

Multi-Core, Main-Memory Joins: Sort vs. Hash Revisited

Presenter: Haonan Wang

Slides Credit: CMU 15-721 Spring 2018

haonanw@mit.edu

March 19, 2019

1 Background

- Sort vs. Hash
- Motivation

2 Merge - Sort Join

- The basic idea
- Sort Phase
- Merge Phase
- Multi-Way Merge

3 Experiment

Section 1

Background

Subsection 1

Sort vs. Hash

Sort vs. Hash

There are two main approaches for the PARALLEL JOIN ALGORITHMS:

- Hash Join
- Sort-Merge Join

History of Hash VS. Sort

- 1970s Sorting
- 1980s Hashing
- 1990s Equivalent
- 2000s Hashing
- 2010s Hashing (Partitioned vs. Non-Partitioned)
- 2020s ???

What Is Merge-Sort Join

Sort-merge join algorithm explained

Customers table

Id	Login
2	User#2
1	User#1
4	User#4
3	User#3

↓

Id	Login
1	User#1
2	User#2
3	User#3
4	User#4

Orders table

Id	User id
1001	2
1002	4
1003	4
1004	1

↓

Id	User id
1004	1
1001	2
1002	4
1003	4

Sorting

Merging

1. Takes the first row in the left - (1, User#1). Does it have a match in the right ? Yes

SIMD?

What is SIMD?

A class of CPU instructions that allow the processor to perform the same operation on multiple data points simultaneously.

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_1+y_1 \\ x_2+y_2 \\ \vdots \\ x_n+y_n \end{bmatrix}$$

Both current AMD and Intel CPUs have ISA and microarchitecture support SIMD operations.

→ MMX, 3DNow!, SSE, SSE2, SSE3, SSE4, AVX

SIMD Makes Sorting Better Than Hashing?



SORT VS. HASH REVISITED: FAST JOIN IMPLEMENTATION ON MODERN MULTI-CORE CPUS
VLDB 2009



- Hashing is faster than Sort-Merge.
- Sort-Merge is faster w/ wider SIMD.



DESIGN AND EVALUATION OF MAIN MEMORY HASH JOIN ALGORITHMS FOR MULTI-CORE CPUS
SIGMOD 2011



- Trade-offs between partitioning & non-partitioning Hash-Join.



MASSIVELY PARALLEL SORT-MERGE JOINS IN MAIN MEMORY MULTI-CORE DATABASE SYSTEMS
VLDB 2012



- Sort-Merge is already faster than Hashing, even without SIMD.



MASSIVELY PARALLEL NUMA-AWARE HASH JOINS
IMDM 2013



- Ignore what we said last year.
- You really want to use Hashing!



MAIN-MEMORY HASH JOINS ON MULTI-CORE CPUS: TUNING TO THE UNDERLYING HARDWARE
ICDE 2013



- New optimizations and results for Radix Hash Join.



AN EXPERIMENTAL COMPARISON OF THIRTEEN RELATIONAL EQUI-JOINS IN MAIN MEMORY
SIGMOD 2016



- Hold up everyone! Let's look at everything for real!

Section 2

Merge - Sort Join

The basic idea for the designing

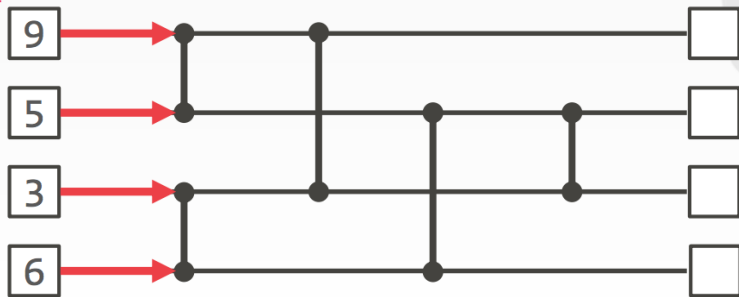
- Partition Phase(Optional)
 - Partition R and assign them to workers / cores.
- Sort Phase
 - Sort the tuples of R and S based on the join key.
- Merge Phase
 - Scan the sorted relations and compare tuples.
 - The outer relation R only needs to be scanned once.

