# Lecture 22:
# Query Execution

Monday, March 6, 2006
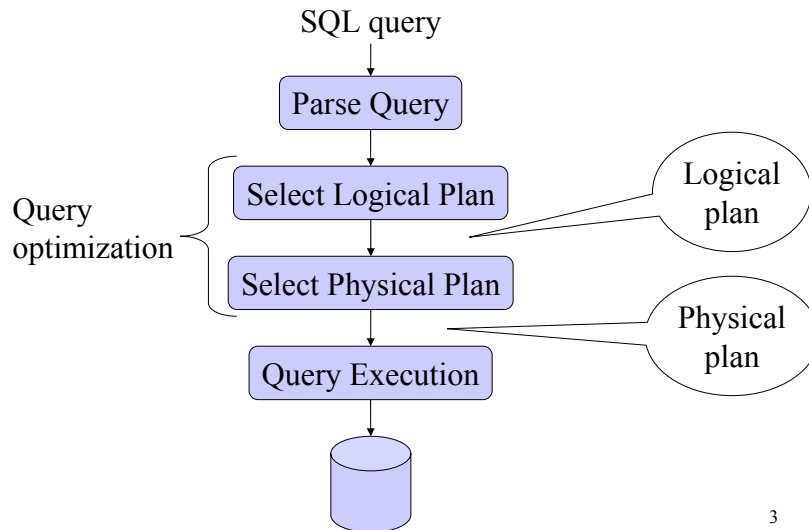
1

# Outline

- Query execution: 15.1 – 15.5

2

# Architecture of a Database Engine
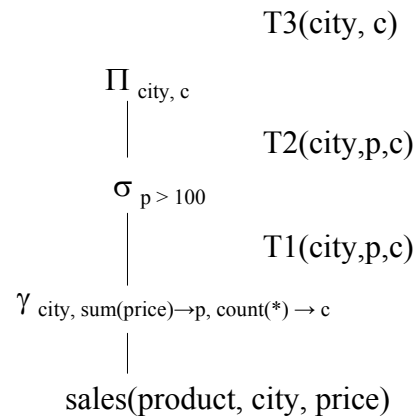
SQL query

↓

Parse Query

↓

Select Logical Plan ——— Logical plan

Query optimization

Select Physical Plan ——— Physical plan

↓

Query Execution

↓

3

# Logical Algebra Operators

- Union, intersection, difference
- Selection $\sigma$
- Projection $\Pi$
- Join $|x|$
- Duplicate elimination $\delta$
- Grouping $\gamma$
- Sorting $\tau$

4

# Logical Query Plan

SELECT city, count(*)
FROM sales
GROUP BY city
HAVING sum(price) > 100

T3(city, c)

$\Pi_{city, c}$

T2(city,p,c)

$\sigma_{p > 100}$

T1(city,p,c)

$\gamma_{city, \, sum(price) \to p, \, count(*) \to c}$
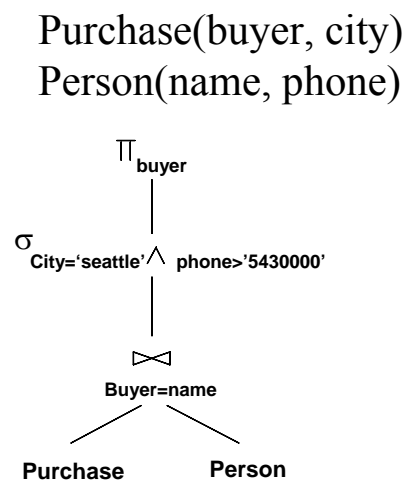
sales(product, city, price)

