# DBs + GPUs:
# Where are we and what's next?

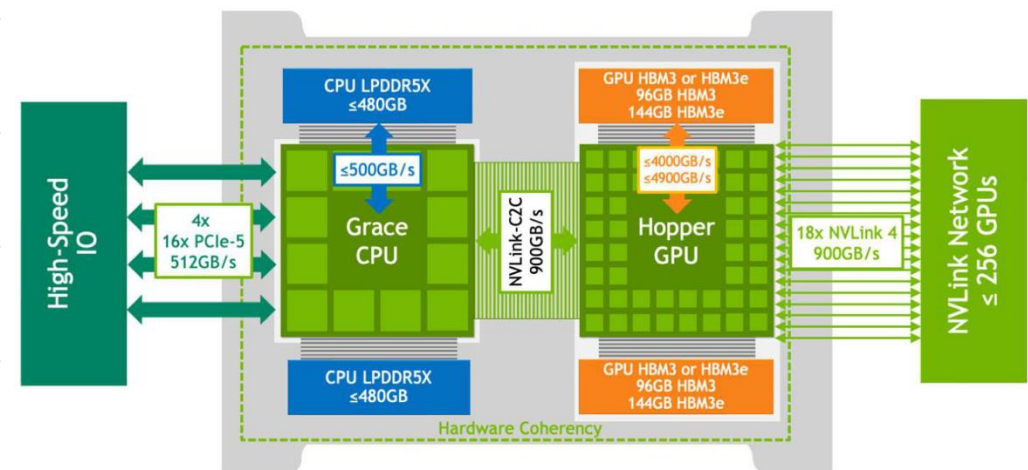Bowen Wu

Systems Group @ ETHz

Lunch Seminar, 31/10/2025

# Why run databases on GPUs?

|  | Active Threads | Memory Bandwidth | Memory Capacity |
|---|---|---|---|
| Nvidia V100 PCIe | 163,840 | 897 GB/s | 16 GB |
| Nvidia A100 SXM | 221,184 | 1,935 GB/s | 80 GB |
| Nvidia H100 PCIe | 233,472 | 2,039 GB/s | 80 GB |
| Nvidia H200 | 233,472$^?$ | 4.8 TB/s | 141 GB |
| AMD  MI300X | ? | 5.3 TB/s | 192 GB |

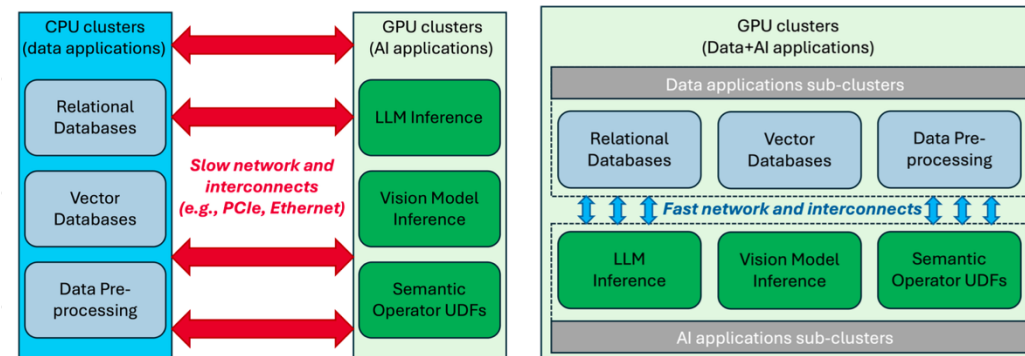Massive *parallelism* and high memory *bandwidth* make GPUs suitable for accelerating databases.



Fast *interconnects* and *networks* make it possible to process *large* datasets at an unprecedented speed.

# Why run databases on GPUs?

**Microsoft inks $33 billion in deals with 'neoclouds' like Nebius, CoreWeave — Nebius deal alone secures 100,000 Nvidia GB300 chips for internal use**

(a) Existing architecture.    (b) Prospective architecture.

Data centers and clouds are being increasingly dominated by GPUs, which cannot be fully utilized by AI.

Making both DBs and AI run on the GPUs enables optimizations to make both systems run more efficiently.

# Current State

- **Research**
  - More than two decades of research
  - Crystal (fully in-GPU)
  - Operator studies
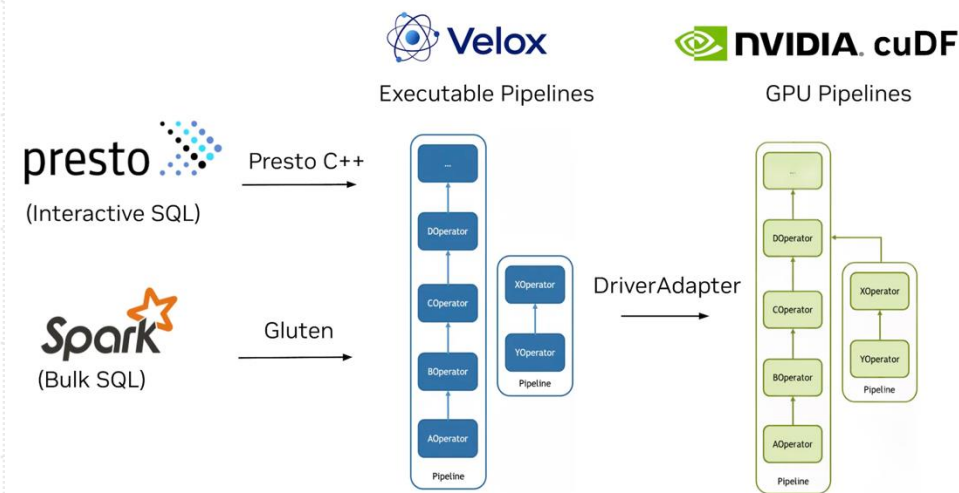  - Maximus/Eiger (SG), Microsoft TQP, SiriusDB

- **Industry**
  - cuDF (NVIDIA), hipDF (AMD)
  - BlazingSQL, Kinetica, HeavyDB, Voltron Data (Startups)
  - Velox + Wave/cuDF (Meta)
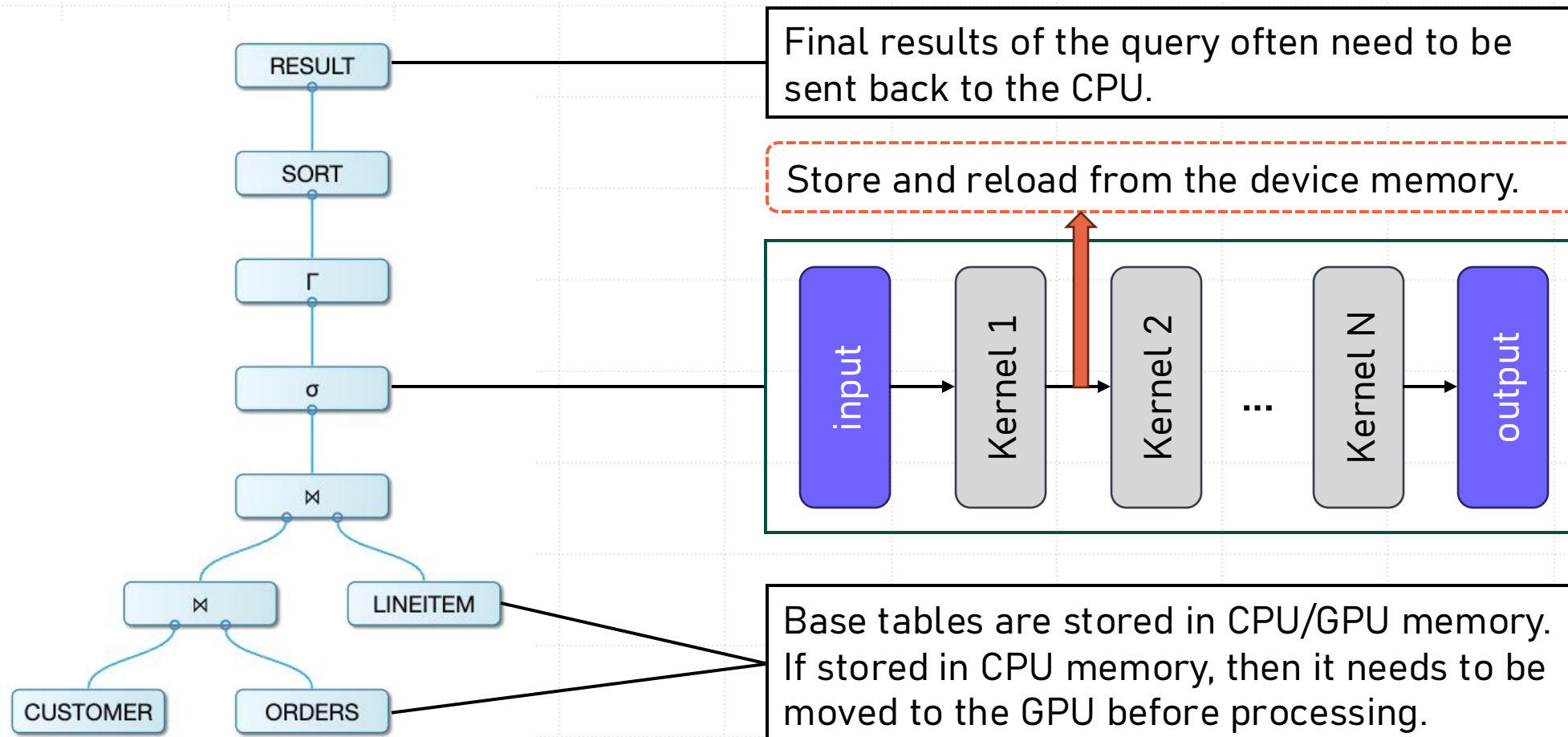  - IBM, Databricks, Microsoft, …



