# Unlocking GPU potential with JIT

Anastasia Ailamaki
with Periklis Chrysogelos and Panos Sioulas

# One hardware does not fit all
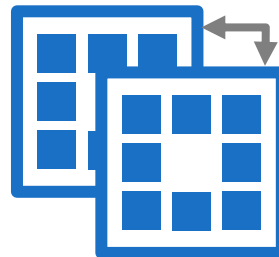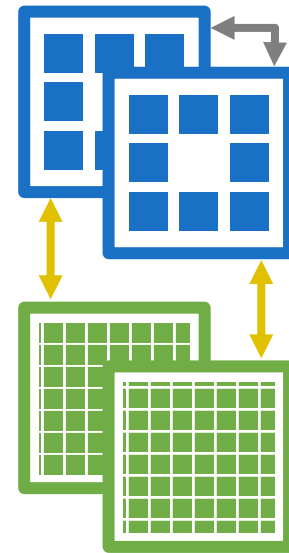


single-core
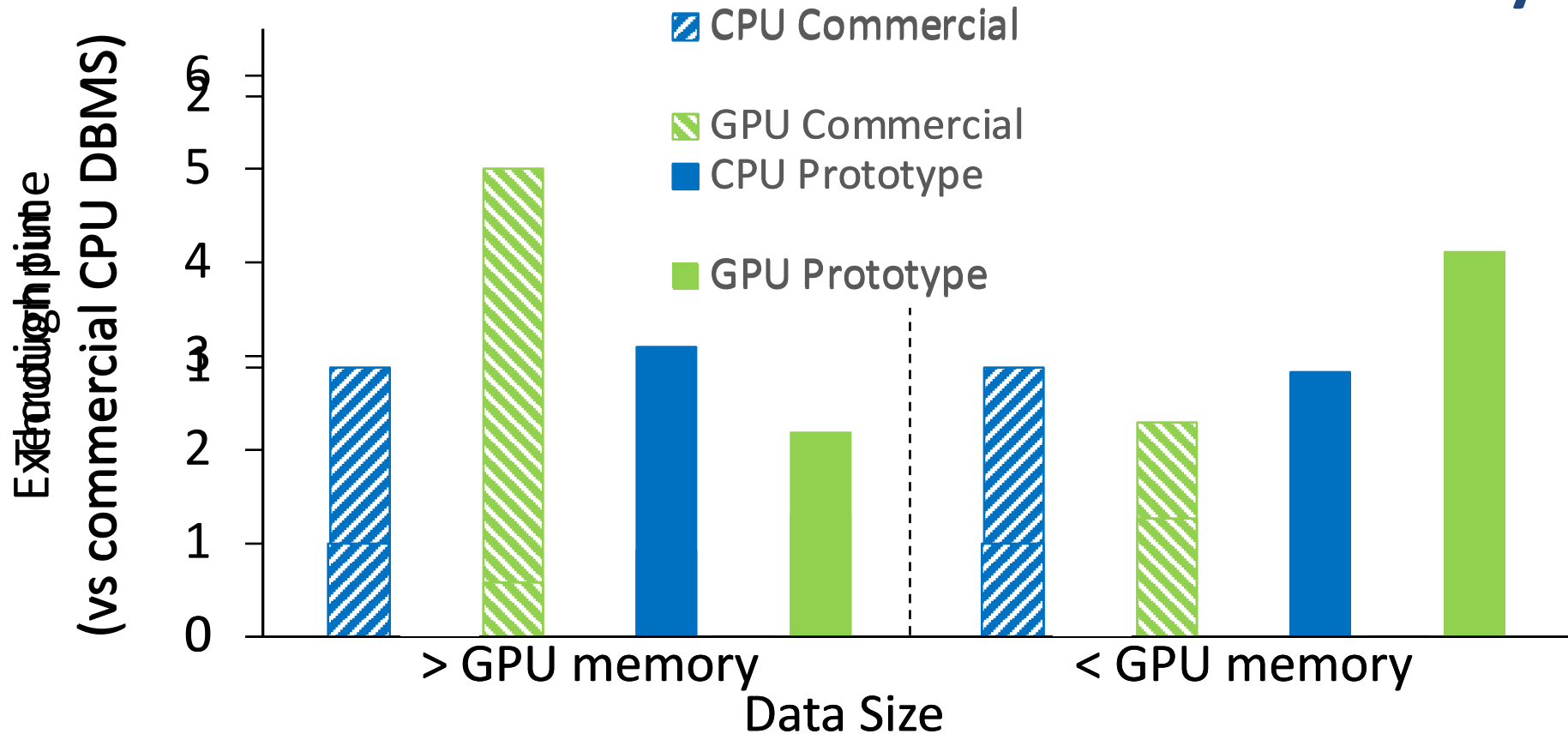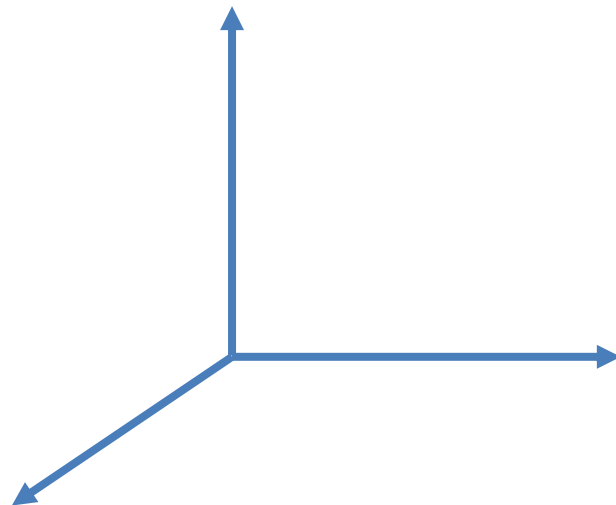single-CPU