Subsection 2

Sort Phase

Sorting Networks(1)

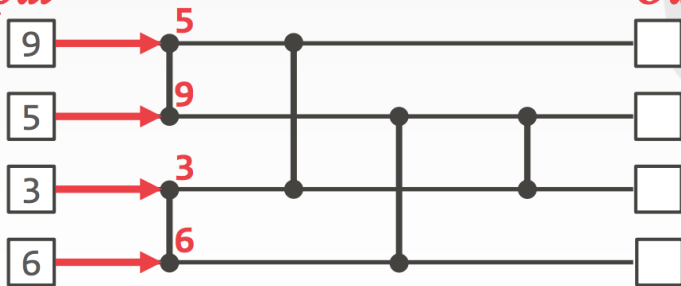
Input



Output

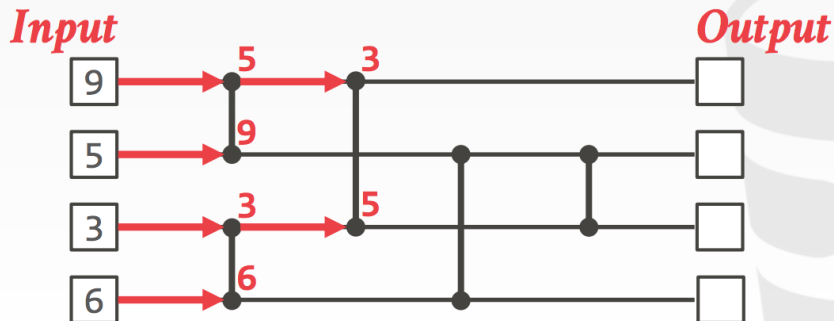
Sorting Networks(2)

Input

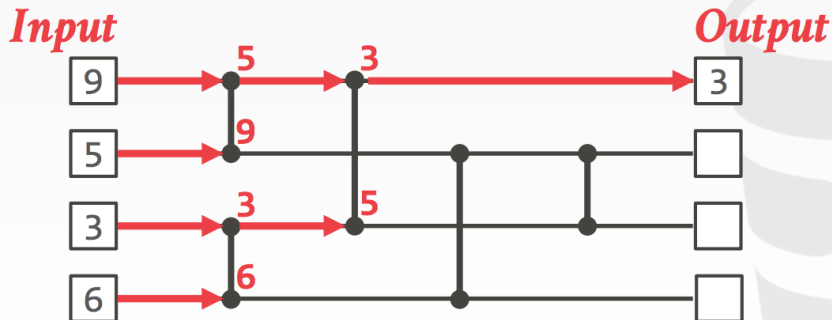


Output

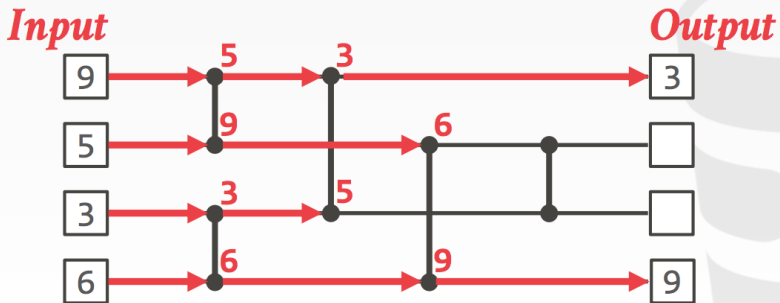
Sorting Networks(3)



Sorting Networks(4)

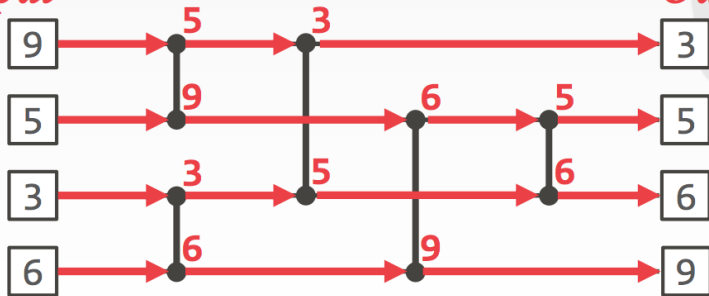


Sorting Networks(5)



