# Crystal GPU Database
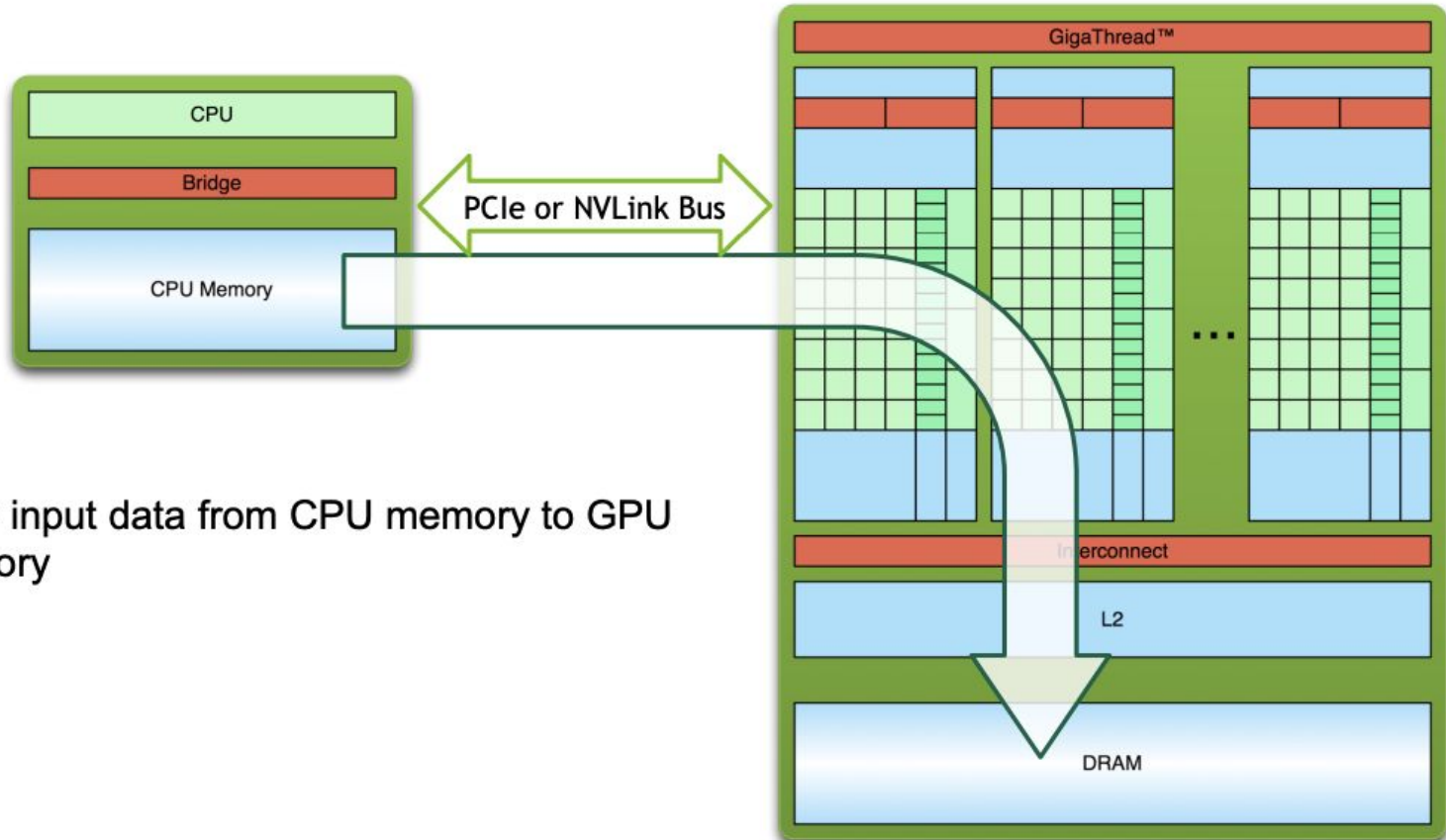
By Devesh
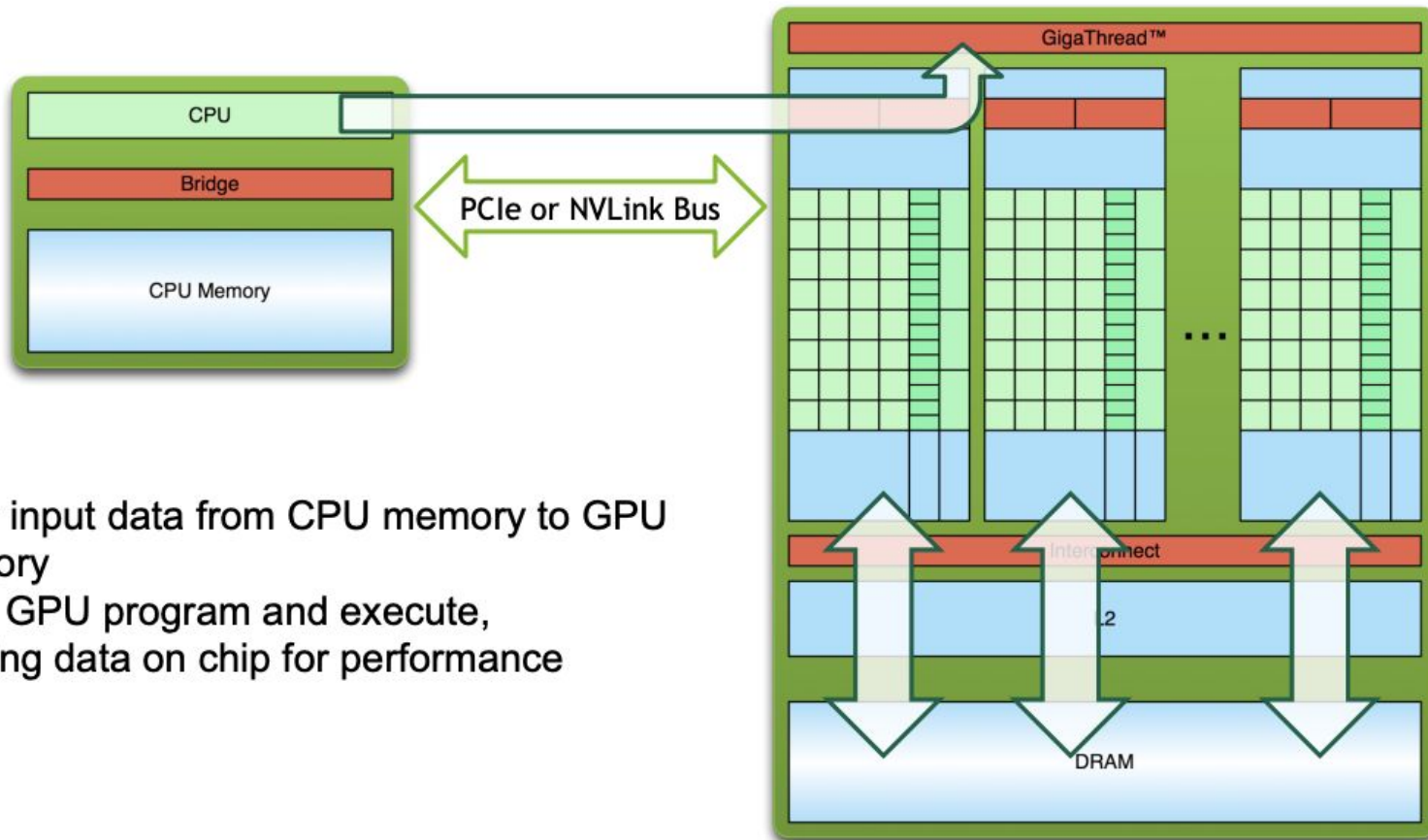
# Intro to GPUS

