# Introduction to Data Management
## CSE 344

Lecture 8-9: Relational Algebra
and Query Evaluation

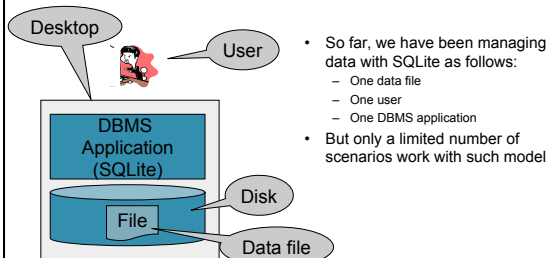Magda Balazinska - CSE 344, Fall 2012          1

---

# Where We Are

- Motivation for using a DBMS for managing data
- SQL, SQL, SQL
  - Declaring the schema for our data (CREATE TABLE)
  - Inserting data one row at a time or in bulk (INSERT/.import)
  - Modifying the schema and updating the data (ALTER/UPDATE)
  - Querying the data (SELECT)
  - Tuning queries (CREATE INDEX)

- Next step: More knowledge of how DBMSs work
  - Client-server architecture
  - Relational algebra and query execution

Magda Balazinska - CSE 344, Fall 2012          2

---

# Data Management with SQLite

Desktop

User

DBMS Application (SQLite)
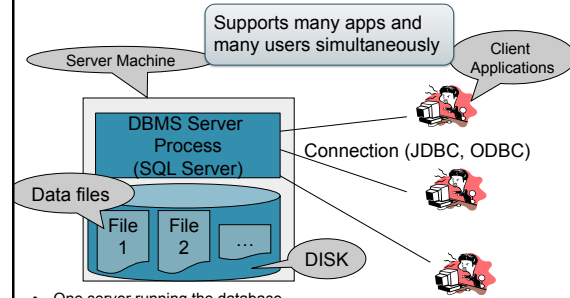
Disk

File

Data file

- So far, we have been managing data with SQLite as follows:
  - One data file
  - One user
  - One DBMS application
- But only a limited number of scenarios work with such model

Magda Balazinska - CSE 344, Fall 2012          3

---

# Client-Server Architecture

Supports many apps and many users simultaneously

Server Machine

Client Applications

DBMS Server Process (SQL Server)

Connection (JDBC, ODBC)

Data files

File 1     File 2     …

DISK

- One server running the database
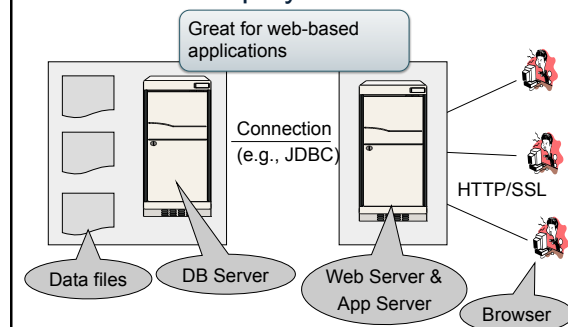- Many clients, connecting via the ODBC or JDBC (Java Database Connectivity) protocol
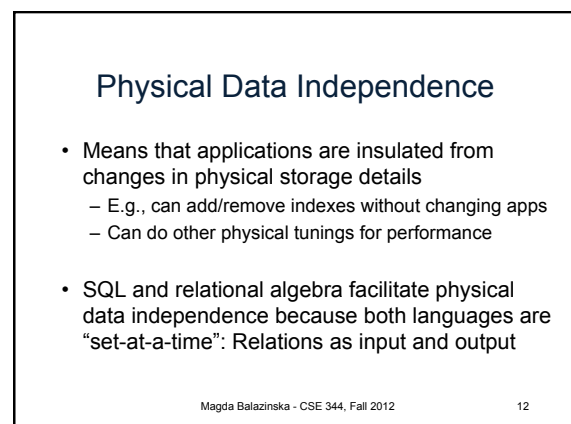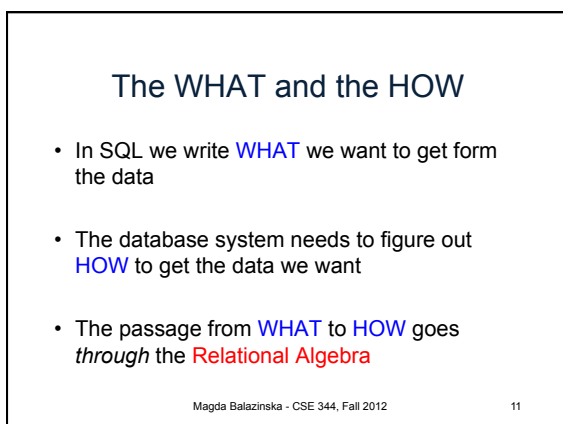
4

