

ADMINISTRIVIA

Project #2:

- \rightarrow Feedback Submission: Saturday April 1st
- \rightarrow Final Submission: Monday May 1st
- \rightarrow Sign up for a system if you haven't yet!

Project #3

- \rightarrow Status Update Presentation: Monday April 3rd
- → Final Presentations: Friday May 5th @ 5:30pm

TODAY'S AGENDA

Background Sorting Algorithms Parallel Sort-Merge Join Evaluation

SORT-MERGE JOIN (R⊳⊲S)

Phase #1: Sort

 \rightarrow Sort the tuples of **R** and **S** based on the join key(s).

Phase #2: Merge

- \rightarrow Maintain two iterators (one per sorted relation) and compare tuples at each position.
- \rightarrow The outer relation **R** only needs to be scanned once.

SORT-MERGE JOIN (R⋈S)





PARALLEL SORT-MERGE JOINS

Sorting is the most expensive part.

Use hardware correctly to speed up the join algorithm as much as possible.

- \rightarrow Utilize as many CPU cores as possible.
- \rightarrow Be mindful of NUMA boundaries.
- \rightarrow Use SIMD instructions where applicable.



15-721 (Spring 2023)

PARALLEL SORT-MERGE JOIN (R S)

Phase #1: Partitioning (optional)

- \rightarrow Partition **R** and assign them to workers / cores.
- \rightarrow Can use the radix partitioning approach discussed last class.

Phase #2: Sort

 \rightarrow Sort the tuples of **R** and **S** based on the join key.

Phase #3: Merge

- \rightarrow Scan the sorted relations and compare tuples.
- \rightarrow The outer relation **R** only needs to be scanned once.



SORT PHASE

Quicksort is probably what most DBMSs will use. Mergesort is good but requires O(N) additional storage for intermediate results.

We will first discuss a mergesort implementation that takes advantage of NUMA and parallel cores for in-memory data.



CACHE-CONSCIOUS SORTING

Level #1: In-Register Sorting → Sort runs that fit into CPU registers.

Level #2: In-Cache Sorting

 \rightarrow Merge Level #1 output into runs that fit into CPU caches.

 \rightarrow Repeat until sorted runs are ½ cache size.

Level #3: Out-of-Cache Sorting

 \rightarrow Used when the runs of Level #2 exceed the size of caches.

CACHE-CONSCIOUS SORTING



CMU·DB 15-721 (Spring 2023)

Abstract model for sorting keys.

- \rightarrow Fixed wiring "paths" for lists with the same # of elements.
- \rightarrow Efficient to execute on modern CPUs because of limited data dependencies and no branches.



orig = [9,5,3,6]	
<pre>wires₁[0] = min(orig[0], orig[1]) wires₁[1] = max(orig[0], orig[1])</pre>	1
<pre>wires₁[2] = min(orig[2], orig[3]) wires₁[3] = max(orig[2], orig[3])</pre>	
<pre>output[0] = min(wires₁[0], wires₁[2]) wires₂[2] = max(wires₁[0], wires₁[2]) wires₂[1] = min(wires₁[1], wires₁[3]) output[3] = max(wires₁[1], wires₁[3])</pre>	2
<pre>output[1] = min(wires₂[1], wires₂[2]) output[2] = max(wires₂[1], wires₂[2])</pre>	3



<64-bit Join Key, 64-bit Tuple Pointer>





Instructions: \rightarrow 4 LOAD





Instructions: \rightarrow 4 LOAD

Instructions: \rightarrow 10 MIN/MAX





Instructions: \rightarrow 4 LOAD

Instructions: \rightarrow 10 MIN/MAX

Instructions: \rightarrow 8 SHUFFLE \rightarrow 4 STORE



LEVEL #2 - BITONIC MERGE NETWORK

17

Similar technique as a Sorting Network but merges two locally-sorted lists into a globally-sorted list.

Can expand network to merge progressively larger lists up to ½ LLC size.

Intel's Measurements \rightarrow 2.25–3.5× speed-up over SISD implementation.



15-721 (Spring 2023)

LEVEL #2 - BITONIC MERGE NETWORK



LEVEL #3 - MULTI-WAY MERGING

- Use the Bitonic Merge Networks but split the process up into tasks.
- \rightarrow Still one worker thread per core.
- \rightarrow Link together tasks with a cache-sized FIFO queue.

A task blocks when either its input queue is empty, or its output queue is full.

Requires more CPU instructions but brings bandwidth and compute into balance.



