# Declarative Languages and Relational Algebra

CSE 4/562: Database Systems | Lecture 2

**DB. Sys.: T.C.B.:** 2.1-2.4, 5.1

What is different about these code snippets?

```sql
SELECT R.A, SUM(S.B)
FROM R, S
WHERE R.B = S.B
```

```python
for r in R:
  lookup[r.B] += [r.A]
for s in S:
  for a in lookup[s.B]:
    result[a] += s.C
```
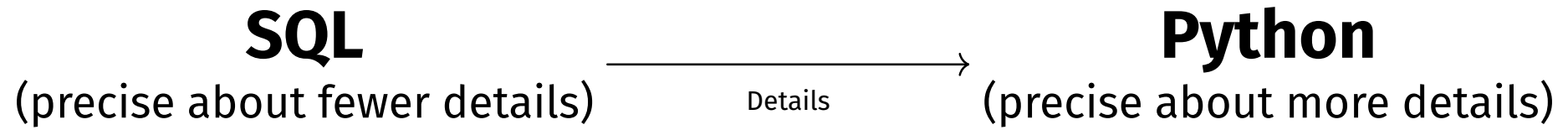
```
SELECT R.A, SUM(S.B)
FROM R, S
WHERE R.B = S.B
```

```
for r in R:
    lookup[r.B] += [r.A]
for s in S:
    for a in lookup[s.B]:
        result[a] += s.C
```

- The python code preprocesses R and scans S (and not visa versa)
- The python code explicitly creates a lookup table on R.
- The python code explicitly uses a dictionary for the lookup table.
- The python code explicitly encodes the result as a dictionary.
- SQL only states sources, constraints, and output format.

**SQL**
(precise about fewer details)

⟶ Details

**Python**
(precise about more details)

# SQL is a **Declarative Language**

## Postgresql

```
SELECT *
FROM ...
```

## Postgresql

```
SELECT *
FROM ...
```
$\longrightarrow$ Query

## Postgresql

```
SELECT *          ⎯⎯⎯⎯→  Query  ⎯⎯⎯⎯→
FROM ...
```

$Path_1,$

$...,$

$Path_n$

## Postgresql

SELECT *
FROM ...  $\longrightarrow$  Query  $\longrightarrow$  $Path_1,$ $...,$ $Path_n$  $\longrightarrow$  $Path_i$ + EvalState

## Postgresql

$$\text{SELECT *}\text{FROM ...} \longrightarrow \text{Query} \longrightarrow \text{Path}_1, ..., \text{Path}_n \longrightarrow \text{Path}_i + \text{EvalState} \longrightarrow \text{Results}$$

## Postgresql

```
SELECT *
FROM ...
```
$\longrightarrow$ Query $\longrightarrow$ $Path_1,$ $...,$ $Path_n$ $\longrightarrow$ $Path_i$ + EvalState $\longrightarrow$ Results

## Apache Spark

```
SELECT *
FROM ...
```

## Postgresql

```
SELECT *
FROM ...
```
$\longrightarrow$ Query $\longrightarrow$ $Path_1, ..., Path_n$ $\longrightarrow$ $Path_i$ + EvalState $\longrightarrow$ Results

## Apache Spark

```
SELECT *
FROM ...
```
$\longrightarrow$ Command

## Postgresql

```
SELECT *                    ⟶  Query  ⟶         Path₁,        ⟶         Pathᵢ         ⟶  Results
FROM ...                                         ...,                     +
                                                 Pathₙ                 EvalState
```

$$\text{SELECT * FROM ...} \longrightarrow \text{Query} \longrightarrow \begin{matrix} \text{Path}_1, \\ ..., \\ \text{Path}_n \end{matrix} \longrightarrow \begin{matrix} \text{Path}_i \\ + \\ \text{EvalState} \end{matrix} \longrightarrow \text{Results}$$

## Apache Spark

$$\text{SELECT * FROM ...} \longrightarrow \text{Command} \longrightarrow \text{LogicalPlan}$$

## Postgresql

```
SELECT *
FROM ...
```
$\longrightarrow$ Query $\longrightarrow$ $Path_1,$ ..., $Path_n$ $\longrightarrow$ $Path_i$ + EvalState $\longrightarrow$ Results

## Apache Spark

```
SELECT *
FROM ...
```
$\longrightarrow$ Command $\longrightarrow$ LogicalPlan

## Postgresql

$$\texttt{SELECT *} \atop \texttt{FROM ...}$$ $\longrightarrow$ Query $\longrightarrow$ $$\text{Path}_1, \atop ..., \atop \text{Path}_n$$ $\longrightarrow$ $$\text{Path}_i \atop + \atop \text{EvalState}$$ $\longrightarrow$ Results

## Apache Spark

$$\texttt{SELECT *} \atop \texttt{FROM ...}$$ $\longrightarrow$ Command $\longrightarrow$ LogicalPlan $\rightarrow$ PhysicalPlan

## Postgresql

SELECT *
FROM ... $\longrightarrow$ Query $\longrightarrow$ $Path_1,$ ..., $Path_n$ $\longrightarrow$ $Path_i$ + EvalState $\longrightarrow$ Results

## Apache Spark

SELECT *
FROM ... $\longrightarrow$ Command $\longrightarrow$ LogicalPlan $\rightarrow$ PhysicalPlan $\longrightarrow$ Results

1. **Parse the query**

2. **Make a Dumb™ plan**

3. **Come up with better plan ideas and pick one**

4. **Fill in the specific algorithms and data structures**

# How do we express a plan?

## SQL is Redundant

```sql
SELECT A, COUNT(DISTINCT B)
FROM R
WHERE C > 5
GROUP BY A
HAVING SUM(C) < 20
```

## SQL is Redundant

```sql
SELECT A, COUNT(DISTINCT B)
FROM R
WHERE C > 5
GROUP BY A
HAVING SUM(C) < 20
```

- Duplication in the language leads to duplication in the code.
- It's harder to reason about all possible combinations.

