

# CMPT 354: Database System I

Lecture 5. Relational Algebra

# What have we learned

- Lec 1. Database History
- Lec 2. Relational Model
- Lec 3-4. SQL

# Why Relational Algebra matter?

- An essential topic to understand how query processing and optimization work
  - What happened when an SQL is issued to a database?
- Help you master the skills to quickly learn a new query language
  - How to quickly learn XML QL and MongoDB QL?

# Relational Query Languages

- Query languages allow the manipulation and retrieval of data from a database
- Traditionally: QL != programming language
  - Doesn't need to be turing complete
  - Not designed for computation
  - Supports easy, efficient access to large databases
- Recent Years:
  - Everything interesting involves a large data set
  - QLs are quite powerful for expressing algorithms at scale

# Formal Query Languages

- Relational Algebra
  - Procedural query language
  - used to represent execution plans
- Relational Calculus
  - Non-procedural (declarative) query language
  - Describe **what** you want, rather than **how** to compute it
  - Foundation for SQL

# Results of a Query

- Query is a function over **relations**

$$Q(R_1, \dots, R_n) = R_{\text{result}}$$

- The schema of the result relation is determined by the input relation and the query
- Because the result of a query is a relation, it can be used as input to another query

$$Q(\text{blue grid}) = \text{orange grid}, Q(\text{orange grid}) = \text{gray grid}, \dots$$

# Sets v.s. Bags

- Sets: {a, b, c}, {a, d, e, f}, {}, ...
- Bags: {a, a, b, c}, {b, b, b, b, b}, ...
- Relational Algebra has two flavors:
  - Set semantics = standard Relational Algebra
  - Bag semantics = extended Relational Algebra
- DB systems implement bag semantics (Why?)

# Sets v.s. Bags

- Sets: {a, b, c}, {a, d, e, f}, {}, ...
- Bags: {a, a, b, c}, {b, b, b, b, b}, ...
- Relational Algebra has two flavors:
  - **Set semantics = standard Relational Algebra**
  - Bag semantics = extended Relational Algebra
- DB systems implement bag semantics (Why?)