# SIMPLE PROCESSING FLOW
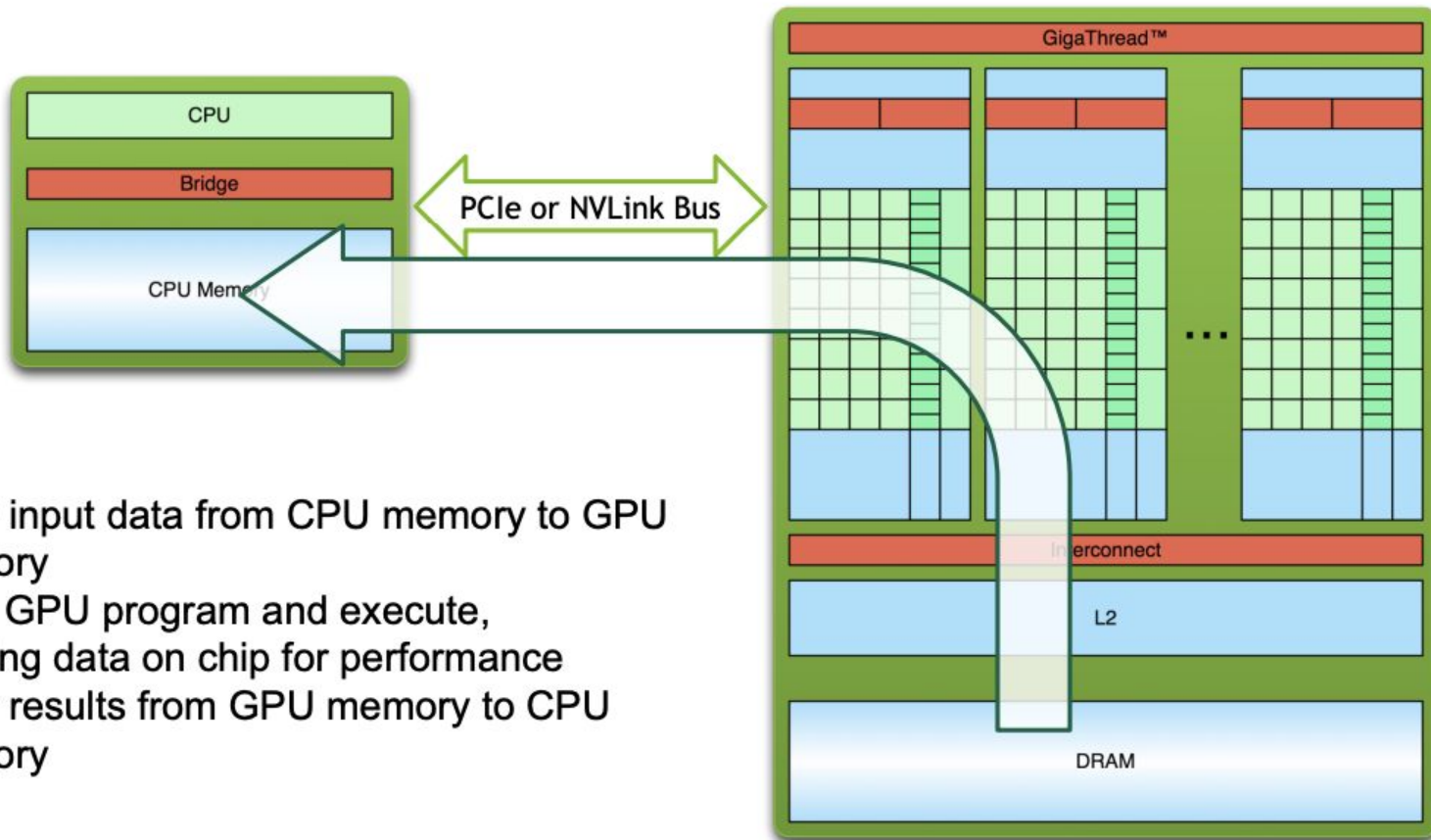


1. Copy input data from CPU memory to GPU memory

PCIe or NVLink Bus

CPU

Bridge

CPU Memory

GigaThread™

Interconnect

L2

DRAM

# SIMPLE PROCESSING FLOW

CPU

Bridge

PCIe or NVLink Bus

CPU Memory

GigaThread™

Interconnect

DRAM

1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance

7

NVIDIA.

