# Hardware Acceleration of Database Operations

Jared Casper and Kunle Olukotun

Pervasive Parallelism Laboratory

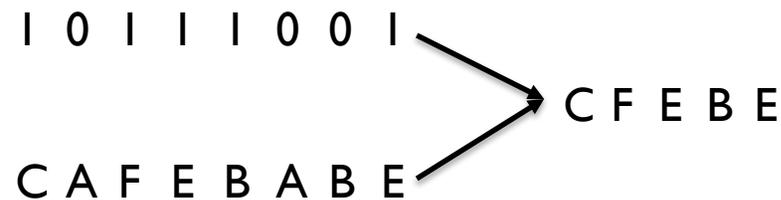Stanford University

# Database machines

- Database machines from late 1970s
  - Put some compute on the disk track/head/unit
  - Processors got faster, I/O performance did not
  - Processor could keep up with disk
    - No performance left on the table
- Today's database machines
  - Made up of general purpose components
  - Massive amounts of memory
  - Very high speed interconnect
  - Tables, even databases, fit entirely within memory
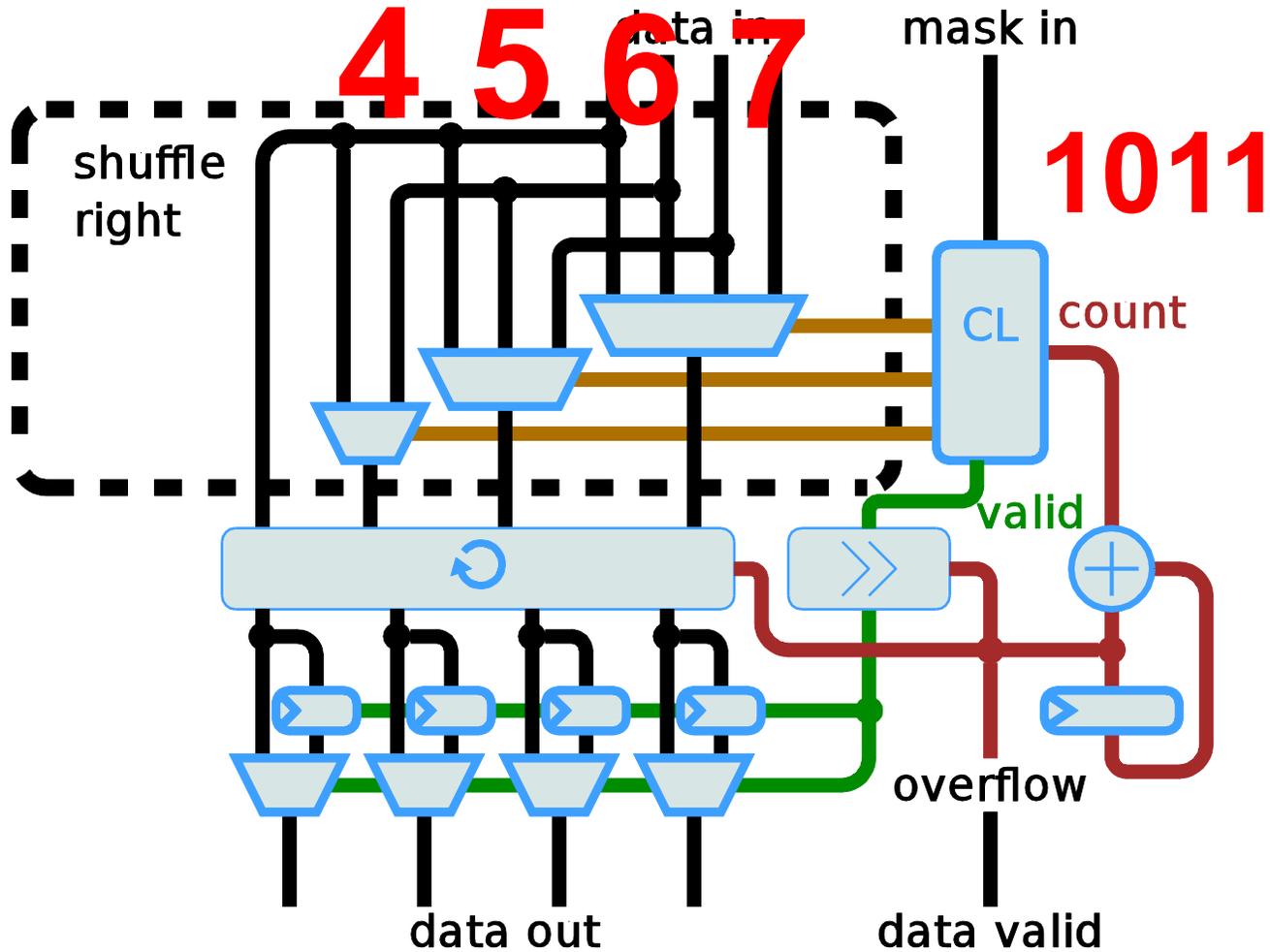
# Database Operation Acceleration

- Processors can not keep up with memory
  - Join performance is at 100s of million tuples per second
  - 64-bit tuples $\rightarrow$ 2-3 GB/s
  - Chips can get over 100 GB/s
  - Performance is being left on the table
- Follow 10x10 rule, build accelerators
- Three acceleration blocks
  - Selection, merge join, sort
  - Combine these to do a sort merge join
  - Goal is to "keep up with memory"

# Select

```
1 0 1 1 1 0 0 1
                    C F E B E
C A F E B A B E
```

- **Software implementation uses SIMD**
  - Read data into SIMD register
  - Use SIMD shuffle operation to move selected data to one end of the register
    - Mask used as index into table for shuffle values
  - Unaligned write to append to output
  - Limited by SIMD width, number of SIMD registers