LEVEL #3 - MULTI-WAY MERGING



ECMU-DB 15-721 (Spring 2023)

IN-PLACE SUPERSCALAR SAMPLESORT

The <u>IPS⁴o algorithm</u> (2017) recursively partition relation by sampling keys to determine partition boundaries.

- \rightarrow Copies data into output buffers during the partitioning phases.
- → When a buffer gets full, the DBMS writes it back into portions of its existing input buffers instead of allocating a new buffer.

This is the sorting algorithm that we used in CMU's NoisePage DBMS (RIP).

15-721 (Spring 2023)

IN-PLACE PARALLEL SUPER SCALAR SAMPLESORT

VECTORIZED QUICKSORT

Google vqsort (2022)

- \rightarrow Use sorting network for less than 256 keys.
- \rightarrow Based on <u>Google Highway</u> library to provide support for different ISAs and SIMD register sizes.
- \rightarrow Claims to be 1.59x faster than IPS⁴o.

Intel x86-simd-sort (2022)

 \rightarrow Aggressive use of AVX512 instructions.



MERGE PHASE

Iterate through the outer table and inner table in lockstep and compare join keys.

May need to backtrack if there are duplicates.

The DBMS can execute this phase in parallel using multiple workers without synchronization if there are separate output buffers.

SORT-MERGE JOIN VARIANTS

Multi-Way Sort-Merge (M-WAY)



MULTI-CORE, MAIN-MEMORY JOINS: SORT VS. HASH REVISITED VLDB 2013

Multi-Pass Sort-Merge (M-PASS)



MULTI-CORE, MAIN-MEMORY JOINS: SORT VS. HASH REVISITED VLDB 2013

Massively Parallel Sort-Merge (MPSM)



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

Outer Table

- \rightarrow Each core sorts in parallel on local data (levels #1/#2).
- \rightarrow Redistribute sorted runs across cores using range partitioning then perform <u>multi-way merge</u> (level #3).

Inner Table

 \rightarrow Same as outer table.

Merge phase is between matching pairs of chunks of outer/inner tables at each core.













CMU·DB 15-721 (Spring 2023)









Same steps as Outer Table



ECMU-DB 15-721 (Spring 2023)



ECMU·DB 15-721 (Spring 2023)

MULTI-PASS SORT-MERGE

Outer Table

- \rightarrow Same level #1/#2 sorting as previous Multi-Way Merge.
- → But instead of redistributing data across cores, perform a global **<u>multi-pass naïve merge</u>** on sorted runs.

Inner Table

 \rightarrow Same as outer table.

Merge phase is between matching pairs of chunks of outer table and inner table.



MULTI-PASS SORT-MERGE



MASSIVELY PARALLEL SORT-MERGE

Outer Table

- \rightarrow Range-partition outer table and redistribute to cores.
- \rightarrow Each core sorts then in parallel on their local partitions.

Inner Table

- \rightarrow Not redistributed like outer table.
- \rightarrow Each core sorts its local data.

Merge phase is between entire sorted run of outer table and a segment of inner table.

MASSIVELY PARALLEL SORT-MERGE



ECMU·DB 15-721 (Spring 2023)

HYPER'S RULES FOR PARALLELIZATION

Rule #1: No random writes to non-local memory

→ Chunk the data, redistribute, and then each core sorts/works on local data.

Rule #2: Only perform sequential reads on nonlocal memory

 $\rightarrow\,$ This allows the hardware prefetcher to hide remote access latency.

Rule #3: No core should ever wait for another

 \rightarrow Avoid fine-grained latching or sync barriers.

EVALUATION

Compare the different join algorithms using a synthetic data set.

- \rightarrow **Sort-Merge:** M-WAY, M-PASS, MPSM
- → **Hash:** Radix Partitioning

Hardware:

- \rightarrow 4 Socket Intel Xeon E4640 @ 2.4GHz
- \rightarrow 8 Cores with 2 Threads Per Core
- \rightarrow 512 GB of DRAM



COMPARISON OF SORT-MERGE JOINS

Workload: 1.6B≈128M (8-byte tuples)



ECMU·DB 15-721 (Spring 2023)



ECMU·DB 15-721 (Spring 2023)

SORT-MERGE JOIN VS. HASH JOIN

Workload: Different Table Sizes (8-byte tuples)



Source: Cagri Balkesen

SECMU.DB





15-721 (Spring 2023)

PARTING THOUGHTS

Hash join is (almost) always the superior choice for a join algorithm on modern hardware. \rightarrow Most enterprise OLAP DBMS support both.

We did not consider the impact of queries where the output needs to be sorted.

We will see sort-merge joins again next class...

NEXT CLASS

Worst-Case Optimal Joins (aka multi-way joins)