multi-core
single-CPU

multi-core
multi-CPU

multi-core
multi-CPU
multi-GPU
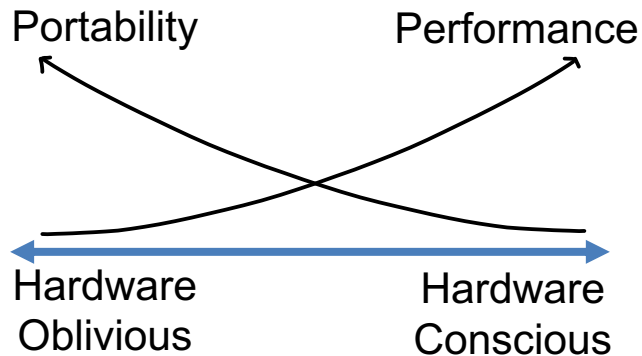
**Rethink query engines for accelerator-level parallelism**

# One hardware fits all: The end of an efficient story

Legend: CPU Commercial, GPU Commercial, CPU Prototype, GPU Prototype

Y-axis: Execution time / Throughput (vs commercial CPU DBMS)

X-axis categories: > GPU memory, < GPU memory — Data Size

**20%-80% throughput loss due to lack of portability**

# Designing query engines for heterogeneous HW

Portability          Performance

Hardware          Hardware
Oblivious          Conscious

**Decomposition of design space to find sweet spot**

# OLAP in heterogeneous servers: design space

[CIDR2019]

**intra-operator**

↳ Performance depends on μ-arch
↳ Tune operators to memory hierarchy specifics

**intra-device**

↳ Portability impacted by specialization
↳ Inject target-specific info using codegen

**inter-device**

↳ Limited device inter-operability
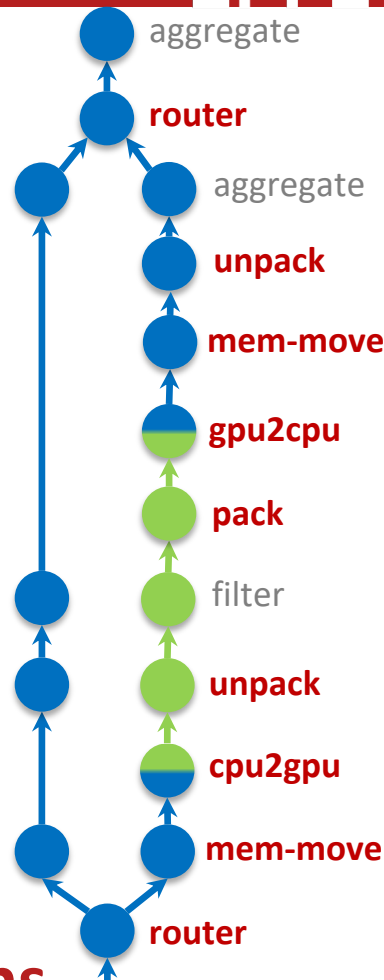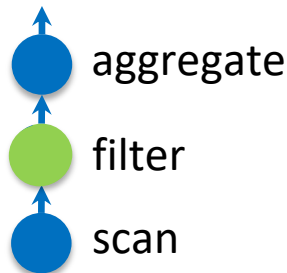↳ Encapsulate heterogeneity and balance load