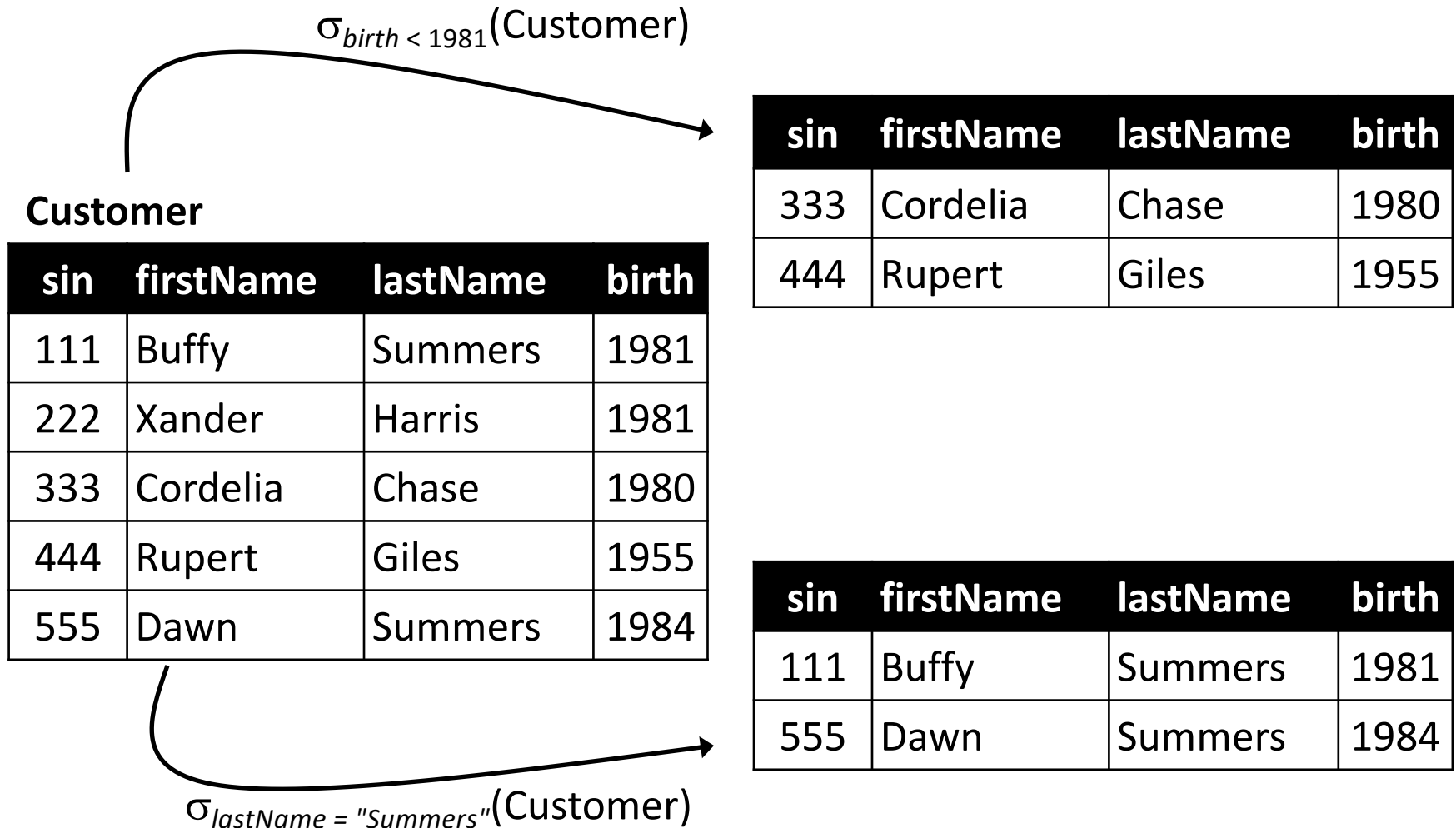
# Relational Algebra Operators

- Core 5 operators
  - Selection ( $\sigma$ )
  - Projection ( $\pi$ )
  - Union ( $\cup$ )
  - Set Difference ( $-$ )
  - Cross product ( $\times$ )
- Additional operators
  - Rename ( $\rho$ )
  - join ( $\bowtie$ )
  - Intersect ( $\cap$ )

# Selection

- The selection operator,  $\sigma$  (sigma), specifies the *rows* to be retained from the input relation
- A selection has the form:  $\sigma_{condition}(relation)$ , where *condition* is a Boolean expression
  - Terms in the condition are comparisons between two fields (or a field and a constant)
  - Using one of the comparison operators:  $<$ ,  $\leq$ ,  $=$ ,  $\neq$ ,  $\geq$ ,  $>$
  - Terms may be connected by  $\wedge$  (and), or  $\vee$  (or),
  - Terms may be negated using  $\neg$  (not)

