# Utilizing fast interconnects on GPUs for data processing

Data Systems Seminars - U Waterloo

March 12, 2024

Prof. Tilmann Rabl

Data Engineering Systems

Hasso Plattner Institute

# Who am I?

- Until 2011: PhD in CS at University of Passau
  - Distributed databases

- Until 2015: Postdoc at University of Toronto
  - Big data systems / benchmarking

- Until 2019
  - Visiting Professor & Research Director at DIMA group, TU Berlin
  - Deputy Director of Department IAM at DFKI
  - Scientific Coordinator of the Berlin Big Data Center

- Since Mai 2019
  - Professor for Data Engineering Systems, Digital Engineering Faculty, HPI, U Potsdam

- Other activities
  - HPI Ombudsperson
  - Director of HPI Data Center

# Hasso Plattner Institute

- Computer Science Institute in Potsdam, Germany
- B.Sc. and M.Sc. degree programs
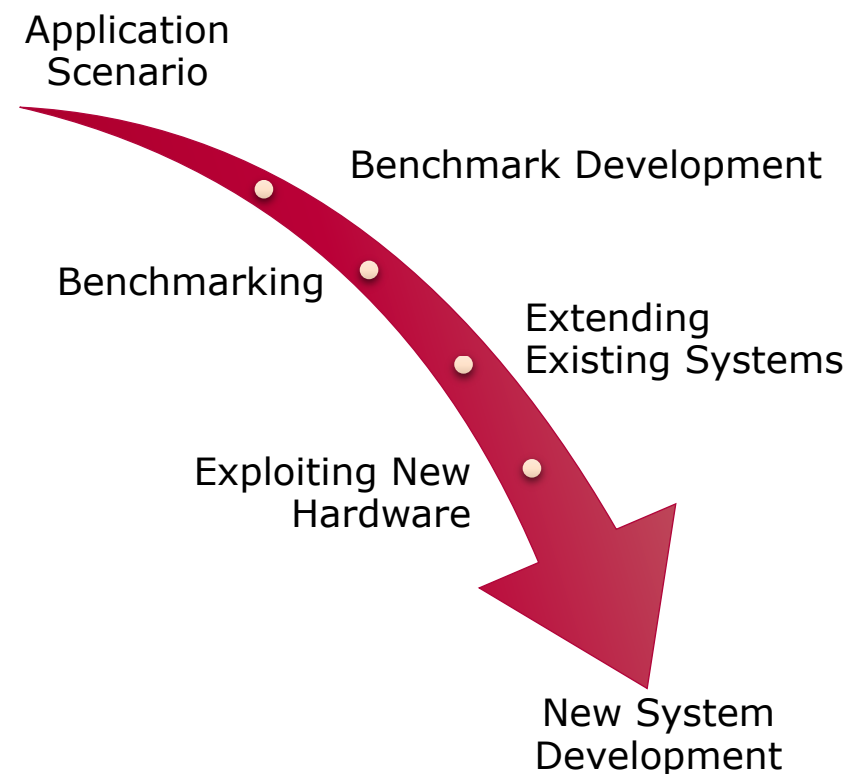- Department of University of Potsdam

# Research Topics

- **Database Systems (on Modern Hardware)**
  - ICDE 21, PVLDB 21, SIGMOD 22, SIGMOD 23

- **Stream Processing / Real-time Analytics**
  - PVLDB 20, TODS 21, SIGMOD 22, EDBT 23

- **Machine Learning Systems**
  - SIGMOD 20, ICDE 21, EDBT 22, EDBT 23

- **Benchmarking**
  - TPCTC 21, SIGMOD 21, PVLDB 22, SIGMOD 23

## Research Approach

Application Scenario

Benchmark Development

Benchmarking

Extending Existing Systems

Exploiting New Hardware

New System Development

# This work



- Clemens Lutz

- Tobias Maltenberger

- Ivan Ilic

# This Talk

1. Quick 101 on data processing on GPU

2. Scalable Joins on a single GPU

3. Mult-GPU Sorting

# Big Fast Data Analysis

- Data is growing

- Messages, tweets, social networks (statuses, check-ins, shared content), blogs, click streams, various logs, …

- Everyone is interested!
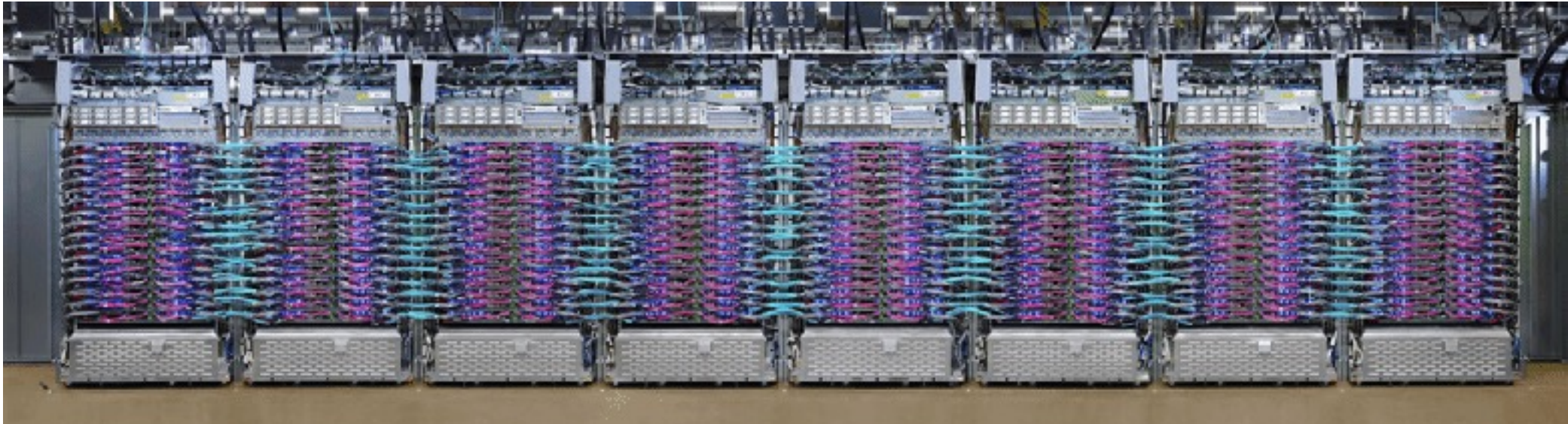
- **The value of data is decreasing with its age!**

# The Hardware Challenge

- Hardware performance is not simply increasing anymore

- Single thread
  –> multi thread

- Multi core
  –> specialization

- It is getting harder to be efficient



42 Years of Microprocessor Trend Data

Transistors (thousands)

Single-Thread Performance (SpecINT x $10^3$)

Frequency (MHz)

Typical Power (Watts)

Number of Logical Cores

Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2017 by K. Rupp

# Hardware – Industry Trend



- TPU cloud, tight coupling, scale up, special hardware

- "Only way to meet growing compute demands" - Amin Vahdat (Google)

- *We need to be hardware-conscious and make systems efficient - Tilmann*

# GPU-accelerated Data Processing

# Graphics Processing Unit (GPU)

- Highly specialized coprocessing chip

- Image rendering -> parallel computations

- Performance increased ~2.4x yearly

- During Moore's law CPU's performance increased 1.8x



Wikipedia

- Programming model was based on DirectX

- Targeting workloads for CPUs required major adaptations

# GPGPUs

- CPU cores != GPU cores -> Fundamental design differences

- CPU cores
  - Minimizing the latency of arithmetic operations
  - Latency-oriented design
  - Tens of cores

- GPU cores
  - Large number of FLOPS and memory accesses
  - Throughput-oriented design (1 555 – 3 000 GB/s)
  - Thousands of cores

- Massive parallelism
  - CPU bandwidth -> ~100s GB/s
  - GPU bandwidth -> ~1555 GB/s (A100) – 3000 GB/s (H100 - 2023)



| Core | Control | Core | Control |
| L1 Cache | | L1 Cache | |
| Core | Control | Core | Control |
| L1 Cache | | L1 Cache | |
| L2 Cache | | L2 Cache | |

L3 Cache

DRAM

CPU

L2 Cache

DRAM

GPU

CUDA Programming Guide

# GPU Architecture

- Compute Units
  - Computation cores
  - Register files
  - L1 cache
  - Shared memory

- Shared L2 cache
  - All compute units can access it

- Global Memory

| Compute Unit (CU or SM) | Compute Unit (CU or SM) | Compute Unit (CU or SM) |
|---|---|---|
| Register | Register | Register |
| L1 Cache ~ 192KB / Local Memory | L1 Cache ~ 192KB / Local Memory | L1 Cache ~ 192KB / Local Memory |

**L2 Cache (up to 400 MB)**

**~1.5 - 3 TB/s**

**GPU Global Memory (up to 188 GB)**

**~ 16 - 900 GB/s**

**Host Memory (DRAM)**

# GPU Interconnects

# System Interconnects

- PCI Express

  - Connects CPU to storage, memory and coprocessors

  - Most commonly used CPU-GPU interconnect

  - Current systems use versions 3.0 and 4.0

  - x1, x4, x8, x16 lanes

  - Not necessarily all lanes are used

  - PCIe 3.0x16 ->16 GB/s, PCIe 4.0x16 -> 32 GB/s

- From PCIe 5.0 also Compute Express Link (CXL)

  - Cache coherent access DtoH and HtoD (like NVLink)



PCI Express x1, x4, x8, x16

# GPU Interconnects

- NVIDIA NVLink
  - Mainly used as an inter-GPU interconnect
  - High bandwidth P2P data transfers
  - NVLink 2.0 -> 150 GB/s, NVLink 3.0 -> 300 GB/s
  - NVLink 4.0 -> 900 GB/s



NVLink Bridge 2.0 (left) and 3.0 (right)

- NVIDIA NVSwitch
  - Switching element connecting up to 18 GPUs
  - Enables all-to-all GPU interconnectivity in a system
  - Non-blocking data transfers
  - GPU-GPU bandwidth -> up to 900 GB/s



NVSwitch

# Single CPU – 2 GPU

- Single socket CPU System

- Direct main memory access

- No NUMA effects

- Single/Shared PCIe x.0 CPU-GPU lane

- Thinkstation P620
  - CPU -> AMD Threadripper PRO 3995WX 64 Core
  - GPU -> 2x NVIDIA A5000 24GB NVLink 3.0 connected (56 GB/s)
  - CPU-GPU -> Single PCIe 4.0x16 lane (32 GB/s)

# Dual CPU – 4 GPU

- Dual socket CPU System
  - CPUs interconnected via a proprietary interconnect
  - NUMA effects

- Dedicated main memory access

- **CPU-GPU connected via NVLink**

- **Equal bandwidth across all processors**

- IBM AC922
  - CPU -> 2x IBM Power 9
  - GPU -> 4 x NVIDIA V100 32GB NVLink 2.0 connected (75 GB/s)
  - CPU-GPU -> NVLink 2.0 interconnect (75 GB/s)
  - CPU-CPU -> IBM XBus 64 GB/s

# Dual CPU – 8 GPU

- Single/Shared PCIe 4.0 CPU-GPU lane

- Switching elements between GPUs

- All to all GPU communication

- Full GPU-GPU Bandwidth

- GPU-GPU bandwidth close to global memory bandwidth



- DGX A100
  - CPU -> 2x AMD EPYC 7742
  - GPU -> 8x NVIDIA A100 40 GB NVLink 3.0 connected (300 GB/s)
  - CPU-GPU -> **Shared** PCIe 4.0x16 lane (32 GB/s)
  - CPU-CPU -> AMD Infinity Fabric (102 GB/s)