T1, T2, T3 = temporary tables

5

---

# Logical Query Plan

SELECT P.buyer
FROM Purchase P, Person Q
WHERE P.buyer=Q.name AND
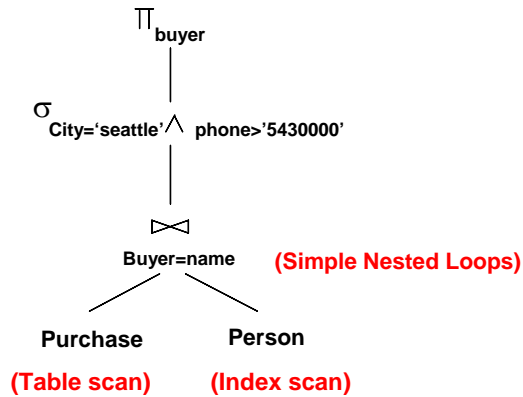    P.city='seattle' AND
    Q.phone > '5430000'

Purchase(buyer, city)
Person(name, phone)

$\Pi_{buyer}$

$\sigma_{City='seattle' \wedge \, phone>'5430000'}$

$\bowtie_{Buyer=name}$

Purchase        Person

6

3

# Physical Query Plan

SELECT  S.buyer
FROM Purchase P, Person Q
WHERE P.buyer=Q.name AND
    Q.city='seattle' AND
    Q.phone > '5430000'

$\Pi_{buyer}$

$\sigma_{City='seattle' \wedge phone>'5430000'}$

$\bowtie$ Buyer=name     **(Simple Nested Loops)**

**Purchase**     **Person**

**(Table scan)**     **(Index scan)**

Query Plan:
• logical tree
• implementation choice at every node
• scheduling of operations.

Some operators are from relational algebra, and others (e.g., scan) are not.

7

---

# Question in Class

Logical operator:
   **Product(pname, cname) |×| Company(cname, city)**

Propose three physical operators for the join, assuming the tables are in main memory:

1.
2.
3.

8

# Question in Class

**Product(pname, cname) |x| Company(cname, city)**

- 1000000 products
- 1000 companies

How much time do the following physical operators take if the data is **in main memory** ?

- Nested loop join          time =
- Sort and merge = merge-join    time =
- Hash join             time =

9

# Cost Parameters

The *cost* of an operation = total number of I/Os

    result assumed to be delivered in main memory

Cost parameters:

- B(R) = number of blocks for relation R
- T(R) = number of tuples in relation R
- V(R, a) = number of distinct values of attribute a
- M = size of main memory buffer pool, in blocks
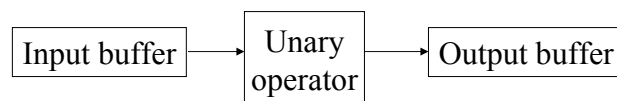
10

5

# Cost Parameters

- *Clustered* table R:
  - Blocks consists only of records from this table
  - $B(R) \ll T(R)$
- *Unclustered* table R:
  - Its records are placed on blocks with other tables
  - $B(R) \approx T(R)$

- When a is a key, $V(R,a) = T(R)$
- When a is not a key, $V(R,a)$

11

# Selection and Projection

Selection $\sigma(R)$, projection $\Pi(R)$
- Both are *tuple-at-a-time* algorithms
- Cost: $B(R)$

Input buffer $\longrightarrow$ Unary operator $\longrightarrow$ Output buffer

12

6

# Main Memory Hash Join

Hash join:  R |x| S
- Scan S, build buckets in main memory
- Then scan R and join

- Cost: B(R) + B(S)
- Assumption: B(S) <= M

13

# Duplicate Elimination

Duplicate elimination $\delta(R)$
- Hash table in main memory

- Cost: B(R)
- Assumption: B($\delta$(R)) <= M

14

# Grouping

Grouping:

Product(name, department, quantity)

$\gamma_{department,\ sum(quantity)}$ (Product) →
Answer(department, sum)

Main memory hash table

Question: How ?

