CSE544 Data Management

Lecture 6: Data Models, Relational Algebra

Announcements

Project teams due Friday

HW2 is posted

Review of "What goes around..." today

Where We Are

 We are done with SQL; Please continue to read and learn on your own

 Today: data models, and why the relational model wins; then RA

Next lectures: DBMS internals

References

 M. Stonebraker and J. Hellerstein. What Goes Around Comes Around. In "Readings in Database Systems" (aka the Red Book). 4th ed.

Outline

- Early data models
 - IMS
 - CODASYL

Relational Model in some detail

Data models that followed the relational model

Early Proposal 1: IMS*

What is it?

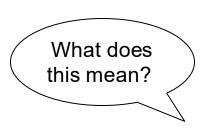
Early Proposal 1: IMS*

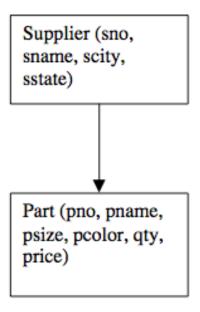
Hierarchical data model

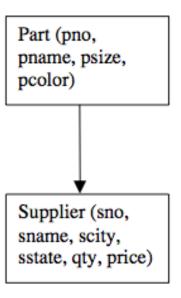
- Record
 - Type: collection of named fields with data types
 - Instance: must match type definition
 - Each instance has a key
 - Record types arranged in a tree
- IMS database is collection of instances of record types organized in a tree

IMS Example

Figure 2 from "What goes around comes around"

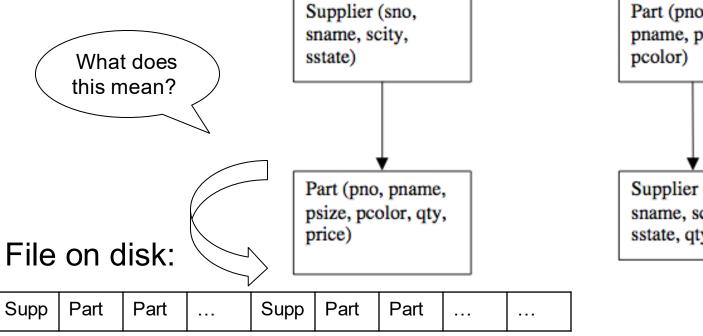


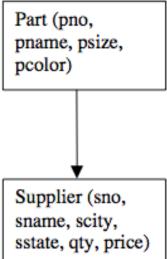




IMS Example

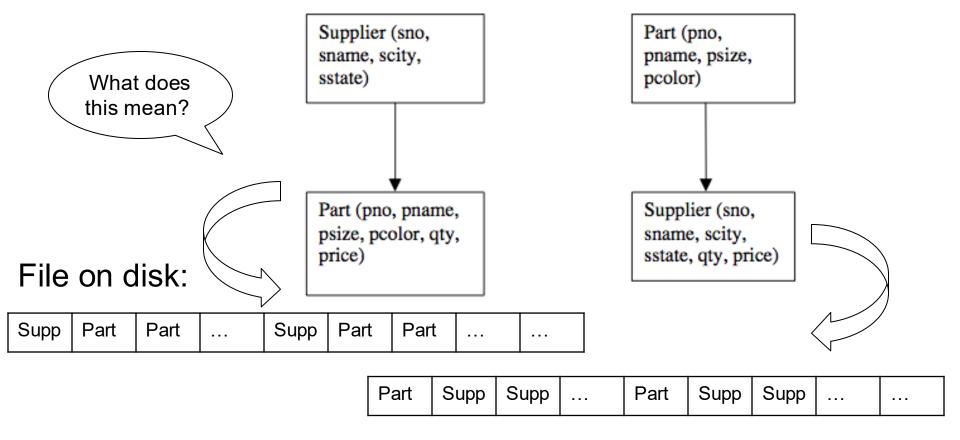
Figure 2 from "What goes around comes around"





IMS Example

Figure 2 from "What goes around comes around"



- Tree-structured data model
 - Redundant data; existence depends on parent

- Tree-structured data model
 - Redundant data; existence depends on parent
- Record-at-a-time user interface
 - User must specify algorithm to access data

- Tree-structured data model
 - Redundant data; existence depends on parent
- Record-at-a-time user interface
 - User must specify algorithm to access data
- Very limited physical independence
 - Phys. organization limits possible operations
 - Application programs break if organization changes
- Some logical independence but limited

Data Manipulation Language: DL/1

How does a programmer retrieve data in IMS?

