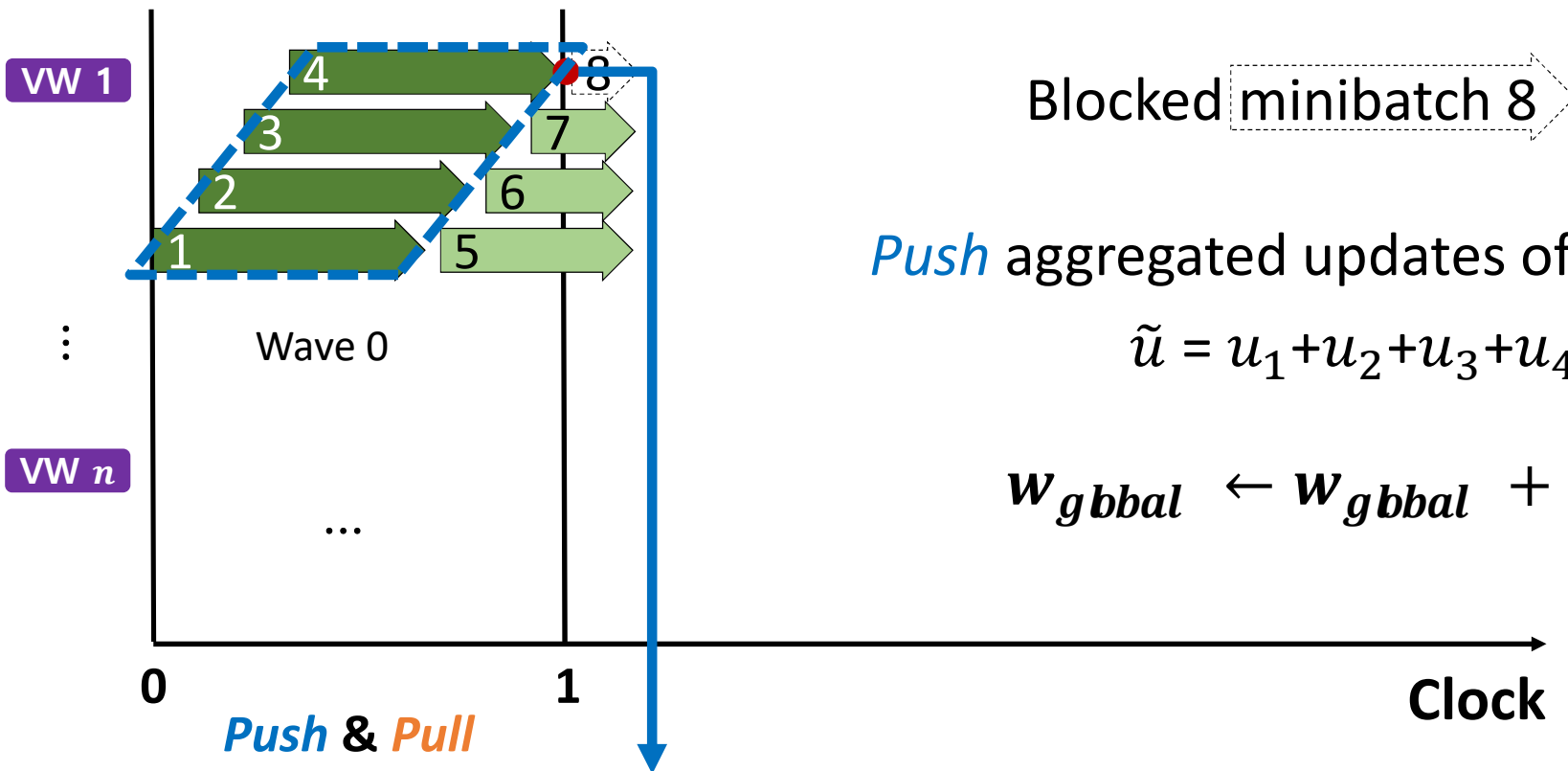- **Data Parallelism with Multiple VWs**

- Evaluation
- Conclusion

**Wave**: Sequence of concurrently executing $N_m$ minibatches



Wave 0   Wave 1

VW 1

Minibatch 4   8
Minibatch 3   7
Minibatch 2   6
Minibatch 1   5

$N_m$  ...

**Progress of minibatch execution**

VW $n$   ...

0   1   2

*Push* & *Pull*

Parameter Server: $w_{global}$

Clock

- ***Push* occurs every clock**

VW 1

4    8

3    7

2    6

1    5

⋮

Wave 0

VW $n$

...