# Select



**4 5 6 7**

**1011**

data in     mask in

shuffle right

CL   count

valid
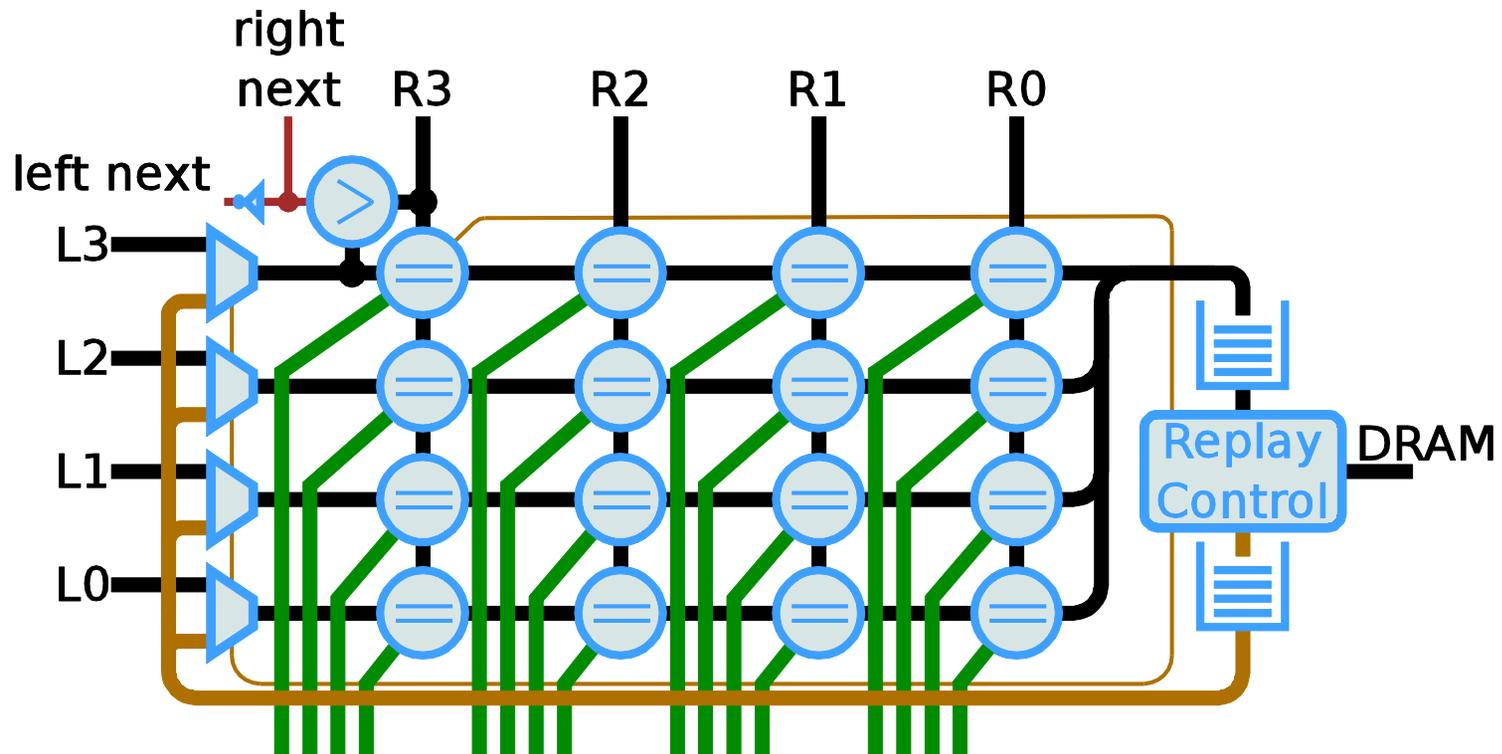
count

overflow

data out     data valid

# Merge Join

- Scan two sorted columns, output matching values
  - Can have associated values or record IDs
  - Output cross product when multiple values
  - Generally viewed as the "free" thing after sorting
    - More an indication of how slow sorting is
- Software implementations have bad branching behaviour
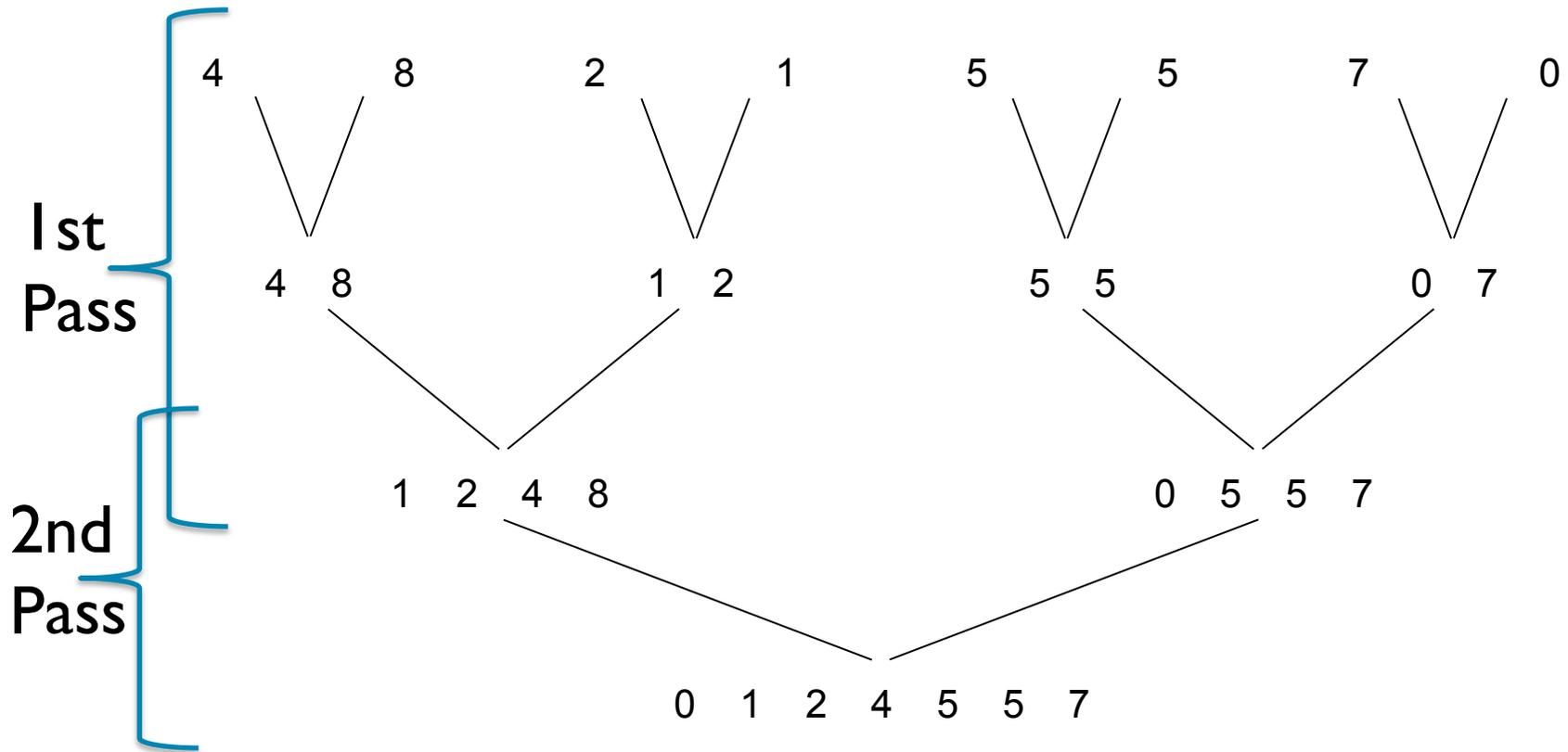  - Limits the IPC $\rightarrow$ hard to keep up with memory

# Merge Join



- Output is bitmask of equal keys with corresponding values
  - Ready for input into the select block