# SIMPLE PROCESSING FLOW

CPU

Bridge

PCIe or NVLink Bus

CPU Memory

GigaThread™

Interconnect

L2

DRAM
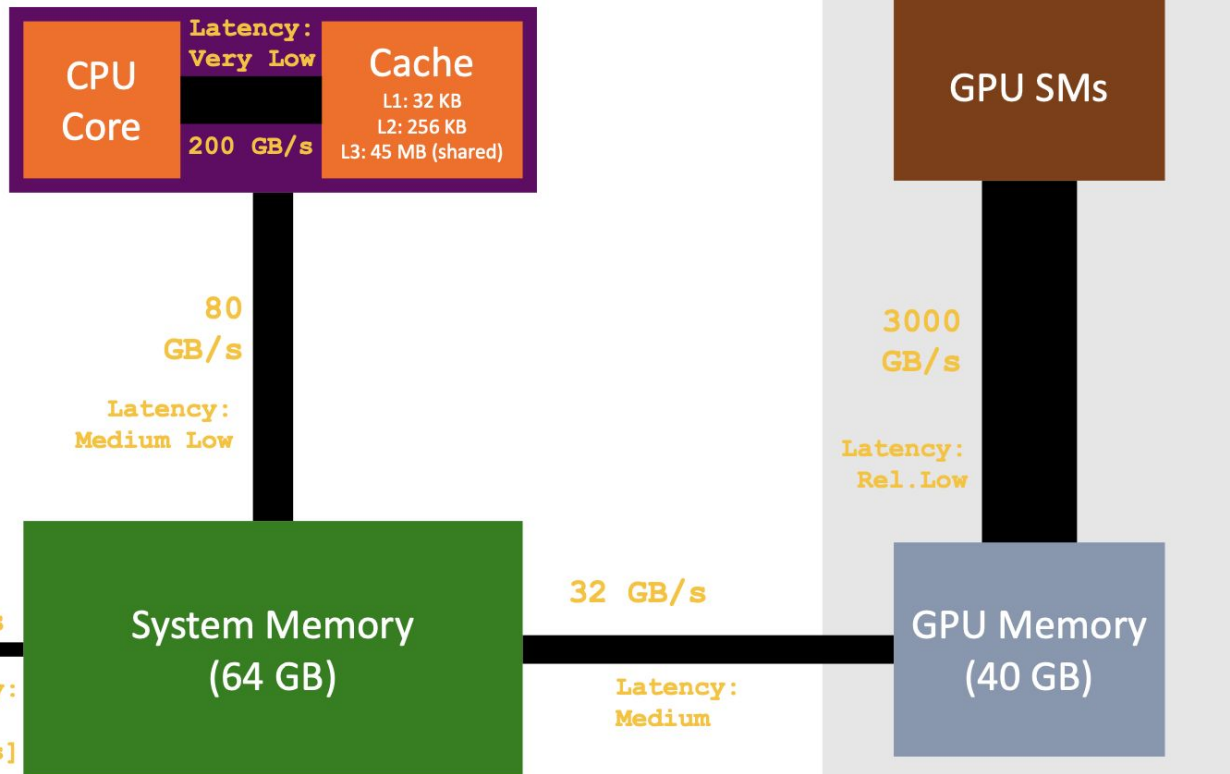
1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance
3. Copy results from GPU memory to CPU memory
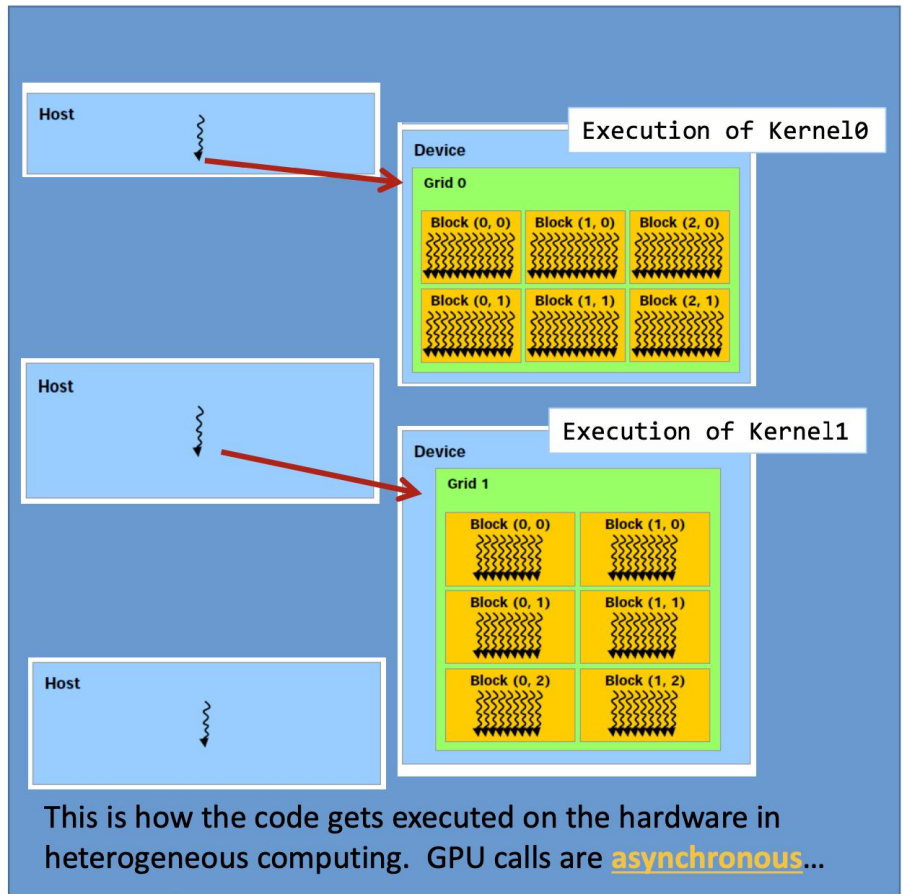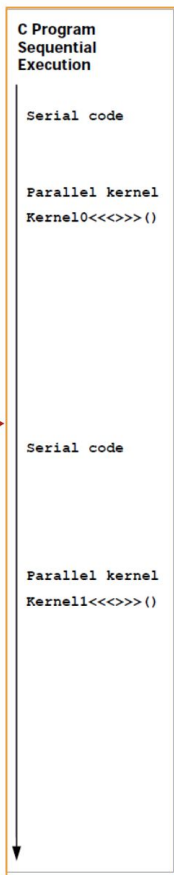
# Latencies, Bandwidths, and Limits

**CPU Core**

**Latency: Very Low**

**Cache**
L1: 32 KB
L2: 256 KB
L3: 45 MB (shared)

**200 GB/s**

**GPU SMs**

NOTE: The **width** of the black lines is proportional to the bandwidth.

**80 GB/s**

Latency: Medium Low

**3000 GB/s**

Latency: Rel.Low

Infiniband to Next Node

**6GB/s**

Latency: High [4-6 us]

**System Memory (64 GB)**

**32 GB/s**

Latency: Medium

**GPU Memory (40 GB)**

# It's all about threads



C Program
Sequential
Execution

Serial code

Parallel kernel
Kernel0<<<>>>()

Serial code

Parallel kernel
Kernel1<<<>>>()

This is how your
C code looks like

Host

Host

Host

**Execution of Kernel0**

Device

Grid 0

Block (0, 0)  Block (1, 0)  Block (2, 0)

Block (0, 1)  Block (1, 1)  Block (2, 1)

**Execution of Kernel1**

Device

Grid 1

Block (0, 0)  Block (1, 0)

Block (0, 1)  Block (1, 1)

Block (0, 2)  Block (1, 2)

This is how the code gets executed on the hardware in
heterogeneous computing.  GPU calls are asynchronous...

