



Introduction to Database Systems

CSE 444



Lecture 18: Relational Algebra

Outline

- ▶ Motivation and sets v.s. bags
- ▶ Relational Algebra
- ▶ Translation from SQL to the Relational Algebra

- ▶ Read Sections 2.4, 5.1, and 5.2
 - ▶ [Old edition: 5.1 through 5.4]
 - ▶ These book sections go over relational operators

The WHAT and the HOW

- ▶ In SQL, we write **WHAT** we want to get from 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**

SQL = WHAT

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

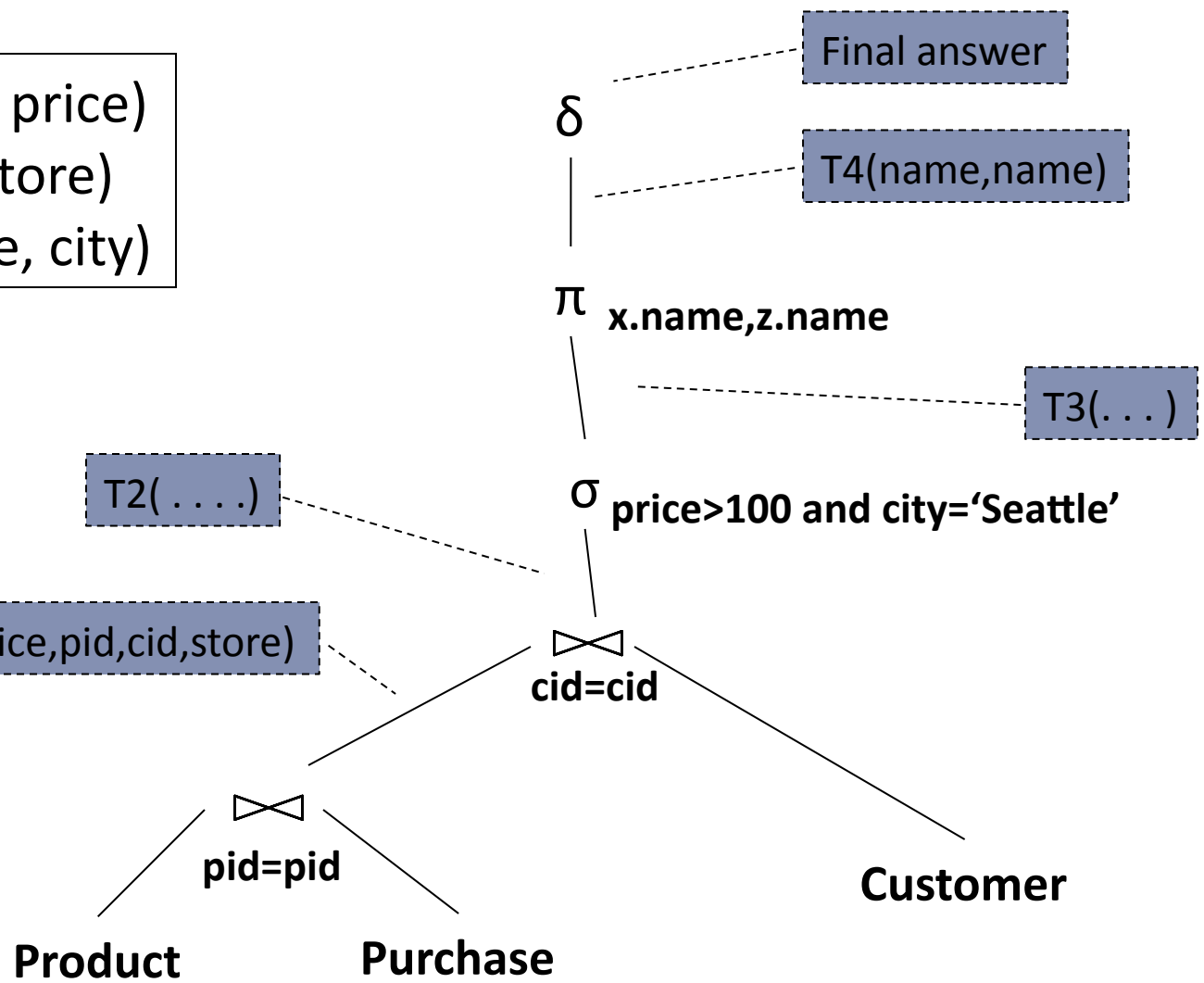
```
SELECT DISTINCT x.name, z.name  
FROM Product x, Purchase y, Customer z  
WHERE x.pid = y.pid and y.cid = z.cid and  
       x.price > 100 and z.city = 'Seattle'
```

