# Introduction to Data Structure (Data Management)
# Lecture 12

Felipe P. Vista IV

# Reminder

- Everybody, make sure that your name in ZOOM is in the following format:
  - University ID Num Name (no "( )")
  - Ex: 202054321 Juan Dela Cruz

  -

  - Not changing your name to this format
    - you might be marked Absent

    * → absent?

- # Relational Algebra (Part 2)

- # Query Evaluation

- # From Logical Plans to Physical Plans

# Relational Algebra (Part 2)

# So which join is it?

- Theta join: $R \bowtie_\theta S = \sigma_\theta (R \times S)$
    - Join of R and S with a join condition $\theta$
    - Cross product followed by a selection $\theta$

# So which join is it?

- Theta join: $R \bowtie_\theta S = \sigma_\theta(R \times S)$
  - Join of R and S with a join condition $\theta$
  - Cross product followed by a selection $\theta$

- Equijoin: $R \bowtie_\theta S = \sigma_\theta(R \times S)$
  - Join condition $\theta$ consist only of equalities

# So which join is it?

- Theta join: $R \bowtie_\theta S = \sigma_\theta(R \times S)$
  - Join of R and S with a join condition $\theta$
  - Cross product followed by a selection $\theta$

- Equijoin: $R \bowtie_\theta S = \sigma_\theta(R \times S)$
  - Join condition $\theta$ consist only of equalities

- Natural join1: $R \bowtie S = \pi_A(\sigma_\theta(R \times S))$
  - Equijoin
  - Equality on all fields with same name in R and in S
  - Projection $\pi_A$ drops all redundant attributes

# So which join is it?

- When we use $R \bowtie S$,

    - we usually mean an equijoin

    - but 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
  - Does not eliminate/remove duplicate columns

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

# Outer Join Example (Review)

**Lecture**

| Id | Name | Room | Class |
|----|------|------|-------|
| 1 | S1 | 406 | Discrete Math |
| 2 | S2 | 408 | Data Struc |
| 3 | S3 | 409 | Programming |

**Research**

| Project | Focus | Room |
|---------|-------|------|
| Brain | Deep Learning | 406 |
| Ships | Sensors | 408 |

$$L \bowtie R$$

|  |  |  |  |  |  |  |
|--|--|--|--|--|--|--|
|  |  |  |  |  |  |  |
|  |  |  |  |  |  |  |
|  |  |  |  |  |  |  |

# Outer Join Example (Review)

**Lecture**

| Id | Name | Room | Class |
|----|------|------|-------|
| 1 | S1 | 406 | Discrete Math |
| 2 | S2 | 408 | Data Struc |
| 3 | S3 | 409 | Programming |

**Research**

| Project | Focus | Room |
|---------|-------|------|
| Brain | Deep Learning | 406 |
| Ships | Sensors | 408 |

$$L \bowtie R$$

| L.Id | L.Name | L.Room | L.Class | R.Project | R.Focus | R.Room |
|------|--------|--------|---------|-----------|---------|--------|
| 1 | S1 | 406 | Discrete Math | Brain | Deep Learning | 406 |
| 2 | S2 | 408 | Data Struc | Ships | Sensors | 408 |
| 3 | S3 | 409 | Programming | Null | Null | Null |

# More Examples

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

a. Name of supplier of parts with size greater than 10

a. Name of supplier of red parts or parts with size greater than 10

# More Examples

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

a. **Name of supplier of parts with size greater than 10**

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

a. **Name of supplier of red parts or parts with size greater than 10**

$$\Pi_{sname}(Supplier \bowtie Supply \bowtie (\sigma_{pcolor='red' OR psize>10}(Part)))$$

# More Examples

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

a.  Name of supplier of parts with size greater than 10

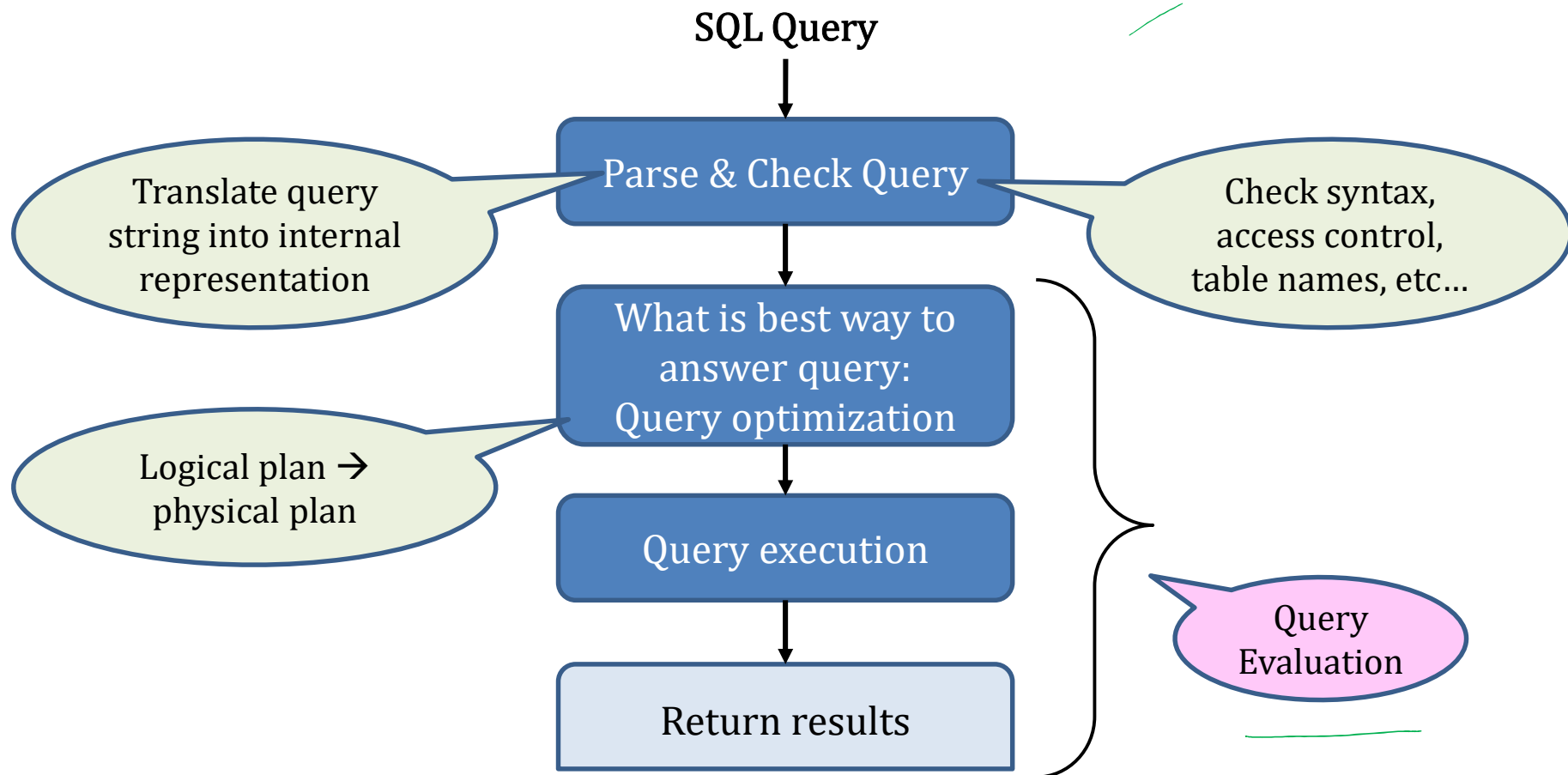$$\Pi_{sname}(\text{Supplier} \bowtie \text{Supply} \bowtie (\sigma_{psize>10})(\text{Part}))$$

a.  Name of supplier of red parts or parts with size greater than 10

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

INTRO TO DATA STRUCTURE

# Query Evaluation

# Query Evaluation Steps

SQL Query

Parse & Check Query

Translate query string into internal representation

Check syntax, access control, table names, etc...

