Intro to Database Systems (15-445/645) Join Algorithms



ADMINISTRIVIA

Homework 3 due Sunday → You may not turn in Homework 3 late

Midterm exam Wednesday, March 1st

- \rightarrow Practice exam coming tomorrow
- \rightarrow Exam accommodations? Schedule with EOS

Project 2 is available

- \rightarrow First checkpoint due Friday, March 3rd (15% of P2 grade)
- \rightarrow Overall due Wednesday, March 22nd (85% of P2 grade)

ECMU·DB 15-445/645 (Spring 2023)

LAST TIME

Finished concurrent B+Trees

Sorting \rightarrow Top-k heap sort \rightarrow External merge sort

Aggregations \rightarrow External hashing

ECMU·DB 15-445/645 (Spring 2023)

RECALL: QUERY PLAN

The operators are arranged in a tree.

Data flows from the leaves of the tree
up towards the root.
→ We will discuss the granularity of the data movement next week.

The output of the root node is the result of the query.

SELECT R.id, S.cdate
FROM R JOIN S
ON R.id = S.id
WHERE S.value > 100



ECMU·DB 15-445/645 (Spring 2023

WHY DO WE NEED TO JOIN?

We normalize tables in a relational database to avoid unnecessary repetition of information.

We then use the join operator to reconstruct the original tuples without any information loss.

JOIN ALGORITHMS

We will focus on performing binary joins (two tables) using <u>inner equijoin</u> algorithms.
→ These algorithms can be tweaked to support other joins.
→ Multi-way joins exist primarily in research literature.

In general, we want the smaller table to always be the left table ("outer table") in the query plan.
→ The optimizer will (try to) figure this out when generating the physical plan.

ECMU-DB 15-445/645 (Spring 2023

JOIN OPERATORS

Decision #1: Output

→ What data does the join operator emit to its parent operator in the query plan tree?

Decision #2: Cost Analysis Criteria

→ How do we determine whether one join algorithm is better than another?

SELECT R.id, S.cdate
FROM R JOIN S
ON R.id = S.id
WHERE S.value > 100



ECMU·DB 15-445/645 (Spring 2023

OPERATOR OUTPUT

For tuple $r \in R$ and tuple $s \in S$ that match on join attributes, concatenate r and s together into a new tuple.

Output contents can vary:

- \rightarrow Depends on processing model
- \rightarrow Depends on storage model
- \rightarrow Depends on data requirements in query





CMU·DB 15-445/645 (Spring 2023

OPERATOR OUTPUT: DATA

Early Materialization:

 \rightarrow Copy the values for the attributes in outer and inner tuples into a new output tuple.

Subsequent operators in the query plan never need to go back to the base tables to get more data.

WHERE S.value > 100 R(id, na id na TT R.id, S.cdate 123 lab R.id=S.id R.id R.name S.id S.value S.cdate 123 123 1000 2/23/23 abc 123 abc 123 2000 2/23/23

SELECT R.id, S.cdate

ON R.id = S.id

FROM R JOIN S

SECMU-DB 15-445/645 (Spring 2023)

OPERATOR OUTPUT: RECORD IDS

Late Materialization:

→ Only copy the joins keys along with the Record IDs of the matching tuples.

Ideal(?) for column stores because the DBMS does not copy data that is not needed for the query.



SELECT R.id, S.cdate

ON R.id = S.id

FROM R JOIN S

10

COST ANALYSIS CRITERIA

Assume:

 $\rightarrow M$ pages in table **R**, *m* tuples in **R** $\rightarrow N$ pages in table **S**, *n* tuples in **S** SELECT R.id, S.cdate
FROM R JOIN S
ON R.id = S.id
WHERE S.value > 100

Cost Metric: # of I/Os to compute join

We ignore overall output costs because it depends on the data and is the same for all algorithms.

ECMU·DB 15-445/645 (Spring 2023

JOIN VS CROSS-PRODUCT

R▷⊲S is the most common operation and thus must be carefully optimized.
R×S followed by a selection is inefficient because the cross-product is large.