# NVIDIA Grace Hopper

- **Grace CPU**
  - ARM
  - 72 Cores
  - 480 GB DDR5x
  - 64x PCIe 5 lanes

- **Hopper GPU**
  - 96 GB HBM3

- **NVLink 4.0**
  - 900 GB/s

# Data Transfers – ThinkStation P620

- Single PCIe 4.0 channel per GPU enables efficient CPU-GPU parallel transfers

- Performance drops for bi-directional CPU-GPU transfers

- Direct communication between the two GPUs

- 2x performance discrepancy between CPU-GPU and P2P data transfers

- Significant performance drops due to NUMA effects

- CPU-CPU interconnect is the main bottleneck

- Using two GPUs often faster than all four

- No performance discrepancy between CPU-GPU and P2P data transfers

# Data Transfers - DGX A100

- Consistent P2P bandwidth

- Main memory transfers limited by the PCIe interconnect

- Parallel transfers limited by shared PCIe bus

- Despite having NVSwitch, parallel P2P transfers are faster between co-located GPUs

- >10x performance discrepancy between CPU-GPU and P2P data transfers

# Topology-aware Algorithm Design

- Pre-allocate memory and reduce data transfers

- Organize the P2P communication

- Maximize the bandwidth according to the system's topology

- When data is located on one socket, keep the computation on the same NUMA node

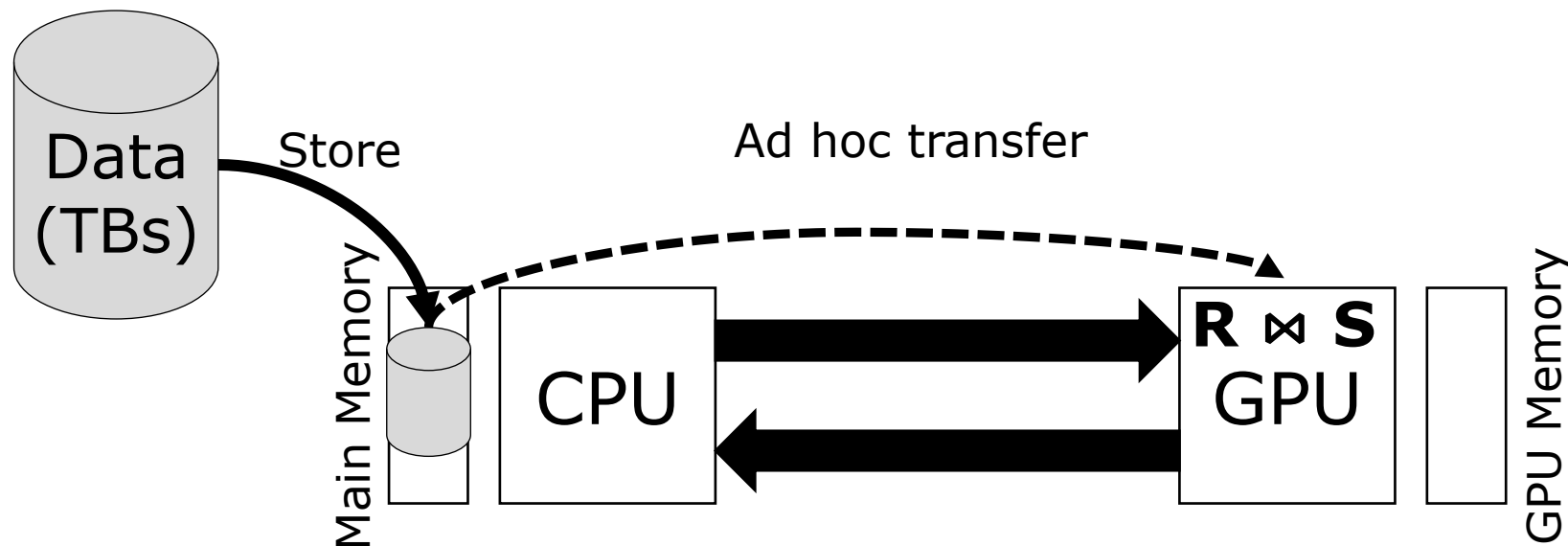



24

# Scalable GPU-based Join

# Goal

*Scale GPU-accelerated data management to **arbitrary data volumes**.*
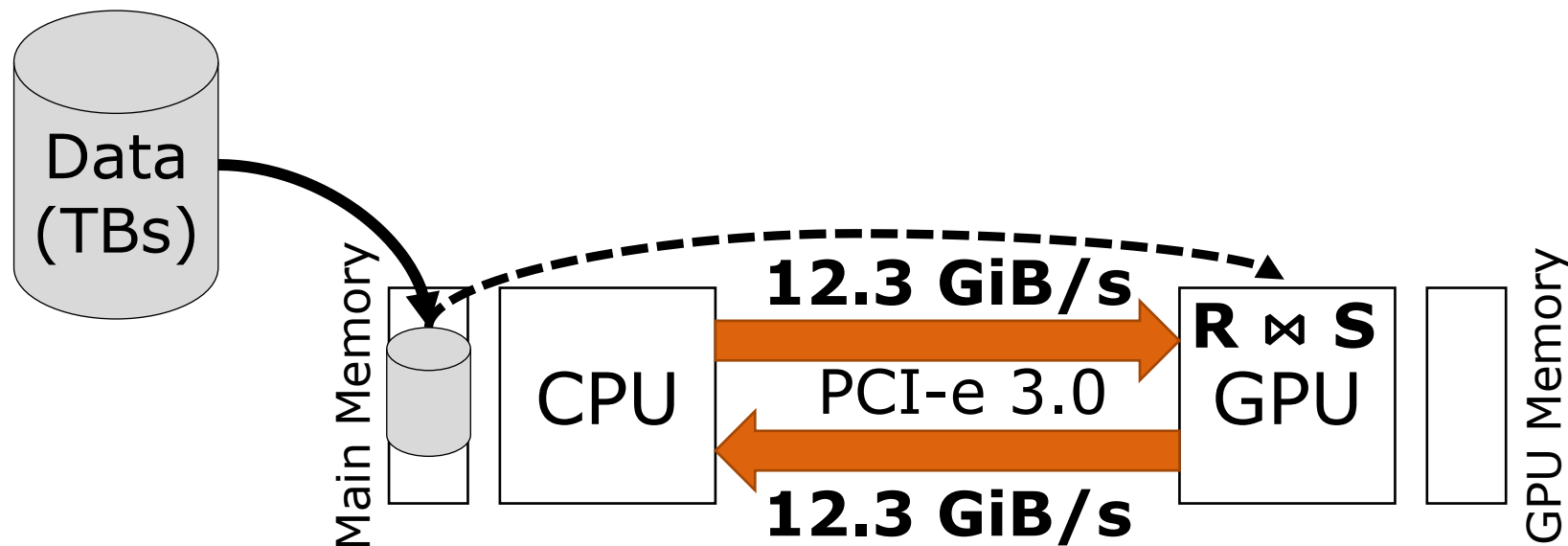
# Problem: Transfer Bandwidth



- Today's GPU databases:

- Store data in main memory
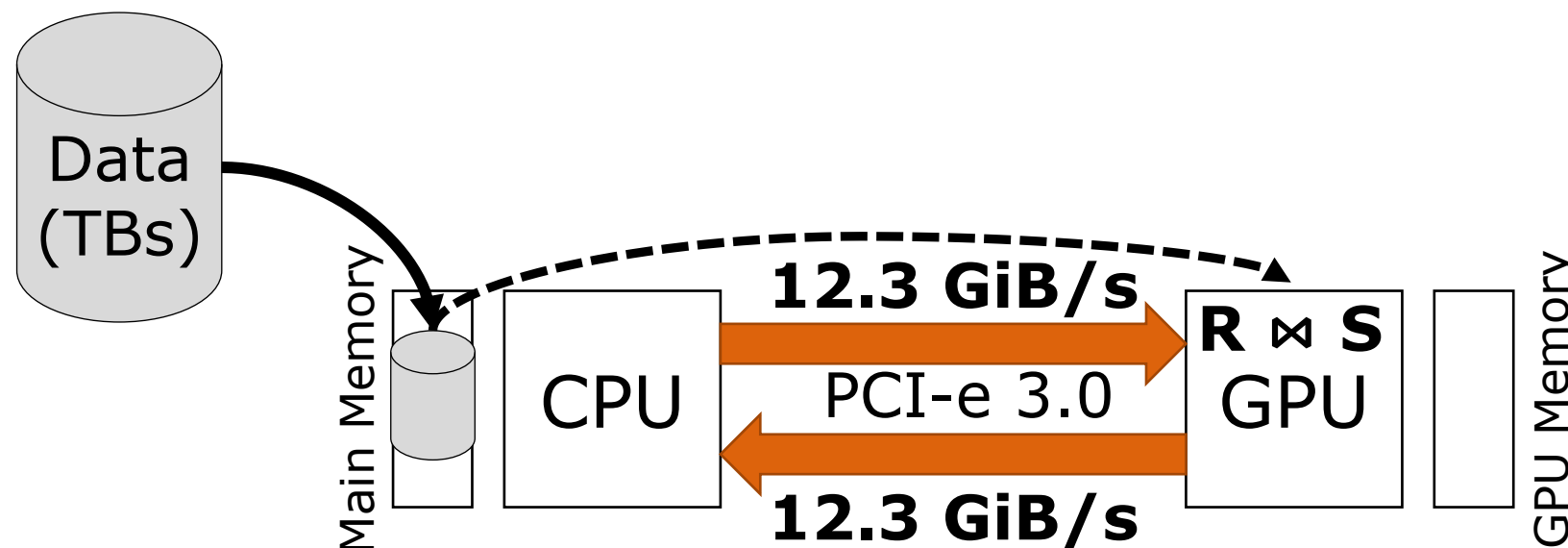
- Perform data processing on GPU

# Problem: Transfer Bandwidth



- Ad hoc data transfer over PCI-e 3.0

- GPU capable of much higher throughput

# Problem: Transfer Bandwidth



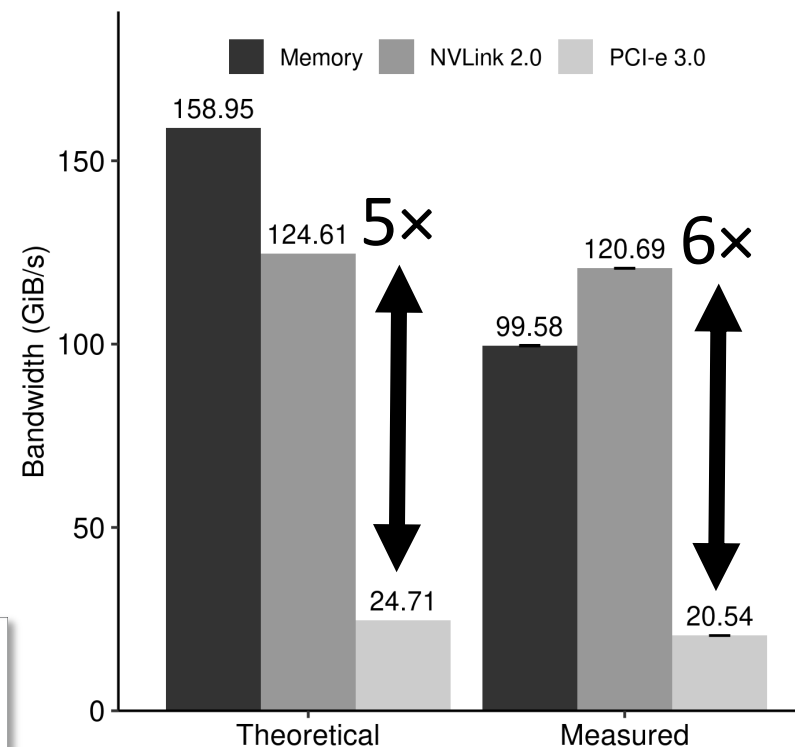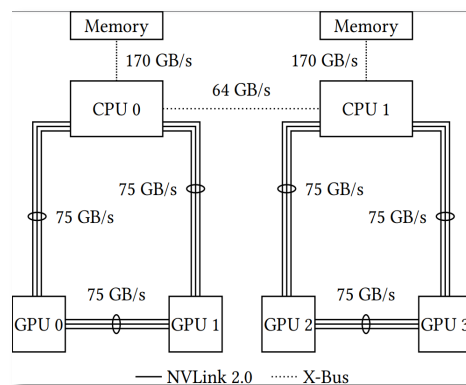**Interconnect bandwidth & GPU memory capacity limit scalability**