## SQL is Redundant

```
SELECT A, COUNT(DISTINCT B)
FROM R
WHERE C > 5
GROUP BY A
HAVING SUM(C) < 20
```

## SQL is Order-Independent

```
SELECT R.A, T.D
FROM R, S, T
WHERE R.B = S.B
  AND S.C = T.C
```

- Duplication in the language leads to duplication in the code.
- It's harder to reason about all possible combinations.

## SQL is Redundant

```
SELECT A, COUNT(DISTINCT B)
FROM R
WHERE C > 5
GROUP BY A
HAVING SUM(C) < 20
```

- Duplication in the language leads to duplication in the code.
- It's harder to reason about all possible combinations.

## SQL is Order-Independent

```
SELECT R.A, T.D
FROM R, S, T
WHERE R.B = S.B
  AND S.C = T.C
```

- Order of operations has a huge effect on performance.
- Reordering operations can create opportunities to "inline" pairs of operators.

## SQL is Composable

```sql
SELECT *
FROM R, S, (
    SELECT C, SUM(D) FROM T GROUP BY C
) T,
WHERE R.B = S.B AND S.C = T.C
```
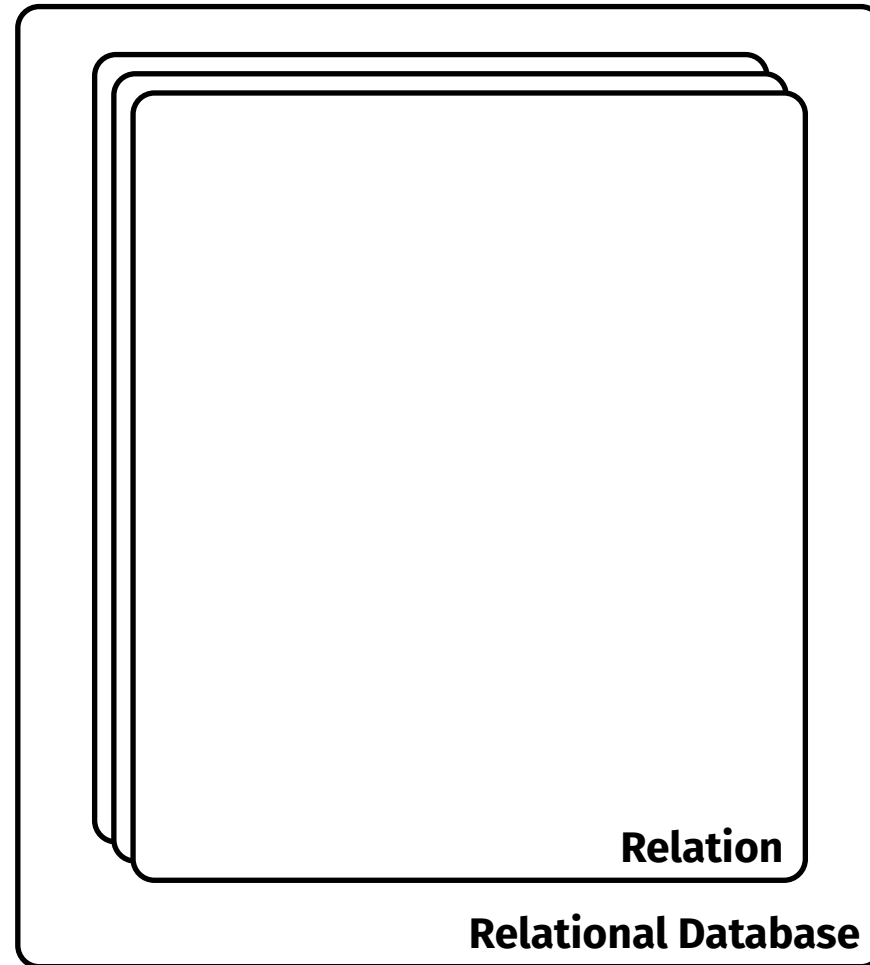
## SQL is Composable

```sql
SELECT *
FROM R, S, (
  SELECT C, SUM(D) FROM T GROUP BY C
) T,
WHERE R.B = S.B AND S.C = T.C
```
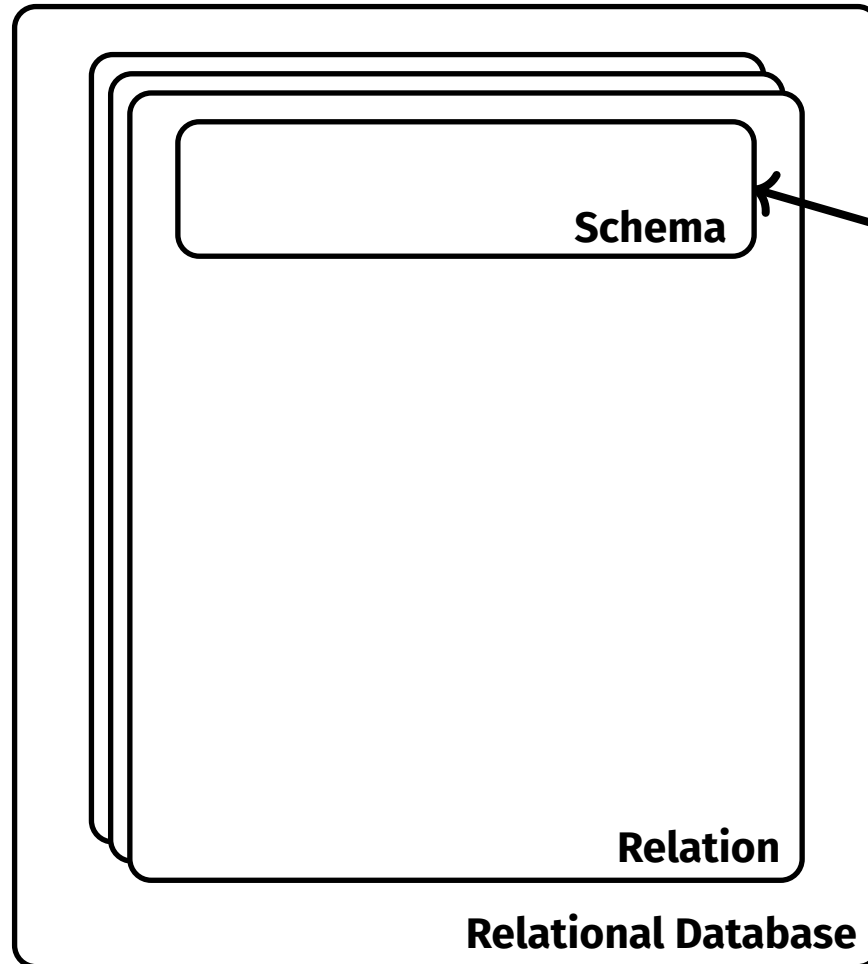
- Queries define relations: You can use a query in place of any relation.