There are many algorithms for reducing join cost, but no algorithm works well in all scenarios.

JOIN ALGORITHMS

- Nested Loop Join
- \rightarrow Naïve
- \rightarrow Block
- \rightarrow Index
- Sort-Merge Join
- Hash Join
- \rightarrow Simple
- \rightarrow GRACE (Externally Partitioned)
- \rightarrow Hybrid

ECMU-DB 15-445/645 (Spring 2023)

NAÏVE NESTED LOOP JOIN

foreach tuple r \in R: Outer foreach tuple s \in S: Inner if r and s match then emit

R(id,name)

id	name
600	MethodMan
200	GZA
100	Andy
300	ODB
500	RZA
700	Ghostface
400	Raekwon

S(id, value, cdate)

id	value	cdate
100	2222	2/23/23
500	7777	2/23/23
400	6666	2/23/23
100	9999	2/23/23
200	8888	2/23/23

ECMU·DB 15-445/645 (Spring 2023)

NAÏVE NESTED LOOP JOIN

Why is this algorithm bad? \rightarrow For every tuple in **R**, it scans **S** once **Cost: M + (m \cdot N)**

R(id,name)

id	name	
600	MethodMan	
200	GZA	
100	Andy	
300	ODB	
500	RZA	
700	Ghostface	
400	Raekwon	

S(id,value,cdate)

id	value	cdate
100	2222	2/23/23
500	7777	2/23/23
400	6666	2/23/23
100	9999	2/23/23
200	8888	2/23/23

N pages n tuples 16

M pages **m** tuples

ECMU·DB 15-445/645 (Spring 2023)

NAÏVE NESTED LOOP JOIN

Example database:

→ Table R: M = 1000, m = 100,000 - 4 KB pages → 6 MB → Table S: N = 500, n = 40,000

Cost Analysis:

 $\rightarrow M + (m \cdot N) = 1000 + (100000 \cdot 500) = 50,001,000 \text{ IOs}$ \rightarrow At 0.1 ms/IO, Total time \approx 1.3 hours

What if smaller table (S) is used as the outer table? $\rightarrow N + (n \cdot M) = 500 + (40000 \cdot 1000) = 40,000,500 \text{ IOs}$ \rightarrow At 0.1 ms/IO, Total time \approx 1.1 hours

SECMU-DB 15-445/645 (Spring 2023



M pages m tuples

ECMU·DB 15-445/645 (Spring 2023) 600MethodMan200GZA100Andy300ODB500RZA700Ghostface400Raekwon

100	2222	2/23/23
500	7777	2/23/23
400	6666	2/23/23
100	9999	2/23/23
200	8888	2/23/23

N pages **n** tuples 18

This algorithm performs fewer disk accesses. \rightarrow For every block in **R**, it scans **S** once.

Cost: $M + ((\# blocks in \mathbf{R}) \cdot N)$

R(id, name)

id name 600 MethodMan GZA 200 100 Andy ODB 300 500 RZA 700 Ghostface 400 Raekwon

S(id, value, cdate)

id	value	cdate
100	2222	2/23/23
500	7777	2/23/23
400	6666	2/23/23
100	9999	2/23/23
200	8888	2/23/23

N pages n tuples

19

M pages *m* tuples

SECMU:DB 15-445/645 (Spring 2023)

The smaller table should be the outer table. We determine size based on the number of pages, not the number of tuples.

6 2 1 3 5 7 4 **M** pages *m* tuples Security CMU.DB 15-445/645 (Spring 2023)

R(id, name)

d	name	/
00	MethodMan	
00	GZA	
00	Andy	
00	ODB	
00	RZA	
00	Ghostface	
.00	Raekwon	/

S(id, value, cdate)

	id	value	cdate
	100	2222	2/23/23
	500	7777	2/23/23
/	400	6666	2/23/23
	100	9999	2/23/23
	200	8888	2/23/23

N pages n tuples

20

If we have **B** buffers available:

- \rightarrow Use **B-2** buffers for each block of the outer table.
- \rightarrow Use one buffer for the inner table, one buffer for output.