# Executing the blocks
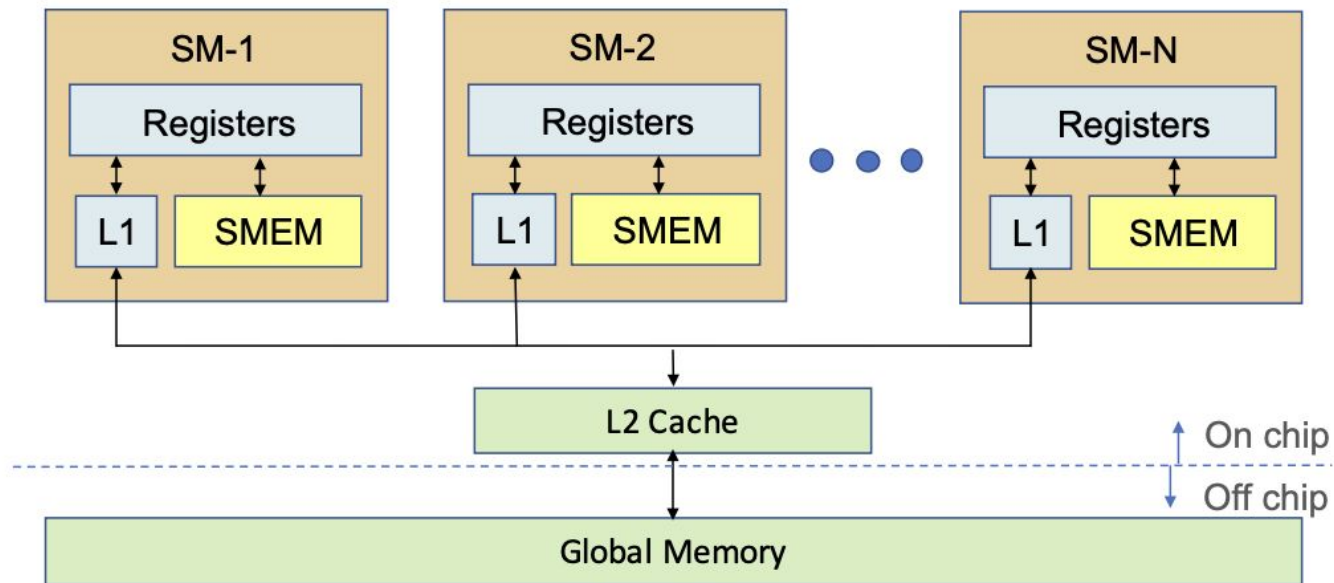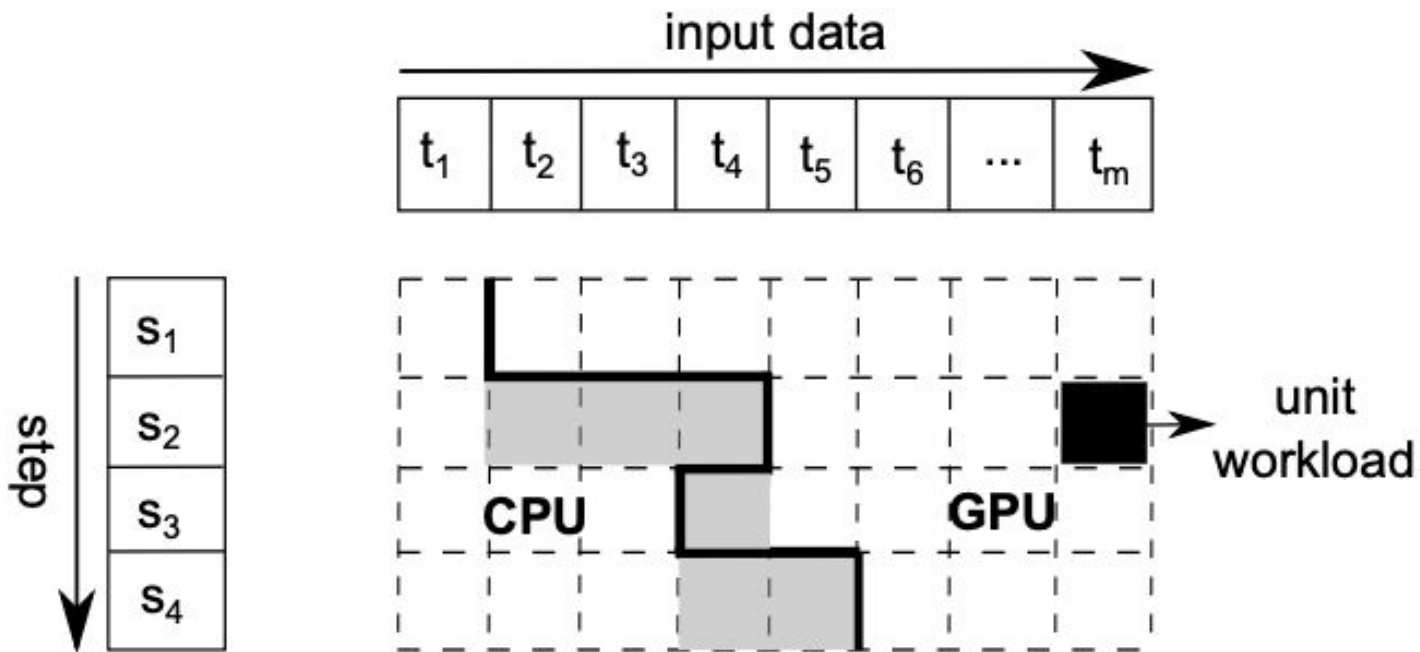
# Memory Hierarchy



