





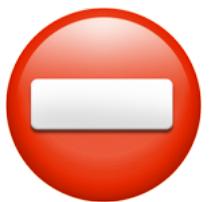
## Visibility



Lack of fine-grained  
visibility  
at queue pair level



Pre-race & race time!



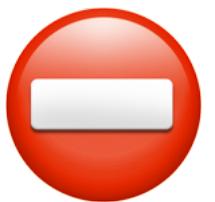
Visibility



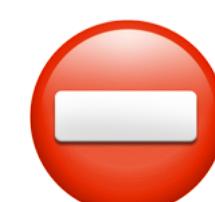
Lack of fine-grained visibility at queue pair level



Pre-race & race time!



Visibility

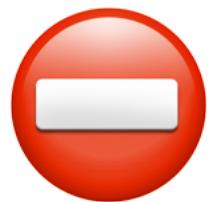


Reliability

Lack of fine-grained visibility at queue pair level



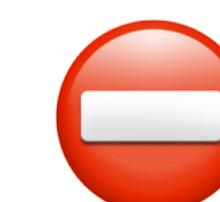
Pre-race & race time!



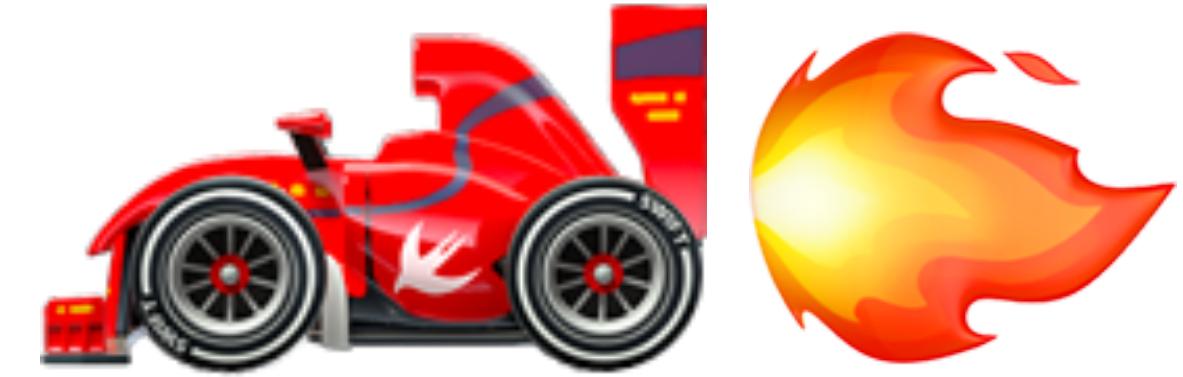
Visibility



NIC/link flapping



Reliability



# Turbocharging AI/ML Workloads

---

Revving up Speed and Resilience !

**Lerna Ekmekcioglu**  
**Sr. Solutions Engineer, Clockwork Systems**

---

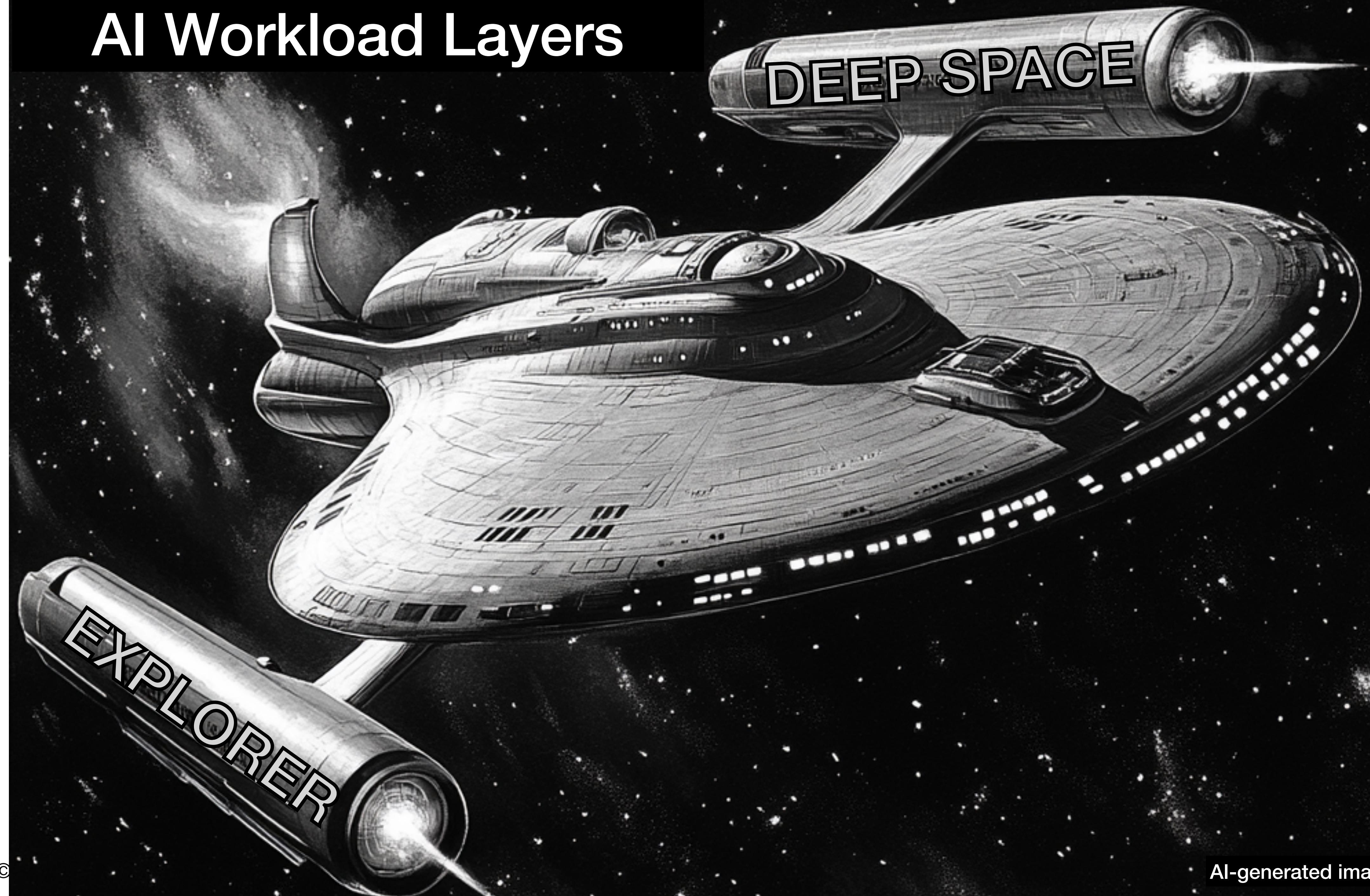
# Turbocharging AI/ML Workloads

---

Revving up Speed  and Resilience  !

**Lerna Ekmekcioglu**  
**Sr. Solutions Engineer, Clockwork Systems**

# AI Workload Layers



# AI Workload Layers

DEEP SPACE

1

DeepSpaceExplorer LLM Model

EXPLORER

# AI Workload Layers

DEEP SPACE

1

DeepSpaceExplorer LLM Model

2

Deepspeed, PyTorch, etc.

EXPLORER

# AI Workload Layers

DEEP SPACE

Large scale  
jobs

1

DeepSpaceExplorer LLM Model

2

Deepspeed, PyTorch, etc.

3

NCCL (Nvidia Collective Communications Library)

Inter-GPU  
communication

EXPLORER

# AI Workload Layers

DEEP SPACE

1

DeepSpaceExplorer LLM Model

2

Deepspeed, PyTorch, etc.

3

NCCL (Nvidia Collective Communications Library)

4

Network Device (Infiniband or RoCE)





Traditional networks



AI networks



Traditional networks



AI networks

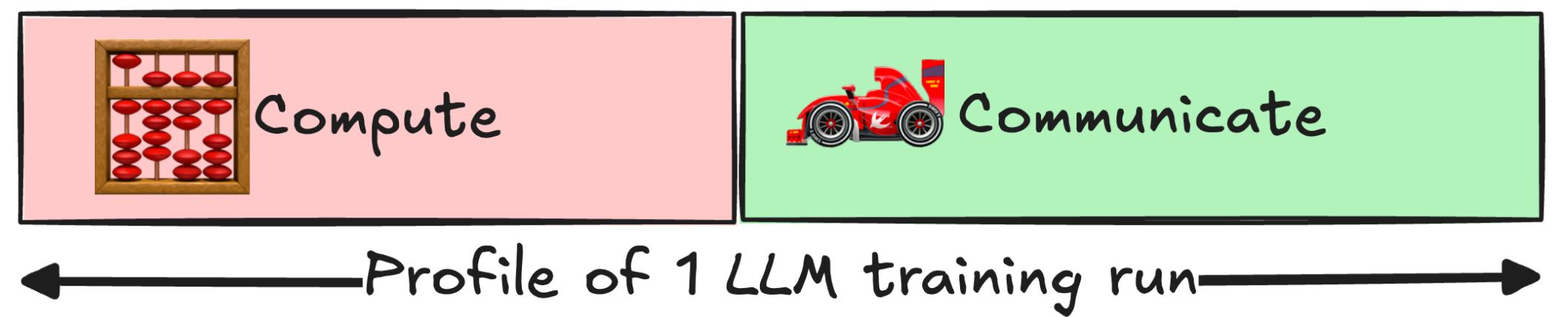


Traditional networks



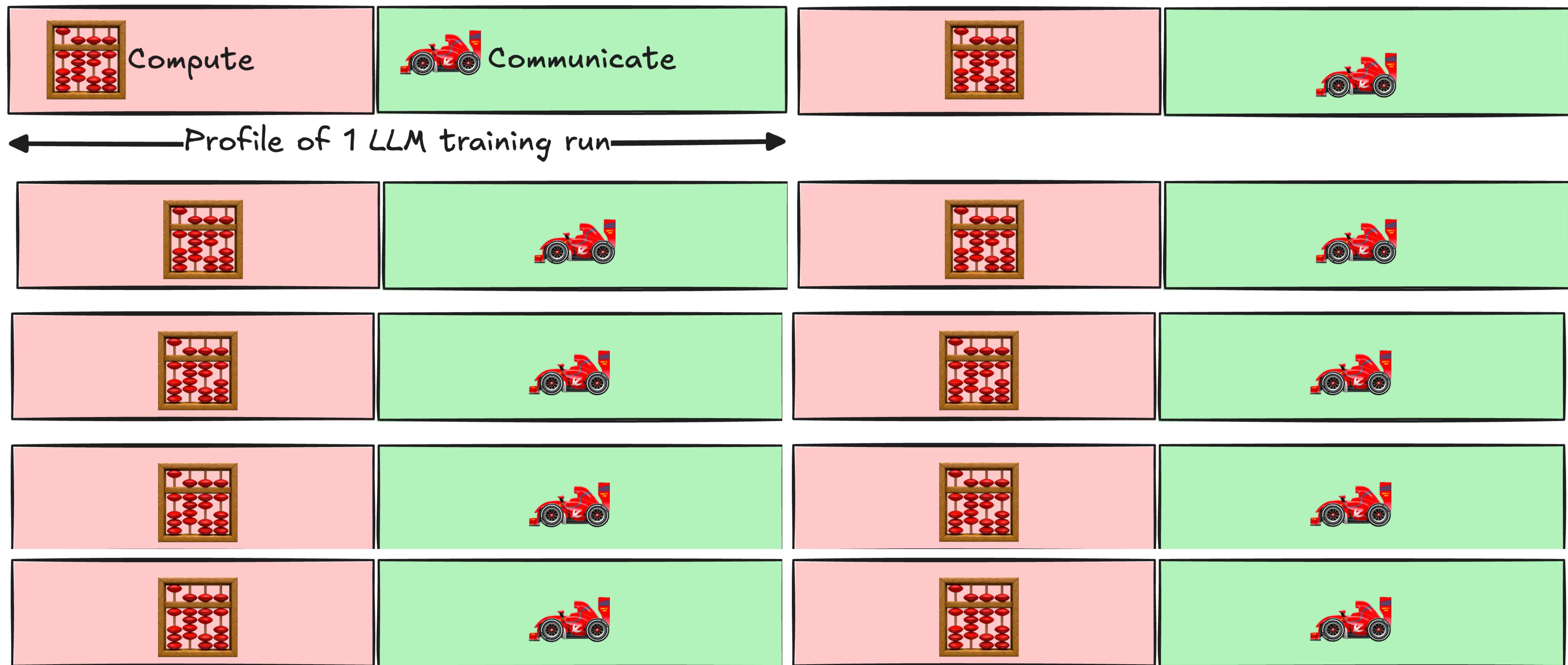
AI networks

# Profile of an LLM training job

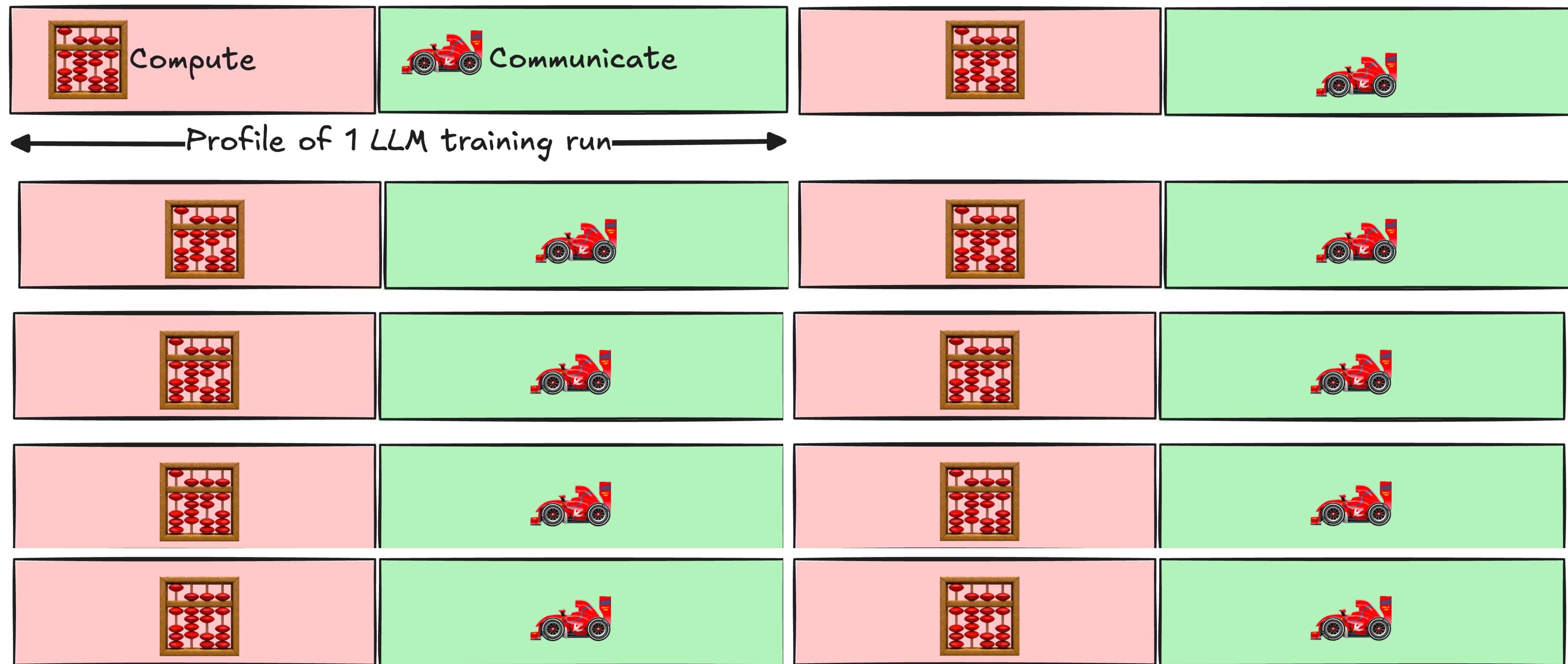


Source: Nvidia

# Profile of an LLM training job



Source: Nvidia

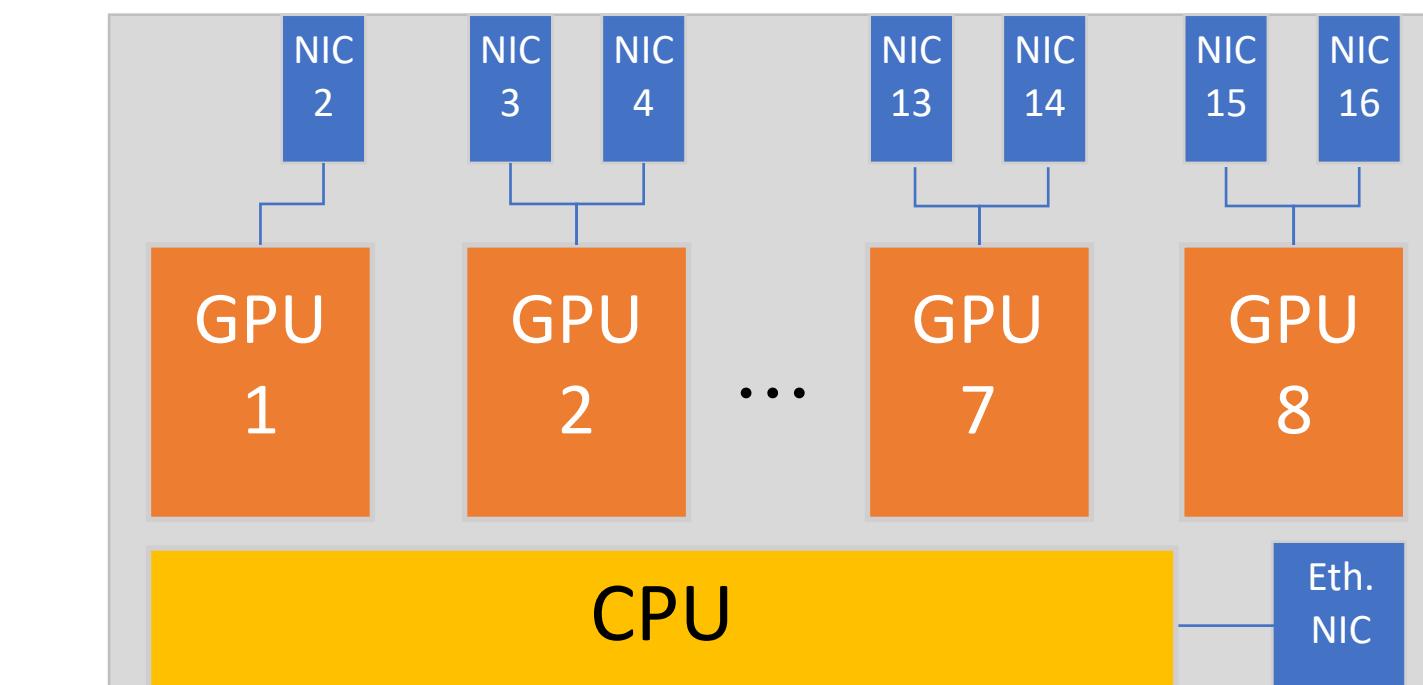
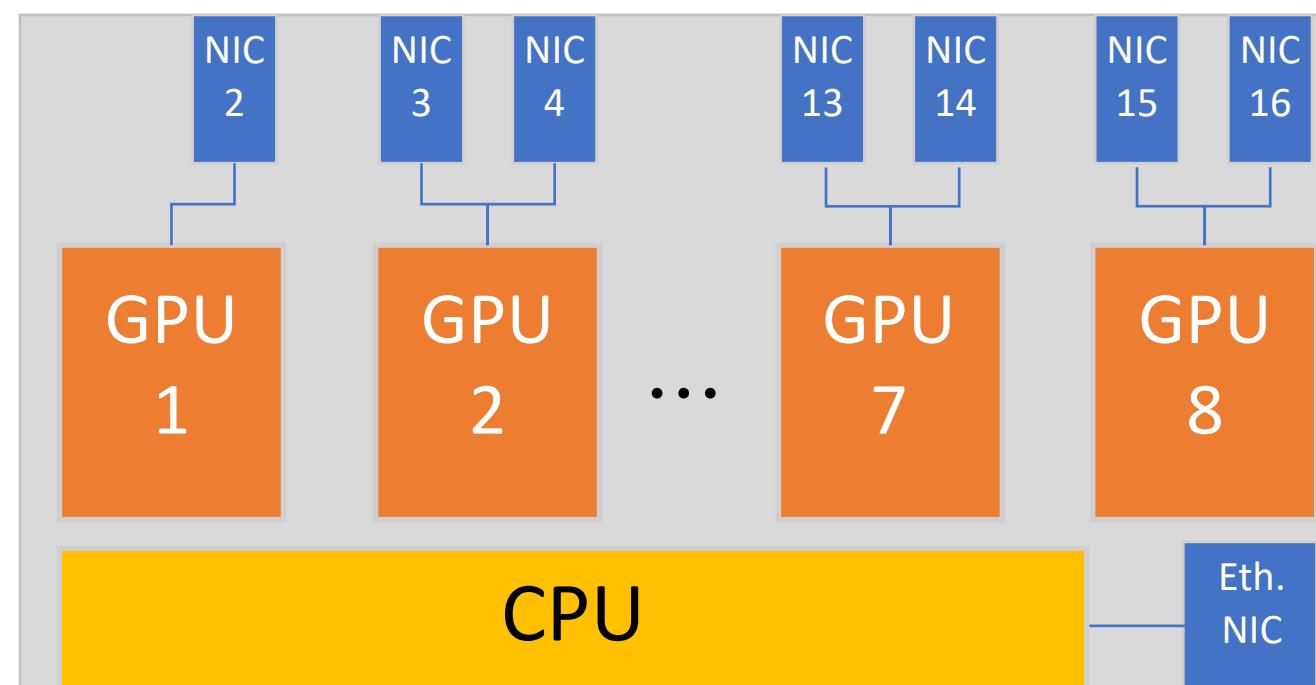


> 30 % of total job completion time  
is spent on networking!