What is best way to answer query: Query optimization

Logical plan → physical plan

Query execution

Query Evaluation

Return results

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

# From SQL to RA

```
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="Jeonju"
```
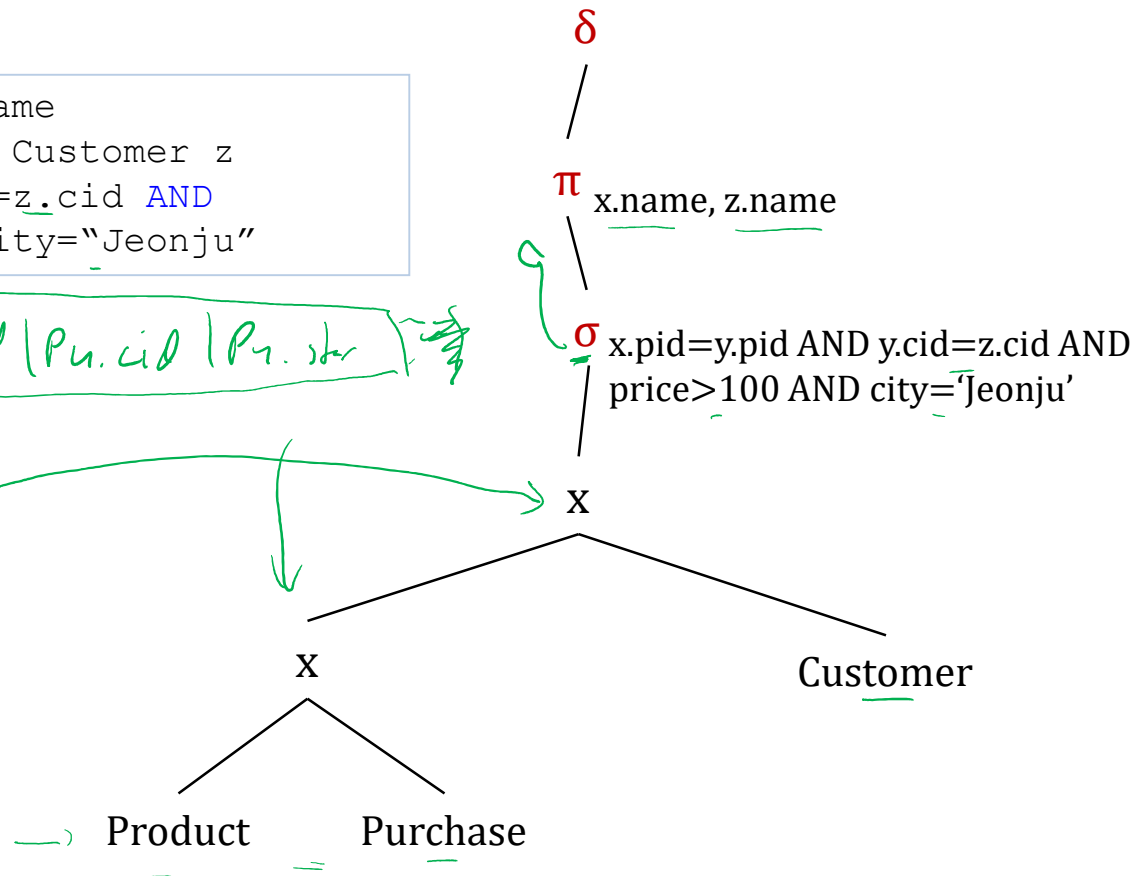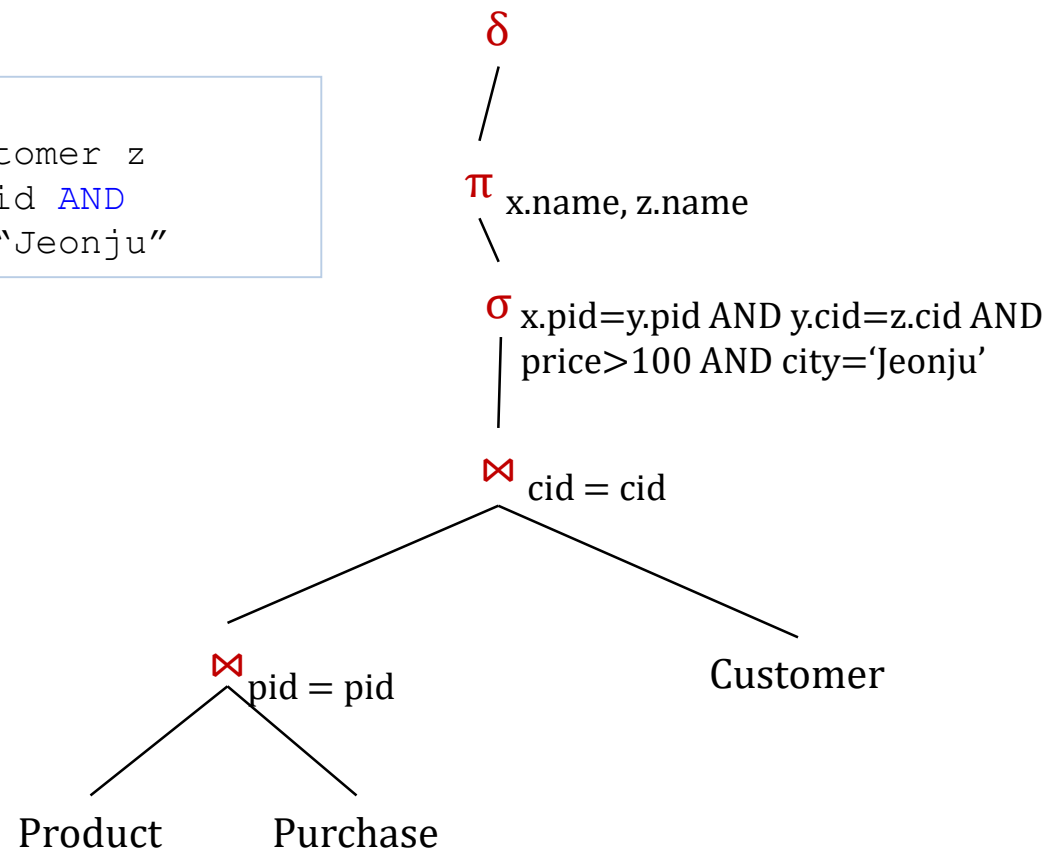
Product(**pid**, name, price)
Purchase(**pid**, **cid**, store)
Customer(**cid**, name, city)

# From SQL to RA

```
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="Jeonju"
```

$\delta$

$\pi$ x.name, z.name

$\sigma$ x.pid=y.pid AND y.cid=z.cid AND
price>100 AND city='Jeonju'

| P.pid | P.name | P.price | Pu.pid | Pu.cid | Pu.stor |
|-------|--------|---------|--------|--------|---------|