15

# Nested Loop Joins

- Tuple-based nested loop R ⋈ S

```
for each tuple r in R do
  for each tuple s in S do
    if r and s join then output (r,s)
```

- Cost: T(R) B(S) when S is clustered
- Cost: T(R) T(S) when S is unclustered

16

# Nested Loop Joins

- We can be much more clever

- *Question*: how would you compute the join in the following cases ? What is the cost ?

  - B(R) = 1000, B(S) = 2, M = 4

  - B(R) = 1000, B(S) = 3, M = 4
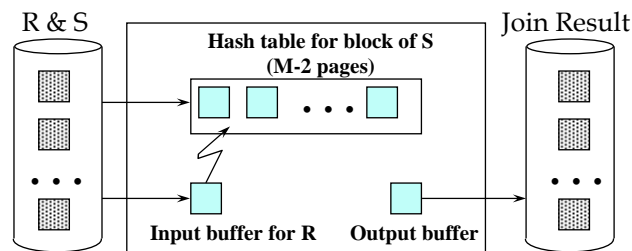
  - B(R) = 1000, B(S) = 6, M = 4

17

# Nested Loop Joins

- Block-based Nested Loop Join

for each (M-2) blocks bs of S do
  for each block br of R do
      for each tuple s in bs
          for each tuple r in br do
              if "r and s join" then output(r,s)

18

# Nested Loop Joins

R & S          **Hash table for block of S**          Join Result

**(M-2 pages)**
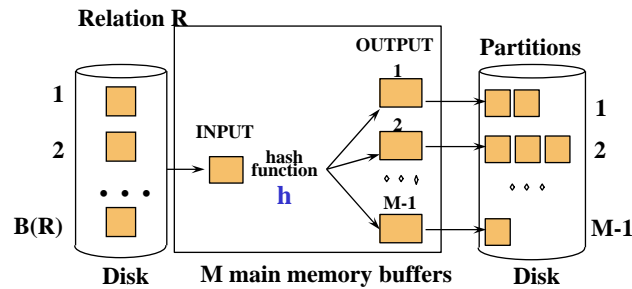
**Input buffer for R**     **Output buffer**

# Nested Loop Joins

- Block-based Nested Loop Join
- Cost:
  - Read S once: cost $B(S)$
  - Outer loop runs $B(S)/(M-2)$ times, and each time need to read R: costs $B(S)B(R)/(M-2)$
  - Total cost: $B(S) + B(S)B(R)/(M-2)$
- Notice: it is better to iterate over the smaller relation first
- R |x| S: R=outer relation, S=inner relation

# Partitioned Hash Algorithms

- Idea: partition a relation R into buckets, on disk
- Each bucket has size approx. B(R)/M

**Relation R**

**OUTPUT** **Partitions**

1

2

B(R)

**INPUT**

**hash function** **h**

**OUTPUT**
1
2
M-1

**Partitions**
1
2
M-1

**Disk**     **M main memory buffers**     **Disk**

- Does each bucket fit in main memory ?
  - Yes if B(R)/M <= M,   i.e. B(R) <= M$^2$

21

# Duplicate Elimination

- Recall:  $\delta(R)$ = duplicate elimination
- Step 1. Partition R into buckets
- Step 2. Apply $\delta$ to each bucket (may read in main memory)

- Cost: 3B(R)
- Assumption: B(R) <= M$^2$

22

# Grouping

- Recall:  $\gamma(R)$ = grouping and aggregation
- Step 1. Partition R into buckets
- Step 2. Apply $\gamma$ to each bucket (may read in main memory)

- Cost: 3B(R)
- Assumption: B(R) <= $M^2$
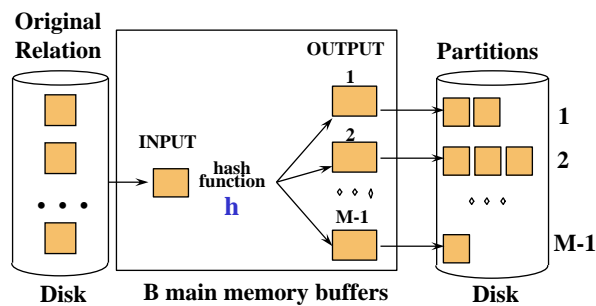
23