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

Relational Algebra = HOW

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

Temporary tables
 T1, T2, ...



Relational Algebra = HOW

The order is now clearly specified:

- ▶ Iterate over PRODUCT...
- ▶ ...join with PURCHASE...
- ▶ ...join with CUSTOMER...
- ▶ ...select tuples with Price>100 and City='Seattle'...
- ▶ ...eliminate duplicates...
- ▶ ...and that's the final answer !

Relations

- ▶ **A relation is a set of tuples**
 - ▶ Sets: {a,b,c}, {a,d,e,f}, { }, . . .
- ▶ But, commercial DBMS's implement **relations that are bags** rather than sets
 - ▶ Bags: {a, a, b, c}, {b, b, b, b, b}, . . .

Sets v.s. Bags

Relational Algebra has two flavors:

- ▶ **Over sets**: theoretically elegant but limited
- ▶ **Over bags**: needed for SQL queries + more efficient
 - ▶ Example: Compute average price of all products

We discuss set semantics

- ▶ We mention bag semantics only where needed

Outline

- ▶ Motivation and sets v.s. bags
- ▶ Relational Algebra
- ▶ Translation from SQL to the Relational Algebra

Relational Algebra

- ▶ **Query language** associated with relational model
- ▶ **Queries specified in an operational manner**
 - ▶ A query gives a step-by-step procedure
- ▶ **Relational operators**
 - ▶ Take one or two relation instances as argument
 - ▶ Return one relation instance as result
 - ▶ Easy to compose into relational algebra expressions

Relational Algebra (1 / 3)

Five basic operators:

- ▶ **Union** (\cup) and **Set difference** ($-$)
- ▶ **Selection**: $\sigma_{\text{condition}}(S)$
 - ▶ Condition is Boolean combination (\wedge, \vee) of terms
 - ▶ Term is: attribute op constant, attr. op attr.
 - ▶ Op is: $<, \leq, =, \neq, \geq,$ or $>$
- ▶ **Projection**: $\pi_{\text{list-of-attributes}}(S)$
- ▶ **Cross-product** or **cartesian product** (\times)

Relational Algebra (2 / 3)

Derived or auxiliary operators:

- ▶ **Intersection** (\cap), **Division** (R/S)
- ▶ **Join**: $R \bowtie_{\theta} S = \sigma_{\theta}(R \times S)$
- ▶ Variations of joins
 - ▶ Natural, equijoin, theta-join
 - ▶ Outer join and semi-join
- ▶ **Rename** $\rho_{B_1, \dots, B_n}(S)$

Relational Algebra (3 / 3)

Extensions for bags

- ▶ **Duplicate elimination:** δ
- ▶ **Group by:** γ [Same symbol as aggregation]
 - ▶ Partitions tuples of a relation into “groups”
- ▶ **Sorting:** τ

Other extensions

- ▶ **Aggregation:** γ (min, max, sum, average, count)

Union and Difference

- ▶ $R1 \cup R2$
- ▶ Example:
 - ▶ $\text{ActiveEmployees} \cup \text{RetiredEmployees}$
- ▶ $R1 - R2$
- ▶ Example:
 - ▶ $\text{AllEmployees} - \text{RetiredEmployees}$

Be careful when applying to bags!

What about Intersection?