R(id, name)

id	name	/
600	MethodMan	
200	GZA	
100	Andy	
300	ODB	
500	RZA	
700	Ghostface	
400	Raekwon	/

S(id, value, cdate)

id	value	cdate
100	2222	2/23/23
500	7777	2/23/23
400	6666	2/23/23
100	9999	2/23/23
200	8888	2/23/23

N pages **n** tuples 22

M pages m tuples

ECMU·DB 15-445/645 (Spring 2023)



This algorithm uses **B-2** buffers for scanning **R**. Cost: $M + (\lceil M \mid (B-2) \rceil \cdot N)$

If the outer relation fits in memory (M < B-2): \rightarrow Cost: M + N = 1000 + 500 = 1500 I/Os \rightarrow At 0.1ms per I/O, Total time \approx 0.15 seconds

If we have B=102 buffer pages: \rightarrow Cost: $M + (\lceil M / (B-2) \rceil \cdot N) = 1000 + 10*500 = 6000 I/Os$ \rightarrow Or can switch inner/outer relations, giving us cost: 500 + 5*1000 = 5500 I/Os

EFCMU·DB 15-445/645 (Spring 2023)

NESTED LOOP JOIN

Why is the basic nested loop join so bad?
→ For each tuple in the outer table, we must do a sequential scan to check for a match in the inner table.

We can avoid sequential scans by using an index to find inner table matches.

 \rightarrow Use an existing index for the join.



SCMU·DB 15-445/645 (Spring 2023)

NESTED LOOP JOIN SUMMARY

27

Key Takeaways

- \rightarrow Pick the smaller table as the outer table.
- \rightarrow Buffer as much of the outer table in memory as possible.
- \rightarrow Loop over the inner table (or use an index).

Algorithms

- \rightarrow Naïve
- \rightarrow Block
- \rightarrow Index

ECMU·DB 15-445/645 (Spring 2023

28

Phase #1: Sort

- \rightarrow Sort both tables on the join key(s).
- \rightarrow You can use any appropriate sort algorithm
- → These phases are distinct from the sort/merge phases of an external merge sort, from the previous class

Phase #2: Merge

- \rightarrow Step through the two sorted tables with cursors and emit matching tuples.
- \rightarrow May need to backtrack depending on the join type.

ECMU-DB 15-445/645 (Spring 2023

```
sort R,S on join keys
cursor<sub>R</sub> < R<sub>sorted</sub>, cursor<sub>S</sub> < S<sub>sorted</sub>
while cursor<sub>R</sub> and cursor<sub>S</sub>:
    if cursor<sub>R</sub> > cursor<sub>S</sub>:
        increment cursor<sub>S</sub>
    if cursor<sub>R</sub> < cursor<sub>S</sub>:
        increment cursor<sub>R</sub> (and possibly
            backtrack cursor<sub>s</sub>)
    elif cursor<sub>R</sub> and cursor<sub>S</sub> match:
        emit
        increment cursor<sub>S</sub>
```

ECMU-DB 15-445/645 (Spring 2023

R(id, name)

id	name	
100	Andy	
200	GZA	
200	GZA	6
300	ODB	1:4
400	Raekwon	
500	RZA	
600	MethodMan	
700	Ghostface	

S(id,value,cdate)				
X	id	value	cdate	\ni
	100	2222	2/23/23	
	100	9999	2/23/23	
	200	8888	2/23/23	
	400	6666	2/23/23	
	500	7777	2/23/23	

```
Sort!
```

SELECT R.id, S.cdate
FROM R JOIN S
ON R.id = S.id
WHERE S.value > 100

Output Buffer

R.id	R.name	S.id	S.value	S.cdate
100	Andy	100	2222	2/23/23
100	Andy	100	9999	2/23/23
200	GZA	200	8888	2/23/23
200	GZA	200	8888	2/23/23
400	Raekwon	200	6666	2/23/23
500	RZA	500	7777	2/23/23

SCMU·DB 15-445/645 (Spring 2023)

Sort!