1. The language should not be redundant.

2. The language should include order of operations

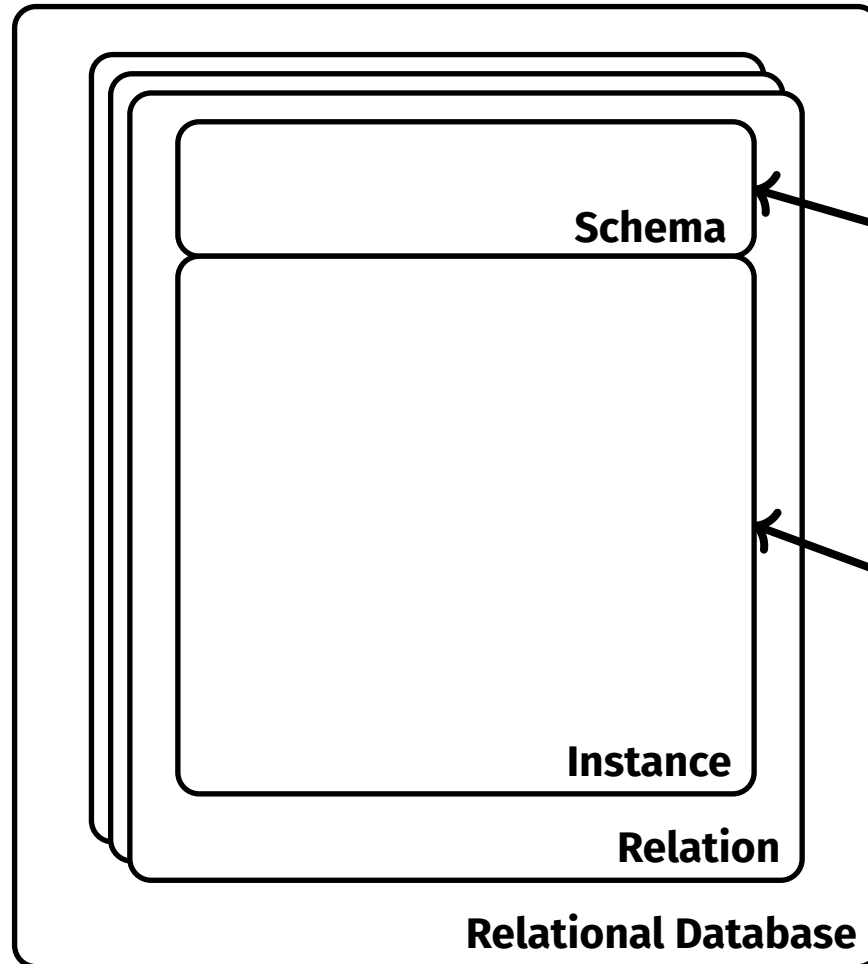3. The language should be made of composable building blocks.