| c.cid | c.na | c.city |
|-------|------|--------|

x

x                                Customer

Product        Purchase

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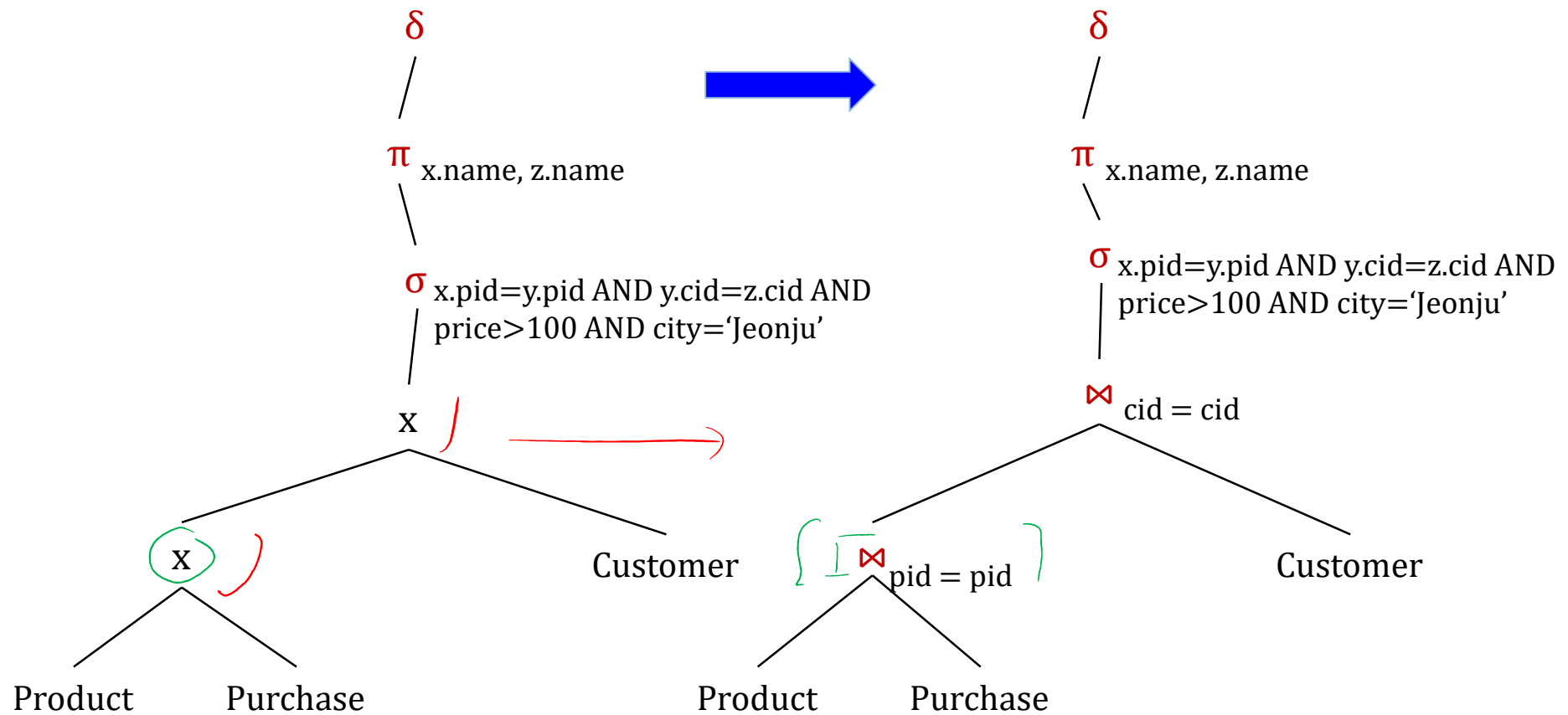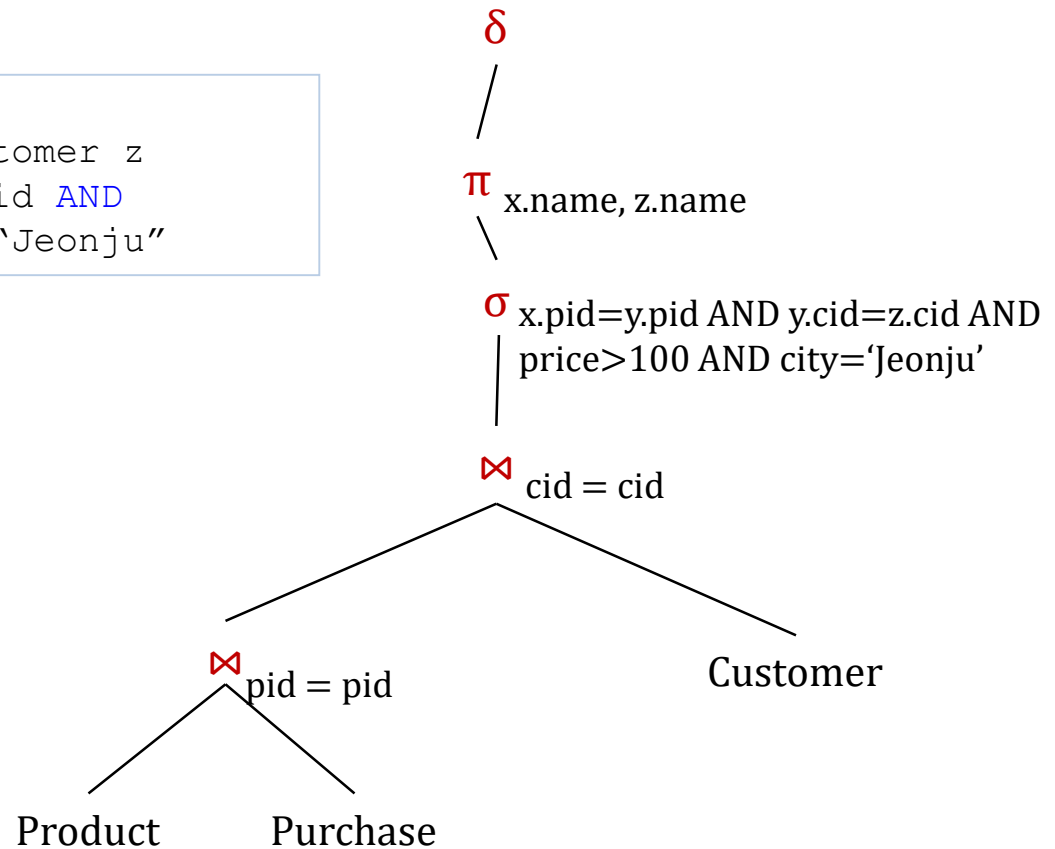Customer(**cid**, name, city)

# From SQL to RA

```
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="Jeonju"
```

$\delta$

$\pi$ x.name, z.name

$\sigma$ x.pid=y.pid AND y.cid=z.cid AND price>100 AND city='Jeonju'

$\bowtie$ cid = cid

$\bowtie$ pid = pid

Customer

Product       Purchase

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

# From SQL to RA

δ

π x.name, z.name

σ x.pid=y.pid AND y.cid=z.cid AND price>100 AND city='Jeonju'

x

⨝ x

Product    Purchase

Customer

δ

π x.name, z.name

σ x.pid=y.pid AND y.cid=z.cid AND price>100 AND city='Jeonju'