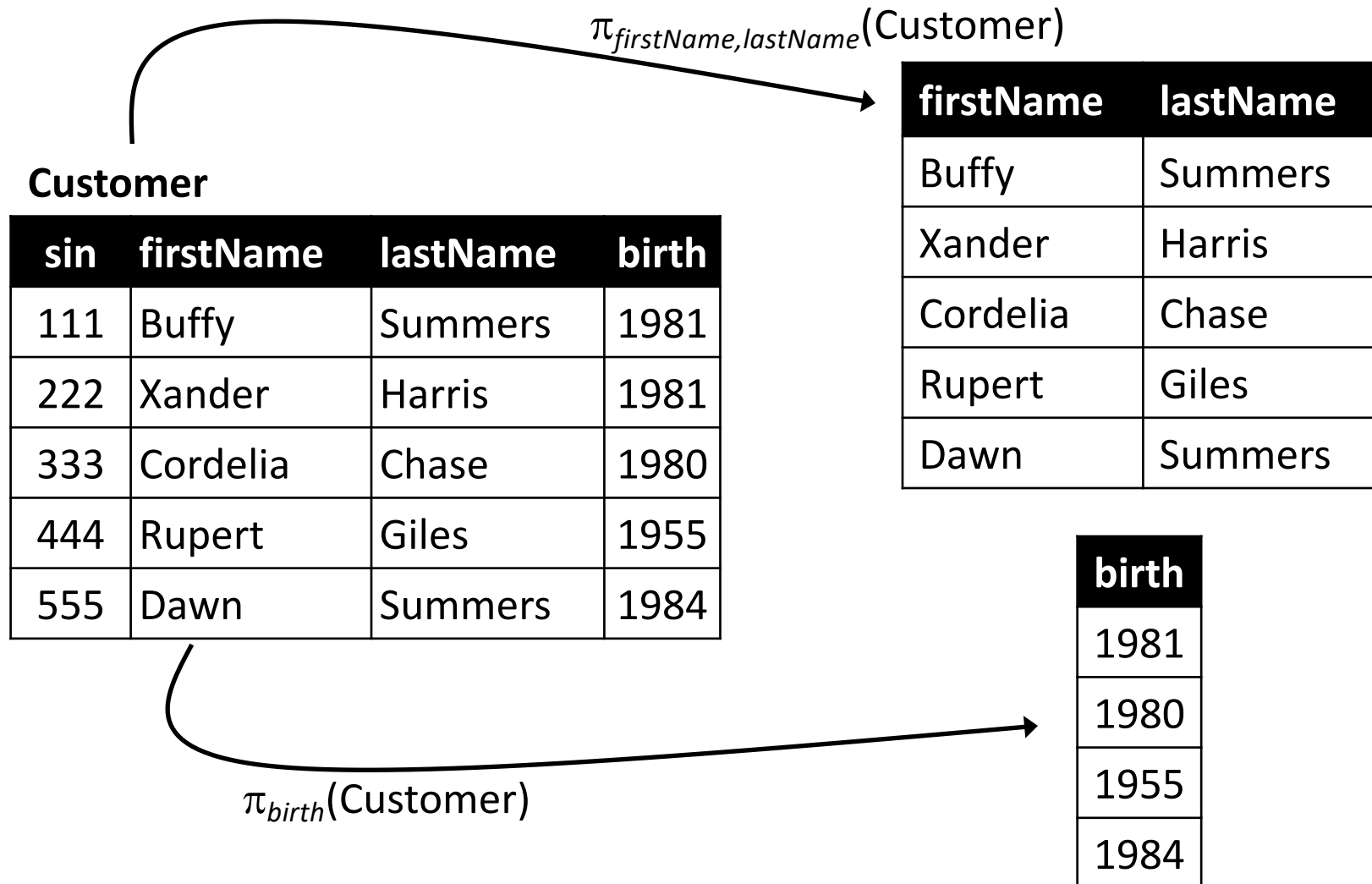
# Selection Example



# Projection

- The projection operator,  $\pi$  (pi), specifies the columns to be retained from the input relation
- A selection has the form:  $\pi_{columns}(relation)$ 
  - Where *columns* is a comma separated list of column names
  - The list contains the names of the columns to be retained in the result relation

# Projection Example



# Selection and Projection Notes

- Selection and projection eliminate duplicates
  - Since relations are sets
- Both operations require one input relation
- The schema of the result of a selection is *the same as* the schema of the input relation
- The schema of the result of a projection contains just those attributes in the projection list

# Composing Selection and Projection

$\pi_{sin, firstName}(\sigma_{birth < 1982 \wedge lastName = "Summers"}(Customer))$