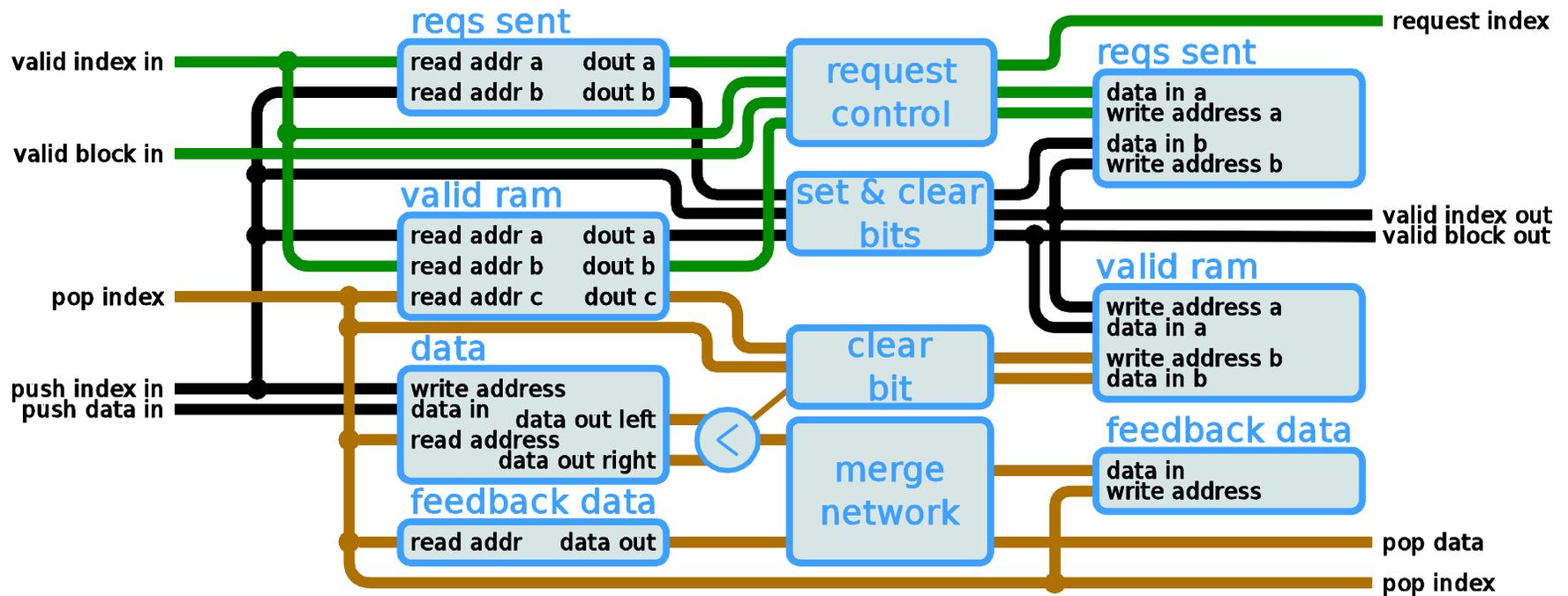1st
Pass

2nd
Pass

| 4 | 8 | 2 | 1 | 5 | 5 | 7 | 0 |

4 8        1 2              5 5        0 7