**Velox** + Follow
940 followers
1w · 🌐

Velox now runs at GPU speed 🚀 IBM and NVIDIA have teamed up to bring cuDF-powered GPU execution to Velox delivering big gains for #Presto and Apache Gluten.

# Introduction to GPU-based DB

RESULT ── Final results of the query often need to be sent back to the CPU.

SORT

Γ

σ ────────────────→ input → Kernel 1 → Kernel 2 ... Kernel N → output

Store and reload from the device memory.

⋈

⋈          LINEITEM ── Base tables are stored in CPU/GPU memory. If stored in CPU memory, then it needs to be moved to the GPU before processing.

CUSTOMER    ORDERS

# Is the query exec bound by *GPU* or *IC*?

- **MaxBench** – Single-GPU DB Benchmark

- Authors: Marko Kabić, Bowen Wu, Jonas Dann, Gustavo Alonso

**Table 1: Hardware Configurations for the empirical evaluation.**

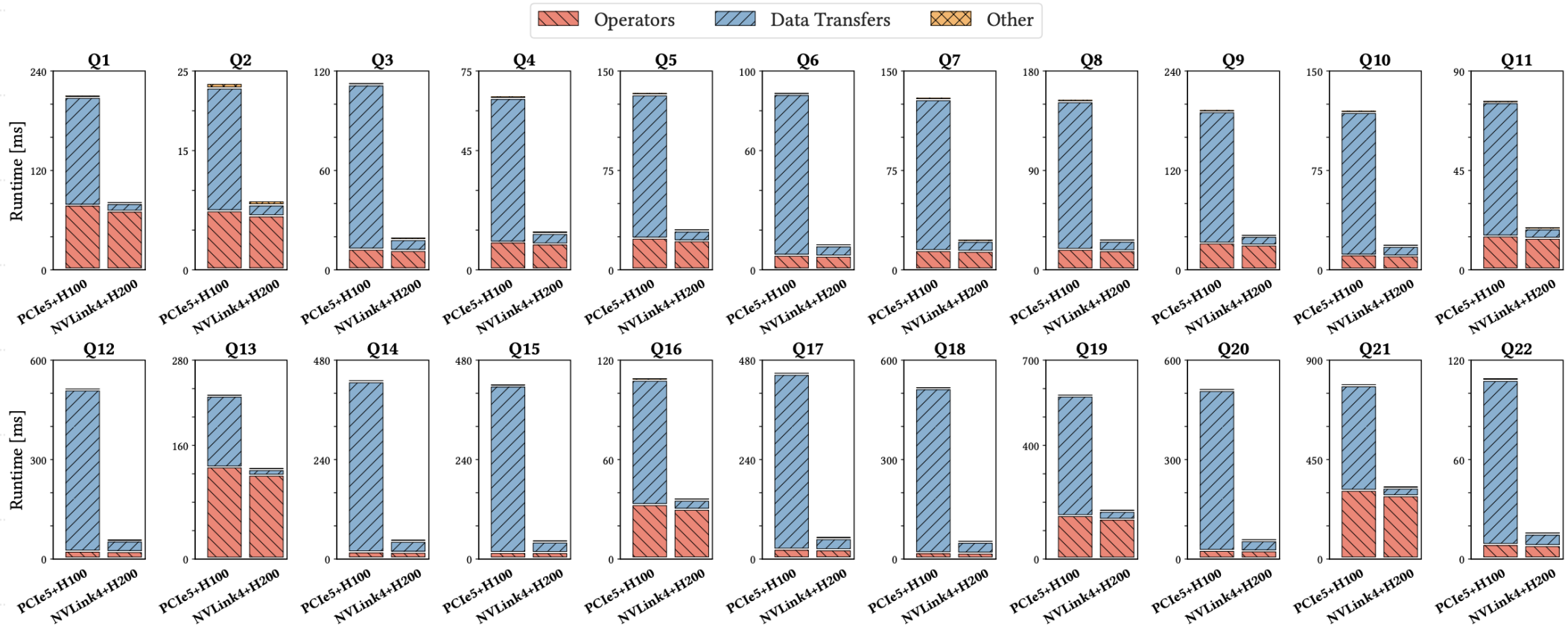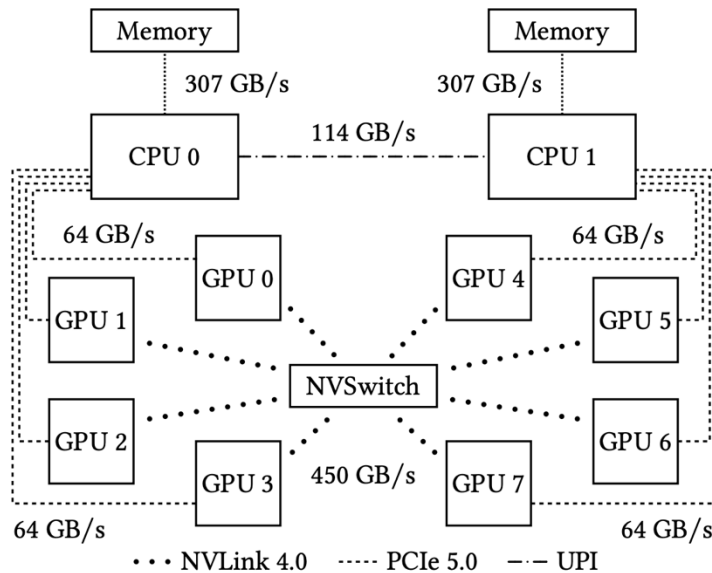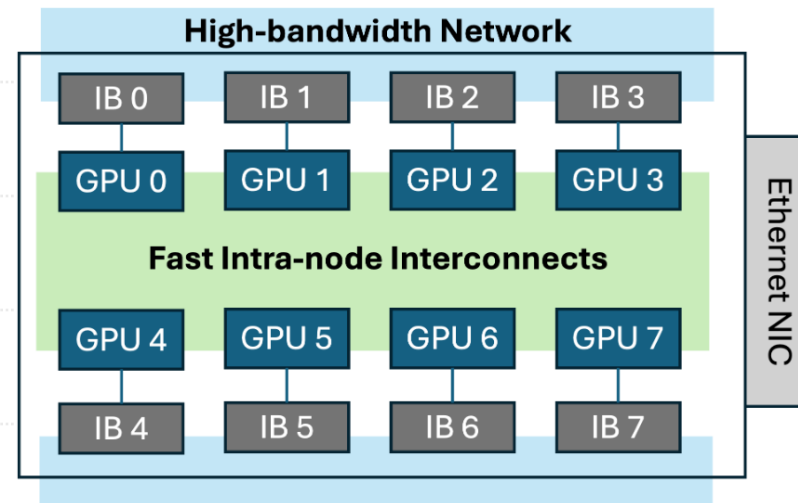| Configuration | $C_1$ = PCIe3+A100 | $C_2$ = PCIe5+H100 | $C_3$ = GH200 (NVLink4+H200) | $C_4$ = PCIe4+RTX3090 |
|---|---|---|---|---|
| CPU | Intel Xeon Platinum 8171M | AMD EPYC 9124 | NVIDIA Grace | AMD EPYC 7313 |
| CPU cores | 2x6 | 2x16 | 72 ARM Neoverse V2 cores | 2x16 |
| GPU | NVIDIA A100 40GB | NVIDIA H100 80GB | NVIDIA H200 96GB | NVIDIA RTX3090 24GB |
| GPU Mem. Bandwidth | 1.55 TB/s | 4 TB/s | 4 TB/s | 0.936 TB/s |
| GPU clock (base-boost) | 765–1410 MHz | 1080–1785 Mhz | 1980–1980 Mhz | 1395-2100 MHz |
| Power cap | 400W | 400W | 624.15W / 900W | 350W |
| Interconnect (IC) | PCIe 3.0 | PCIe 5.0 | NVLink 4.0 | PCIe 4.0 |
| IC Bandwidth (1-way) | 16 GB/s | 64 GB/s | 450 GB/s | 32 GB/s |