Customer

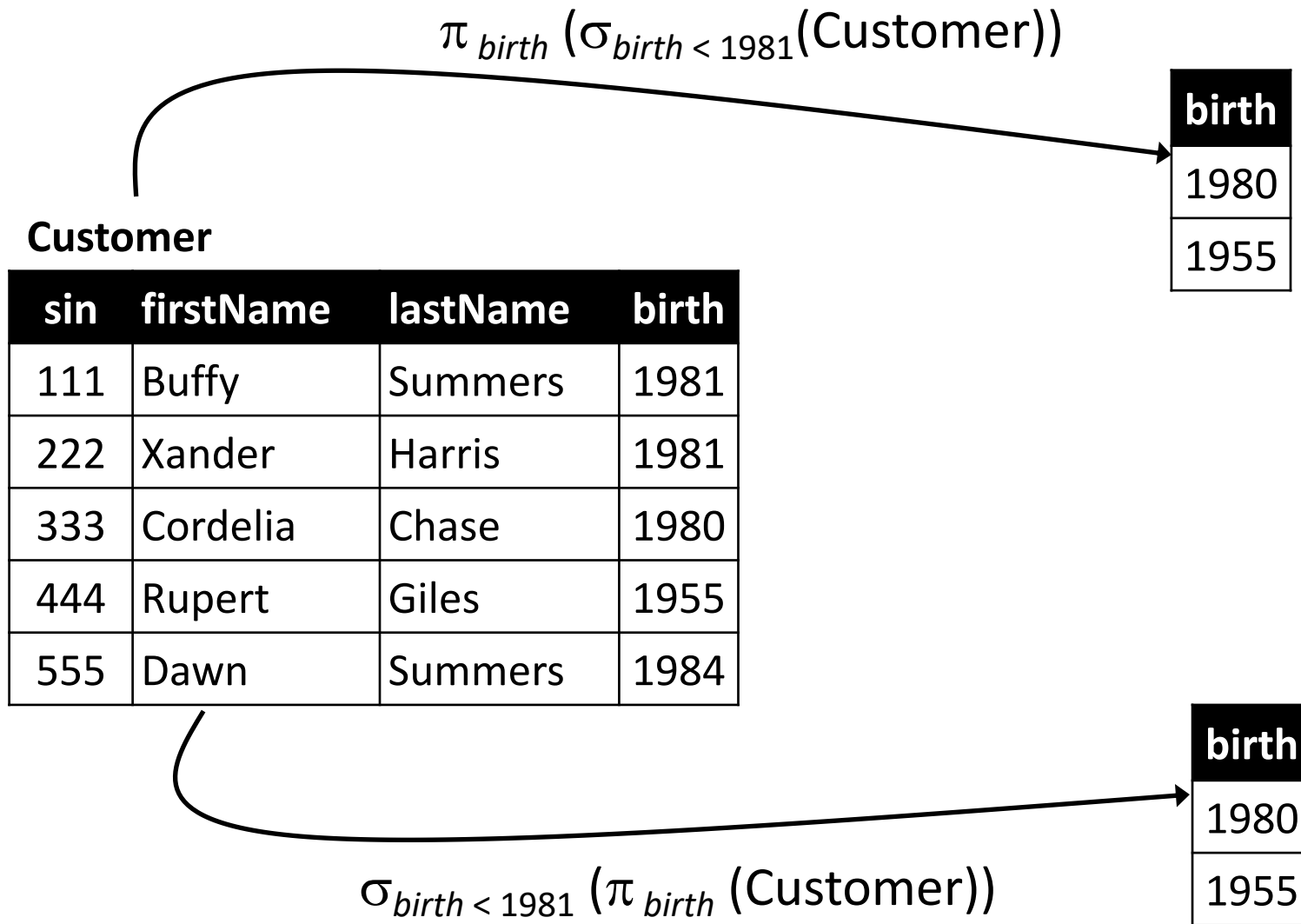
sin	firstName	lastName	birth
111	Buffy	Summers	1981
222	Xander	Harris	1981
333	Cordelia	Chase	1980
444	Rupert	Giles	1955
555	Dawn	Summers	1984

