# Databases and Information Systems

## Relational Algebra

Dmitriy Traytel
slides partly by Marcos Vaz Salles

UNIVERSITY OF COPENHAGEN

# Do-It-Yourself Recap:
# Fill in the blanks

P

| x | y |
|---|---|
| 1 | 10 |
| 2 | 12 |
| 42 | 10 |

$$\{x \mapsto 42\} \quad \neg P(x,10)$$

$$\{x \mapsto 12\} \quad \neg P(x,10)$$

$$[\![\neg P(x,10)]\!] =$$

$$[\![P(x,10) \lor P(2,x)]\!] =$$

$$[\![\neg P(x,10) \lor P(2,x)]\!] =$$

$$[\![P(x,10) \lor P(2,y)]\!] =$$

# What should we learn today?

- Understand the notion of domain independence and be able to argue whether a given query is domain independent

- Understand the relational algebra normal form (RANF) and be able to determine whether a relational calculus query is in RANF

- Explain the operators of the relational algebra

- Formulate queries in the relational algebra

- Explain limitations of the relational algebra in terms of query expressiveness

- Formulate queries with the main extensions of the relational algebra including bag semantics, grouping, aggregation, sorting, and outer join

# Finite vs Infinite

- Fundamental problem with relational calculus:
  $[\![ \phi ]\!]$ is not always a finite relation

- Some examples for such "**unsafe**" queries
  - $\phi$ = Ships(n,cl,l) $\lor$ Outcomes(n,b,r)
  - $\phi$ = P(x) $\lor$ Q(y)
  - $\phi$ = $\neg$ P(x)
  - $\phi$ = x $\approx$ y

- But: query evaluation works with finite tables (why?)

# Domain Independence

- All tables (predicates) are finite

- Q: Where does the infiniteness of $[\![\phi]\!]$ come from?

# Domain Independence

- All tables (predicates) are finite

- Q: Where does the infiniteness of $[\![\phi]\!]$ come from?

- A: The domain $\mathbb{D}$

# Domain Independence

- All tables (predicates) are finite

- Q: Where does the infiniteness of $[\![\phi]\!]$ come from?

- A: The domain $\mathbb{D}$

- $\mathbb{D}$ can be seen as a parameter of query evaluation: $[\![\phi]\!]_{\mathbb{D}}$

# Domain Independence

- All tables (predicates) are finite

- Q: Where does the infiniteness of $[\![\phi]\!]$ come from?

- A: The domain $\mathbb{D}$

- $\mathbb{D}$ can be seen as a parameter of query evaluation: $[\![\phi]\!]_{\mathbb{D}}$

- $\phi$ is <u>domain-independent</u> if for all $\mathbb{D}, \mathbb{E}$: $[\![\phi]\!]_{\mathbb{D}} = [\![\phi]\!]_{\mathbb{E}}$

# Domain Independence

- All tables (predicates) are finite

- Q: Where does the infiniteness of $[\![\phi]\!]$ come from?

- A: The domain $\mathbb{D}$

- $\mathbb{D}$ can be seen as a parameter of query evaluation: $[\![\phi]\!]_{\mathbb{D}}$

- $\phi$ is domain-independent if for all $\mathbb{D}, \mathbb{E}$: $[\![\phi]\!]_{\mathbb{D}} = [\![\phi]\!]_{\mathbb{E}}$

- For example: $P(x) \wedge Q(y)$ is domain-independent: $[\![P(x) \wedge Q(y)]\!]_{\mathbb{D}} = DB(P) \times DB(Q)$

AND

# Unsafe $\implies$ Not Domain Independent

$$[\![\phi]\!]_{\mathbb{D}}$$

$\phi = P(x) \lor Q(y)$    {(x,y). x ∈ DB(P) and y ∈ $\mathbb{D}$ or x ∈ $\mathbb{D}$ and y ∈ DB(Q)}

$\phi = \neg\, P(x)$    {(x). x ∈ $\mathbb{D}$ and x ∉ DB(P)}

$\phi = x \approx y$    {(x,x). x ∈ $\mathbb{D}$}

# Safe $\implies$ Domain Independent?

- Here, safe means "finite query result"

- Implication does **not** hold!

# Safe $\implies$ Domain Independent?

- Here, safe means "finite query result"

- Implication does **not** hold!

- Any formula without free variables (= closed) has a finite query result

# Safe $\Longrightarrow$ Domain Independent?

- Here, safe means "finite query result"

- Implication does **not** hold!

- Any formula without free variables (= closed) has a finite query result

- $[\![\forall x.\ P(x)]\!]_{\mathbb{D}} = \{()\}$ if $\mathbb{D} = DB(P)$

# Safe $\Longrightarrow$ Domain Independent?

- Here, safe means "finite query result"

- Implication does **not** hold!

- Any formula without free variables (= closed) has a finite query result

- $[\![\forall x.\ P(x)]\!]_{\mathbb{D}} = \{()\}$ if $\mathbb{D} = DB(P)$

- $[\![\forall x.\ P(x)]\!]_{\mathbb{D}} = \{\}$ if $DB(P) \subsetneq \mathbb{D}$

# When is a formula domain-independent?

- Undecidable problem
  (= there exists no algorithm that answers the above question precisely)

- Resort to **syntactic** overapproximations:
  - under easy-to-check conditions a formula is domain-independent
    - e.g., the formula is a conjunction of n predicates (conjunctive queries)
  - conditions not met $\implies$ the formula may or may not be domain-independent

# Relational Algebra Normal Form

- A particular syntactic overapproximation

- RANF $\implies$ domain-independent

- Even better: RANF $\implies$ each "subformula" evaluates to a finite relation

- Has something to do with Relational Algebra (coming soon)

# RANF (continued)

```
ranf(P(t1, …, tn))      ⟺  true
ranf(t1 ≈ t2)           ⟺  false
ranf(¬ φ)               ⟺  fv(φ)={} and ranf(φ)
ranf(φ ∨ ψ)             ⟺  ranf(φ) and ranf(ψ) and fv(φ)=fv(ψ)
ranf(φ ∧ ψ)             ⟺  …
ranf(∀x. φ)             ⟺  false
ranf(∃x. φ)             ⟺  ranf(φ)
```