31

Sort Cost (R): $2M \cdot (1 + \lceil \log_{B-1} \lceil M / B \rceil \rceil)$ Sort Cost (S): $2N \cdot (1 + \lceil \log_{B-1} \lceil N / B \rceil \rceil)$ Merge Cost: (M + N)

Total Cost: Sort + Merge

ECMU·DB 15-445/645 (Spring 2023)

Example database:

- \rightarrow Table R: M = 1000, m = 100,000
- → Table S: N = 500, n = 40,000

With B=100 buffer pages, both R and S can be sorted in two passes:

- \rightarrow Sort Cost (**R**) = 2000 · (1 + $[\log_{99} 1000 / 100]) = 4000 \text{ I/Os}$
- \rightarrow Sort Cost (S) = 1000 · (1 + $\lceil \log_{99} 500 / 100 \rceil)$ = 2000 I/Os
- \rightarrow Merge Cost = (1000 + 500) = **1500 I/Os**
- → Total Cost = 4000 + 2000 + 1500 = **7500 I/Os**
- \rightarrow At 0.1 ms/IO, Total time \approx 0.75 seconds

ECMU·DB 15-445/645 (Spring 2023)

The worst case for the merging phase is when the join attribute of all the tuples in both relations contains the same value.

Cost: $(M \cdot N)$ + (sort cost)

ECMU·DB 15-445/645 (Spring 2023

WHEN IS SORT-MERGE JOIN USEFUL?

One or both tables are already sorted on join key. Output must be sorted on join key.

The input relations may be sorted either by an explicit sort operator, or by scanning the relation using an index on the join key.

HASH JOIN

If tuple $r \in R$ and a tuple $s \in S$ satisfy the join condition, then they have the same value for the join attributes.

If that value is hashed to some partition \mathbf{i} , the R tuple must be in $\mathbf{r}_{\mathbf{i}}$ and the S tuple in $\mathbf{s}_{\mathbf{i}}$.

Therefore, **R** tuples in \mathbf{r}_i need only to be compared with **S** tuples in \mathbf{s}_i .

ECMU-DB 15-445/645 (Spring 2023)

SIMPLE HASH JOIN ALGORITHM

Phase #1: Build

 \rightarrow Scan the outer relation and populate a hash table using the hash function h_1 on the join attributes.

Phase #2: Probe

 \rightarrow Scan the inner relation and use h_1 on each tuple to jump to a location in the hash table and find a matching tuple.



HASH TABLE CONTENTS

Key: The attribute(s) that the query is joining on
→ The hash table needs to store the key to verify that we have a correct match, in case of hash collisions.

Value: It varies

- → Depends on what the next query operators will do with the output from the join
- \rightarrow Early vs. Late Materialization

OPTIMIZATION: PROBE FILTER

Create a probe filter (such as a <u>Bloom Filter</u>) during the build phase if the key is likely to not exist in the inner relation

 \rightarrow Check the filter before probing the hash table

 \rightarrow Fast because the filter fits in CPU cache



ECMU·DB 15-445/645 (Spring 2023

BLOOM FILTERS

Uses a bitmap to probabilistically answer set membership queries

 \rightarrow False negatives will never occur

 \rightarrow False positives can sometimes occur

Insert(x):

 \rightarrow Use k hash functions to set bits in the filter to 1

Lookup(x):

 \rightarrow Check whether the bits are 1 for each hash function

See the **Bloom Filter Calculator** if you build one

ECMU-DB 15-445/645 (Spring 2023)

BLOOM FILTERS



Insert('RZA')

Insert('GZA')

 $Lookup(RZA') \rightarrow TRUE$

Lookup('Raekwon') → FALSE

Lookup('ODB')**→ TRUE**

ECMU-DB 15-445/645 (Spring 2023)

HASH JOINS OF LARGE RELATIONS

What happens if we do not have enough memory to fit the entire hash table?

We do not want to let the buffer pool manager swap out the hash table pages at random.