intermediate relation

sin	firstName	lastName	birth
111	Buffy	Summers	1981

sin	firstName
111	Buffy

# Composing Selection and Projection





# Commutative property

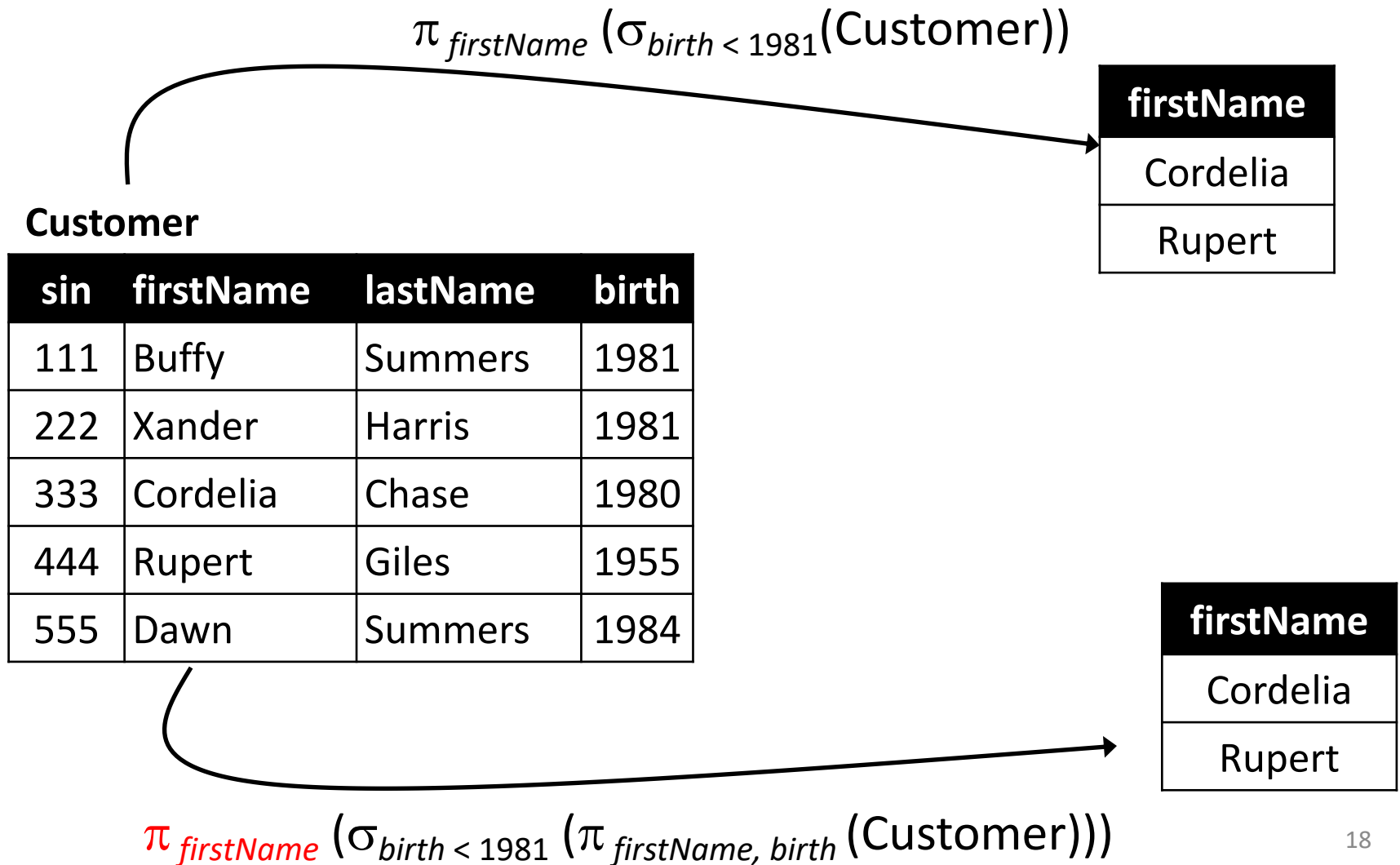
- For example:
  - $x + y = y + x$
  - $x * y = y * x$
- Does it hold for projection and selection?

