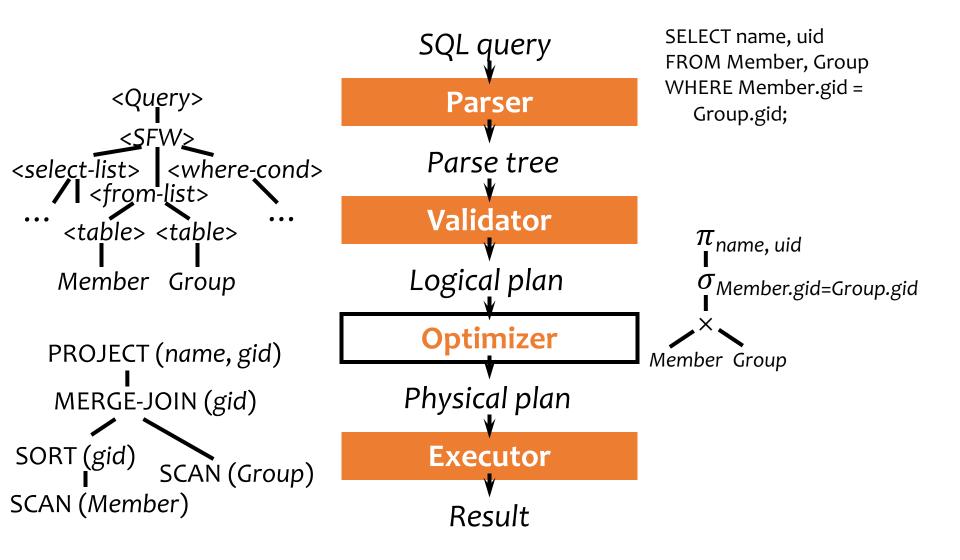
Lecture 18: Query Processing & Optimization

CS348 Spring 2025: Introduction to Database Management

Instructor: Xiao Hu

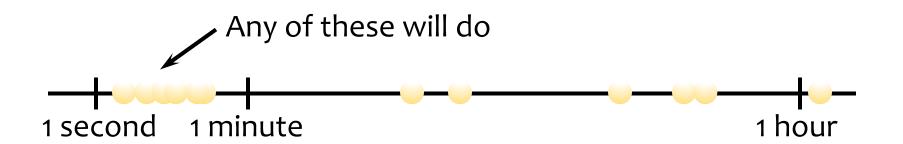
Sections: 001, 002, 003

A query's trip through the DBMS



(Recap) Query optimization

- Why query optimization?
- Search space
 - What are the possible equivalent logical plans? (lecture 17)
 - What are the possible physical plans? (lectures 16-17)
- Search strategy
 - Rule-based strategy
 - Cost-estimation-based strategy



