

CS 764: Topics in Database Management Systems Lecture 2: Join

Xiangyao Yu 9/9/2020 Piazza

Sample reviews and exam questions

Lectures after the exam: state-of-the-art research in database systems

Email me if you have problems submitting the review

Today's Paper: Join

Join Processing in Database Systems with Large Main Memories

LEONARD D. SHAPIRO North Dakota State University

We study algorithms for computing the equijoin of two relations in a system with a standard architecture but with large amounts of main memory. Our algorithms are especially efficient when the main memory available is a significant fraction of the size of one of the relations to be joined; but they can be applied whenever there is memory equal to approximately the square root of the size of one relation. We present a new algorithm which is a hybrid of two hash-based algorithms and which dominates the other algorithms we present, including sort-merge. Even in a virtual memory environment, the hybrid algorithm dominates all the others we study.

Finally, we describe how three popular tools to increase the efficiency of joins, namely filters, Babb arrays, and semijoins, can be grafted onto any of our algorithms.

Categories and Subject Descriptors: H.2.0 [Database Management]: General; H.2.4 [Database Management]: Systems—query processing; H.2.6 [Database Management]: Database Machines

General Terms: Algorithms, Performance

Additional Key Words and Phrases: Hash join, join processing, large main memory, sort-merge join

ACM Transactions on Database Systems, 1986

Agenda

System architecture and assumptions Notations

- Join algorithms
 - Sort merge join
 - Simple hash join
 - GRACE hash join
 - Hybrid hash join

Partition overflow and additional techniques

System Architecture and Assumptions



CPU: uniprocessor

- Avoids sync complexity
- Could be built on systems of the day

Memory

Tens of Megabytes

Focus only on equi-join

Notation

Relations: R, S (| R | < | S |) Join: S \bowtie R Memory: M

I R I: number of blocks in relation R (similar for S and M) **F**: hash table for R occupies I R I * F blocks

Join Algorithms

Sort Merge Join

Phase 1: Produce sorted runs of S and R

Phase 2: Merge runs of S and R, output join result



Sort Merge Join

Phase 1: Produce sorted runs of S and R

Phase 2: Merge runs of S and R, output join result



Unsorted R and S

Sorted runs of R and S

Sort Merge Join

Phase 1: Produce sorted runs of S and R

Phase 2: Merge runs of S and R, output join result



Unsorted R and S

Sorted runs of R and S

Find matches in sorted runs ¹⁰

Sort Merge Join – Phase 1

Phase 1: Produce sorted runs of S and R

- Each run of S will be 2 × I M I average length
- Q: Where does 2 come from? A: Replacement selection

Memory

Priority queue (heap)



Memory layout in Phase 1



Naïve solution:

Each run contains | M | blocks

- Load I M I blocks
- Sort
- Output I M I blocks



Replacement selection:

load I M I blocks and sort

While heap is not empty

Output one tuple and load one tuple from input buffer If the new tuple < any tuple in output save the tuple for next run (heap size reduces) else

heap reorder



Replacement selection:

Each run contains 2 × I M I blocks

load I M I blocks and sort

https://opendsa-server.cs.vt.edu/ODSA/Books/Everything/html/ExternalSort.html

While heap is not empty

Output one tuple and load one tuple from input buffer If the new tuple < any tuple in output save the tuple for next run (heap size reduces) else

heap reorder



Replacement selection:

Each run contains 2 × I M I blocks

load | M | blocks and sort

https://opendsa-server.cs.vt.edu/ODSA/Books/Everything/html/ExternalSort.html

While heap is not empty

Output one tuple and load one tuple from input buffer If the new tuple < any tuple in output save the tuple for next run (heap size reduces) else

heap reorder

Total number of runs
=
$$\frac{|S|}{2 \times |M|} + \frac{|R|}{2 \times |M|} \le \frac{|S|}{|M|}$$

Sort Merge Join – Phase 2

Phase 2: Merge runs of S and R, output join result

• One input buffer required for each run



Memory layout in Phase 2

Sort Merge Join – Phase 2

Phase 2: Merge runs of S and R, output join result

• One input buffer required for each run

Requirement

 $|M| \ge total number runs$

Satisfied if
$$|M| \ge \frac{|S|}{|M|}$$

namely $|M| \ge \sqrt{|S|}$



Memory layout in Phase 2

Hash Join

Build a hash table on the smaller relation (**R**) and probe with larger (**S**) Hash tables have overhead, call it **F**

When **R** doesn't fit fully in memory, partition hash space into ranges



Simple Hash Join

Build a hash table on R



Simple Hash Join – 1st pass

- Build a hash table on R
- If **R** does not fit in memory, find a subset of buckets that fit in memory