0                    1

Push & Pull

Parameter Server: $w_{gbbal}$

Blocked minibatch 8

*Push* aggregated updates of wave0 ($\tilde{u}$)

$$\tilde{u} = u_1 + u_2 + u_3 + u_4$$

$$w_{gbbal} \leftarrow w_{gbbal} + \tilde{u}$$

Clock

20

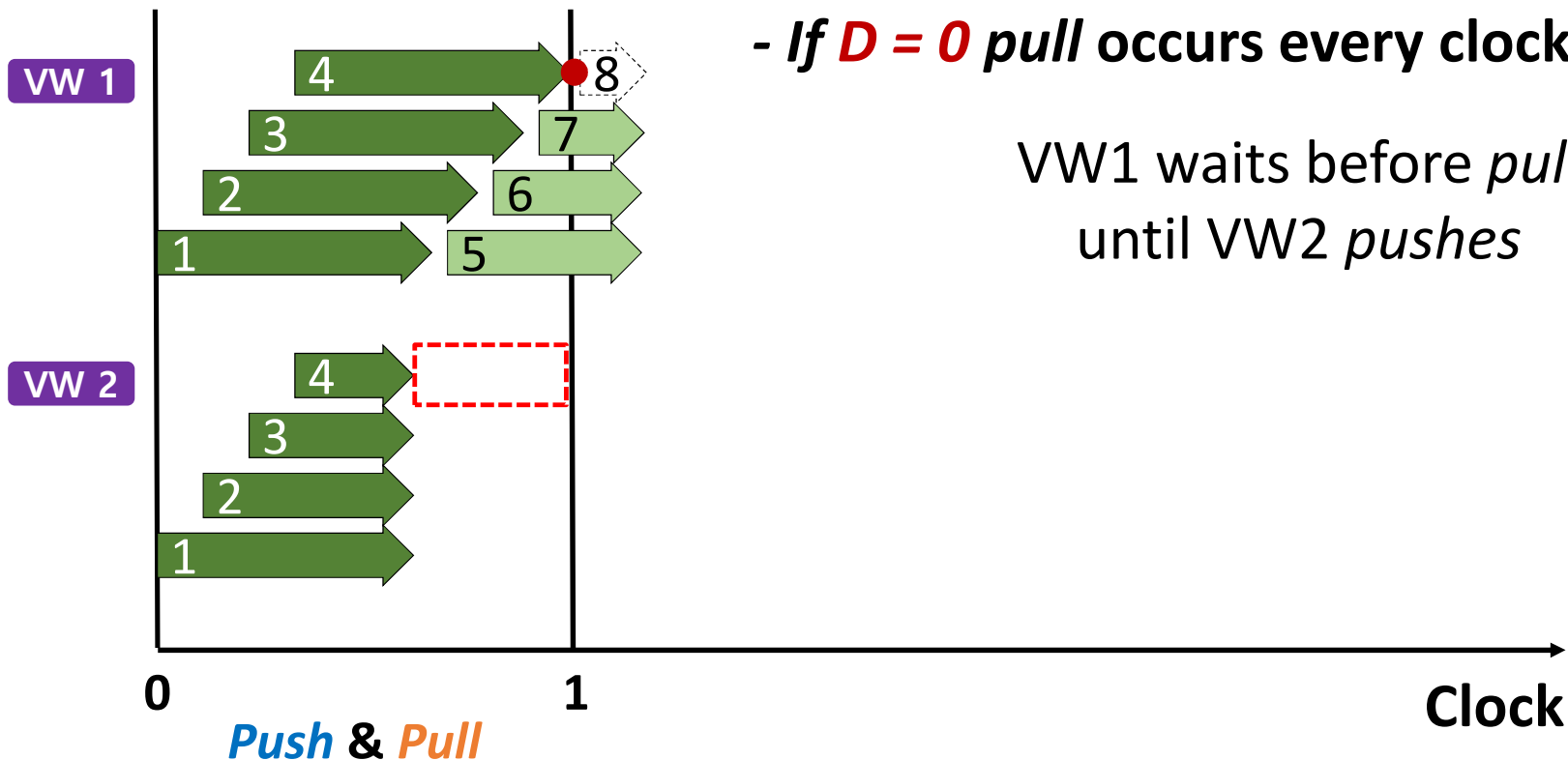- *Pull* occurs intermittently - Depending on user defined *clock distance D*