# RANF (continued)

ranf(φ) and ranf(ψ) or

ranf...
ranf...
ranf(¬φ)            ⟺  ...(φ) ... and ranf(φ)
ranf(φ ∨ ψ)         ⟺  ...anf(φ) and ranf(ψ) and fv(φ)=fv(ψ)
ranf(φ ∧ ψ)         ⟺  ...
ranf(∀x. φ)         ⟺  false
ranf(∃x. φ)         ⟺  ranf(φ)

# RANF (continued)

ranf(φ) and ranf(ψ) or
ranf(φ) and ψ=¬χ and ranf(χ) and fv(χ)⊆fv(φ) or

ranf($\quad$)
ranf($\quad$)
ranf($\quad$φ)
ranf(φ ∨ ψ) $\iff$ ranf(φ) and ranf(ψ) and fv(φ)=fv(ψ)
ranf(φ ∧ ψ) $\iff$ ...
ranf(∀x. φ) $\iff$ false
ranf(∃x. φ) $\iff$ ranf(φ)

# RANF (continued)

ranf(φ) and ranf(ψ) or
ranf(φ) and ψ=¬χ and ranf(χ) and fv(χ)⊆fv(φ) or
ranf(φ) and ψ=t1 ≈ t2 and (fv(t1)⊆fv(φ) or fv(t2)⊆fv(φ)) or

ranf($\varphi \lor \psi$)        ⟺  ranf(φ) and ranf(ψ) and fv(φ)=fv(ψ)
ranf($\varphi \land \psi$)        ⟺  ...
ranf($\forall x.\ \varphi$)        ⟺  false
ranf($\exists x.\ \varphi$)        ⟺  ranf(φ)

# RANF (continued)

ranf(φ)  and  ranf(ψ)  or
ranf(φ)  and  ψ=¬χ  and  ranf(χ)  and  fv(χ)⊆fv(φ)  or
ranf(φ)  and  ψ=t1 ≈ t2  and  (fv(t1)⊆fv(φ)  or  fv(t2)⊆fv(φ))  or
ranf(φ)  and  ψ=¬t1 ≈ t2  and  fv(t1)⊆fv(φ)  and  fv(t2)⊆fv(φ)

ranf(φ ∨ ψ)            ⟺ ranf(φ)  and  ranf(ψ)  and  fv(φ)=fv(ψ)
ranf(φ ∧ ψ)            ⟺ …
ranf(∀x. φ)            ⟺ false
ranf(∃x. φ)            ⟺ ranf(φ)

# Quiz: RANF or Not?

$P(x) \lor Q(y)$

$P(x) \land Q(y)$

$P(x,y) \land Q(y,z)$

$P(x,y) \land \neg Q(y,z)$

$P(x,y,z) \land \neg Q(y,z)$

$P(x,y,z) \land \neg y \approx z$

$P(x,y,z) \land \neg (Q(y,z) \lor R(x,z))$

$P(x,y,z) \land \neg Q(y,z) \land \neg R(x,z)$

# Quiz: RANF or Not?

P(x) ∨ Q(y)    ✘

P(x) ∧ Q(y)

P(x,y) ∧ Q(y,z)

P(x,y) ∧ ¬ Q(y,z)

P(x,y,z) ∧ ¬ Q(y,z)

P(x,y,z) ∧ ¬ y ≈ z

P(x,y,z) ∧ ¬ (Q(y,z) ∨ R(x,z))

P(x,y,z) ∧ ¬ Q(y,z) ∧ ¬ R(x,z)

# Quiz: RANF or Not?

$P(x) \lor Q(y)$ ✖

$P(x) \land Q(y)$ ✔

$P(x,y) \land Q(y,z)$

$P(x,y) \land \lnot Q(y,z)$

$P(x,y,z) \land \lnot Q(y,z)$

$P(x,y,z) \land \lnot y \approx z$

$P(x,y,z) \land \lnot (Q(y,z) \lor R(x,z))$

$P(x,y,z) \land \lnot Q(y,z) \land \lnot R(x,z)$

# Quiz: RANF or Not?

P(x) ∨ Q(y)    ✘

P(x) ∧ Q(y)    ✔

P(x,y) ∧ Q(y,z)    ✔

P(x,y) ∧ ¬ Q(y,z)

P(x,y,z) ∧ ¬ Q(y,z)

P(x,y,z) ∧ ¬ y ≈ z

P(x,y,z) ∧ ¬ (Q(y,z) ∨ R(x,z))

P(x,y,z) ∧ ¬ Q(y,z) ∧ ¬ R(x,z)

# Quiz: RANF or Not?

$P(x) \vee Q(y)$ ❌

$P(x) \wedge Q(y)$ ✅

$P(x,y) \wedge Q(y,z)$ ✅

$P(x,y) \wedge \neg Q(y,z)$ ❌

$P(x,y,z) \wedge \neg Q(y,z)$

$P(x,y,z) \wedge \neg y \approx z$

$P(x,y,z) \wedge \neg (Q(y,z) \vee R(x,z))$

$P(x,y,z) \wedge \neg Q(y,z) \wedge \neg R(x,z)$

# Quiz: RANF or Not?

P(x) ∨ Q(y)          ✗

P(x) ∧ Q(y)          ✓

P(x,y) ∧ Q(y,z)      ✓

P(x,y) ∧ ¬ Q(y,z)    ✗

P(x,y,z) ∧ ¬ Q(y,z)  ✓

P(x,y,z) ∧ ¬ y ≈ z

P(x,y,z) ∧ ¬ (Q(y,z) ∨ R(x,z))

P(x,y,z) ∧ ¬ Q(y,z) ∧ ¬ R(x,z)

# Quiz: RANF or Not?

P(x) ∨ Q(y)    ✗

P(x) ∧ Q(y)    ✓

P(x,y) ∧ Q(y,z)    ✓

P(x,y) ∧ ¬ Q(y,z)    ✗

P(x,y,z) ∧ ¬ Q(y,z)    ✓

P(x,y,z) ∧ ¬ y ≈ z    ✓

P(x,y,z) ∧ ¬ (Q(y,z) ∨ R(x,z))

P(x,y,z) ∧ ¬ Q(y,z) ∧ ¬ R(x,z)