Data Manipulation Language: DL/1

How does a programmer retrieve data in IMS?

- Each record has a hierarchical sequence key (HSK)
- HSK defines semantics of commands:
 - get_next; get_next_within_parent
- DL/1 is a record-at-a-time language
 - Programmers construct algorithm, worry about optimization

Data storage

How is data physically stored in IMS?

Data storage

How is data physically stored in IMS?

- Root records
 - Stored sequentially (sorted on key)
 - Indexed in a B-tree using the key of the record
 - Hashed using the key of the record
- Dependent records
 - Physically sequential
 - Various forms of pointers
- Selected organizations restrict DL/1 commands
 - No updates allowed due to sequential organization
 - No "get-next" for hashed organization

Data Independence

What is it?

Data Independence

What is it?

 Physical data independence: Applications are insulated from changes in physical storage details

 Logical data independence: Applications are insulated from changes to logical structure of the data

Lessons from IMS

- Physical/logical data independence needed
- Tree structure model is restrictive

 Record-at-a-time programming forces user to do optimization

Early Proposal 2: CODASYL

What is it?

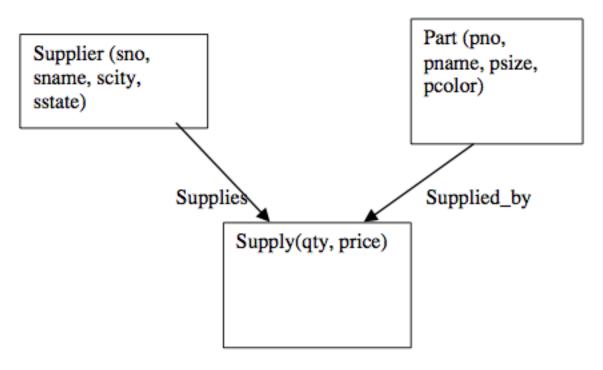
Early Proposal 2: CODASYL

What is it?

- Networked data model
- Primitives are also record types with keys
- Record types are organized into network
- Multiple parents; arcs = "sets"
- More flexible than hierarchy
- Record-at-a-time data manipulation language

CODASYL Example

Figure 5 from "What goes around comes around"



CODASYL Limitations

 No data independence: application programs break if organization changes

Record-at-a-time: "navigate the hyperspace"

The Programmer as Navigator





Outline

Early data models

Relational Model in some detail

Data models that followed the relational model

Relational Model Overview

Ted Codd 1970

What was the motivation? What is the model?





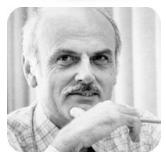
Relational Model Overview

Ted Codd 1970

- Motivation: logical and physical data independence
- Store data in a simple data structure (table)
- Access data through set-at-a-time language
- No need for physical storage proposal

Relational Database: A Practical Foundation for Productivity





Great Debate

- Pro relational
 - What were the arguments?

- Against relational
 - What were the arguments?

How was it settled?

Great Debate

- Pro relational
 - CODASYL is too complex
 - No data independence
 - Record-at-a-time hard to optimize
 - Trees/networks not flexible enough
- Against relational
 - COBOL programmers cannot understand relational languages
 - Impossible to implement efficiently
- Ultimately settled by the market place

Outline

Early data models

Relational Model in some detail

Data models that followed the relational model

Entity-relationship

Object-relational

Semistructured

Entity-relationship

Abandoned as data model.