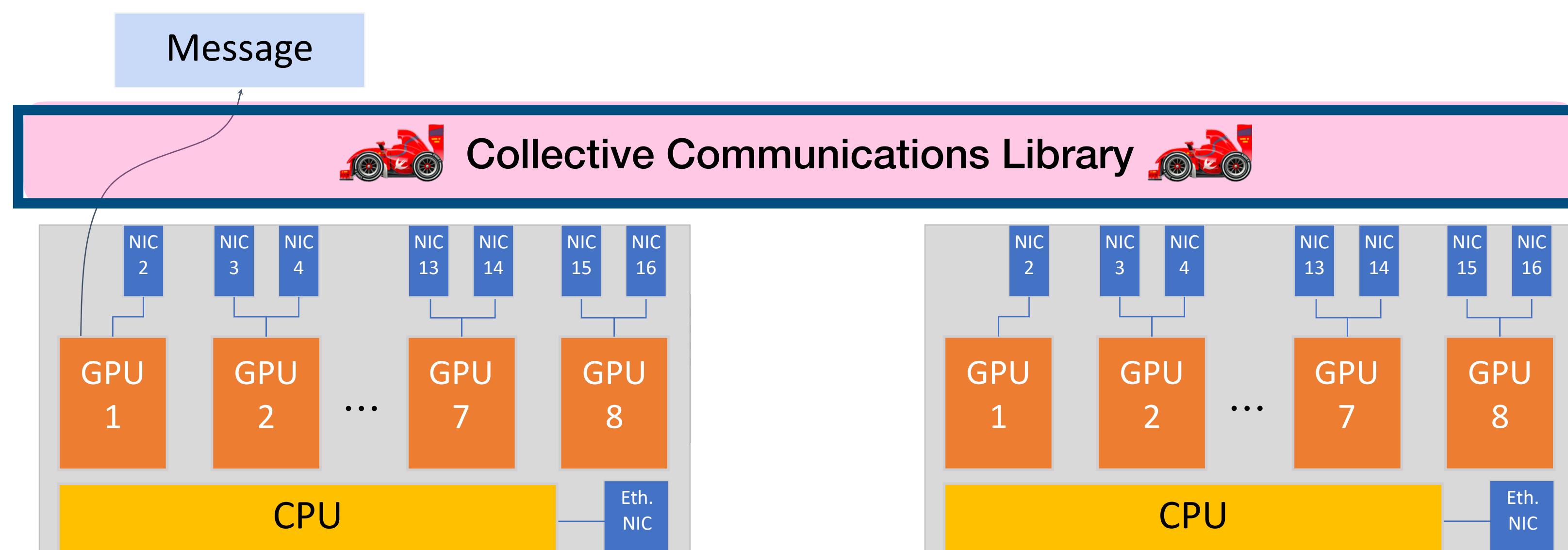
# Under the hood: NCCL or other communication library

Inter-GPU  
communication



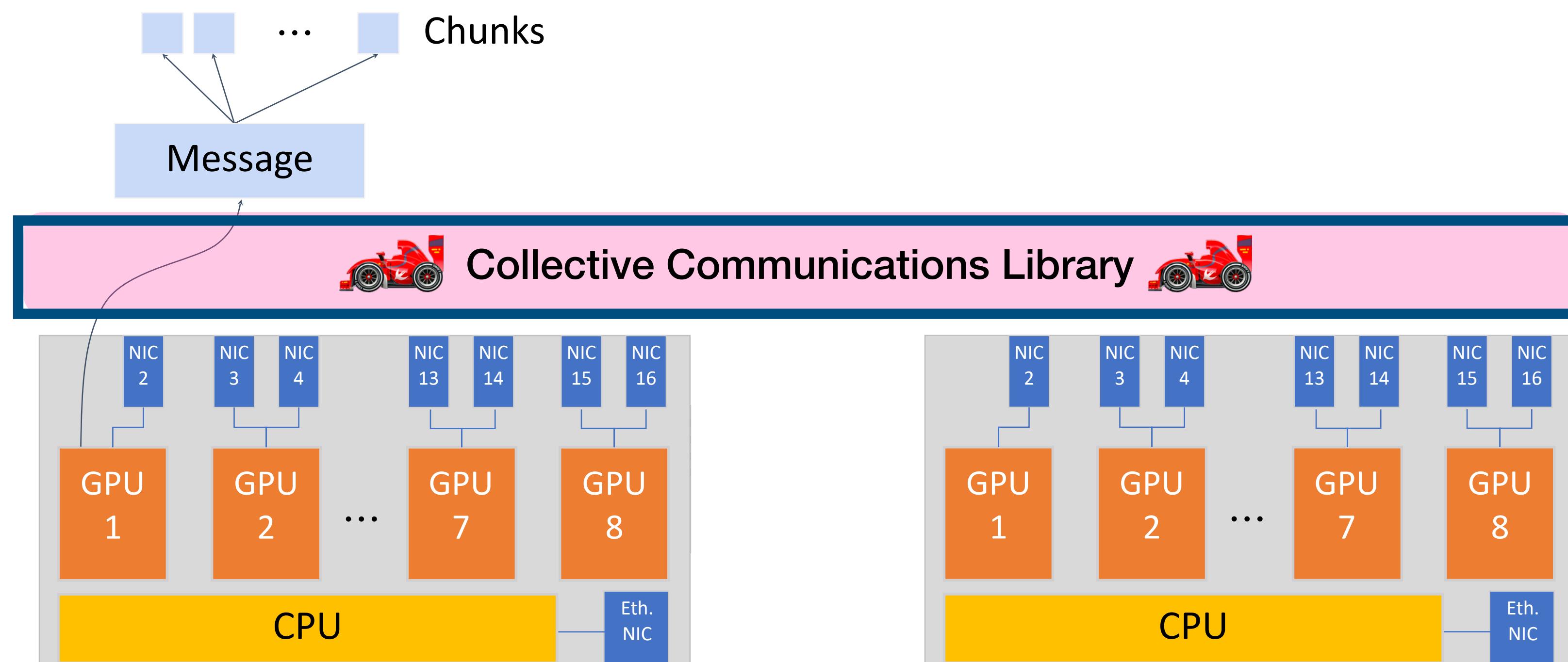
# Under the hood: NCCL or other communication library

Inter-GPU  
communication



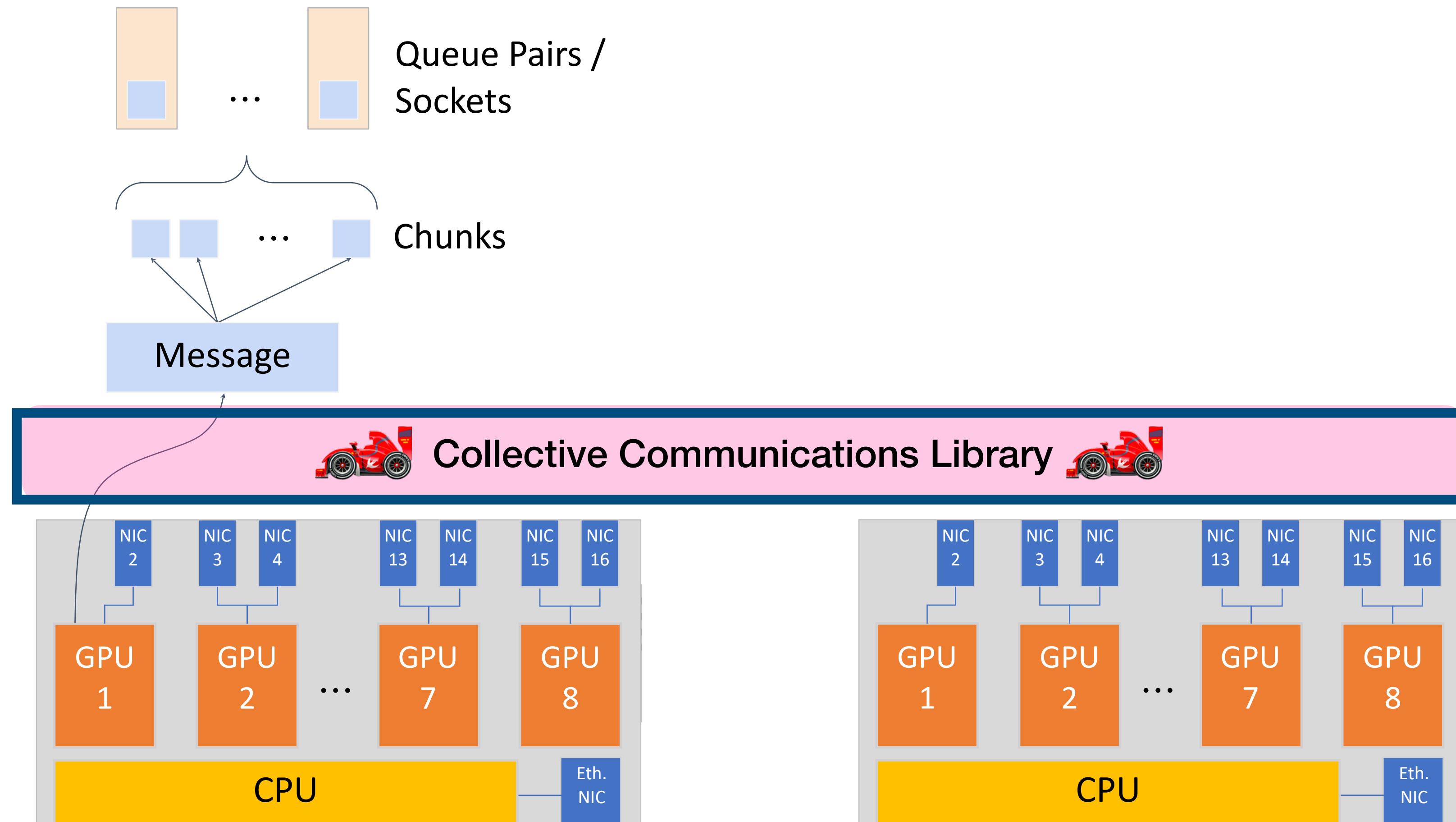
# Under the hood: NCCL or other communication library