⨝ cid = cid

⨝ pid = pid

Product    Purchase

Customer

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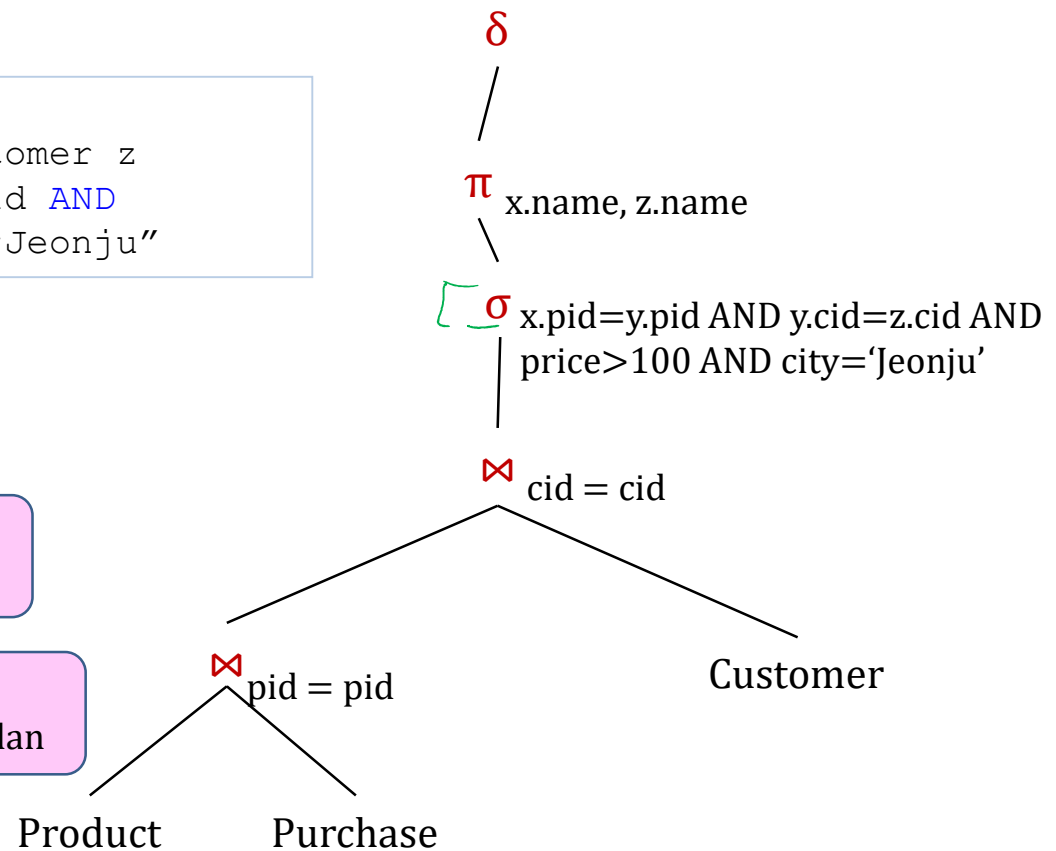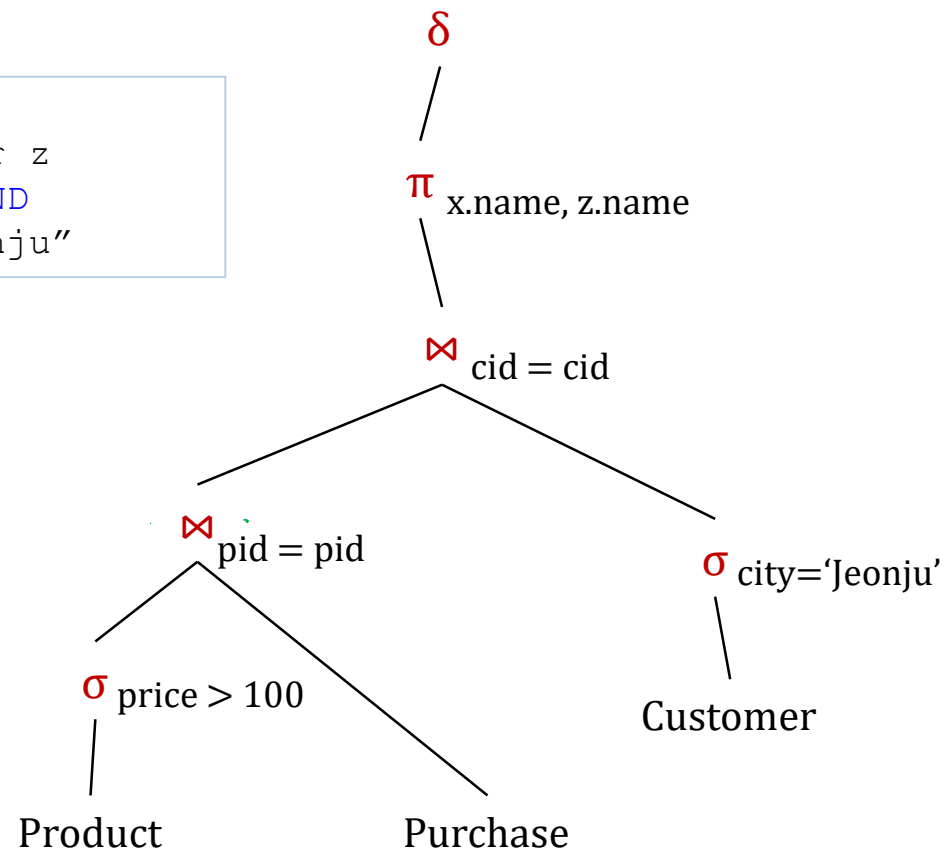Customer(**cid**, name, city)

# From SQL to RA

```
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="Jeonju"
```

Is there another way

$\delta$

$\pi$ x.name, z.name

$\sigma$ x.pid=y.pid AND y.cid=z.cid AND price>100 AND city='Jeonju'

$\bowtie$ cid = cid

$\bowtie$ pid = pid

Customer

Product        Purchase

# From SQL to RA

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

```sql
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="Jeonju"
```

We can do this way!

Push selections down the query plan

Query Optimization: Find an equivalent optimal plan

$\delta$

$\pi$ x.name, z.name

$\sigma$ x.pid=y.pid AND y.cid=z.cid AND price>100 AND city='Jeonju'

$\bowtie$ cid = cid

$\bowtie$ pid = pid

Customer

Product　　Purchase

# From SQL to RA

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

```sql
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="Jeonju"
```

We can do this way!

Push selections down the query plan

Query Optimization:
Find an equivalent optimal plan

$\delta$

$\pi_{x.name, z.name}$

$\bowtie cid = cid$

$\bowtie pid = pid$

$\sigma city='Jeonju'$

$\sigma price > 100$

Product

Purchase

Customer

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

# From SQL to RA

(A)

δ

π x.name, z.name

σ x.pid=y.pid AND y.cid=z.cid AND
  price>100 AND city='Jeonju'

⋈ cid = cid

⋈ pid = pid

Product    Purchase

Customer

(B)

δ

π x.name, z.name

⋈ cid = cid

⋈ pid = pid

σ price > 100

Product    Purchase

σ city='Jeonju'

Customer

# Extended RA: Operator on Bags

- Duplicate elimination ($\delta$)

- Grouping & Aggregation ($\gamma$)

- Sorting ($\tau$)

# Logical Query Plan

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

# Logical Query Plan

SGL

```
SELECT city, count(*)
FROM sales
GROUP BY city
WHERE sum(price) > 100
```
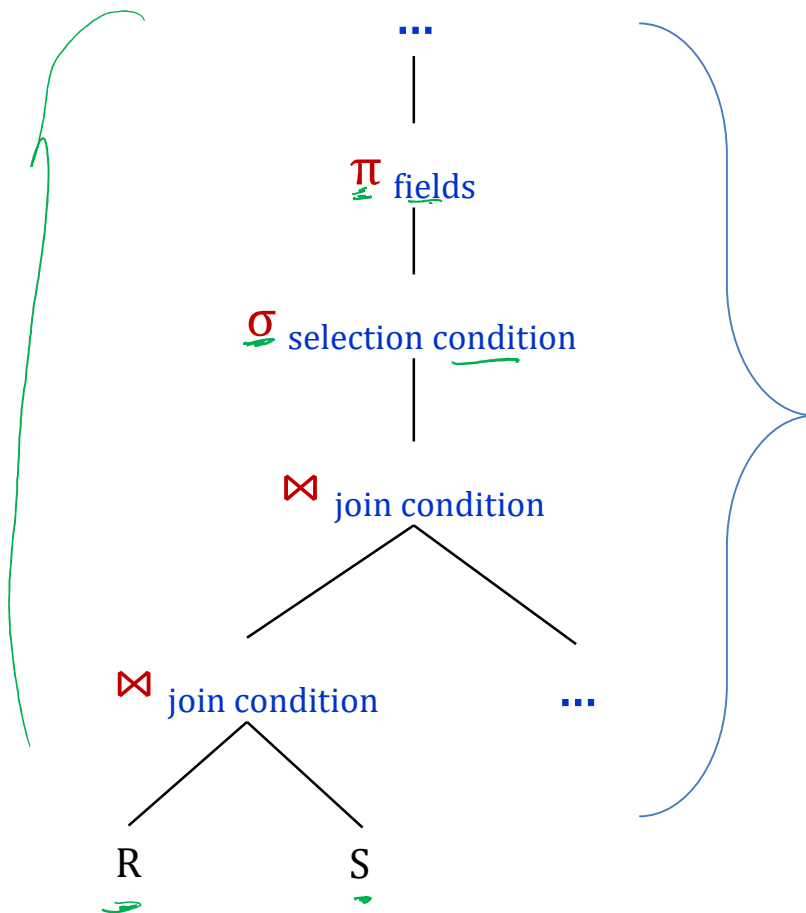
$\pi_{city, c}$

$\sigma_{p > 100}$

$\gamma_{city, sum(price) \rightarrow p, count(*) \rightarrow c}$

Sales(product, city, price)

T3(city, c)

T2(city, p, c)

$T_1$(city, p, c)

$T_1, T_2, T_3$ = temporary tables

# Typical for Block (1/2)

```
SELECT fields
FROM R, S
WHERE condition
```

**SELECT-PROJECT-JOIN**
**Query**

# Typical for Block (1/2)

...

