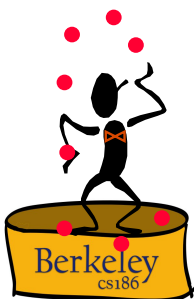


Join Algorithms

R&G 14.4



Rendezvous



- Grouping/Aggregation is one kind of rendezvous
 - groups of matching items within a single table
- Join is the other main kind of rendezvous
 - combinations of items from multiple tables

Cross Product



Given two collections R and S:

$R \times S \Rightarrow$ all pairs $\{r, s\}$ of items in R, S
– a.k.a Cartesian product

Can always do this to create a single “supertable” and then use a single-table query
– correct but rarely best performing

“Theta” Join



- $R \bowtie_{\theta} S$: all pairs $\{r, s\}$ where $\theta(r, s)$
 - e.g. FriendRequests \bowtie_{θ} Users
 - θ is “friendID = ID”
 - e.g. Family \bowtie_{θ} Family
 - θ is “age < age”
- A common case: EquiJoin
 - i.e., θ is an equality test
 - special case: one side of $=$ is a “key”
 - E.g. Enrolled.studentID = Students.ID
 - This is like doing “lookups” into the Students table

Schema for Examples



Sailors (sid: integer, sname: string, rating: integer, age: real)
Reserves (sid: integer, bid: integer, day: dates, name: string)

- Sailors:
 - Each tuple is 50 bytes long, 80 tuples per page, 500 pages.
 - [S]=500, $p_S=80$.
- Reserves:
 - Each tuple is 40 bytes, 100 tuples per page, 1000 pages.
 - [R]=1000, $p_R=100$.

Joins

```
SELECT *
FROM   Reserves R1, Sailors S1
WHERE  R1.sid=S1.sid
```



- Joins are very common.
- $R \times S$ is large
 - so, $R \times S$ followed by a “filter” is inefficient.
- Many approaches to reduce join cost.
- Join techniques we will cover today:
 - Nested-loops join
 - Index-nested loops join
 - Sort-merge join
 - Hash Joins

Some Cost Notation



- $[R]$: the number of pages to store R
- p_R : number of records per page of R
- $|R|$: the number of records in R
– cardinality
- Note: $p_R * [R] = |R|$

Simple Nested Loops Join



$R \bowtie S$: foreach **record** r in R do
 foreach **record** s in S do
 if $\theta(r_i, s_j)$ then add $\langle r, s \rangle$ to result

- Cost = $(p_R * [R]) * [S] + [R] = 100 * 1000 * 500 + 1000$ IOs
– At 10ms/IO, Total time: ???
- What if smaller relation (S) was “outer”?
- What assumptions are being made here?
- What is cost if one relation can fit entirely in memory?

Page-Oriented NestLoop Join



$R \bowtie S$: foreach **page** b_R in R do
 foreach **page** b_S in S do
 foreach **record** r in b_R do
 foreach **record** s in b_S do
 if $\theta(r_i, s_j)$ then add $\langle r, s \rangle$ to result

- Cost = $[R] * [S] + [R] = 1000 * 500 + 1000$
- If smaller relation (S) is outer, cost = $500 * 1000 + 500$
- Much better than naïve per-tuple approach!

Page-Oriented NestLoop Join



$R \bowtie S$: foreach **page** b_R in R do
 foreach **page** b_S in S do
 foreach **record** r in b_R do
 foreach **record** s in b_S do
 if $\theta(r_i, s_j)$ then add $\langle r, s \rangle$ to result

I/O for S

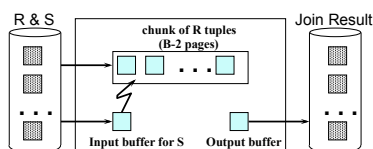
I/O for R

- Cost = $[R] * [S] + [R] = 1000 * 500 + 1000$
- If smaller relation (S) is outer, cost = $500 * 1000 + 500$
- Much better than naïve per-tuple approach!