  - If *D = 0* *pull* occurs every clock

VW1 waits before *pull*
until VW2 *pushes*

VW 1

| 4 | 8 |
| 3 | 7 |
| 2 | 6 |
| 1 | 5 |

VW 2

| 4 |
| 3 |
| 2 |
| 1 |

0      1      **Clock**

*Push* & *Pull*

Parameter Server: $w_{gbbal}$

- *Pull* occurs intermittently - Depending on user defined *clock distance D*



**If $D = 0$**
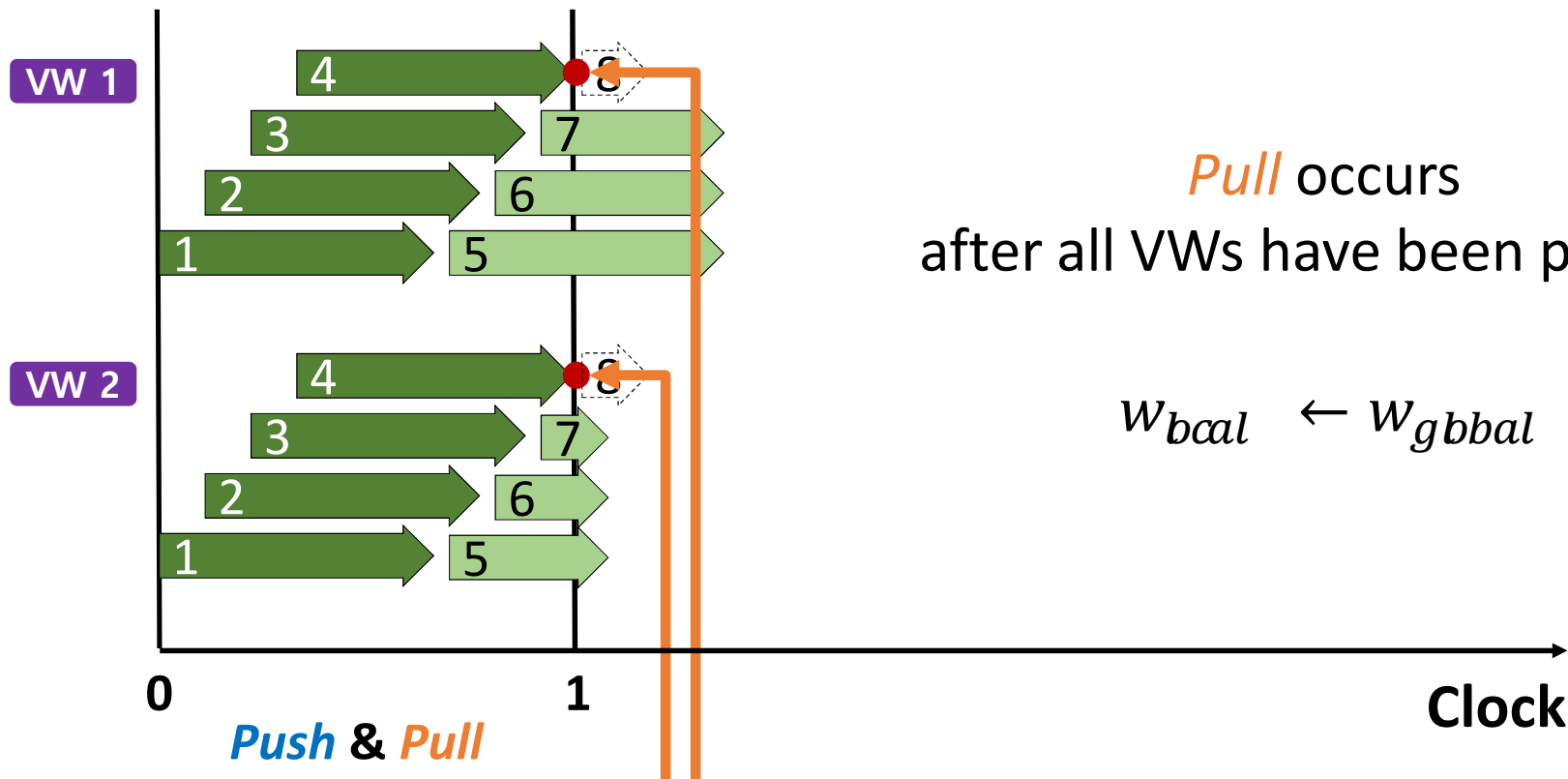
VW1 waits before *pull*
until VW2 *pushes*

VW2 *Push* aggregated updates ($\tilde{u}$)

$$w_{gbbal} \leftarrow w_{gbbal} + \tilde{u}$$

VW 1

VW 2

Push & Pull

0          1          Clock

Parameter Server: $w_{gbbal}$

- *Pull* occurs intermittently - Depending on user defined *clock distance D*



*If D = 0*

*Pull* occurs
after all VWs have been pushed

$$w_{bcal} \leftarrow w_{gbbal}$$

**VW 1**

**VW 2**

0

1

**Clock**

*Push* & *Pull*

Parameter Server: $w_{gbbal}$

- *Pull* **occurs intermittently - Depending on user defined *clock distance D***