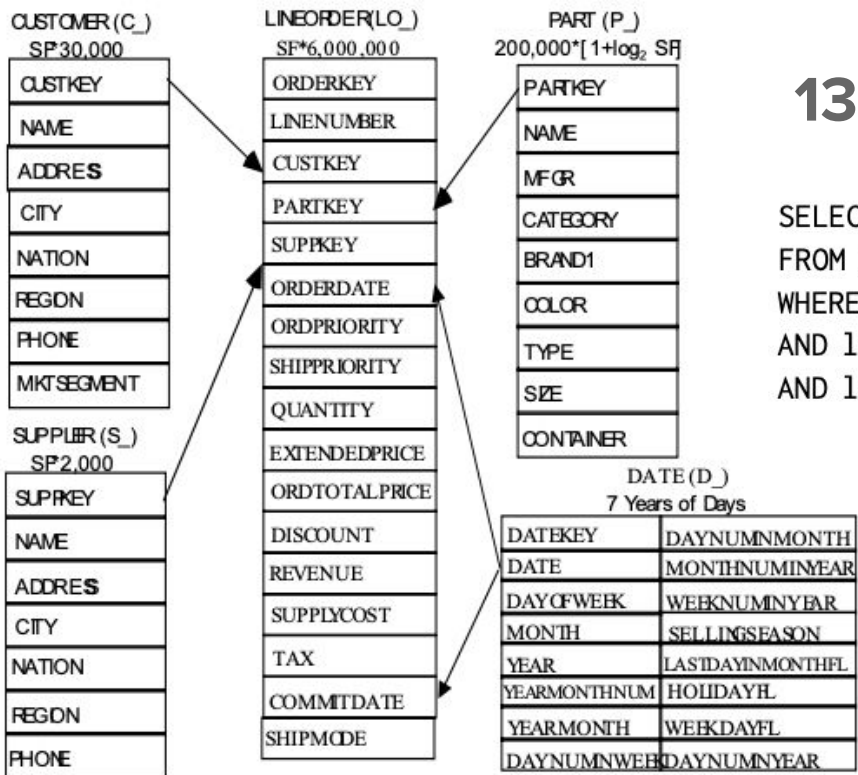**Figure 1: GPU Memory Hierarchy**

# Or as the Mythbusters explained

# Previous Works

# GPU as a Co Processor



**Data resides in CPU memory and is moved to the GPU during query execution**
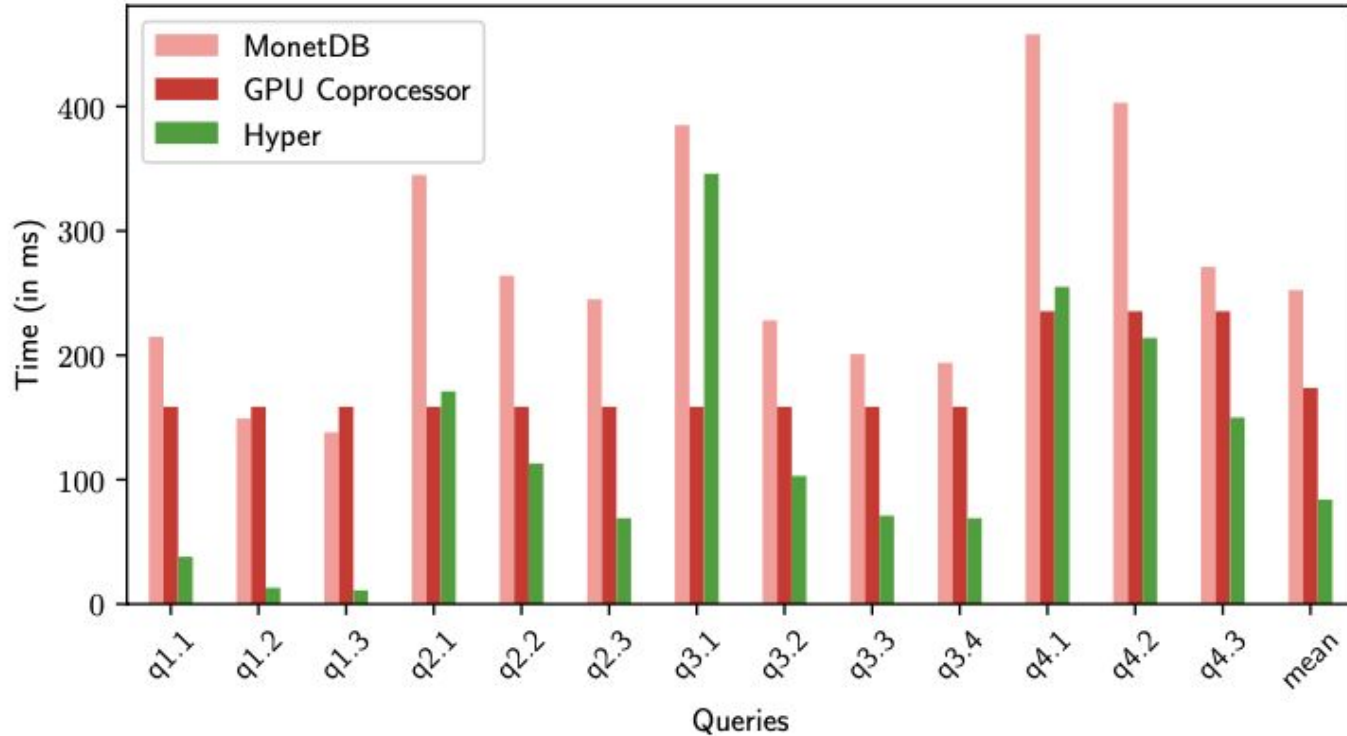
# SSB Benchmark

## CUSTOMER (C_)
### SP 30,000
| |
|---|
| CUSTKEY |
| NAME |
| ADDRESS |
| CITY |
| NATION |
| REGION |
| PHONE |
| MKTSEGMENT |

## SUPPLIER (S_)
### SP 2,000
| |
|---|
| SUPPKEY |
| NAME |
| ADDRESS |
| CITY |
| NATION |
| REGION |
| PHONE |

## LINEORDER(LO_)
### SF*6,000,000
| |
|---|
| ORDERKEY |
| LINENUMBER |
| CUSTKEY |
| PARTKEY |
| SUPPKEY |
| ORDERDATE |
| ORDPRIORITY |
| SHIPPRIORITY |
| QUANTITY |
| EXTENDEDPRICE |
| ORDTOTALPRICE |
| DISCOUNT |
| REVENUE |
| SUPPLYCOST |
| TAX |
| COMMITDATE |
| SHIPMODE |

## PART (P_)
### $200,000*[1+\log_2 SF]$
| |
|---|
| PARTKEY |
| NAME |
| MFGR |
| CATEGORY |
| BRAND1 |
| COLOR |
| TYPE |
| SIZE |
| CONTAINER |

## DATE (D_)
### 7 Years of Days
| | |
|---|---|
| DATEKEY | DAYNUMNMONTH |
| DATE | MONTHNUMINYEAR |
| DAYOFWEEK | WEEKNUMINYEAR |
| MONTH | SELLINGSEASON |
| YEAR | LASTDAYINMONTHFL |
| YEARMONTHNUM | HOLIDAYFL |
| YEARMONTH | WEEKDAYFL |
| DAYNUMNWEEK | DAYNUMNYEAR |

## 13 different queries like this

```
SELECT SUM(lo_extendedprice * lo_discount) AS revenue
FROM lineorder
WHERE lo_quantity < 25
AND lo_orderdate >= 19930101 AND lo_orderdate <= 19940101
AND lo_discount >= 1 AND lo_discount <= 3;
```

## Using SF = 20

# Comparison against other OLAP databases



**On average GPU co processing is 1.4x times slower than Hyper**

Crystal

# Solution

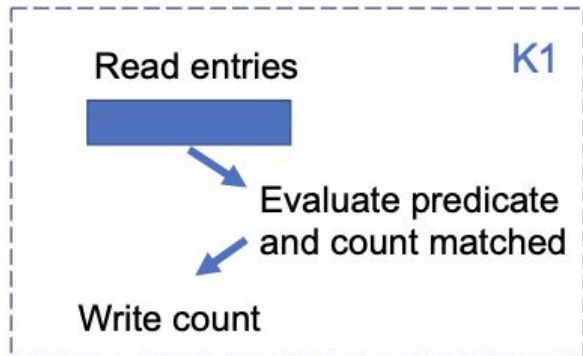Store the working set in the GPU memory itself as done by



Crystal builds on these using a tile based execution model

# Problem with current GPU approach

**Motivating query**

```
Q0: SELECT y FROM R WHERE y > v;
```
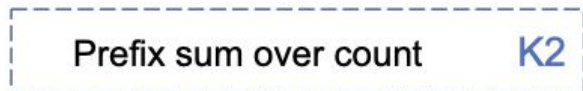
# How current systems do it



"Each thread strides through the dataset and determines number of matching rows"
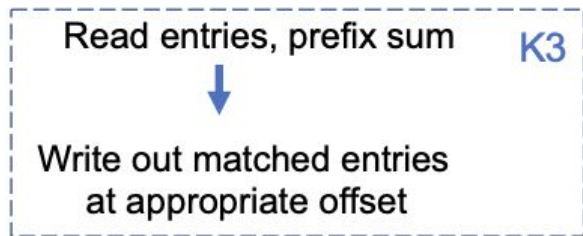
If 100 total threads then
Thread 0 checks rows 0, 100, 200, ...
Thread 1 checks rows 1, 101, 201, ...

Example count arr: [1, 3, 4, 2, 5]

Prefix sum: [0, 1, 4, 8, 10, 15]
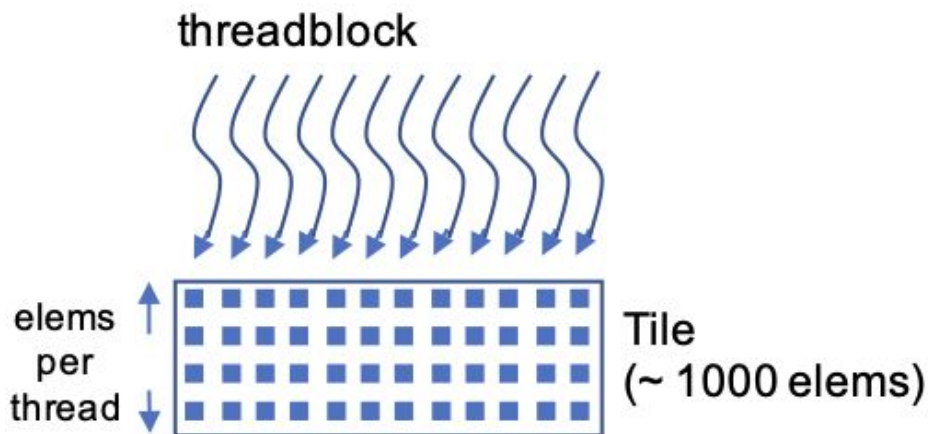
Thread 0 will write result in indices [0, 1)
Thread 1 will write result in indices [1, 4)
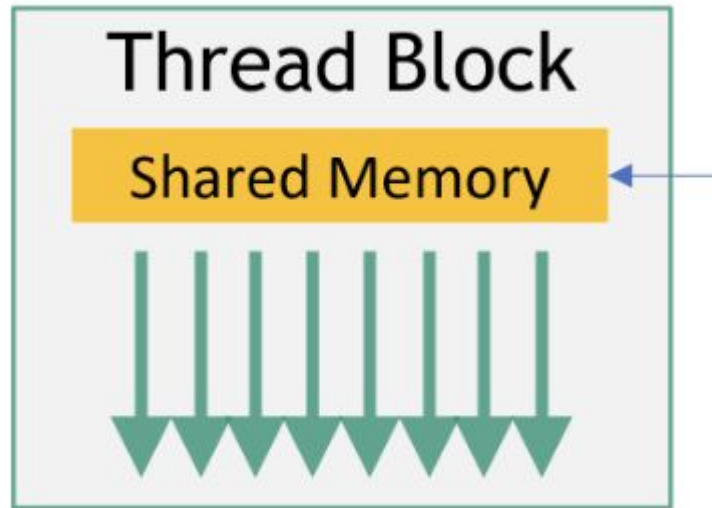Thread 2 will write result in [4, 8)
...