---

# Client-Server Architecture

- One *server* that runs the DBMS (or RDBMS):
  - Your own desktop, or
  - Some beefy system, or
  - A cloud service (SQL Azure)
- Many *clients* run apps and connect to DBMS
  - Microsoft's Management Studio (for SQL Server), or
  - psql (for postgres)
  - Some Java program (HW5) or some C++ program
- Clients "talk" to server using JDBC/ODBC protocol

Magda Balazinska - CSE 344, Fall 2012          5

---

# DBMS Deployment: 3 Tiers

Great for web-based applications

Connection (e.g., JDBC)

HTTP/SSL

Data files     DB Server     Web Server & App Server     Browser

---

## DBMS Deployment: Cloud

Great for web-based applications too

Data Files

HTTP/SSL

DB Server

Web & App Server

Developers

Users

Magda Balazinska - CSE 344, Fall 2012

## Using a DBMS Server

1. Client application establishes connection to server
2. Client must authenticate self
3. Client submits SQL commands to server
4. Server executes commands and returns results

DBMS

Magda Balazinska - CSE 344, Fall 2012    8

## Query Evaluation Steps Review

SQL query

Parse & Check Query

Translate query string into internal representation

Check syntax, access control, table names, etc.

Decide how best to answer query: query optimization

Query Execution

Query Evaluation

Return Results

Magda Balazinska - CSE 344, Fall 2012    9

## Question: How does Query Evaluation Work?

Magda Balazinska - CSE 344, Fall 2012    10

## The WHAT and the HOW

- In SQL we write WHAT we want to get form the data

- The database system needs to figure out HOW to get the data we want

- The passage from WHAT to HOW goes *through* the Relational Algebra

Magda Balazinska - CSE 344, Fall 2012    11

## Physical Data Independence

- Means that applications are insulated from changes in physical storage details
  – E.g., can add/remove indexes without changing apps
  – Can do other physical tunings for performance

- SQL and relational algebra facilitate physical data independence because both languages are "set-at-a-time": Relations as input and output

Magda Balazinska - CSE 344, Fall 2012    12

## Overview: SQL = WHAT

Product(pid, name, price)
Purchase(pid, cid, store)
Customer(cid, name, city)

SELECT DISTINCT x.name, z.name
FROM Product x, Purchase y, Customer z
WHERE x.pid = y.pid and y.cid = y.cid and
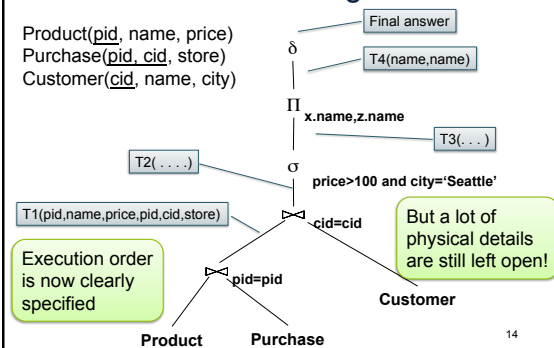        x.price > 100 and z.city = 'Seattle'

It's clear WHAT we want, unclear HOW to get it

Magda Balazinska - CSE 344, Fall 2012    13

---

## Overview: Relational Algebra = HOW

Product(pid, name, price)
Purchase(pid, cid, store)
Customer(cid, name, city)

Final answer

$\delta$    T4(name,name)

$\Pi$ x.name,z.name

T2( . . . . )    T3(. . . )

$\sigma$ price>100 and city='Seattle'

T1(pid,name,price,pid,cid,store)

⋈ cid=cid

But a lot of physical details are still left open!