**VW 1**

| 4 | 8 |
| 3 | 7 |
| 2 | 6 |
| 1 | 5 |

**If $D = 0$**

Minibatch 8 starts with $w_8$

$$w_8 = w_0 + (u_1 + u_2 + u_3 + u_4)_{vw1,vw2}$$

**VW 2**

| 4 | 8 |
| 3 | 7 |
| 2 | 6 |
| 1 | 5 |

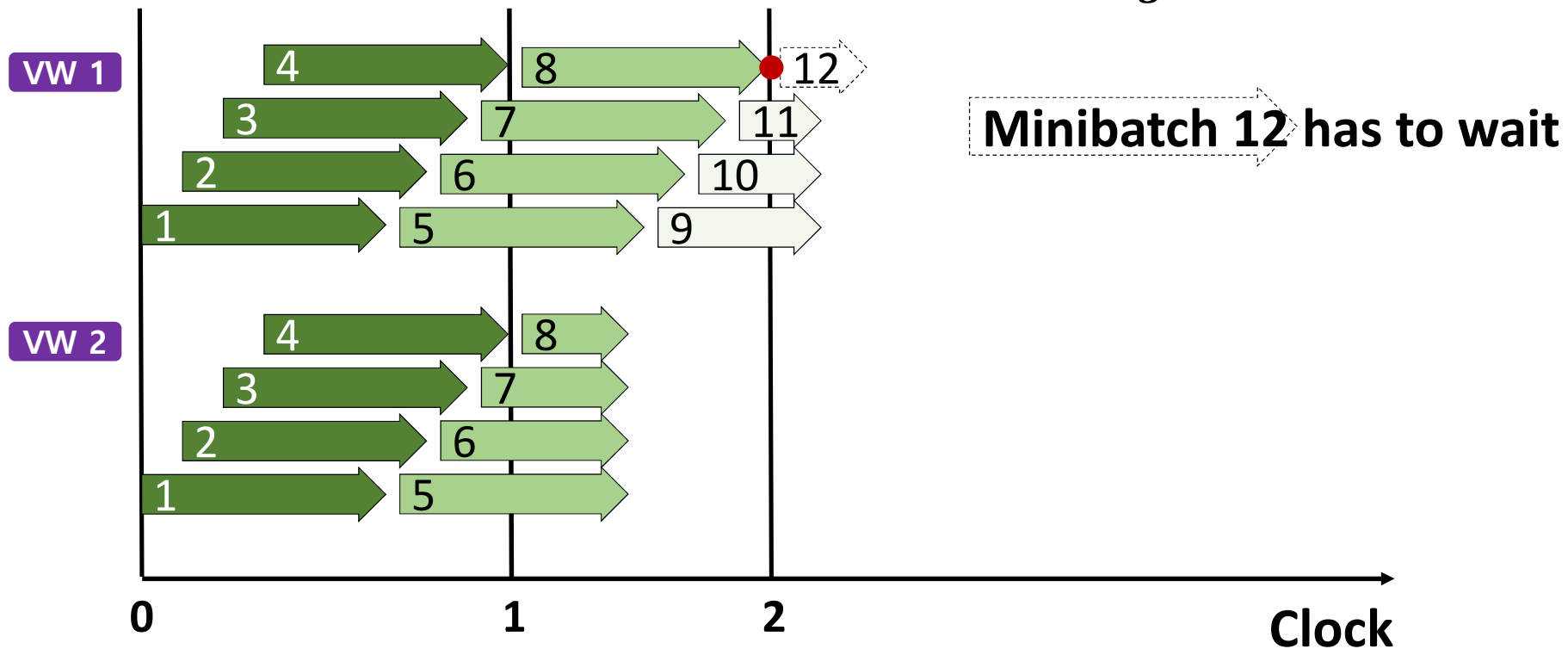0                    1

*Push* **&** *Pull*

**Clock**

Parameter Server: $w_{gbbal}$

- Local staleness ($S_{bcal}$) and global staleness ($S_{gbbal}$) with WSP



$$w_{11} = w_0 + (u_1 + u_2 + u_3 + u_4)_{vw1, vw2}$$
$$+ (u_5 + u_6 + u_7)_{vw1}$$

$$S_{gbbal} \begin{cases} (u_8 + u_9 + u_{10})_{vw1} \rightarrow S_{bcal} \\ (u_5 + u_6 + u_7)_{vw2} \end{cases}$$

$$w_{gbbal} = w_0 + (u_1 + u_2 + u_3 + u_4)_{vw1, vw2}$$

- **Local staleness ($S_{local}$) and global staleness ($S_{global}$) with WSP**



**Minibatch 12 has to wait**

Parameter Server: $w_{global}$

26

## Outline

- Motivation & Background

- HetPipe in a Nutshell

- Our System: HetPipe

- **Evaluation**
  - **Setup**
  - **Resource Allocation for Virtual Workers**
  - **Results**