## Selective obliviousness

# Inter-device: HetExchange

[VLDB2019]

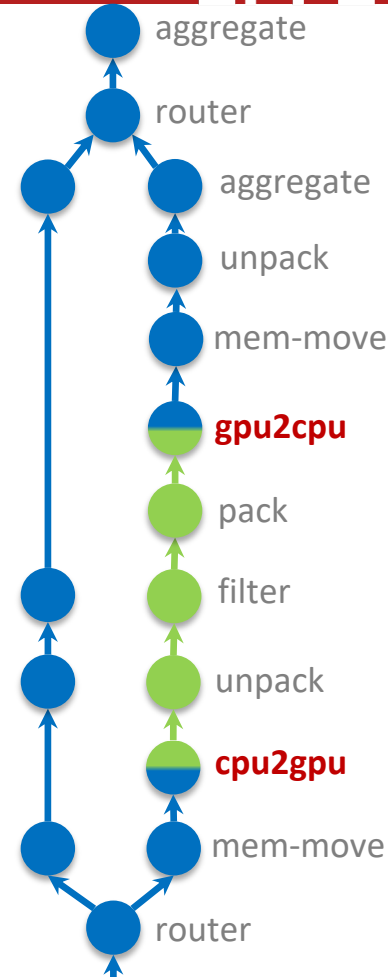- Decouple data- from control-flow
- Operators encapsulate trait conversions

| Flow | Scope | Trait |
|------|-------|-------|
| Control | Delegation | Heterogeneous Parallelism |
| | Routing | Homogeneous Parallelism |
| Data | Transfer | Data Locality |
| | Granularity | Execution Granularity |

aggregate

filter

scan

aggregate
router
aggregate
unpack
mem-move
gpu2cpu
pack
filter
unpack
cpu2gpu
mem-move
router

**Optimizer can produce cross-device plans**
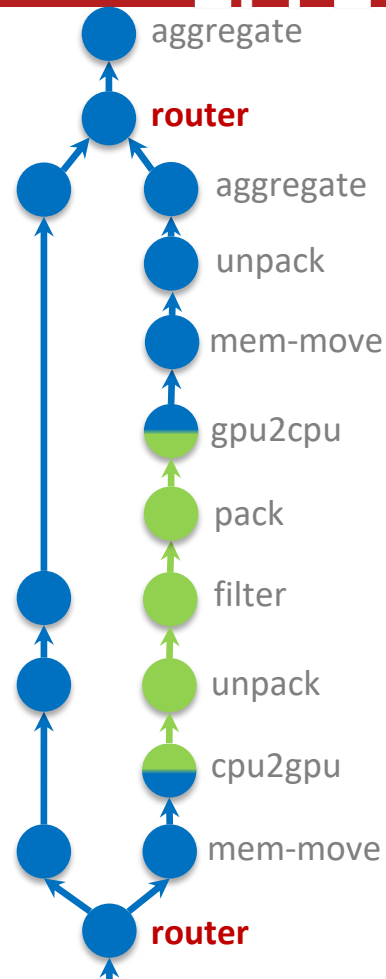
6

# Device Boundary Crossings

- Cross-device pipelined execution

- Hand-over execution to next device

- Launch kernels/threads, synchronize, backpressure

- Only operators aware of device heterogeneity

**Encapsulate heterogeneous parallelism**



aggregate

router

aggregate

unpack

mem-move

**gpu2cpu**

pack

filter

unpack

**cpu2gpu**

mem-move

router

# Concurrent Execution

- Horizontal & Vertical parallelism

- Instantiate pipelines multiple times

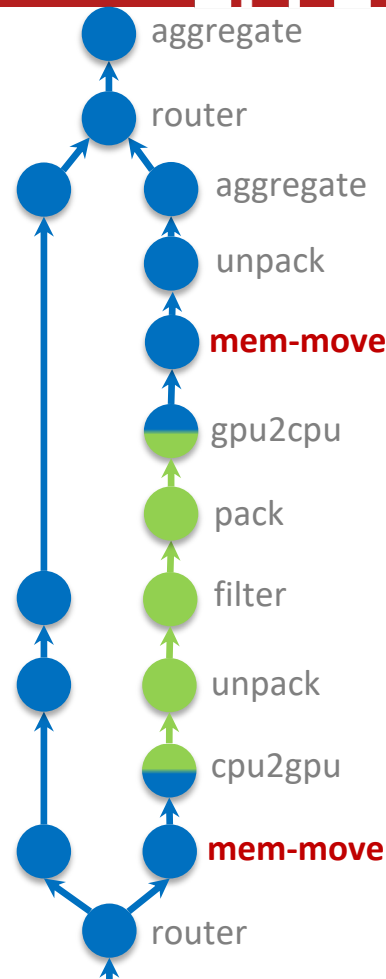- Routing policies: load-balance, partition, locality