Sorting Networks(6)

Input



Output

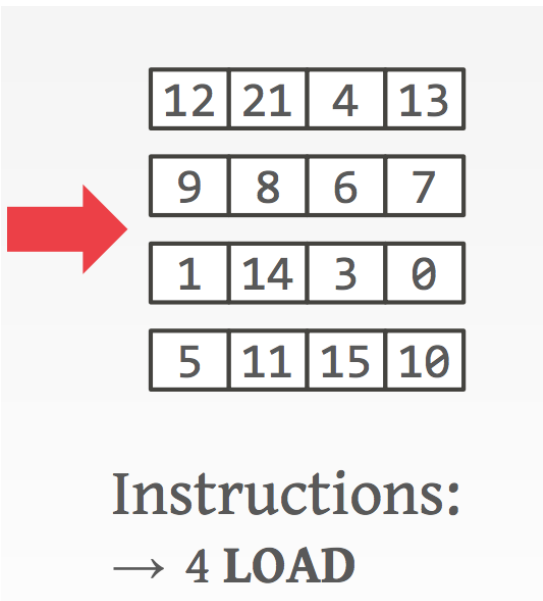
Sorting Networks Summary(1)

```
e = min (a, b)
f = max (a, b)
g = min (c, d)
h = max (c, d)
i = max (e, g)
j = min (f, h)
w = min (e, g)
x = min (i, j)
y = max (i, j)
z = max (f, h)
```

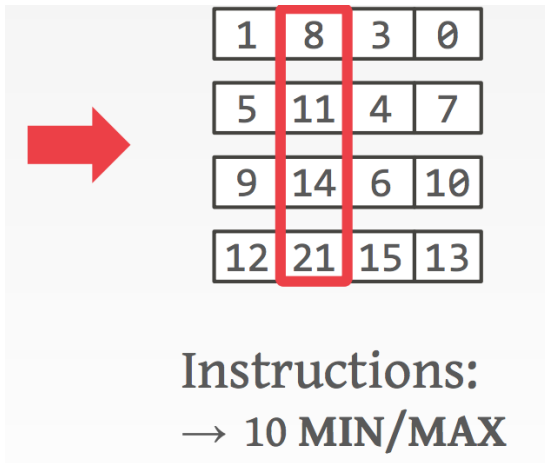
Sorting Networks Summary(2)

- Always has fixed wiring paths for lists with the same number of elements.
- Efficient to execute on modern CPUs because of limited data dependencies and no branches.

Sorting Network Speed Up With SIMD(1)



Sorting Network Speed Up With SIMD(2)

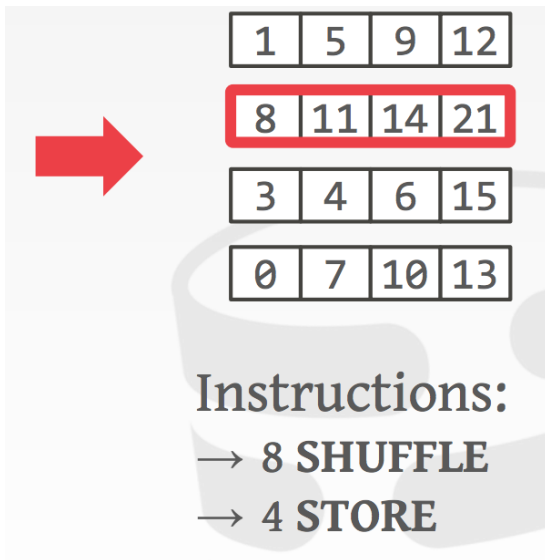


A 4x4 grid of numbers is shown. A red arrow points from the left towards the grid. A red vertical box highlights the second column of the grid, which contains the numbers 8, 11, 14, and 21.