# MaxBench – Single–GPU DB Benchmark



**Figure 3: The full TPC-H benchmark (SF=10) run on hardware configurations $C_2$ = PCIe5 + H100 and $C_3$ = NVLink4+H200.**

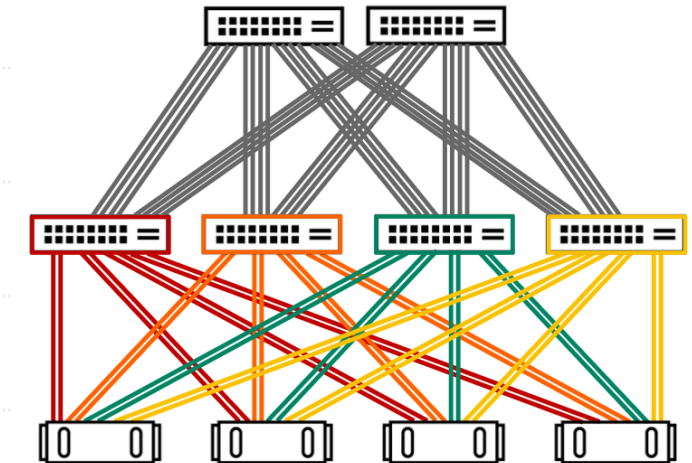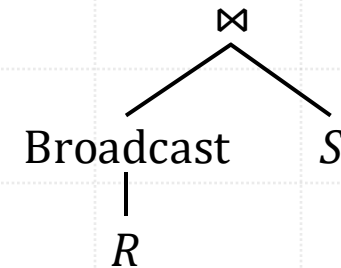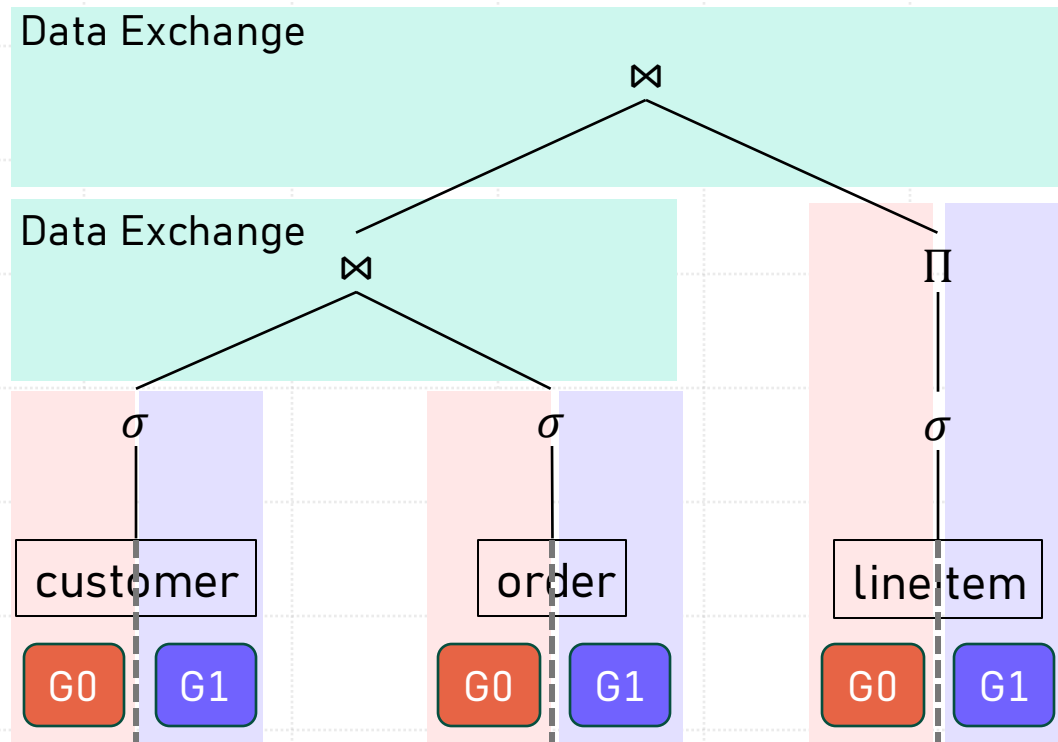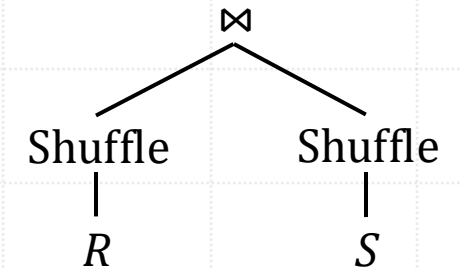# Distributed GPU cluster



8x H100 SXM5 80 GB[1]



GPUs+NICs[2]



Switches[3]