- ▶ It is a derived operator
- ▶ $R1 \cap R2 = R1 - (R1 - R2)$
- ▶ Also expressed as a join (will see later)
- ▶ Example
 - ▶ `UnionizedEmployees` \cap `RetiredEmployees`

Relational Algebra (1 / 3)

Five basic operators:

- ▶ **Union** (\cup) and **Set difference** ($-$)
- ▶ **Selection**: $\sigma_{\text{condition}}(S)$
 - ▶ Condition is Boolean combination (\wedge, \vee) of terms
 - ▶ Term is: attribute op constant, attr. op attr.
 - ▶ Op is: $<, \leq, =, \neq, \geq,$ or $>$
- ▶ **Projection**: $\pi_{\text{list-of-attributes}}(S)$
- ▶ **Cross-product** or **cartesian product** (\times)

Selection

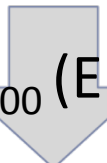
- ▶ Returns all tuples that satisfy a condition
- ▶ Notation: $\sigma_c(R)$
- ▶ Examples
 - ▶ $\sigma_{\text{Salary} > 40000}(\text{Employee})$
 - ▶ $\sigma_{\text{name} = \text{"Smith"}}(\text{Employee})$
- ▶ The condition c can be
 - ▶ Boolean combination (\wedge, \vee) of terms
 - ▶ Term is: attribute op constant, attr. op attr.
 - ▶ Op is: $<, \leq, =, \neq, \geq,$ or $>$

Maps to the **WHERE** clause in SQL!

Selection example

SSN	Name	Salary
1234545	John	20000
5423341	Smith	60000
4352342	Fred	50000

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



SSN	Name	Salary
5423341	Smith	60000
4352342	Fred	50000

Projection

- ▶ Eliminates columns
- ▶ Notation: $\pi_{A_1, \dots, A_n}(R)$
- ▶ Example: project social-security number and names:
 - ▶ $\pi_{SSN, Name}(Employee)$
 - ▶ Output schema: Answer(SSN, Name)

Semantics differs over set or over bags

Projection example

Set semantics: duplicate elimination automatic

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

$\pi_{\text{Name,Salary}}(\text{Employee})$

Name	Salary
John	20000
John	60000

Example

Bag semantics: no duplicate elimination;
need explicit δ

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

$\pi_{\text{Name,Salary}}(\text{Employee})$

Name	Salary
John	20000
John	60000
John	20000

Selection & Projection Examples

Patient

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

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

zip	disease
98125	flu
98125	heart
98120	lung
98120	heart

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

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

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

zip
98120
98125

Relational Algebra (1 / 3)

Five basic operators:

- ▶ **Union** (\cup) and **Set difference** ($-$)
- ▶ **Selection**: $\sigma_{\text{condition}}(S)$
 - ▶ Condition is Boolean combination (\wedge, \vee) of terms
 - ▶ Term is: attribute op constant, attr. op attr.
 - ▶ Op is: $<$, \leq , $=$, \neq , \geq , or $>$
- ▶ **Projection**: $\pi_{\text{list-of-attributes}}(S)$
- ▶ **Cross-product** or **cartesian product** (\times)

Cartesian Product

- ▶ Each tuple in R1 with each tuple in R2
- ▶ Notation: $R1 \times R2$
- ▶ Example:
 - ▶ Employee \times Dependents
- ▶ Rare in practice; mainly used to express joins

Cartesian Product Example

Employee

Name	SSN
John	999999999
Tony	777777777

Dependents

EmployeeSSN	Dname
999999999	Emily
777777777	Joe

Employee x Dependents

Name	SSN	EmployeeSSN	Dname
John	999999999	999999999	Emily
John	999999999	777777777	Joe
Tony	777777777	999999999	Emily
Tony	777777777	777777777	Joe

Relational Algebra (2 / 3)

Derived or auxiliary operators:

- ▶ **Intersection** (\cap), **Division** (R/S)
- ▶ **Join**: $R \bowtie_{\theta} S = \sigma_{\theta}(R \times S)$
- ▶ Variations of joins
 - ▶ Natural, equijoin, theta-join
 - ▶ Outer join and semi-join
- ▶ **Rename** $\rho_{B_1, \dots, B_n}(S)$