# Quiz: RANF or Not?

P(x) ∨ Q(y)     ✘

P(x) ∧ Q(y)     ✔

P(x,y) ∧ Q(y,z)     ✔

P(x,y) ∧ ¬ Q(y,z)     ✘

P(x,y,z) ∧ ¬ Q(y,z)     ✔

P(x,y,z) ∧ ¬ y ≈ z     ✔

P(x,y,z) ∧ ¬ (Q(y,z) ∨ R(x,z))     ✘

P(x,y,z) ∧ ¬ Q(y,z) ∧ ¬ R(x,z)

# Quiz: RANF or Not?

$P(x) \lor Q(y)$ ❌

$P(x) \land Q(y)$ ✅

$P(x,y) \land Q(y,z)$ ✅

$P(x,y) \land \lnot Q(y,z)$ ❌

$P(x,y,z) \land \lnot Q(y,z)$ ✅

$P(x,y,z) \land \lnot y \approx z$ ✅

$P(x,y,z) \land \lnot (Q(y,z) \lor R(x,z))$ ❌

$P(x,y,z) \land \lnot Q(y,z) \land \lnot R(x,z)$ ✅

# Competition Task: Smells Like Cheese Spirit

```
loves(who:string,what:string)
smells(what:string)
```

```
e.g. loves(Bill,cheese)
     loves(Bill,jackfruit)
     loves(Bill,tomato)
     loves(Elon,cheese)
     loves(Elon,fish)
     loves(Elon,jackfruit)
     loves(Jeff,cheese)
     loves(Jeff,Jeff)
     smells(cheese)
     smells(jackfruit)
     smells(fish)
```

# Competition Task: Smells Like Cheese Spirit

```
loves(who:string,what:string)
smells(what:string)
```

Compute the lovers of cheese
that love all things that smell.

```
e.g. loves(Bill,cheese)
     loves(Bill,jackfruit)
     loves(Bill,tomato)
     loves(Elon,cheese)
     loves(Elon,fish)
     loves(Elon,jackfruit)
     loves(Jeff,cheese)
     loves(Jeff,Jeff)
     smells(cheese)
     smells(jackfruit)
     smells(fish)
```

# Competition Task: Smells Like Cheese Spirit

loves(who:string,what:string)
smells(what:string)

Compute the lovers of cheese
that love all things that smell.

The most efficient RC query in RANF wins.

e.g. loves(Bill,cheese)
loves(Bill,jackfruit)
loves(Bill,tomato)
loves(Elon,cheese)
loves(Elon,fish)
loves(Elon,jackfruit)
loves(Jeff,cheese)
loves(Jeff,Jeff)
smells(cheese)
smells(jackfruit)
smells(fish)

# Competition Task: Smells Like Cheese Spirit

```
loves(who:string,what:string)
smells(what:string)
```

Compute the lovers of cheese
that love all things that smell.

The most efficient RC query in RANF wins.

**Hint #1**: any feature supported by RC-eval is allowed

```
e.g. loves(Bill,cheese)
     loves(Bill,jackfruit)
     loves(Bill,tomato)
     loves(Elon,cheese)
     loves(Elon,fish)
     loves(Elon,jackfruit)
     loves(Jeff,cheese)
     loves(Jeff,Jeff)
     smells(cheese)
     smells(jackfruit)
     smells(fish)
```

# Competition Task: Smells Like Cheese Spirit

```
loves(who:string,what:string)
smells(what:string)
```

Compute the lovers of cheese
that love all things that smell.

The most efficient RC query in RANF wins.

**Hint #1**: any feature supported by RC-eval is allowed
**Hint #2**: RC-eval uses VeriMon to evaluate RANF queries

https://doi.org/10.1007/978-3-031-17715-6_1

```
e.g. loves(Bill,cheese)
     loves(Bill,jackfruit)
     loves(Bill,tomato)
     loves(Elon,cheese)
     loves(Elon,fish)
     loves(Elon,jackfruit)
     loves(Jeff,cheese)
     loves(Jeff,Jeff)
     smells(cheese)
     smells(jackfruit)
     smells(fish)
```

# Codd's Theorem

For every domain-independent relational calculus query there exists an equivalent query in RANF

# RANF $\implies$ each "subformula" evaluates to a <u>finite relation</u>

$$(P(x,y,z) \land \lnot Q(y,z)) \land \lnot R(x,z)$$

# RANF $\implies$ each "subformula" evaluates to a finite relation

$$\wedge \neg$$
$$\wedge \neg$$
$$P(x,y,z) \qquad Q(y,z) \qquad R(x,z)$$

# RANF $\implies$ each "subformula" evaluates to a finite relation



P(x,y,z)

| x | y | z |
|---|---|---|
| 1 | 10 | 4 |
| 2 | 12 | 4 |
| 42 | 10 | 5 |

Q(y,z)

| y | z |
|---|---|
| 10 | 4 |
| 13 | 4 |

R(x,z)

14

# RANF $\implies$ each "subformula" evaluates to a finite relation

| x | y | z |
|---|---|---|
| 2 | 12 | 4 |
| 42 | 10 | 5 |

$\wedge \neg$

$\wedge \neg$

P(x,y,z)        Q(y,z)        R(x,z)

| x | y | z |
|---|---|---|
| 1 | 10 | 4 |
| 2 | 12 | 4 |
| 42 | 10 | 5 |

| y | z |
|---|---|
| 10 | 4 |
| 13 | 4 |

# RANF $\implies$ each "subformula" evaluates to a finite relation

| x | y | z |
|---|---|---|
| 2 | 12 | 4 |
| 42 | 10 | 5 |

$\wedge \neg$

$\wedge \neg$

P(x,y,z)          Q(y,z)          R(x,z)

| x | y | z |
|---|---|---|
| 1 | 10 | 4 |
| 2 | 12 | 4 |
| 42 | 10 | 5 |

| y | z |
|---|---|
| 10 | 4 |
| 13 | 4 |

| x | z |
|---|---|
| 12 | 4 |
| 42 | 5 |

# RANF $\implies$ each "subformula" evaluates to a finite relation

| x | y | z |
|---|---|---|
| 2 | 12 | 4 |

| x | y | z |
|---|---|---|
| 2 | 12 | 4 |
| 42 | 10 | 5 |

$\wedge \neg$

$\wedge \neg$

P(x,y,z)