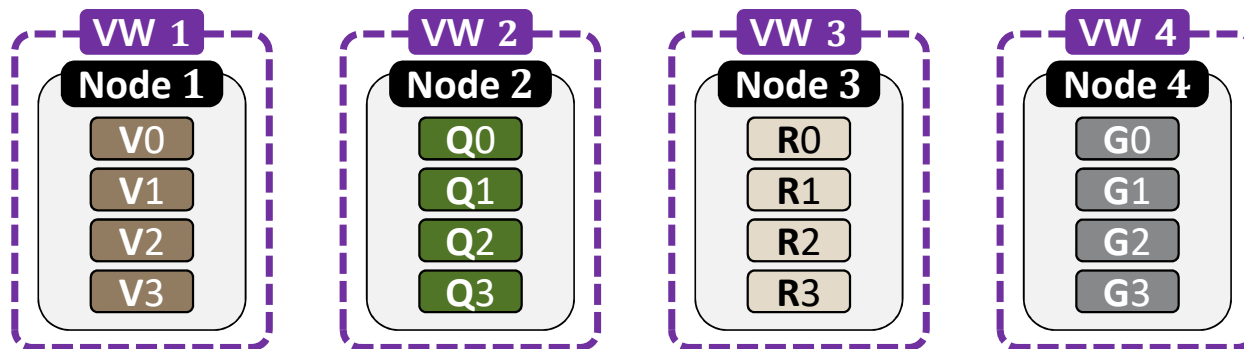- Conclusion

27

- **Cluster setup** - 4 heterogeneous GPU nodes



InfiniBand (56 Gbps)

| Node 1 | Node 2 | Node 3 | Node 4 |
| --- | --- | --- | --- |
| V0 | Q0 | R0 | G0 |
| V1 | Q1 | R1 | G1 |
| V2 | Q2 | R2 | G2 |
| V3 | Q3 | R3 | G3 |
| TITAN V | Quadro P4000 | TITAN RTX | GeForce RTX 2060 |

Computation power
V > R > G > Q

Memory size
R > V > Q > G

- **Two DNN models**

| | ResNet-152 | VGG-19 |
| --- | --- | --- |
| Dataset, minibatch size | ImageNet, 32 | |
| Model parameter size | 230 MB | 548 MB |
| Characteristic | Large activation output | Large parameter size |

- **NP (Node Partition)**

| VW 1 | VW 2 | VW 3 | VW 4 |
|------|------|------|------|
| Node 1 | Node 2 | Node 3 | Node 4 |
| V0 | Q0 | R0 | G0 |
| V1 | Q1 | R1 | G1 |
| V2 | Q2 | R2 | G2 |
| V3 | Q3 | R3 | G3 |

- Minimum communication overhead within VW

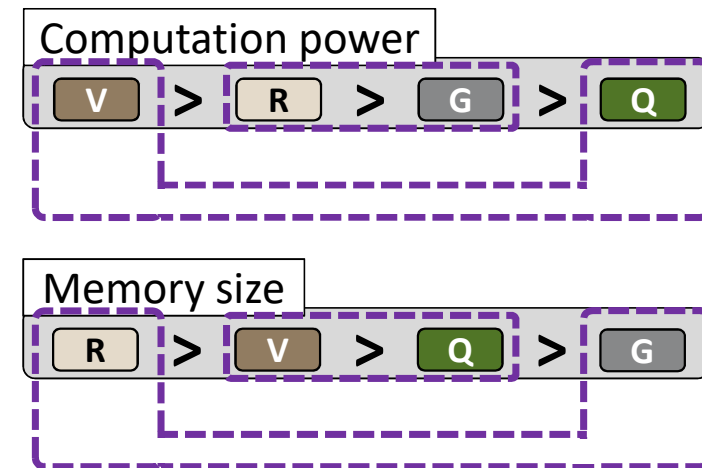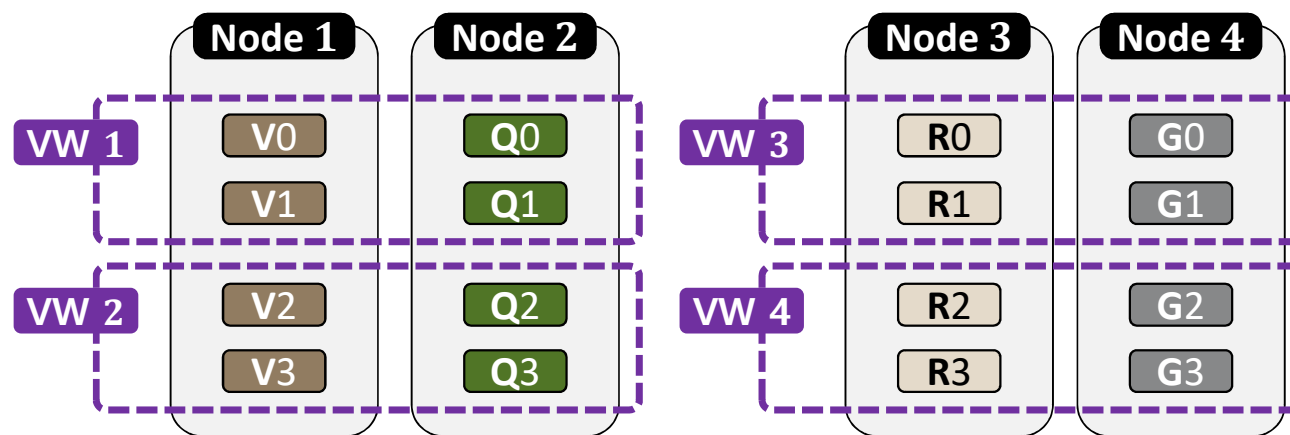- Performance of each virtual worker varies
- Straggler may degrade performance with DP