Block Nested Loops Join



- Page-oriented NL doesn't exploit extra buffers :(
- Idea to use memory efficiently:

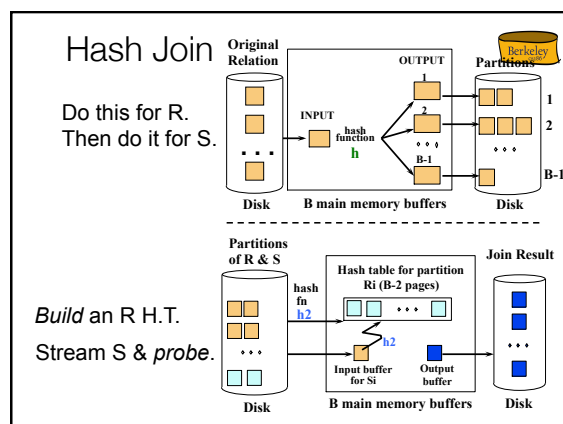
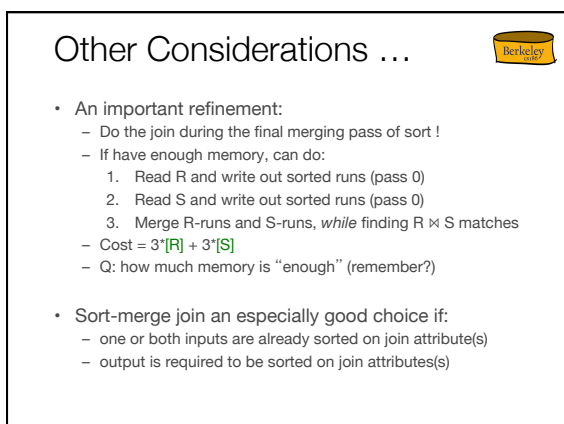
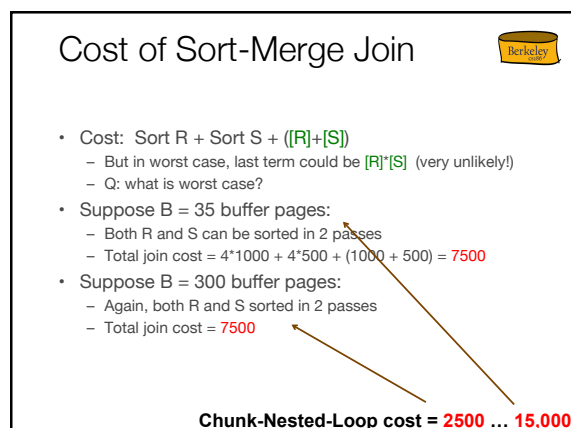
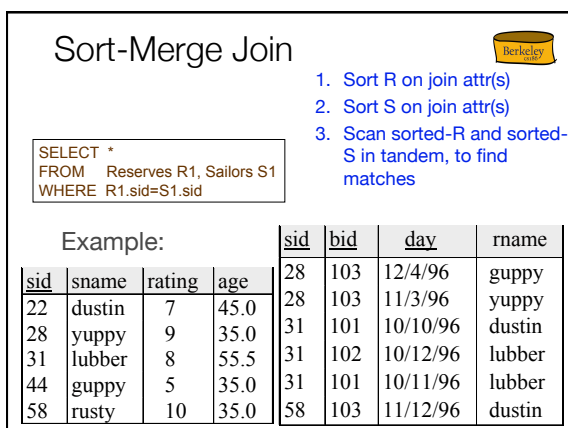
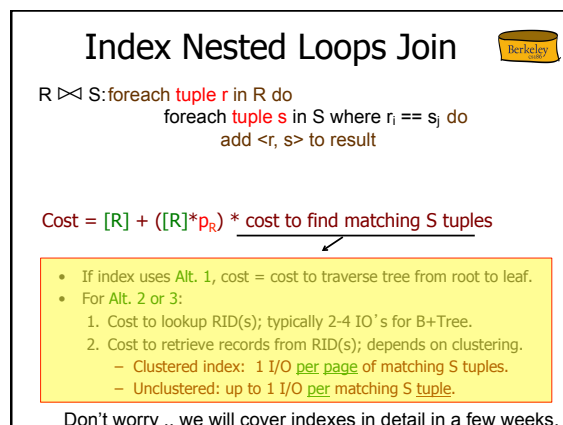
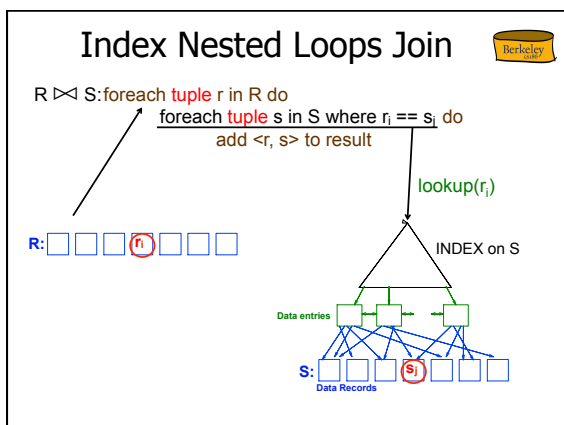