# Distributed DB



Data Exchange

Data Exchange

$\bowtie$

$\bowtie$

$\Pi$

$\sigma$

$\sigma$

$\sigma$

customer

order

lineitem

G0  G1

G0  G1

G0  G1

$\bowtie$

Broadcast        $S$

$R$

$\bowtie$
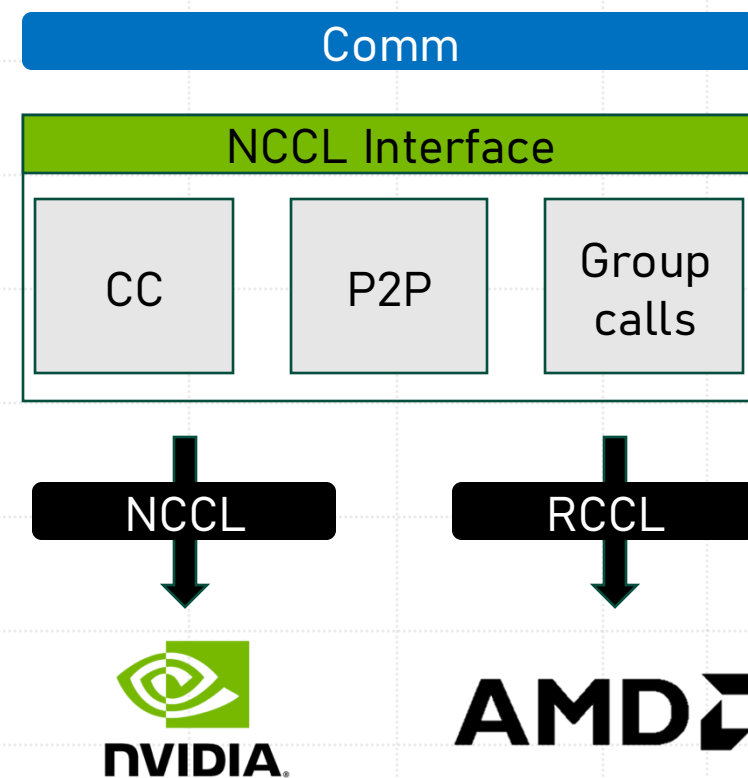
Shuffle        Shuffle
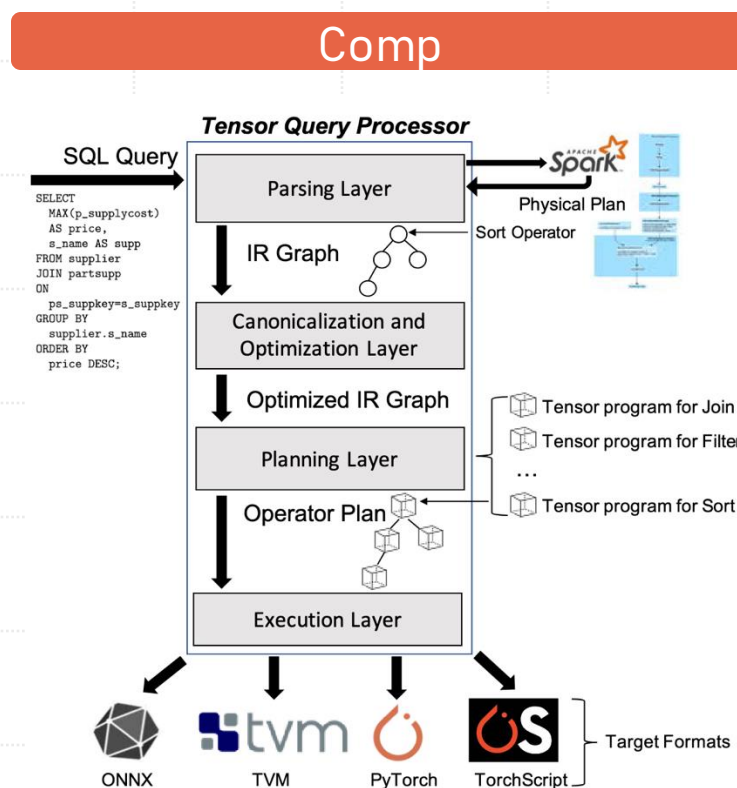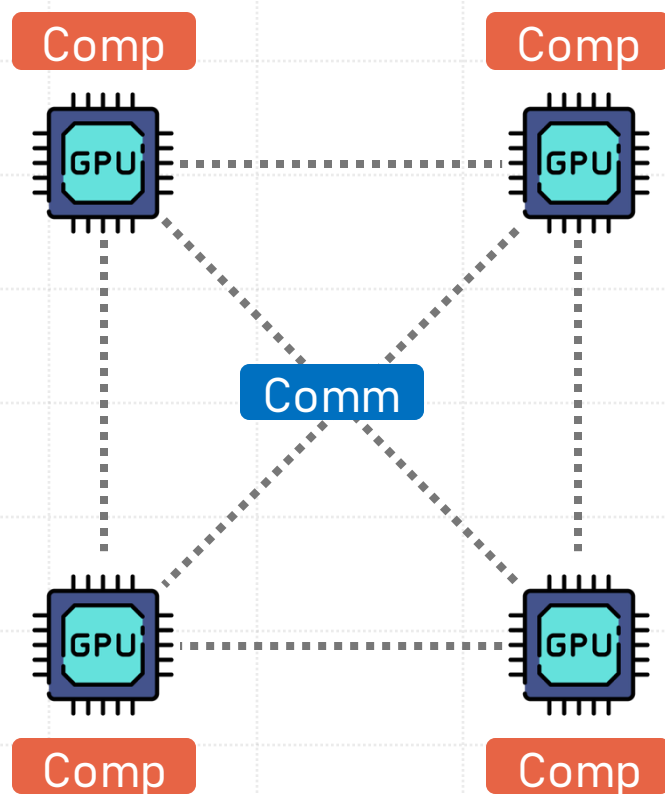
$R$        $S$

**Broadcast**

- Works well when the data is small.
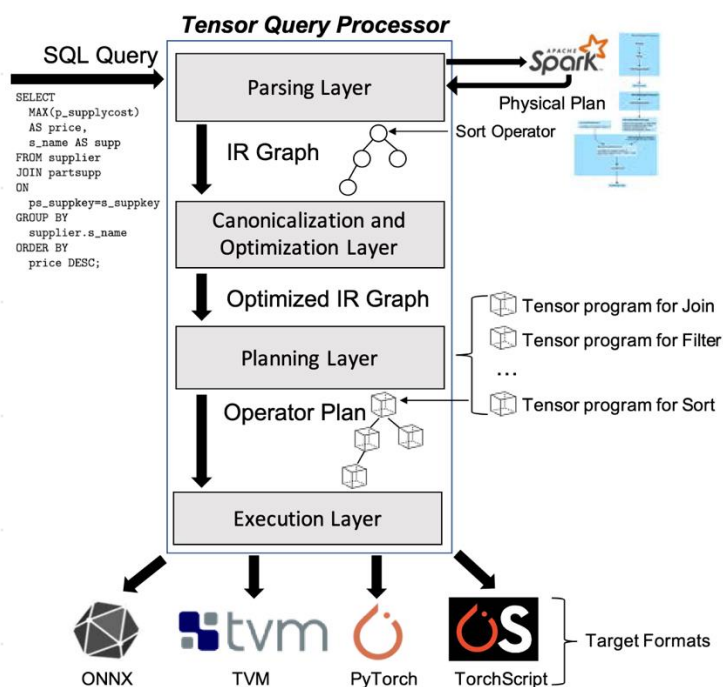- No extra overhead.
- Resistant to skew (later)

**Shuffle**

- Works well when the data is large.
- Extra overhead of partitioning data.
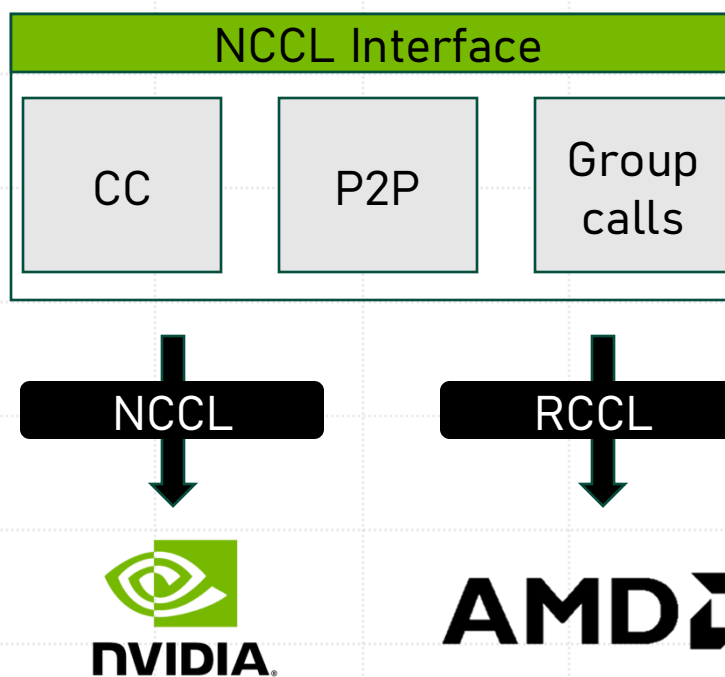- Not very resistant to skew (later)

# Our System

# Our System



**Comp**

**Comm**

NCCL Interface

| CC | P2P | Group calls |

NCCL → NVIDIA

RCCL → AMD

Using existing ML libraries gives
- Good out-of-box performance
- Portability
- Easy-to-develop

But it may not fully suit the DB workload, which often exhibits lots of *irregularity*.

***Research question:***
How efficiently can we build a DB with GPU acceleration using ML libraries?

# Experiment Setup

- We report the total time of all 22 TPC–H queries at SF=1000 or 3000.

- Data are already partitioned and loaded into the GPUs' memory.