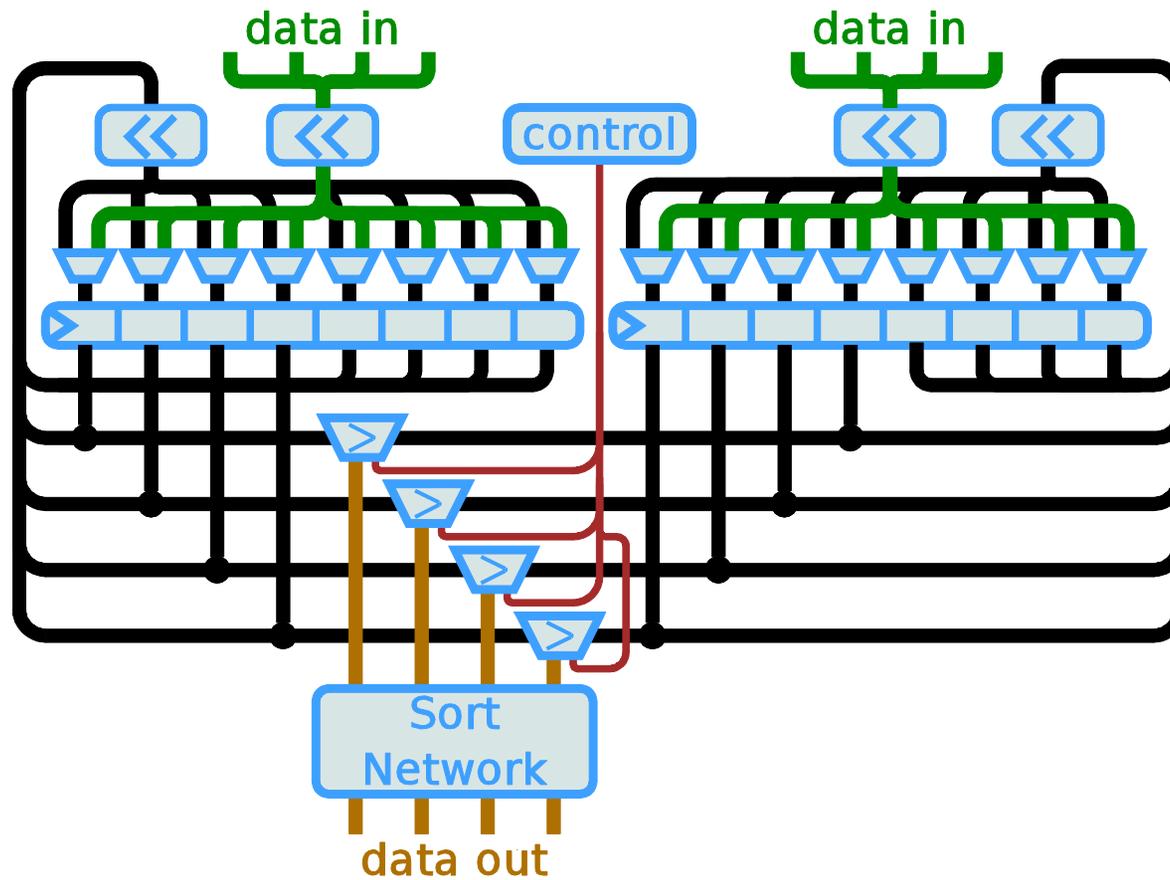1  2  4  8                    0  5  5  7

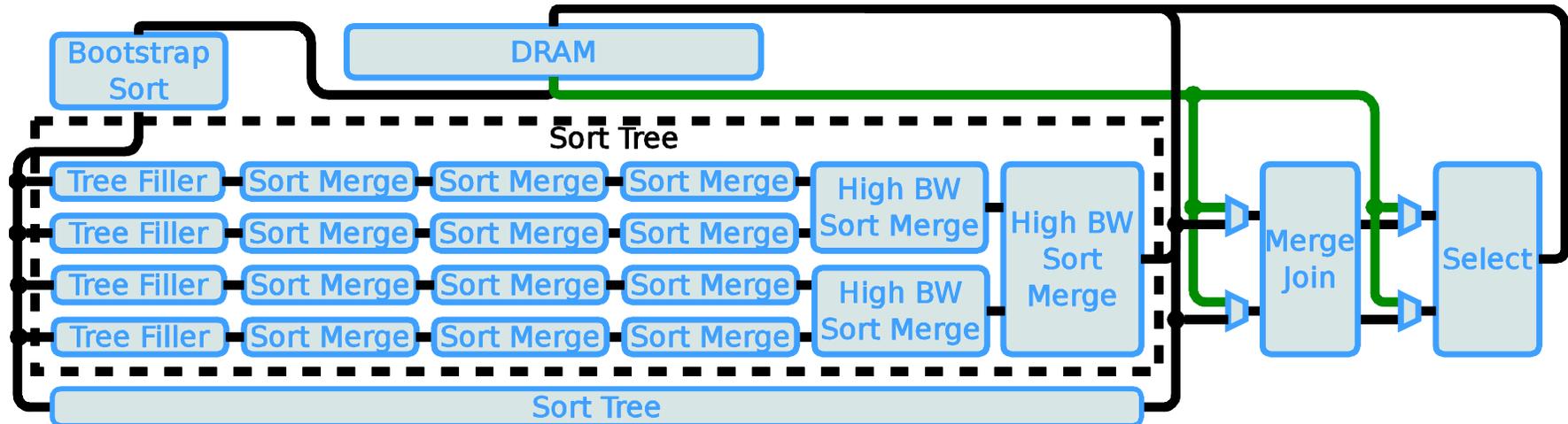0  1  2  4  5  5  7

# Merge Sort Level