Execution order is now clearly specified

⋈ pid=pid

Customer

Product    Purchase

Magda Balazinska - CSE 344, Fall 2012    14

---

## Relational Algebra

Magda Balazinska - CSE 344, Fall 2012    15

---

## Sets v.s. Bags

- Sets: {a,b,c}, {a,d,e,f}, { }, . . .
- Bags: {a, a, b, c}, {b, b, b, b, b}, . . .
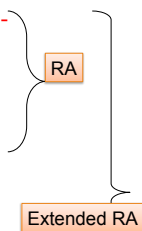
Relational Algebra has two semantics:
- Set semantics = standard Relational Algebra
- Bag semantics = extended Relational Algebra

Magda Balazinska - CSE 344, Fall 2012    16

---

## Relational Algebra Operators

- Union ∪, intersection ∩, difference -
- Selection $\sigma$
- Projection $\Pi$
- Cartesian product ×, join ⋈
- Rename $\rho$

RA

- Duplicate elimination $\delta$
- Grouping and aggregation $\gamma$
- Sorting $\tau$

Extended RA

Magda Balazinska - CSE 344, Fall 2012    17

---

## Union and Difference

R1 ∪ R2
R1 – R2

What do they mean over bags ?

Magda Balazinska - CSE 344, Fall 2012    18

3

## What about Intersection ?

- Derived operator using minus

$$R1 \cap R2 = R1 - (R1 - R2)$$

- Derived using join (will explain later)

$$R1 \cap R2 = R1 \bowtie R2$$

## Selection

- Returns all tuples which satisfy a condition

$$\sigma_c(R)$$

- Examples
  - $\sigma_{Salary > 40000}$ (Employee)
  - $\sigma_{name = \text{"Smith"}}$ (Employee)
- The condition c can be =, <, ≤, >, ≥, <>

Employee

| SSN | Name | Salary |
|---|---|---|
| 1234545 | John | 200000 |
| 5423341 | Smith | 600000 |
| 4352342 | Fred | 500000 |

$\sigma_{Salary > 40000}$ (Employee)

| SSN | Name | Salary |
|---|---|---|
| 5423341 | Smith | 600000 |
| 4352342 | Fred | 500000 |

## Projection

- Eliminates columns

$$\Pi_{A1,\ldots,An}(R)$$

- Example: project social-security number and names:
  - $\Pi_{SSN, Name}$ (Employee)
  - Answer(SSN, Name)

Different semantics over sets or bags!  Why?

Employee

| SSN | Name | Salary |
|---|---|---|
| 1234545 | John | 20000 |
| 5423341 | John | 60000 |
| 4352342 | John | 20000 |

$\Pi_{Name,Salary}$ (Employee)

| Name | Salary |
|---|---|
| John | 20000 |
| John | 60000 |
| John | 20000 |

Bag semantics

| Name | Salary |
|---|---|
| John | 20000 |
| John | 60000 |

Set semantics

Which is more efficient?

23

## Composing RA Operators

Patient

| no | name | zip | disease |
|---|---|---|---|
| 1 | p1 | 98125 | flu |
| 2 | p2 | 98125 | heart |
| 3 | p3 | 98120 | lung |
| 4 | p4 | 98120 | heart |

$\pi_{zip,disease}$(Patient)

| zip | disease |
|---|---|
| 98125 | flu |
| 98125 | heart |
| 98120 | lung |
| 98120 | heart |

$\sigma_{disease='heart'}$(Patient)

| no | name | zip | disease |
|---|---|---|---|
| 2 | p2 | 98125 | heart |
| 4 | p4 | 98120 | heart |

$\pi_{zip}(\sigma_{disease='heart'}$(Patient))

| zip |
|---|
| 98120 |
| 98125 |

## Cartesian Product

- Each tuple in R1 with each tuple in R2

$$R1 \times R2$$

- Rare in practice; mainly used to express joins

## Cross-Product Example

**Employee**

| Name | SSN |
|------|-----|
| John | 999999999 |
| Tony | 777777777 |

**Dependent**

| EmpSSN | DepName |
|--------|---------|
| 999999999 | Emily |
| 777777777 | Joe |

**Employee ✕ Dependent**

| Name | SSN | EmpSSN | DepName |
|------|-----|--------|---------|
| John | 999999999 | 999999999 | Emily |
| John | 999999999 | 777777777 | Joe |
| Tony | 777777777 | 999999999 | Emily |
| Tony | 777777777 | 777777777 | Joe |

## Renaming

- Changes the schema, not the instance

$$\rho_{B1,\ldots,Bn}(R)$$

- Example:
  - $\rho_{N, S}(\text{Employee}) \rightarrow \text{Answer}(N, S)$

Not really used by systems, but needed on paper

## Natural Join

$$R1 \bowtie R2$$

- Meaning: $R1 \bowtie R2 = \Pi_A(\sigma(R1 \times R2))$

- Where:
  - Selection $\sigma$ checks equality of all common attributes
  - Projection eliminates duplicate common attributes