Cost: Scan outer + (#outer chunks * scan inner)
#outer chunks = $\lceil \text{outer} / \text{chunksize} \rceil$

Block NestLoop Examples



- Say we have $B = 100 + 2$ memory buffers
- Join cost = $\text{[outer]} + (\text{outer-chunks} * \text{[inner]})$
– #outer chunks = $\text{[outer]} / 100$
- With R as outer ($[R] = 1000$):
– Scanning R costs 1000 IO's (done in 10 chunks)
– Per chunk of R , we scan S ; costs $10 * 500$ IOs
– Total = $1000 + 10 * 500$.
- With S as outer ($[S] = 500$):
– Scanning S costs 500 IO's (done in 5 chunks)
– Per chunk of S , we scan R ; costs $5 * 1000$ IO's
– Total = $500 + 5 * 1000$.



Cost of Hash Join



- Partitioning phase: read+write both relations
 $\Rightarrow 2([R]+[S])$ I/Os
- Matching phase: read both relations, write output
 $\Rightarrow [R]+[S] + [\text{output}]$ I/Os
- Total cost of 2-pass hash join = $3([R]+[S])+[\text{output}]$

Q: what is cost of 2-pass *sort-merge join*?

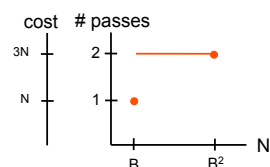
Q: how much memory needed for 2-pass *sort-merge join*?

Q: how much memory needed for 2-pass *hash join*?

Exploit excess memory



- Have B memory buffers
- Want to hash relation of size N blocks

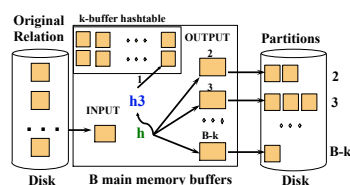


If $B < N < B^2$, will have unused memory ...

Hybrid Hashing

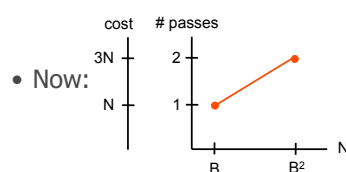


Idea: keep one of the hash buckets in memory!



Q: how do we choose the value of k ?

Cost savings: hybrid hashing



• Now:

Hash Join vs. Sort-Merge Join



- Sorting pros:
 - Good if input already sorted, or need output sorted
 - Not sensitive to data skew or bad hash functions
- Hashing pros:
 - Can be cheaper due to hybrid hashing
 - For join: # passes depends on *size of smaller relation*
 - Good if input already hashed, or need output hashed

Recap



- Nested Loops Join
 - Works for arbitrary Θ
 - Make sure to utilize memory in "chunks"
- Index Nested Loops
 - For equi-joins
 - When you already have an index on one side
- Sort/Hash
 - For equi-joins
 - No index required
- No clear winners – may want to implement them all
- Be sure you know the cost model for each
 - You will need it for query optimization!