| x | y | z |
|---|---|---|
| 1 | 10 | 4 |
| 2 | 12 | 4 |
| 42 | 10 | 5 |

Q(y,z)

| y | z |
|---|---|
| 10 | 4 |
| 13 | 4 |

R(x,z)

| x | z |
|---|---|
| 12 | 4 |
| 42 | 5 |

14

# Relational Algebra

Give names to the supported table operations
Use these as basic operators

| x | y | z |
|---|---|---|
| 1 | 10 | 4 |
| 2 | 12 | 4 |
| 42 | 10 | 5 |

$\wedge \neg$

| y | z |
|---|---|
| 10 | 4 |
| 13 | 4 |

$=$

# Relational Algebra

Give names to the supported table operations
Use these as basic operators

| x | y | z |
|----|----|---|
| 1 | 10 | 4 |
| 2 | 12 | 4 |
| 42 | 10 | 5 |

$\wedge \neg$

| y | z |
|----|---|
| 10 | 4 |
| 13 | 4 |

$=$

| x | y | z |
|----|----|---|
| 2 | 12 | 4 |
| 42 | 10 | 5 |

# Relational Algebra

| x | z |
|---|---|
| 12 | 4 |
| 42 | 5 |

$\wedge$

| y | z |
|---|---|
| 10 | 4 |
| 13 | 4 |

=

Give names to the supported table operations
Use these as basic operators

| x | y | z |
|---|---|---|
| 1 | 10 | 4 |
| 2 | 12 | 4 |
| 42 | 10 | 5 |

$\wedge \neg$

| y | z |
|---|---|
| 10 | 4 |
| 13 | 4 |

=

| x | y | z |
|---|---|---|
| 2 | 12 | 4 |
| 42 | 10 | 5 |

# Relational Algebra

| x | z |
|---|---|
| 12 | 4 |
| 42 | 5 |

inner join (42, 5, not printed)

$\wedge$

| y | z |
|---|---|
| 10 | 4 |
| 13 | 4 |

$=$

| x | y | z |
|---|---|---|
| 12 | 10 | 4 |
| 12 | 13 | 4 |

## Give names to the supported table operations
## Use these as basic operators

| x | y | z |
|---|---|---|
| 1 | 10 | 4 |
| 2 | 12 | 4 |
| 42 | 10 | 5 |

$\wedge \neg$

| y | z |
|---|---|
| 10 | 4 |
| 13 | 4 |

$=$

| x | y | z |
|---|---|---|
| 2 | 12 | 4 |
| 42 | 10 | 5 |

# Relational Algebra

- Basic operations:
  - Selection $\sigma$ Selects a subset of rows from relation
  - Projection $\pi$   Deletes unwanted columns from relation
  - Cross-product $\times$   Allows us to combine two relations
  - Set-difference $-$   Tuples in relation 1, but not in relation 2
  - Union $\cup$ Tuples in relation 1 and in relation 2


- Additional operations:
  - Intersection $\cap$, join $\bowtie$, antijoin $\triangleright$, division $\div$, renaming $\rho$:  Not essential, but (very!) useful.


- Each operation returns a finite relation, i.e., operations can be composed! (Algebra is "closed".)

# Example Instances

| sid | bid | day |
|-----|-----|-----|
| 22 | 101 | 10.10 |
| 58 | 103 | 11.12 |

R1

| sid | sname | rating | age |
|-----|-------|--------|-----|
| 22 | dustin | 7 | 45.0 |
| 31 | lubber | 8 | 55.5 |
| 58 | rusty | 10 | 35.0 |

S1

- "Sailors" and "Reserves" relations for our examples.
- In relational calculus: formula's variables = column names
- In relational algebra: column names (= attributes) fixed
- Names in results "inherited" from names in input relations.

| sid | sname | rating | age |
|-----|-------|--------|-----|
| 28 | yuppy | 9 | 35.0 |
| 31 | lubber | 8 | 55.5 |
| 44 | guppy | 5 | 35.0 |
| 58 | rusty | 10 | 35.0 |

S2

# Projection

| sid | sname | rating | age |
|-----|-------|--------|------|
| 28 | yuppy | 9 | 35.0 |
| 31 | lubber | 8 | 55.5 |
| 44 | guppy | 5 | 35.0 |
| 58 | rusty | 10 | 35.0 |

S2

- Deletes attributes that are not in projection list.
- Schema of result contains exactly the fields in the projection list, with the same names that they had in the (only) input relation.
- Projection operator has to eliminate duplicates!  (Why??)
  - Note: real systems typically don't do duplicate elimination unless the user explicitly asks for it. (Why not?)

$\pi_{\text{sname,rating}}(S2)$

| sname | rating |
|-------|--------|
| yuppy | 9 |
| lubber | 8 |
| guppy | 5 |
| rusty | 10 |

$\pi_{\text{age}}(S2)$

| age |
|------|
| 35.0 |
| 55.5 |

# Projection

| sid | sname | rating | age |
|-----|-------|--------|------|
| 28 | yuppy | 9 | 35.0 |
| 31 | lubber | 8 | 55.5 |
| 44 | guppy | 5 | 35.0 |
| 58 | rusty | 10 | 35.0 |

S2

- Deletes attributes that are not in projection list.
- Schema of result contains exactly the fields in the projection list, with the same names that they had in the (only) input relation.
- Projection operator has to eliminate duplicates!  (Why??)
  - Note: real systems typically don't do duplicate elimination unless the user explicitly asks for it. (Why not?)

$\pi_{\text{sname,rating}}(\text{S2})$

$\exists \text{sid,age.}$
$\text{S2}_{(\text{sid,sname,rating,age})}$

| sname | rating |
|-------|--------|
| yuppy | 9 |
| lubber | 8 |
| guppy | 5 |
| rusty | 10 |

$\pi_{\text{age}}(\text{S2})$

| age |
|------|
| 35.0 |
| 55.5 |

Source: Ramakrishnan & Gehrke   18

# Projection

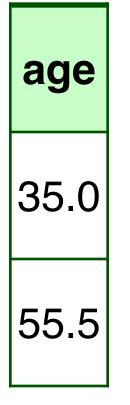| sid | sname | rating | age |
|-----|-------|--------|------|
| 28 | yuppy | 9 | 35.0 |
| 31 | lubber | 8 | 55.5 |
| 44 | guppy | 5 | 35.0 |
| 58 | rusty | 10 | 35.0 |

S2

- Deletes attributes that are not in projection list.
- Schema of result contains exactly the fields in the projection list, with the same names that they had in the (only) input relation.
- Projection operator has to eliminate duplicates!  (Why??)
  - Note: real systems typically don't do duplicate elimination unless the user explicitly asks for it. (Why not?)