$$\pi_{columns}(\sigma_{condition}(R)) = \pi_{condition}(\sigma_{columns}(R)) ?$$

- What about

$$\pi_{firstName}(\sigma_{birth < 1981}(\text{Customer}))?$$

# Commutative property



# Set Operations Review

$$A = \{1, 3, 6\}$$

$$B = \{1, 2, 5, 6\}$$

Union ( $\cup$ )

$$A \cup B \equiv B \cup A$$

$$A \cup B = \{1, 2, 3, 5, 6\}$$

Intersection ( $\cap$ )

$$A \cap B \equiv B \cap A$$

$$A \cap B = \{1, 6\}$$

Set Difference ( $-$ )

$$A - B \neq B - A$$

$$A - B = \{3\}$$

$$B - A = \{2, 5\}$$

# Union Compatible Relations

$$A \text{ op } B = R_{\text{result}}$$

- where  $\text{op} = \cup, \cap$ , or  $-$
- $A$  and  $B$  must be **union compatible**
  - Same number of fields
  - Field  $i$  in each schema have the same type

# Union Compatible Relations

Intersection of the Employee and Customer relations

Customer

sin	firstName	lastName	birth
111	Buffy	Summers	1981
222	Xander	Harris	1981
333	Cordelia	Chase	1980
444	Rupert	Giles	1955
555	Dawn	Summers	1984

Employee

sin	firstName	lastName	salary
208	Clark	Kent	80000.55
111	Buffy	Summers	22000.78
412	Carol	Danvers	64000.00

The two relations are not union compatible as birth is a DATE and salary is a REAL

We can carry out preliminary operations to make the relations union compatible

$$\pi_{sin, firstName, lastName}(\text{Customer}) \cap \pi_{sin, firstName, lastName}(\text{Employee})$$

# Union Compatible Relations

$$A \text{ op } B = R_{\text{result}}$$

- where  $\text{op} = \cup, \cap$ , or  $-$
- $A$  and  $B$  must be **union compatible**
  - Same number of fields
  - Field  $I$  in each schema have the same type
- Result schema borrowed from  $A$

$$A(\text{age int}) \cup B(\text{num int}) = R_{\text{result}} (\text{age int})$$

# Union

A

sin	firstName	lastName
111	Buffy	Summers
222	Xander	Harris
333	Cordelia	Chase
444	Rupert	Giles
555	Dawn	Summers

B

sin	firstName	lastName
208	Clark	Kent
111	Buffy	Summers
412	Carol	Danvers

$A \cup B$

sin	firstName	lastName
111	Buffy	Summers
222	Xander	Harris
333	Cordelia	Chase
444	Rupert	Giles
555	Dawn	Summers
208	Clark	Kent
412	Carol	Danvers

# Set Difference

A

sin	firstName	lastName
111	Buffy	Summers
222	Xander	Harris
333	Cordelia	Chase
444	Rupert	Giles
555	Dawn	Summers

A – B

sin	firstName	lastName
222	Xander	Harris
333	Cordelia	Chase
444	Rupert	Giles
555	Dawn	Summers

B

sin	firstName	lastName
208	Clark	Kent
111	Buffy	Summers
412	Carol	Danvers

B – A

sin	firstName	lastName
208	Clark	Kent
412	Carol	Danvers



# Note on Set Difference

- Notice that most operators are monotonic
  - Increasing size of inputs  $\rightarrow$  outputs grow
- Set Difference is non-monotonic
  - Example:  $A - B$
  - Increasing the size of  $B$  could decrease output size
- Set difference is blocking:
  - For  $A - B$ , must wait for all  $B$  tuples before any results