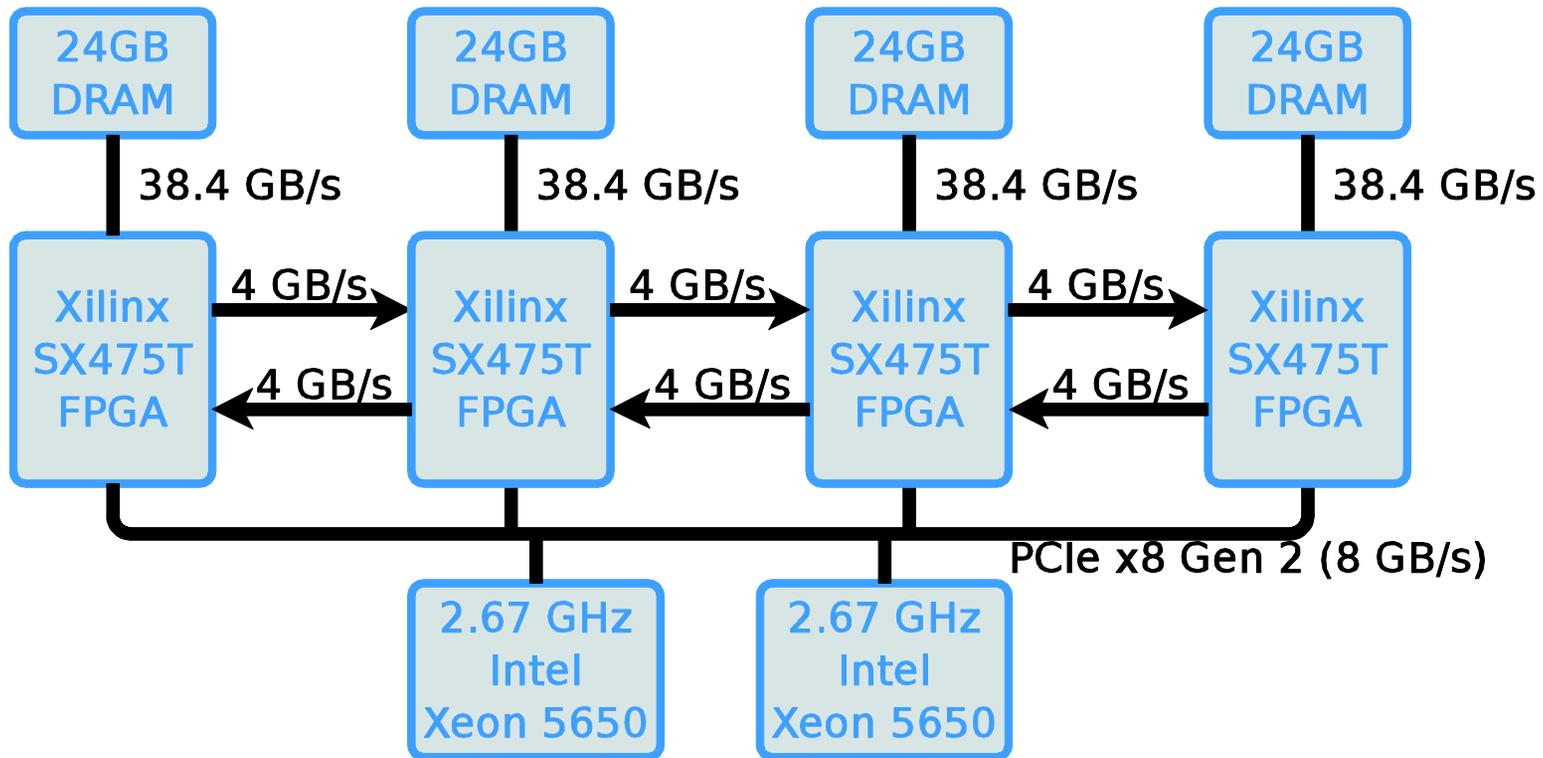
# High Bandwidth Sort Merge Node

# Sort Merge Join



- Sort, merge join, and select blocks are combined to perform