## Natural Join Example

**R**

| A | B |
|---|---|
| X | Y |
| X | Z |
| Y | Z |
| Z | V |

**S**

| B | C |
|---|---|
| Z | U |
| V | W |
| Z | V |

**R ⋈ S =**
$\Pi_{ABC}(\sigma_{R.B=S.B}(R \times S))$

| A | B | C |
|---|---|---|
| X | Z | U |
| X | Z | V |
| Y | Z | U |
| Y | Z | V |
| Z | V | W |

## Natural Join Example 2

AnonPatient P

| age | zip | disease |
|-----|-----|---------|
| 54 | 98125 | heart |
| 20 | 98120 | flu |

Voters V

| name | age | zip |
|------|-----|-----|
| p1 | 54 | 98125 |
| p2 | 20 | 98120 |

P ⋈ V

| age | zip | disease | name |
|-----|-----|---------|------|
| 54 | 98125 | heart | p1 |
| 20 | 98120 | flu | p2 |

## Natural Join

- Given schemas R(A, B, C, D), S(A, C, E), what is the schema of R ⋈ S ?

- Given R(A, B, C),  S(D, E), what is R ⋈ S ?

- Given R(A, B),  S(A, B),  what is  R ⋈ S ?

---

## Theta Join

- A join that involves a predicate

$$R1 \bowtie_\theta R2 \;=\; \sigma_\theta (R1 \times R2)$$

- Here θ can be any condition
- For our voters/disease example:

$P \bowtie_{\text{P.zip = V.zip and P.age < V.age + 5 and P.age > V.age - 5}} V$

---

## Equijoin

- A theta join where θ is an equality

$$R1 \bowtie_{A=B} R2 \;=\; \sigma_{A=B} (R1 \times R2)$$

- This is by far the most used variant of join in practice

---

## Equijoin Example

AnonPatient P

| age | zip | disease |
|-----|-------|---------|
| 54 | 98125 | heart |
| 20 | 98120 | flu |

Voters V

| name | age | zip |
|------|-----|-------|
| p1 | 54 | 98125 |
| p2 | 20 | 98120 |

$P \bowtie_{\text{P.age=V.age}} V$

| age | P.zip | disease | name | V.zip |
|-----|-------|---------|------|-------|
| 54 | 98125 | heart | p1 | 98125 |
| 20 | 98120 | flu | p2 | 98120 |

---

## Join Summary

- **Theta-join**: $R \bowtie_\theta S = \sigma_\theta(R \times S)$
  - Join of R and S with a join condition θ
  - Cross-product followed by selection θ
- **Equijoin**: $R \bowtie_\theta S = \pi_A (\sigma_\theta(R \times S))$
  - Join condition θ consists only of equalities
  - Projection $\pi_A$ drops all redundant attributes
- **Natural join**: $R \bowtie S = \pi_A (\sigma_\theta(R \times S))$
  - Equijoin
  - Equality on **all** fields with same name in R and in S

---

## So Which Join Is It ?

- When we write R ⋈ S we usually mean an equijoin, but we often omit the equality predicate when it is clear from the context

## More Joins

- **Outer join**
  - Include tuples with no matches in the output
  - Use NULL values for missing attributes

- Variants
  - Left outer join
  - Right outer join
  - Full outer join

## Outer Join Example

AnonPatient P

| age | zip | disease |
|-----|-----|---------|
| 54 | 98125 | heart |
| 20 | 98120 | flu |
| 33 | 98120 | lung |

AnnonJob J

| job | age | zip |
|-----|-----|-----|
| lawyer | 54 | 98125 |
| cashier | 20 | 98120 |

$P \bowtie V$

| age | zip | disease | job |
|-----|-----|---------|-----|
| 54 | 98125 | heart | lawyer |
| 20 | 98120 | flu | cashier |
| 33 | 98120 | lung | null |

## Some Examples

```
Supplier(sno,sname,scity,sstate)
Part(pno,pname,psize,pcolor)
Supply(sno,pno,qty,price)
```

Q2: Name of supplier of parts with size greater than 10
$\pi_{sname}(Supplier \bowtie Supply \bowtie (\sigma_{psize>10} (Part))$

Q3: Name of supplier of red parts or parts with size greater than 10
$\pi_{sname}(Supplier \bowtie Supply \bowtie (\sigma_{psize>10} (Part) \cup \sigma_{pcolor='red'} (Part) ) )$