$\pi$ fields

$\sigma$ selection condition

$\bowtie$ join condition

$\bowtie$ join condition          ...

R          S

```
SELECT fields
FROM R, S
WHERE condition
```

**SELECT-PROJECT-JOIN
Query**

# Typical for Block (2/2)

```
SELECT fields
FROM R, S
WHERE condition
GROUP BY fields
HAVING condition
```

# Typical for Block (2/2)

$\pi$ fields

$\sigma$ HAVING condition

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

$\sigma$ WHERE condition

$\bowtie$ join condition

R          S

...          ...

```
SELECT fields
FROM R, S
WHERE condition
GROUP BY fields
HAVING condition
```

# How About Subqueries?

Supplier(**sno**, sname, scity, sprov)
Part(**pno**, pname, psize, pcolor)
Supply(**sno**, **pno**, qty, price)

```
SELECT Q.sno
FROM Supplier Q
WHERE Q.sprov = 'CAPIZ'
   AND NOT EXISTS
   (SELECT *
    FROM Supply P
    WHERE P.sno = Q.sno
        AND P.price > 100;
```

Correlation!

# How About Subqueries?

Supplier(**sno**, sname, scity, sprov)
Part(**pno**, pname, psize, pcolor)
Supply(**sno**, **pno**, qty, price)

```
SELECT Q.sno
FROM Supplier Q
WHERE Q.sprov = 'CAPIZ'
   AND NOT EXISTS
   (SELECT *
    FROM Supply P
    WHERE P.sno = Q.sno
        AND P.price > 100;
```

De-Correlation

# How About Subqueries?

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

```sql
SELECT Q.sno
FROM Supplier Q
WHERE Q.sprov = 'CAPIZ'
   AND NOT EXISTS
   (SELECT *
    FROM Supply P
    WHERE P.sno = Q.sno
        AND P.price > 100;
```

De-Correlation

```sql
SELECT Q.sno
FROM Supplier Q
WHERE Q.sprov = 'CAPIZ'
   AND Q.sno NOT IN
   (SELECT P.sno
    FROM Supply P
    WHERE P.price > 100;
```

# How About Subqueries?

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

```sql
SELECT Q.sno
FROM Supplier Q
WHERE Q.sprov = 'CAPIZ'
    EXCEPT
   (SELECT P.sno
    FROM Supply P
    WHERE P.price > 100;
```

**Un-nesting**

**EXCEPT** = set difference

# How About Subqueries?

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

```
SELECT Q.sno
FROM Supplier Q
WHERE Q.sprov = 'CAPIZ'
   EXCEPT
  (SELECT P.sno
   FROM Supply P
   WHERE P.price > 100;
```

$n = 100$

**EXCEPT** = set difference

Un-nesting

```
SELECT Q.sno
FROM Supplier Q
WHERE Q.sprov = 'CAPIZ'
   AND Q.sno NOT IN
   (SELECT P.sno
    FROM Supply P
    WHERE P.price > 100;
```

# How About Subqueries?

Supplier(**sno,** sname, scity, sprov)
Part(**pno,** pname, psize, pcolor)
Supply(**sno, pno,** qty, price)
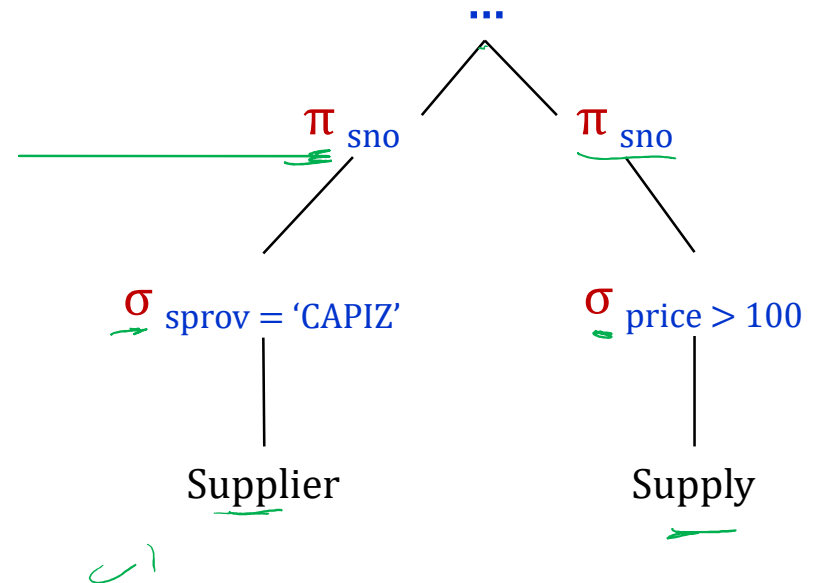
```
(SELECT Q.sno
FROM Supplier Q
WHERE Q.sprov = 'CAPIZ'
   EXCEPT
   (SELECT P.sno
    FROM Supply P
    WHERE P.price > 100;
```

Finally...

# How About Subqueries?

Supplier(**sno**, sname, scity, sprov)
Part(**pno**, pname, psize, pcolor)
Supply(**sno**, **pno**, qty, price)

```
(SELECT Q.sno
 FROM Supplier Q
 WHERE Q.sprov = 'CAPIZ'
     EXCEPT
     (SELECT P.sno
      FROM Supply P
      WHERE P.price > 100;
```

Finally...
Logical Plan!

...

$\pi_{sno}$          $\pi_{sno}$

$\sigma_{sprov = 'CAPIZ'}$          $\sigma_{price > 100}$

Supplier          Supply

# Reminder

- Everybody, make sure that your name in ZOOM is in the following format:
  - University ID Num Name (no "( )")
  - Ex: 202054321 Juan Dela Cruz
  -
  - Not changing your name to this format
    - you might be marked Absent
    - * → absent?

# From Logical Plans to Physical Plans

# Physical Operators

- Each of the logical operators may have one or more implementations = physical operations

- Discuss some basic physical operators, paying special attention to join

# Main Memory Algorithms

- ## Logical operator

```
Product(pid, name, price)
Purchase(pid, cid, store)
```

Product(pid, name, price) $\bowtie_{pid=pid}$ Purchase(pid, cid, store)

Propose three physical operators for the join,
   assuming the tables are in main memory:
1. Nested Loop Join              O( ?? )
2. Merge Join                     O( ?? )
3. Hash Join                      O( ?? )

Take note that **pid** is a key.