*"Transfer bottleneck"*

# Game Changer

- **Fast interconnects**
  - e.g., NVLink 2.0, Infinity Fabric, CXL

- **High bandwidth (124 GB/s total)**

- **System-wide cache-coherence**
  - data-dependent memory access
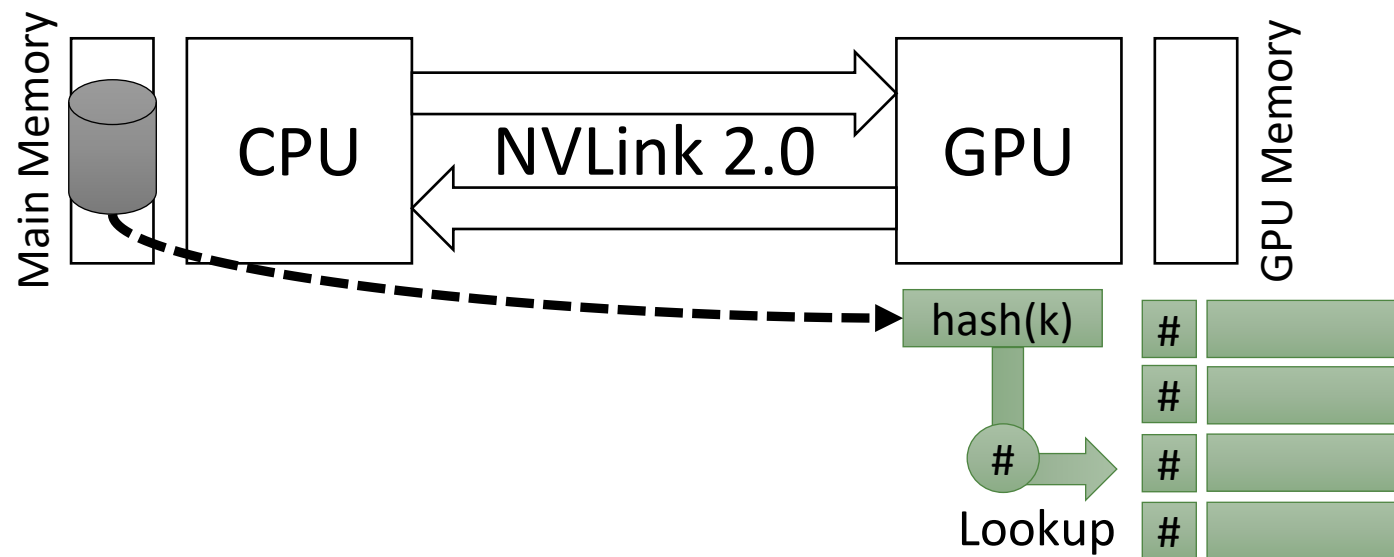  - fine-grained CPU+GPU cooperation

- **E.g., AC 922**

# Solution

Hash Join
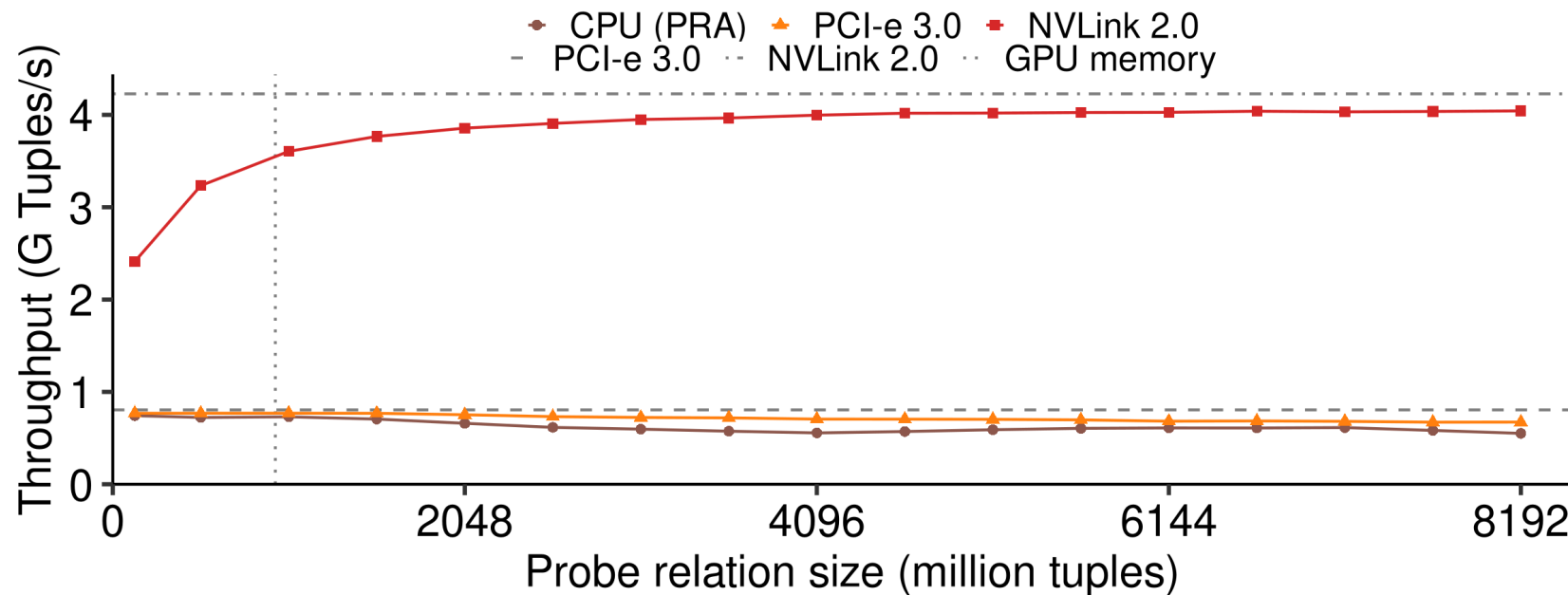- Probe-side scaling
- Build-side scaling
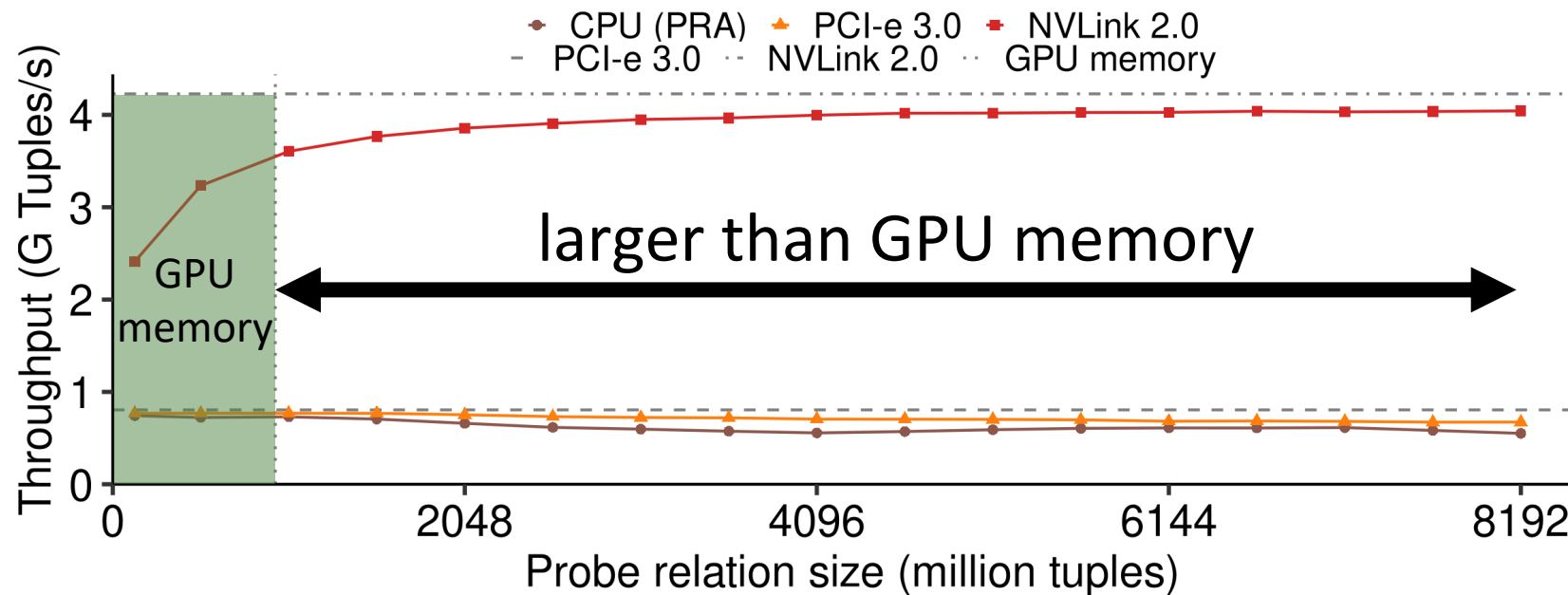- GPU+CPU cooperation

**Interconnect feature: High bandwidth**
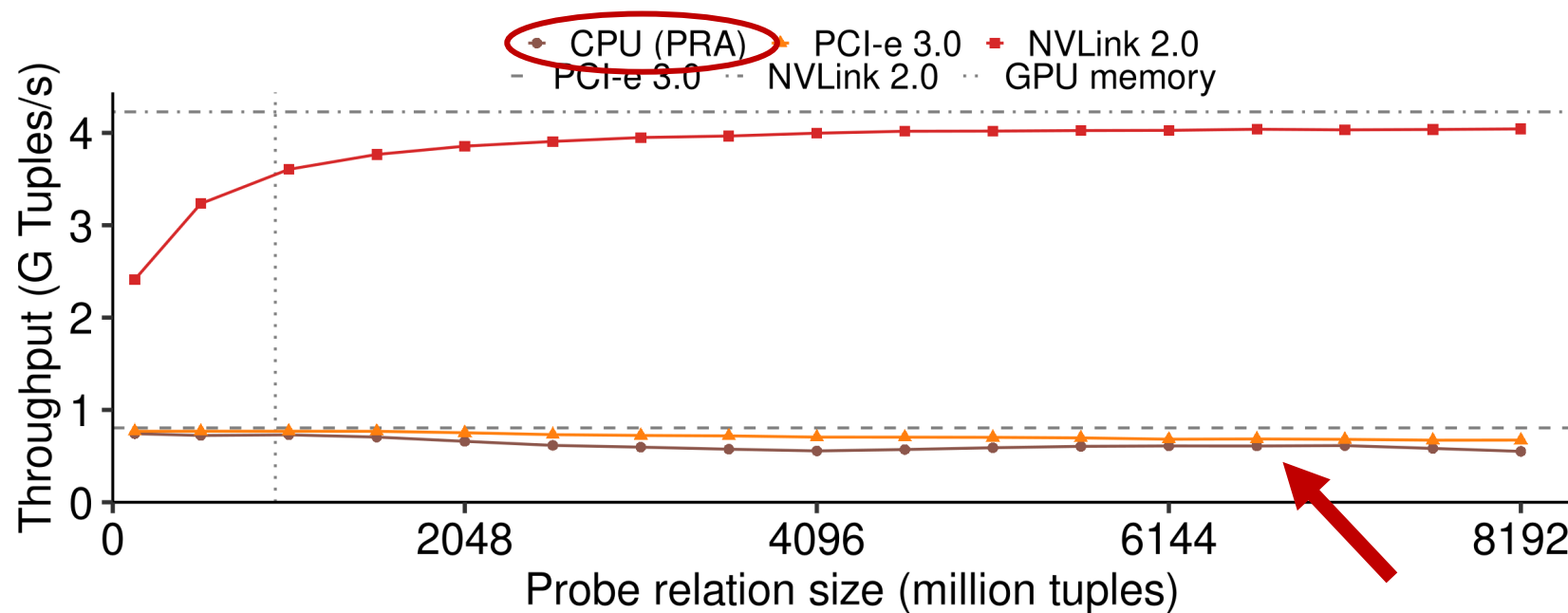
# Probe-Side Scaling



- Up to 2⋈122 GB

# Probe-Side Scaling



- Up to 2⋈122 GB

# Probe-Side Scaling

Throughput (G Tuples/s) vs. Probe relation size (million tuples)