$\pi_{\text{sname,rating}}(S2)$

$\exists$sid,age.
S2$_{\text{(sid,sname,rating,age)}}$

| sname | rating |
|-------|--------|
| yuppy | 9 |
| lubber | 8 |
| guppy | 5 |
| rusty | 10 |

$\pi_{\text{age}}(S2)$

$\exists$sid,sname,rating.
S2$_{\text{(sid,sname,rating,age)}}$

| age |
|------|
| 35.0 |
| 55.5 |

# Selection

| sid | sname | rating | age |
|-----|-------|--------|------|
| 28 | yuppy | 9 | 35.0 |
| 31 | lubber | 8 | 55.5 |
| 44 | guppy | 5 | 35.0 |
| 58 | rusty | 10 | 35.0 |

S2

- Selects rows that satisfy selection condition.

- No duplicates in result!  (Why?)

- Schema of result identical to schema of (only) input relation.

- Result relation can be the input for another relational algebra operation! (Operator composition.)

$\sigma_{\text{rating}>8}(S2)$

| sid | sname | rating | age |
|-----|-------|--------|------|
| 28 | yuppy | 9 | 35.0 |
| 58 | rusty | 10 | 35.0 |

$\pi_{\text{sname,rating}}(\sigma_{\text{rating}>8}(S2))$

| sname | rating |
|-------|--------|
| yuppy | 9 |
| rusty | 10 |

# Selection

| sid | sname | rating | age |
|-----|-------|--------|------|
| 28 | yuppy | 9 | 35.0 |
| 31 | lubber | 8 | 55.5 |
| 44 | guppy | 5 | 35.0 |
| 58 | rusty | 10 | 35.0 |

S2

- Selects rows that satisfy selection condition.
- No duplicates in result!  (Why?)
- Schema of result identical to schema of (only) input relation.
- Result relation can be the input for another relational algebra operation! (Operator composition.)

$\sigma_{\text{rating}>8}(S2)$

| sid | sname | rating | age |
|-----|-------|--------|------|
| 28 | yuppy | 9 | 35.0 |
| 58 | rusty | 10 | 35.0 |

$\pi_{\text{sname,rating}}(\sigma_{\text{rating}>8}(S2))$

| sname | rating |
|-------|--------|
| yuppy | 9 |
| rusty | 10 |

S2(sid,sname,rating,age)

$\wedge$ rating > 8

# Selection

| sid | sname | rating | age |
|-----|-------|--------|------|
| 28 | yuppy | 9 | 35.0 |
| 31 | lubber | 8 | 55.5 |
| 44 | guppy | 5 | 35.0 |
| 58 | rusty | 10 | 35.0 |

S2

- Selects rows that satisfy selection condition.

- No duplicates in result!  (Why?)

- Schema of result identical to schema of (only) input relation.

- Result relation can be the input for another relational algebra operation! (Operator composition.)

$\sigma_{\text{rating}>8}(\text{S2})$

| sid | sname | rating | age |
|-----|-------|--------|------|
| 28 | yuppy | 9 | 35.0 |
| 58 | rusty | 10 | 35.0 |

S2(sid,sname,rating,age)

$\wedge$ rating > 8

$\pi_{\text{sname,rating}}(\sigma_{\text{rating}>8}(\text{S2}))$

| sname | rating |
|-------|--------|
| yuppy | 9 |
| rusty | 10 |

$\exists$sid,age.

S2(sid,sname,rating,age)

$\wedge$ rating > 8

# Union, Intersection, Set-Difference

- These operations take two input relations that must be union-compatible:
- Same number of fields.
- "Corresponding" fields have the same type.
- What is the schema of result?

| sid | sname | rating | age |
|-----|-------|--------|------|
| 22 | dustin | 7 | 45.0 |
| 28 | yuppy | 9 | 35.0 |
| 31 | lubber | 8 | 55.5 |
| 44 | guppy | 5 | 35.0 |
| 58 | rusty | 10 | 35.0 |

S1 ∪ S2

| sid | sname | rating | age |
|-----|-------|--------|------|
| 31 | lubber | 8 | 55.5 |
| 58 | rusty | 10 | 35.0 |

S1 ∩ S2

| sid | sname | rating | age |
|-----|-------|--------|------|
| 22 | dustin | 7 | 45.0 |

S1 − S2

# Union, Intersection, Set-Difference

- These operations take two input relations that must be union-compatible:
- Same number of fields.
- "Corresponding" fields have the same type.
- What is the schema of result?

| sid | sname | rating | age |
|-----|-------|--------|------|
| 22 | dustin | 7 | 45.0 |
| 28 | yuppy | 9 | 35.0 |
| 31 | lubber | 8 | 55.5 |
| 44 | guppy | 5 | 35.0 |
| 58 | rusty | 10 | 35.0 |

S1 ∪ S2

| sid | sname | rating | age |
|-----|-------|--------|------|
| 31 | lubber | 8 | 55.5 |
| 58 | rusty | 10 | 35.0 |

S1 ∩ S2

| sid | sname | rating | age |
|-----|-------|--------|------|
| 22 | dustin | 7 | 45.0 |

S1 − S2

$S1_{(sid,sname,rating,age)} \lor S2_{(sid,sname,rating,age)}$

# Union, Intersection, Set-Difference

- These operations take two input relations that must be union-compatible:
- Same number of fields.
- "Corresponding" fields have the same type.
- What is the schema of result?