Inter-GPU  
communication

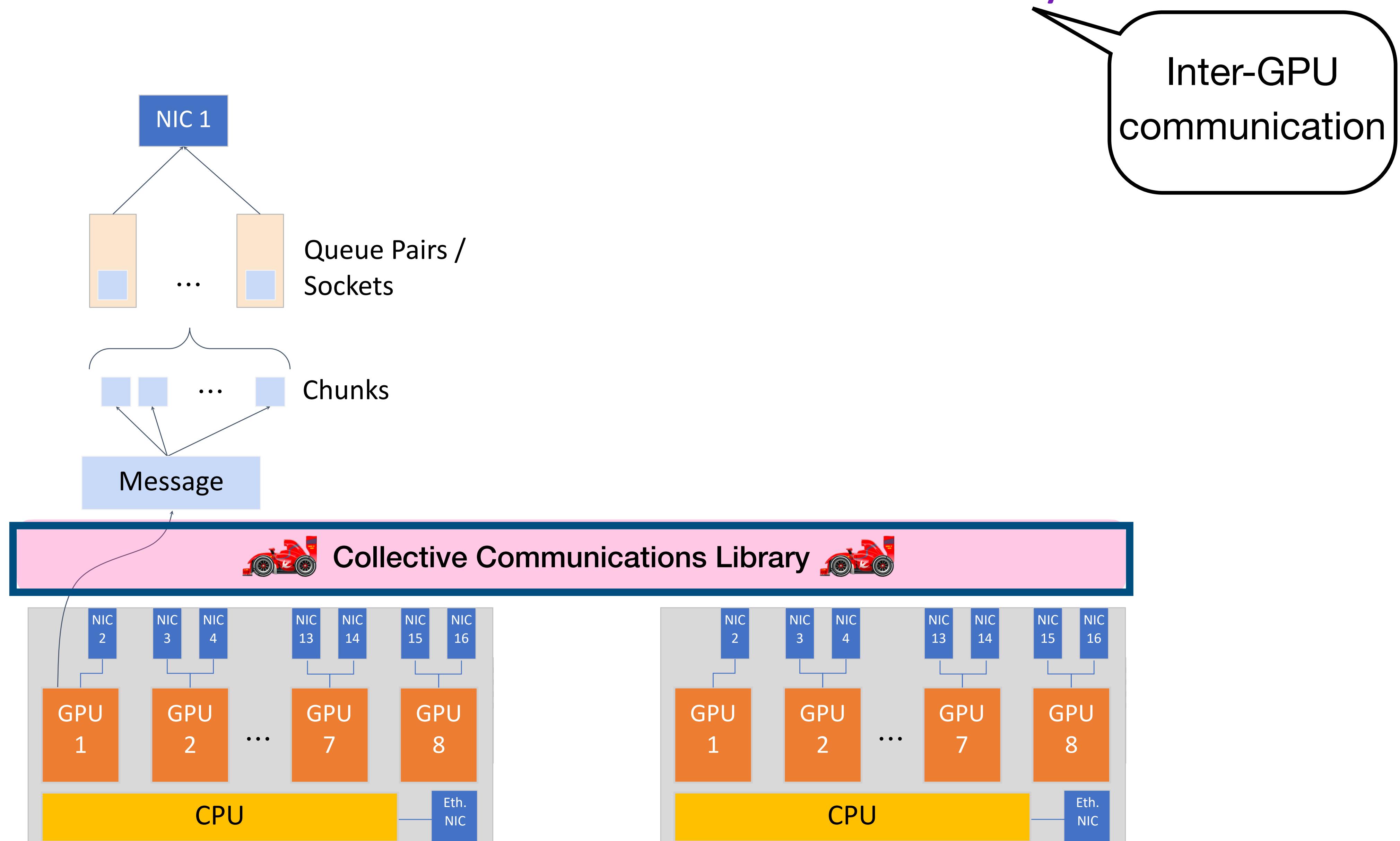


# Under the hood: NCCL or other communication library

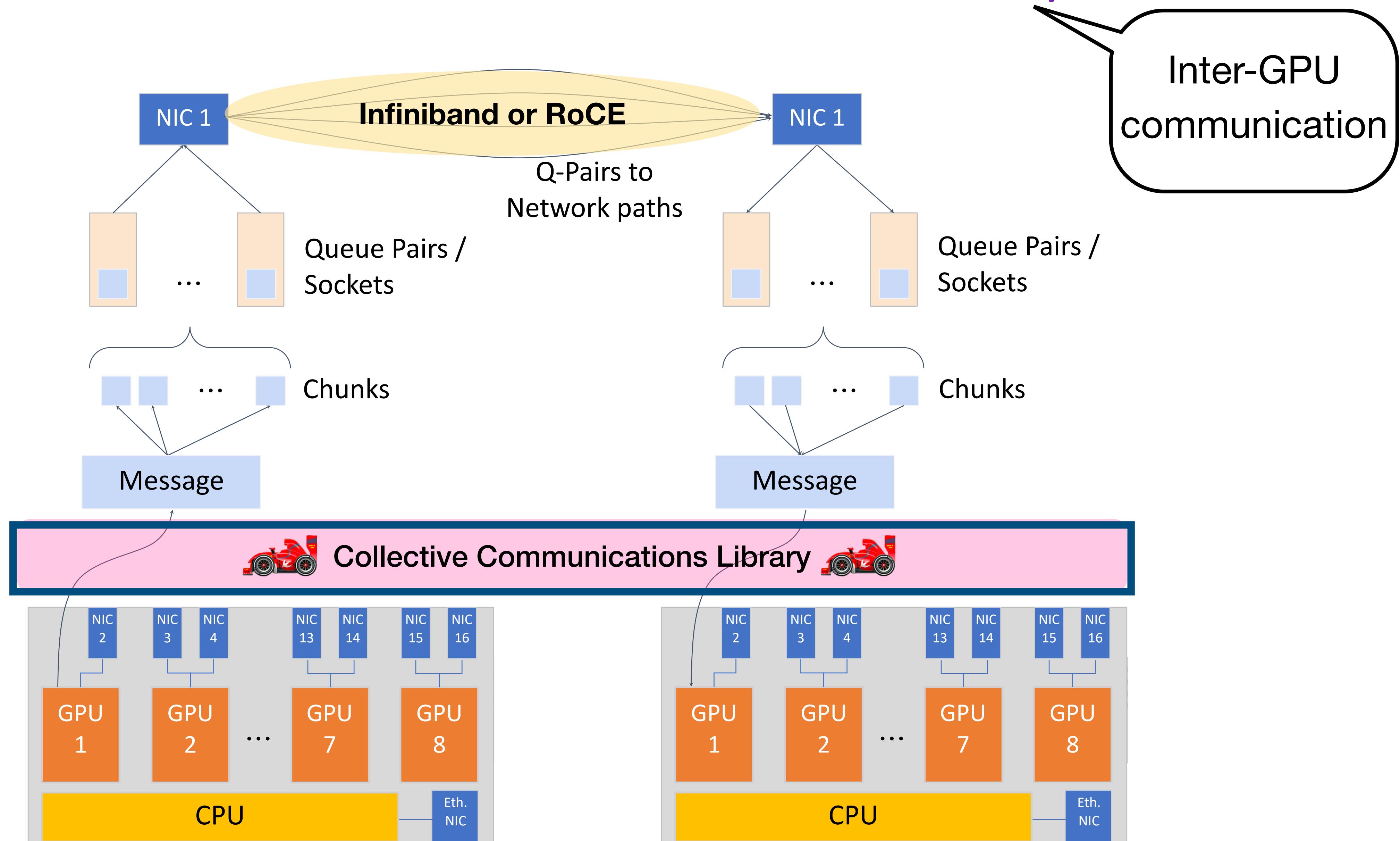
Inter-GPU  
communication

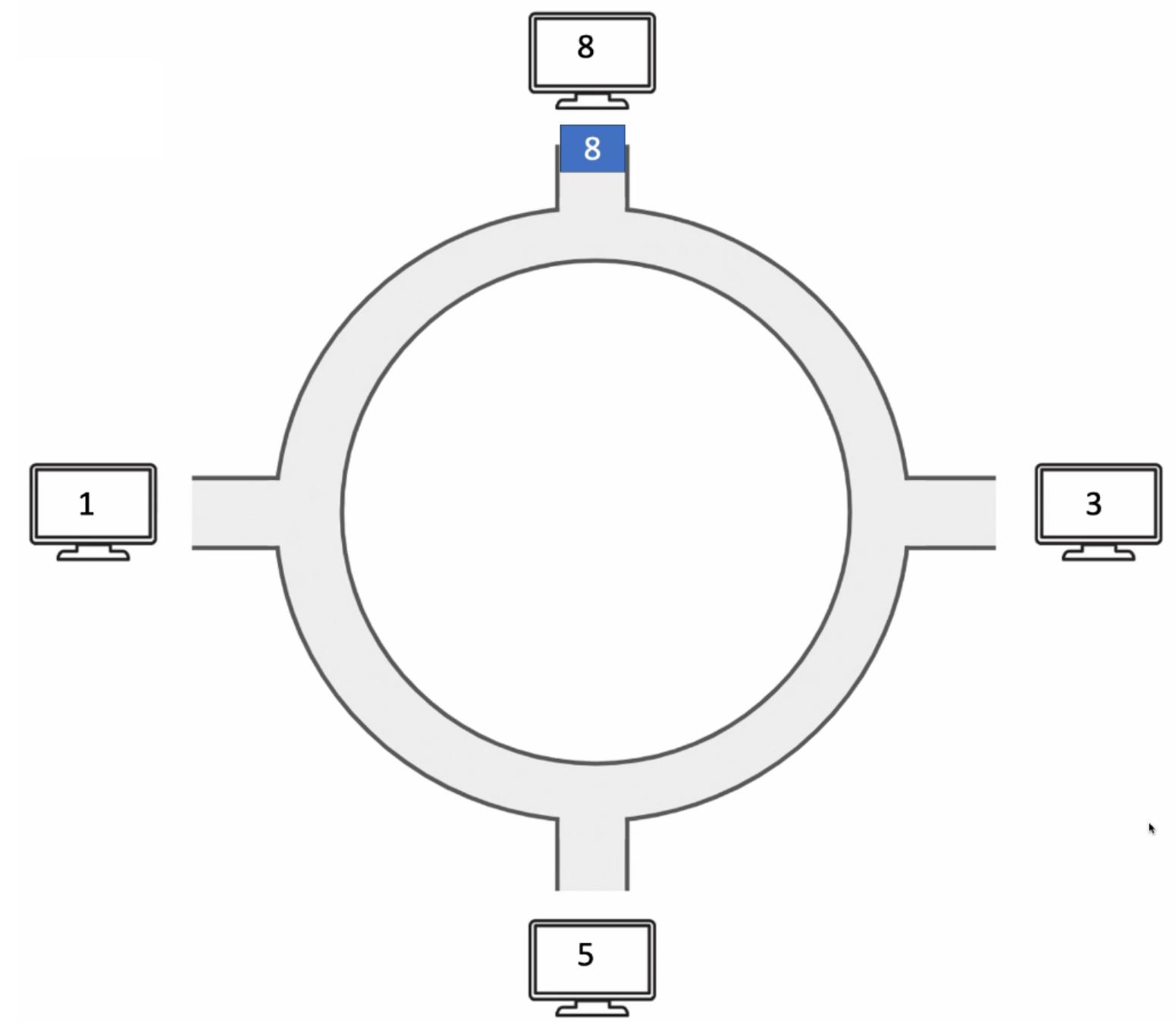


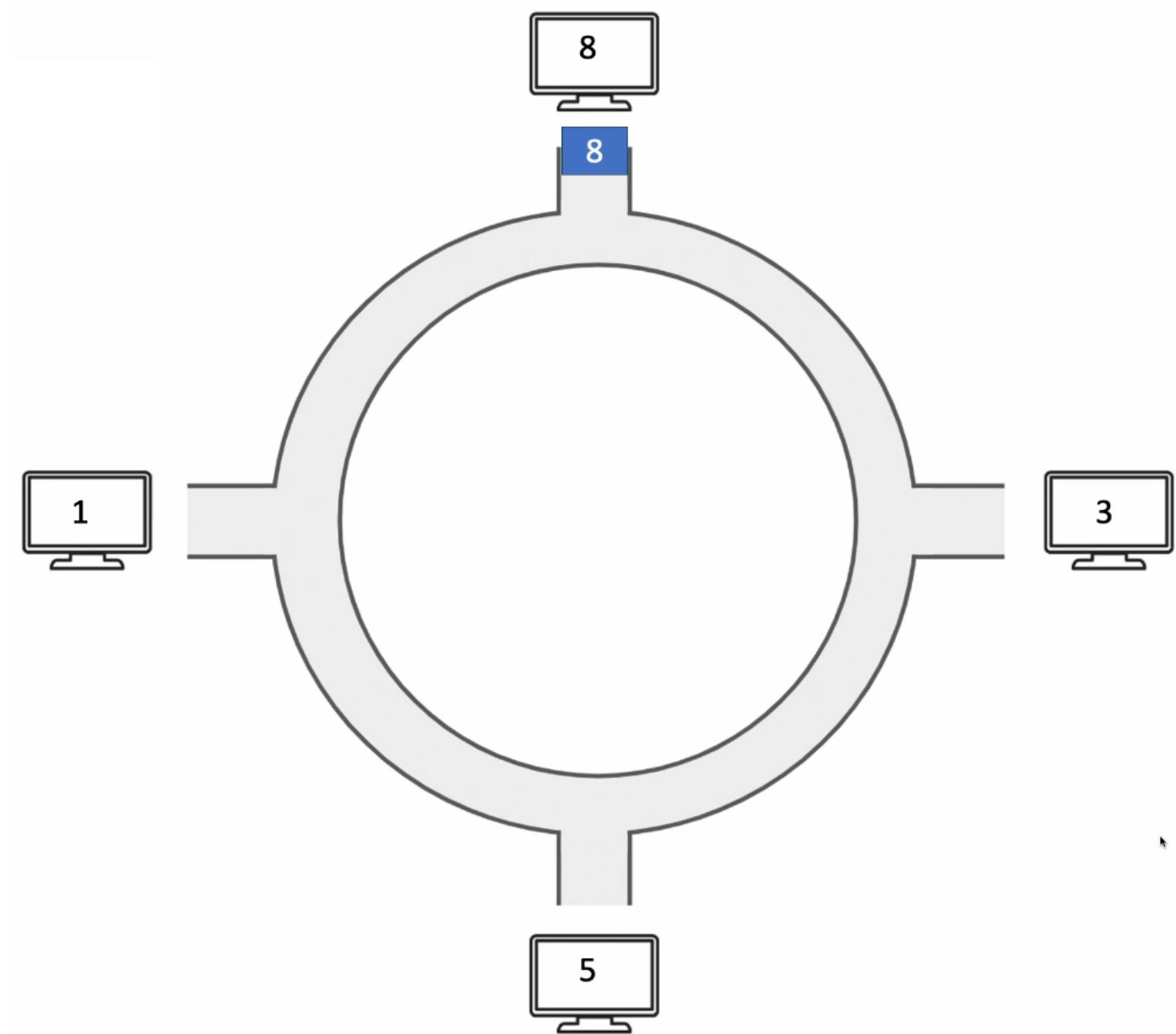
# Under the hood: NCCL or other communication library



# Under the hood: NCCL or other communication library



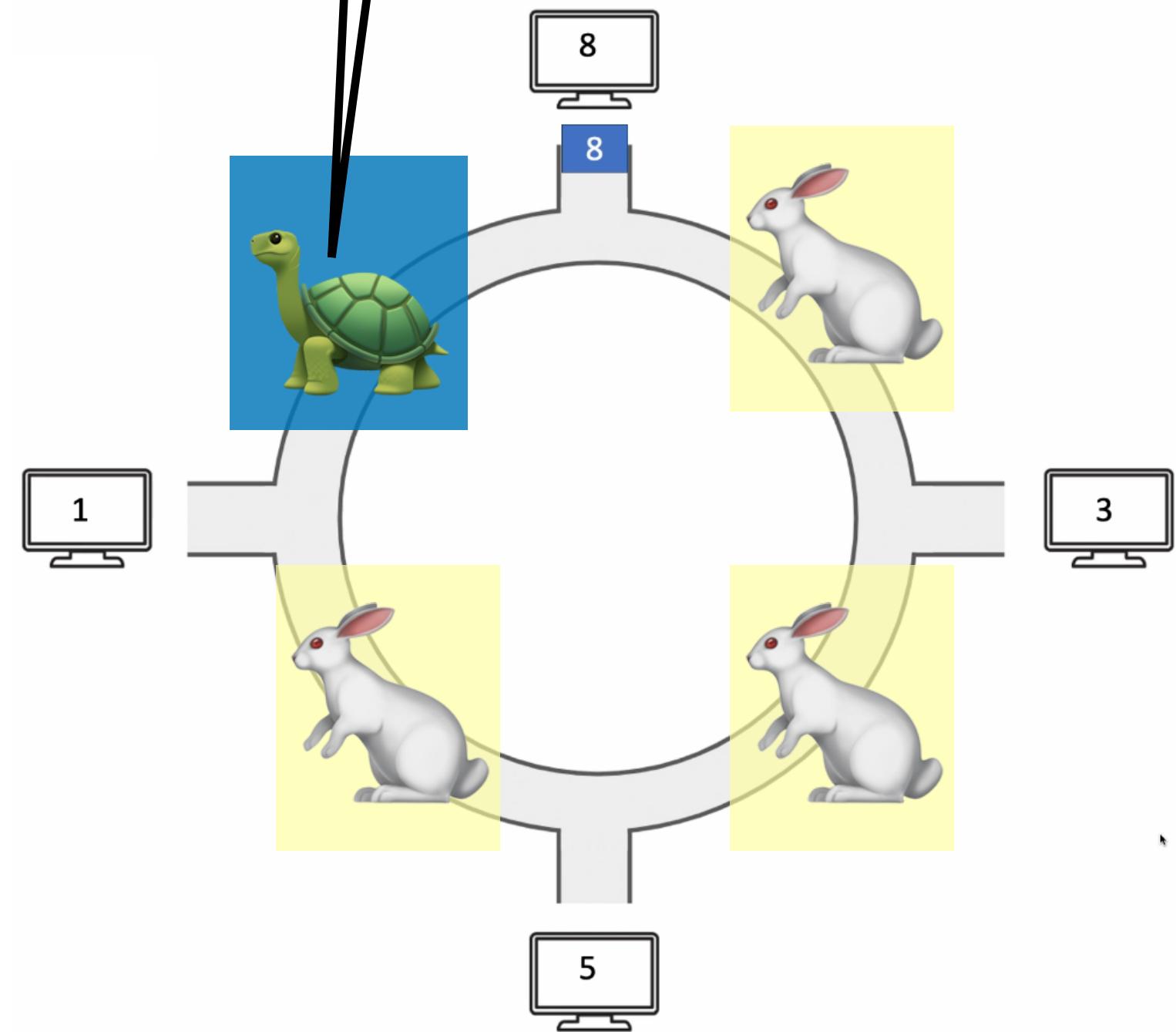




Slowest team member  
determines speed.



Slowest flow (or node)  
determines throughput.



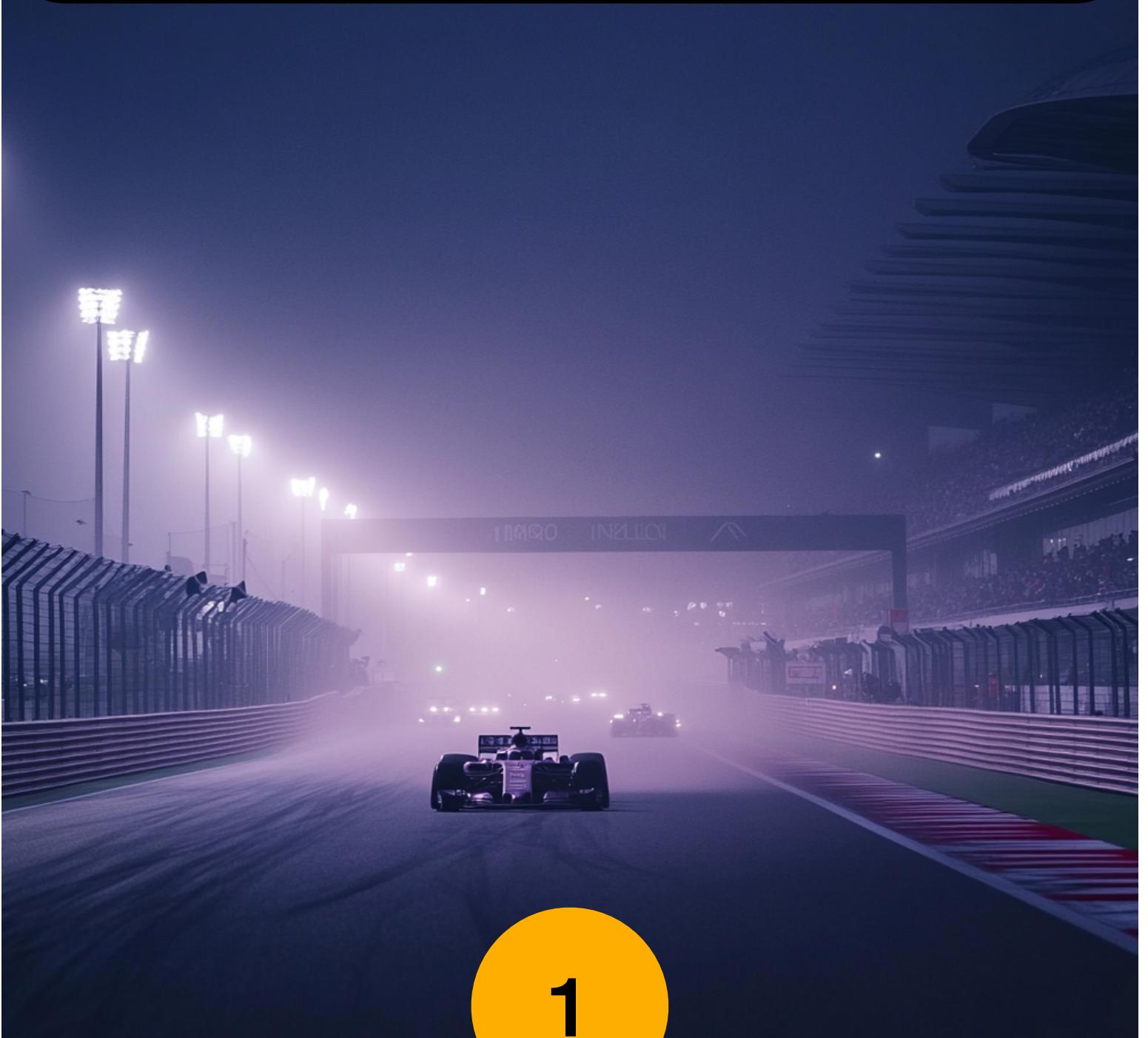
Slowest team member  
determines speed.