**Relational Database**

Relation

Relational Database

**Specifies the name of the relation, name and type of each column, and any other constraints**

```
Officers(
    firstname string,
    lastname string,
    id int
)
```

Schema

Relation

Relational Database

**Specifies the name of the relation, name and type of each column, and any other constraints**

```
Officers(
    firstname string,
    lastname string,
    id int
)
```

**The Data**

```
[Jean Luc, Picard, 2360]
[Benjamin, Sisko, 2365]
[Kathryn, Janeway, 2370]
```

Schema

Instance

Relation

Relational Database

Schema

Instance

Relation

Relational Database

Columns (degree/arity)

**Specifies the name of the relation, name and type of each column, and any other constraints**

```
Officers(
    firstname string,
    lastname string,
    id int
)
```
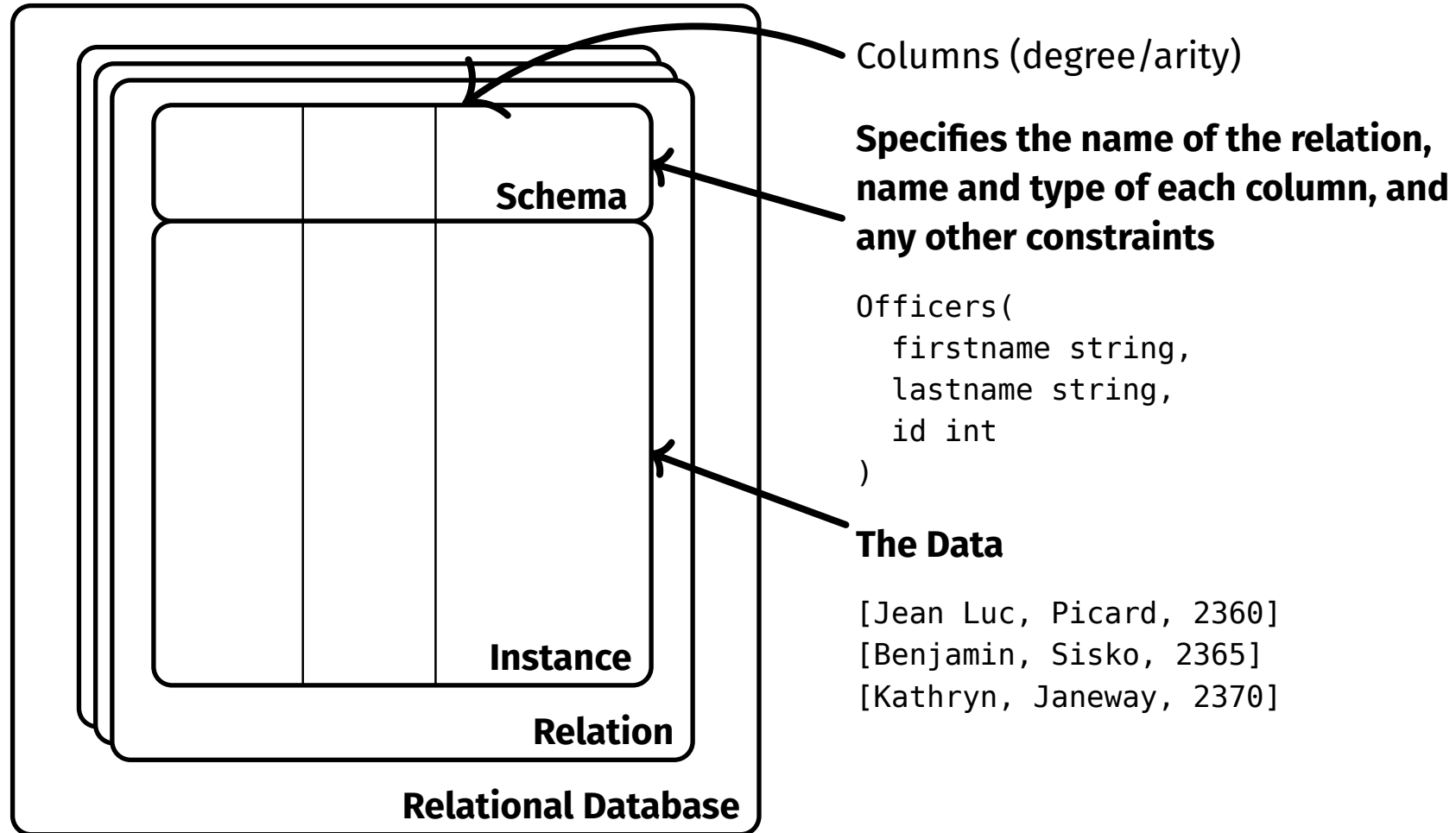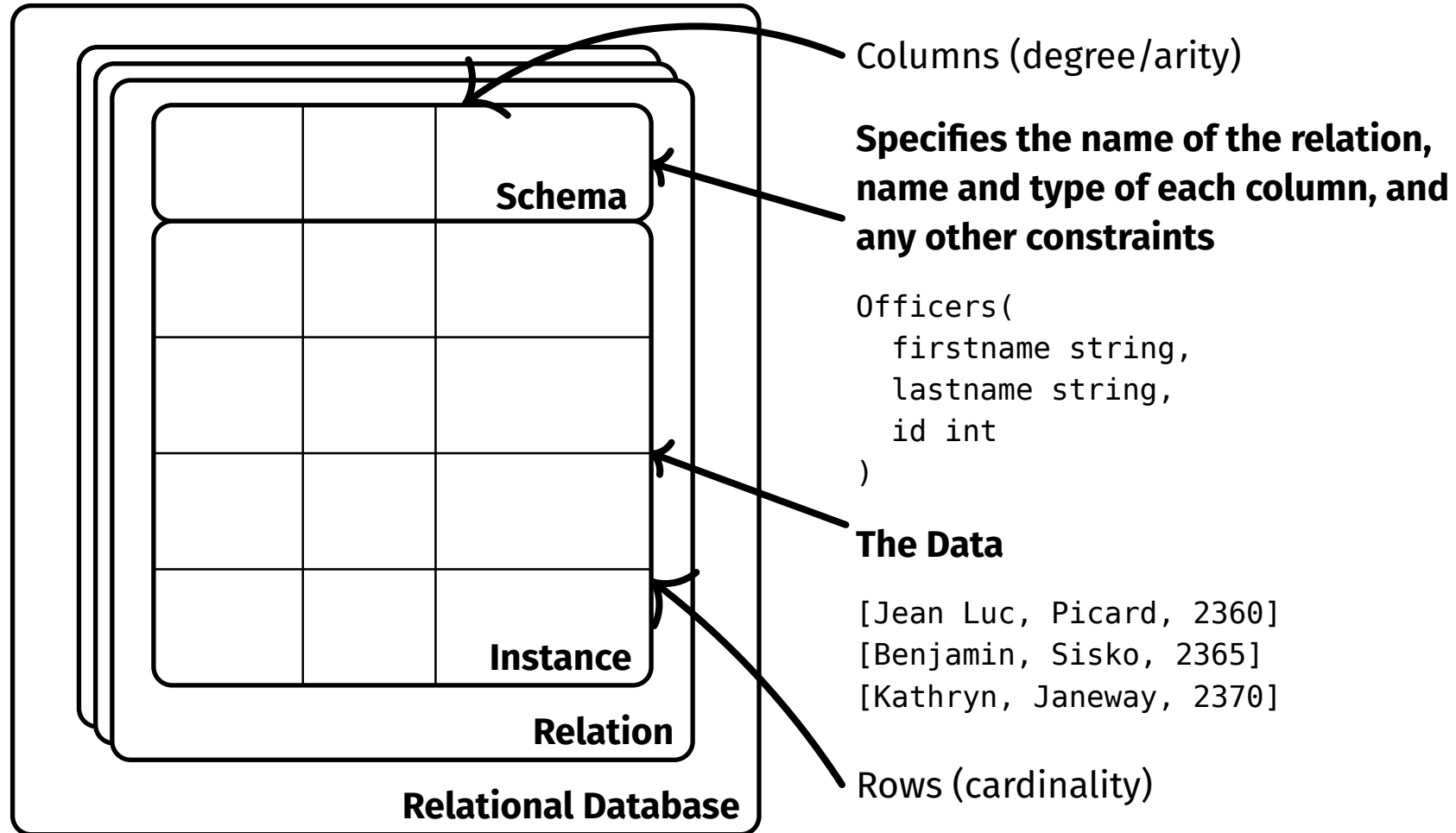
**The Data**

```
[Jean Luc, Picard, 2360]
[Benjamin, Sisko, 2365]
[Kathryn, Janeway, 2370]
```