**Problem: Requires 3 kernels and 2 iterations over the dataset**
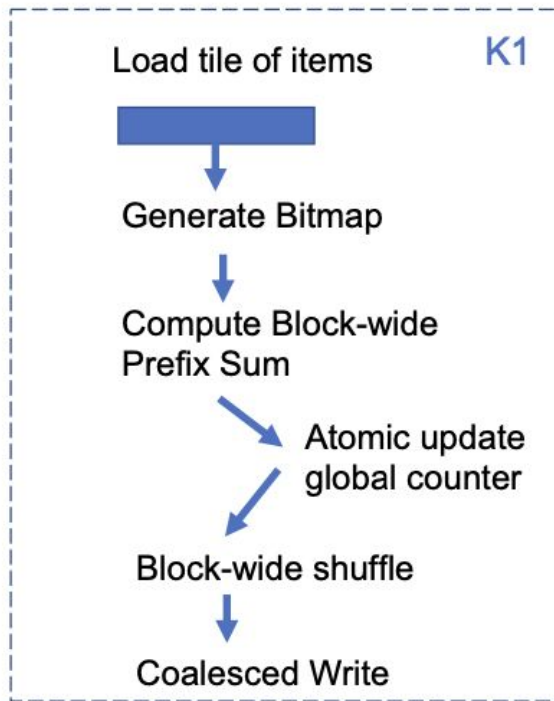
# Introducing Tile Based Execution Model



threadblock

elems per thread

Tile (~ 1000 elems)

Thread Block

Shared Memory

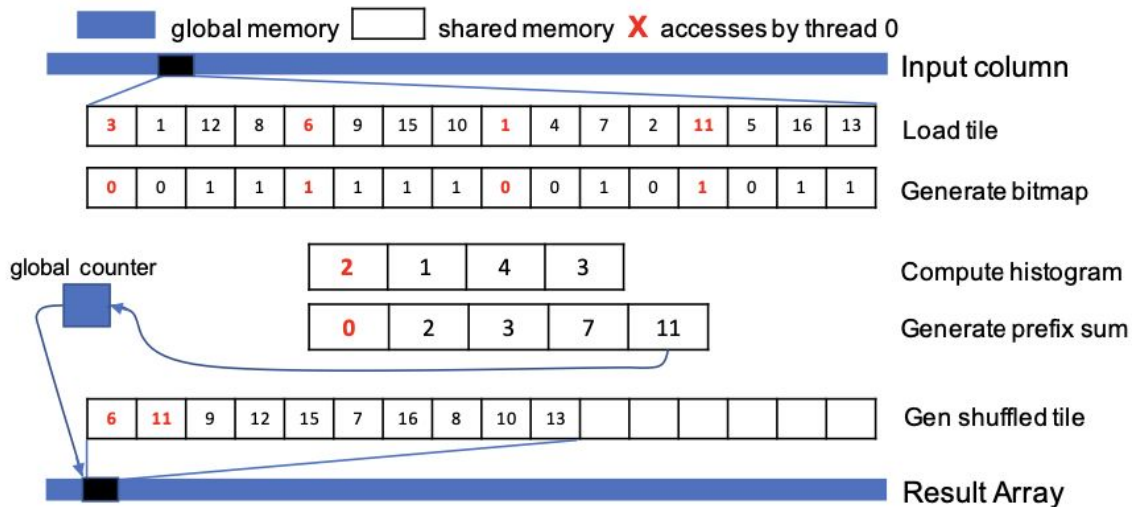**Tiles rather than threads are basic units of execution**

**Load all of tile's data into shared memory once and reuse it**

# Going back to the example



(b) With Tile-based processing
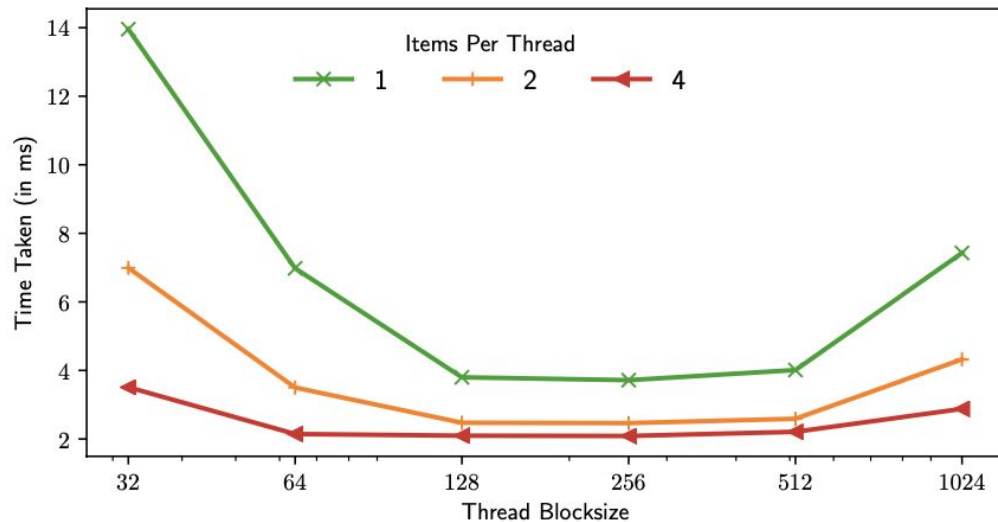
**SELECT Y FROM R WHERE Y > 5**

# Converting this to code

| Primitive | Description |
|---|---|
| BlockLoad | Copies a tile of items from global memory to shared memory. Uses vector instructions to load full tiles. |
| BlockLoadSel | Selectively load a tile of items from global memory to shared memory based on a bitmap. |
| BlockStore | Copies a tile of items in shared memory to device memory. |
| BlockPred | Applies a predicate to a tile of items and stores the result in a bitmap array. |
| BlockScan | Co-operatively computes prefix sum across the block. Also returns sum of all entries. |
| BlockShuffle | Uses the thread offsets along with a bitmap to locally rearrange a tile to create a contiguous array of matched entries. |
| BlockLookup | Returns matching entries from a hash table for a tile of keys. |
| BlockAggregate | Uses hierarchical reduction to compute local aggregate for a tile of items. |