# Networking challenges of AI/ML workloads



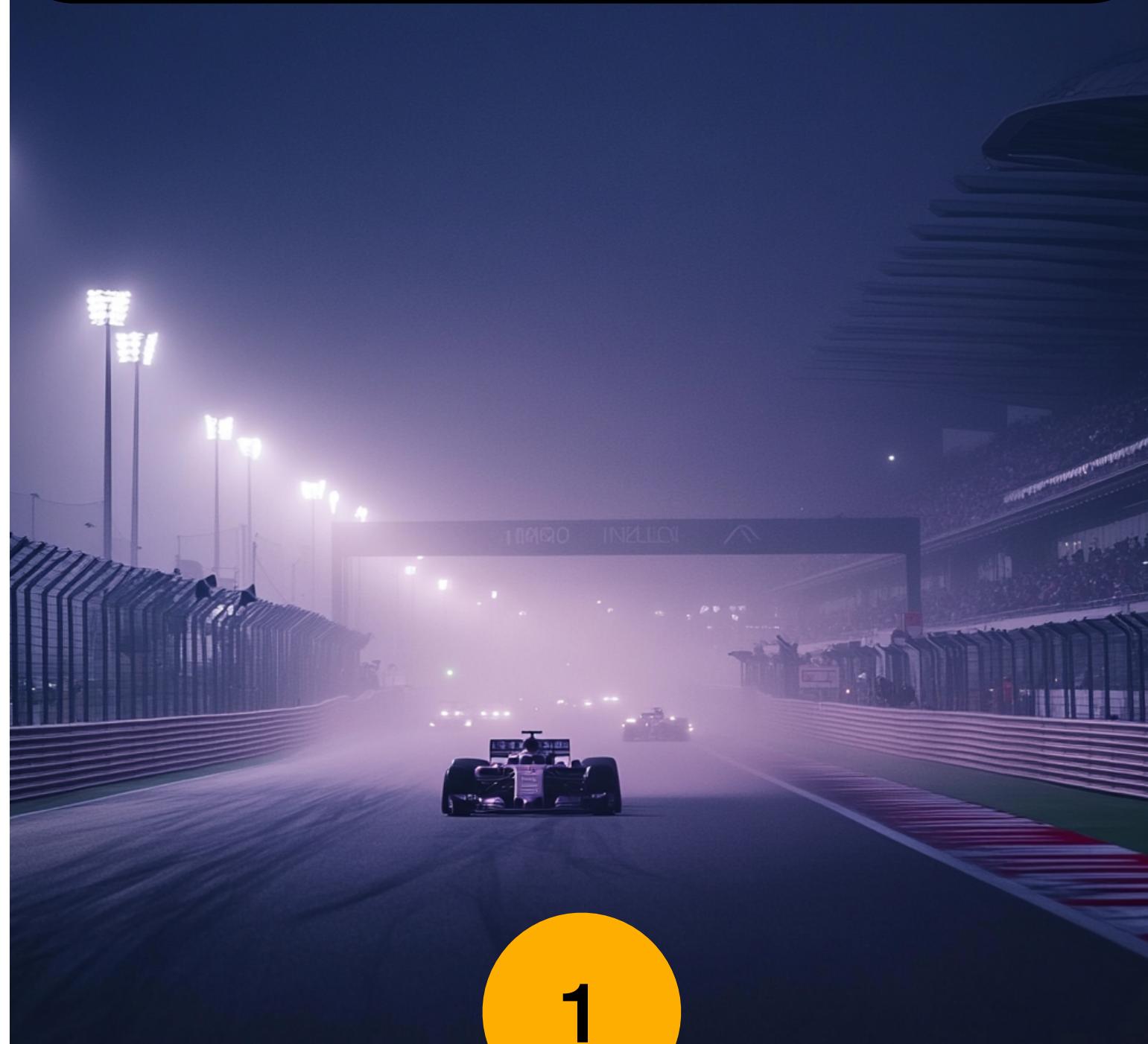
## Visibility



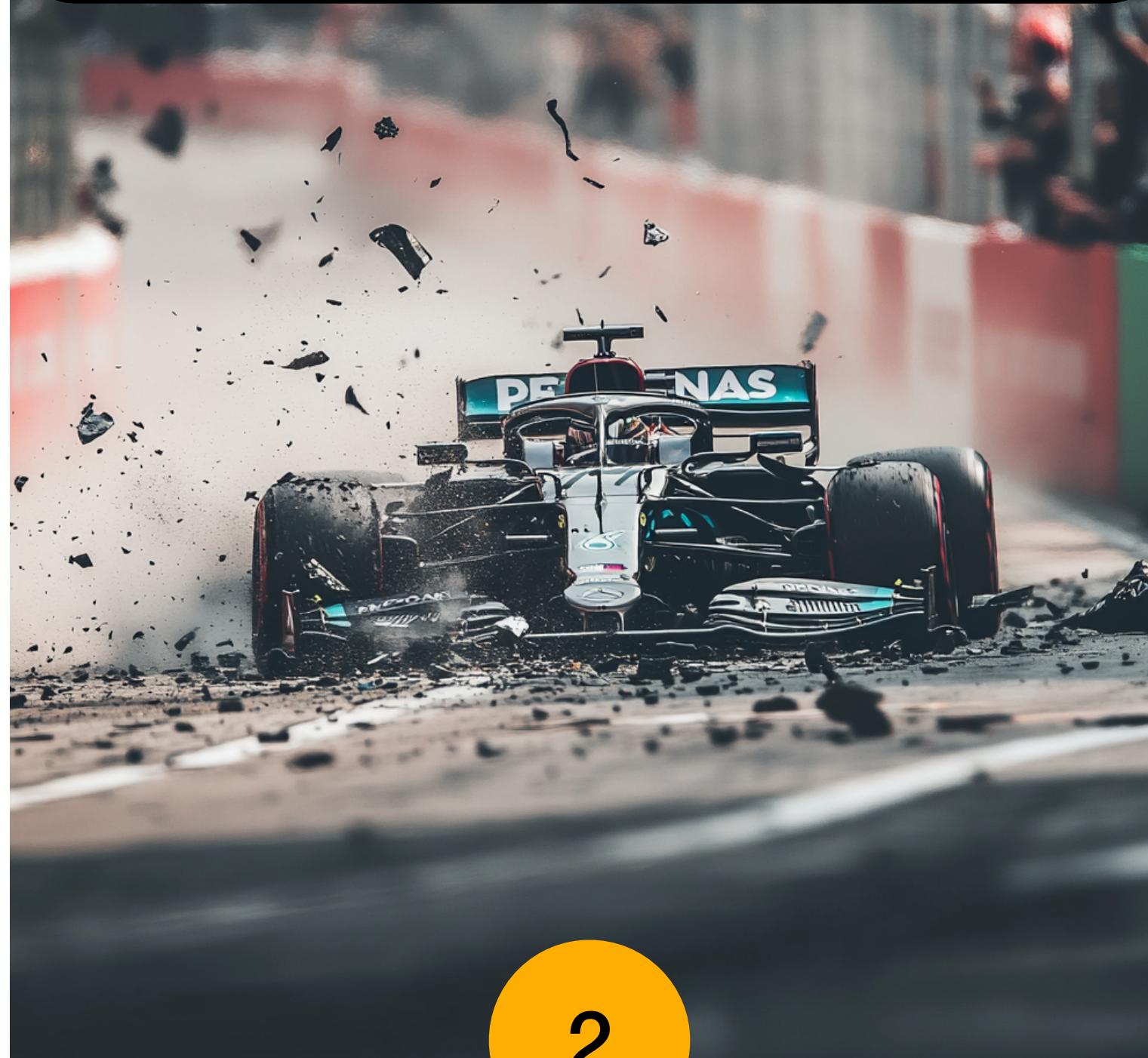
# Networking challenges of AI/ML workloads



Visibility



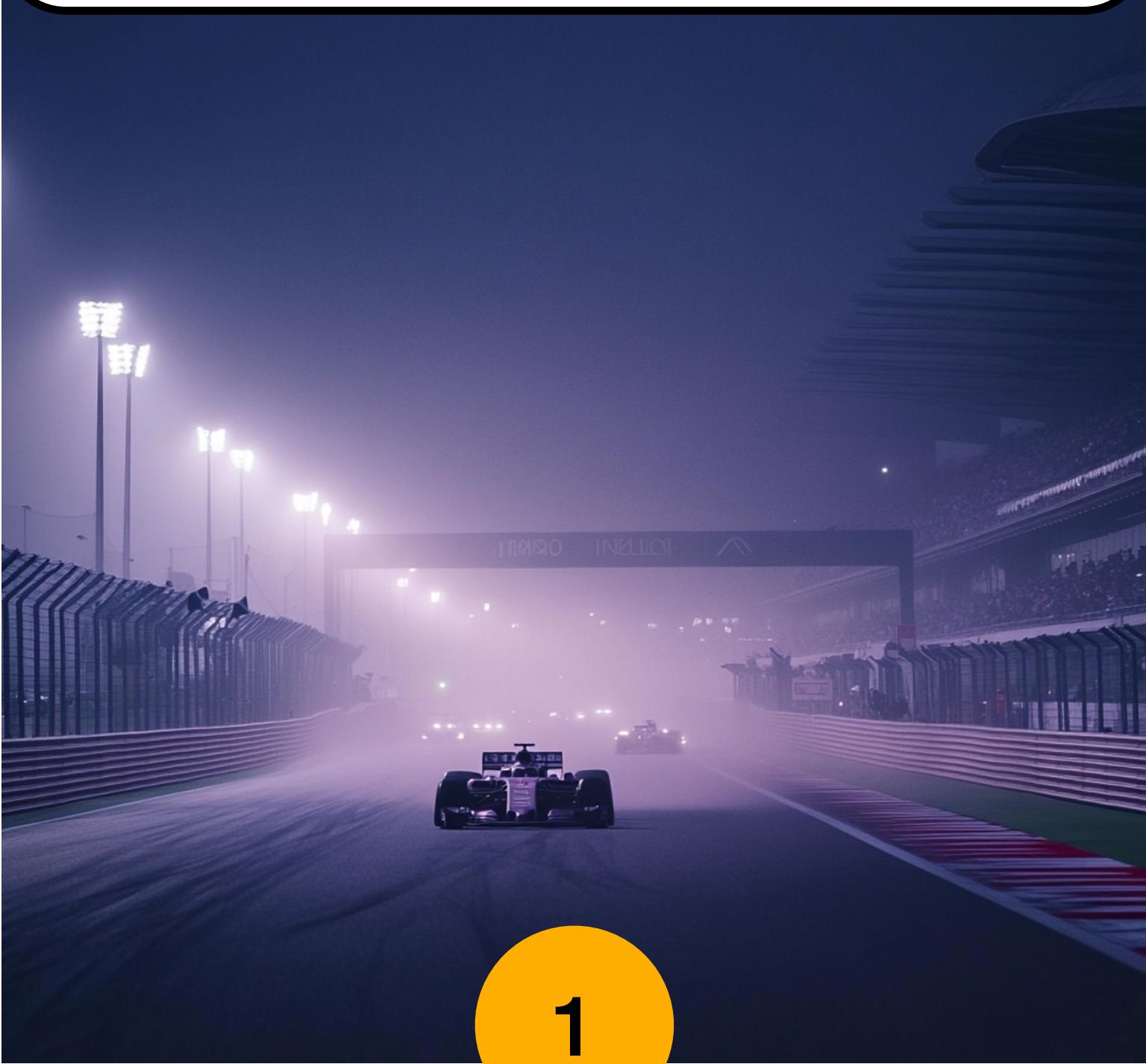
Reliability



# Networking challenges of AI/ML workloads



Visibility



1



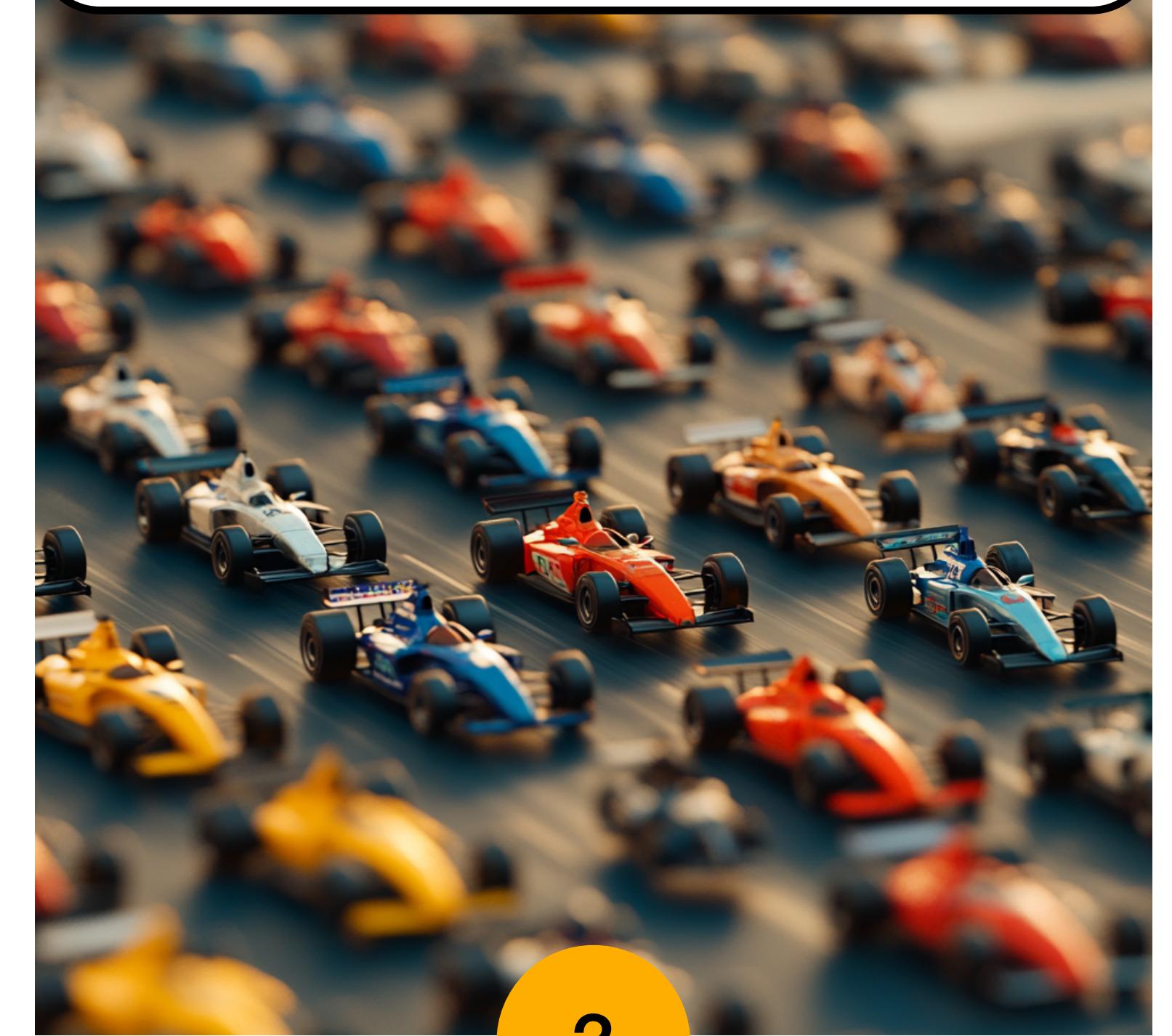
Reliability



2



Performance



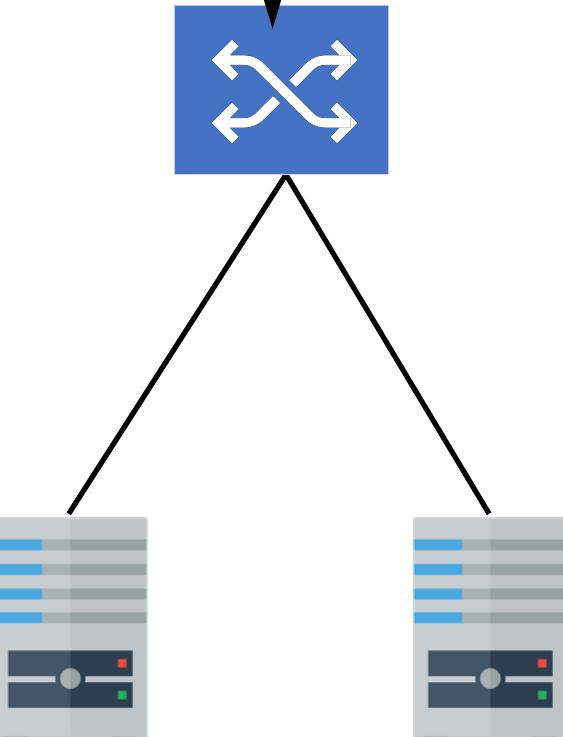
3

1

# Lack of visibility



The ToR has  
3 down uplinks  
since 6 months ago  
but we are not sure why!

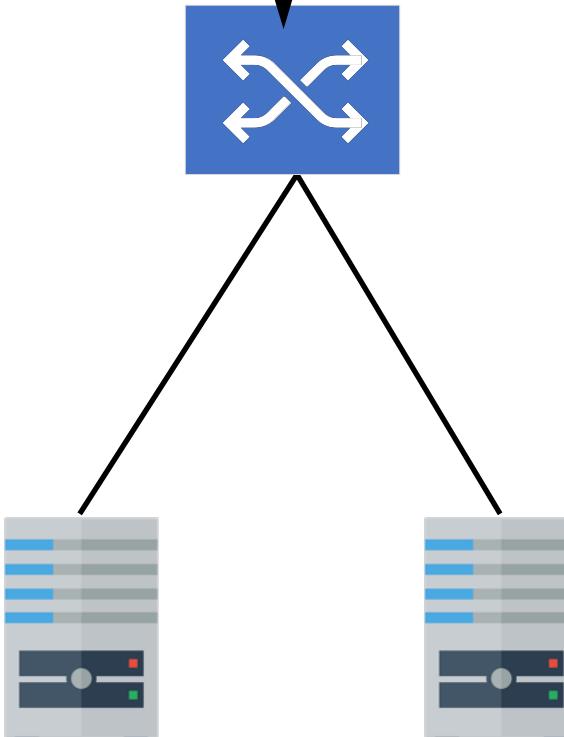


1

# Lack of visibility



The ToR has  
3 down uplinks  
since 6 months ago  
but we are not sure why!



## Why?

- ▶ Aggregate level metrics
- ▶ No fine grained instrumentation
- ▶ No visibility down to the queue pair level

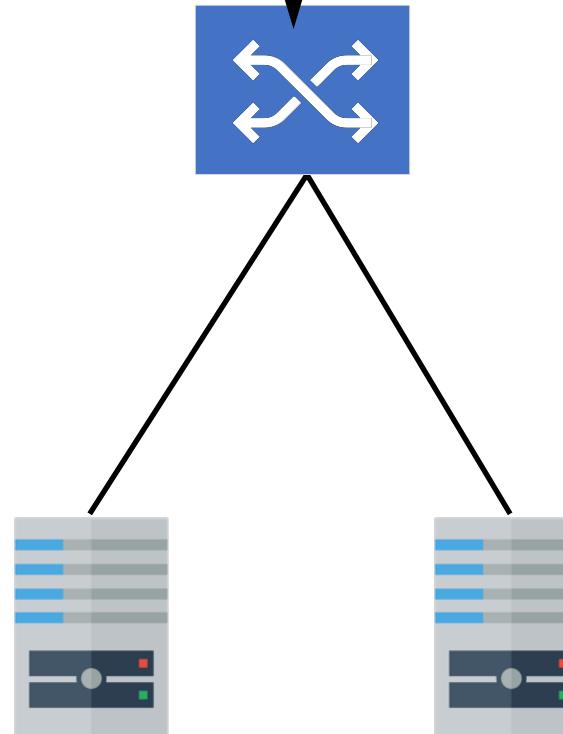


1

# Lack of visibility



The ToR has  
3 down uplinks  
since 6 months ago  
but we are not sure why!



## Why?

- ▶ Aggregate level metrics
- ▶ No fine grained instrumentation
- ▶ No visibility down to the queue pair level



## So what?

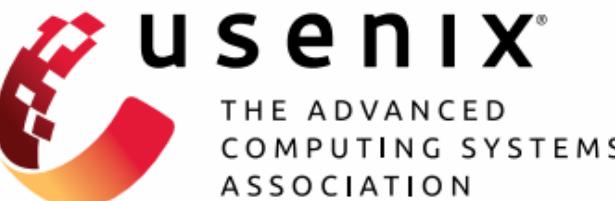
- ▶ Unable to identify root cause to resolve issues quickly.
- ▶ Not only at fleet level but at workload level

# Visibility



# Huygens clock synchronization

- ▶ Software based
- ▶ High precision at scale 



The Advanced Computing Systems Association

## Exploiting a Natural Network Effect for Scalable, Fine-grained Clock Synchronization

Yilong Geng, Shiyu Liu, and Zi Yin, Stanford University; Ashish Naik, Google Inc.;  
Balaji Prabhakar and Mendel Rosenblum, Stanford University; Amin Vahdat, Google Inc.