# Partitioned Hash Join

R |x| S
- Step 1:
  - Hash S into M buckets
  - send all buckets to disk
- Step 2
  - Hash R into M buckets
  - Send all buckets to disk
- Step 3
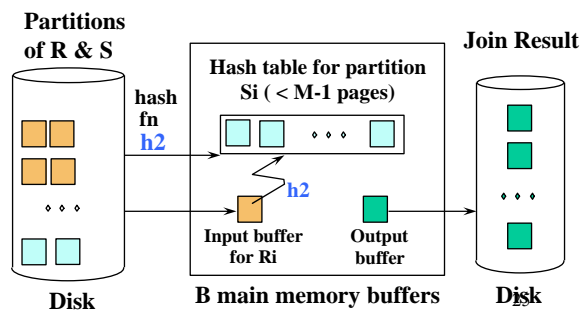  - Join every pair of buckets

24

# Hash-Join

**Original Relation**

- Partition both relations using hash fn **h**: R tuples in partition i will only match S tuples in partition i.

**OUTPUT**

**Partitions**

**INPUT**

**hash function h**

1
2
M-1

**Disk**

**B main memory buffers**

**Disk**

1
2
M-1

- ❖ Read in a partition of R, hash it using **h2 (<> h!)**. Scan matching partition of S, search for matches.

**Partitions of R & S**

**hash fn h2**

**Hash table for partition Si ( < M-1 pages)**

**h2**

**Input buffer for Ri**

**Output buffer**

**Join Result**

**Disk**

**B main memory buffers**

**Disk**

---

# Partitioned Hash Join

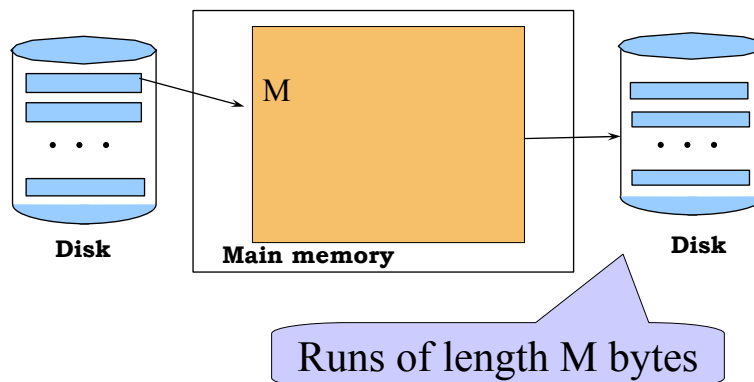- Cost: $3B(R) + 3B(S)$
- Assumption: $\min(B(R), B(S)) <= M^2$

26

# External Sorting

- Problem:
- Sort a file of size B with memory M
- Where we need this:
  - ORDER BY in SQL queries
  - Several physical operators
  - Bulk loading of B+-tree indexes.
- Will discuss only 2-pass sorting, for when $B < M^2$
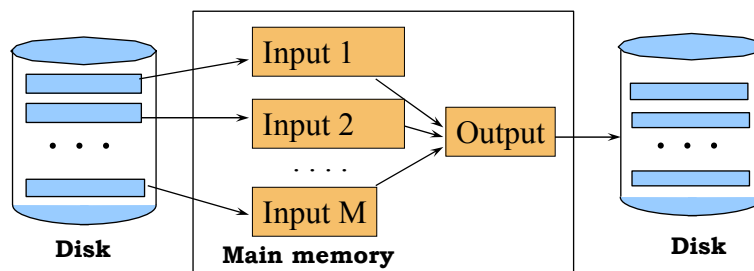
27

# External Merge-Sort: Step 1

- Phase one: load M bytes in memory, sort



Runs of length M bytes

28

# External Merge-Sort: Step 2

- Merge M – 1 runs into a new run
- Result: runs of length M (M – 1)$\approx$ M$^2$



If B <= M$^2$  then we are done

# Cost of External Merge Sort

- Read+write+read = 3B(R)

- Assumption: B(R) <= M$^2$

# Duplicate Elimination

Duplicate elimination $\delta(R)$
- Idea: do a two step merge sort, but change one of the steps

- Question in class: which step needs to be changed and how ?