1	8	3	0
5	11	4	7
9	14	6	10
12	21	15	13

Instructions:
→ 10 MIN/MAX

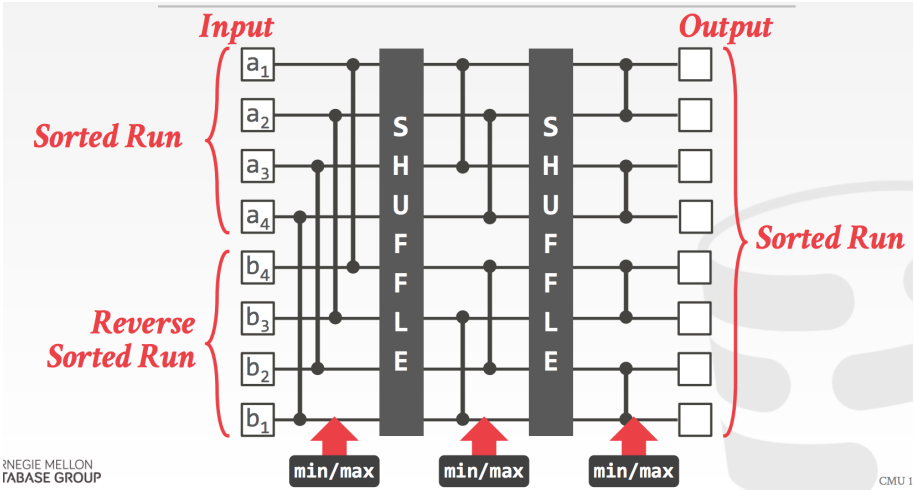
Sorting Network Speed Up With SIMD(3)



Subsection 3

Merge Phase

Bitonic Merge Networks



Merging Larger Lists using Bitonic Merge

Algorithm 1: Merging larger lists with help of bitonic merge kernel `bitonic_merge4()` ($k = 4$).

```
1 a ← fetch4(in1); b ← fetch4(in2);
2 repeat
3   ⟨a, b⟩ ← bitonic_merge4(a, b);
4   emit a to output;
5   if head(in1) < head(in2) then
6     | a ← fetch4(in1);
7   else
8     | a ← fetch4(in2);
9 until eof(in1) or eof(in2);
10 ⟨a, b⟩ ← bitonic_merge4(a, b);
11 emit4(a); emit4(b);
12 if eof(in1) then
13   | emit rest of in2 to output;
14 else
15   | emit rest of in1 to output;
```

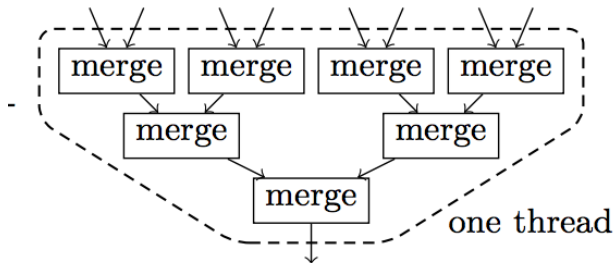


Figure 3: Multi-way merging.

Merging-Sort Hierarchy(Summary)

- *in-register sorting*, with runs that fit into (SIMD) CPU registers;
- *in-cache sorting*, where runs can still be held in a CPU-local cache;
- *out-of-cache sorting*, once runs exceed cache sizes.

Subsection 4

Multi-Way Merge

Impact Of Numa

- In practice, at least some merging passes will inevitably cross NUMA boundaries.
- multisocket systems show an increasing asymmetry, where the NUMA interconnect bandwidth stays further and further behind the aggregate memory bandwidth that the individual memory controllers could provide.

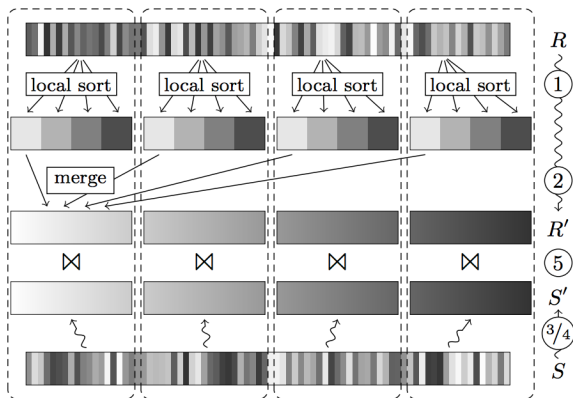


Figure 4: m -way: NUMA-aware sort-merge join with multi-way merge and SIMD.