* time complexity : the computational complexity that describes the
amount of time it takes to run an algorithm
** '$n$' – is the input size

# Main Memory Algorithms

- ## Logical operator

Product(**pid**, name, price)
Purchase(**pid**, **cid**, store)

Product(**pid**, name, price) $\bowtie_{pid=pid}$ Purchase(**pid**, **cid**, store)

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

1. Nested Loop Join          → $O( n^2 )$          Two nested loops
2. Merge Join                $O( ?? )$
3. Hash Join                 $O( ?? )$

\* time complexity : the computational complexity that describes the amount of time it takes to run an algorithm
\*\* '$n$' – is the input size

# Main Memory Algorithms

- ## Logical operator

```
Product(pid, name, price)
Purchase(pid, cid, store)
```

Product(**pid**, name, price) $\bowtie_{pid=pid}$ Purchase(**pid**, **cid**, store)

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

1. Nested Loop Join               $O(n^2)$
2. Merge Join                   $O(n \log n)$
3. Hash Join                    $O(??)$

> Sort both: O(n log n)
> Merge: O(n)

\* time complexity : the computational complexity that describes the amount of time it takes to run an algorithm
\*\* '$n$' – is the input size

# Main Memory Algorithms

- Logical operator

> Product(**pid**, name, price) $\bowtie_{pid=pid}$ Purchase(**pid**, **cid**, store)

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

1. Nested Loop Join      $O(n^2)$
2. Merge Join      $O(n \log n)$
3. Hash Join      $O(n) \dots O(n^2)$

> Add n to hash: $O(n)$?
> Lookup n in hash: $O(n)$

\* time complexity : the computational complexity that describes the amount of time it takes to run an algorithm
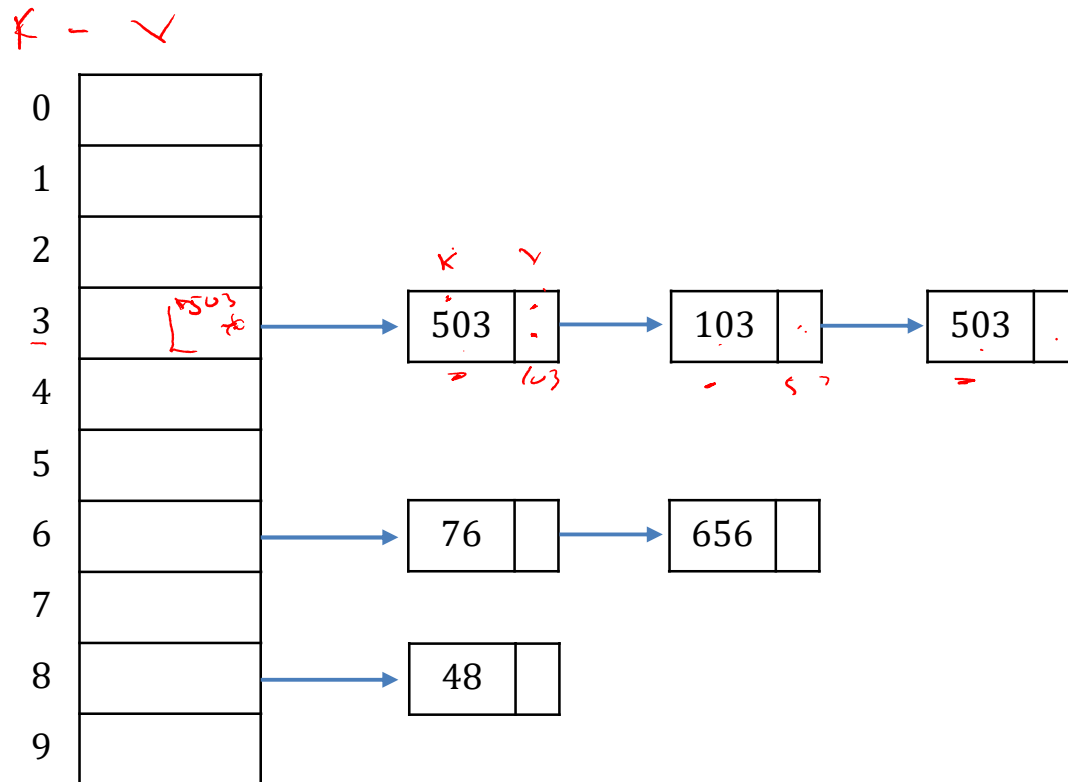\*\* '*n*' – is the input size

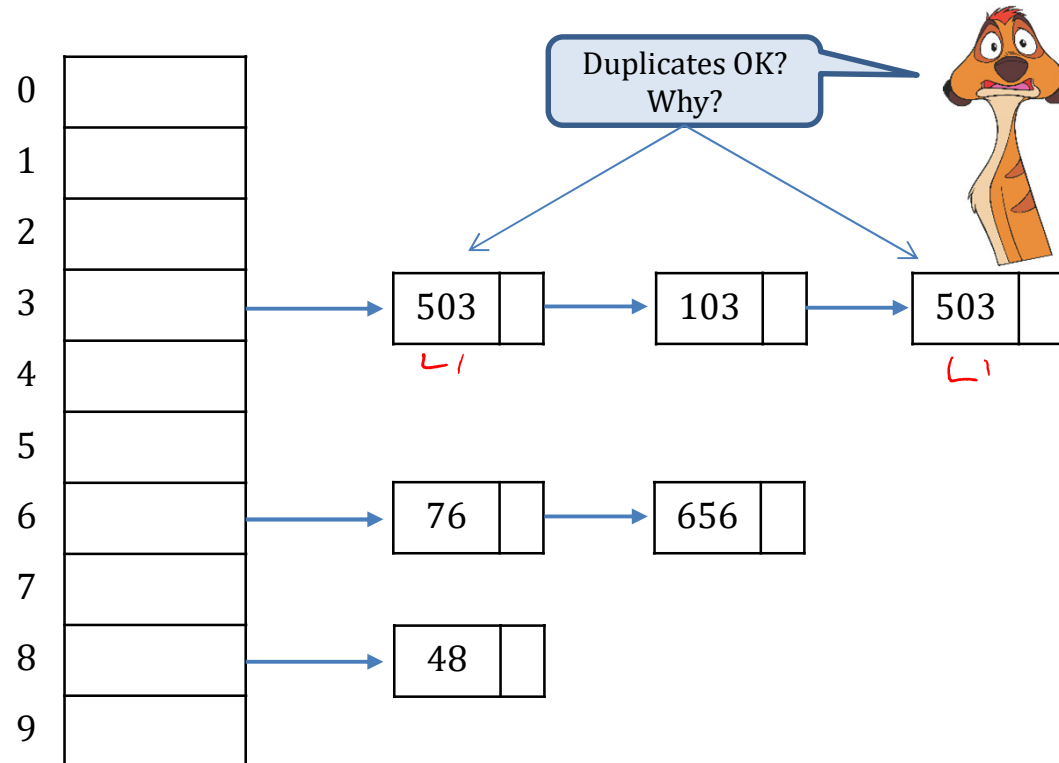# Brief Review of Hash Tables

A naive hash function:

$$h(x) = x \bmod 10$$

Operations:

find(103) = ??
insert(488) = ??

# Brief Review of Hash Tables

A naive hash function:

$$h(x) = x \bmod 10$$

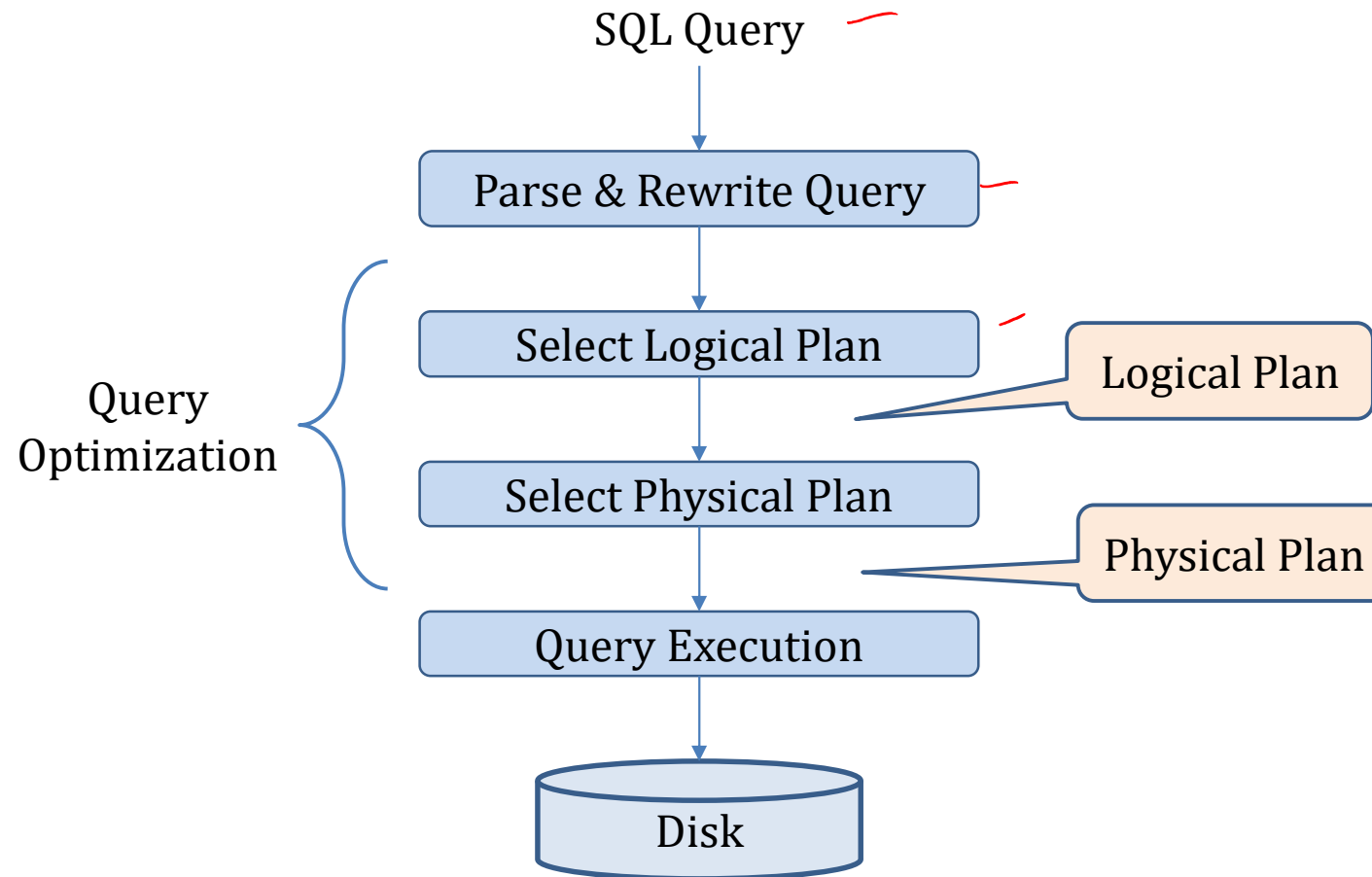Operations:

find(103) = ??
insert(488) = ??

# Brief Review of Hash Tables

Values can have different tuple IDs!

A naive hash function:

$$h(x) = x \bmod 10$$

Duplicates OK? Why?

Operations:

find(103) = ??
insert(488) = ??

# Query Evaluation Steps Review

SQL Query

↓

Parse & Rewrite Query

↓

Select Logical Plan ← Logical Plan

↓

Query Optimization

Select Physical Plan ← Physical Plan

↓

Query Execution

↓

Disk

# Relational Algebra

Supplier(**sid**, sname, scity, sprov)
Supply(**sid**, **pno**, quantity)

```
SELECT sname
FROM Supplier x, Supply y
WHERE x.sid = y.sid
  AND y.pno = 2
  AND x.scity = 'Jeonju'
  AND x.prov = 'CAPIZ';
```

*(handwritten: "∧" = AND)*

Give a relational algebra expression for this query.

$$\pi_{xxxx} (\sigma_{xxxx='xxxx' \land xxxx='xxxx' \land xxxx=xxxx} ( Supplier \bowtie_{xxxx=xxxx} Supply))$$

*(handwritten annotations: sname; y.pno = 2; x.scity = 'Jeonju'; x.prov = 'CAPIZ'; x.sid = y.sid)*

# Relational Algebra

Supplier(**sid**, sname, scity, sprov)
Supply(**sid**, **pno**, quantity)

```
SELECT sname
FROM Supplier x, Supply y
WHERE x.sid = y.sid
   AND y.pno = 2
   AND x.scity = 'Jeonju'
   AND x.prov = 'CAPIZ';
```
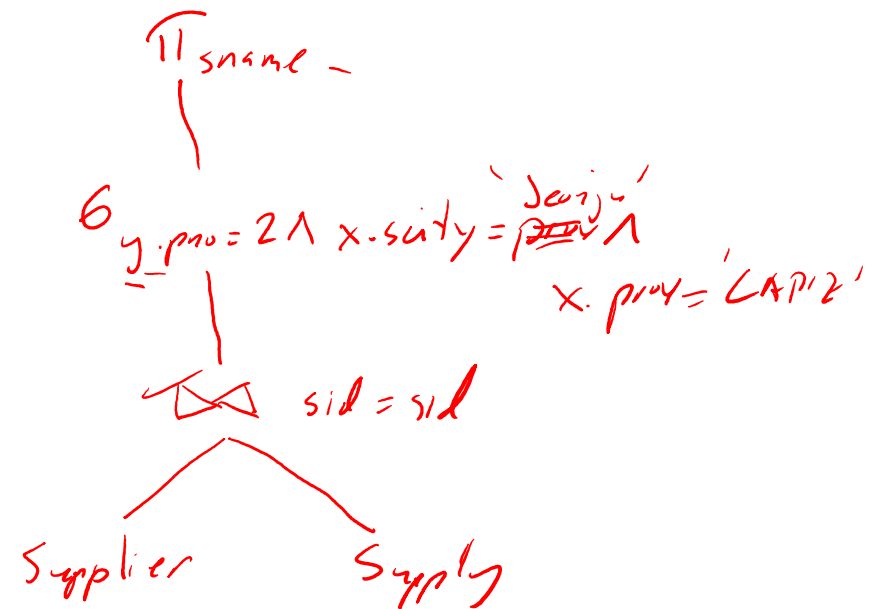
Give a relational algebra expression for this query.

$$\pi_{sname}(\sigma_{scity='Jeonju' \wedge sprov='CAPIZ' \wedge pno = 2}(\text{Supplier} \bowtie_{sid=sid} \text{Supply}))$$

# Relational Algebra

Supplier(**sid**, sname, scity, sprov)
Supply(**sid**, **pno**, quantity)

```
SELECT sname
FROM Supplier x, Supply y
WHERE x.sid = y.sid
  AND y.pno = 2
  AND x.scity = 'Jeonju'
  AND x.prov = 'CAPIZ';
```

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

$\Pi_{sname}$

$\sigma_{y.pno=2 \land x.scity='Jeonju' \land x.prov='CAPIZ'}$

$\bowtie_{sid=sid}$

Supplier      Supply

# Relational Algebra

```
Supplier(sid, sname, scity, sprov)
Supply(sid, pno, quantity)
```

```
SELECT sname
FROM Supplier x, Supply y
WHERE x.sid = y.sid
  AND y.pno = 2
  AND x.scity = 'Jeonju'
  AND x.prov = 'CAPIZ';
```

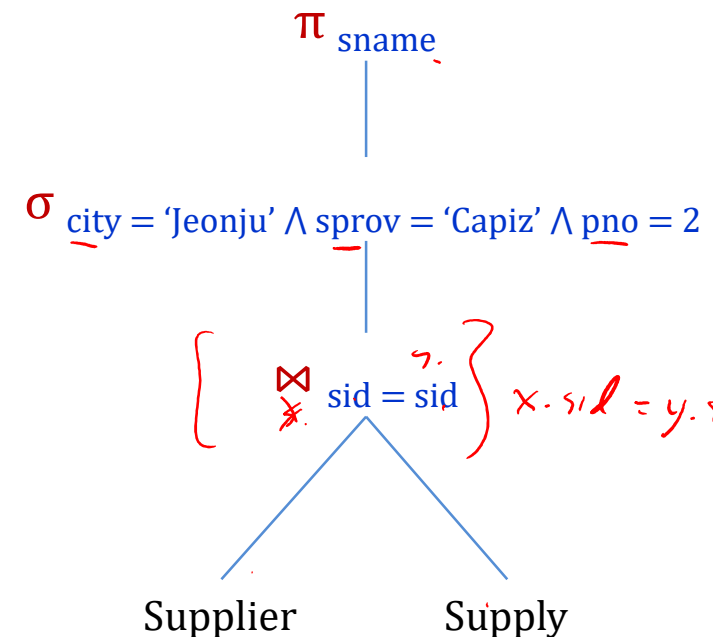Relational algebra expression is also called the "logical query plan"