Legend: CPU (PRA), PCI-e 3.0, NVLink 2.0, PCI-e 3.0, NVLink 2.0, GPU memory

- Up to 2⋈122 GB
- CPU baseline: Radix-partitioned hash join

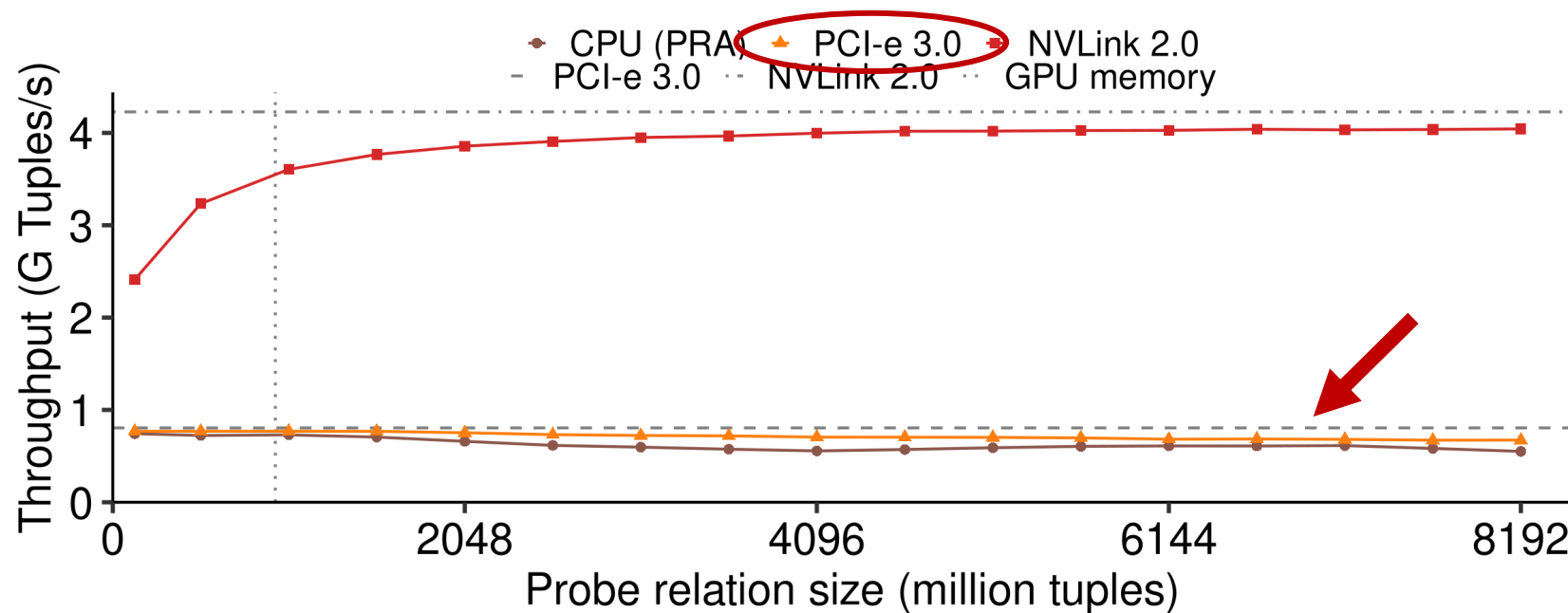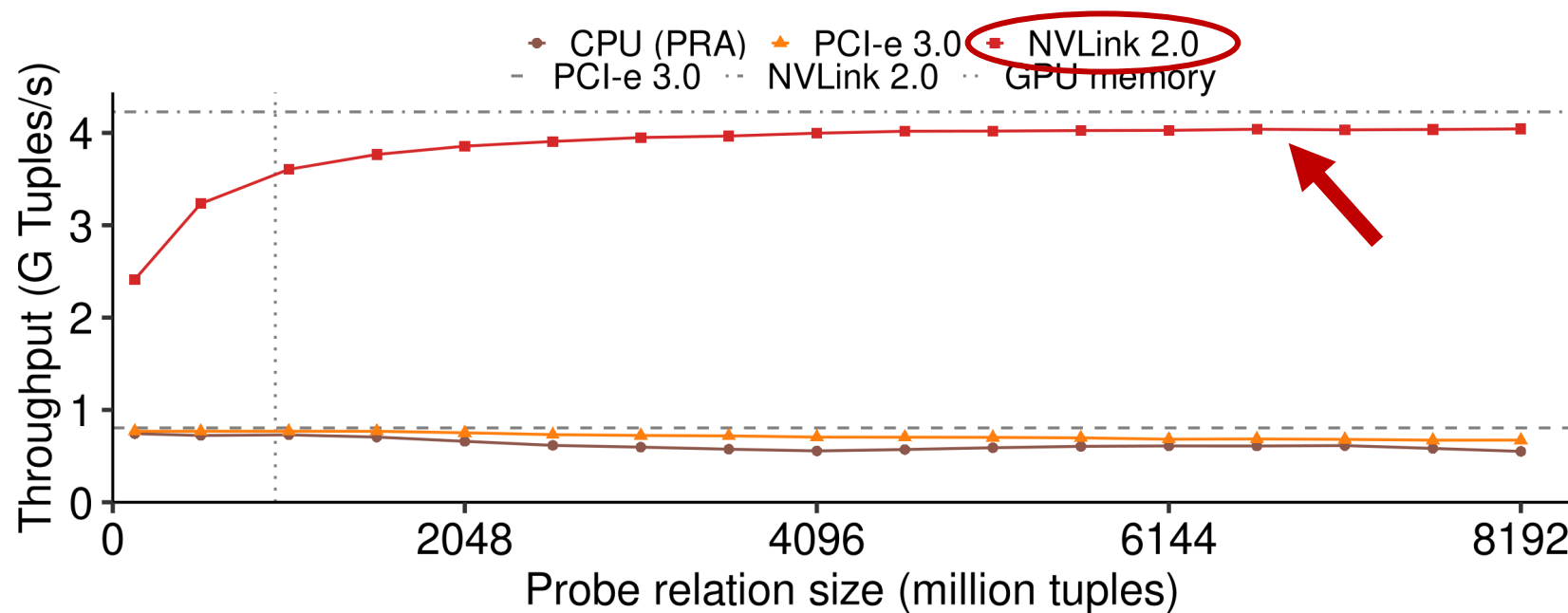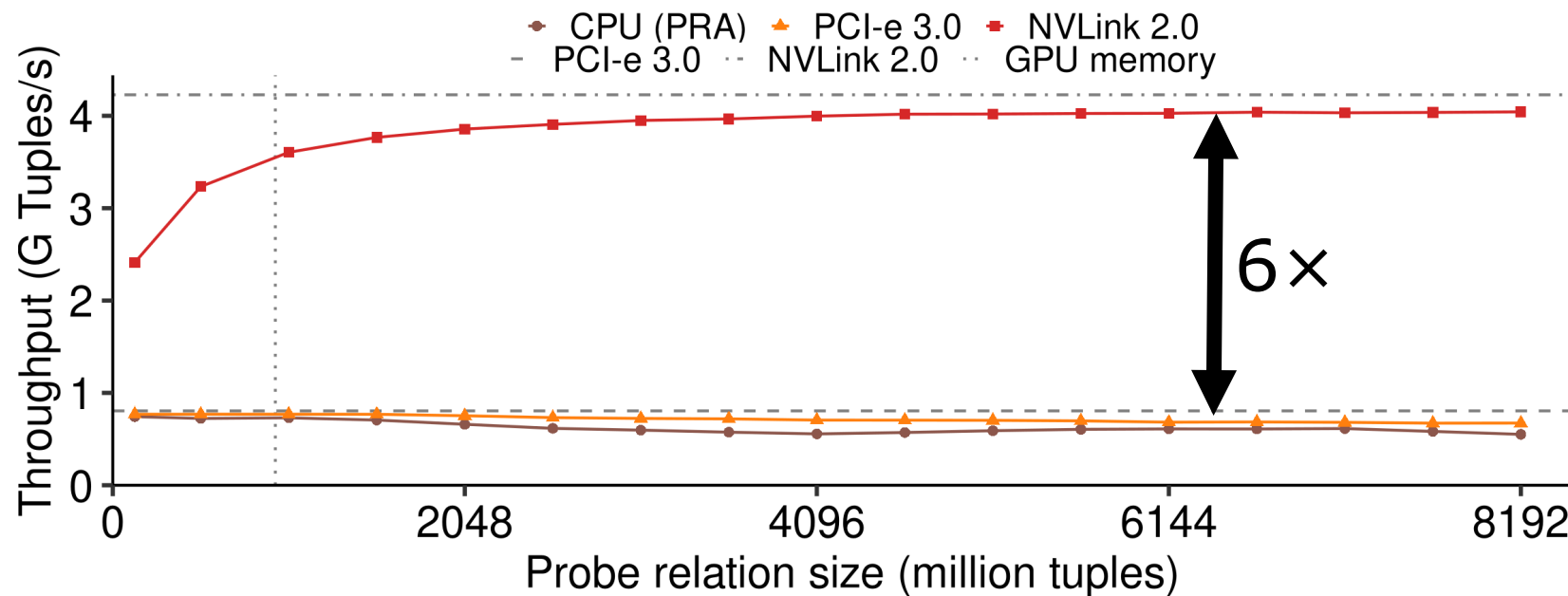# Probe-Side Scaling