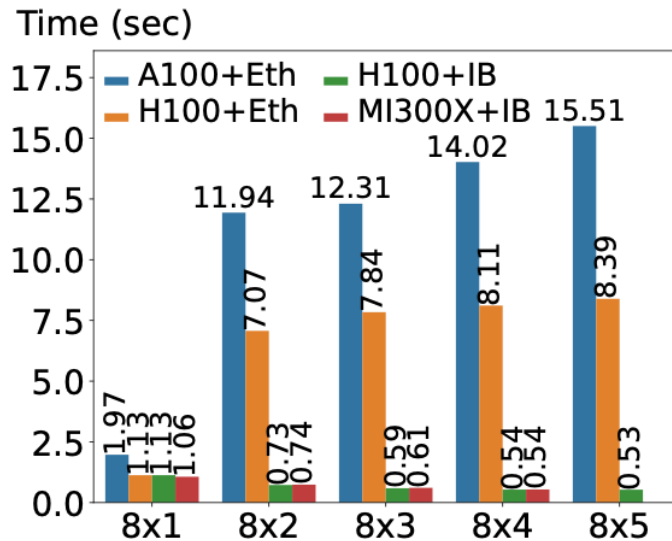Table 3: Cluster Configurations. Eth: Ethernet. IB: InfiniBand.

| Cluster | GPU Type | HBM (GiB) | k | V | GPU Interconnect | | CPU type | CPU Cores | CPU Mem (GiB) | Price/hour (USD) |
| | | | | | Intra-VM | Inter-VM | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | NVIDIA A100 | 80 | | 7 | NVLink 300 GB/sec | Eth: 1x50 Gbits/sec | AMD EPYC 7V12 | | 1800 | 32.77* |
| 2 | NVIDIA H100 | 79.6 | 8 | 5 | NVLink 450 GB/sec | Eth: 1x100 Gbits/sec / IB: 8x400 Gbits/sec | Intel Xeon Platinum 8480C | 96 | 1900 | 98.32 |
| 3 | AMD MI300X | 191.5 | | 4 | Infinity Fabric 448 GB/sec | Eth: 1x100 Gbits/sec / IB: 8x400 Gbits/sec | Intel Xeon Platinum 8480C | | 1850 | 63.6 |

* This is the price for 8x200 Gbits/sec Infiniband. The Eth version, where we run our experiment, is not publicly listed.
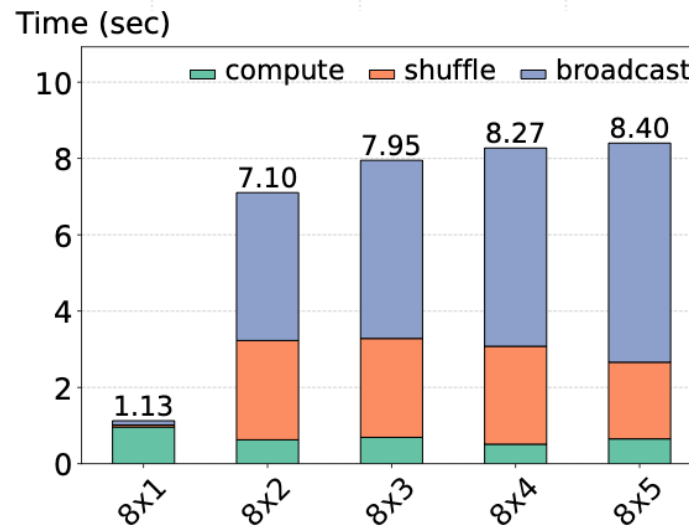
# Key Results

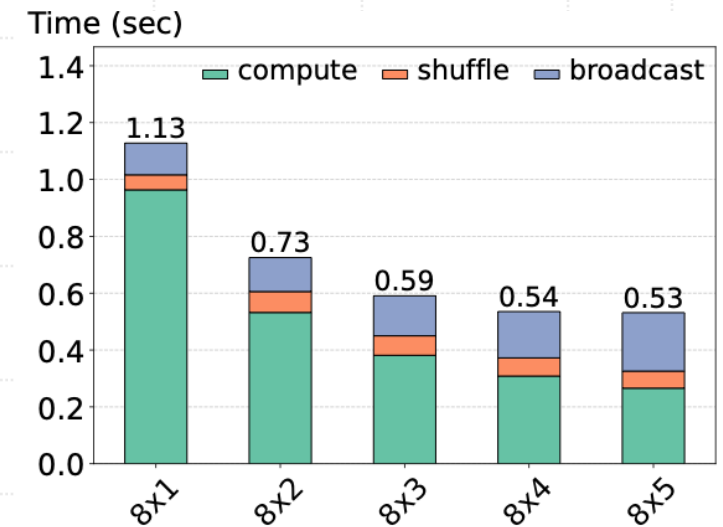*How will the performance change if my cluster changes or my workload changes?*
- Scale-out (more machines)?
- Scale-up (more GPUs/machine, faster NVLink and network)?
- Workload sensitivity (skew, data placement, etc.)

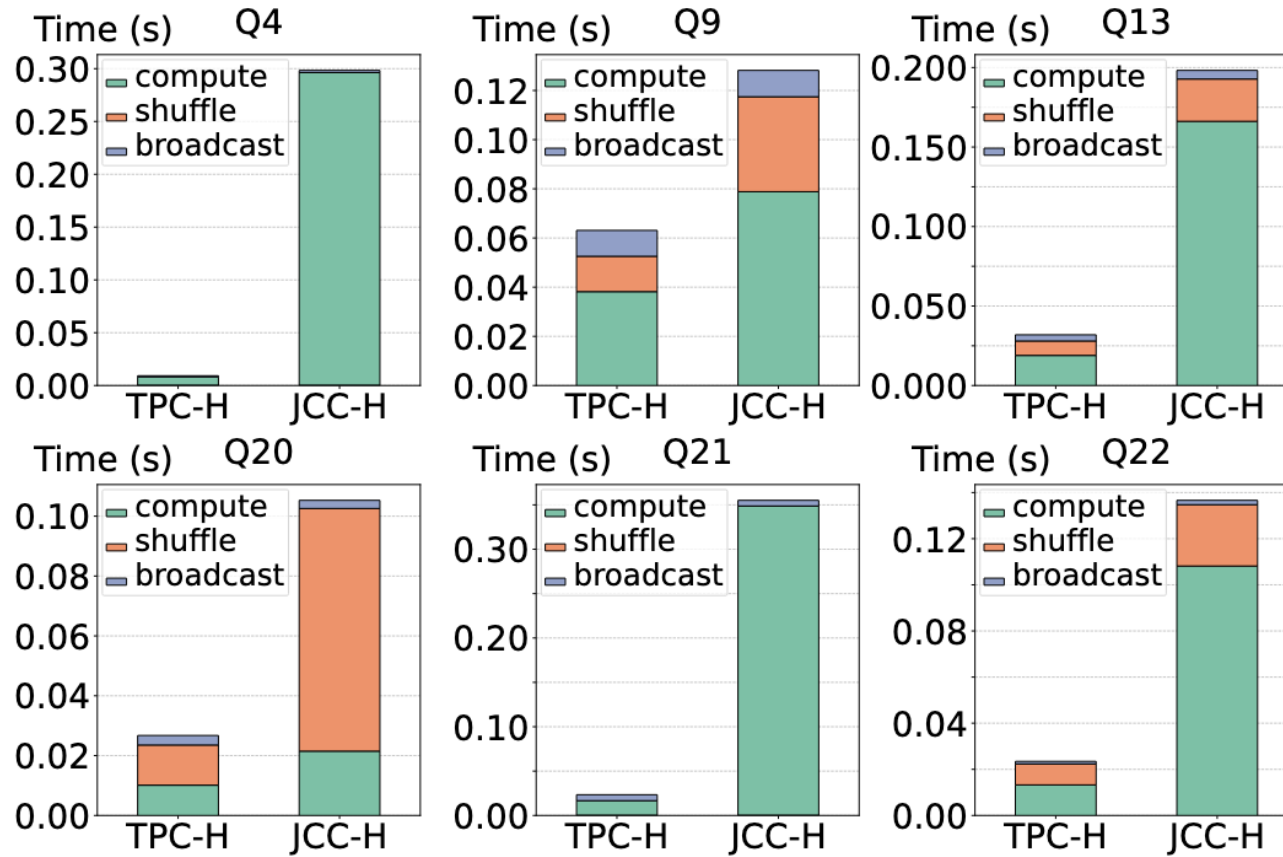- GPU-based DBs can be 1-2 orders of magnitude faster than CPU-based ones.



- A fast network is necessary for good scalability.



- With a fast network, the query exec transitions from GPU-bound to network-bound as the cluster grows.