- **ED (Equal Distribution)**



- Performance will be the same across the VWs
- Mitigates the straggler problem

- High communication overhead within each VW

- **HD (Hybrid Distribution)**

| | Node 1 | Node 2 | | Node 3 | Node 4 |
|---|---|---|---|---|---|
| VW 1 | V0 / V1 | Q0 / Q1 | VW 3 | R0 / R1 | G0 / G1 |
| VW 2 | V2 / V3 | Q2 / Q3 | VW 4 | R2 / R3 | G2 / G3 |

Computation power: $V > R > G > Q$

Memory size: $R > V > Q > G$

- Mitigates the straggler problem
- Reduces communication overhead within each VW

# Parameter Placement

- **Round-robin policy (default)**
  - Can be used in all three policies: **NP**, **ED**, and **HD**

Parameters of each layer:

Example: ED

| Node 1 | Node 2 | Node 3 | Node 4 |
|--------|--------|--------|--------|

VW 1

VW 4 | V3 | Q3 | R3 | G3 |

- **Local placement policy**
  - **ED-local** Parameters of each layer:



- Significantly reduces communication overhead
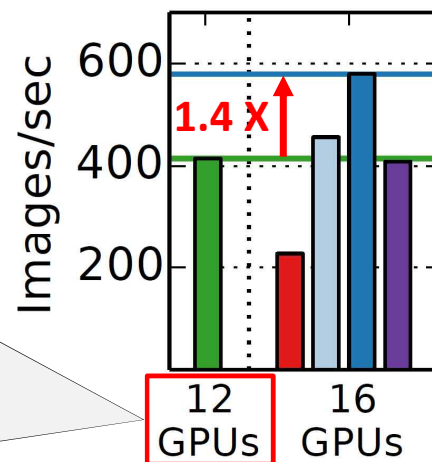
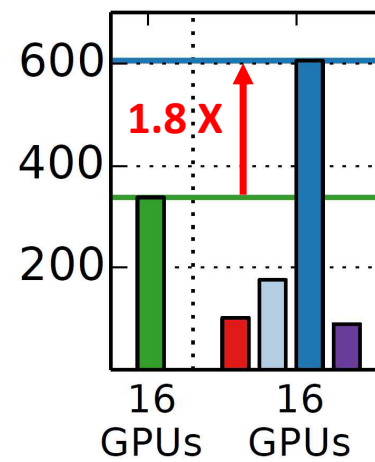- Parameter communication occurs