- Up to 2⋈122 GB
- CPU baseline: Radix-partitioned hash join

# Probe-Side Scaling



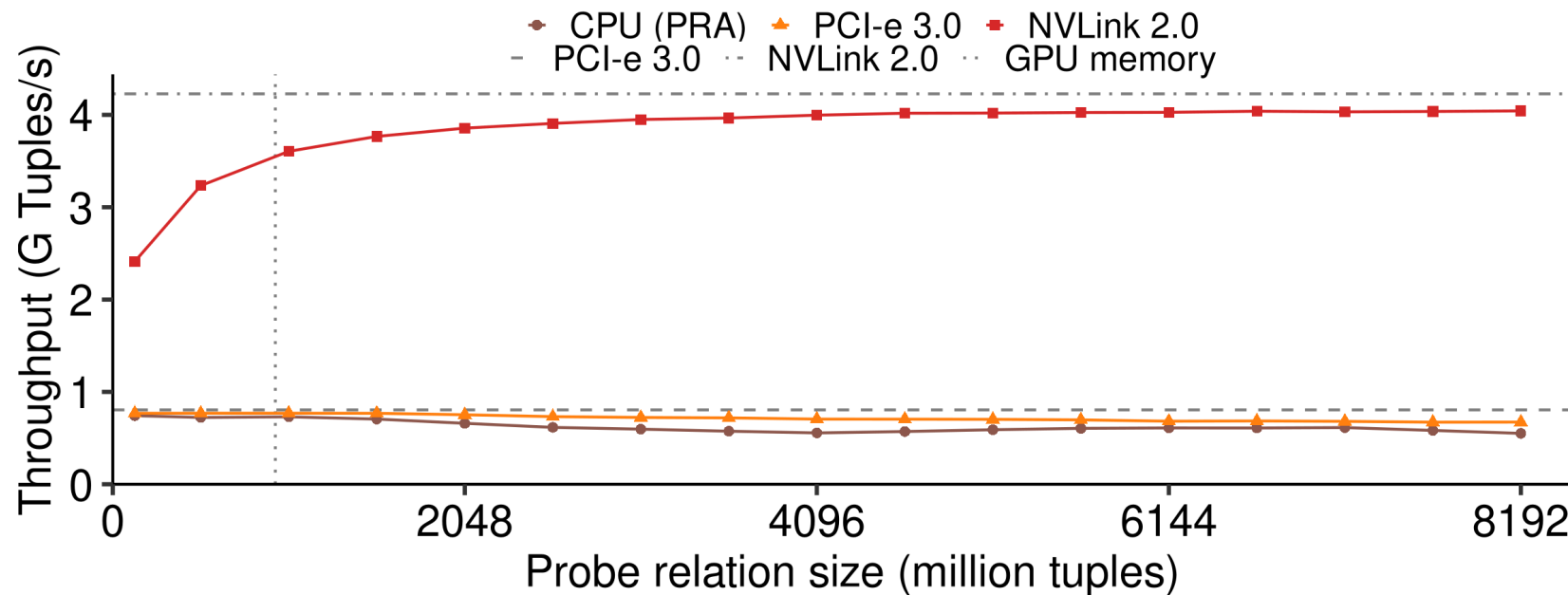- Up to 2⋈122 GB
- CPU baseline: Radix-partitioned hash join

# Probe-Side Scaling



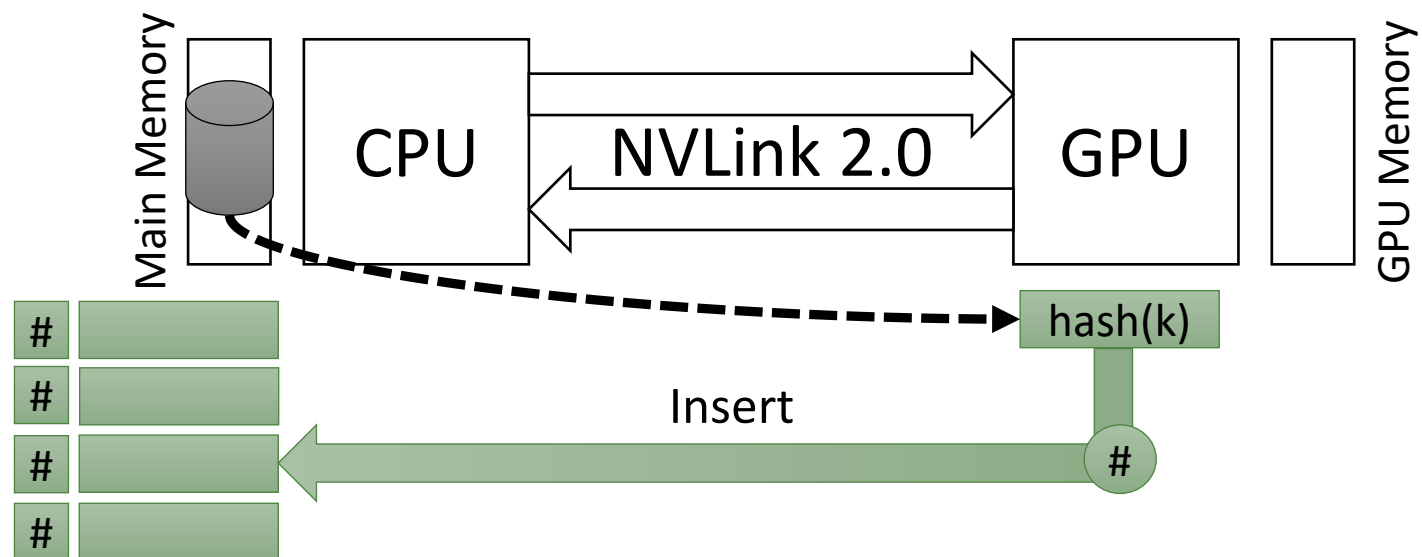- Up to 2⋈122 GB
- CPU baseline: Radix-partitioned hash join

# Probe-Side Scaling
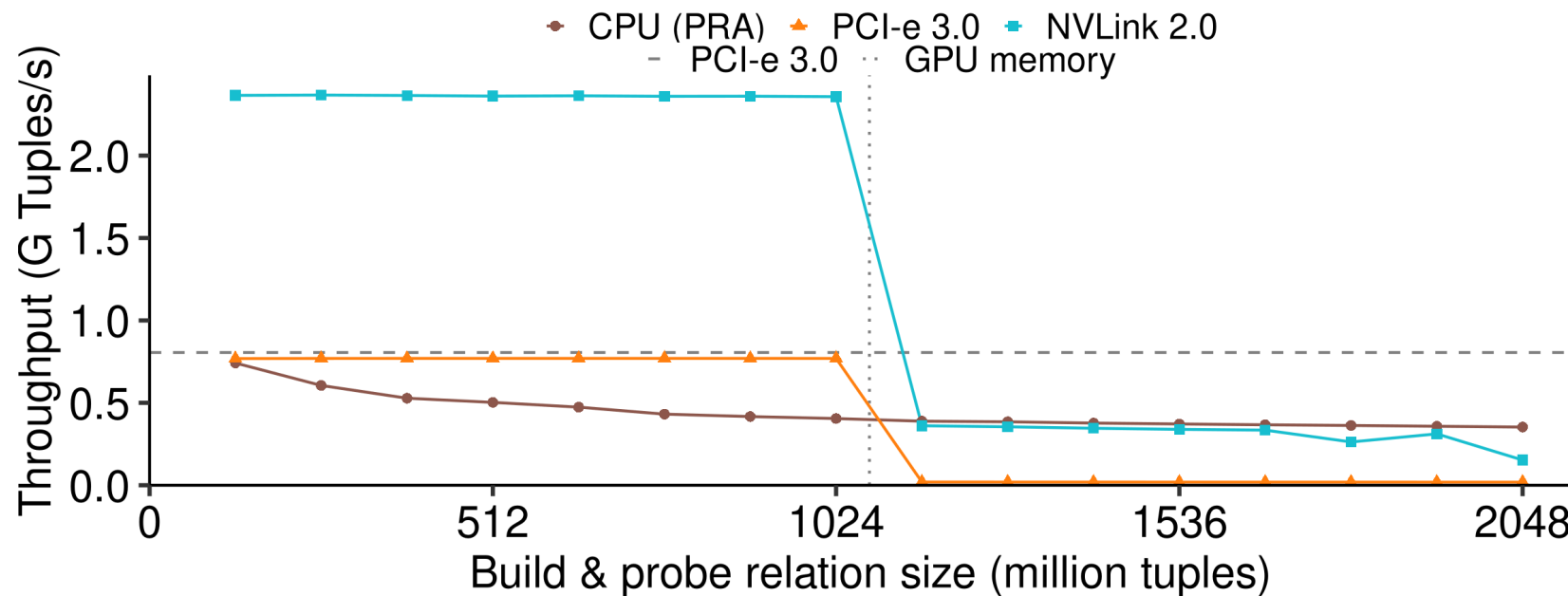


GPUs can efficiently process large, out-of-core data

# Build-Side Scaling



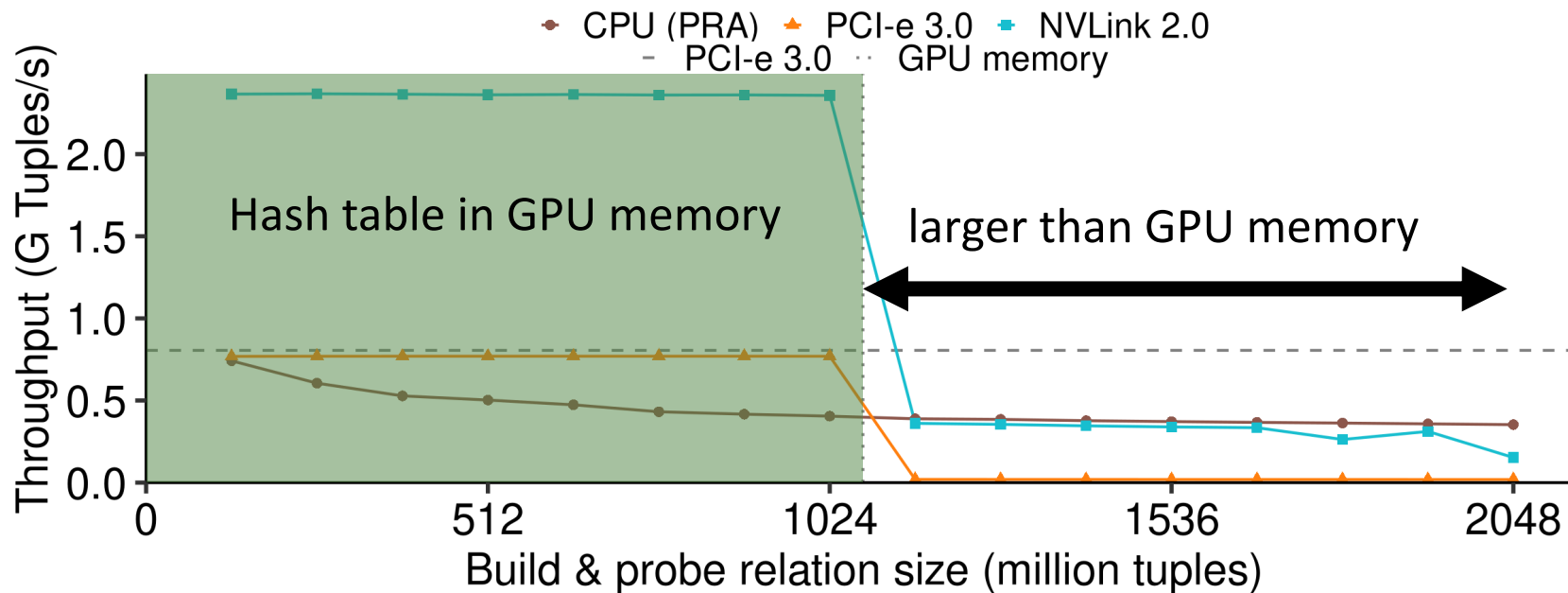**Interconnect feature: Data-dependent memory access**

# Build-Side Scaling



- Up to 30⋈30 GB with a 30 GB hash table **= 90 GB**

# Build-Side Scaling



Throughput (G Tuples/s) vs Build & probe relation size (million tuples)