Search strategy

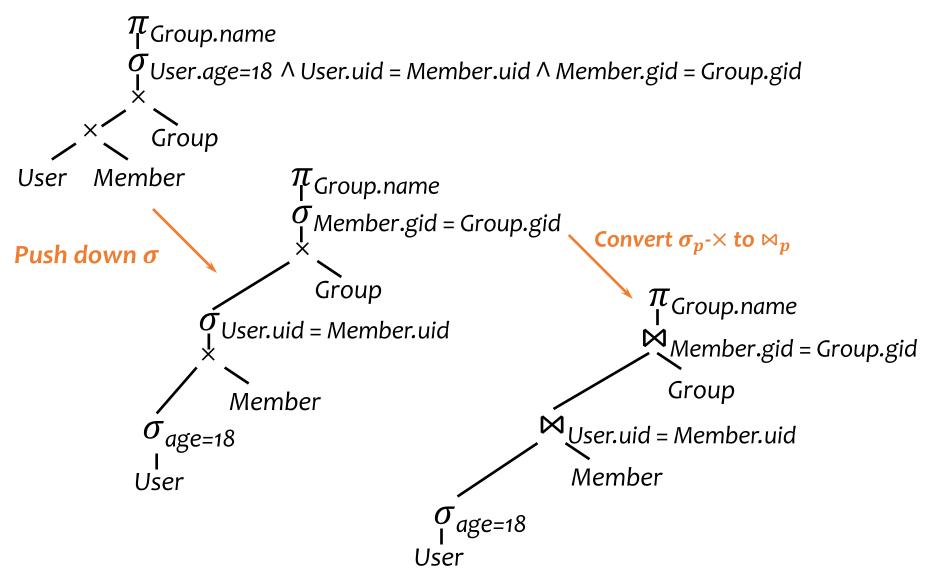


Rule-based query optimization

- Push $\sigma/\pi/-$ down as much as possible
 - $\sigma_{p \wedge p_R \wedge p_S}(R \bowtie_{p'} S) = (\sigma_{p_R} R) \bowtie_{p \wedge p'} (\sigma_{p_S} S)$
 - $\pi_L(\sigma_p R) = \pi_L(\sigma_p(\pi_{L \cup L'} R))$
 - $\pi_L(R \bowtie_p S) = \pi_L((\pi_{L_R}R) \bowtie_p (\pi_{L_S}S))$
 - $(R \bowtie S) (T \bowtie S) = (R T) \bowtie S$
 - $(R \bowtie S) (T \bowtie W) = ((R T) \bowtie S) \cup (R \bowtie (S W))$
- Join smaller relations first and avoid cross product
- (Many other rules to be further exploited)

Why? Reduce the size of intermediate results

An Example



From rule-based to cost-based opt.

- Rule-based optimization
 - Apply algebraic equivalence to rewrite plans into cheaper ones
- Cost-based optimization
 - Rewrite logical plan to combine "blocks" as much as possible
 - Optimize query block by block
 - Enumerate logical plans (already covered)
 - Estimate the cost of the plans
 - Pick a plan with an acceptable cost
 - Focus: select-project-join blocks



"Selinger"-style query optimization \leftarrow

Patricia Selinger

Cost estimation

Physical plan example:

INDEX-NESTED-LOOP-JOIN (gid)

Index on Group(gid)

Index on Member(uid)

What is its input size?
How many users with age 18?

INDEX-SCAN (age = 18)
Index on User(age)

• We have I/O cost formula for each operator

Lectures 16-17

- Example: INDEX-NESTED-LOOP-JOIN (uid) takes $O(B(R) + |R| \cdot lookup + fetch)$
- We need the size of the intermediate results

Cardinality Estimation



Selections with equality predicates

Consider $\sigma_{A=v}R$

- DBMSs typically store the following in the catalog
 - Size of *R*: |*R*|
 - Number of distinct A values in R: $|\pi_A R|$
- Assumption of uniformity: A-values are uniformly distributed in tuples from *R*
- $|\sigma_{A=v}R| \approx \frac{|R|}{|\pi_{A}R|}$
 - Selectivity factor of (A = v) is $\frac{1}{|\pi_A R|}$
 - Selectivity: the probability that any row will satisfy a predicate

Conjunctive predicates

Consider $\sigma_{A=u \land B=v}R$

- Assumption of selection independence: (A = u) and (B = v) independently select tuple in R
 - Counterexample: major and advisor, or *A* is the key
- $|\sigma_{A=u \wedge B=v}R| \approx \frac{|R|}{|\pi_{A}R| \cdot |\pi_{B}R|}$
 - Selectivity factor of (A = u) is $\frac{1}{|\pi_A R|}$
 - Selectivity factor of (B = v) is $\frac{1}{|\pi_B R|}$
 - Selectivity factor of $(A = u) \wedge (B = v)$ is $\frac{1}{|\pi_A R| \cdot |\pi_B R|}$
 - Reduce total size by all selectivity factors

Negated and disjunctive predicates

Consider $\sigma_{A\neq v}R$

- $|\sigma_{A\neq v}R| \approx |R| \cdot (1 \frac{1}{|\pi_A R|})$
 - Selectivity factor of $\neg p$ is (1 selectivity factor of p)