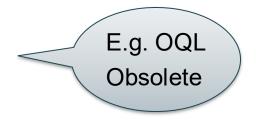
Adopted for schema design

Object-relational

Semistructured

Entity-relationship

Object-relational



Semistructured

Entity-relationship

Object-relational

Semistructured

XML, Json, Protobuf Nested data; tree-like

Entity-relationship

Object-relational

Semistructured

Key-value pairs

NoSQL:

- GET(K), PUT(K,V)
- · Great scalability
- Poor functionality

Data Independence in the Relational Model

Data Independence

- Logical data independence:
 - SQL views

- Physical data independence:
 - 1. Declarative query: FO or SQL
 - 2. Query plan: Relational Algebra
 - 3. Query optimization / execution

Data Independence

- Logical data independence:
 - SQL views

- Physical data independence:
 - 1. Declarative query: FO or SQL
 - 2. Query plan: Relational Algebra
 - 3. Query optimization / execution

SQL Views

CREATE VIEW RedSuppliers AS ...

Creates a new relation name

 The content of that relation is defined by a SELECT-FROM-WHERE query

SQL Views

CREATE VIEW RedSuppliers AS
SELECT DISTINCT x.*
FROM Supplier x, Supply y, Part z
WHERE x.sno=y.sno and y.pno=z.pno and z.pcolor='Red';

SQL Views

```
CREATE VIEW RedSuppliers AS

SELECT DISTINCT x.*

FROM Supplier x, Supply y, Part z

WHERE x.sno=y.sno and y.pno=z.pno and z.pcolor='Red';
```

RedSuppliers is added to the database schema

SQL Views

CREATE VIEW RedSuppliers AS

SELECT DISTINCT x.*

FROM Supplier x, Supply y, Part z

WHERE x.sno=y.sno and y.pno=z.pno and z.pcolor='Red';

RedSuppliers is added to the database schema Later we can use it, like any table:

SELECT DISTINCT x.*
FROM RedSupplier x, Supply y, Part z
WHERE x.sno=y.sno and y.pno=z.pno and z.pcolor='Blue';

SQL Views

CREATE VIEW RedSuppliers AS

SELECT DISTINCT x.*

FROM Supplier x, Supply y, Part z

WHERE x.sno=y.sno and y.pno=z.pno and z.pcolor='Red';

RedSuppliers is added to the database schema Later we can use it, like any table:

What does this query compute?

SELECT DISTINCT x *

FROM RedSupplier x, Supply y, Part z
WHERE x.sno=y.sno and y.pno=z.pno and z.pcolor='Blue';

Discussion

- CREATE VIEW...
 - Add view name permanently to the schema
 - To delete, type DROP VIEW ...
- CREATE TEMPORAY VIEW...
 - Supported in postgres
 - Add view name to the schema...
 - ...remove at the end of the session
- WITH Only local to one query

Virtual v.s. Materialized Views

- Virtual view the default
 - The query defining the view is evaluated every time when it is used
 - The data is always up to date!
 - But we re-compute it every time: inefficient
- Materialized view
 - CREATE MATERALIZED VIEW ...
 - Supported by some systems
 - Computed only once: efficient
 - But may not be up to date

Data Independence

- Logical data independence:
 - SQL views

- Physical data independence:
 - 1. Declarative query: FO or SQL
 - 2. Query plan: Relational Algebra
 - 3. Query optimization / execution

Data Independence

- Logical data independence:
 - SQL views

- Physical data independence:
 - 1. Declarative query: FO or SQL
 - 2. Query plan: Relational Algebra
 - 3. Query optimization / execution

We saw this

Next

Next lectures

Relational Algebra

Executing SQL Queries

User writes in SQL

System: SQL → Relational Algebra

Query Plan: is optimized, the executed

Relational Algebra

Five operators:

- Selection σ
- Projection Π
- Join or cartesian product ⋈, ×
- Union U
- Difference –

 $\sigma_{condition}(T)$

Returns tuples in T that satisfy the condition

SELECT *
FROM T
WHERE condition;

$$\sigma_{condition}(T)$$

Returns tuples in T that satisfy the condition

$$\sigma_{C=20 \wedge D \geq 66}(S) =$$

S=

SELECT *
FROM T
WHERE condition;