Renaming

- ▶ Changes the schema, not the instance
- ▶ Notation: $\rho_{B_1, \dots, B_n}(R)$
- ▶ Example:
 - ▶ $\rho_{\text{LastName}, \text{SocSecNo}}(\text{Employee})$
 - ▶ Output schema:
Answer(LastName, SocSecNo)

Renaming Example

Employee

Name	SSN
John	999999999
Tony	777777777

$\rho_{LastName, SocSecNo}(\text{Employee})$

Employee

LastName	SocSecNo
John	999999999
Tony	777777777

Relational Algebra (2 / 3)

Derived or auxiliary operators:

- ▶ **Intersection** (\cap), **Division** (R/S)
- ▶ **Join**: $R \bowtie_{\theta} S = \sigma_{\theta}(R \times S)$
- ▶ Variations of joins
 - ▶ Natural, equijoin, theta-join
 - ▶ Outer join and semi-join
- ▶ **Rename** $\rho_{B_1, \dots, B_n}(S)$

Different Types of Join

- ▶ **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 π_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

Theta-Join Example

AnonPatient P

age	zip	disease
54	98125	heart
20	98120	flu

AnonJob J

job	age	zip
lawyer	54	98125
cashier	20	98120

$P \bowtie_{P.age=J.age \wedge P.zip=J.zip \wedge P.age < 50} J$

P.age	P.zip	disease	job	J.age	J.zip
20	98120	flu	cashier	20	98120

Equijoin Example

AnonPatient P

age	zip	disease
54	98125	heart
20	98120	flu

AnonJob J

job	age	zip
lawyer	54	98125
cashier	20	98120

$P \bowtie_{P.age=J.age} J$

age	P.zip	disease	job	J.zip
54	98125	heart	lawyer	98125
20	98120	flu	cashier	98120

Natural Join Example

AnonPatient P

age	zip	disease
54	98125	heart
20	98120	flu

AnonJob J

job	age	zip
lawyer	54	98125
cashier	20	98120

$P \bowtie J$

age	P.zip	disease	job
54	98125	heart	lawyer
20	98120	flu	cashier

So which join is it ?

- ▶ When we write $R \bowtie S$ 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

AnonJob 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

Semijoin

- ▶ $R \bowtie S = \Pi_{A_1, \dots, A_n} (R \bowtie S)$
- ▶ Where A_1, \dots, A_n are the attributes in R
- ▶ Example:
 - ▶ Employee \bowtie Dependents

Example of Algebra Queries

Q1: Jobs of patients who have heart disease

$\pi_{\text{job}}(\text{AnonJob} \bowtie (\sigma_{\text{disease}='heart'}(\text{AnonPatient})))$

More 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

$\pi_{\text{sname}}(\text{Supplier} \bowtie \text{Supply} \bowtie (\sigma_{\text{psize} > 10}(\text{Part})))$

Q3: Name of supplier of red parts or parts with size greater than 10

$\pi_{\text{sname}}(\text{Supplier} \bowtie \text{Supply} \bowtie (\sigma_{\text{psize} > 10}(\text{Part}) \cup \sigma_{\text{pcolor} = \text{'red'}}(\text{Part})))$

RA Expressions v.s. Programs

- ▶ An Algebra Expression is like a program
 - ▶ Several operations
 - ▶ Strictly specified order

- ▶ But Algebra expressions have limitations

RA and Transitive Closure

- ▶ Cannot compute “transitive closure”

Name1	Name2	Relationship
Fred	Mary	Father
Mary	Joe	Cousin
Mary	Bill	Spouse
Nancy	Lou	Sister

- ▶ Find all direct and indirect relatives of Fred
- ▶ Cannot express in RA!!! Need to write Java program