Consider $\sigma_{A=u \vee B=v}R$

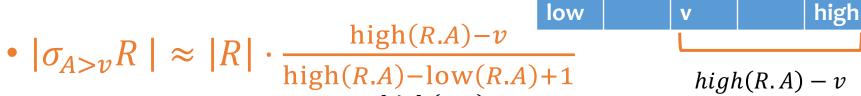
- $|\sigma_{A=u \vee B=v}R| \approx |R| \cdot (1/|\pi_{AR}| + 1/|\pi_{BR}|)$?
 - Tuples satisfying (A = u) and (B = v) are counted twice!
- $|\sigma_{A=u \vee B=v}R| \approx |R| \cdot (1/|\pi_{AR}| + 1/|\pi_{BR}| 1/|\pi_{AR}||\pi_{BR}|)$
 - Inclusion-exclusion principle

Range predicates

Consider $\sigma_{A>v}R$

- DBMSs typically store the following in the catalog
 - Largest R.A value: high(R.A)
 - Smallest R.A value: low(R.A)

high(R.A) - low(R.A)



• Selectivity factor is $\frac{\text{high}(R.A) - v}{\text{high}(R.A) - \text{low}(R.A)}$

Two-way natural join

- $Q = R(A, B) \bowtie S(A, C)$
- Assumption of containment of value sets: every tuple in the "smaller" relation (one with fewer distinct values for the join attribute) joins with some tuple in the other relation
 - That is, if $|\pi_A R| \leq |\pi_A S|$ then $\pi_A R \subseteq \pi_A S$
 - Certainly not true in general
 - But holds many practical cases
- $|Q| \approx \frac{|R| \cdot |S|}{\max(|\pi_A R|, |\pi_A S|)}$
 - Selectivity factor of R.A = S.A is $\frac{1}{\max(|\pi_A R|, |\pi_A S|)}$

Multiway natural join

- $Q: R(A,B) \bowtie S(B,C) \bowtie T(C,D)$
- What is the number of distinct *C* values in the join of *R* and *S*?
- Assumption of preservation of value sets
 - A non-join attribute does not lose values from its set of possible values
 - That is, if C is in S but not R, then $\pi_C(R \bowtie S) = \pi_C S$
 - Certainly not true in general
 - But holds many practical cases

Multiway natural join (cont'd)

- $Q: R(A,B) \bowtie S(B,C) \bowtie T(C,D)$
- Reduce the total size by the selectivity factor of each join predicate
 - $R.B = S.B: \frac{1}{\max(|\pi_B R|, |\pi_B S|)}$
 - $|R \bowtie S| = \frac{|R| \cdot |S|}{\max(|\pi_B R|, |\pi_B S|)}$
 - $(R \bowtie S).C = T.C: \frac{1}{\max(|\pi_C(R \bowtie S)|, |\pi_C T|)} = \frac{1}{\max(|\pi_C S|, |\pi_C T|)}$

•
$$|Q| \approx \frac{|R| \cdot |S| \cdot |T|}{\max(|\pi_B R|, |\pi_B S|) \cdot \max(|\pi_C S|, |\pi_C T|)}$$

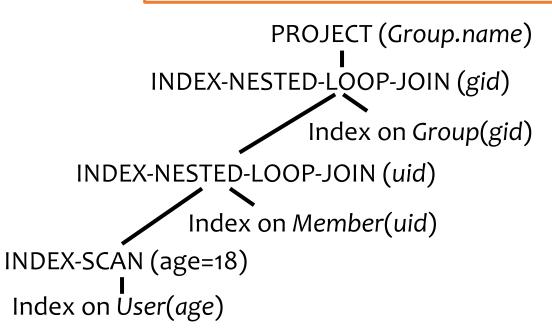
Summary of Cardinality Estimation

- Lots of assumptions and very rough estimation
 - An accurate estimator is not needed
 - Maybe okay if we overestimate or underestimate, since it may not change the query plan selection
- (MUCH) Better techniques
 - Histograms
 - Sketches
 - Machine learning approaches (one of the most popular topics in database research now)

User (<u>uid</u> int, name string, age int, pop float) Group (<u>gid</u> string, name string) Member (<u>uid</u> int, <u>gid</u> string)

Data statistics:

- |User| = 1000
- |Member| = 50000
- |Group| = 100
- $|\pi_{name}User| = 50$
- $|\pi_{uid}Member| = 500$
- $|\pi_{gid}Member| = 100$
- $|\pi_{qid}Group| = 100$



What are the sizes of intermediate join results?