## From SQL to RA

## From SQL to RA

Product(pid, name, price)
Purchase(pid, cid, store)
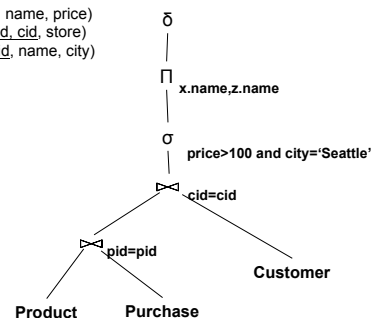Customer(cid, name, city)

```
SELECT DISTINCT x.name, z.name
FROM Product x, Purchase y, Customer z
WHERE x.pid = y.pid and y.cid = y.cid and
      x.price > 100 and z.city = 'Seattle'
```

## From SQL to RA

Product(pid, name, price)
Purchase(pid, cid, store)
Customer(cid, name, city)

δ
|
Π **x.name,z.name**
|
σ **price>100 and city='Seattle'**
|
⋈ **cid=cid**
|
⋈ **pid=pid** — **Customer**
|
**Product**   **Purchase**

42

## An Equivalent Expression

Query optimization = finding cheaper, equivalent expressions

$\delta$

$\Pi$ x.name,z.name

$\bowtie$ cid=cid

$\bowtie$ pid=pid

$\sigma$ city='Seattle'

$\sigma$ price>100

Customer

Product    Purchase

43

---

## Extended RA: Operators on Bags

- Duplicate elimination $\delta$
- Grouping $\gamma$
- Sorting $\tau$

Magda Balazinska - CSE 344, Fall 2012      44

---

## Logical Query Plan

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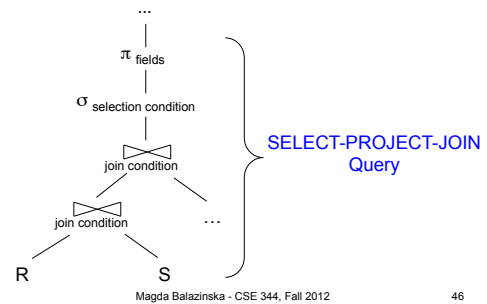FROM sales
GROUP BY city
HAVING sum(price) > 100

T3(city, c)

$\Pi$ city, c

T2(city,p,c)

$\sigma$ p > 100

T1(city,p,c)

$\gamma$ city, sum(price)→p, count(*) → c

T1, T2, T3 = temporary tables      sales(product, city, price)

Magda Balazinska - CSE 344, Fall 2012      45

---

## Typical Plan for Block (1/2)

...

$\pi$ fields

$\sigma$ selection condition

$\bowtie$ join condition

SELECT-PROJECT-JOIN Query

$\bowtie$ join condition      ...

R          S

Magda Balazinska - CSE 344, Fall 2012      46

---

## Typical Plan For Block (2/2)

having condition

$\gamma$ fields, sum/count/min/max(fields)

$\pi$ fields

$\sigma$ selection condition

$\bowtie$ join condition

…          …

Magda Balazinska - CSE 344, Fall 2012      47

---

## How about Subqueries?

Supplier(sno,sname,scity,sstate)
Part(pno,pname,psize,pcolor)
Supply(sno,pno,price)

SELECT  Q.sno
FROM Supplier Q
WHERE  Q.sstate = 'WA'
  and not exists
   (SELECT *
    FROM Supply P
    WHERE P.sno = Q.sno
      and P.price > 100)

Magda Balazinska - CSE 344, Fall 2012      48

**How about Subqueries?**

Supplier(sno,sname,scity,sstate)
Part(pno,pname,psize,pcolor)
Supply(sno,pno,price)

```
SELECT  Q.sno
FROM Supplier Q
WHERE  Q.sstate = 'WA'
  and not exists
     (SELECT *
      FROM Supply P
      WHERE P.sno = Q.sno
          and P.price > 100)
```

Correlation !