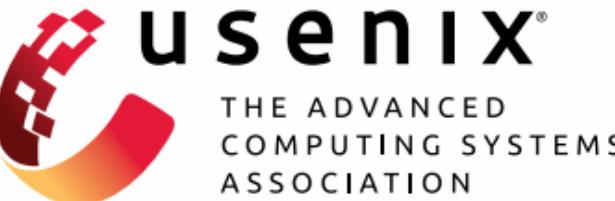
<https://www.usenix.org/conference/nsdi18/presentation/geng>



# Huygens clock synchronization

- ▶ Software based
- ▶ High precision at scale 

Measures  
one way delays



Exploiting a Natural Network Effect for Scalable,  
Fine-grained Clock Synchronization

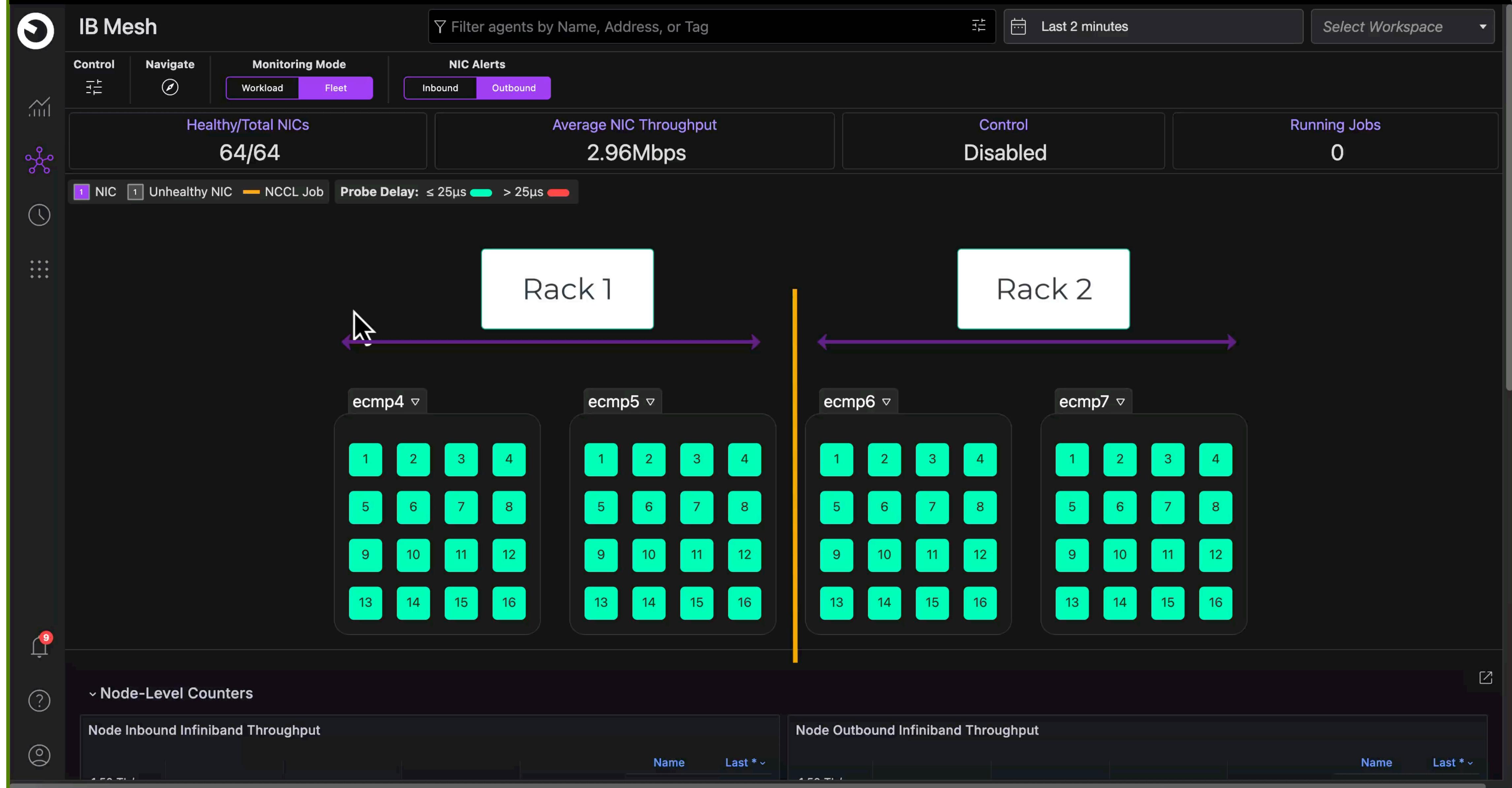
Yilong Geng, Shiyu Liu, and Zi Yin, Stanford University; Ashish Naik, Google Inc.;  
Balaji Prabhakar and Mendel Rosenblum, Stanford University; Amin Vahdat, Google Inc.

<https://www.usenix.org/conference/nsdi18/presentation/geng>

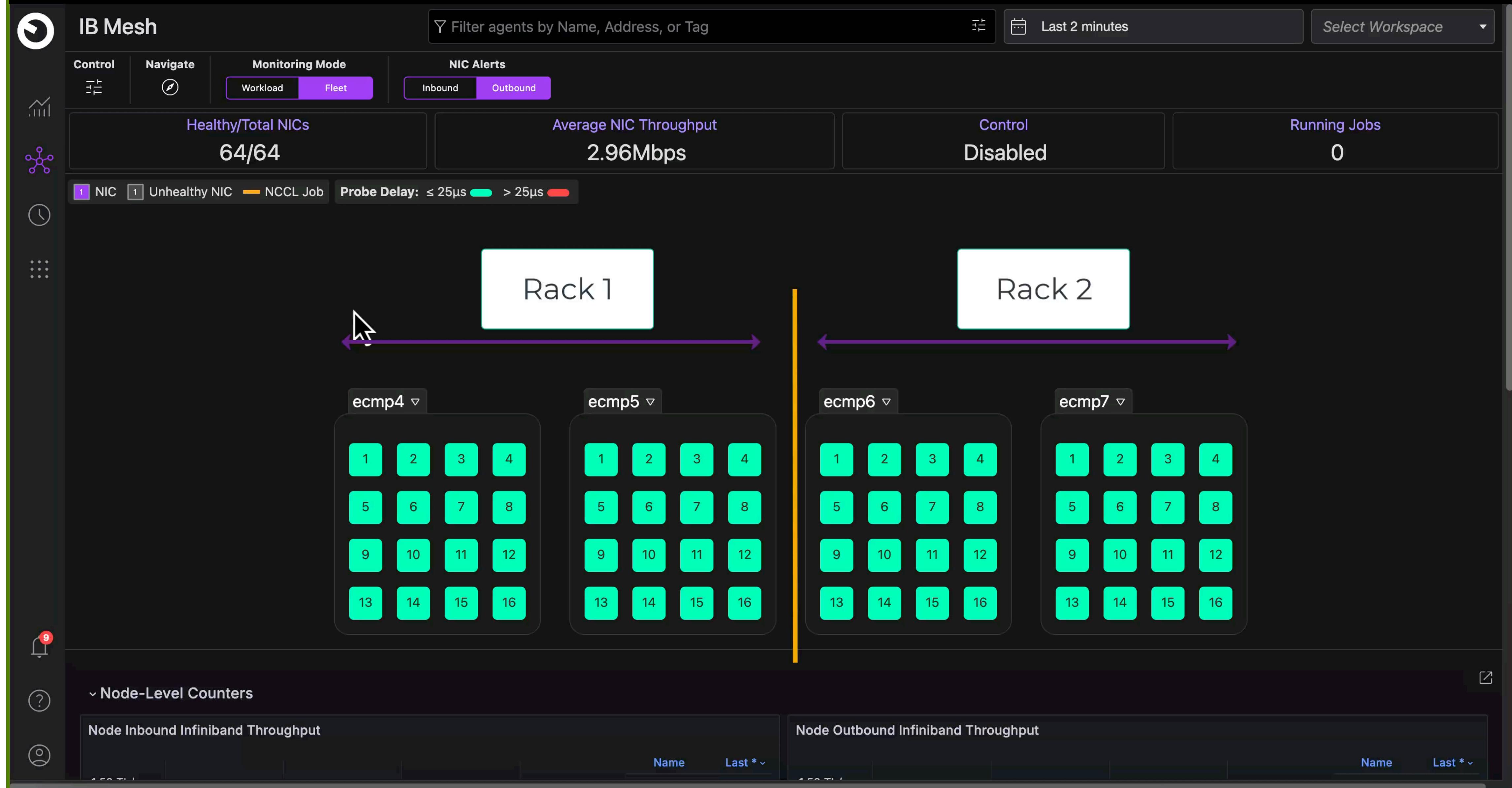


AI-generated image

# Visibility demo



# Visibility demo



2

## Lack of reliability



**Alert**



**Link flapped 10 times!**  
Send technicians  
to swap cables!



2

## Lack of reliability



Link flapped 10 times!  
Send technicians  
to swap cables!



Common in large clusters

Why?



Source: [semianalysis](#)

AI-generated images

2

## Lack of reliability



Link flapped 10 times!  
Send technicians  
to swap cables!



### Why?

Common in large clusters



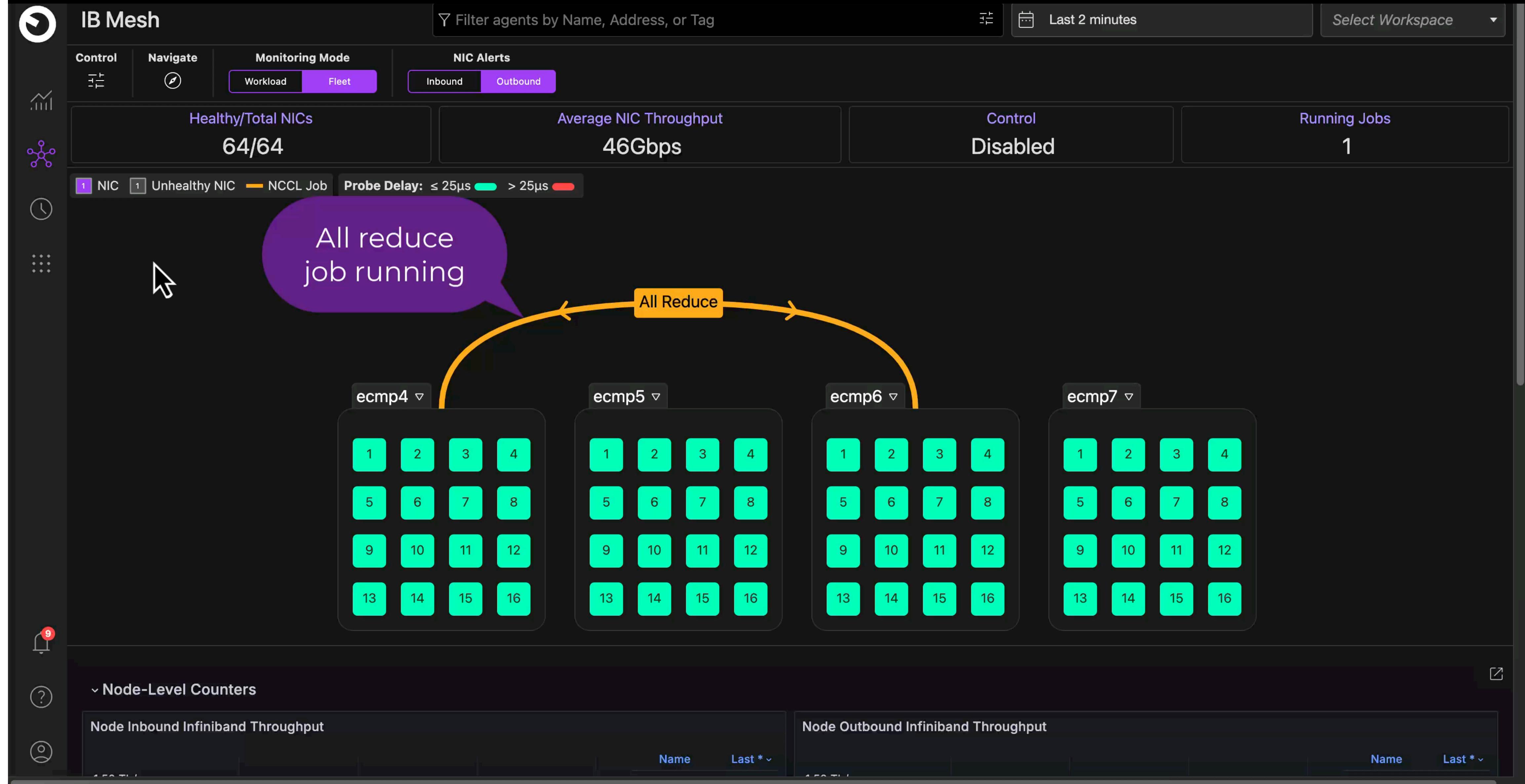
### So what?

- ▶ Causes job to crash.
- ▶ Forces to restart from previous checkpoint.
- ▶ Lost ⏳💰

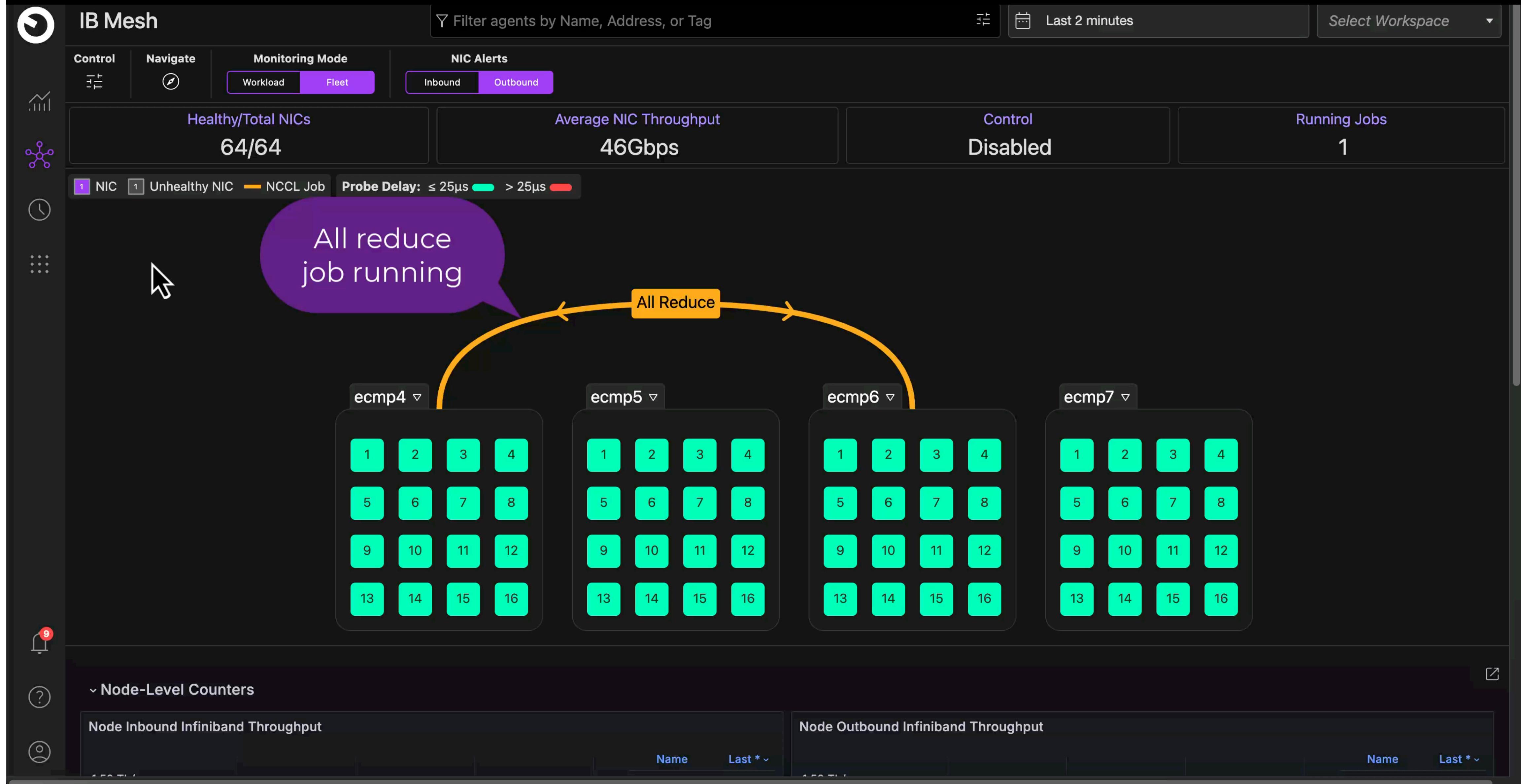
Source: [semanalysis](#)