User (<u>uid</u> int, name string, age int, pop float) Group (<u>gid</u> string, name string) Member (<u>uid</u> int, <u>gid</u> string)

Data statistics:

- |User| = 1000
- |Member| = 50000
- |Group| = 100
- $|\pi_{age}User| = 50$
- $|\pi_{uid}Member| = 500$
- $|\pi_{aid}Member| = 100$
- $|\pi_{gid}Group| = 100$

PROJECT (Group.name)

INDEX-NESTED-LOOP-JOIN (gid)

Index on Group(gid)

INDEX-NESTED-LOOP-JOIN (uid)

Index on Member(uid)

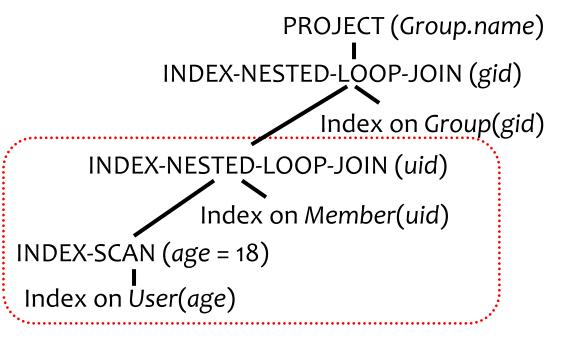
INDEX-SCAN (age = 18)

Index on User(age)

- What is the intermediate size?
 - Assumption of uniformity
 - $|\sigma_{age=18}(User)| \approx \frac{1000}{50} = 20$

User (<u>uid</u> int, name string, age int, pop float) Group (<u>gid</u> string, name string) Member (<u>uid</u> int, <u>gid</u> string)

- |User| = 1000
- |Member| = 50000
- |Group| = 100
- $|\pi_{age}User| = 50$
- $|\pi_{uid}Member| = 500$
- $|\pi_{aid}Member| = 100$
- $|\pi_{qid}Group| = 100$



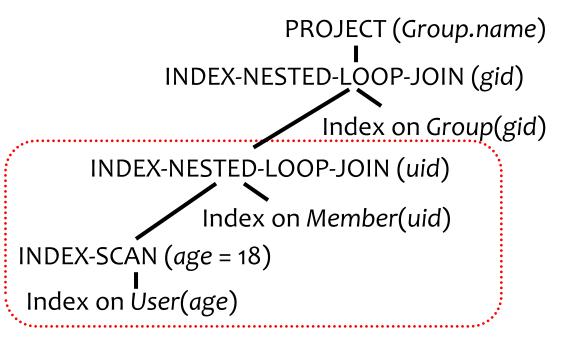
- What is the intermediate join size?
 - Assume uniformity of age in User
 - Assume containment of value sets in User and Member on uid

•
$$|\sigma_{age=18}(User)| \bowtie Member| \approx \frac{|\sigma_{age=18}(User)| \cdot |Member|}{\max(|\pi_{uid}\sigma_{age=18}(User)|, |\pi_{uid}Member|)}$$

$$= \frac{20.50000}{\max(20,500)} = 2000$$

User (<u>uid</u> int, name string, age int, pop float) Group (<u>gid</u> string, name string) Member (<u>uid</u> int, <u>gid</u> string)

- |User| = 1000
- |Member| = 50000
- |Group| = 100
- $|\pi_{age}User| = 50$
- $|\pi_{uid}Member| = 500$
- $|\pi_{aid}Member| = 100$
- $|\pi_{qid}Group| = 100$



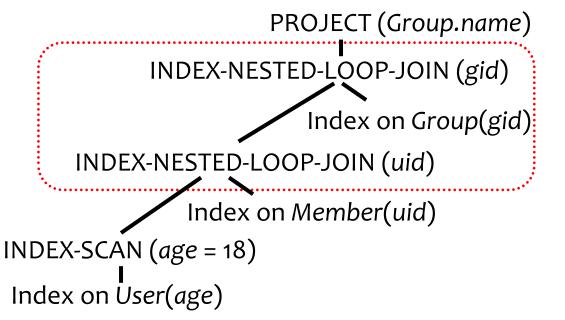
- What is the intermediate join size? (an alternative estimate)
 - Assume the instance follows the foreign-key constraint
 - Assume that within each group, users of different ages join with equal probability

