# Relational Algebra

# DBMS Architecture

How does a SQL engine work ?

- SQL query $\rightarrow$ relational algebra plan
- Relational algebra plan $\rightarrow$ Optimized plan
- Execute each operator of the plan

# Relational Algebra

- Five operators:
  - Union: $\cup$
  - Difference: -
  - Selection: $\sigma$
  - Projection: $\Pi$
  - Cartesian Product: $\times$
- Derived or auxiliary operators:
  - Intersection, complement
  - Joins (natural,equi-join, theta join, semi-join)
  - Renaming: $\rho$

# 1. Union and 2. Difference

- R1 $\cup$ R2
- Example:
  - ActiveEmployees $\cup$ RetiredEmployees

- R1 – R2
- Example:
  - AllEmployees -- RetiredEmployees

# 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

# 3. Selection

- Returns all tuples which satisfy a condition
- Notation: $\sigma_c(R)$
- Examples
  - $\sigma_{Salary > 40000}$ (Employee)
  - $\sigma_{name = \text{"Smith"}}$ (Employee)
- The condition c can be $=, <, \leq, >, \geq, <>$

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

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

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

# 4. Projection

- Eliminates columns, then removes duplicates

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

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

| SSN | Name | Salary |
|---|---|---|
| 1234545 | John | 200000 |
| 5423341 | John | 600000 |
| 4352342 | John | 200000 |

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

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

# 5. Cartesian Product

- Each tuple in R1 with each tuple in R2
- Notation: R1 $\times$ R2
- Example:
  - Employee $\times$ Dependents
- Very 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 |

# Renaming

- Changes the schema, not the instance
- Notation: $\rho_{B1,\ldots,Bn}(R)$
- Example:
  - $\rho_{LastName,\ SocSocNo}(Employee)$
  - Output schema:
    Answer(LastName, SocSocNo)

# Renaming Example

**Employee**

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

$$\rho_{LastName,\ SocSocNo}\ (\textbf{Employee})$$

| LastName | SocSocNo |
|----------|----------|
| John | 999999999 |
| Tony | 777777777 |

# Natural Join

- Notation: R1 $|\times|$ R2

- Meaning:  R1 $|\times|$ R2 = $\Pi_A(\sigma_C(R1 \times R2))$

- Where:
  - The selection $\sigma_C$ checks equality of all common attributes
  - The projection eliminates the duplicate common attributes

# Natural Join Example

## Employee

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

## Dependents

| SSN | Dname |
|-----|-------|
| 999999999 | Emily |
| 777777777 | Joe |

**Employee $\bowtie$ Dependents =**

$$\Pi_{\text{Name, SSN, Dname}}(\sigma_{\text{SSN=SSN2}}(\text{Employee x } \rho_{\text{SSN2, Dname}}(\text{Dependents})))$$

| Name | SSN | Dname |
|------|-----|-------|
| John | 999999999 | Emily |
| Tony | 777777777 | Joe |

# Natural Join

- R=

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

S=

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

- R |×| S=

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

# Natural Join

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

- Given R(A, B, C),  S(D, E), what is R $\bowtie$ S ?

- Given R(A, B),  S(A, B),  what is  R $\bowtie$ S ?