Legend: CPU (PRA), PCI-e 3.0, NVLink 2.0, PCI-e 3.0, GPU memory

Hash table in GPU memory

larger than GPU memory

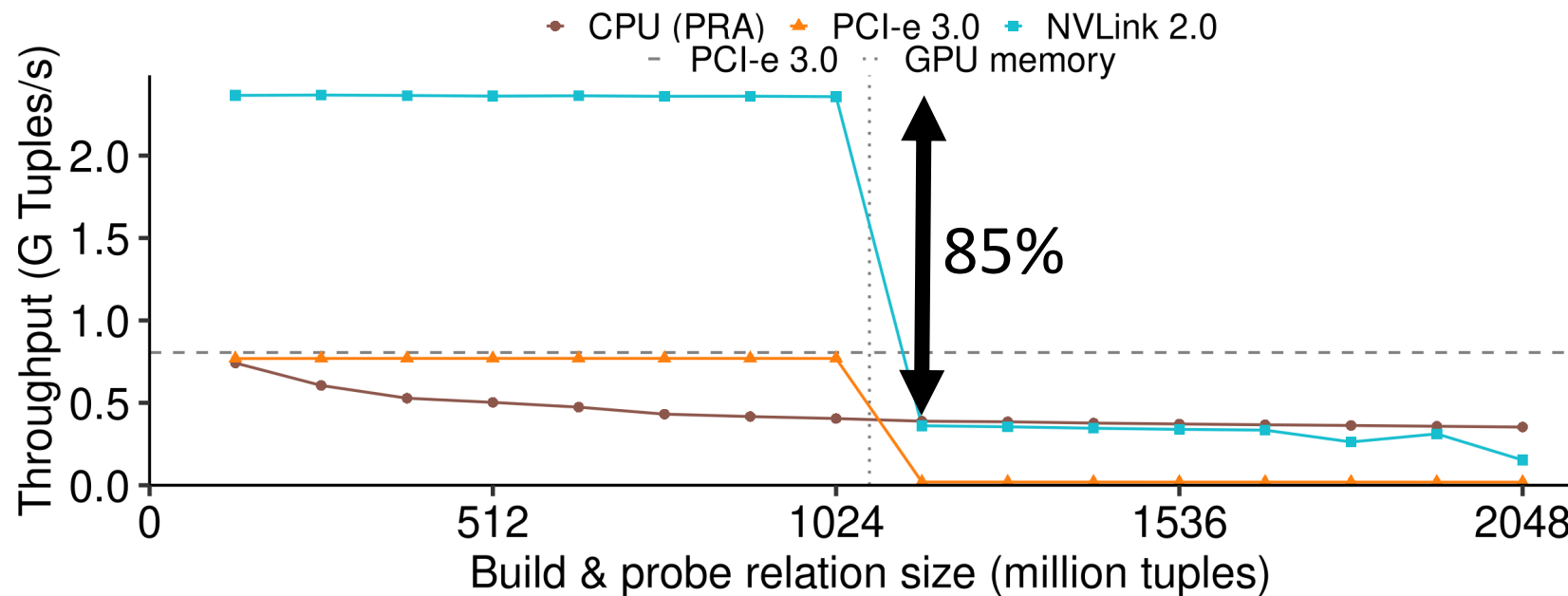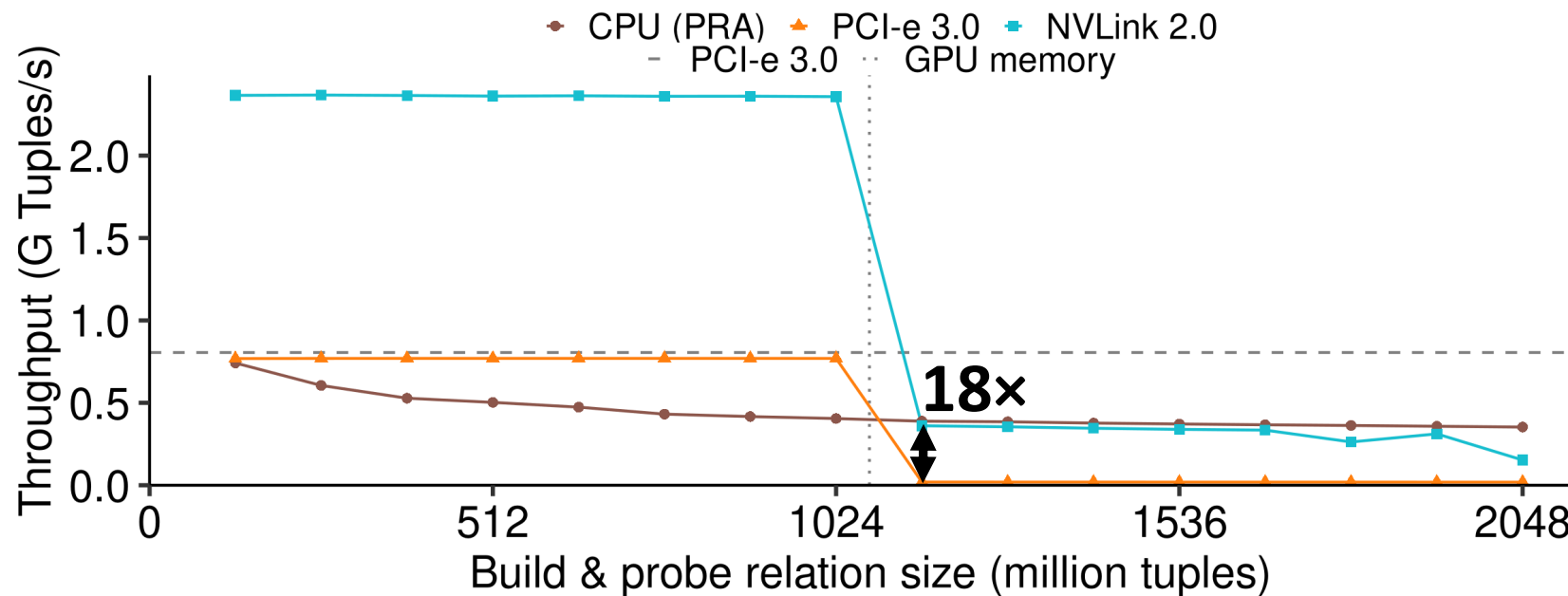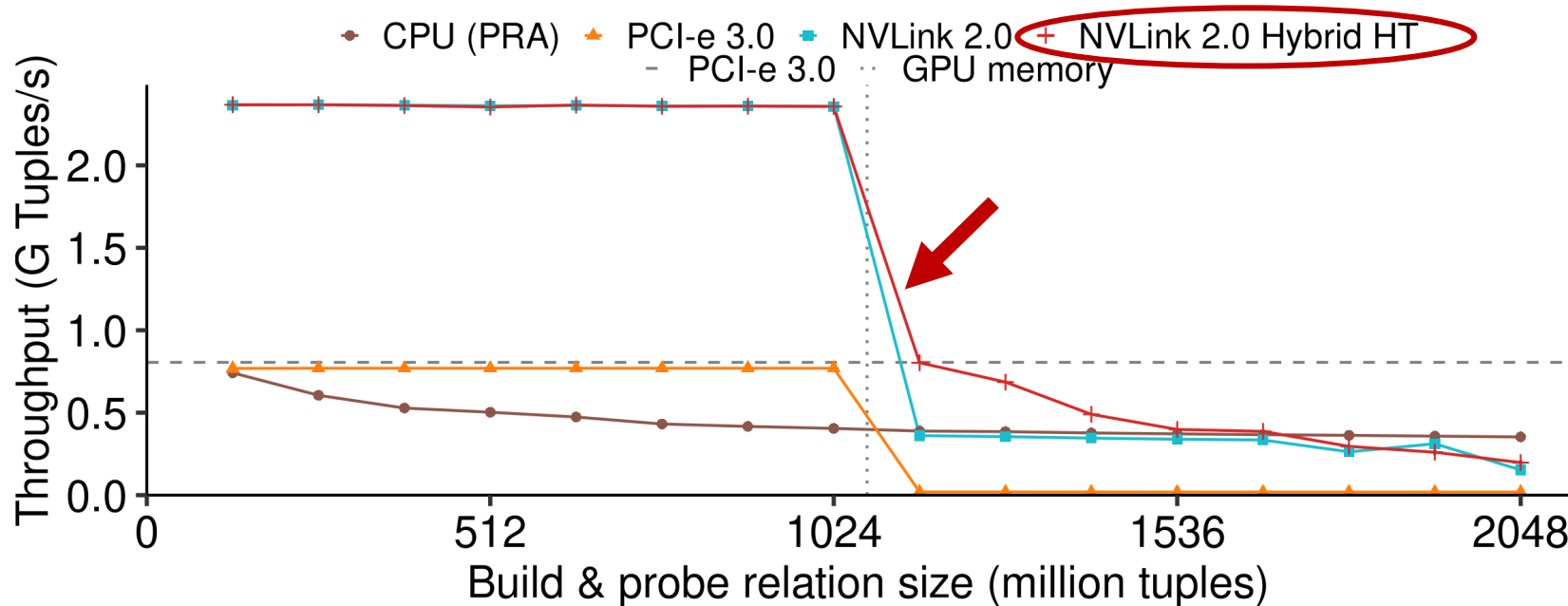- Up to 30⋈30 GB with a 30 GB hash table **= 90 GB**