an full sort merge join in hardware
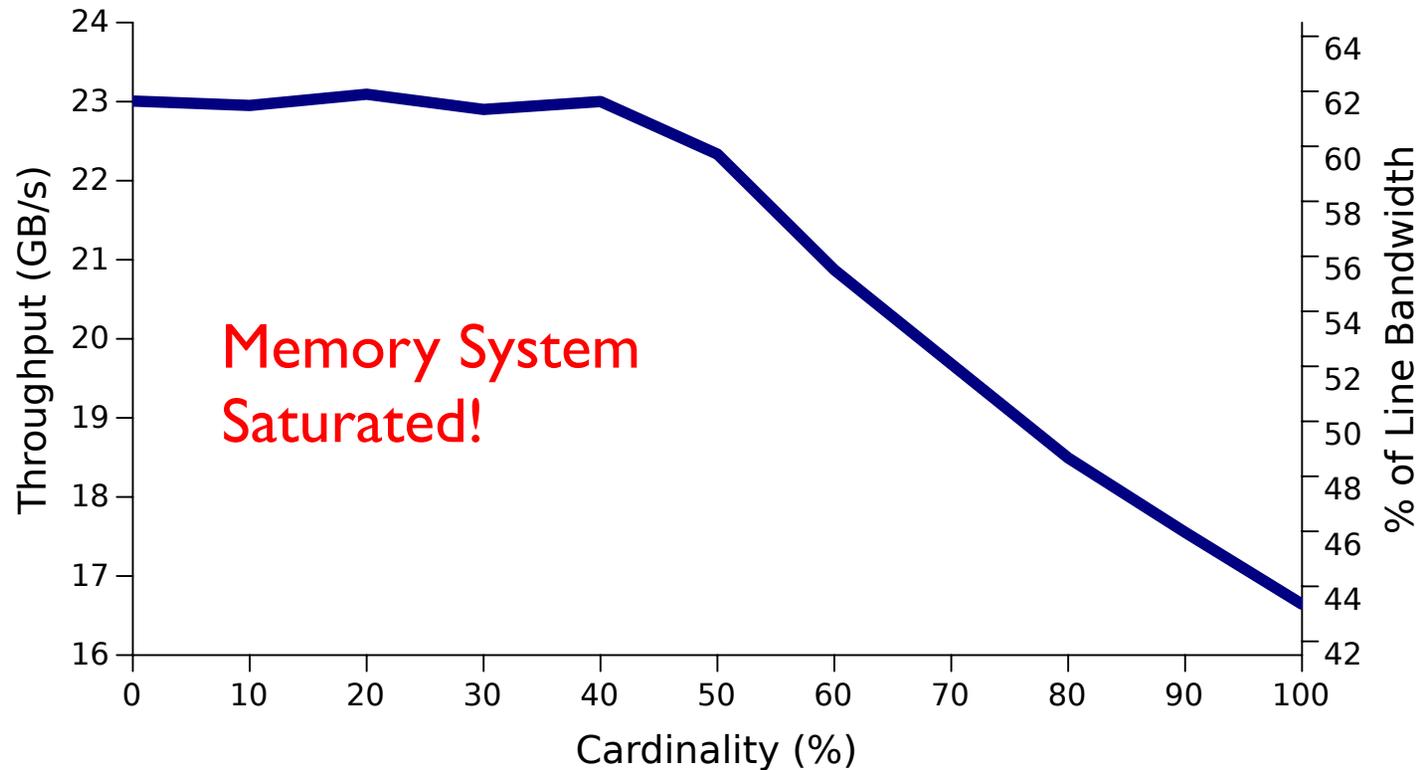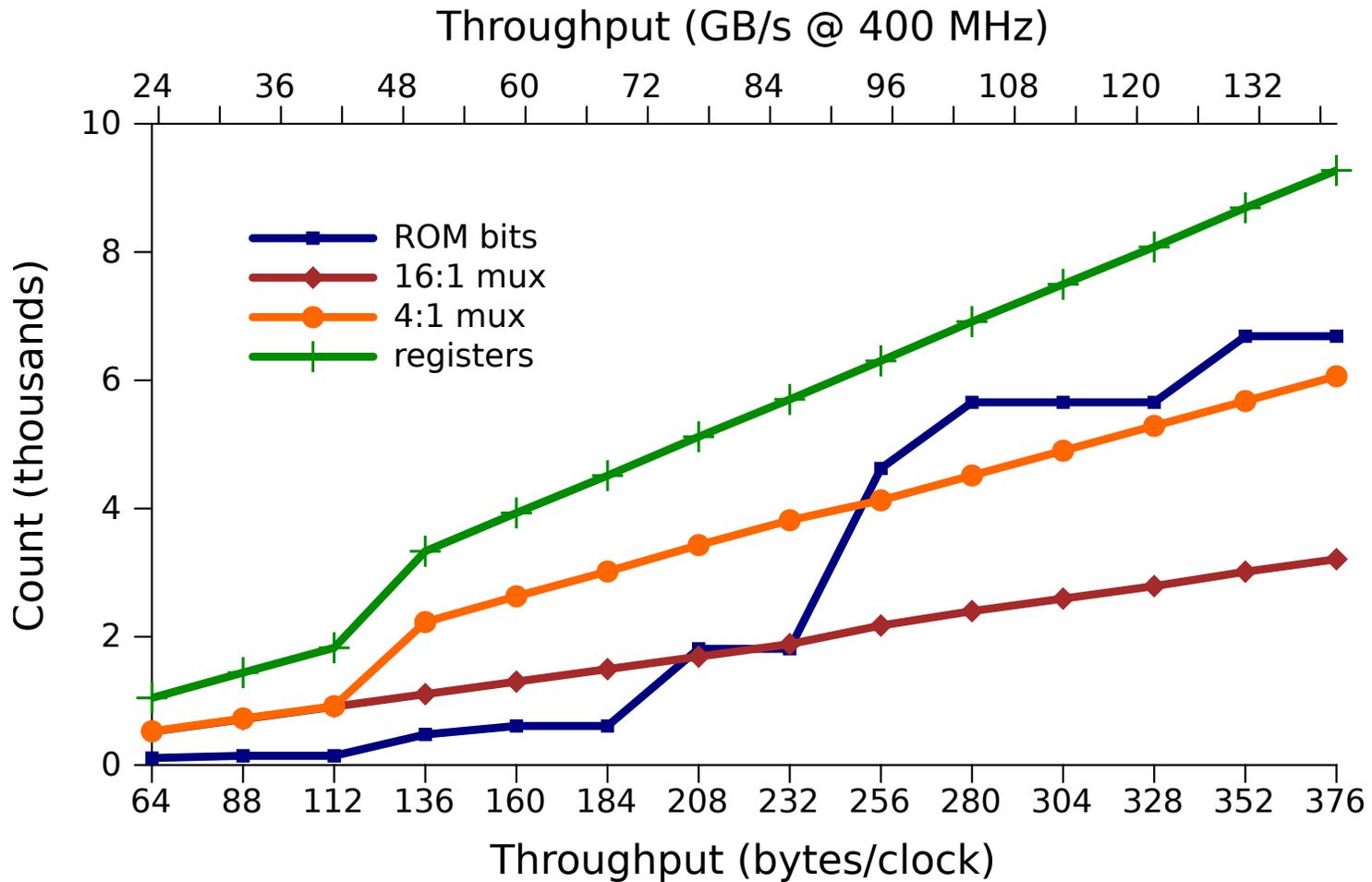
# Prototyping Platform - Maxeler