•
$$|\sigma_{age=18}(User) \bowtie Member| = |\sigma_{age=18}(User \bowtie Member)|$$

 $\approx \frac{1}{50} \cdot 50000 = 1000$

User (<u>uid</u> int, name string, age int, pop float) Group (<u>gid</u> string, name string) Member (<u>uid</u> int, <u>gid</u> string)

- |User| = 1000
- |Member| = 50000
- |Group| = 100
- $|\pi_{age}User| = 50$
- $|\pi_{uid}Member| = 500$
- $|\pi_{qid}Member| = 100$
- $|\pi_{gid}Group| = 100$

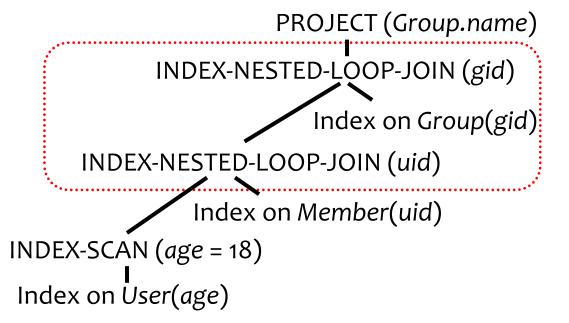


- What is the intermediate join size?
 - Assumption of preservation of value sets
 - $Q' = \sigma_{age=18}(User) \bowtie Member$ (suppose we use estimator 2000)

•
$$|Q' \bowtie Group| \approx \frac{|Q'| \cdot |Group|}{\max(|\pi_{gid}Q'|, |\pi_{gid}Group|)} = \frac{2000 \cdot 100}{\max(100, 100)} = 2000$$

User (<u>uid</u> int, name string, age int, pop float) Group (<u>gid</u> string, name string) Member (<u>uid</u> int, <u>gid</u> string)

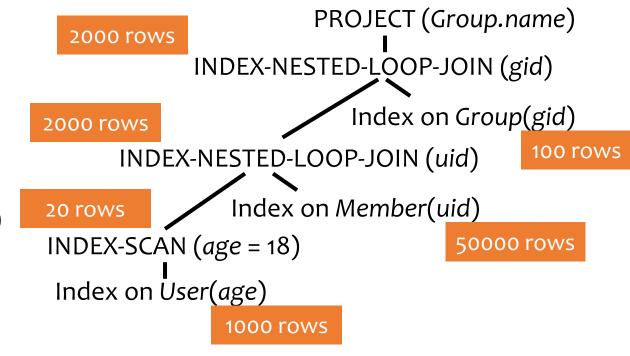
- |User| = 1000
- |Member| = 50000
- |Group| = 100
- $|\pi_{age}User| = 50$
- $|\pi_{uid}Member| = 500$
- $|\pi_{qid}Member| = 100$
- $|\pi_{qid}Group| = 100$



- What is the intermediate join size? (an alternative estimate)
 - Assume the instance follows the foreign-key constraint
 - $Q' = \sigma_{age=18}(User) \bowtie Member$ (suppose we use estimator 2000)
 - $|Q' \bowtie Group| \approx |Q'| = 2000$

User (<u>uid</u> int, name string, age int, pop float) Group (<u>gid</u> string, name string) Member (<u>uid</u> int, <u>gid</u> string)

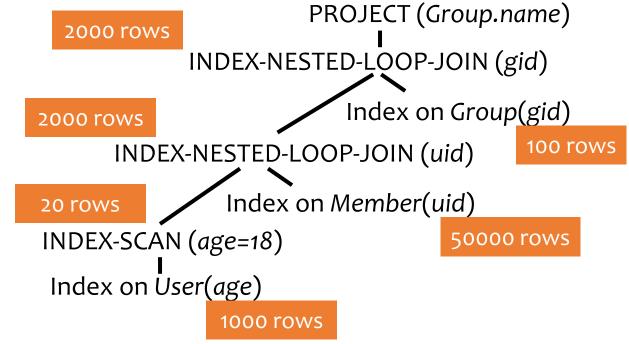
- |User| = 1000
- |Member| = 50000
- |Group| = 100
- $|\pi_{age}User| = 50$
- $|\pi_{uid}Member| = 500$
- $|\pi_{qid}Member| = 100$
- $|\pi_{qid}Group| = 100$