# Effect of *data skew*



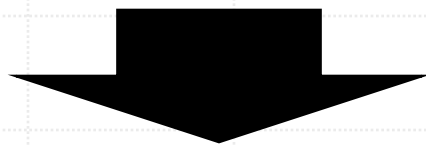Figure 21: Time breakdown comparison. (V=5)

JCC-H 1TB
Same schema + same query + skewed data

# Key Contributions

**How will the performance change if my cluster changes or my workload changes?**
- Scale-out (more machines)?
- Scale-up (more GPUs/machine, faster NVLink and network)?
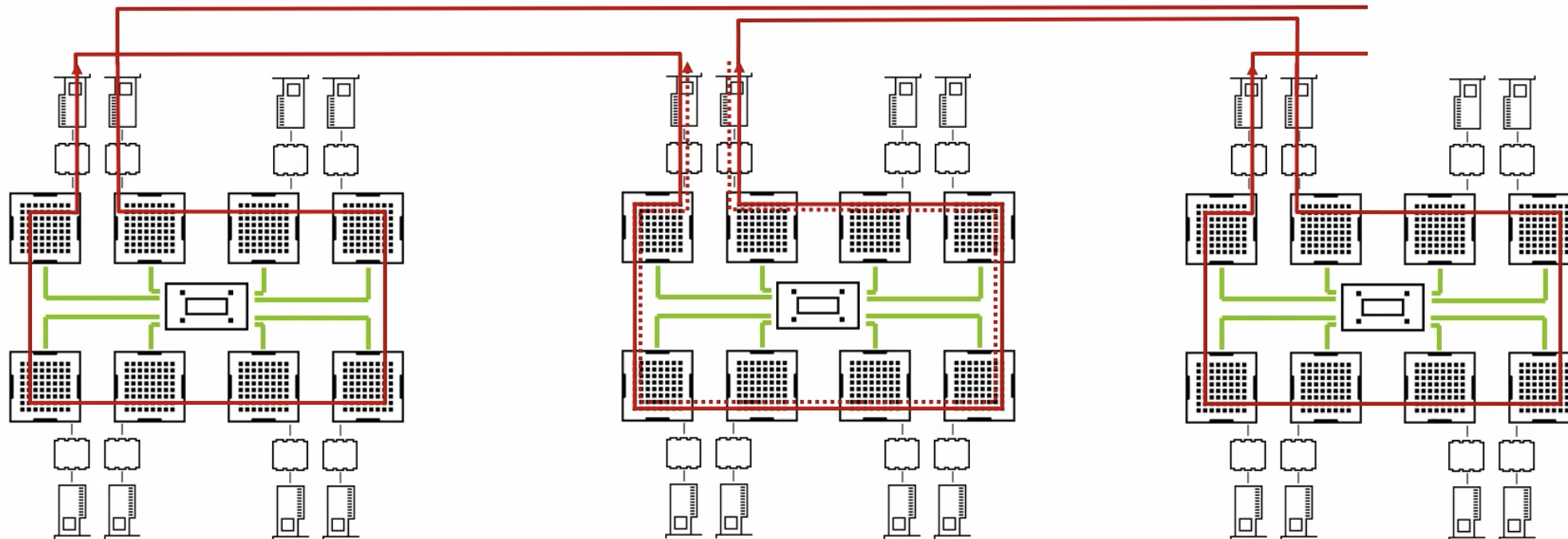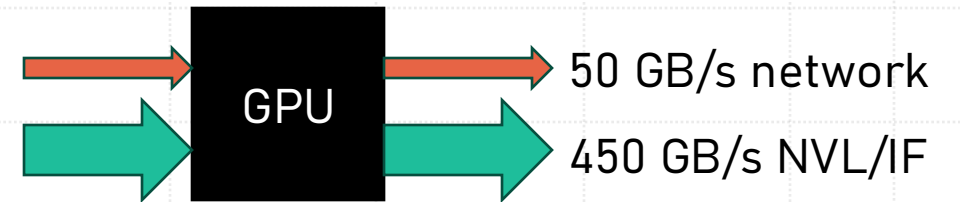- Workload sensitivity (skew, data placement, etc.)

**Model the computation and communication of query execution**

**Challenges**:
- The distributed GPU cluster has a very heterogeneous interconnect.
- The algorithm used by NCCL is obscure.
- Previous effort main focused on modeling ML workloads.
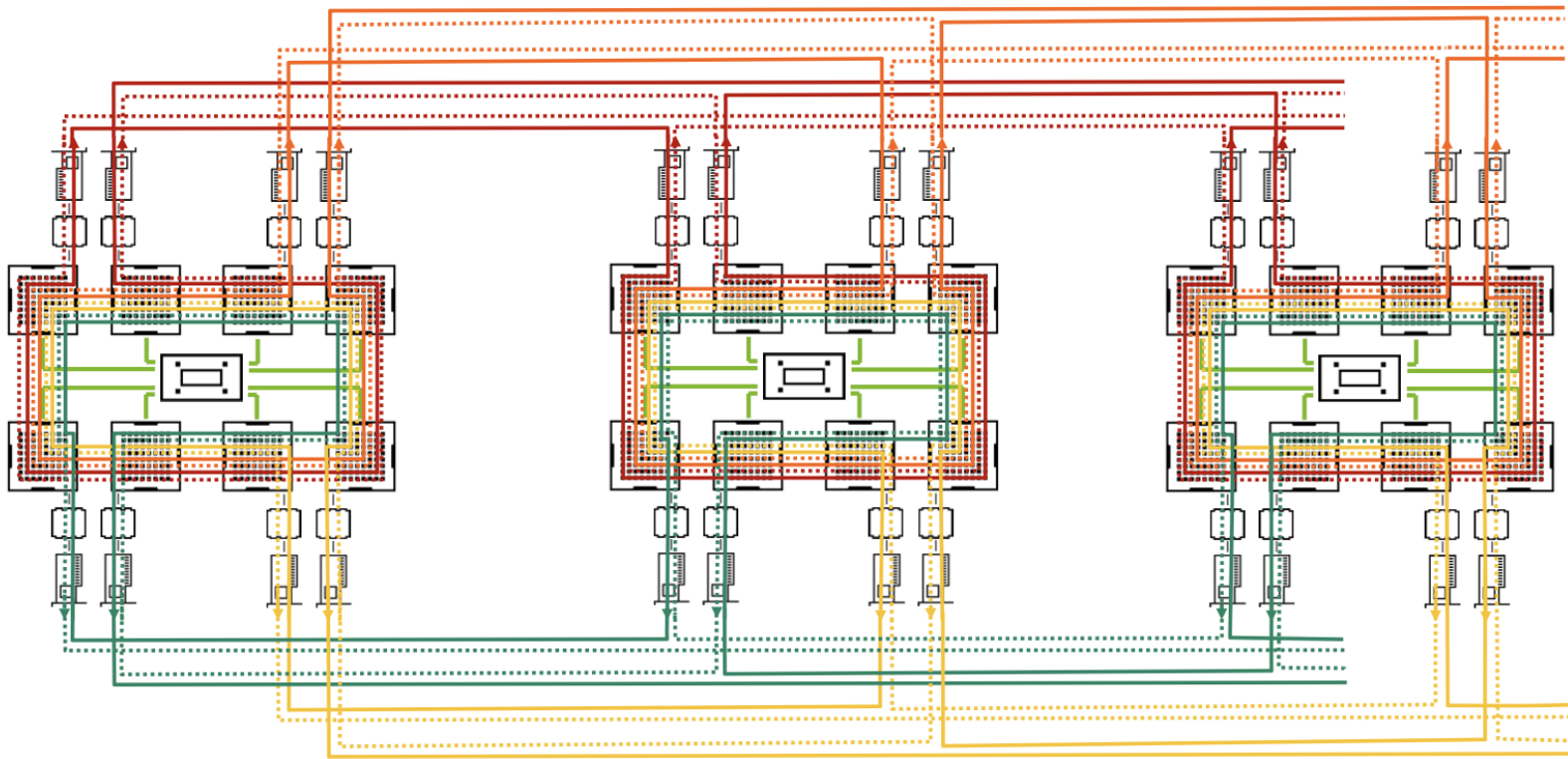
# Modeling Broadcast



GPU

50 GB/s network

450 GB/s NVL/IF

Formation of a single ring

# Modeling Broadcast

Formation of all rings

| k | Number of GPUs per machine (or node, VM) |
|---|---|
| V | Number of machines (or nodes, VMs) |
| N | Total number of GPUs = $k \times V$ |
| $B_g$ | Unidirectional inter-GPU bandwidth within each machine |
| $B_n$ | Unidirectional network bandwidth at each machine |
| S | Total dataset size processed by the GPU cluster. |
| $G_{ij}$ | The $i$-th GPU in the $j$-th machine. |
| $m_{ij \to pq}$ | The message sent from $G_{ij}$ to $G_{pq}$. |

# Modeling Broadcast



$\frac{S/N}{\min(B_n, B_g)}$ time/step

In total, you need $(N-1)$ steps.