# Theta Join

- A join that involves a predicate
- R1 $|\!\times\!|_\theta$ R2  =  $\sigma_\theta$ (R1 $\times$ R2)
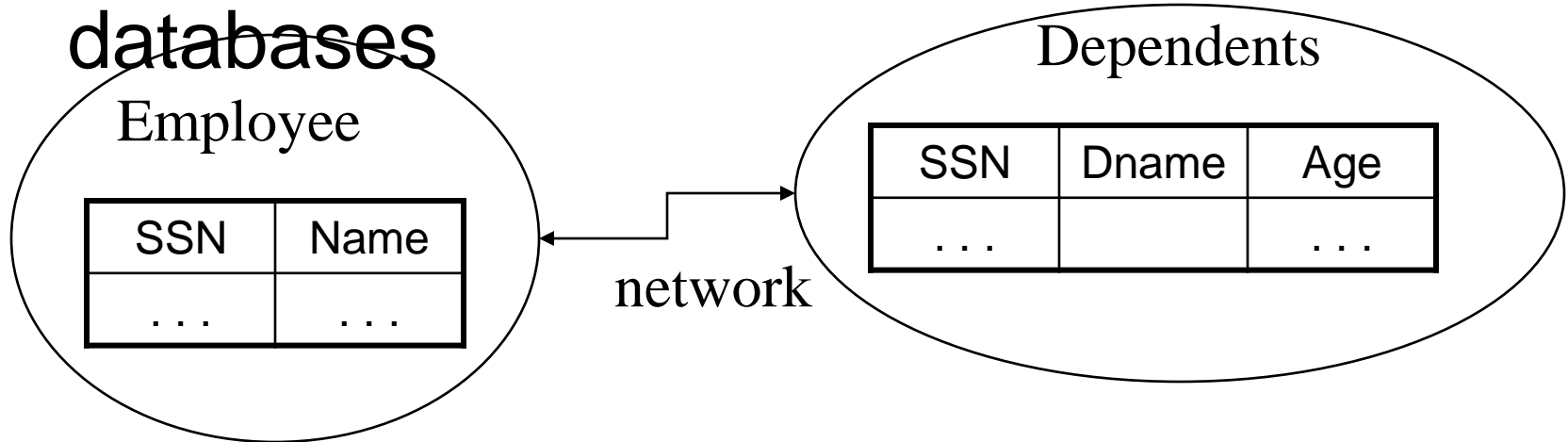- Here $\theta$ can be any condition

# Eq-join

- A theta join where $\theta$ is an equality
- R1 $|\times|_{A=B}$ R2 $= \sigma_{A=B}$ (R1 $\times$ R2)
- Example:
  - Employee $|\times|_{SSN=SSN}$ Dependents


- Most useful join in practice

# Semijoin

- $R \ltimes S = \Pi_{A1,\ldots,An} (R \bowtie S)$
- Where $A_1, \ldots, A_n$ are the attributes in R
- Example:
  - Employee $\ltimes$ Dependents