Outline

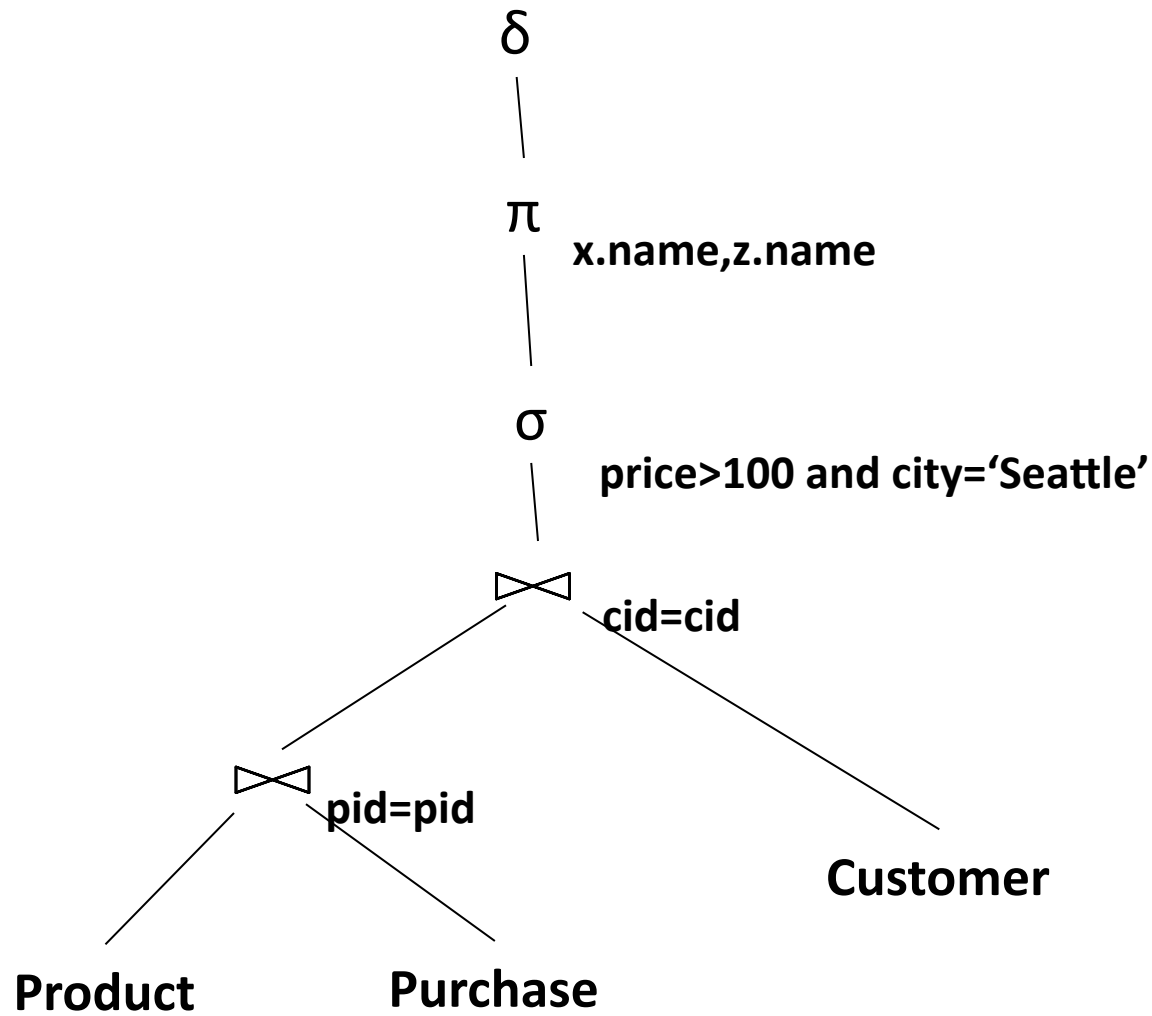
- ▶ Motivation and sets v.s. bags
- ▶ Relational Algebra
- ▶ Translation from SQL to the Relational Algebra