R= A B
1 10
1 20
2 20

$$\sigma_{condition}(T)$$

Returns tuples in T that satisfy the condition

$$\sigma_{C=20 \land D \ge 66}(S) = \begin{array}{|c|c|} \hline C \\ \hline 20 \\ \hline 20 \\ \hline \end{array}$$

S=

SELECT *
FROM T
WHERE condition;

С	D
10	33
10	44
20	55
20	66
20	77
30	66

D

66

77

$$\sigma_{condition}(T)$$

$$\sigma_{D \ge 66}(S) =$$

С	D
20	66
20	77
30	66

Returns tuples in T that satisfy the condition

$$\sigma_{C=20 \wedge D \geq 66}(S) =$$

С	D
20	66
20	77

SELECT *
FROM T
WHERE condition;

А	В
1	10
1	20
2	20

С	D
10	33
10	44
20	55
20	66
20	77
30	66

Projection

 $\Pi_{col1,col2,...}(T)$

Return columns 1, 2, ... from T; all rows

SELECT T.col1, T.col2,...
FROM T
WHERE condition;

R=

Α	В
1	10
1	20
2	20

С	D
10	33
10	44
20	55
20	66
20	77
30	66

Projection

$$\Pi_{col1,col2,...}(T)$$

Return columns 1, 2, ... from T; all rows

	ن
$\Pi_C(S) = \frac{1}{2}$	10
	10
	20
	20
	20
	30

SELECT T.col1, T.col2,...
FROM T
WHERE condition;

Α	В
1	10
1	20
2	20

C	D
10	33
10	44
20	55
20	66
20	77
30	66

Bag semantics...

Projection

 $\Pi_{col1,col2,...}(T)$

Return columns 1, 2, ... from T; all rows

	С
$\Pi_C(S) =$	10
	10
	20
	20
	20
	30

SELECT T.col1, T.col2,...
FROM T
WHERE condition;

R=

Α	В
1	10
1	20
2	20

С	D
10	33
10	44
20	55
20	66
20	77
30	66

Bag semantics...

Projection

 $\Pi_{col1,col2,...}(T)$

Return columns 1, 2, ... from T; all rows

SELECT T.col1, T.col2,...
FROM T
WHERE condition;

R=

Α	В
1	10
1	20
2	20

C	D
10	33
10	44
20	55
20	66
20	77
30	66

Join

 $T_1 \bowtie_{cond} T_2$

Joins the two tables

SELECT *
FROM T1, T2
WHERE cond;

Join

 $T_1 \bowtie_{cond} T_2$

$$R \bowtie_{B=C} S =$$

Joins the two tables

Α	В	С	D
1	10	10	33
1	10	10	44
1	20	20	55
1	20	20	66
1	20	20	77
2	20	20	55
2	20	20	66
2	20	20	77

SELECT *
FROM T1, T2
WHERE cond;

Join

 $T_1 \bowtie_{cond} T_2$

 $R x \bowtie_{x.B=y.B} R y =$

y.A x.Ax.B y.B

Joins the two tables

Disambiguate attributes: various conventions...

SELECT *
FROM T1, T2
WHERE cond;

Many Variants of Joins

Eq-join:

$$R\bowtie_{B=C} S$$

Theta-join:

$$R\bowtie_{B\leq C\wedge A*D\leq 20} S$$

Cartesian product:

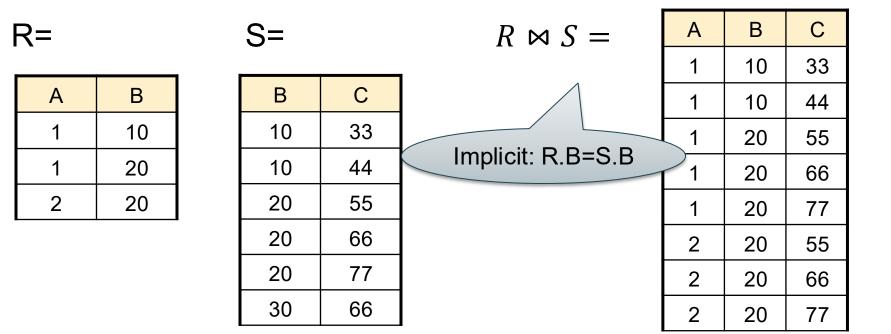
$$R \times S$$

Natural join:

next!