AI-generated images

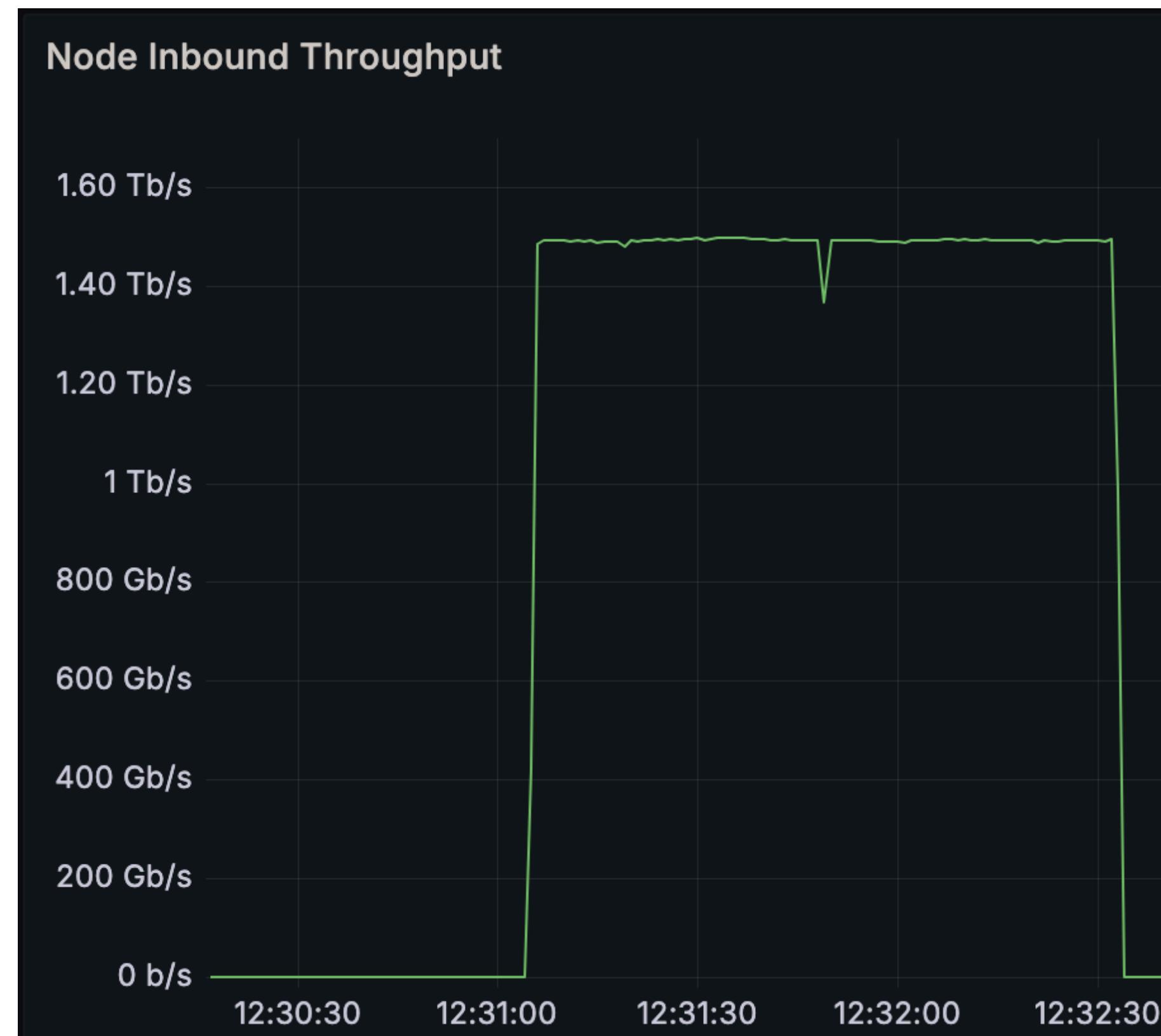
# NIC flapping demo



# NIC flapping demo



# Resilience to NIC Flaps in Oracle Cloud



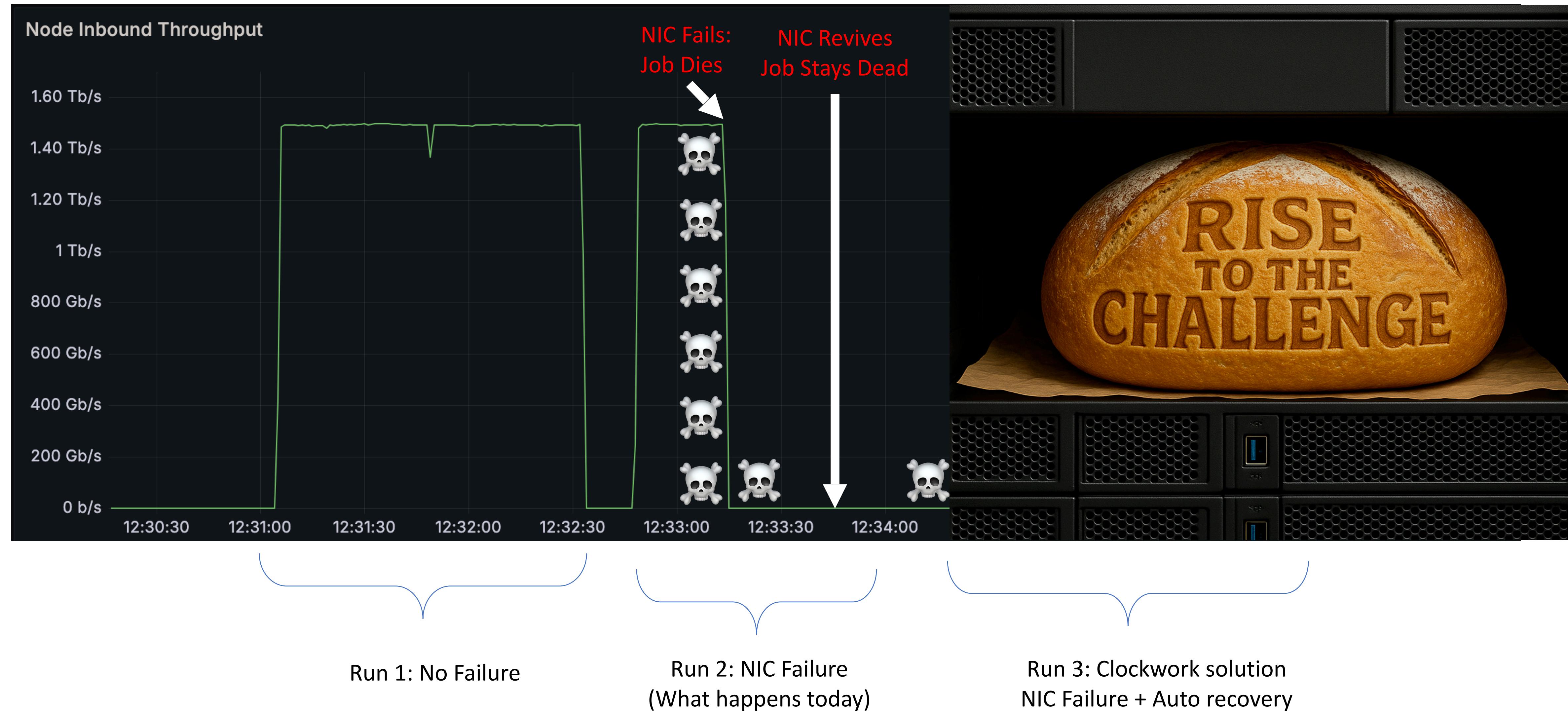
Run 1: No Failure

Run 2: NIC Failure  
(What happens today)

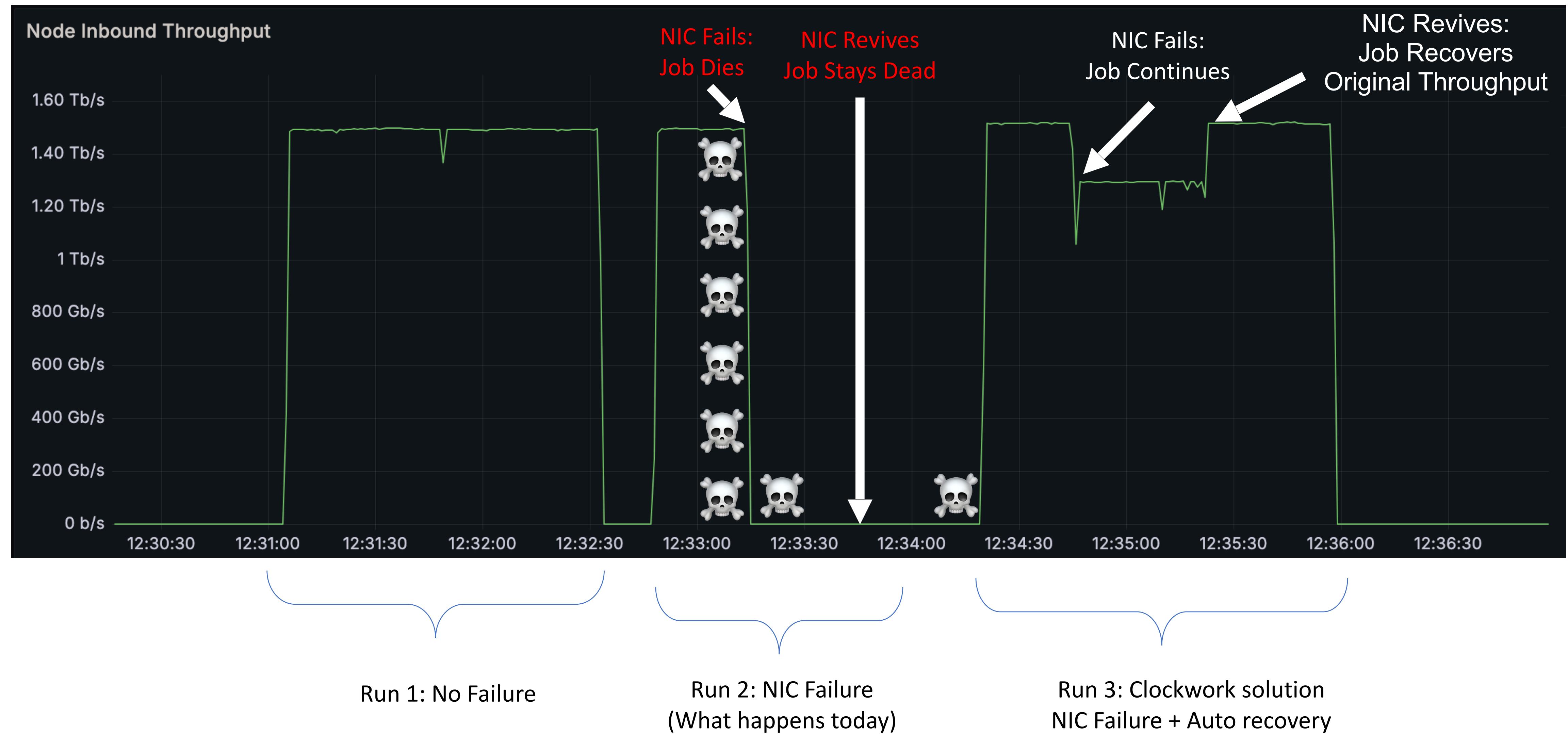
Run 3: Clockwork solution  
NIC Failure + Auto recovery



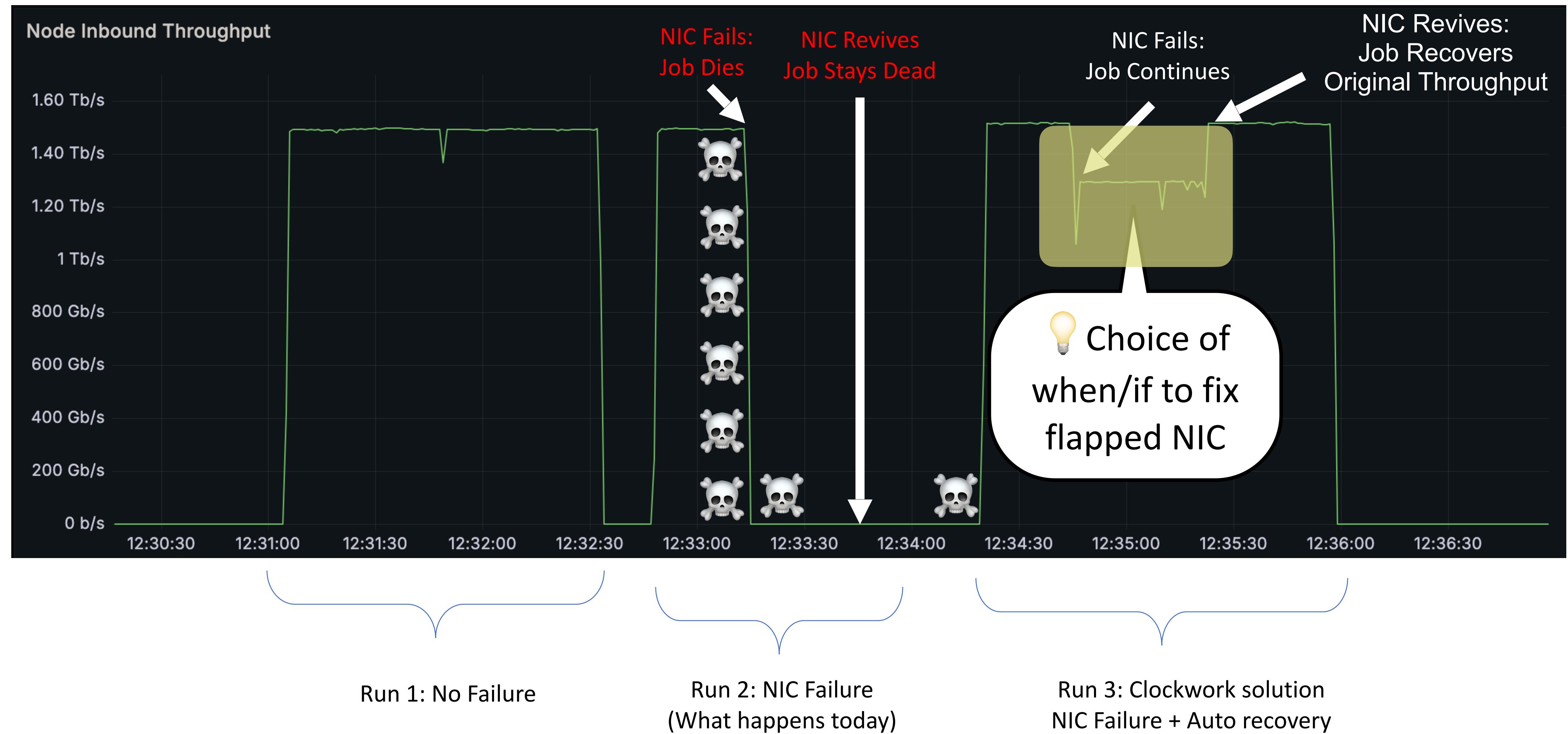
# Resilience to NIC Flaps in Oracle Cloud



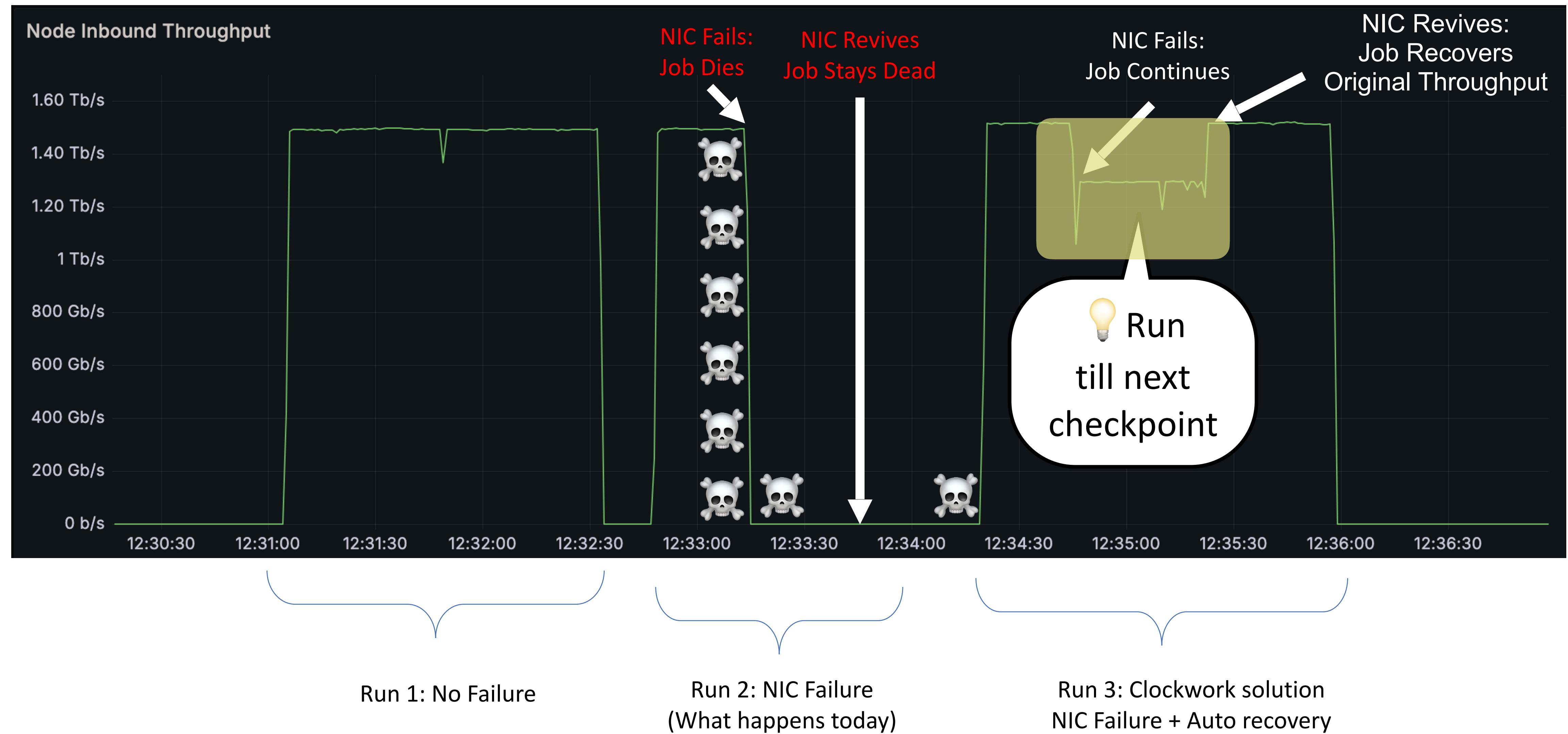
# Resilience to NIC Flaps in Oracle Cloud



# Resilience to NIC Flaps in Oracle Cloud



# Resilience to NIC Flaps in Oracle Cloud

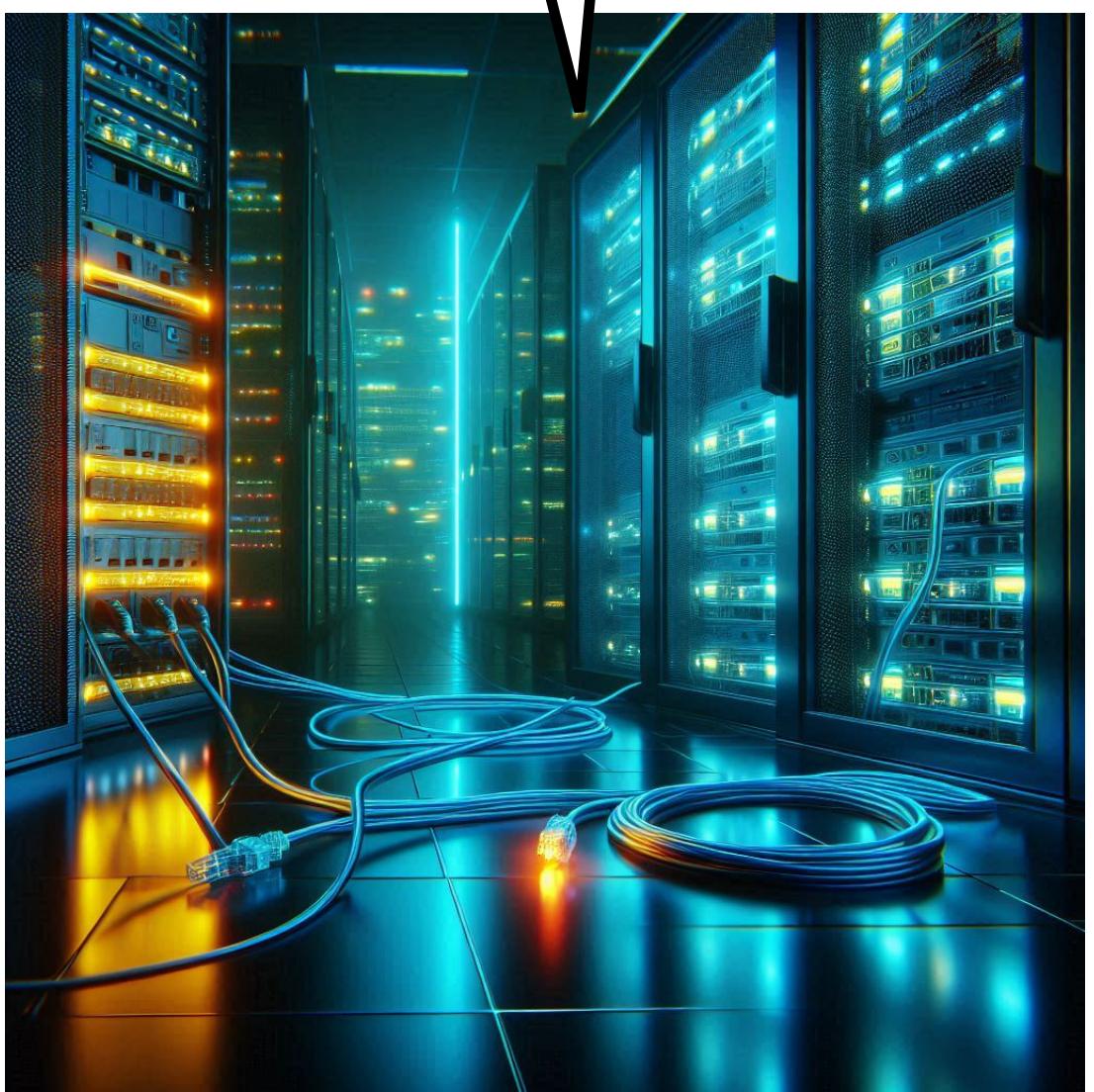


3

## Lack of predictable performance



The app is slow!