# Build-Side Scaling



- Up to 30⋈30 GB with a 30 GB hash table **= 90 GB**

# Build-Side Scaling
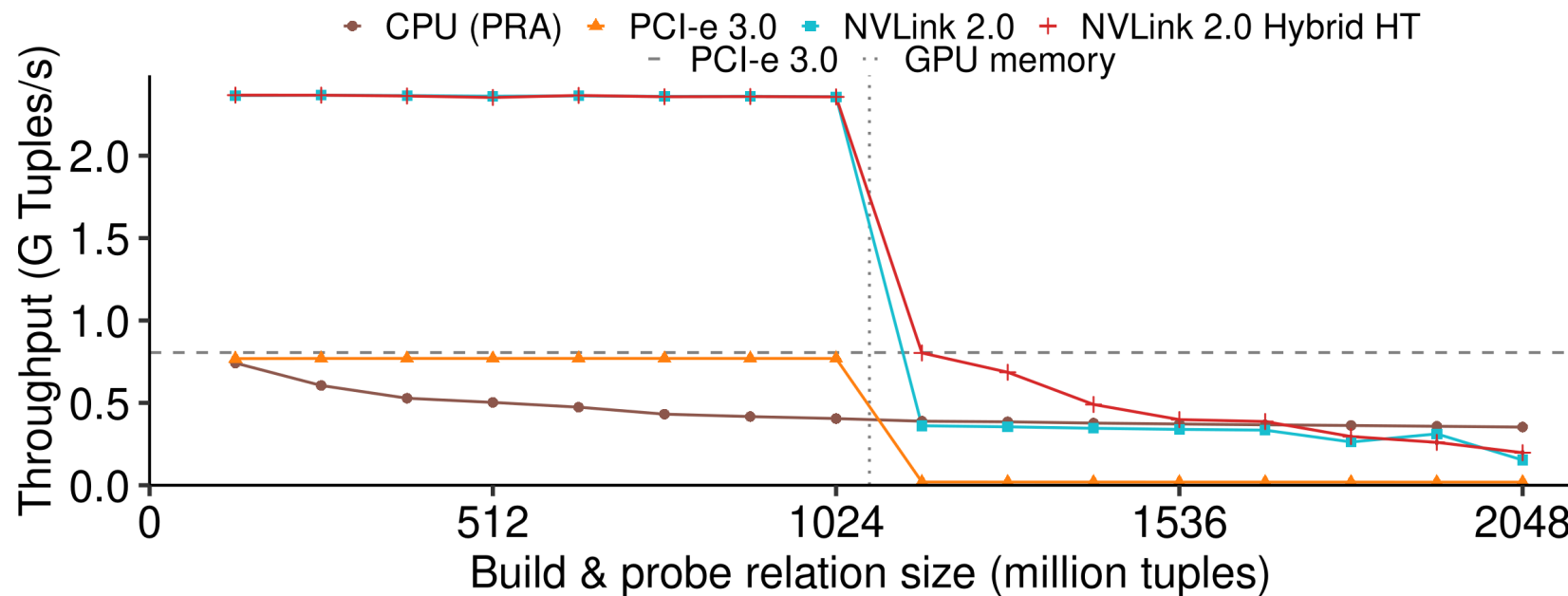


- Up to 30⋈30 GB with a 30 GB hash table **= 90 GB**

- Up to 30⋈30 GB with a 30 GB hash table **= 90 GB**

- Hybrid hash table spills to CPU memory

# Build-Side Scaling



GPUs are able to operate on large, out-of-core data structures

… but should cache data structures in GPU memory

# Conclusion

*We explore **in which ways** fast interconnects **benefit databases**:*

- Out-of-core **data sets**

- Out-of-core **data structures**
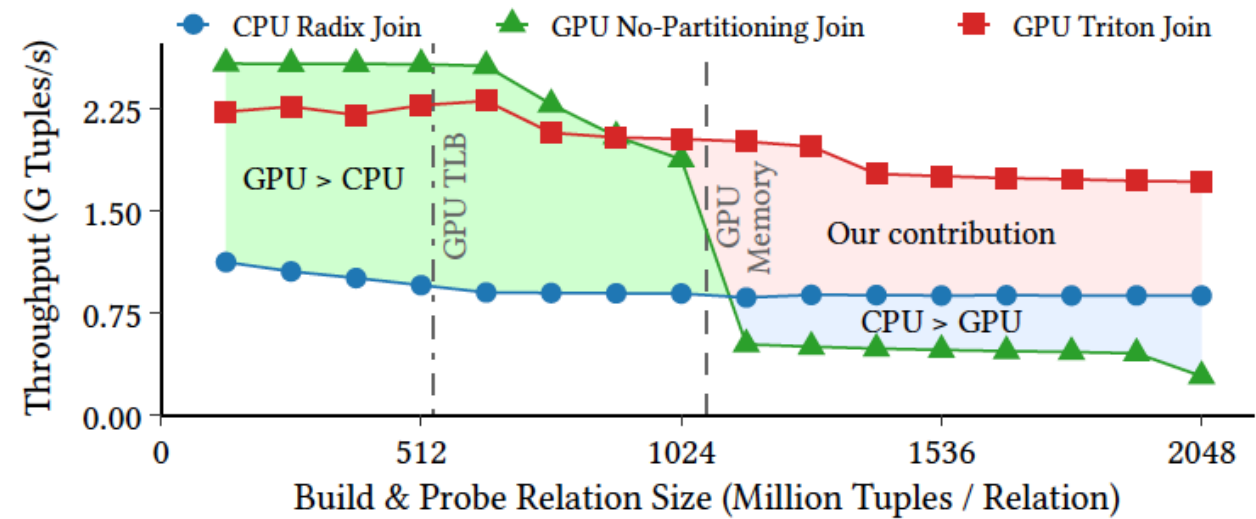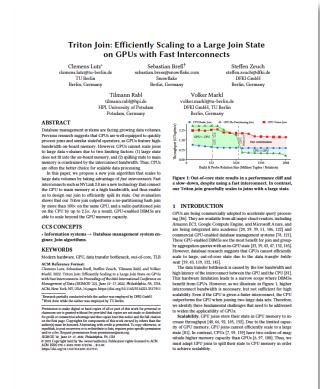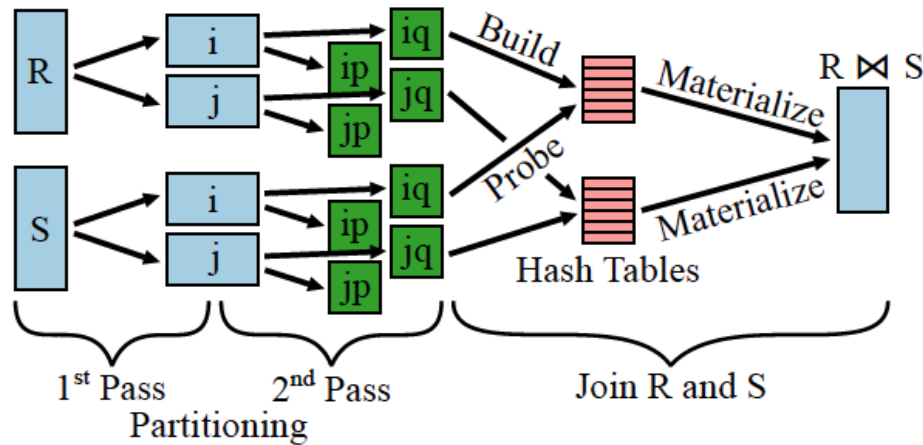
- **Fine-grained** cooperative co-processing

*Key technology enablers:*

- Hardware accelerators

- Fast interconnects

- Radix-partitioned GPU hash join

- Hierarchical partitioning



*Triton Join: Efficiently Scaling to a Large Join State on GPUs with Fast Interconnects – Clemens Lutz et al. SIGMOD 2022*

# Multi-GPU Sorting

**Evaluating Multi-GPU Sorting with Modern Interconnects @ SIGMOD 2022**

*Tobias Maltenberger, Ivan Ilic, Ilin Tolovski, Tilmann Rabl*

**RMG Sort: A Radix-Partitioning-Based Multi-GPU Sorting Algorithm @ BTW 2023**

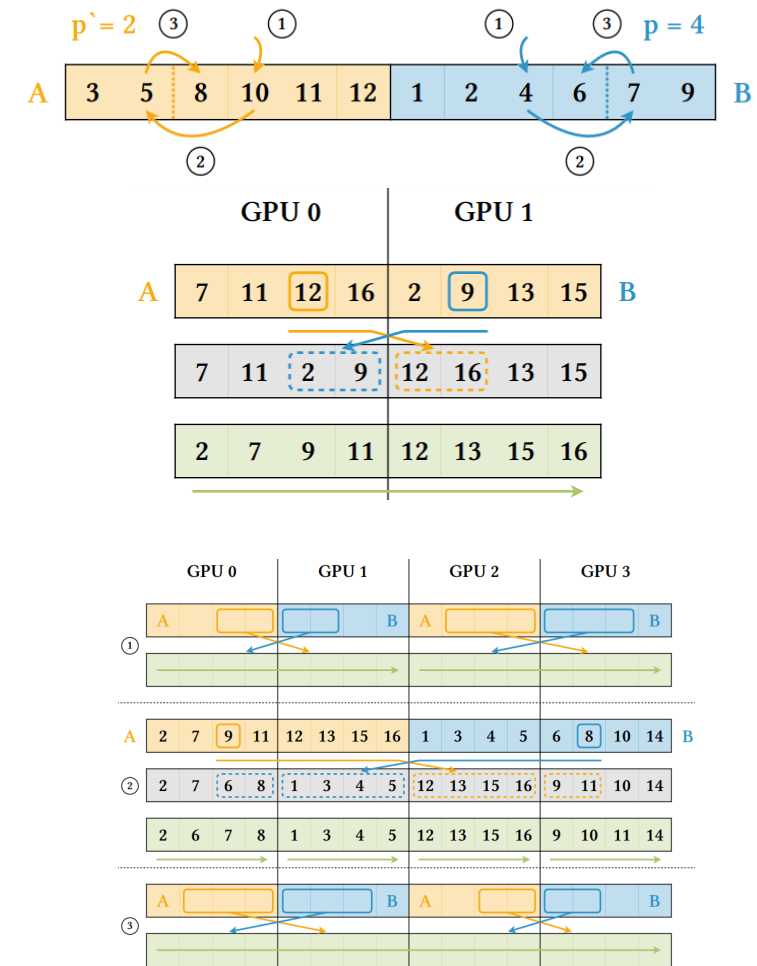*Ivan Ilic, Ilin Tolovski, Tilmann Rabl*

# Peer-to-Peer-based Sorting

- Algorithm stages:
  1. Split data across the GPUs
  2. Find the pivot on each GPU
  3. Exchange data between GPUs
  4. Sorted subarrays

- In-place sorting algorithm

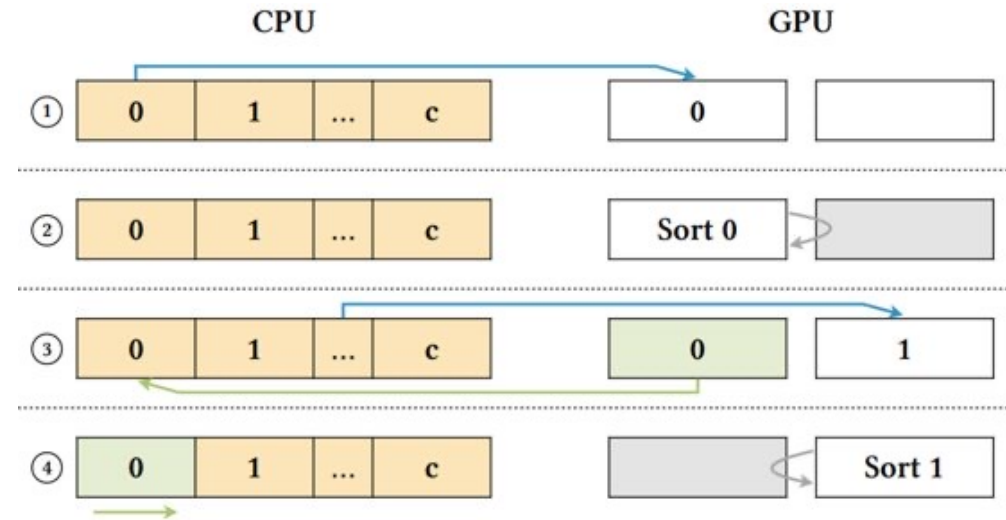- Algorithm limited by the total GPU memory

- Max. data to sort: ~50% of total GPU memory

# Heterogenous CPU-GPU Sorting

- Algorithm stages:
  1. Fill a single GPU with data
  2. Execute an in-place sort
  3. Bi-directional CPU-GPU data exchange
  4. CPU Merge & GPU Sort executed in parallel
  5. Final merge on CPU



- Sorted data can be larger than combined GPU capacity

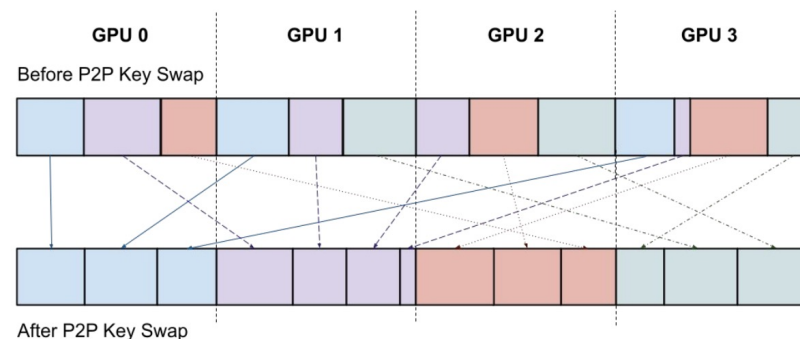- Main memory bound algorithm

# Radix-Partitioning-based Sorting

- Algorithm stages:
  1. Partition keys based on radix value (MSB)
  2. Exchange buckets between all GPUs via P2P transfers
  3. Sort buckets on each GPU using a single-GPU sorting primitive

- In-place sorting algorithm

- Algorithm limited by the total GPU memory

- Max. data to sort: ~50% of total GPU memory

```
00000010 00111010 01110001 01100010  →  Bucket [2]
...
00000001 10101011 01011000 00110101  →  Bucket [1]
```