Columns (degree/arity)

**Specifies the name of the relation, name and type of each column, and any other constraints**

```
Officers(
    firstname string,
    lastname string,
    id int
)
```

**The Data**

```
[Jean Luc, Picard, 2360]
[Benjamin, Sisko, 2365]
[Kathryn, Janeway, 2370]
```

Rows (cardinality)

Schema

Instance

Relation

Relational Database

**Everything is a relation**

- Q(R)

**Everything is a relation**

- Q(R)
- Q(Q(R))

## Everything is a relation

- Q(R)
- Q(Q(R))
- Q(Q(Q(R)))

**Everything is a relation**

- Q(R)
- Q(Q(R))
- Q(Q(Q(R)))
- Q Q moar.

## Everything is a relation

- Q(R)
- Q(Q(R))
- Q(Q(Q(R)))
- Q Q moar.

A query language with this property is **closed**.

## Everything is a relation

- Q(R)
- Q(Q(R))
- Q(Q(Q(R)))
- Q Q moar.

A query language with this property is **closed**.

## Simplifying Assumptions

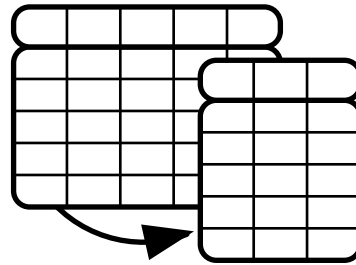- All attributes have **unique** names.
- Each instance is a Bag, Set or List
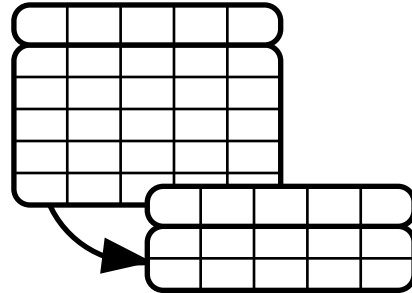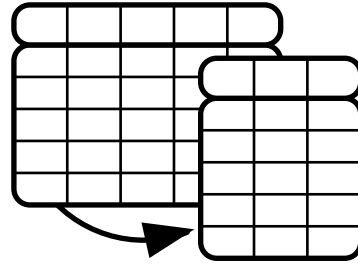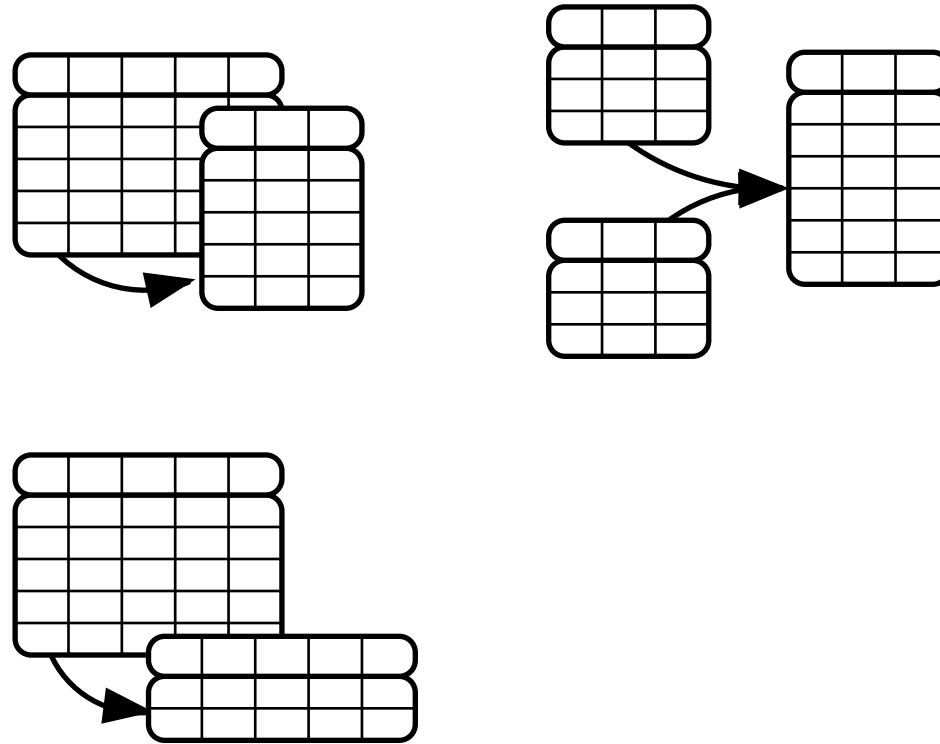
## Everything is a relation

- Q(R)
- Q(Q(R))
- Q(Q(Q(R)))
- Q Q moar.

A query language with this property is **closed**.

## Simplifying Assumptions

- All attributes have **unique** names.
- Each instance is a Bag, ~~Set, or List~~

|  | Filter | Merge | Derive |
| --- | --- | --- | --- |
| **Column** | | | |
| **Row** | | | |

| SQL Field | Operation |
|---|---|
| SELECT | |
| FROM | |
| JOIN | |
| WHERE | |
| GROUP, BY/Aggregate | |
| HAVING | |

| SQL Field | Operation |
|---|---|
| SELECT | Filter Column + Derive Column |
| FROM | |
| JOIN | |
| WHERE | |
| GROUP, BY/Aggregate | |
| HAVING | |

| SQL Field | Operation |
|---|---|
| SELECT | Filter Column + Derive Column |
| FROM | Merge Column |
| JOIN | |
| WHERE | |
| GROUP, BY/Aggregate | |
| HAVING | |

| SQL Field | Operation |
|---|---|
| SELECT | Filter Column + Derive Column |
| FROM | Merge Column |
| JOIN | Merge Column + Filter Row |
| WHERE | |
| GROUP, BY/Aggregate | |
| HAVING | |

| SQL Field | Operation |
|---|---|
| SELECT | Filter Column + Derive Column |
| FROM | Merge Column |
| JOIN | Merge Column + Filter Row |
| WHERE | Filter Row |
| GROUP, BY/Aggregate | |
| HAVING | |

| SQL Field | Operation |
| --- | --- |
| SELECT | Filter Column + Derive Column |
| FROM | Merge Column |
| JOIN | Merge Column + Filter Row |
| WHERE | Filter Row |
| GROUP, BY/Aggregate | Derive Row |
| HAVING | |

| SQL Field | Operation |
|---|---|
| SELECT | Filter Column + Derive Column |
| FROM | Merge Column |
| JOIN | Merge Column + Filter Row |
| WHERE | Filter Row |
| GROUP, BY/Aggregate | Derive Row |
| HAVING | Derive Row + Filter Row |

# Columns and rows are different

## Columns

• There is a "small" number (10,000s at most).

• Query planning knows everything about columns.

• Identified explicitly (by name or position)

## Rows

• There is a "large" number (Millions, Billions, More)

• Query planning has only statistics (if it has anything).

• Sets/Bags have no explicit identity (Each tuple's attributes identify it)[1]

---

[1]Technically false: Many DB systems have RowIDs... but these are not intended for the user

## Design Questions

- How do we indicate which columns/rows to keep/discard?