| sid | sname | rating | age |
|-----|-------|--------|------|
| 22 | dustin | 7 | 45.0 |
| 28 | yuppy | 9 | 35.0 |
| 31 | lubber | 8 | 55.5 |
| 44 | guppy | 5 | 35.0 |
| 58 | rusty | 10 | 35.0 |

| sid | sname | rating | age |
|-----|-------|--------|------|
| 31 | lubber | 8 | 55.5 |
| 58 | rusty | 10 | 35.0 |

| sid | sname | rating | age |
|-----|-------|--------|------|
| 22 | dustin | 7 | 45.0 |

$$S1_{(sid,sname,rating,age)} \wedge S2_{(sid,sname,rating,age)}$$

and - суворий

## S1 ∪ S2    S1 ∩ S2    S1 − S2

$$S1_{(sid,sname,rating,age)} \vee S2_{(sid,sname,rating,age)}$$

or

# Union, Intersection, Set-Difference

- These operations take two input relations that must be union-compatible:
- Same number of fields.
- "Corresponding" fields have the same type.
- What is the schema of result?

| sid | sname | rating | age |
|-----|-------|--------|------|
| 22 | dustin | 7 | 45.0 |
| 28 | yuppy | 9 | 35.0 |
| 31 | lubber | 8 | 55.5 |
| 44 | guppy | 5 | 35.0 |
| 58 | rusty | 10 | 35.0 |

| sid | sname | rating | age |
|-----|-------|--------|------|
| 31 | lubber | 8 | 55.5 |
| 58 | rusty | 10 | 35.0 |

| sid | sname | rating | age |
|-----|-------|--------|------|
| 22 | dustin | 7 | 45.0 |

$S1_{(sid,sname,rating,age)} \wedge \neg\, S2_{(sid,sname,rating,age)}$

$S1_{(sid,sname,rating,age)} \wedge S2_{(sid,sname,rating,age)}$

## S1 ∪ S2           ## S1 ∩ S2           ## S1 − S2

$S1_{(sid,sname,rating,age)} \vee S2_{(sid,sname,rating,age)}$

# Cross-Product

- Each row of S1 is paired with each row of R1.
- Result schema has one field per field of S1 and R1, with field names "inherited" if possible.
- Conflict: Both S1 and R1 have a field called sid.

R1

| sid | bid | day |
|-----|-----|-------|
| 22 | 101 | 10.10 |
| 58 | 103 | 11.12 |

S1

| sid | sname | rating | age |
|-----|-------|--------|------|
| 22 | dustin | 7 | 45.0 |
| 31 | lubber | 8 | 55.5 |
| 58 | rusty | 10 | 35.0 |

S1 × R1

| sid | sname | rating | age | sid | bid | day |
|-----|-------|--------|------|-----|-----|-------|
| 22 | dustin | 7 | 45.0 | 22 | 101 | 10.10 |
| 22 | dustin | 7 | 45.0 | 58 | 103 | 11.12 |
| 31 | lubber | 8 | 55.5 | 22 | 101 | 10.10 |
| 31 | lubber | 8 | 55.5 | 58 | 103 | 11.12 |
| 58 | rusty | 10 | 35.0 | 22 | 101 | 10.10 |
| 58 | rusty | 10 | 35.0 | 58 | 103 | 11.12 |

# Cross-Product

- Each row of S1 is paired with each row of R1.

- Result schema has one field per field of S1 and R1, with field names "inherited" if possible.

- Conflict: Both S1 and R1 have a field called sid.

R1

| sid | bid | day |
|-----|-----|-----|
| 22 | 101 | 10.10 |
| 58 | 103 | 11.12 |

S1

| sid | sname | rating | age |
|-----|-------|--------|-----|
| 22 | dustin | 7 | 45.0 |
| 31 | lubber | 8 | 55.5 |
| 58 | rusty | 10 | 35.0 |

S1 × R1

| sid | sname | rating | age | sid | bid | day |
|-----|-------|--------|-----|-----|-----|-----|
| 22 | dustin | 7 | 45.0 | 22 | 101 | 10.10 |
| 22 | dustin | 7 | 45.0 | 58 | 103 | 11.12 |
| 31 | lubber | 8 | 55.5 | 22 | 101 | 10.10 |
| 31 | lubber | 8 | 55.5 | 58 | 103 | 11.12 |
| 58 | rusty | 10 | 35.0 | 22 | 101 | 10.10 |
| 58 | rusty | 10 | 35.0 | 58 | 103 | 11.12 |

S1$_{(sid1,sname,rating,age)}$ $\wedge$ R1$_{(sid2,bid,day)}$

# Renaming

- The $\rho$ operator gives a new schema to a relation.

- $\rho_{\text{R1(A1,...,An)}}$(R2) makes R1 be a relation with attributes A1,…,An and the same tuples as R2.

- Simplified notation: R1(A1,…,An) := R2.

- In our previous example:  Res(sid1,sname,rating,age,sid2,bid,day) := S1 $\times$ R1

# Joins

- Condition Join:  $R \bowtie_C S = \sigma_C (R \times S)$

- $S1 \bowtie_{S1.sid<R1.sid} R1$
  where sid1 < sid2 (condition)

| sid | sname | rating | age | sid | bid | day |
|-----|-------|--------|------|-----|-----|-------|
| 22 | dustin | 7 | 45.0 | 58 | 103 | 11.12 |
| 31 | lubber | 8 | 55.5 | 58 | 103 | 11.12 |

- Result schema same as that of cross-product.

- Fewer tuples than cross-product,
  might be able to compute more efficiently

- Sometimes called a theta-join.

23

# Joins

- Condition Join: $R \bowtie_C S = \sigma_C (R \times S)$

- S1 $\bowtie_{S1.sid<R1.sid}$ R1

| sid | sname | rating | age | sid | bid | day |
|-----|-------|--------|------|-----|-----|-------|
| 22 | dustin | 7 | 45.0 | 58 | 103 | 11.12 |
| 31 | lubber | 8 | 55.5 | 58 | 103 | 11.12 |

$S1_{(sid1,sname,rating,age)} \wedge R1_{(sid2,bid,day)}$

$\wedge$ sid1 $<$ sid2

- Result schema same as that of cross-product.

- Fewer tuples than cross-product, might be able to compute more efficiently

- Sometimes called a theta-join.

# Joins

- <span style="color:red">Equijoin</span>:  A special case of condition join where the condition contains only equalities

- S1 $\bowtie_{sid}$ R1

| sid | sname | rating | age | bid | day |
|-----|-------|--------|------|-----|-------|
| 22 | dustin | 7 | 45.0 | 101 | 10.10 |
| 58 | rusty | 10 | 35.0 | 103 | 11.12 |

- Result schema similar to cross-product, but only one copy of fields for which equality is specified.