Natural Join

- Joins on common column names
- Retains only one copy of the column



What do these natural joins compute?

• $R(A,B) \bowtie S(B,C)$

• $R(A,B) \bowtie S(C,D)$

• $R(A,B)\bowtie S(A,B)$

What do these natural joins compute?

• $R(A,B) \bowtie S(B,C)$

Joins on B=B

• $R(A,B) \bowtie S(C,D)$

• $R(A,B)\bowtie S(A,B)$

What do these natural joins compute?

• $R(A,B) \bowtie S(B,C)$

• $R(A,B) \bowtie S(C,D)$

Joins on B=B

Cartesian product $R \times S$

• $R(A,B) \bowtie S(A,B)$

What do these natural joins compute?

•
$$R(A,B) \bowtie S(B,C)$$

• $R(A,B) \bowtie S(C,D)$

• $R(A,B) \bowtie S(A,B)$

Joins on B=B

Cartesian product $R \times S$

Intersection $R \cap S$

Even More Joins

Inner joins:
⋈ (all variations)

Left outer join:

Semi-join:

Anti-semi-join:

People use various symbols

X

Semi-Join

 $R \ltimes S$

Α	В
1	20
1	30

Tuples in R that join with S

SELECT DISTINCT R.*
FROM R, S
WHERE join-cond;

R=

Α	В
1	20
1	30
2	40

В	С
10	33
10	44
20	55
20	66
20	77
30	66

Anti Semi-Join

$$R \triangleright S = R - R \bowtie S$$

Α	В
2	40

S=

Tuples in R that do not join with S

Finally: Union and Difference

Set operations:

$$R \cup S$$
, $R - S$

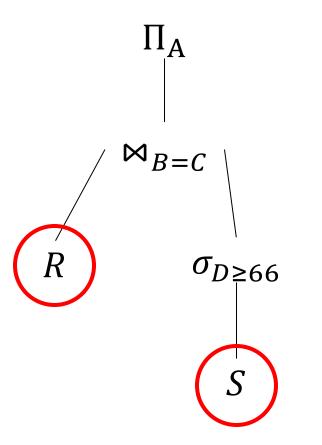
R, S must have the same schema!