PARTITIONED HASH JOIN

Hash join when tables do not fit in memory.

- → **Partition Phase:** Hash both tables on the join attribute into partitions.
- → **Probe Phase:** Compares tuples in corresponding partitions for each table.

Sometimes called **GRACE Hash** Join.

 \rightarrow Named after the GRACE <u>database</u> <u>machine</u> from Japan in the 1980s.



GRACE University of Tokyo

ECMU·DB 15-445/645 (Spring 2023



PARTITIONED HASH JOIN PARTITION PHASE

Hash R into k buckets.

Hash S into k buckets with same hash function.

Write buckets to disk when they get full.



PARTITIONED HASH JOIN PROBE PHASE

Read corresponding partitions into memory one pair at a time, hash join their contents.

PARTITIONED HASH JOIN EDGE CASES

If a partition does not fit in memory, recursively partition it with a different hash function

- \rightarrow Repeat as needed
- \rightarrow Eventually hash join the corresponding (sub-)partitions

If a single join key has so many matching records that they don't fit in memory, use a block nested loop join for that key

SCMU·DB 15-445/645 (Spring 2023)

RECURSIVE PARTITIONING

49

ANALYSIS OF PARTITIONED HASH JOIN

50

How big a table can be joined without recursive partitioning?

 \rightarrow Up to **B-1** partitions

 \rightarrow Each could be about as big as *B***-2** pages

Answer: About (B-1) · (B-2) pages

- \rightarrow If the partitions are approximately equal size, a table of
 - **N** pages needs about **sqrt(N)** buffers
- \rightarrow In practice, use a "fudge factor" f > 1: sqrt($f \cdot N$)
- \rightarrow Only partitions of the outer table need to fit in memory

ECMU·DB 15-445/645 (Spring 2023)

COST OF PARTITIONED HASH JOIN

If we don't need recursive partitioning: \rightarrow Cost: 3(M + N)

Partition phase:

- \rightarrow Read+write both tables
- → **2(M+N)** I/Os

Probe phase:

→ Read both tables (in total, one partition at a time) → M+N I/Os

ECMU-DB 15-445/645 (Spring 2023)

PARTITIONED HASH JOIN

Example database: $\rightarrow M = 1000, m = 100,000$ $\rightarrow N = 500, n = 40,000$

Cost Analysis: \rightarrow 3 · (M + N) = 3 · (1000 + 500) = 4,500 IOs \rightarrow At 0.1 ms/IO, Total time \approx 0.45 seconds

ECMU·DB 15-445/645 (Spring 2023)

OPTIMIZATION: HYBRID HASH JOIN Use some buckets for a simple in-memory hash join, have some buckets spill to disk. R(id, name) S(id,value,cdate) 0 2 k-1 9 **SECMU-DB** 15-445/645 (Spring 2023)

53

HASH JOIN OBSERVATIONS

The inner table can be any size. \rightarrow Only outer table (or its partitions) need to fit in memory

If we know the size of the outer table, then we can use a static hash table.

 \rightarrow Less computational overhead

If we do not know the size, then we must use a dynamic hash table or allow for overflow pages.

ECMU-DB 15-445/645 (Spring 2023

JOIN ALGORITHMS: SUMMARY

	Algorithm	IO Cost	Example
	Naïve Nested Loop Join	<i>M</i> + (<i>m</i> ⋅ <i>N</i>)	1.3 hours
	Block Nested Loop Join	<i>M</i> + (0.55 seconds
	Index Nested Loop Join	<i>M</i> + (<i>m</i> ⋅ C)	Variable
)	Sort-Merge Join	M + N + (sort cost)	0.75 seconds
	Hash Join	3 · (M + N)	0.45 seconds

CONCLUSION

Hashing is almost always better than sorting for operator execution.

Caveats:

- \rightarrow Sorting is better on non-uniform data.
- \rightarrow Sorting is better when result needs to be sorted.

Good DBMSs use either (or both).

57 **NEXT CLASS** Composing operators together to execute queries. SECMU.DB 15-445/645 (Spring 2023)