# Compare Throughput with Horovod

- **Baseline Horovod**
  - State-of-the-art DP using AllReduce



- For ResNet-152, the whole model is too large to be loaded into a single *G* type GPU (batch size = 32)
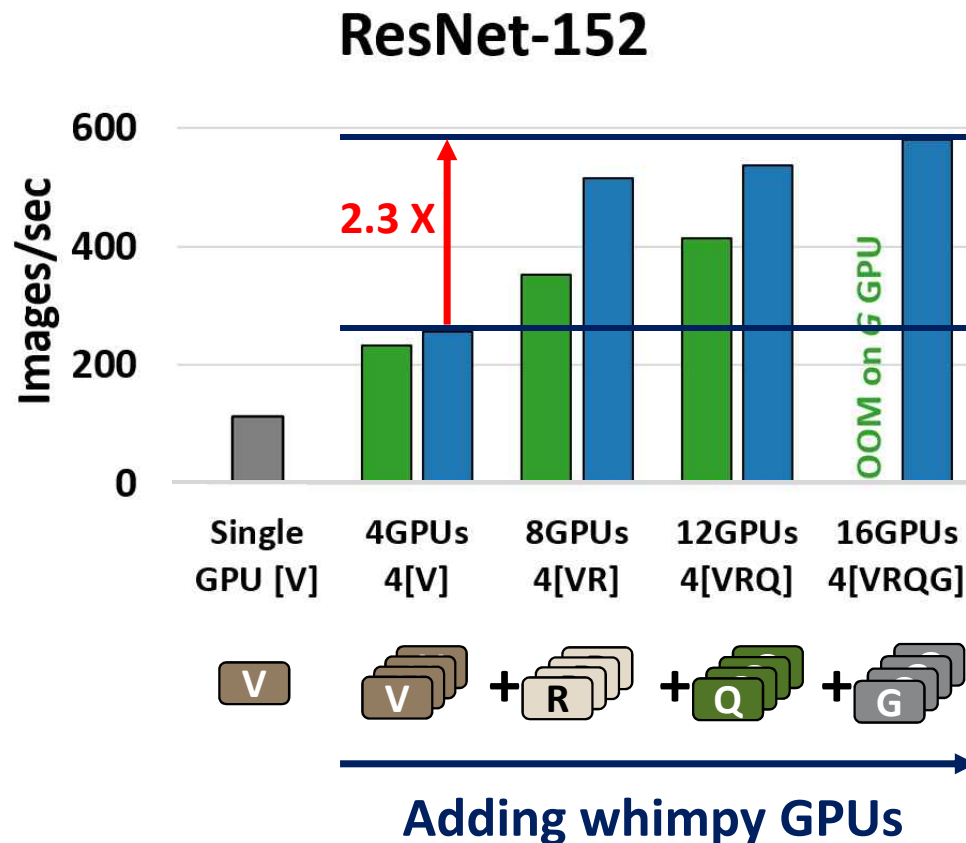
**ResNet-152**

**VGG-19**

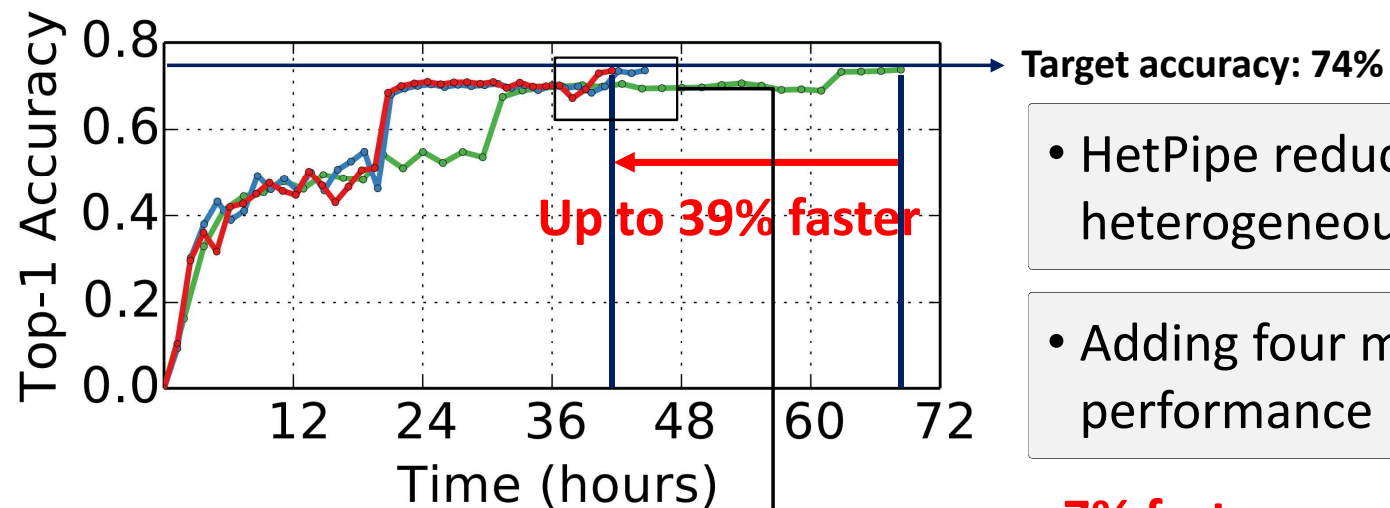- ED: reduces the straggler problem
- ED-local: significantly reduces communication overhead