- What are the estimated I/O costs? (Some rough calculation)
- System requirements:
 - Each disk/memory block can hold up to 10 rows (from any table);
 - All tables are stored compactly on disk (10 rows per block);
 - 8 memory blocks are available for query processing: M=8

User (<u>uid</u> int, name string, age int, pop float) Group (<u>gid</u> string, name string) Member (<u>uid</u> int, <u>gid</u> string)

- |User| = 1000
- |Member| = 50000
- |Group| = 100
- $|\pi_{age}User| = 50$
- $|\pi_{uid}Member| = 500$
- $|\pi_{qid}Member| = 100$
- $|\pi_{qid}Group| = 100$



- INDEX-SCAN on User
 - Index lookup: typically 4 I/Os
 - Suppose tuples in the User table are clustered by uid but not age
 - In the worst case, we need 20 I/Os to retrieve tuples
 - Cache the intermediate results in main memory using 2 blocks

User (<u>uid</u> int, name string, age int, pop float) Group (<u>gid</u> string, name string) Member (<u>uid</u> int, <u>gid</u> string)

Data statistics:

- |User| = 1000
- |Member| = 50000
- |Group| = 100
- $|\pi_{age}User| = 50$
- $|\pi_{uid}Member| = 500$
- $|\pi_{aid}Member| = 100$
- $|\pi_{aid}Group| = 100$

PROJECT (Group.name)

INDEX-NESTED-LOOP-JOIN (gid)

Index on Group(gid)

INDEX-NESTED-LOOP-JOIN (uid)

100 rows

Index on Member(uid)

INDEX-SCAN (age=18)

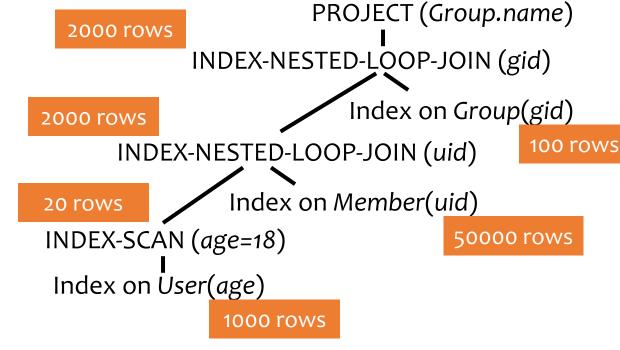
Index on User(age)

1000 rows

- INDEX-NESTED-LOOP-JOIN on $\sigma_{a,qe=18}(User)$ and Member
 - Suppose tuples in Member table clustered by uid
 - Index lookup: typically 4 I/Os
 - Probing takes $20 \cdot lookup + fetch = 20 * 4 + \frac{2000}{10} = 280 I/Os$
 - Writing the intermediate results back to disk takes $\frac{2000}{10} = 200 \text{ I/Os}$

User (<u>uid</u> int, name string, age int, pop float) Group (<u>gid</u> string, name string) Member (<u>uid</u> int, <u>gid</u> string)

- |User| = 1000
- |Member| = 50000
- |Group| = 100
- $|\pi_{age}User| = 50$
- $|\pi_{uid}Member| = 500$
- $|\pi_{aid}Member| = 100$
- $|\pi_{qid}Group| = 100$



- INDEX-NESTED-LOOP-JOIN on Q^\prime and Group
 - $Q' = \sigma_{age=18}(User) \bowtie Member$
 - For each tuple in Q', probe the index on Group(gid)
 - Suppose tuples in Group table are clustered by gid
 - It takes $B(Q') + |Q'| \cdot lookup + fetch = 200 + 2000 * 4 + 2000 I/Os$
 - When a join result is generated in memory, compute the projection on the fly

From rule-based to cost-based opt.