| 24GB DRAM | | 24GB DRAM | | 24GB DRAM | | 24GB DRAM |
|---|---|---|---|---|---|---|

38.4 GB/s　　38.4 GB/s　　38.4 GB/s　　38.4 GB/s

| Xilinx SX475T FPGA | 4 GB/s → ← 4 GB/s | Xilinx SX475T FPGA | 4 GB/s → ← 4 GB/s | Xilinx SX475T FPGA | 4 GB/s → ← 4 GB/s | Xilinx SX475T FPGA |
|---|---|---|---|---|---|---|

| 2.67 GHz Intel Xeon 5650 | 2.67 GHz Intel Xeon 5650 |
|---|---|

PCIe x8 Gen 2 (8 GB/s)

# Select Throughput



**Software achieved 7 GB/s (33%)**

- STREAM achieved 12 GB/s (57%)
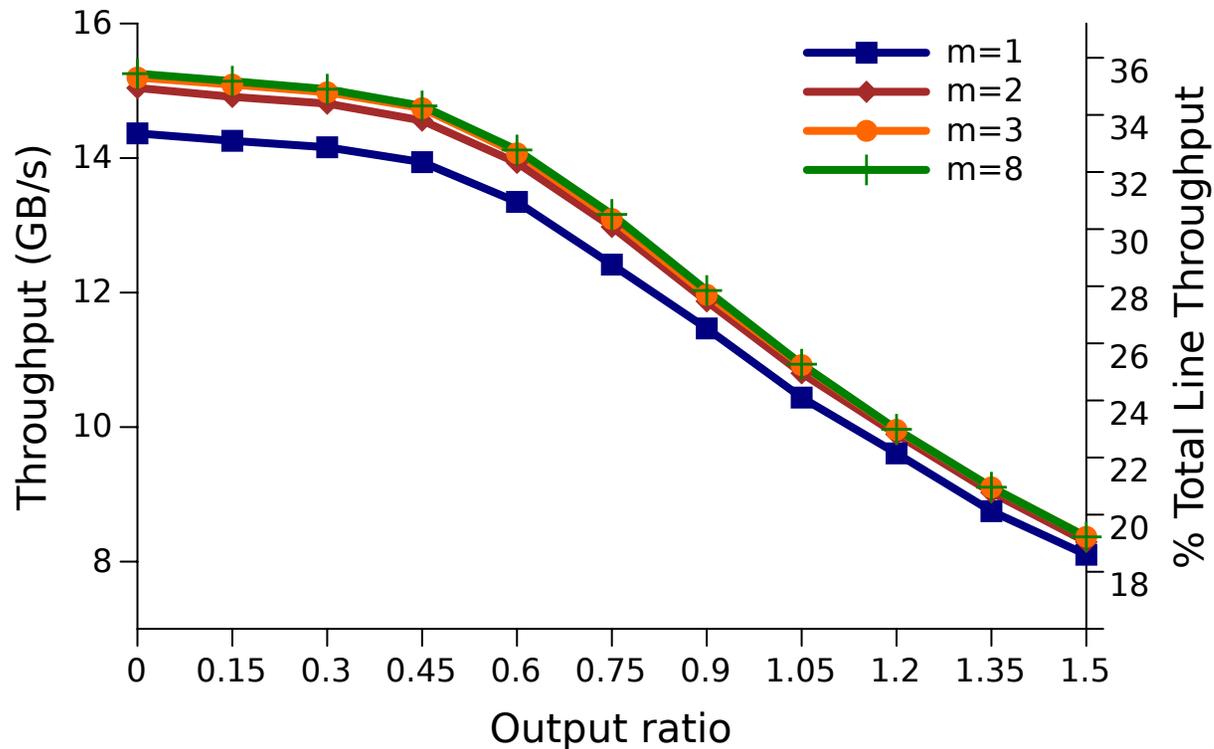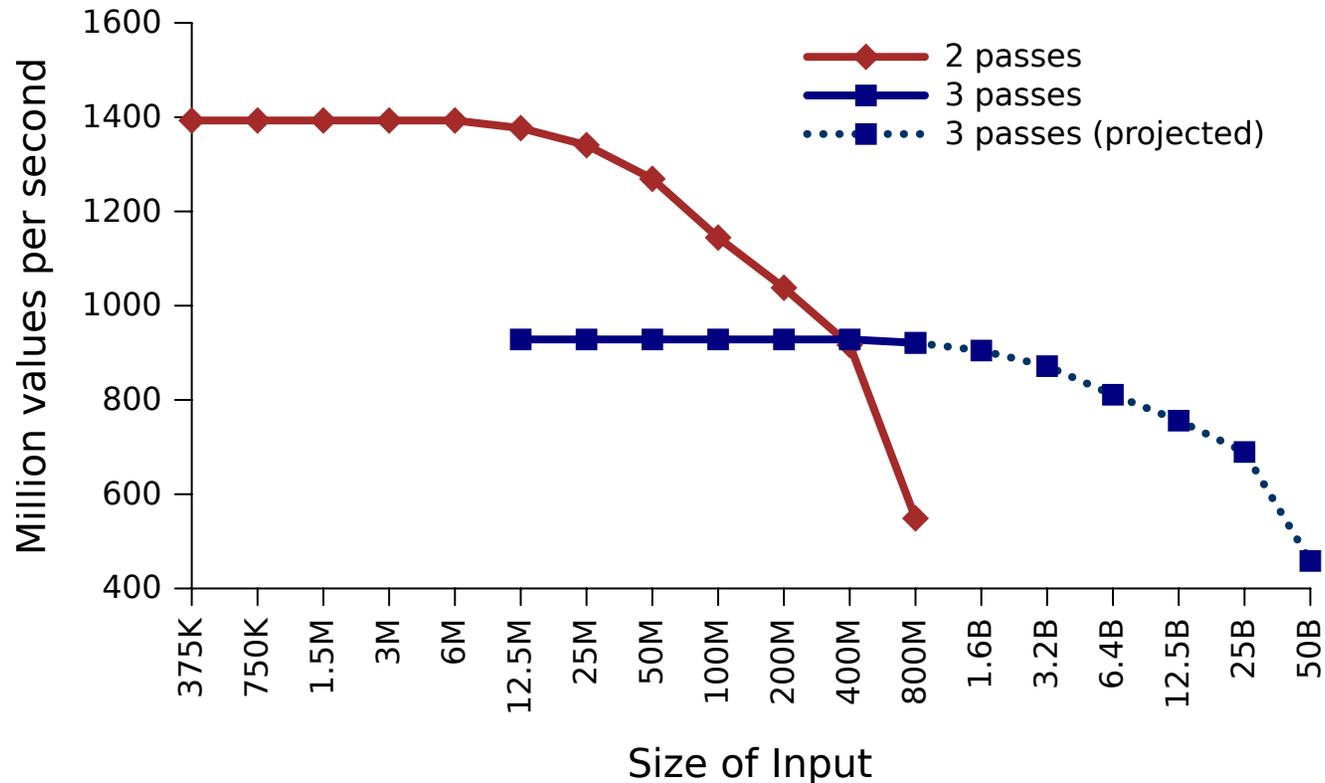
# Merge Join Throughput