**Encapsulate homogeneous parallelism**

aggregate

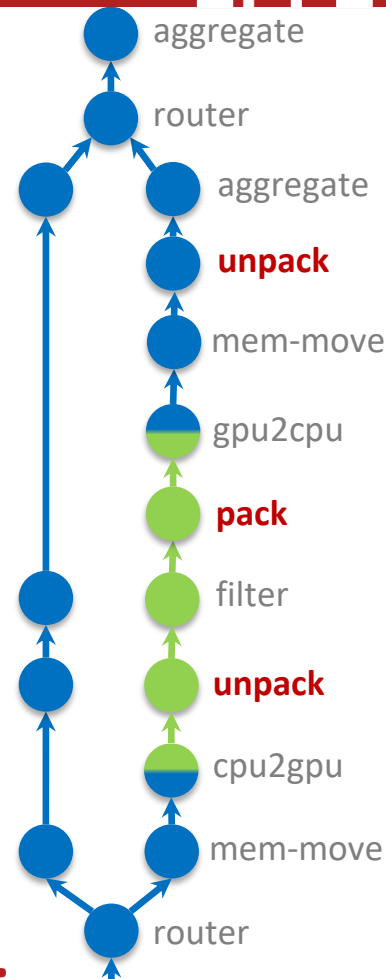**router**

aggregate

unpack

mem-move

gpu2cpu

pack

filter

unpack

cpu2gpu

mem-move

**router**

# Data Transfers

- Handle memory transfers/prefetching

- Hide memory topology

- Overlap transfers with execution

**Hide memory heterogeneity**

aggregate

router

aggregate

unpack

**mem-move**

gpu2cpu

pack

filter

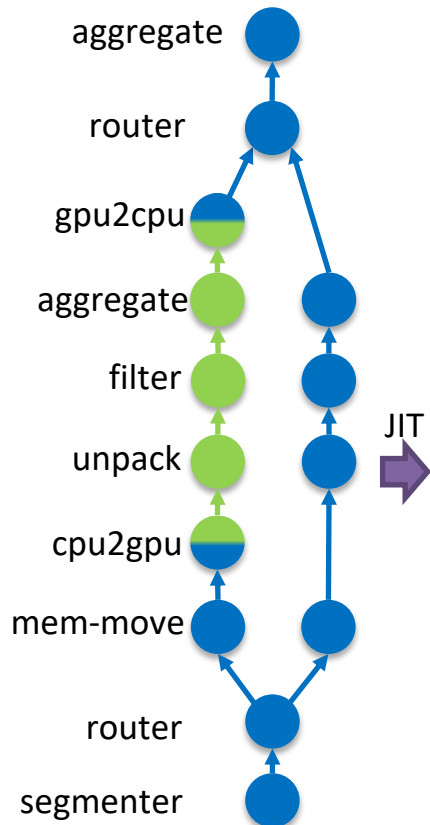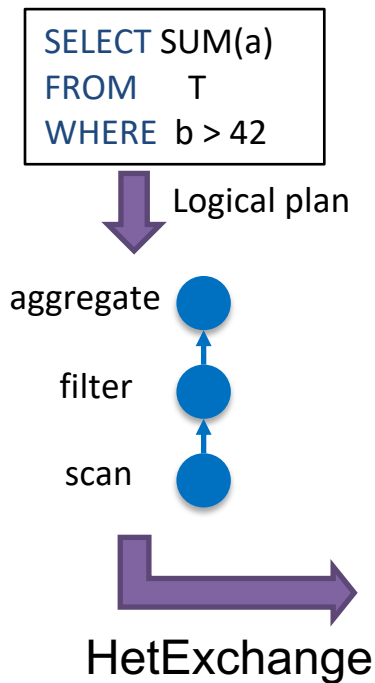unpack