From SQL to RA

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

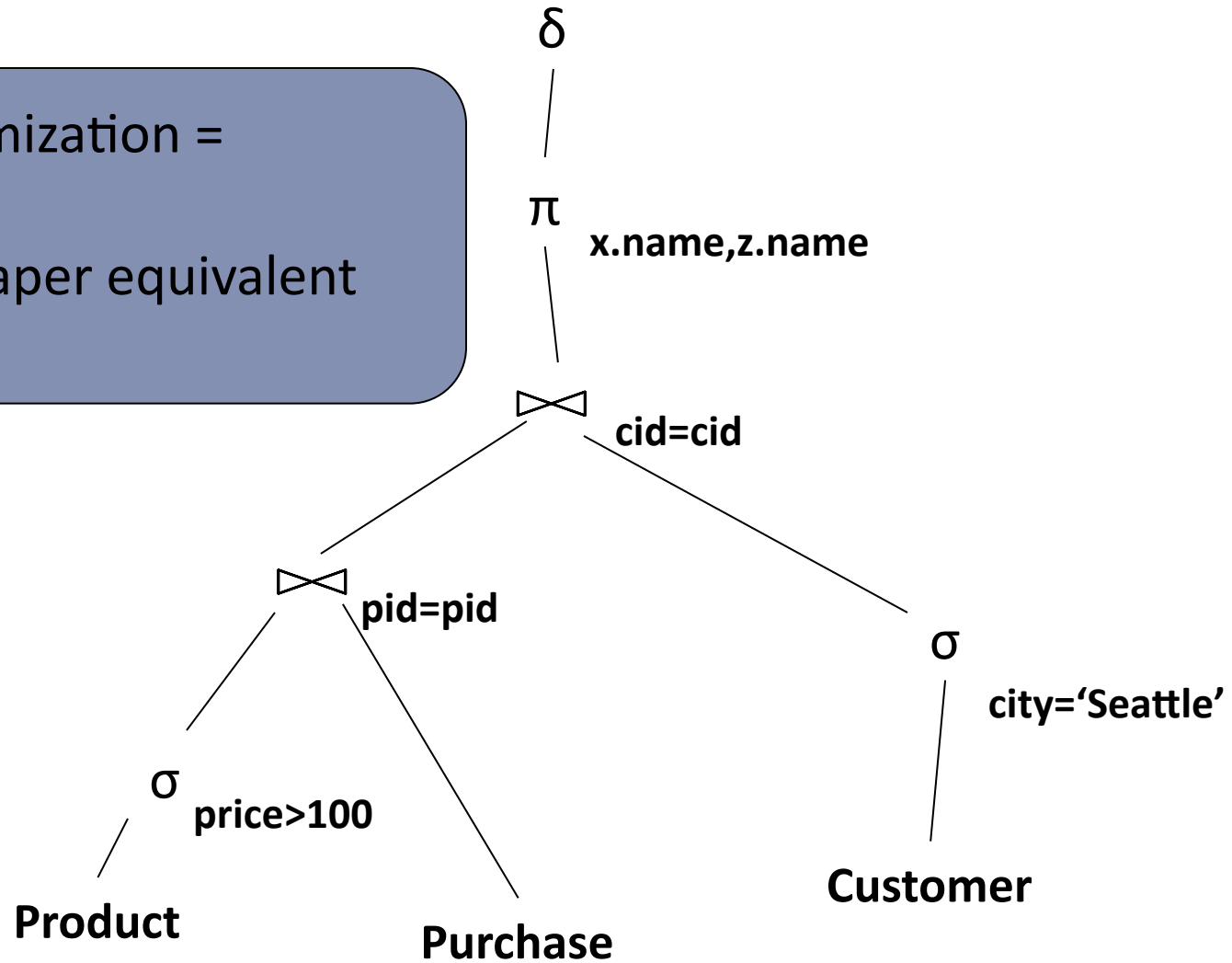
```
SELECT DISTINCT x.name, z.name  
FROM Product x, Purchase y, Customer z  
WHERE x.pid = y.pid and y.cid = z.cid and  
       x.price > 100 and z.city = 'Seattle'
```

From SQL to RA



An Equivalent Expression

Query optimization =
finding cheaper equivalent
expressions



Operators on Bags

- ▶ Duplicate elimination δ
- ▶ Grouping γ
- ▶ Sorting τ

Logical Query Plan

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

T1, T2, T3 = temporary tables