- With additional GPUs, HetPipe achieves up to 2.3X speed up

- Additional whimpy systems allow for faster training

- **ResNet-152**



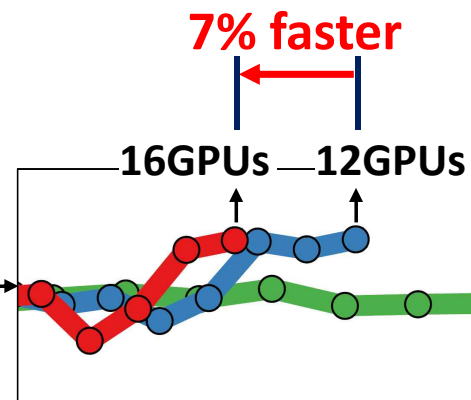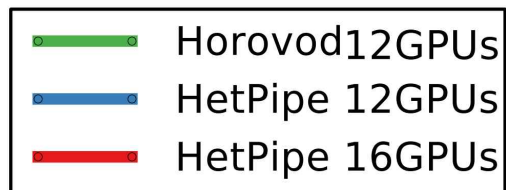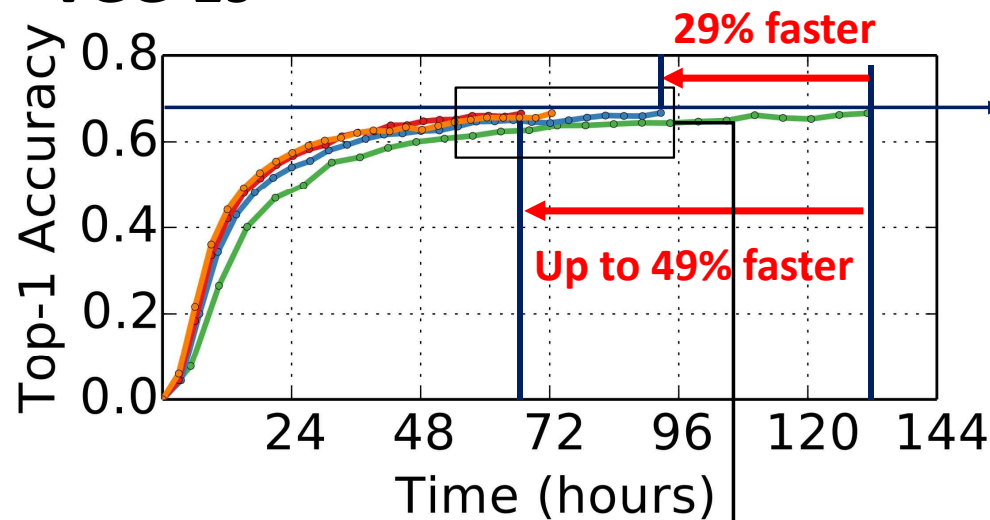Target accuracy: 74%

Up to 39% faster

- HetPipe reduces straggler problem in heterogeneous environment

- Adding four more whimpy *G* GPUs, performance improves even more

7% faster

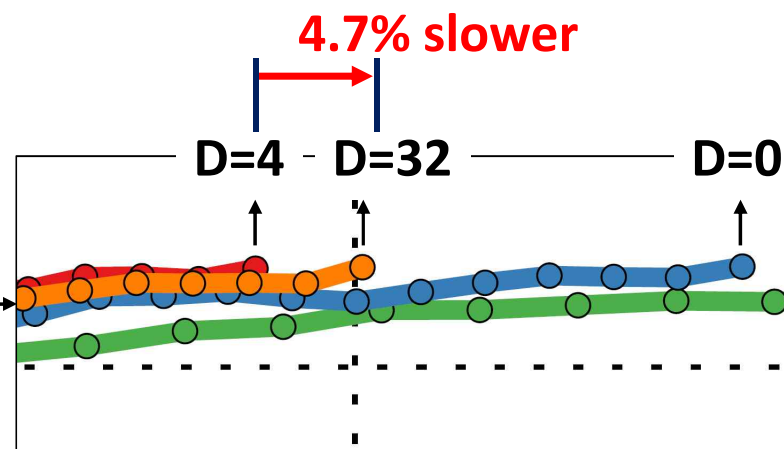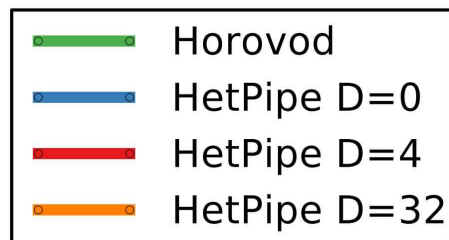16GPUs    12GPUs

Horovod12GPUs
HetPipe 12GPUs
HetPipe 16GPUs

- **VGG-19**



- HetPipe (D=0) is 29% faster than Horovod
- Higher global staleness (i.e., 32) can degrade convergence performance

37

## Conclusion

- HetPipe makes it possible to efficiently **train large DNN models with heterogeneous GPUs**

- **Integrate pipelined model parallelism with data parallelism**

- Propose a novel parameter synchronization model: **WSP**

- DNN models converge up to **49% faster with HetPipe**

# Thank you!



39