cpu2gpu

**mem-move**

router

# Execution Granularity

- Processing: in-registers => tuple-at-a-time

- Memory transfers: packets => block-at-a-time

- Transition between execution granularities

- Create homogeneous (reg. policy) packets

**Transition between execution granularities**

aggregate

router

aggregate

**unpack**

mem-move

gpu2cpu

**pack**

filter

**unpack**

cpu2gpu

mem-move

router

# Heterogeneity-aware plans



```
SELECT SUM(a)
FROM    T
WHERE  b > 42
```

Logical plan

HetExchange

JIT

Efficiency
&
Operator portability

Generate pipelines and instantiate

# HetExchange in a JITed engine



aggregate

router

gpu2cpu

aggregate

filter

unpack

cpu2gpu

mem-move

router

segmenter

JIT

device crossing

9

device crossing

# Device providers

CPU
Provider

```
def unpack_filter_reduce(data_block, N,
    state)
local_acc ← 0
for i = threadIdInWorker to N − 1 with
    step #threadsInWorker
    t ← data_block[i]
    if t.a > 42
        local_acc ← local_acc + t.b
nh_acc ← neighborhood_reduce(local_acc
if thread neighborhood leader
    atomic_add(state.acc, nh_acc)
```
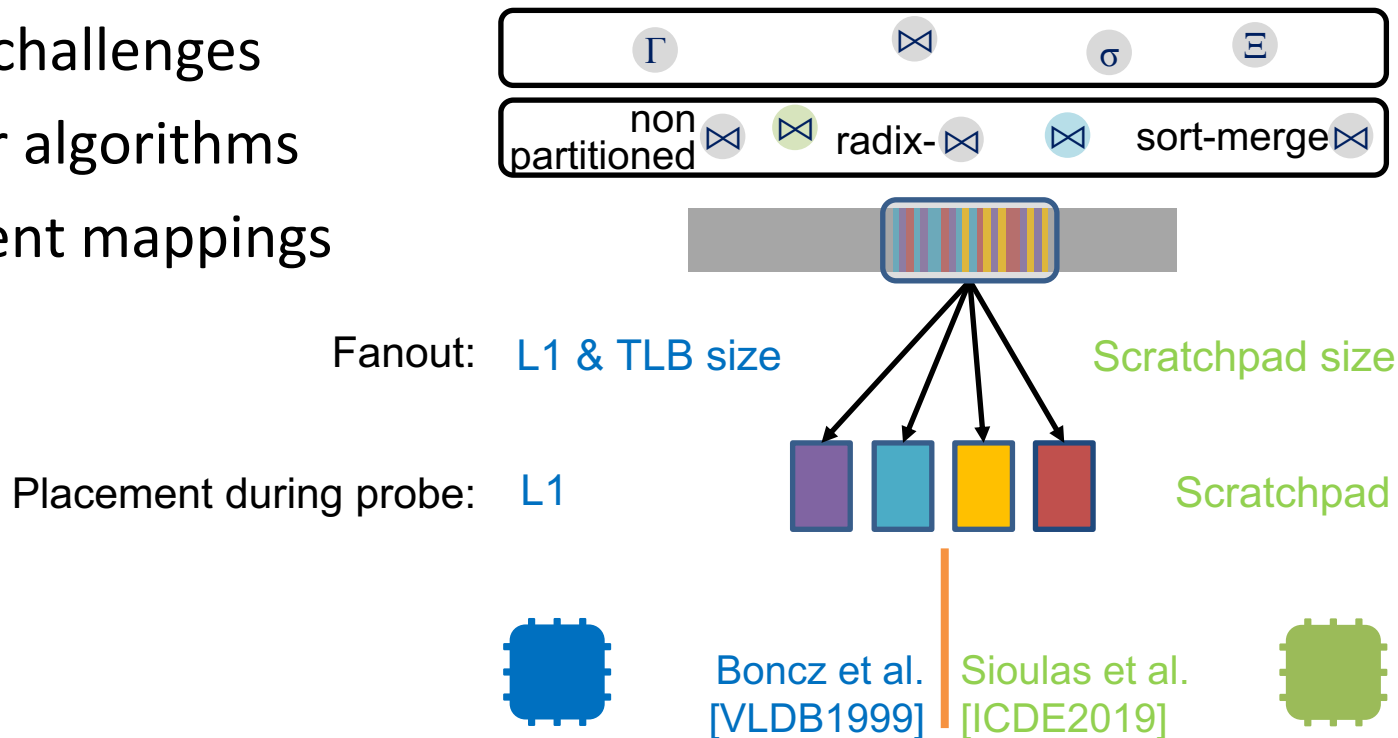
```
function ufr_cpu(data_block, N, state)
local_acc ← 0
for i = 0 to N − 1
    t ← data_block[i]
    if t.a > 42
        local_acc ← local_acc + t.b
nh_acc ← local_acc
state.acc ← state.acc + nh_acc
```

```
gpu_kernel ufr_gpu(data_block, N, state)
local_acc ← 0
for i = threadIdInGrid to N − 1 with
    step gridSize
    t ← data_block[i]
    if t.a > 42
        local_acc ← local_acc + t.b
nh_acc ← threadblock_reduce(local_acc)
if threadId = 0
    atomicAdd(state.acc, nh_acc)
```