# Semijoins in Distributed Databases

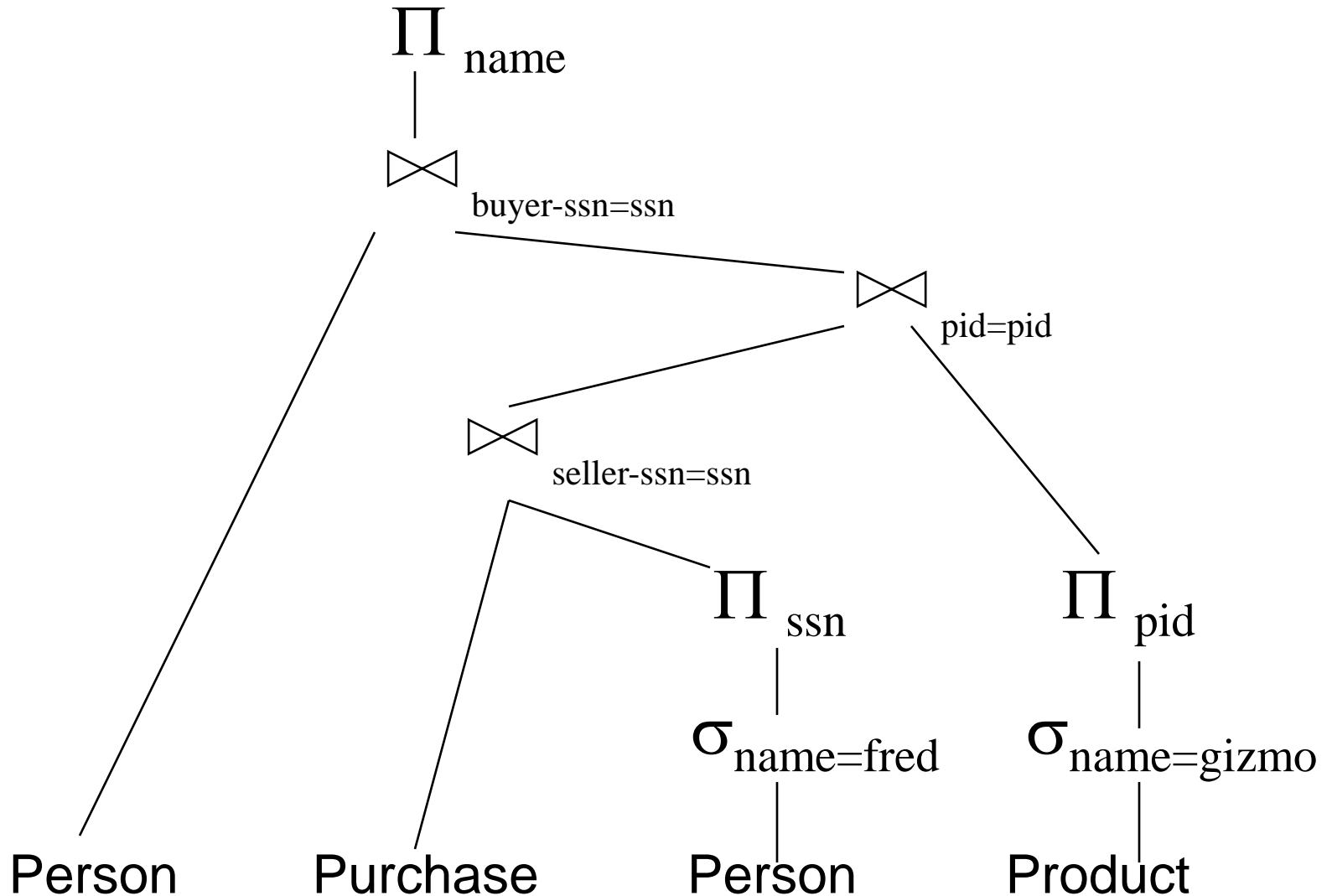- Semijoins are used in distributed databases

Employee

| SSN | Name |
|-----|------|
| . . . | . . . |

network

Dependents

| SSN | Dname | Age |
|-----|-------|-----|
| . . . | | . . . |

$$\text{Employee} \ltimes_{ssn=ssn} (\sigma_{age>71} (\text{Dependents}))$$

$T = \Pi_{SSN} \sigma_{age>71} (\text{Dependents})$

$R = \text{Employee} \ltimes T$

$\text{Answer} = R \bowtie \text{Dependents}$

# Complex RA Expressions

$\Pi_{\text{name}}$

$\bowtie_{\text{buyer-ssn=ssn}}$

$\bowtie_{\text{pid=pid}}$

$\bowtie_{\text{seller-ssn=ssn}}$

$\Pi_{\text{ssn}}$

$\Pi_{\text{pid}}$

$\sigma_{\text{name=fred}}$

$\sigma_{\text{name=gizmo}}$

Person          Purchase          Person          Product

# Operations on Bags

A **bag** = a set with repeated elements

All operations need to be defined carefully on bags

- $\{a,b,b,c\} \cup \{a,b,b,b,e,f,f\} = \{a,a,b,b,b,b,b,c,e,f,f\}$

- $\{a,b,b,b,c,c\} - \{b,c,c,c,d\} = \{a,b,b,d\}$

- $\sigma_C(R)$: preserve the number of occurrences

- $\Pi_A(R)$: no duplicate elimination

- Cartesian product, join: no duplicate elimination

Important ! Relational Engines work on bags, not sets !

# Finally: RA has Limitations !

- 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 C program