## Crystal provides primitives for each of these steps

# Tuning some key parameters

- **Thread Block Size: Number of threads per thread block**
- **Items per thread: Number of rows each thread will process**



- **R has 2^29 rows**
- **Selectivity factor: 0.5**

# Projection Queries

Q1:  SELECT  $ax_1 + bx_2$  FROM  R;
Q2:  SELECT  $\sigma(ax_1 + bx_2)$  FROM  R;

**Note: Q2 is performing a sigmoid operation**

# Will Crystal against

- CPU based multi threaded implementation

- CPU-Opt: CPU + Non temporal writes (write out cache line to main memory) + SIMD optimizations

- Modeling: Calculating expected runtimes given hardware specifications

# How does modeling work

**For query** $\text{SELECT } ax_1 + bx_2 \text{ FROM R;}$

$$runtime = \frac{2 \times 4 \times N}{B_r} + \frac{4 \times N}{B_w}$$

**Time to load both columns**

**Time to write result back**

- N = # of rows

- Br = Read bandwidth

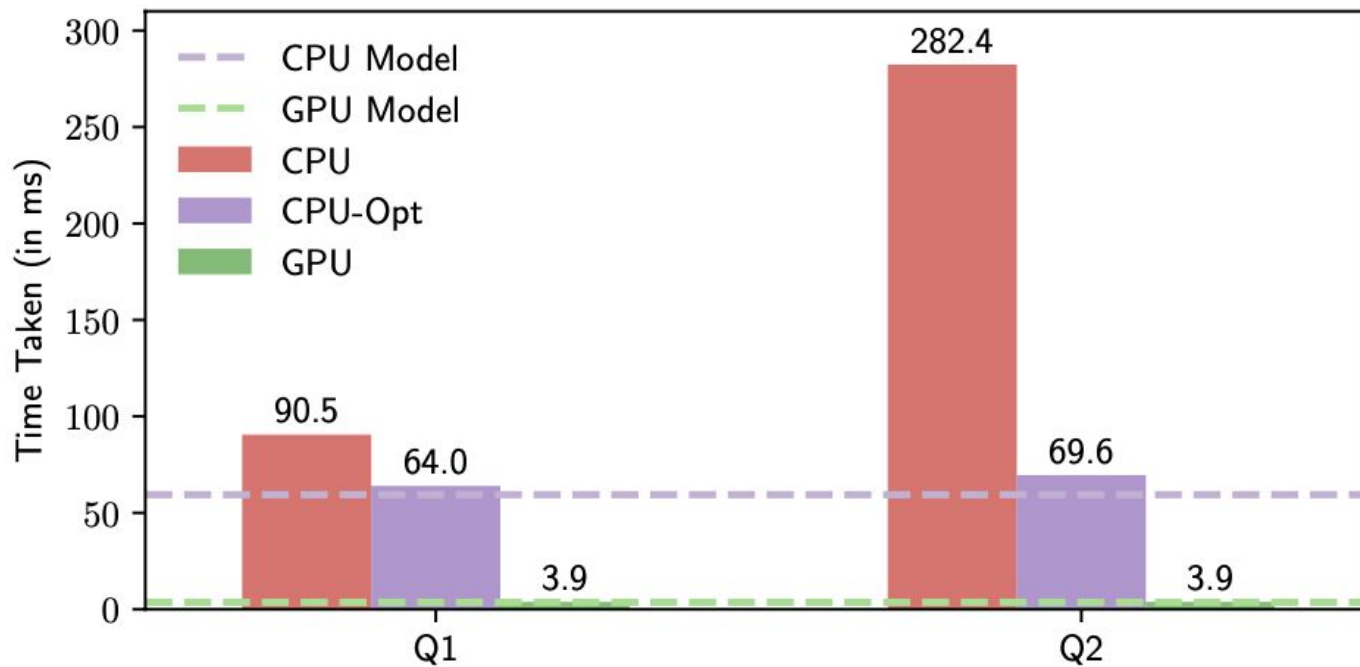- Bw = Write bandwidth

# Projection Performance



Table with 2^29 rows

# Selection

Q3: SELECT y FROM R WHERE y < v;

**Two possible approaches:**

for each y in R:
  if y > v:
    output[i++] = v
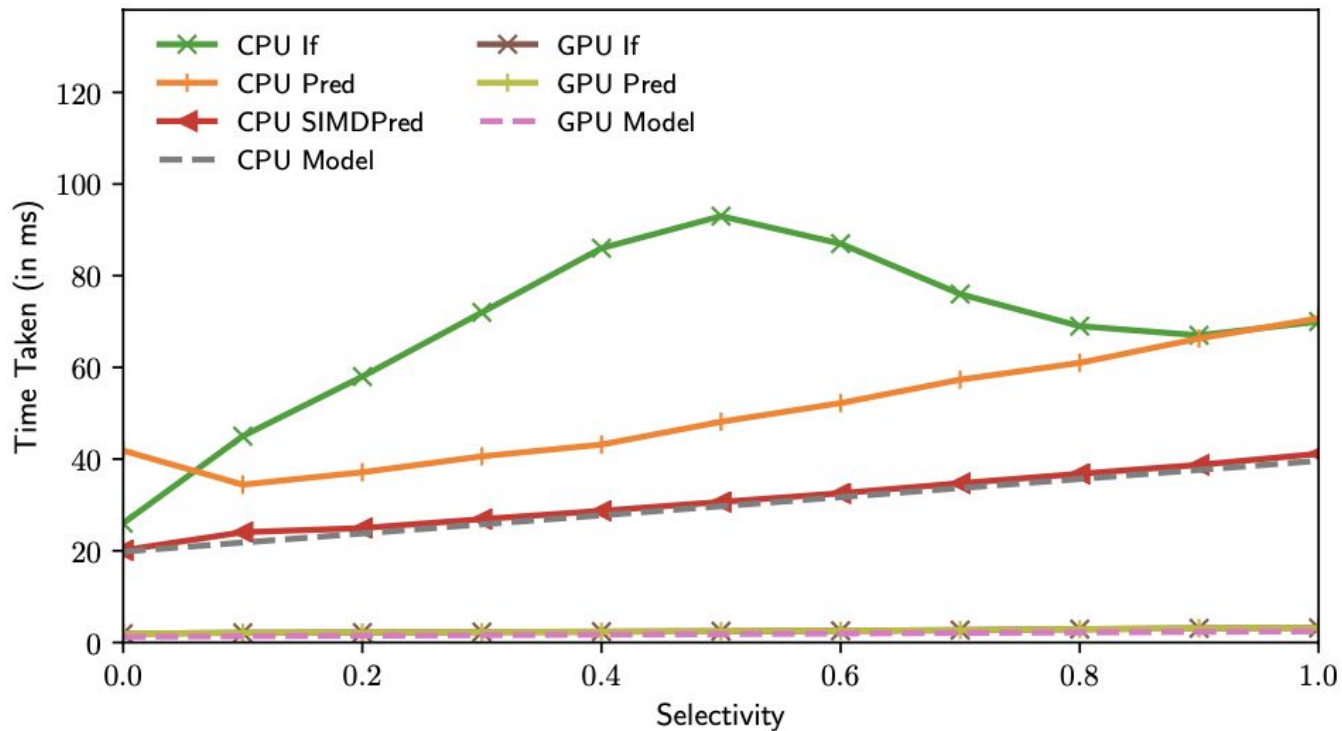
**(a) With branching**