$$\Pi_{A}(R \bowtie_{B=C} \sigma_{D \geq 66}(S))$$

$$\Pi_{A}(R \bowtie_{B=C} \sigma_{D \ge 66}(S))$$

$$\Pi_{A} \qquad \text{RA Plan,}$$
or Query Plan
$$R \qquad \sigma_{D \ge 66}$$

$$\Pi_{A}(R \bowtie_{B=C} \sigma_{D \geq 66}(S))$$

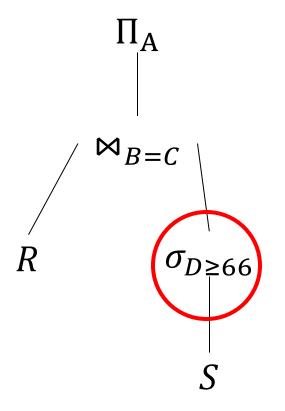


R=	Α	В
	1	10
	1	20
	2	20

С	D
10	33
10	44
20	55
20	66
20	77
30	66

S=

$$\Pi_{A}(R \bowtie_{B=C} \sigma_{D \geq 66}(S))$$



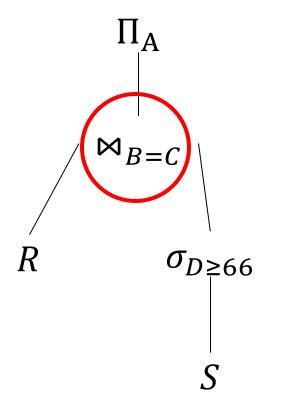
₹=	Α	В
	1	10
	1	20
	2	20

С	D
20	66
20	77
30	66

С	D
10	33
10	44
20	55
20	66
20	77
30	66

R=

$$\Pi_{A}(R\bowtie_{B=C}\sigma_{D\geq 66}(S))$$



Α	В	С	D
1	20	20	66
1	20	20	77
2	20	20	66
2	20	20	77

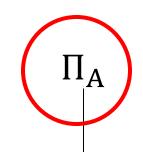
Α	В
1	10
1	20
2	20

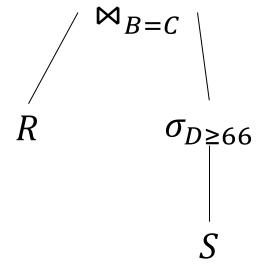
С	D
20	66
20	77
30	66

С	D
10	33
10	44
20	55
20	66
20	77
30	66

S=

$\Pi_{\mathbf{A}}(R)$	$\bowtie_{B=C}$	$\sigma_{D \ge 66}$	(S)
-----------------------	-----------------	---------------------	-----





Α	В	С	D
1	20	20	66
1	20	20	77
2	20	20	66
2	20	20	77

output

Final

С	D
20	66
20	77
30	66

Α

2

С	D
10	33
10	44
20	55
20	66
20	77
30	66

S=

R	_
1 /	

	/
Α	В
1	10
1	20
2	20

Relational Algebra

Five operators:

- Selection σ
- Projection Π
- Join or cartesian product ⋈, ×
- Union U
- Difference –

Relational Algebra

Five operators:

- Selection σ
- Projection Π
- Join or cartesian product ⋈, ×
- Union U
- Difference –

Which operations are monotone?

Relational Algebra

Five operators:

Which operations are monotone?

- Selection σ
- Projection Π
- Join or cartesian product ⋈, ×
- Union U
- Difference –

Monotone

Non-monotone

Extended Relational Algebra

Group-by and aggregate:

Y We only discuss this

• Duplicate elimination: δ

• Sorting: au

Group-by and Aggregates

 $\gamma_{col1,col2,...,agg1,...}(T)$

Standard group-by:

 $\gamma_{C,sum(D)}(S) =$

С	D
10	77
20	196
30	66

SELECT col1,...,agg1(..),agg2(..) FROM T GROUP-BY condition;

С	D
10	33
10	44
20	55
20	66
20	77
30	66

Translation

Every SQL query can be translated into an expression in the Extended RA

Product (<u>pid</u>, name, price)
Purchase (<u>pid</u>, <u>cid</u>, store)
Customer (<u>cid</u>, name, city)

SQL...

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'



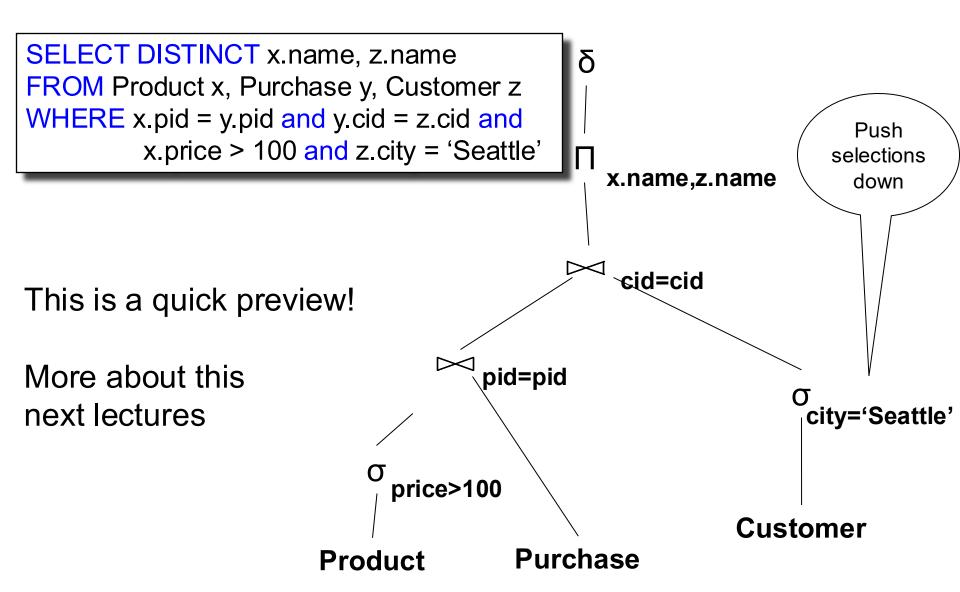
Product (<u>pid</u>, name, price)
Purchase (<u>pid</u>, <u>cid</u>, store)
Customer (<u>cid</u>, name, city)