$\pi_{sname}$

$\sigma_{city = 'Jeonju' \wedge sprov = 'Capiz' \wedge pno = 2}$

$\bowtie_{sid = sid}$   x.sid = y.sid

Supplier          Supply

# Physical Query Plan 1

Supplier(**sid**, sname, scity, sprov)
Supply(**sid**, **pno**, quantity)

(On the fly)    $\pi$ sname

(On the fly)    $\sigma$ city = 'Jeonju' $\wedge$ sprov = 'Capiz' $\wedge$ pno = 2

(Nested loop)    $\bowtie$ sid = sid

Supplier      Supply
(File scan)    (File scan)

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

```
SELECT sname
FROM Supplier x, Supply y
WHERE x.sid = y.sid
  AND y.pno = 2
  AND x.scity = 'Jeonju'
  AND x.prov = 'CAPIZ';
```
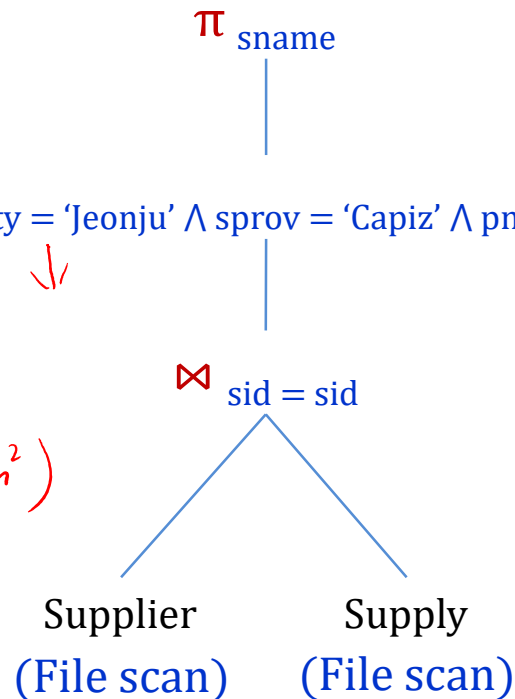
# Physical Query Plan 2

Supplier(**sid**, sname, scity, sprov)
Supply(**sid**, **pno**, quantity)

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

(On the fly)   $\sigma_{city = 'Jeonju' \land sprov = 'Capiz' \land pno = 2}$

PP₂

(Hash Join)    $\bowtie_{sid = sid}$

$O(n) \lll \; \sigma(n^2)$

Supplier       Supply
(File scan)    (File scan)

Same logical query plan
Different physical plan

```
SELECT sname
FROM Supplier x, Supply y
WHERE x.sid = y.sid
  AND y.pno = 2
  AND x.scity = 'Jeonju'
  AND x.prov = 'CAPIZ';
```

SQL

Vanilla Loop    $O(n^2)$

# Physical Query Plan 3

$n^2 = ?$

Sort-merge    $O(n \log(n))$

Supplier(**sid**, sname, scity, sprov)
Supply(**sid**, **pno**, quantity)

Hash Loop Join    $O(n)$

$n = 25,000$

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

Same logical query plan
Different physical plan

(Sort-merge join)    $\bowtie_{\text{sid} = \text{sid}}$    $LP_L \longleftarrow \longrightarrow SQL$

~~(Scan & write to $T_1$)~~
(Scan & write to $T_1$)        (Scan & write to $T_2$)

```
SELECT sname
FROM Supplier x, Supply y
WHERE x.sid = y.sid
   AND y.pno = 2
   AND x.scity = 'Jeonju'
   AND x.prov = 'CAPIZ';
```
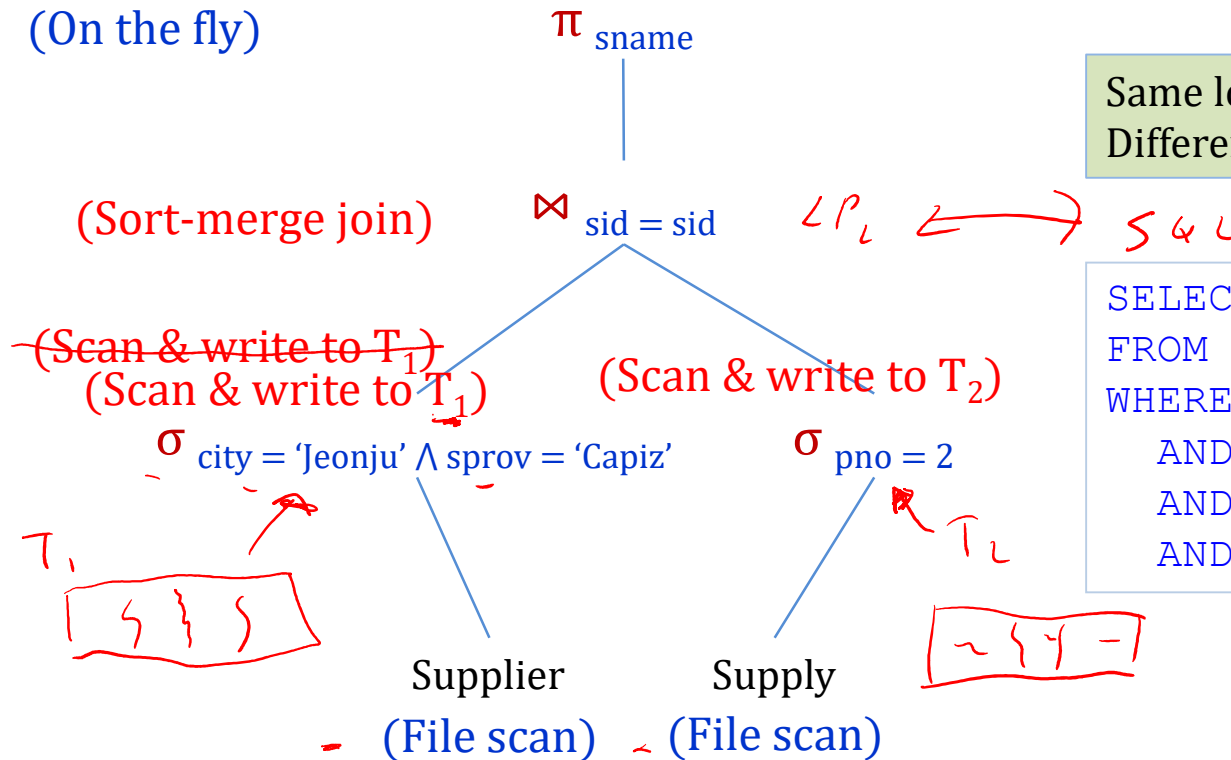
$\sigma_{\text{city} = \text{'Jeonju'} \wedge \text{sprov} = \text{'Capiz'}}$        $\sigma_{\text{pno} = 2}$

$T_1$        $T_2$

Supplier        Supply

(File scan)    (File scan)

# Query Optimization Problem

- For each SQL query … many logical plans

  *Ex: 1 SQL query ⇒ 2 Logical Plans ⇒ LP_A ⇒ 2 PP*
  *LP_A ⇒ 1 PP*

- For each logical plan … many physical plans

- How to find a fast physical plan?
  - Will discuss in a few lectures

# Thank you.