# Intersection

A

sin	firstName	lastName
111	Buffy	Summers
222	Xander	Harris
333	Cordelia	Chase
444	Rupert	Giles
555	Dawn	Summers

B

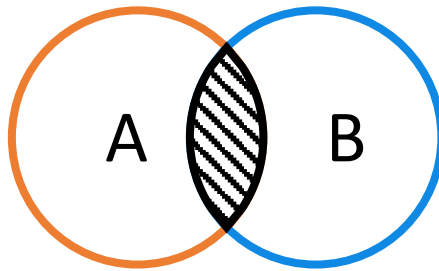
sin	firstName	lastName
208	Clark	Kent
111	Buffy	Summers
412	Carol	Danvers

$A \cap B$

sin	firstName	lastName
111	Buffy	Summers

# Note on Intersect

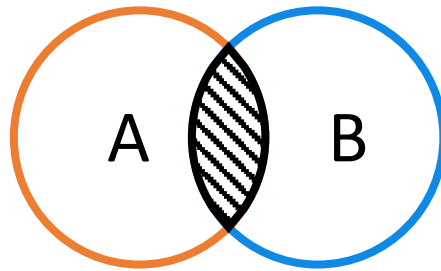
- $A \cap B = R_{\text{result}}$



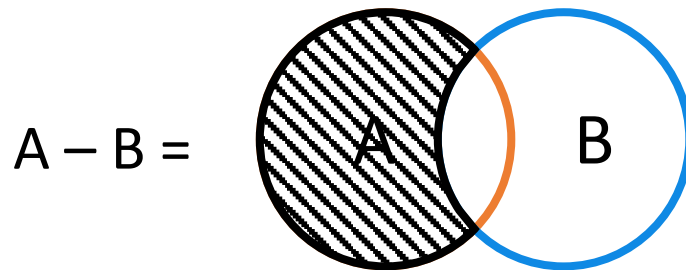
- Can we express using other operators?
  - $A \cap B = ?$

# Note on Intersect

- $A \cap B = R_{\text{result}}$



- Can we express using other operators?
  - $A \cap B = A - (A - B)$



# Cartesian Product

$$A(a_1, \dots, a_n) \times B(a_{n+1}, \dots, a_m) = R_{\text{result}}(a_1, \dots, a_m)$$

- Each row of A paired with each row of B
  - Result schema concatenates A and B's fields
  - Names are inherited if possible (i.e. if not duplicated)
    - If two field names are the same (i.e., a *naming conflict* occurs) and the affected columns are referred to by position
  - If *R* contains *m* records, and *S* contains *n* records, the result relation will contain *m* \* *n* records

# Cartesian Product Example

$\sigma_{lastName = "Summers"}(\text{Customer})$

sin	firstName	lastName	birth
111	Buffy	Summers	1981
555	Dawn	Summers	1984

Account

acc	type	balance	sin
01	CHQ	2101.76	111
02	SAV	11300.03	333
03	CHQ	20621.00	444

$\sigma_{lastName = "Summers"}(\text{Customer}) \times \text{Account}$

<b>1</b>	firstName	lastName	birth	acc	type	balance	<b>8</b>
111	Buffy	Summers	1981	01	CHQ	2101.76	111
111	Buffy	Summers	1981	02	SAV	11300.03	333
111	Buffy	Summers	1981	03	CHQ	20621.00	444
555	Dawn	Summers	1984	01	CHQ	2101.76	111
555	Dawn	Summers	1984	02	SAV	11300.03	333
555	Dawn	Summers	1984	03	CHQ	20621.00	444

# Renaming

- It is sometimes useful to assign names to the results of a relational algebra query
- The rename operator,  $\rho$  (rho)
  - $\rho_S(R)$  renames a relation
  - $\rho_{S(a_1, a_2, \dots, a_n)}(R)$  renames a relation and its attributes
  - $\rho_{\text{new/old}}(R)$  renames specified attributes