## Filtering Rows

- ???

## Filtering Columns

- ???

## Design Questions

- How do we indicate which columns/rows to keep/discard?

## Filtering Rows

- ???

## Filtering Columns

- Explicit list of columns to keep

## Design Questions

- How do we indicate which columns/rows to keep/discard?

## Filtering Rows

- Condition, or "Predicate" for which rows to keep

## Filtering Columns

- Explicit list of columns to keep

## Design Questions

- How do we indicate which columns/rows to keep/discard?

## Filtering Rows

- Condition, or "Predicate" for which rows to keep

## Filtering Columns

- Explicit list of columns to keep

## Language

| Code | Meaning |
|------|---------|
| `Filter($a = b$, In)` | Relation `In`, keeping only rows where $a = b$ |
| `Project([a, b, c], In)` | Relation `In`, keeping only columns a, b, c |

## Design Questions

- How do we indicate which columns/rows to keep/discard?

## Filtering Rows

- Condition, or "Predicate" for which rows to keep

## Filtering Columns

- Explicit list of columns to keep

## Language

| Code | Shorthand | Meaning |
|------|-----------|---------|
| `Filter($a = b$, In)` | $\sigma_{a=b}(\text{In})$ | Relation `In`, keeping only rows where $a = b$ |
| `Project([a, b, c], In)` | $\pi_{a,b,c}(\text{In})$ | Relation `In`, keeping only columns a, b, c |

## Design Questions

- Which rows/columns pair with which other rows/columns?

## Merging Rows

- ???

## Design Questions

• Which rows/columns pair with which other rows/columns?

## Merging Rows

• Pair columns of the same name

## Design Questions

- Which rows/columns pair with which other rows/columns?

## Merging Rows

- Pair columns of the same name
  - ‣ Relations must be "Union-compatible" (same schema) to be paired.
  - ‣ Some engines (e.g., Apache Spark) implicitly add columns of nulls.

## Design Questions

• Which rows/columns pair with which other rows/columns?

## Merging Rows

• Pair columns of the same name
  ‣ Relations must be "Union-compatible" (same schema) to be paired.
  ‣ Some engines (e.g., Apache Spark) implicitly add columns of nulls.

## Language

| Code | Meaning |
|---|---|
| `Union(In1, In2)` | All rows from both relations $In_1$ and $In_2$ |

## Design Questions

- Which rows/columns pair with which other rows/columns?

## Merging Rows

- Pair columns of the same name
  - ‣ Relations must be "Union-compatible" (same schema) to be paired.
  - ‣ Some engines (e.g., Apache Spark) implicitly add columns of nulls.

## Language

| Code | Shorthand | Meaning |
|------|-----------|---------|
| `Union(In1, In2)` | $In_1 \cup In_2$ | All rows from both relations $In_1$ and $In_2$ |

## Design Questions

- Which rows/columns pair with which other rows/columns?

## Merging Columns

- ???

## Design Questions

- Which rows/columns pair with which other rows/columns?

## Merging Columns

- Pair all rows

## Design Questions

- Which rows/columns pair with which other rows/columns?

## Merging Columns

- Pair all rows
  - ‣ Use filter to keep only the rows you want

## Design Questions

- Which rows/columns pair with which other rows/columns?

## Merging Columns

- Pair all rows
  - ‣ Use filter to keep only the rows you want

## Language

| Code | Meaning |
|------|---------|
| `Product(In1, In2)` | Every possible pair of rows from $\text{In}_1, \text{In}_2$ |

## Design Questions

- Which rows/columns pair with which other rows/columns?

## Merging Columns

- Pair all rows
  - ▸ Use filter to keep only the rows you want

## Language

| Code | Shorthand | Meaning |
| --- | --- | --- |
| `Product(In1, In2)` | $In_1 \times In_2$ | Every possible pair of rows from $In_1, In_2$ |

## Design Questions

- Which rows/columns pair with which other rows/columns?

## Merging Columns

- Pair all rows
  - ‣ Use filter to keep only the rows you want

## Language

| Code | Shorthand | Meaning |
|------|-----------|---------|
| `Product(In1, In2)` | $\text{In}_1 \times \text{In}_2$ | Every possible pair of rows from $\text{In}_1, \text{In}_2$ |
| `Join($a = b$, In1, In2)` | $\text{In}_1 \underset{a=b}{\bowtie} \text{In}_2$ | Shorthand for $\sigma_{a=b}(\text{In}_1 \times \text{In}_2)$ |

## Cartesian Product

- $R \times S$: Every pair of one tuple from $R$ and $S$

## Join

- $R \underset{a>b}{\bowtie} S$: Only include pairs where the predicate $a > b$ is true

## Equi-Join

- $R \underset{a=b}{\bowtie} S$: A join that only uses equality predicates (e.g., $a = b \wedge c = d$)
- $R \underset{a}{\bowtie} S$: If the equi-join columns have the same name, we just write the name(s)

## Natural Join

- $R \bowtie S$: If the join predicate is omitted, assume an equi-join between all columns with the same name.

## Design Questions

- How do we specify which columns/rows to combine?
- How do we specify how to define the new row/column?

## Deriving Columns

- ???

## Design Questions

- How do we specify which columns/rows to combine?
- How do we specify how to define the new row/column?