Section 3

Experiment

- Intel Sandy Bridge with a 256-bit AVX instruction set.
- Four-socket configuration, with each CPU socket containing 8 physical cores and 16 thread contexts by the help of the hyper-threading.
- Cache sizes are 32 KiB for L1, 256 KiB for L2, and 20 MiB L3 (the latter shared by the 16 threads within the socket). The cache line size of the system is 64 bytes. TLB1 contains 64/32 entries when using 4 KiB/2 MiB pages (respectively) and 512 TLB2 entries (page size 4 KiB). Total memory available is 512 GiB (DDR3 at 1600 MHz).

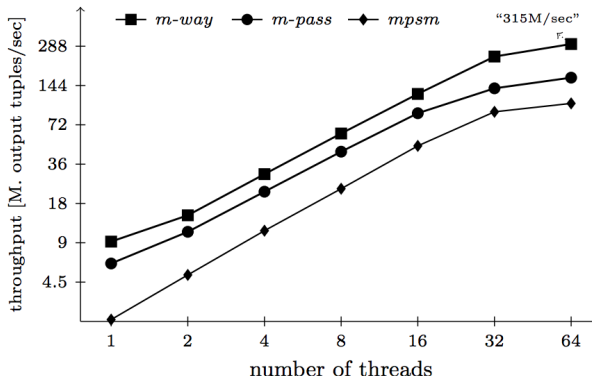
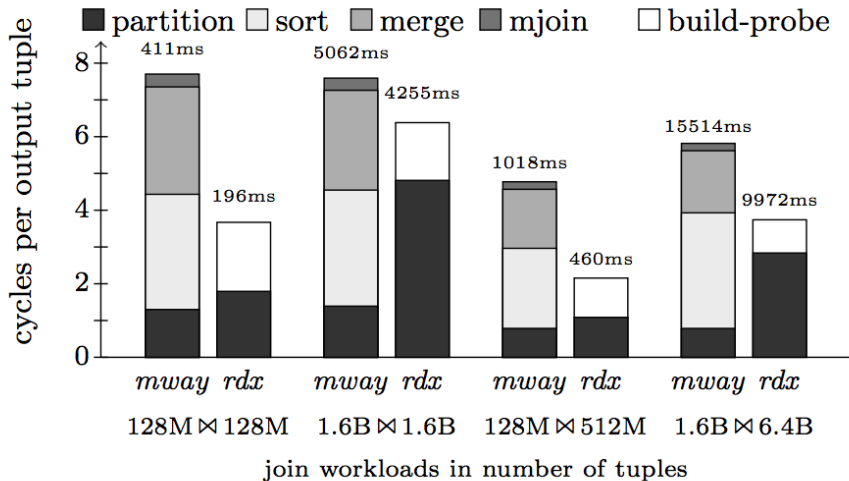


Figure 13: Scalability of sorting-based joins. Workload A, (11.92 GiB \bowtie 11.92 GiB). Throughput metric is output tuples per second, *i.e.* $|S|/\text{execution time}$.

Result(1)



Result(2)

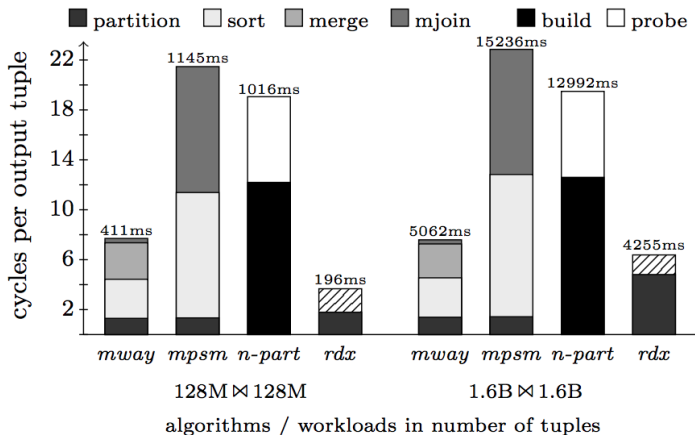


Figure 18: Sort vs. hash join comparison with extended set of algorithms. All using 64 threads.

The End