...to RA

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' x.name,z.name price>100 and city='Seattle' **Specifies** operation order cid=cid pid=pid Customer **Purchase Product**

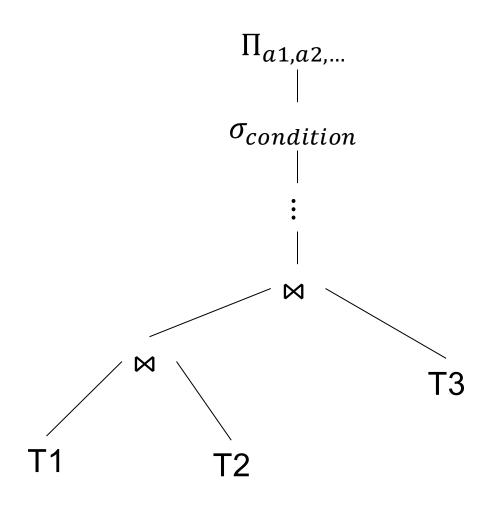
Product(<u>pid</u>, name, price)
Purchase(<u>pid</u>, <u>cid</u>, store)
Customer(<u>cid</u>, name, city)

Customer (cid, name, city) Optimization



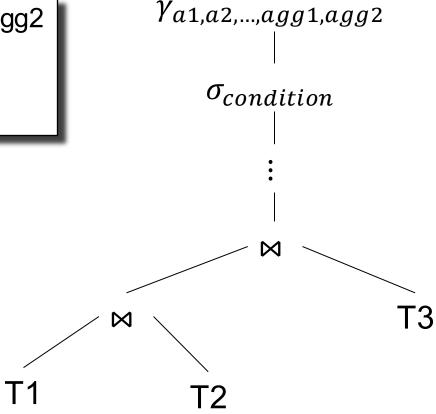
Simple SFW

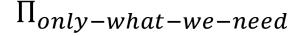
SELECT a1,a2,...
FROM T1, T2, ...
WHERE condition



...add GROUP-BY

```
SELECT a1,a2,...,agg1,agg2
FROM T1, T2, ...
WHERE condition
GROUP BY a1,a2,...
```

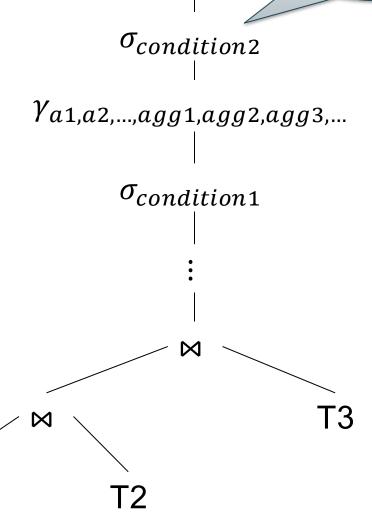




Both HAVING and WHERE use σ

...add HAVING

SELECT a1,a2,...,agg1,agg2 FROM T1, T2, ... WHERE condition1 GROUP BY a1,a2,... HAVING condition2



 SQL queries without subqueries are straightforward to translate

Subqueries may need to be flattened,
 then translated to SQL -- next

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

```
WITH Tmp AS (
                                                        \gamma_{avg(psize)}
         SELECT DISTINCT z.pno, z.psize
         FROM Supplier x, Supply y, Part z
                                                                          Duplicate
         WHERE x.scity = 'Seattle'
                                                                         elimination
           and x.sno=y.sno and y.pno=z.pno)
SELECT avg(psize)
FROM Tmp;
                                                       \Pi_{z.pno,z.psize}
                                                       \bowtie_{y.pno=z.pno}
                                          \bowtie_{x.sno=y.sno}
                                                                       Part z
                              Supplier x
                                                         Supply y
```