- Natural Join $\bowtie$:  Equijoin on all common fields.

# Joins

- Equijoin:  A special case of condition join where the condition contains only equalities

- S1 ⋈$_{sid}$ R1

| sid | sname | rating | age | bid | day |
|-----|-------|--------|-----|-----|-----|
| 22 | dustin | 7 | 45.0 | 101 | 10.10 |
| 58 | rusty | 10 | 35.0 | 103 | 11.12 |

S1$_{(sid,sname,rating,age)}$ ∧ R1$_{(sid,bid,day)}$

- Result schema similar to cross-product, but only one copy of fields for which equality is specified.

- Natural Join ⋈:  Equijoin on all common fields.

Source: Ramakrishnan & Gehrke

# Find names of sailors who've reserved boat #103

$\pi_{\text{sname}}(\sigma_{\text{bid}=103}(R1) \bowtie S1)$

$\text{Tmp1} := \sigma_{\text{bid}=103}(R1)$

$\text{Tmp2} := \text{Tmp1} \bowtie S1$

$\text{Tmp3} := \pi_{\text{sname}}(\text{Tmp2})$

$\pi_{\text{sname}}(\sigma_{\text{bid}=103}(R1 \bowtie S1))$

Different ways to state the same query

R1

| sid | bid | day |
|-----|-----|-------|
| 22 | 101 | 10.10 |
| 58 | 103 | 11.12 |

S1

| sid | sname | rating | age |
|-----|-------|--------|------|
| 22 | dustin | 7 | 45.0 |
| 31 | lubber | 8 | 55.5 |
| 58 | rusty | 10 | 35.0 |

# Relational Query Languages

- Query Languages **!=** programming languages!
  - QLs not expected to be "Turing complete".
  - QLs not intended to be used for complex calculations.
  - QLs support easy, efficient access to large data sets.

- Relational Algebra has limited expressibility
  - What queries are hard to answer?
  - What queries can't we answer?

# Hard, but possible

- Relational Schema:
  - Employees(<u>ssn</u>: integer; sal: real; mgr_ssn: integer).

- **Find the employees with the highest salary:**
  - How would you do it?

- Intuition:
  - Create "dominance" relation: e1 dominates e2 if e1 has salary higher than or equal to e2
  - Divide "dominance" relation by original table: Finds employees who dominate all others

- What about the second highest salary?

# Impossible, but very useful

- Relational Schema:
  - Employees(<u>ssn</u>: integer; sal: real; mgr_ssn: integer).
  - Works_In(<u>ssn</u>: integer; <u>did</u>: integer)

- Find the total amount paid in salaries by department

- Problem: **We do not know how to count!**

- What about finding employees that work in exactly three departments?

# Impossible, but very useful

- Relational Schema:
  - Employees(<u>ssn</u>: integer; sal: real; mgr_ssn: integer).
  - Works_In(<u>ssn</u>: integer; <u>did</u>: integer)

- Find all the superiors of employee with SSN 123-22-3666

- We want a **Transitive Closure**: Comes in handy when you query graphs

- **Problem: We can't write arbitrary loops/recursion!**
  - RA operators only implement implicit "foreach element in relation" loops

# Extensions to Relational Algebra

- Relational query languages restrict expressiveness to obtain ease of use and of optimization.

- There are some very useful queries we cannot express in the relational algebra.

- Many extensions proposed to handle those queries made into SQL.

<span style="color:red">We will study an extended relational algebra next!</span>

# Extensions to Relational Algebra

- Relational query languages restrict expressiveness to obtain ease of use and of optimization.

- There are some very useful queries we cannot express in the relational algebra.

- Many extensions proposed to handle those queries made into SQL.

<span style="color:red">We will study an extended relational algebra next!</span>

- Expressions in projections
- Bags vs sets
- Duplicate elimination
- Grouping/aggregation

- Sorting
- Outer joins
- ~~Recursion~~

# Extending Projection with Expressions

- Using the same $\pi_L$ operator, we allow the list L to contain arbitrary expressions involving attributes:
  - Arithmetic on attributes, e.g., A+B→C.
  - Duplicate occurrences of the same attribute.

$\pi_{rating+sid \rightarrow rs,age,age}(S1)$

| rs | age | age |
|----|-----|-----|
| 29 | 45.0 | 45.0 |
| 39 | 55.5 | 55.5 |
| 68 | 35.0 | 35.0 |

S1

| sid | sname | rating | age |
|-----|-------|--------|-----|
| 22 | dustin | 7 | 45.0 |
| 31 | lubber | 8 | 55.5 |
| 58 | rusty | 10 | 35.0 |

# Relational Algebra on Bags

- A bag (or multiset) is an unordered collection where duplicates are allowed
  - Like a set, but an element may appear more than once.

- Example: {1,2,1,3} is a bag.

- Example: {1,2,3} is a bag that happens to also be a set.

- SQL uses bag semantics

$\pi_{age}(S2)$

| age |
|-----|
| 35.0 |
| 55.5 |
| 35.0 |
| 35.0 |

$\sigma_{age<40.0}(\pi_{age}(S2))$

| age |
|-----|
| 35.0 |
| 35.0 |
| 35.0 |

S2

| sid | sname | rating | age |
|-----|-------|--------|------|
| 28 | yuppy | 9 | 35.0 |
| 31 | lubber | 8 | 55.5 |
| 44 | guppy | 5 | 35.0 |
| 58 | rusty | 10 | 35.0 |

# Bag Product

| age |
|-----|
| 35.0 |
| 35.0 |
| 35.0 |

$\sigma_{\text{age}<40.0}(\pi_{\text{age}}(\text{S2}))$

$\pi_{\text{sname}}(\sigma_{\text{sid}<40}(\text{S2}))$

| sname |
|-------|
| yuppy |
| lubber |

$\sigma_{\text{age}<40.0}(\pi_{\text{age}}(\text{S2})) \times \pi_{\text{sname}}(\sigma_{\text{sid}<40}(\text{S2}))$

| age | sname |
|-----|-------|
| 35.0 | yuppy |
| 35.0 | yuppy |
| 35.0 | yuppy |
| 35.0 | lubber |
| 35.0 | lubber |
| 35.0 | lubber |

S2

| sid | sname | rating | age |
|-----|-------|--------|-----|
| 28 | yuppy | 9 | 35.0 |
| 31 | lubber | 8 | 55.5 |
| 44 | guppy | 5 | 35.0 |
| 58 | rusty | 10 | 35.0 |

Source: Bulskov   33

# Bag Union, Intersection, and Difference

- An element appears in the union of two bags the sum of the number of times it appears in each bag.