## Deriving Columns

- `A + B AS C` or `C = A + B`

## Design Questions

- How do we specify which columns/rows to combine?
- How do we specify how to define the new row/column?

## Deriving Columns

- `A + B AS C` or `C = A + B`

## Language

| Code | Meaning |
|------|---------|
| `Project([A, B, C = $A+B$], In)` | As Project, but derive C |

## Design Questions

- How do we specify which columns/rows to combine?
- How do we specify how to define the new row/column?

## Deriving Columns

- `A + B AS C` or `C = A + B`

## Language

| Code | Shorthand | Meaning |
| --- | --- | --- |
| `Project([A, B, C = $A+B$], In)` | $\pi_{A,B,C=A+B}(\text{In})$ | As Project, but derive C |

## Design Questions

- How do we specify which columns/rows to combine?
- How do we specify how to define the new row/column?

## Deriving Rows

- ???
- ???

## Design Questions

- How do we specify which columns/rows to combine?
- How do we specify how to define the new row/column?

## Deriving Rows

- "Group By" attributes[2]
- ???

---

[2]We'll eventually talk about other grouping strategies (e.g., Window Functions)

## Design Questions

- How do we specify which columns/rows to combine?
- How do we specify how to define the new row/column?

## Deriving Rows

- "Group By" attributes[3]
- "Aggregate functions"
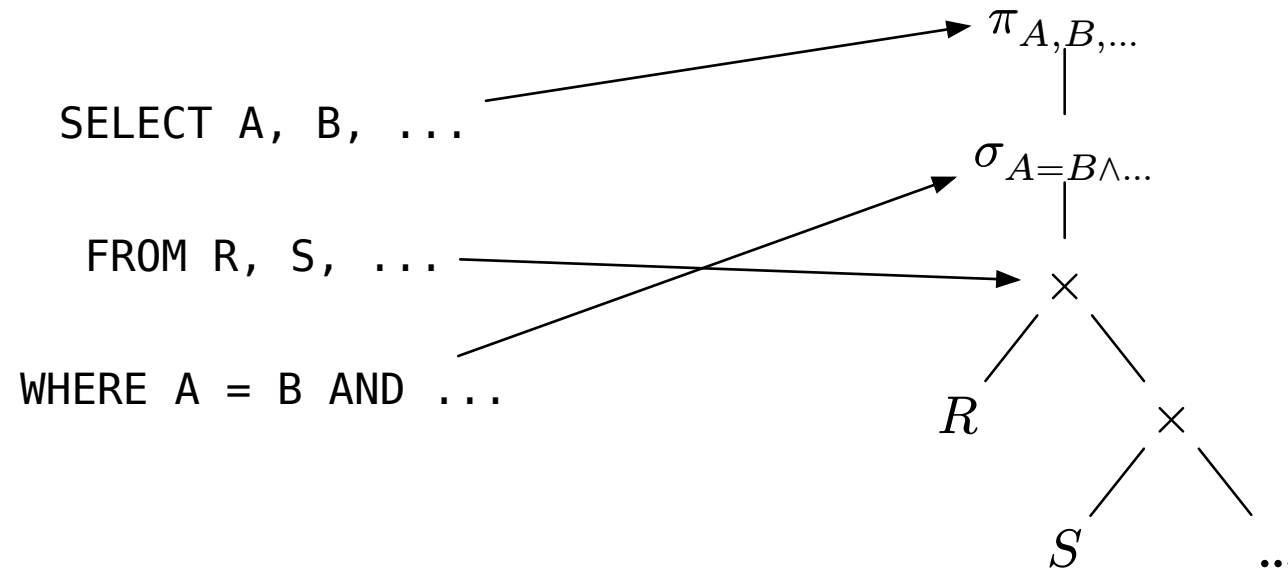
### Code

```
GroupBy([A, B], [SUM(C) AS C, AVG(D) AS D], In)
```

Group tuples by $A$ and $B$; For each group compute $\text{SUM}(C)$, $\text{AVG}(D)$

---

[3]We'll eventually talk about other grouping strategies (e.g., Window Functions)

## Design Questions

- How do we specify which columns/rows to combine?
- How do we specify how to define the new row/column?

## Deriving Rows

- "Group By" attributes[4]
- "Aggregate functions"

### Code                                                    Shorthand

```
GroupBy([A, B], [SUM(C) AS C, AVG(D) AS D], In)
```
$\Sigma_{A,B,C=\text{SUM}(C),D=\text{AVG}(D)}(\text{In})$

Group tuples by $A$ and $B$; For each group compute $\text{SUM}(C)$, $\text{AVG}(D)$

---

[4]We'll eventually talk about other grouping strategies (e.g., Window Functions)

|        | Filter    | Merge        | Derive     |
|--------|-----------|--------------|------------|
| Column | $\pi$     | $\times, \bowtie$ | $\pi$  |
| Row    | $\sigma$  | $\cup$       | $\Sigma$   |

$$\pi_{A,B,\dots}$$

SELECT A, B, ...

$$\sigma_{A=B\wedge\dots}$$

FROM R, S, ...

$$\times$$

WHERE A = B AND ...

$R$    $\times$

$S$    ...

- If you have not yet formed a group **contact me!**.
- Finish the AI quiz ASAP.
- Checkpoint 1 posted and available.