➔ Total time = $(N-1) \frac{S/N}{\min(B_n, B_g)}$

➔ Throughput = $\frac{N}{N-1} \min(B_n, B_g)$

Takeaways:

- Throughput **_decreases_** with #GPUs.

- Faster network won't improve the performance.

# Modeling Broadcast



(a) Broadcast (H100+Eth).

(c) Broadcast (H100+IB).

(e) Broadcast (MI300X+IB).

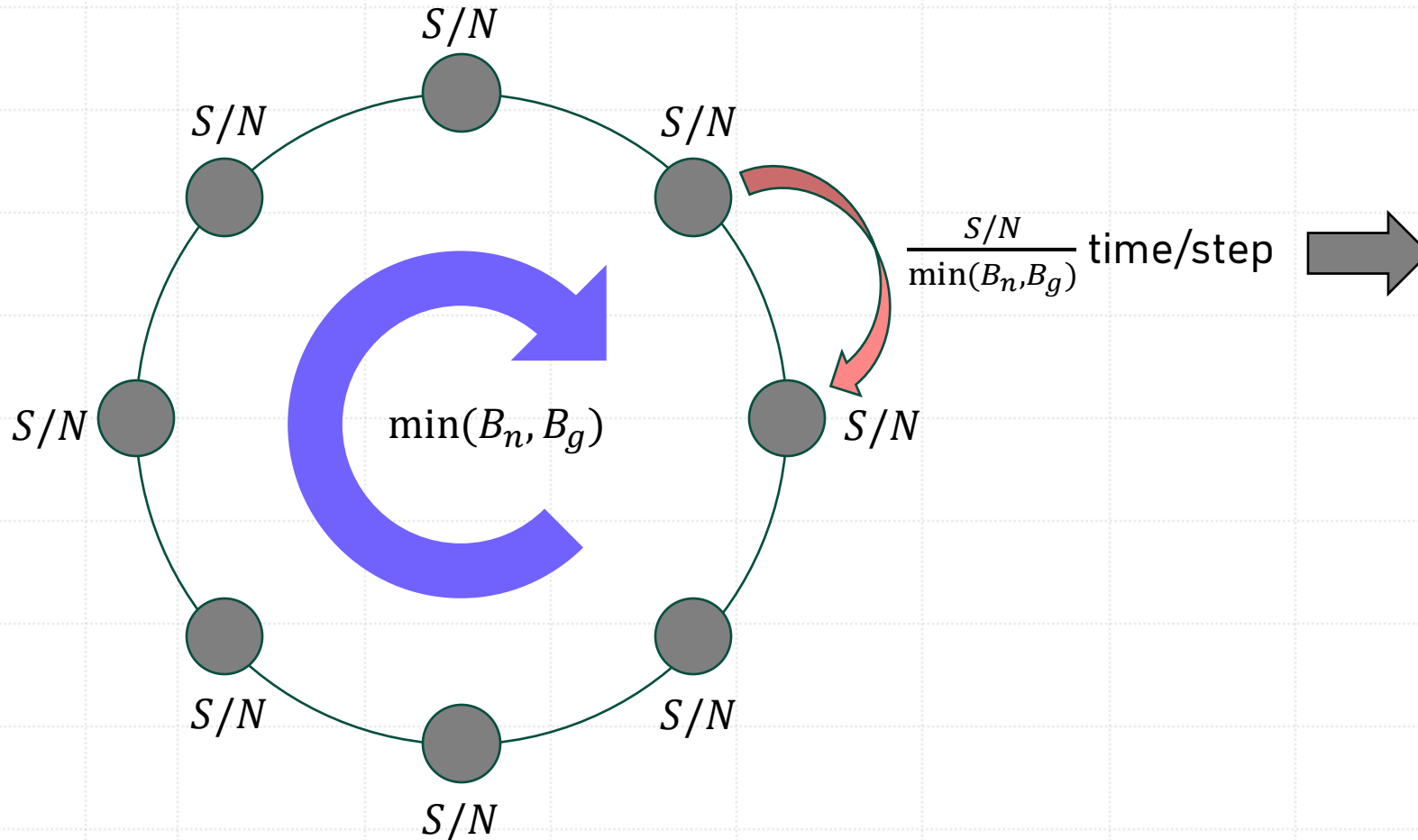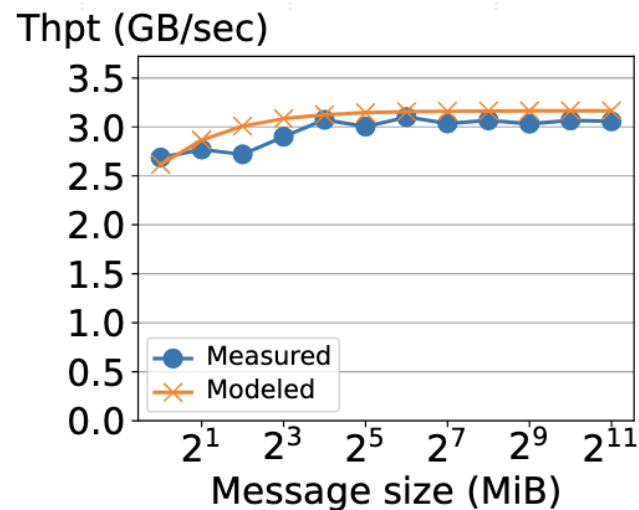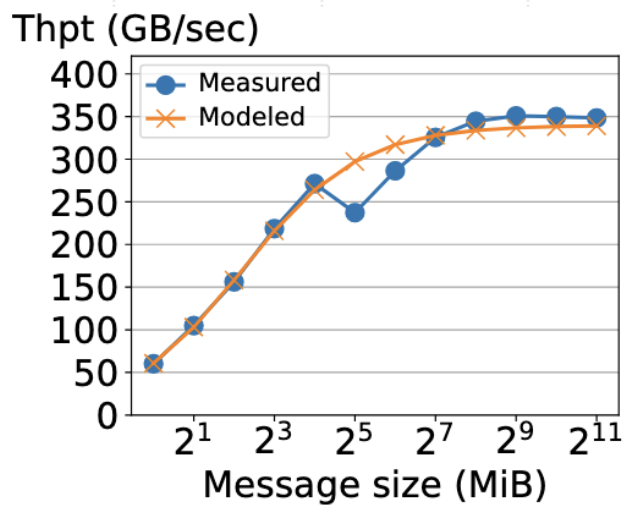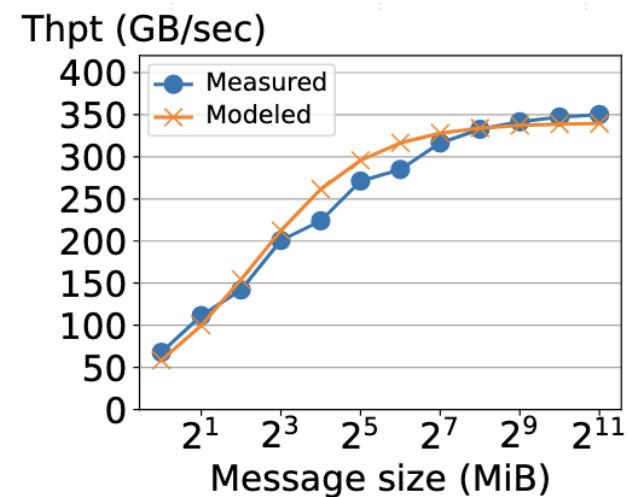# Modeling Broadcast with Skew

| | |
|---|---|
| k | Number of GPUs per machine (or node, VM) |
| V | Number of machines (or nodes, VMs) |
| N | Total number of GPUs = $k \times V$ |
| $B_g$ | Unidirectional inter-GPU bandwidth within each machine |
| $B_n$ | Unidirectional network bandwidth at each machine |
| S | Total dataset size processed by the GPU cluster. |
| $G_{ij}$ | The $i$-th GPU in the $j$-th machine. |
| $m_{ij \rightarrow pq}$ | The message sent from $G_{ij}$ to $G_{pq}$. |