GPU
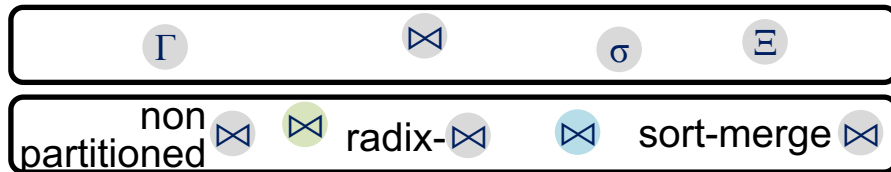Provider

**Inject target-specific info using the JIT infrastructure**

# Device-optimized operators

- Same challenges
- Similar algorithms
- Different mappings

$\Gamma$      $\bowtie$     $\sigma$     $\Xi$

non-partitioned $\bowtie$   $\bowtie$ radix-$\bowtie$   $\bowtie$ sort-merge $\bowtie$

Fanout:    L1 & TLB size      Scratchpad size

Placement during probe:    L1      Scratchpad

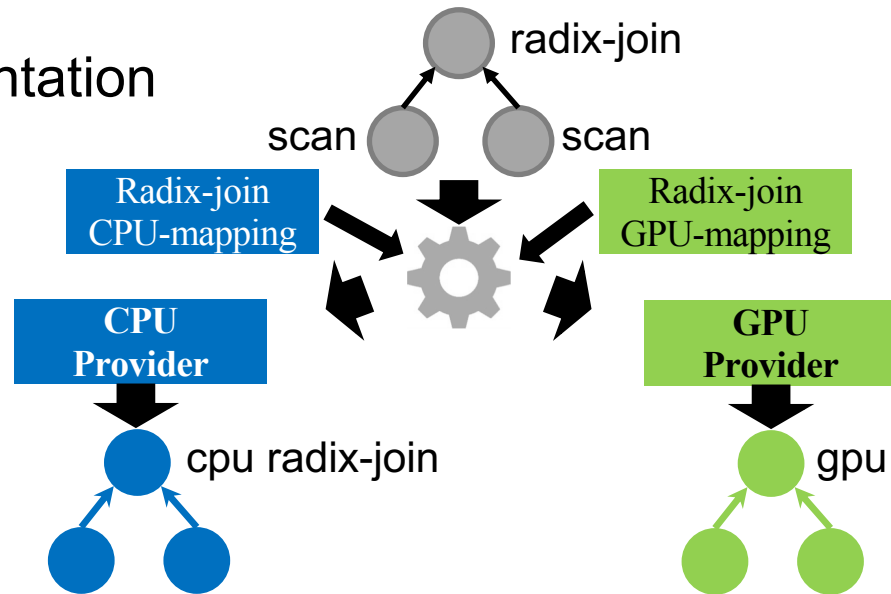Boncz et al. [VLDB1999] | Sioulas et al. [ICDE2019]