Rule-based optimization

Apply algebraic equivalence to rewrite plans into cheaper ones

- Cost-based optimization
 - Rewrite logical plan to combine "blocks" as much as possible
 - Optimize the query block by block
 - Enumerate logical plans (already covered)
 - Estimate the cost of the plans
 - Pick a plan with an acceptable cost
 - Focus: select-project-join blocks



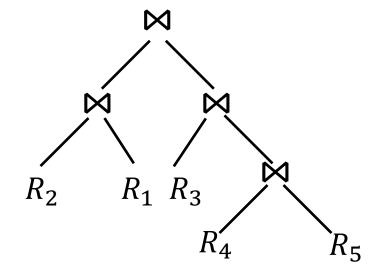
"Selinger"-style query optimization \leftarrow

Patricia Selinger

Search Good Join Ordering

Consider $R_1 \bowtie R_2 \bowtie \cdots \bowtie R_n$:

- Search space is Huge!
 - "Bushy" plan example:



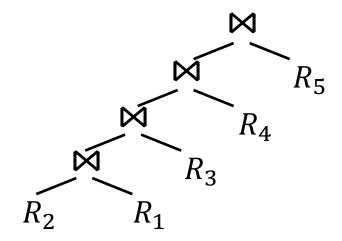
- Just considering different join orders, there are 30240 bushy plans for joining 6 tables!
- (there are more if we consider using different join algorithms)

It is a matter of intermediate join size

- For simplicity, we use the sort-merge join algorithm for demonstration
- $R \bowtie S$ takes $O\left(B(R) + B(S) + \frac{B(R\bowtie S)}{MB}\right)$ I/Os (see Slide 44 in Lecture 16)
- Intermediate join results are needed to be written back to disk, and serve as an input table for the subsequent join

How to minimize the intermediate join size?

Left-deep plans

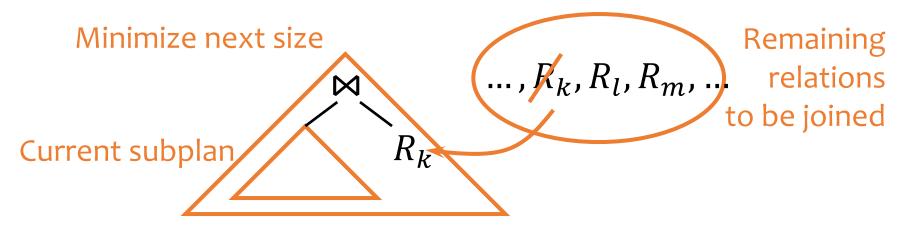


- Heuristic: consider only "left-deep" plans, in which only the left child can be a join
- How many left-deep plans are there for $R_1 \bowtie \cdots \bowtie R_n$?
 - Significantly fewer, but still lots—n! (720 for n=6)

A greedy algorithm for deep-left plans

Consider $R_1 \bowtie R_2 \bowtie \cdots \bowtie R_n$

- Start with the pair R_i , R_j with the smallest estimated size for $R_i \bowtie R_j$
- Repeat until no relation is left: Pick R_k from the remaining relations such that the join of R_k and the current result yields an intermediate result of the smallest estimated size



A DP algorithm for bushy plans

- Dynamic Programming for computing the optimal bushy plans in a bottom-up way:
 - Pass 1: Find the best single-table plans (for each table)
 - Pass 2: Find the best two-table plans (for each pair of tables) by combining best single-table plans
 - •
 - Pass k: Find the best k-table plans (for each combination of k tables) by combining two smaller best plans found in previous passes
 - ...
- Rationale: Any subplan of an optimal plan must also be optimal (otherwise, just replace the subplan to get a better overall plan)

Summary of Search Strategy

- Rule-based strategy
 - Apply algebraic equivalence to rewrite plans
 - But, we need more rewrite rules!
 - Existing rewrite plan is blind to the data statistics
- Cost-based strategy
 - Rewrite logical plan to combine "blocks" as much as possible
 - Optimize the query block by block
 - Enumerate logical plans (already covered)
 - Estimate the cost of the plans
 - Cardinality estimation is still challenging for multiple tables!
 - We need better end-to-end cost estimators
 - Pick a plan with an acceptable cost

A query's trip through the DBMS