Subqueries

Find all suppliers in 'WA' that supply *only* parts under \$100

Subqueries

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

Find all suppliers in 'WA' that supply *only* parts under \$100

Subqueries

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

Find all suppliers in 'WA' that supply *only* parts under \$100

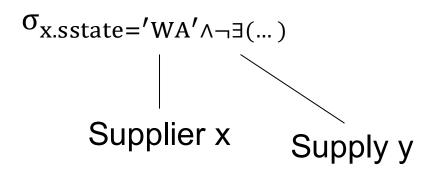
Translate to RA

Subqueries

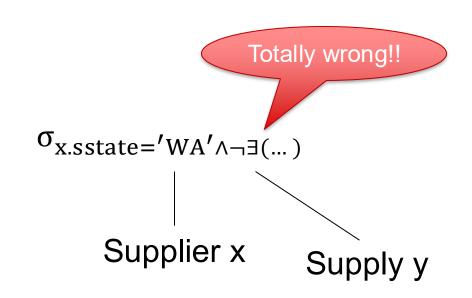
```
SELECT x.sno
FROM Supplier x
WHERE x.sstate = 'WA'
and not exists
(SELECT *
FROM Supply y
WHERE x.sno = y.sno
and y.price > 100)
```

Supplier x

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

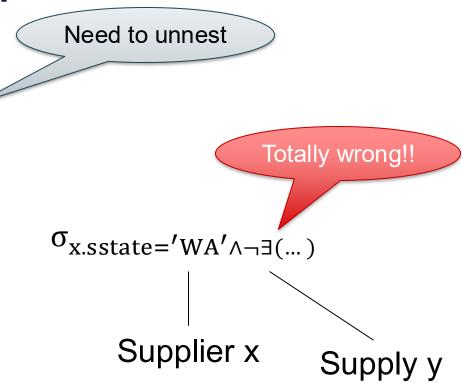


```
SELECT x.sno
FROM Supplier x
WHERE x.sstate = 'WA'
and not exists
(SELECT *
FROM Supply y
WHERE x.sno = y.sno
and y.price > 100)
```



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

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



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

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

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

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

De-Correlation

SELECT x.sno
FROM Supplier x
WHERE x.sstate = 'WA'
and x.sno not in
(SELECT y.sno
FROM Supply y
WHERE y.price > 100)

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

Un-nest

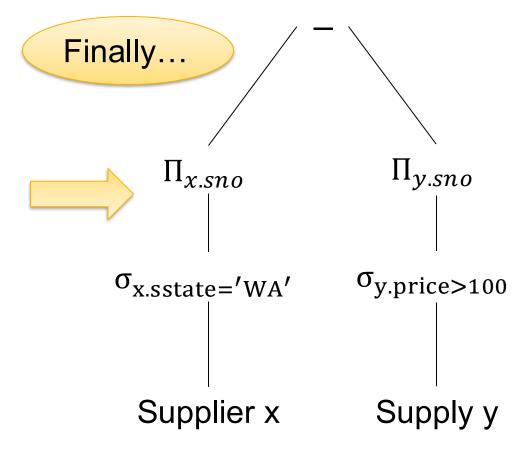
```
(SELECT x.sno
FROM Supplier x
WHERE x.sstate = 'WA')
EXCEPT
(SELECT y.sno
FROM Supply y
WHERE y.price > 100)
```

EXCEPT = set difference

SELECT x.sno
FROM Supplier x
WHERE x.sstate = 'WA'
and x.sno not in
(SELECT y.sno
FROM Supply y
WHERE y.price > 100)

Subqueries

(SELECT x.sno
FROM Supplier x
WHERE x.sstate = 'WA')
EXCEPT
(SELECT y.sno
FROM Supply y
WHERE y.price > 100)



Summary

- User writes in SQL
 - Declarative language
 - Users say WHAT they want
- System: SQL → Relational Algebra
 - Explicit operation order: HOW to compute
 - RA expression a.k.a. Query Plan
- Query Plan: is optimized, then executed

Next lectures