Magda Balazinska - CSE 344, Fall 2012    49

---

**How about Subqueries?**

Supplier(sno,sname,scity,sstate)
Part(pno,pname,psize,pcolor)
Supply(sno,pno,price)

```
SELECT  Q.sno
FROM Supplier Q
WHERE  Q.sstate = 'WA'
  and not exists
     (SELECT *
      FROM Supply P
      WHERE P.sno = Q.sno
          and P.price > 100)
```

De-Correlation

```
SELECT  Q.sno
FROM Supplier Q
WHERE  Q.sstate = 'WA'
  and Q.sno not in
     (SELECT P.sno
      FROM Supply P
      WHERE P.price > 100)
```

Magda Balazinska - CSE 344, Fall 2012    50

---

**How about Subqueries?**

Supplier(sno,sname,scity,sstate)
Part(pno,pname,psize,pcolor)
Supply(sno,pno,price)

Un-nesting

```
(SELECT  Q.sno
 FROM Supplier Q
 WHERE  Q.sstate = 'WA')
  EXCEPT
(SELECT P.sno
 FROM Supply P
 WHERE P.price > 100)
```

EXCEPT = set difference

```
SELECT  Q.sno
FROM Supplier Q
WHERE  Q.sstate = 'WA'
  and Q.sno not in
     (SELECT P.sno
      FROM Supply P
      WHERE P.price > 100)
```

Magda Balazinska - CSE 344, Fall 2012    51

---

**How about Subqueries?**

Supplier(sno,sname,scity,sstate)
Part(pno,pname,psize,pcolor)
Supply(sno,pno,price)

Finally…

```
(SELECT  Q.sno
 FROM Supplier Q
 WHERE  Q.sstate = 'WA')
  EXCEPT
(SELECT P.sno
 FROM Supply P
 WHERE P.price > 100)
```

$-$
$\sigma_{sstate='WA'}$  $\sigma_{Price > 100}$
Supplier    Supply

Magda Balazinska - CSE 344, Fall 2012    52

---

**From Logical Plans
to Physical Plans**

Magda Balazinska - CSE 344, Fall 2012    53

---

**Example**

Supplier(sid, sname, scity, sstate)
Supply(sid, pno, quantity)

```
SELECT sname
FROM Supplier x, Supply y
WHERE x.sid = y.sid
   and  y.pno = 2
   and x.scity = 'Seattle'
   and x.sstate = 'WA'
```

Give a relational algebra expression for this query

Magda Balazinska - CSE 344, Fall 2012    54

## Relational Algebra

Supplier(sid, sname, scity, sstate)
Supply(sid, pno, quantity)

$\pi_{sname}(\sigma_{scity=\text{'Seattle'} \wedge sstate=\text{'WA'} \wedge pno=2}(Supplier \bowtie_{sid=sid} Supply))$
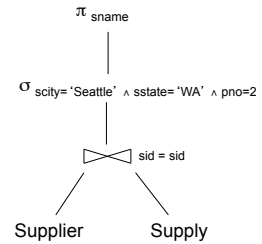
Magda Balazinska - CSE 344, Fall 2012        55

---

## Relational Algebra

Supplier(sid, sname, scity, sstate)
Supply(sid, pno, quantity)

Relational algebra expression is also called the "logical query plan"

$\pi_{sname}$
|
$\sigma_{scity=\text{'Seattle'} \wedge sstate=\text{'WA'} \wedge pno=2}$
|
$\bowtie_{sid=sid}$
/          \
Supplier       Supply

Magda Balazinska - CSE 344, Fall 2012        56

---

## Physical Query Plan 1

(On the fly)        $\pi_{sname}$

A physical query plan is a logical query plan annotated with physical implementation details

(On the fly)
$\sigma_{scity=\text{'Seattle'} \wedge sstate=\text{'WA'} \wedge pno=2}$

(Block-nested loop)
$\bowtie_{sid=sid}$
/              \
Supplier              Supply
(File scan)            (File scan)

Magda Balazinska - CSE 344, Fall 2012        57

---

## Physical Query Plan 2

(On the fly)        $\pi_{sname}$   (d)

(Sort-merge join)        $\bowtie_{sid=sid}$   (c)

(Scan
write to T1)                              (Scan
                                          write to T2)

(a) $\sigma_{scity=\text{'Seattle'} \wedge sstate=\text{'WA'}}$       (b) $\sigma_{pno=2}$
|                                    |
Supplier                             Supply
(File scan)                          (File scan)

Magda Balazinska - CSE 344, Fall 2012        58

---

## Physical Query Plan 3

(On the fly)  (d)  $\pi_{sname}$

(On the fly)
(c)  $\sigma_{scity=\text{'Seattle'} \wedge sstate=\text{'WA'}}$
|
(b)  $\bowtie_{sid=sid}$   (Index nested loop)
/          \
(Use index)
(a) $\sigma_{pno=2}$                     Supplier
|
Supply
(Index lookup on pno )  (Index lookup on sid)
Assume: clustered       Doesn't matter if clustered or not  59