Simple Hash Join – 1st pass

- Build a hash table on R
- If **R** does not fit in memory, find a subset of buckets that fit in memory
- Read in ${\boldsymbol{S}}$ to join with the subset of ${\boldsymbol{R}}$



Simple Hash Join – 1st pass

- Build a hash table on R
- If **R** does not fit in memory, find a subset of buckets that fit in memory
- Read in S to join with the subset of R
- The remaining tuples of ${\bf S}$ and ${\bf R}$ are written back to disk



Simple Hash Join – 2nd pass

- Build a hash table on R
- If **R** does not fit in memory, find a subset of buckets that fit in memory
- Read in S to join with the subset of R
- The remaining tuples of ${\bf S}$ and ${\bf R}$ are written back to disk



Simple Hash Join – 3rd pass

- Build a hash table on R
- If **R** does not fit in memory, find a subset of buckets that fit in memory
- Read in S to join with the subset of R
- The remaining tuples of ${\bf S}$ and ${\bf R}$ are written back to disk



Phase 1: Partition both R and S into pairs of shards Phase 2: Separately join each pairs of partitions



Phase 1: Partition both R and S into pairs of shards

Phase 2: Separately join each pairs of partitions





Memory layout in Phase 1

Phase 1: Partition both R and S into pairs of shards

Phase 2: Separately join each pairs of partitions



1	

Memory layout in Phase 2



Assume k partitions for R and S

In phase 1, needs one output buffer (i.e., block) for each partition

 $k \leq |M|$

Assume **k** partitions for **R** and **S** In phase 1, needs one output buffer (i.e., block) for each partition $k \le |M|$

In phase 2, the hash table of each shard of **R** must fit in memory

$$\frac{\mid R \mid}{k} \times F \le \mid M \mid$$

Assume **k** partitions for **R** and **S** In phase 1, needs one output buffer (i.e., block) for each partition $k \le |M|$

In phase 2, the hash table of each shard of **R** must fit in memory

$$\frac{|R|}{k} \times F \le |M|$$

The maximum size of **R** to perform Grace hash join:

$$|R| \le \frac{|M|}{F} k \le \frac{|M|^2}{F} \qquad |M| \ge \sqrt{|R| \times F}$$

GRACE vs. Simple Hash Join

When I R I × F < I M I

- Simple hash join incurs no IO traffic
- GRACE hash join writes and reads each table (i.e., the partitions) once

When IRI×F>>IMI

- Simple hash join incurs significant IO traffic
- GRACE hash join writes and reads each table (i.e., the partitions) once

When you have two algorithms that are good in different settings, create a hybrid!

When you have two algorithms that are good in different settings, create a hybrid!



Memory layout in Phase 1 of GRACE hash join

When you have two algorithms that are good in different settings, create a hybrid!



Memory layout in Phase 1 of hybrid hash join

When you have two algorithms that are good in different settings, create a hybrid! Memory

Case 1: | R | × F < | M |

Identical to simple hash join

Hash table for R₀

Memory layout in Phase 1 of hybrid hash join

When you have two algorithms that are good in different settings, create a hybrid!

Case 1: | R | × F < | M |

Identical to simple hash join

Case 2: | R | × F >> | M | Similar to GRACE hash join



Memory layout in Phase 1 of hybrid hash join

Evaluation



- Conclusion 1: Hash join is generally better than sort-merge join
- Conclusion 2: Hybrid hash join is strictly better than simple and GRACE hash joins

So far we assume uniform random distribution for **R** and **S**

What if we guess wrong on size required for R hash table and a partition does not fit in memory?

Solution: further divide into smaller partitions range

Additional Techniques

Babb array (or bitmap filter)

- Set a bit for each R tuple
- Use to filter S during initial scan, discard tuple if missing in array

Semijoin

- Project join attributes from R, join to S, then join that result back to R
- Useful if full R tuples won't fit into memory, but join will be selective and filter many S tuples
- Can be added to any join algorithm above

Join – Q/A

Why sqrt(325 MB) is 4 MB?

sqrt(325MB / block_size) = 4 MB / block_size

Modern systems using Babb array?

Join in in-memory database?

Evaluation on real, parallel systems?

Babb filter vs. Bloom filter

Is it possible to make GRACE hash join work when $|M| < \sqrt{|R| \times F}$? For example, |M| = 10, F = 1, |R| = 1000. You may modify the GRACE hash join algorithm as described in the paper.

Is it possible for a sort-merge join algorithm to outperform a hash-based join algorithm? If yes, when can this happen?

Before Next Lecture

Submit discussion summary to https://wisc-cs764-f20.hotcrp.com

- Title: Lecture 2 discussion. group ##
- Authors: Names of students who joined the discussion

Deadline: Thursday 11:59pm

Submit review for

Hong-Tai Chou, David J. DeWitt, *An Evaluation of Buffer Management Strategies for Relational Database Systems*. Algorithmica 1986.