$S$

0            0

$\min(B_n, B_g)$

0            0

0            0

0

### TABLE I (from https://arxiv.org/abs/2507.04786)
### COMPARISON OF NCCL COMMUNICATION PROTOCOLS

| | Simple | LL | LL128 |
|---|---|---|---|
| Design Goal | High bandwidth | Low latency | Low latency and high bandwidth |
| Synchronization Mechanism | Memory fences (high overhead) | Flag-based synchronization | Flag-based synchronization |
| Payload | Data chunks | 4B data + 4B flag | 120B data + 8B flag |
| Bandwidth Utilization | Near peak | 25 ~ 50% of peak [11] | ~ 95% of peak [11] |
| Latency Per-hop | ~ $6\mu s$ | ~ $1\mu s$ | ~ $2\mu s$ |

Due to pipelining and ring formation, the skewed data distribution has a minimal impact on performance, except for a slight increase in latency.
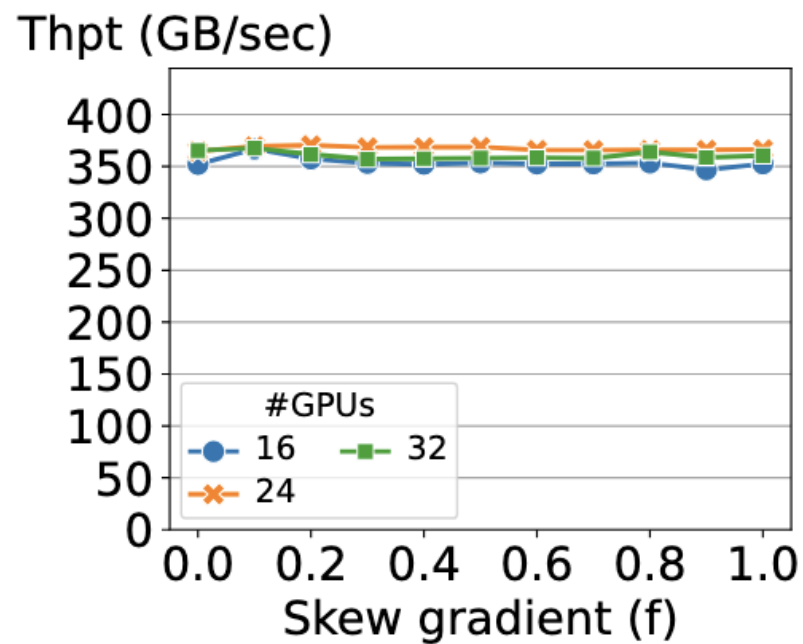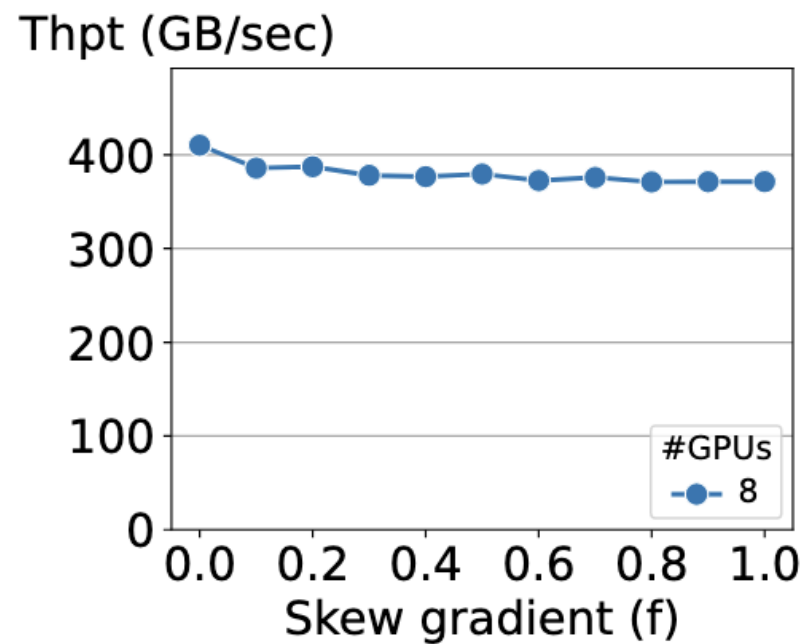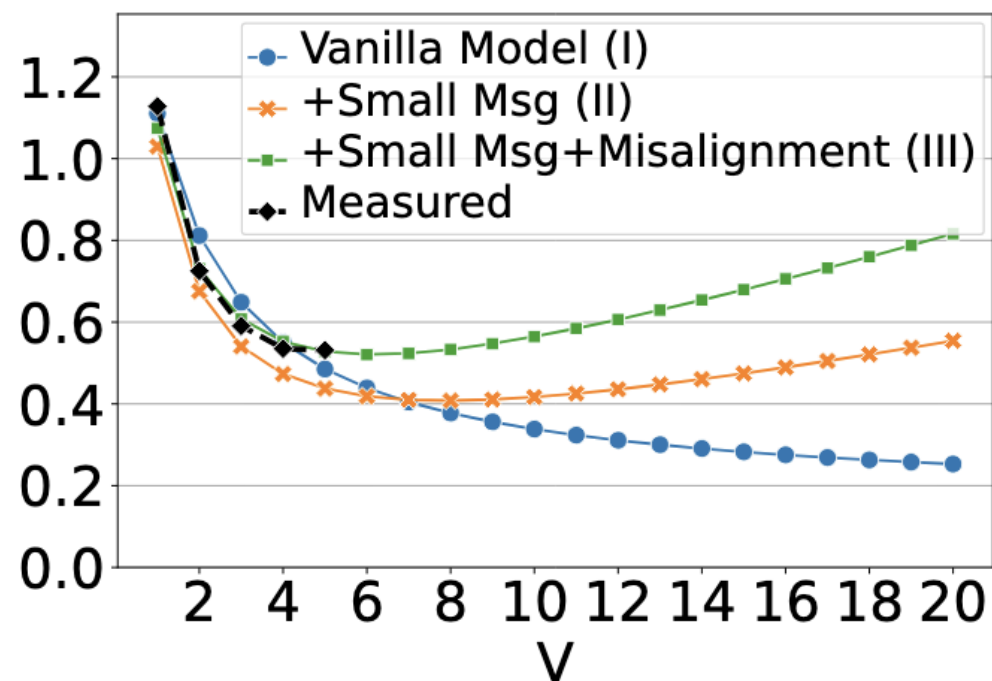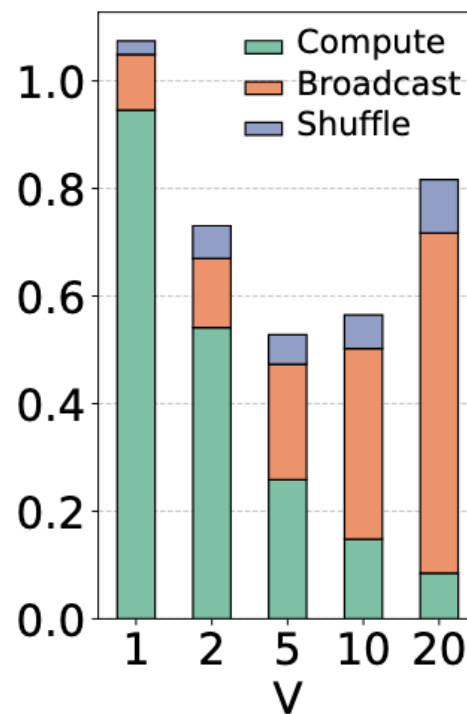
# Model Skew



**Figure 8: Broadcast + data skew.**

# Model TPC–H Queries



(a) Models for TPC-H projections.

(b) Project breakdown.

# Key Takeaways

- Using ML-based software to implement DBs gives good performance when the problem size is large enough.

- Small message sizes, skew, and buffer misalignment are the biggest sources of inefficiencies.

- Together with MaxBench, we see that with fast CPU-GPU, GPU-GPU, and GPU-network interconnects, the query performance often becomes GPU-bound.
  ➔ This calls back a previous work we published: *"Efficiently Processing Joins and Grouped Aggregations on GPUs"*, SIGMOD 2025.

# Next Steps