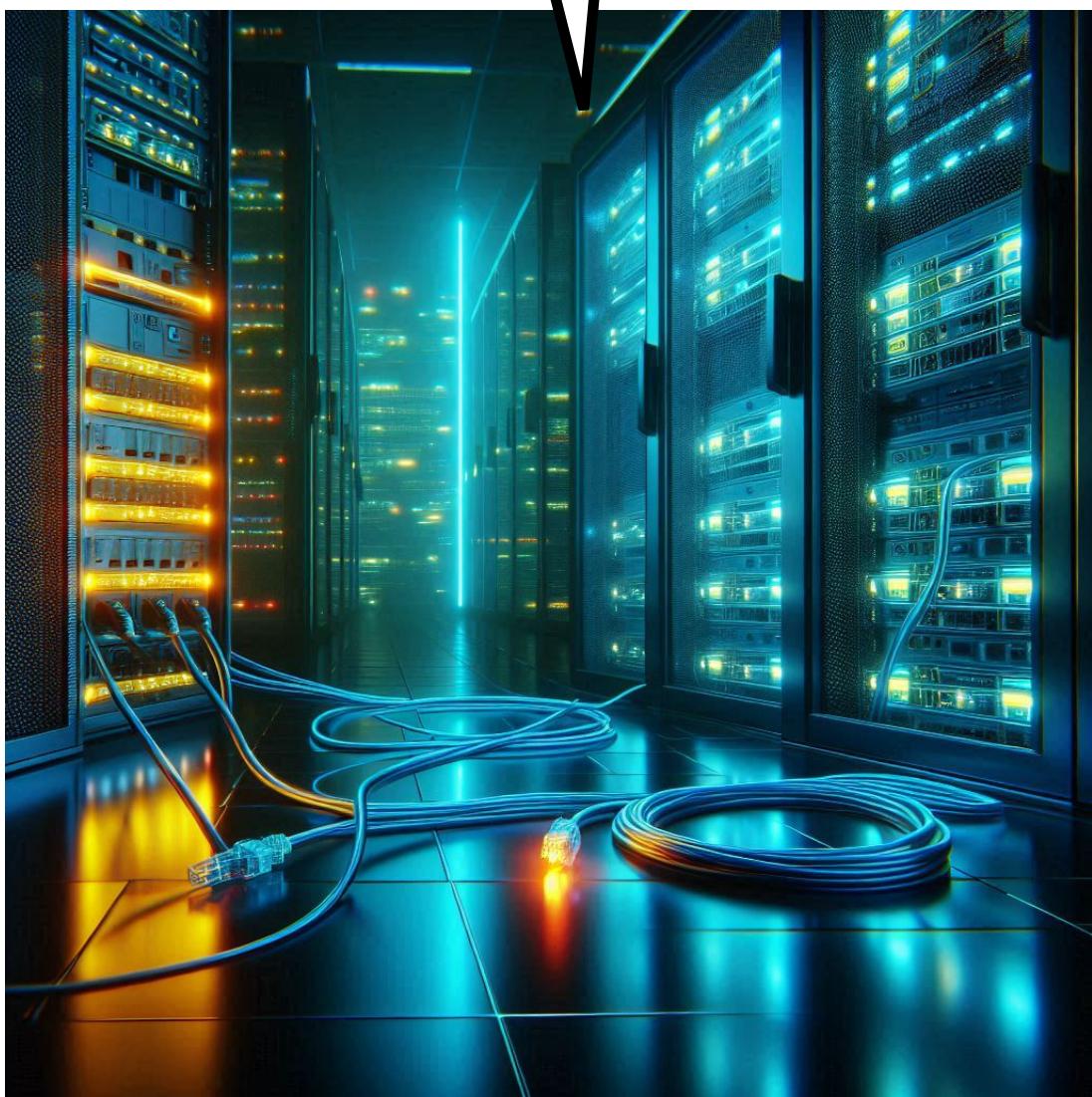


3

## Lack of predictable performance

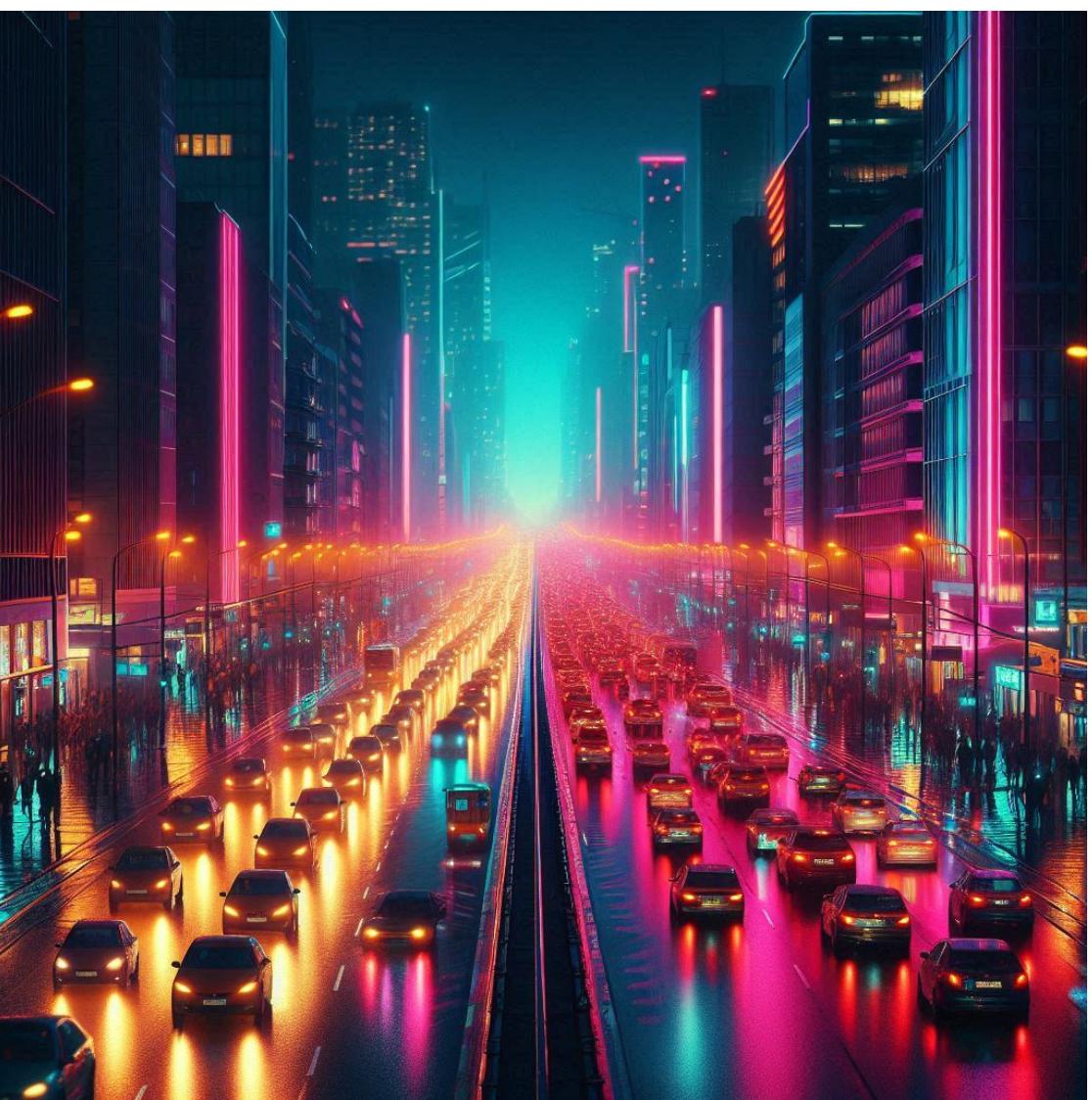


The app is slow!



Why?

Flows contending  
for bandwidth

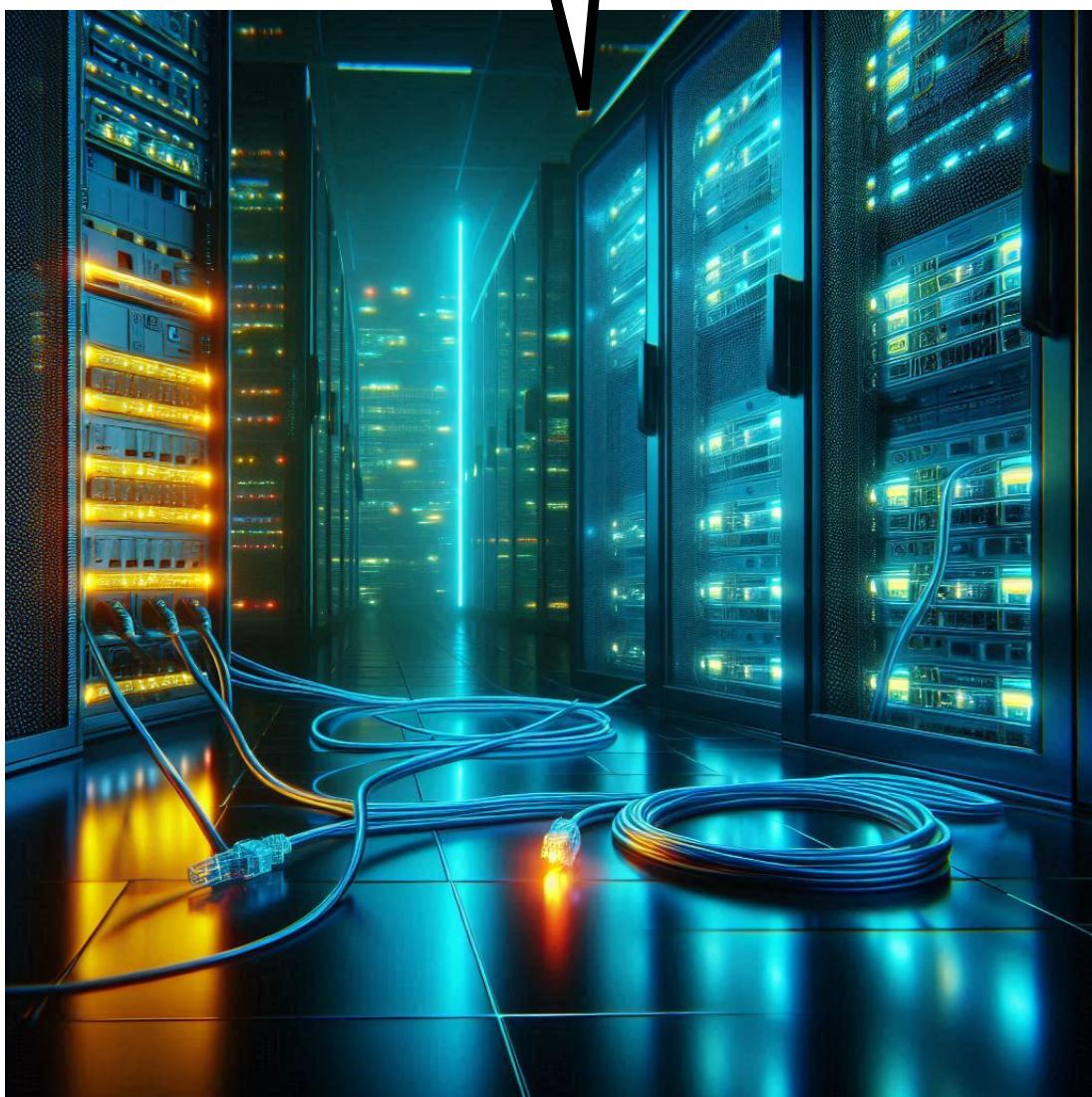


3

## Lack of predictable performance

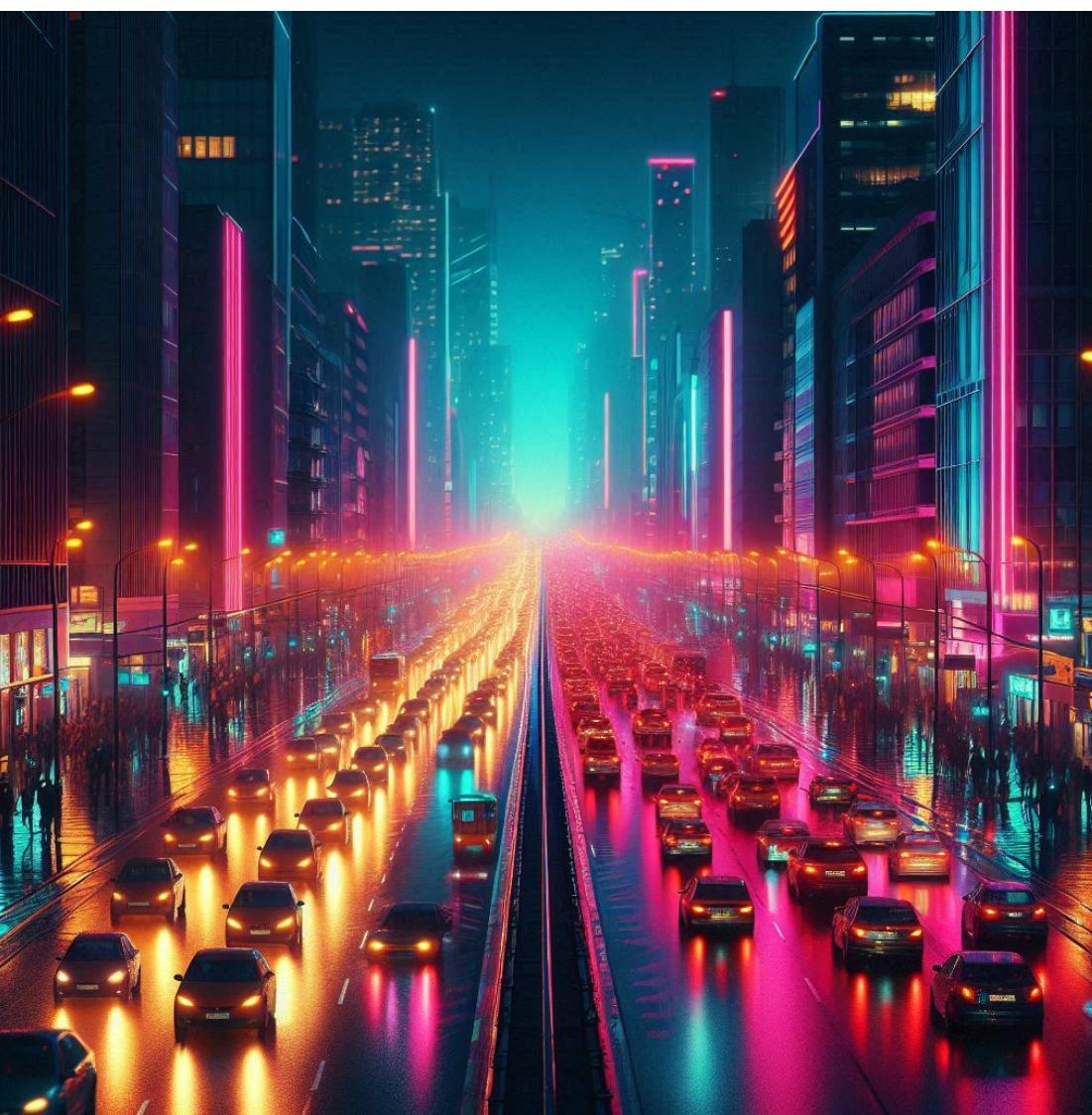


The app is slow!



Why?

Flows contending  
for bandwidth

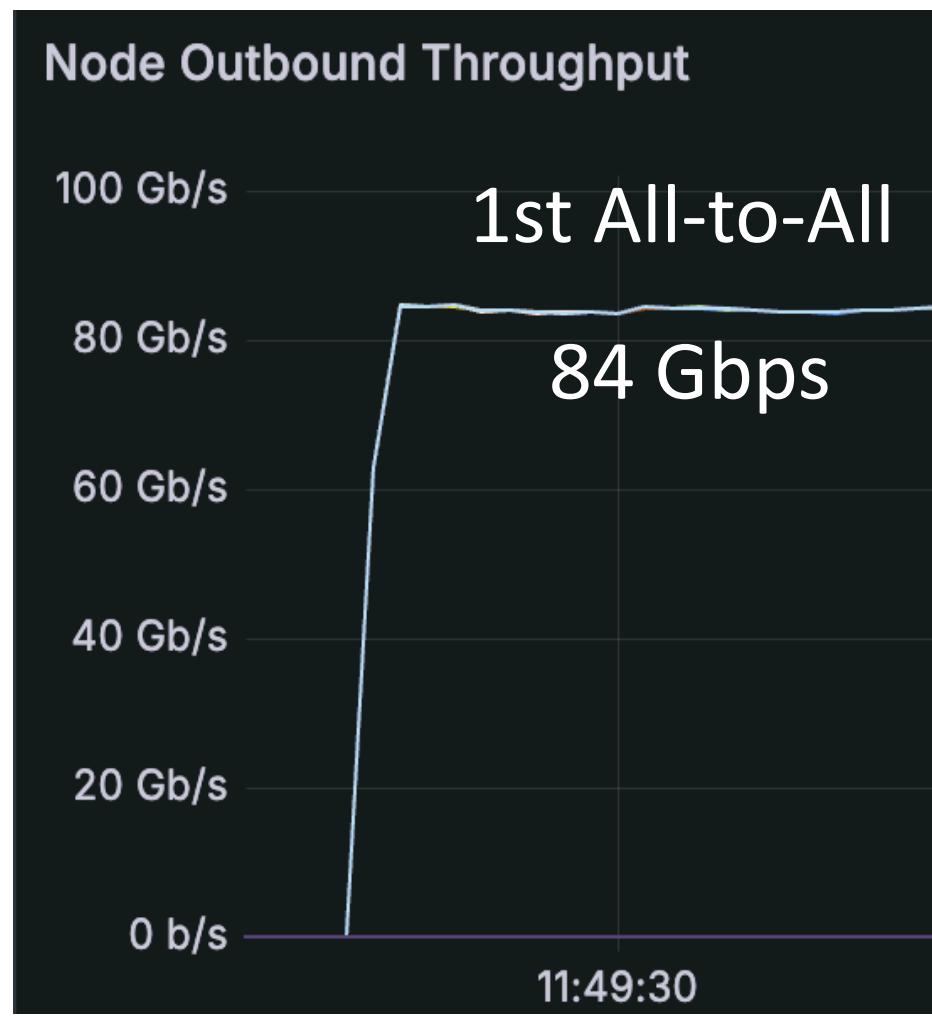


So what?