- Resources required is a quadratic function of desired bandwidth
  - All in comparison logic, routing was the limiting factor
- Above 1.5x output, write bandwidth dominates
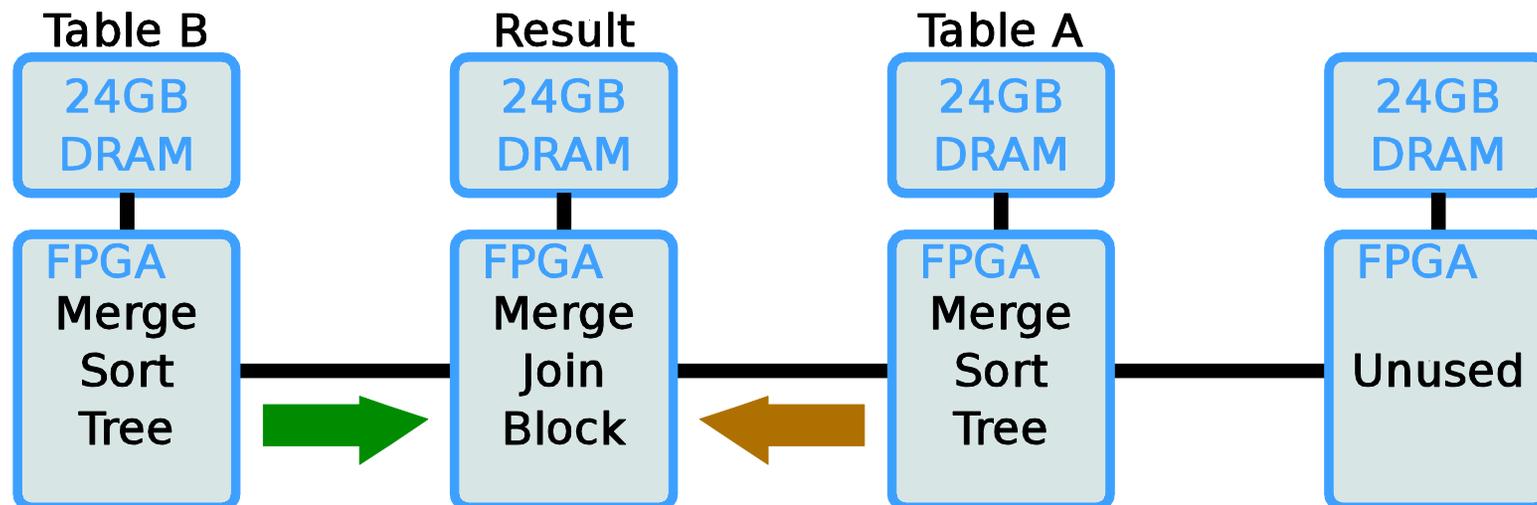  - Throughput above is input consumed

# Sort throughput



- Resources required is a linear function of desired input size
  - Dominated by the memory required to hold working sets
- Recent CPU/GPU numbers ~300M 32-bit values per second

# Sort Merge Join

| Table B | Result | Table A | |
|---|---|---|---|
| 24GB DRAM | 24GB DRAM | 24GB DRAM | 24GB DRAM |
| FPGA Merge Sort Tree | FPGA Merge Join Block | FPGA Merge Sort Tree | FPGA Unused |

- Performance limited by intra-FPGA link
- Total throughput is 800 million tuples/second
    - ~6.5 GB/s
    - 8x previous work on software joins

- FPGAs can be used to saturate memory bandwidth in ways that processors can not
  - Make the most of every byte read
  - In some cases, address bandwidth is just as important as raw data bandwidth
- Scaling your design to high bandwidths can greatly influence the architecture
  - Think streaming
- Next step is to interact with the rest of the system

Questions?