Before P2P Key Swap

After P2P Key Swap
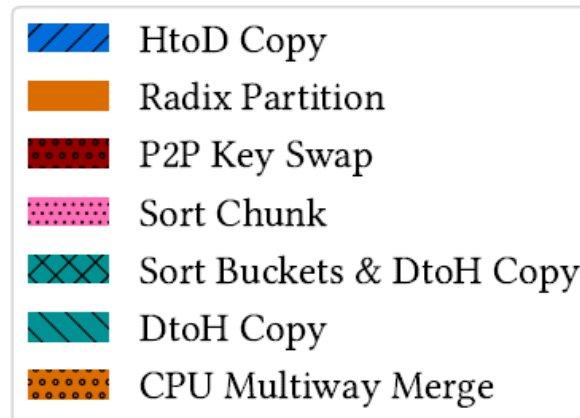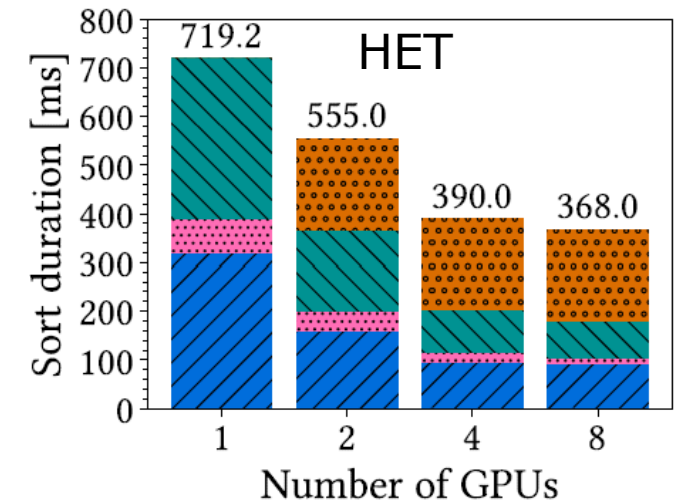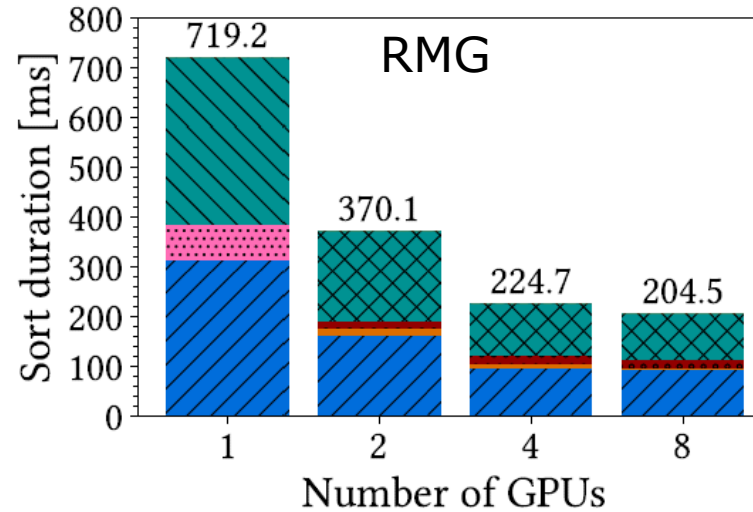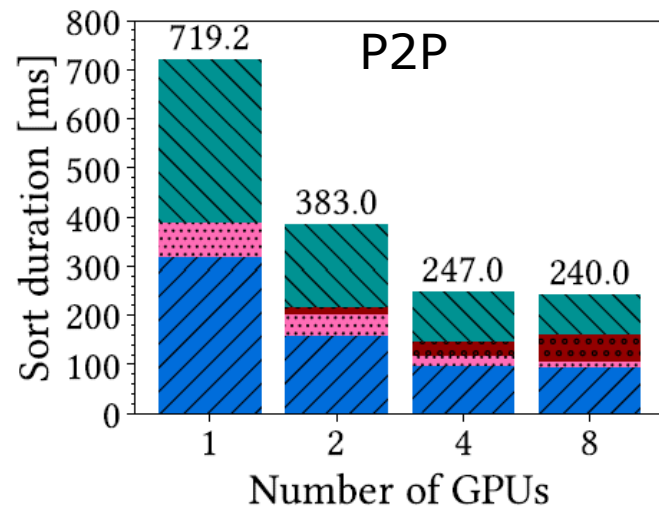
```
Bucket [1]:      00000001 00000000 00000000 00011010
                 ...
                 00000001 10101011 01011000 00110101
                 ...
                 00000001 11111111 11111111 11110101
```
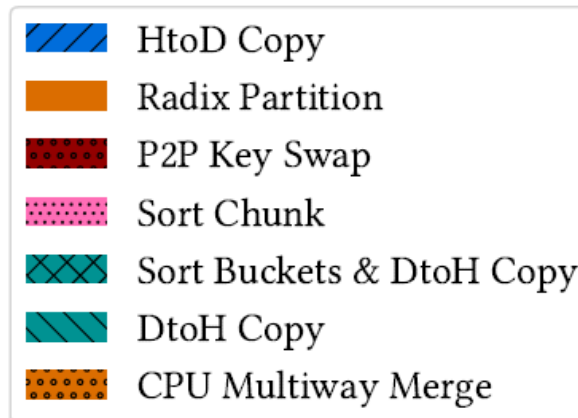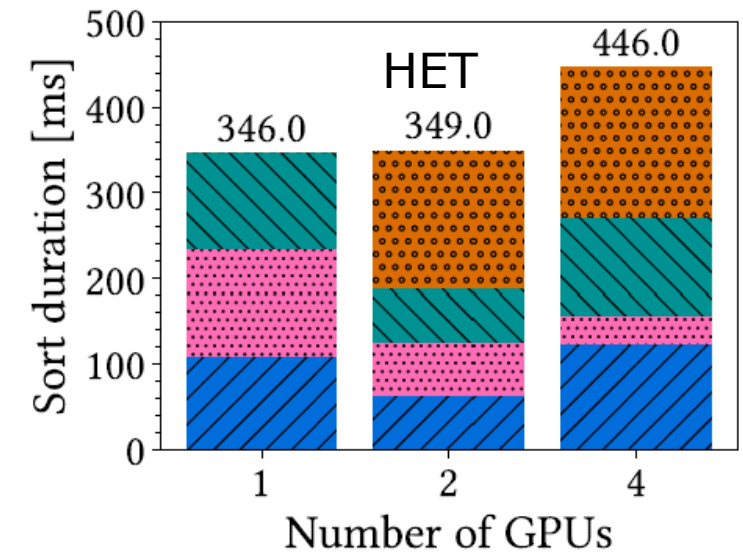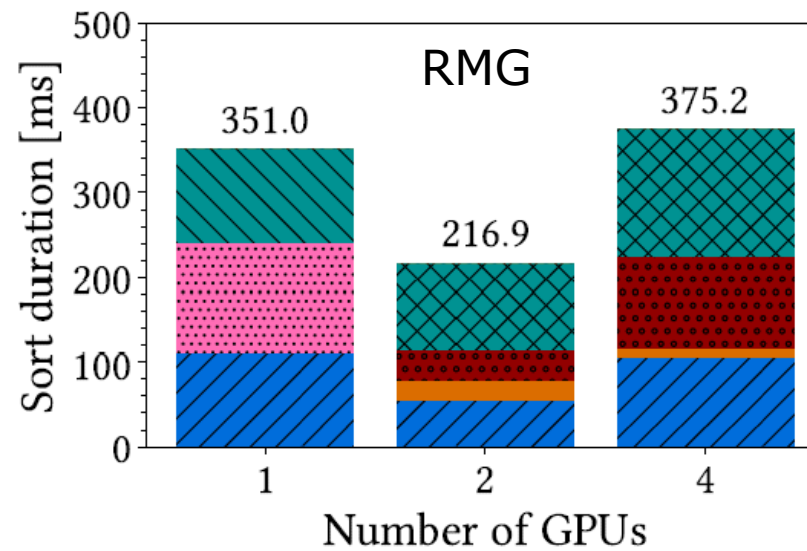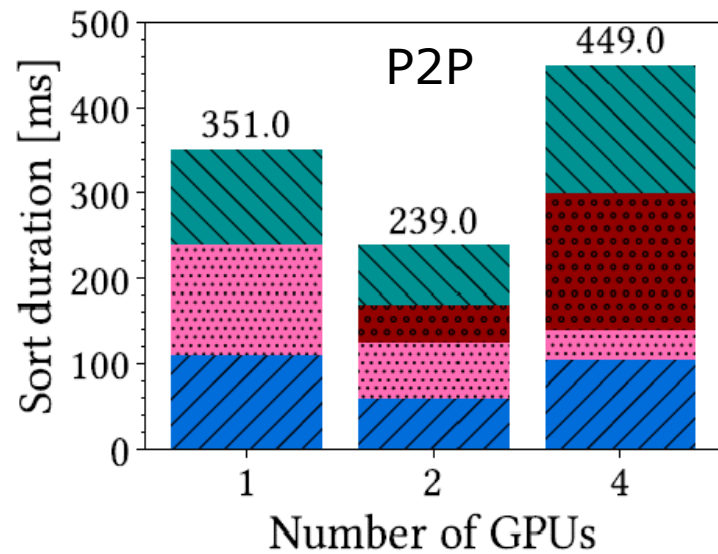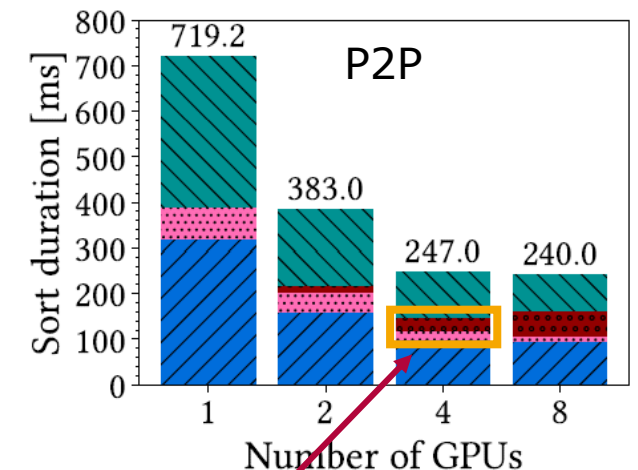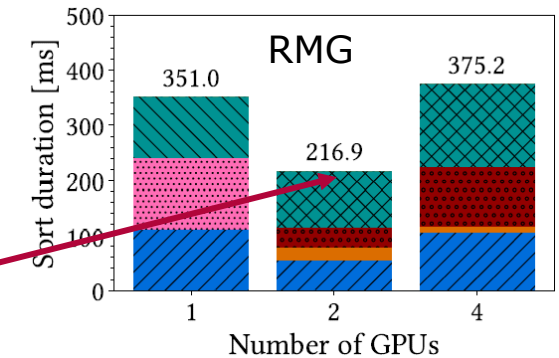
# Takeaways

- Scaling operators on Multi-GPU systems has limited benefits

- Often the best performance is achieved on 2 or 4 GPUs

  - **Overall fastest solution – 2 GPU mode on AC922**

- OLAP workloads rely on data transfers to GPU Memory

  - Short computation time – majority spent on data transfers

- Bound by the interconnect bandwidth

  - **DGX A100 -> P2P throughput: 279 GB/s, HtoD: 24 GB/s**

- Topology awareness and algorithm adaptation is essential



RMG



P2P

**Sorting + merging is only 10-15% of the total runtime**

# Summary

- GPU-accelerated data processing

- Scalable GPU Joins

- Multi-GPU Sorting

- Coming next: Multi-GPU Joins

- Questions?
  - tilmann.rabl@hpi.de