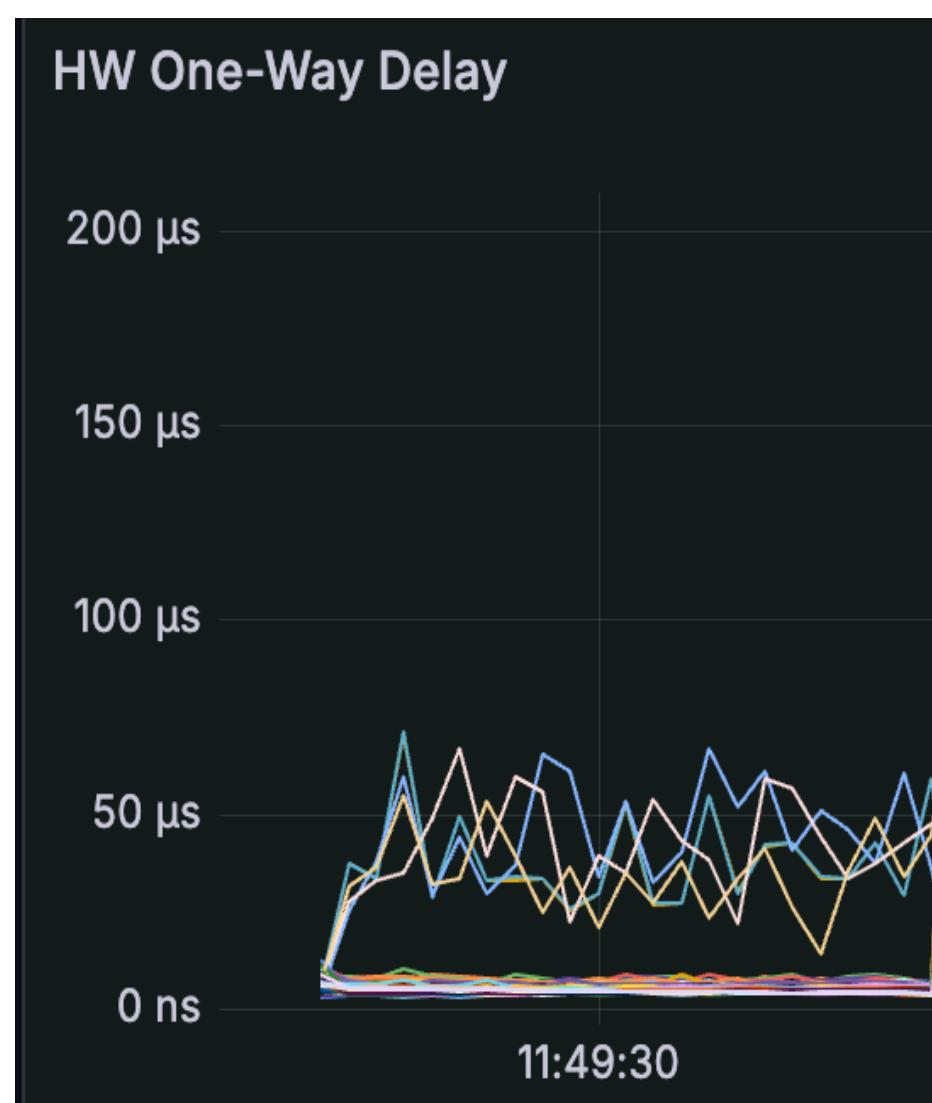
Low throughput,  
high latency

# Detecting and eliminating contention on ECMP Test Bed

Throughput



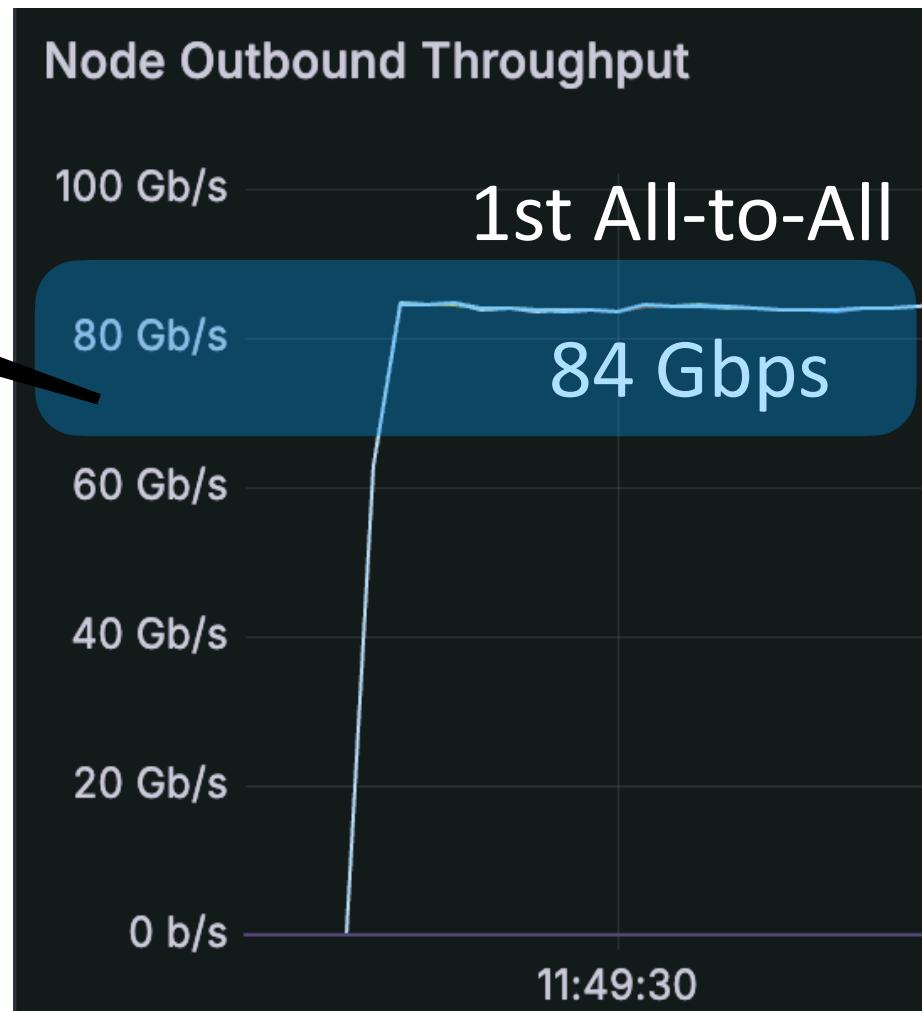
One-way Delays



# Detecting and eliminating contention on ECMP Test Bed

Throughput

⬇️ Lower  
Throughput!



🐢 Queue pairs with  
high one way delays!

One-way Delays



# Detecting and eliminating contention on ECMP Test Bed



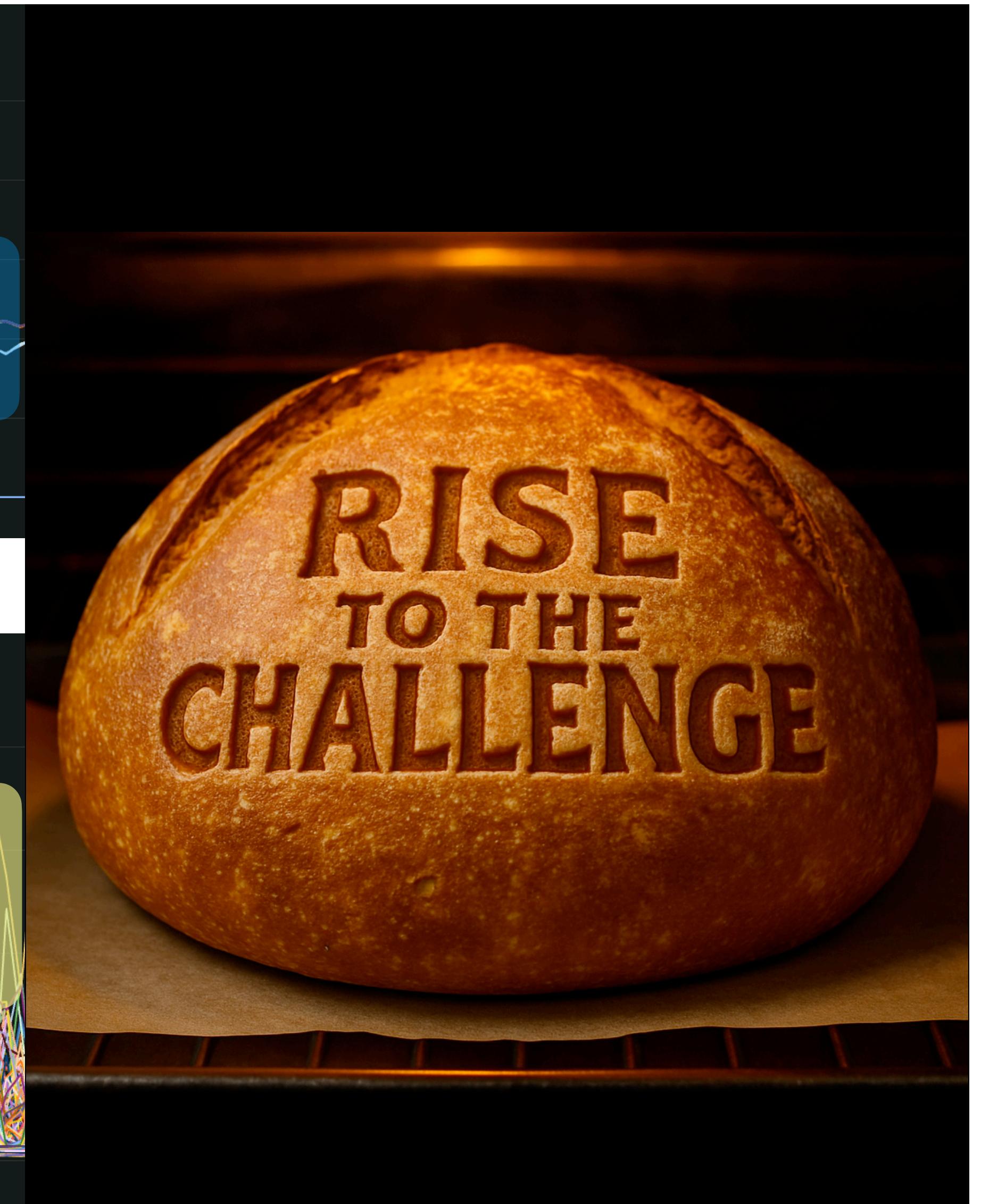
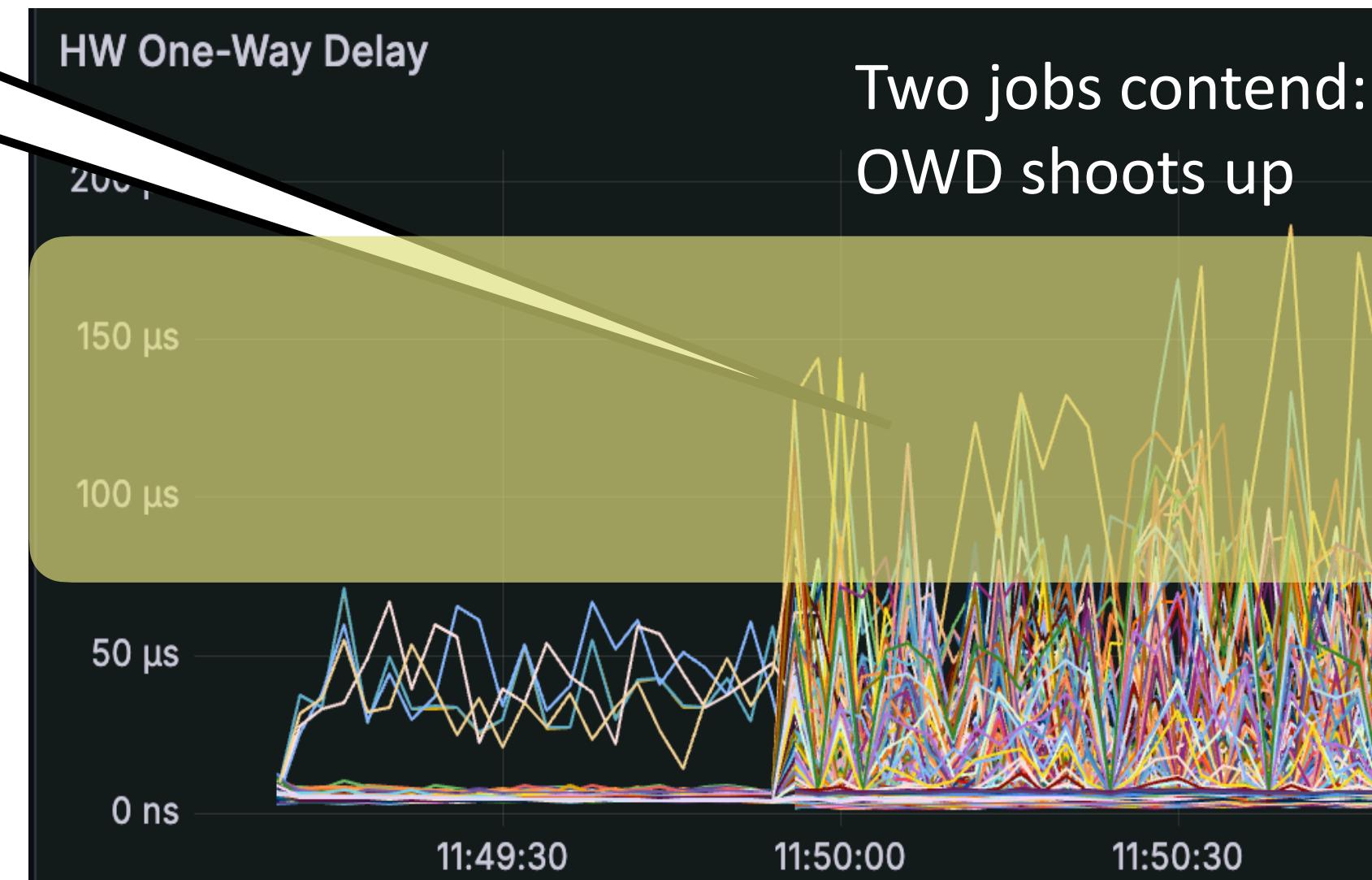
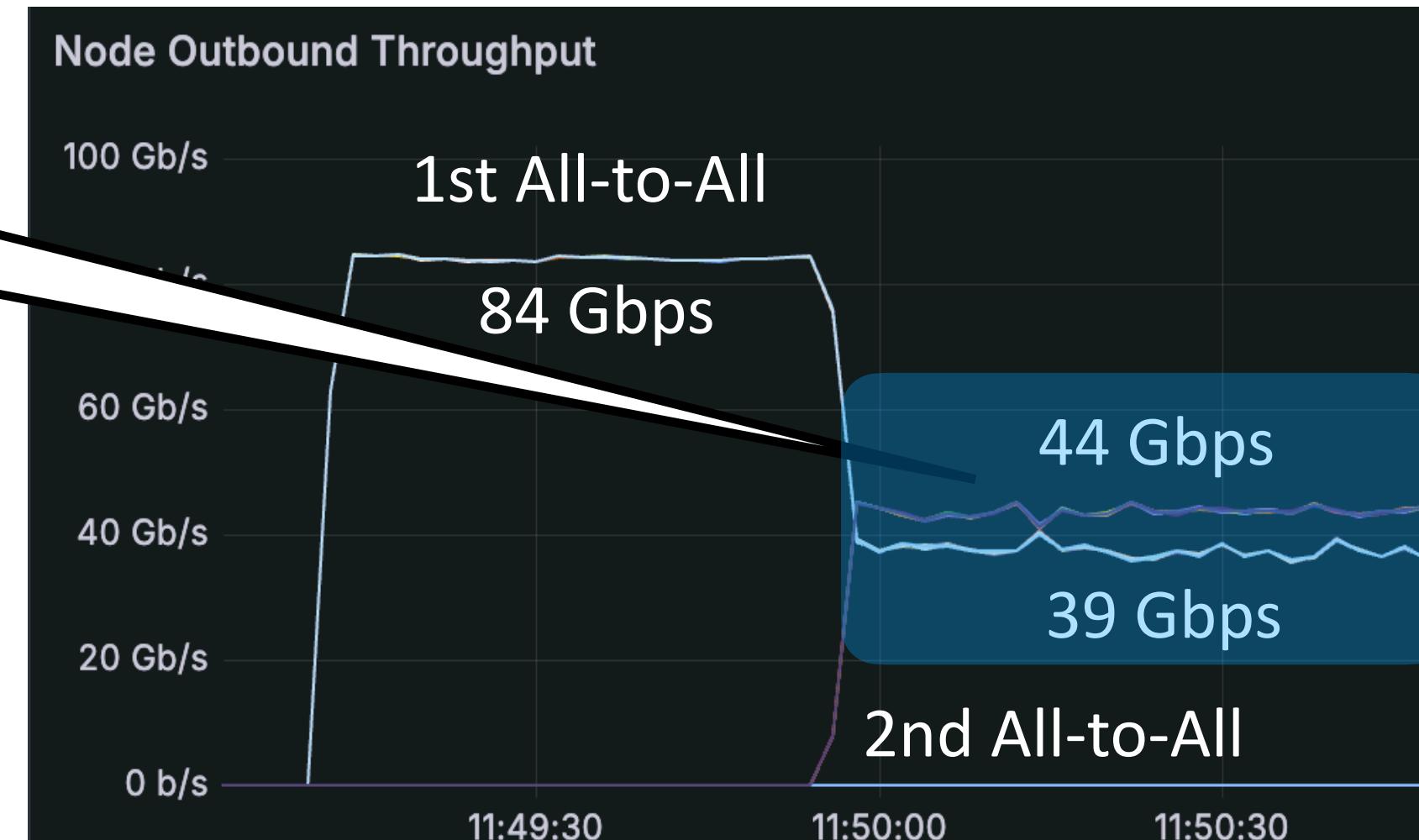
Lower  
Throughput!

Throughput



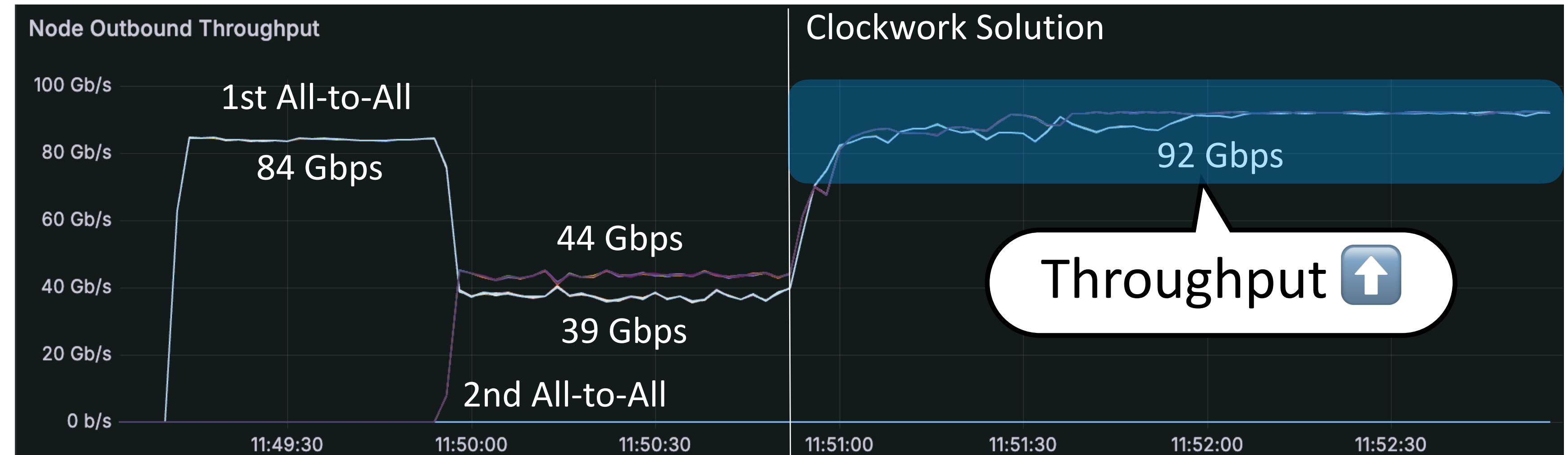
Contention!  
Queue pairs with high  
one way delays

One-way Delays



# Detecting and eliminating contention on ECMP Test Bed

Throughput



One-way Delays



# Key Takeaways

Fine-grained visibility into queue pairs speeds up diagnosis.



Visibility

# Key Takeaways

 Fine-grained visibility into queue pairs speeds up diagnosis.



 Visibility

Checkpoint to recover from crashes due to NIC/link flapping.



 Reliability

# Key Takeaways

 Fine-grained visibility into queue pairs speeds up diagnosis.



 Visibility

Checkpoint to recover from crashes due to NIC/link flapping.



 Reliability

 Minimize delays to slow flows to optimize throughput.



 Performance

# Thank You

Huygens Usenix Paper  
& Slides



[lerna@clockwork.io](mailto:lerna@clockwork.io)



# Thank You

CHECK OUT  
MY FULL  TALK!

Huygens Usenix Paper  
& Slides



lerna@clockwork.io