- Cost = 3B(R)
- Assumption: $B(\delta(R)) <= M^2$

31

# Grouping

Grouping: $\gamma_{a, sum(b)}(R)$
- Same as before: sort, then compute the sum(b) for each group of a's
- Total cost: 3B(R)
- Assumption: $B(R) <= M^2$

32

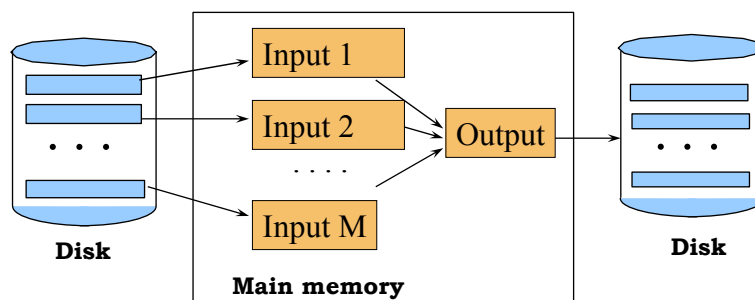# Merge-Join

Join R |x| S
- Step 1a: initial runs for R
- Step 1b: initial runs for S
- Step 2: merge and join

33

# Merge-Join



$M_1 = B(R)/M$ runs for R
$M_2 = B(S)/M$ runs for S
If $B <= M^2$ then we are done

34

# Two-Pass Algorithms Based on Sorting

Join R |x| S

- If the number of tuples in R matching those in S is small (or vice versa) we can compute the join during the merge phase
- Total cost: 3B(R)+3B(S)
- Assumption: $B(R) + B(S) <= M^2$

35

# Index Based Selection

- Selection on equality: $\sigma_{a=v}(R)$
- Clustered index on a:  cost B(R)/V(R,a)
- Unclustered index on a: cost T(R)/V(R,a)

36

# Index Based Selection

- Example:

  $B(R) = 2000$
  $T(R) = 100,000$
  $V(R, a) = 20$

  cost of $\sigma_{a=v}(R) = ?$

- Table scan (assuming R is clustered):
  - $B(R) = 2,000$ I/Os
- Index based selection:
  - If index is clustered: $B(R)/V(R,a) = 100$ I/Os
  - If index is unclustered: $T(R)/V(R,a) = 5,000$ I/Os

- Lesson: don't build unclustered indexes when $V(R,a)$ is small !

37

# Index Based Join

- $R \bowtie S$
- Assume S has an index on the join attribute
- Iterate over R, for each tuple fetch corresponding tuple(s) from S
- Assume R is clustered. Cost:
  - If index is clustered:  $B(R) + T(R)B(S)/V(S,a)$
  - If index is unclustered: $B(R) + T(R)T(S)/V(S,a)$

38

# Index Based Join

- Assume both R and S have a sorted index (B+ tree) on the join attribute
- Then perform a merge join
  - called zig-zag join
- Cost: $B(R) + B(S)$

# Summary of External Join Algorithms

- Block Nested Loop: $B(S) + B(R)*B(S)/M$

- Partitioned Hash: $3B(R)+3B(S)$;
  - $\min(B(R),B(S)) <= M^2$

- Merge Join: $3B(R)+3B(S$
  - $B(R)+B(S) <= M^2$

- Index Join: $B(R) + T(R)B(S)/V(S,a)$