GPU IF

for each y in R:
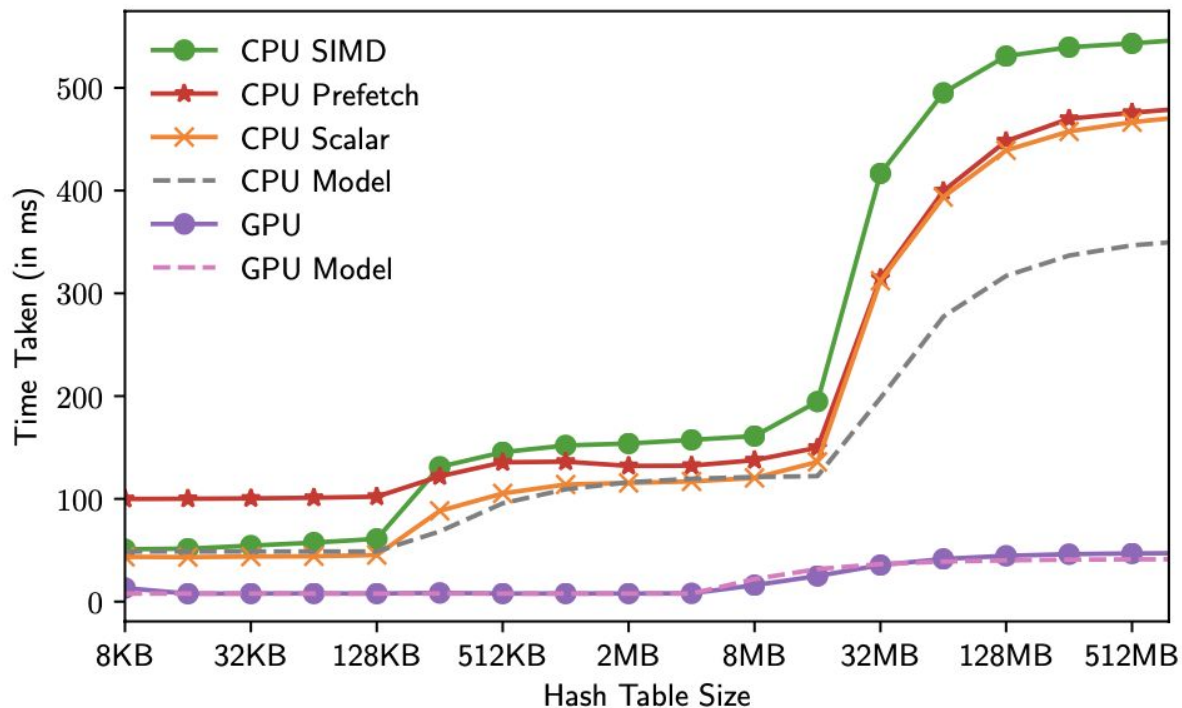  output[i] = y
  i += (y > v)

**(b) With predication**

GPU Pred

# Selection Performance

# Hash Join

Q4: SELECT SUM(A.v + B.v) AS checksum
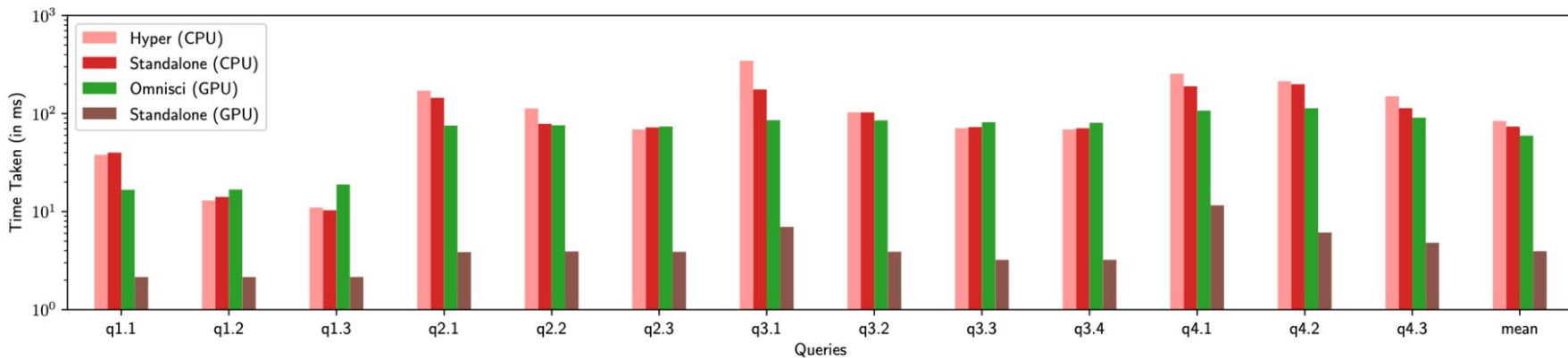    FROM A,B WHERE A.k = B.k

## Important points:

- Probe table: 256 million rows, 50% fill rate, linear probing
- Crystal approach:
  1. Load dating using BlockLoad
  2. Each thread finds matching entries and maintains a local sum
  3. Get overall sum using BlockAggergate

# Hash Join Performance



**Build hash table size**

# SSB Comparison



**Takeaway: Crystal is around 25x than SOTA OLAP**

# The $$$ effect

| | Purchase Cost | Renting Cost |
|---|---|---|
| CPU | $2-5K | $0.504 per hour |
| GPU | $CPU + 8.5K | $3.06 per hour |

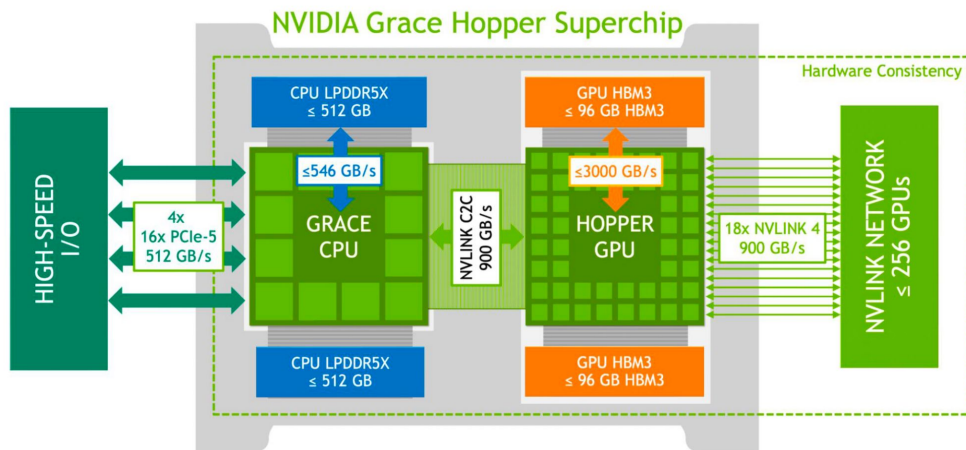**GPUS** are **6x more expensive but 25x faster**

# Limitations

Limited amount of GPU memory:

- Use multiple GPUs

- Bit compression to store more data

Only supports numeric formats

- Strings, Dates, etc..

# Thoughts?