**Reuse algorithms, specialize mappings to hardware**

# Hardware-dependent JIT code



Hardware-aware algorithm

Device-independent implementation

Device-optimized implementation

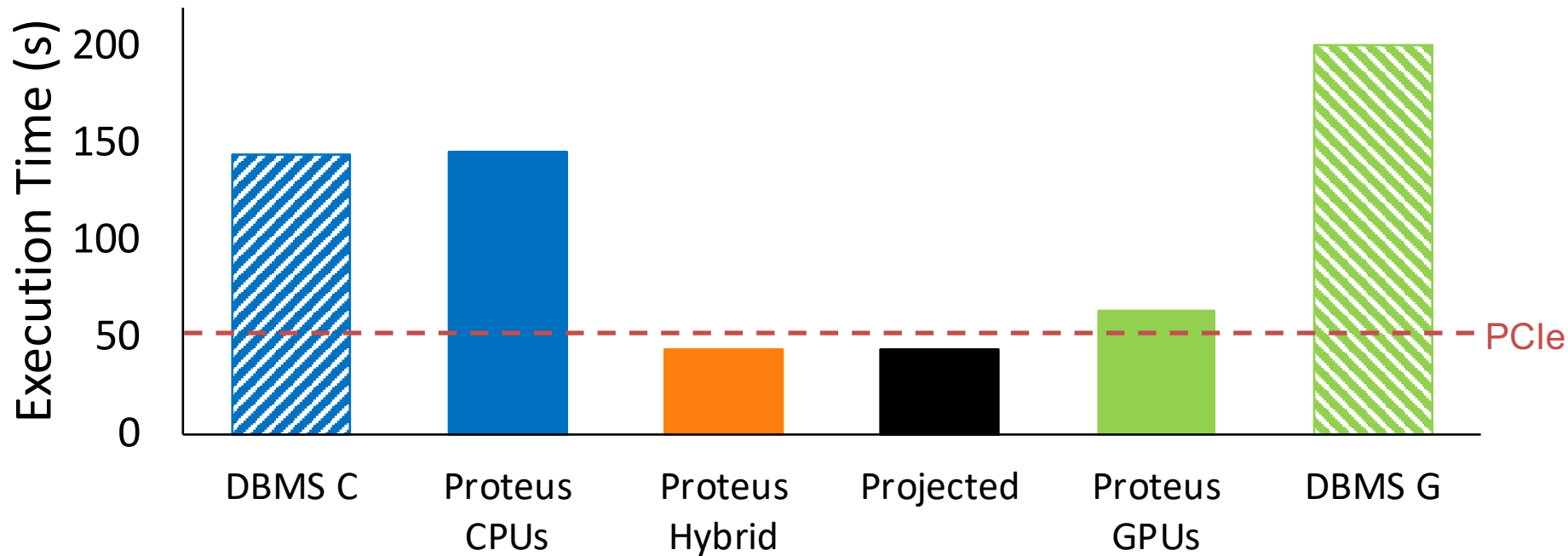**Lower generic description to device-specific code**

# Experimental Setup

- 2x Intel Xeon E5-2650L v3 12-core @ 1.80GHz, 256GB RAM

- 2x NVIDIA GeForce GTX1080, 8GB, PCIe3 x16 per GPU

- DBMS C/G: state-of-the-art commercial DBMS
  - DBMS C: CPU-based, vector-at-a-time, SIMD, based on MonetDB/X100
  - DBMS G: GPU-based, JIT engine

# Performance on CPU-resident data

SSB SF1000, 600GB CSV
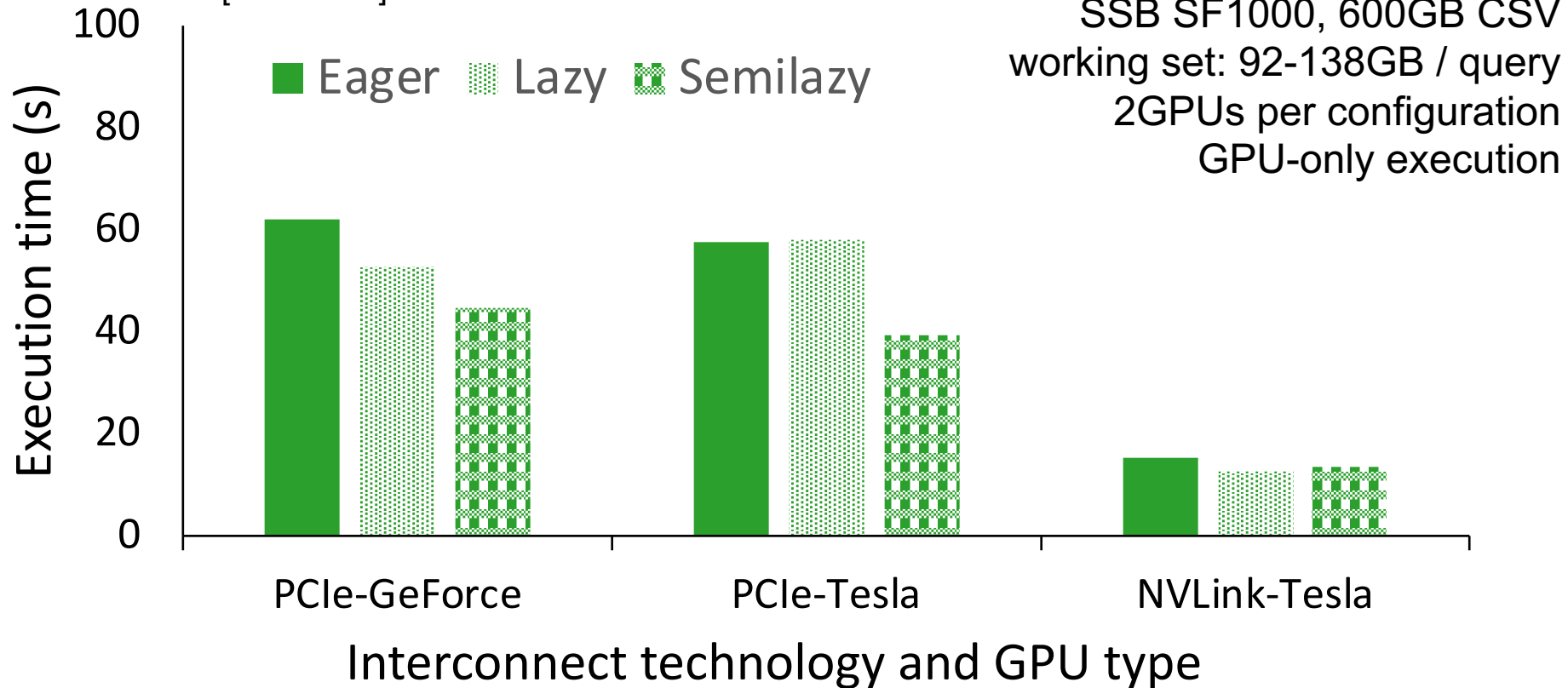working set: 92-138GB / query