  - Example: {1,2,1} ∪ {1,1,2,3,1} = {1,1,1,1,1,2,2,3}

- An element appears in the intersection of two bags the minimum of the number of times it appears in either.

  - Example: {1,2,1,1} ∩ {1,2,1,3} = {1,1,2}.

- An element appears in the difference A − B of bags as many times as it appears in A, minus the number of times it appears in B, but never less than 0 times.

  - Example: {1,2,1,1} − {1,2,3} = {1,1}.

# Beware: Bag Laws != Set Laws

- Some, but not all algebraic laws that hold for sets also hold for bags.

- Examples
  - The commutative law for union $R \cup S = S \cup R$ holds for bags, since addition is commutative

  - However, set union is idempotent, meaning that $S \cup S = S$.
  - For bags, if x appears n times in S, then it appears 2n times in $S \cup S$.
  - Thus $S \cup S \mathrel{!=} S$ in general, e.g., $\{1\} \cup \{1\} = \{1,1\} \mathrel{!=} \{1\}$.

# Duplicate Elimination

- R1 := $\delta$(R2).

- R1 consists of one copy of each tuple that appears in R2 one or more times.

$\pi_{\text{age}}(S2)$

| age |
|-----|
| 35.0 |
| 55.5 |
| 35.0 |
| 35.0 |

$\delta(\pi_{\text{age}}(S2))$

| age |
|-----|
| 35.0 |
| 55.5 |

# Aggregation Operators

- Aggregation operators are not operators of relational algebra.

- Rather, they apply to entire columns of a table and produce a single result.

- The most important examples: SUM, AVG, COUNT, MIN, and MAX.

SUM(rating) = 32
COUNT(sid) = 4
MAX(age) = 55.5
AVG(rating) = 8

S2

| sid | sname | rating | age |
|-----|-------|--------|------|
| 28 | yuppy | 9 | 35.0 |
| 31 | lubber | 8 | 55.5 |
| 44 | guppy | 5 | 35.0 |
| 58 | rusty | 10 | 35.0 |

# Grouping Operator

- $\gamma_L(R)$ where L is a list of elements that are either:
  - Individual (grouping) attributes.
  - AGG(A), where AGG is one of the aggregation operators and A is an attribute.
    - An arrow and a new attribute name renames the component.

- Semantics
  - Form one group for each distinct list of values in R for grouping attributes in list L.
  - Within each group, compute AGG(A) for each aggregation on list L.
  - Result has one tuple for each group:
    - The grouping attributes and
    - Their group's aggregations.

$\gamma_{age,COUNT(rating) \rightarrow cr}(S2)$

| cr | age |
|----|-----|
| 3 | 35.0 |
| 1 | 55.5 |

$\gamma_{age,MAX(rating) \rightarrow mr}(S2)$

| mr | age |
|----|-----|
| 10 | 35.0 |
| 8 | 55.5 |

S2

| sid | sname | rating | age |
|-----|-------|--------|-----|
| 28 | yuppy | 9 | 35.0 |
| 31 | lubber | 8 | 55.5 |
| 44 | guppy | 5 | 35.0 |
| 58 | rusty | 10 | 35.0 |

Source: Bulskov (partial)    38

# Sorting

- $\tau_L(R)$ is the list of tuples of R sorted first on the value of the first attribute on L, then on the second attribute of L, and so on.
  - Break ties arbitrarily.

- $\tau$ is relation itself, but with tuples sorted.
  - Semantics in SQL are murky: result should be a sequence; however, result is treated as a bag if fed to any other operators
  - Some operators can be made order-preserving

$\tau_{\text{rating}}(S2)$

| sid | sname | rating | age |
|-----|-------|--------|------|
| 44  | guppy | 5      | 35.0 |
| 31  | lubber| 8      | 55.5 |
| 28  | yuppy | 9      | 35.0 |
| 58  | rusty | 10     | 35.0 |

S2

| sid | sname | rating | age |
|-----|-------|--------|------|
| 28  | yuppy | 9      | 35.0 |
| 31  | lubber| 8      | 55.5 |
| 44  | guppy | 5      | 35.0 |
| 58  | rusty | 10     | 35.0 |

# Outer Join

- Suppose we join R ⋈ S.

- A tuple of R that has no tuple of S with which it joins is called dangling.
  - Similarly for a tuple of S.

- Outer join ⟗ preserves dangling tuples by padding them with ⊥ (NULL in SQL).

- Variants that preserve only left/right dangling tuples: ⟕, ⟖

$$S1 \overset{\circ}{\bowtie}_L R1 = S1 \overset{\circ}{\bowtie} R1 \qquad S1 \overset{\circ}{\bowtie}_R R1 = S1 \bowtie R1$$

| sid | bid | day |
|-----|-----|-----|
| 22  | 101 | 10.10 |
| 58  | 103 | 11.12 |

R1 (labels the above table)

| sid | sname | rating | age | bid | day |
|-----|-------|--------|-----|-----|-----|
| 22  | dustin | 7 | 45.0 | 101 | 10.10 |
| 58  | rusty  | 10 | 35.0 | 103 | 11.12 |
| 31  | lubber | 8 | 55.5 | ⊥ | ⊥ |

| sid | sname | rating | age |
|-----|-------|--------|-----|
| 22  | dustin | 7 | 45.0 |
| 31  | lubber | 8 | 55.5 |
| 58  | rusty  | 10 | 35.0 |

S1 (labels the above table)

# What should we learn today?

- Understand the notion of domain independence and be able to argue whether a given query is domain independent

- Understand the relational algebra normal form (RANF) and be able to determine whether a relational calculus query is in RANF

- Explain the operators of the relational algebra

- Formulate queries in the relational algebra

- Explain limitations of the relational algebra in terms of query expressiveness

- Formulate queries with the main extensions of the relational algebra including bag semantics, grouping, aggregation, sorting, and outer join