**Hybrid throughput = 88.5% (CPU-only + GPU-only), on average**

# A glimpse into the future

- Effect of interconnects and GPU compute power
- Access to high-throughput network

# Adapting access method to query

[CIDR2020]

SSB SF1000, 600GB CSV
working set: 92-138GB / query
2GPUs per configuration
GPU-only execution

**Execution time (s)** (y-axis: 0, 20, 40, 60, 80, 100)

Legend: ■ Eager ▦ Lazy ▦ Semilazy

X-axis categories: PCIe-GeForce, PCIe-Tesla, NVLink-Tesla
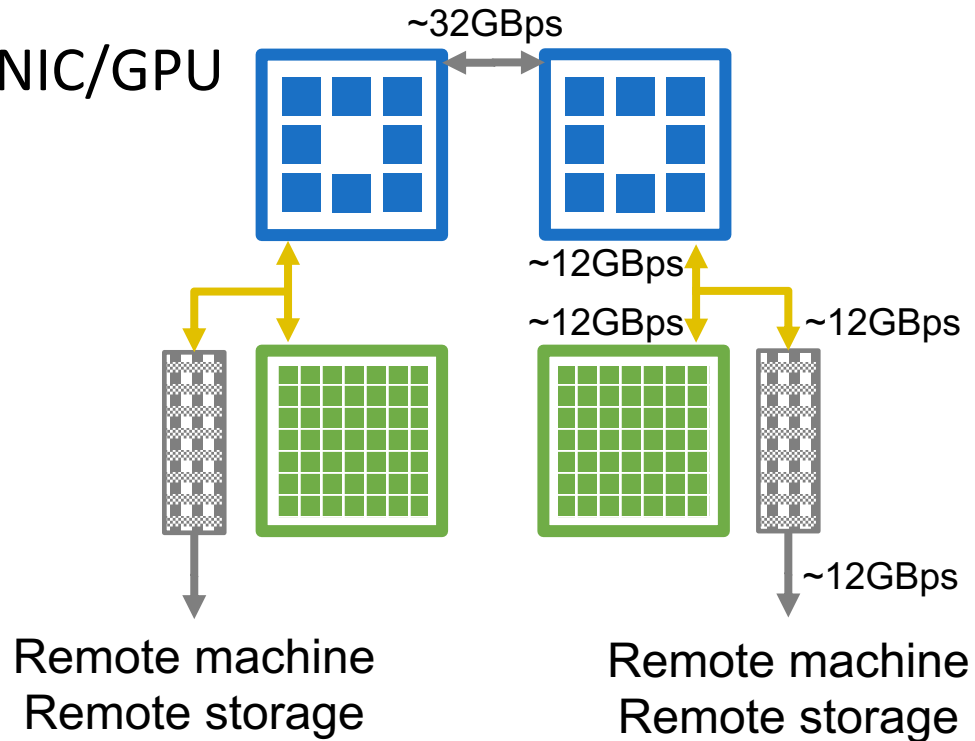
**Interconnect technology and GPU type**

**Up to 45% speed-up by tuning access method to hardware**

# Towards placing the CPU on the side

- Shared & limited PCIe buses to NIC/GPU

- Similar intra/inter-server BW

- Direct NIC-GPU access (RDMA)

~32GBps

~12GBps

~12GBps

~12GBps

~12GBps

Remote machine
Remote storage

Remote machine
Remote storage

**Avoid CPU bottleneck => Device-centric OLAP engines**

# JIT unleashes ALP

- Run on all available devices

- Relational operators oblivious to heterogeneity

- Fast: Inject target-specific information through codegen

- Result: 5x-10x versus CPU-/GPU-specialized systems

# Thank you!