R

sid1	firstName	lastName	birth	acc	type	balance	sid2
111	Buffy	Summers	1981	01	CHQ	2101.76	111
111	Buffy	Summers	1981	02	SAV	11300.03	333
111	Buffy	Summers	1981	03	CHQ	20621.00	444
555	Dawn	Summers	1984	01	CHQ	2101.76	111
555	Dawn	Summers	1984	02	SAV	11300.03	333
555	Dawn	Summers	1984	03	CHQ	20621.00	444

$\rho_{\text{sid1/1, sid2/8}}(R)$

# Largest Balance

- Find the account with the largest balance; return *accNumber*
  - Find accounts which are less than some other account

$$\sigma_{account.balance < d.balance} (Account \times \rho_d (Account))$$

- Use set difference to find the account with the largest balance

$$\pi_{accNumber} (Account) - \pi_{account.accNumber} (\sigma_{account.balance < d.balance} (Account \times \rho_d (Account)))$$



# Relational Algebra Operators

- Core 5 operations
  - Selection ( $\sigma$ )
  - Projection ( $\pi$ )
  - Union ( $\cup$ )
  - Set Difference ( $-$ )
  - Cross product ( $\times$ )
- Additional operations
  - Rename ( $\rho$ )
  - Intersect ( $\cap$ )
  - Join ( $\bowtie$ )

# Relational Algebra Exercises

- **Student** (sID, lastName, firstName, cgpa)
  - 101, Jordan, Michael, 3.8
- **Offering** (oID, dept, cNum, term, instructor)
  - abc, CMPT, 354, Fall 2018, Jiannan
- **Took** (sID, oID, grade)
  - 101, abc, 95

1. sID of all students who have earned some grade over 80 and some grade below 50.

$$\pi_{sID}(\sigma_{grade > 80}(\text{Took})) \cap \pi_{sID}(\sigma_{grade < 50}(\text{Took}))$$

# Relational Algebra Exercises

- **Student** (sID, lastName, firstName, cgpa)
  - 101, Jordan, Michael, 3.8
- **Offering** (oID, dept, cNum, term, instructor)
  - abc, CMPT, 354, Fall 2018, Jiannan
- **Took** (sID, oID, grade)
  - 101, abc, 95

2. Student number of all students who have taken CMPT 354

$$\pi_{sID} (\sigma_{Offering.oID = Took.oID \wedge dept = 'CMPT' \wedge cNum = 354} (Offering \times Took))$$

# (Inner) Joins

- **Motivation**

- Simplify some queries that require a Cartesian product

- **Natural Join:**  $R \bowtie S = \pi_A(\sigma_\theta(R \times S))$

- **Theta Join:**  $R \bowtie_\theta S = \sigma_\theta(R \times S)$

- **Equijoin:**  $R \bowtie_\theta S = \sigma_\theta(R \times S)$

- Join condition  $\theta$  consists only of equalities

# Natural Join

- There is often a natural way to join two relations
  - Join based on common attributes
  - Eliminate duplicate common attributes from the result

**Customer**

sin	firstName	lastName	birth
111	Buffy	Summers	1981
222	Xander	Harris	1981
333	Cordelia	Chase	1980
444	Rupert	Giles	1955

**Employee**

sin	firstName	lastName	salary
208	Clark	Kent	80000.55
111	Buffy	Summers	22000.78
396	Dawn	Allen	41000.21

**Customer ⋈ Employee**

sin	firstName	lastName	birth	salary
111	Buffy	Summers	1981	22000.78

# Natural Join

$$R \bowtie S$$

- Meaning:  $R \bowtie S = \pi_A(\sigma_\theta (R \times S))$
- Where:
  - Selection  $\sigma_\theta$  checks equality of all common attributes (i.e., attributes with same names)
  - Projection  $\pi_A$  eliminates duplicate common attributes
- The natural join of two tables with *no fields in common* is the Cartesian product
  - Not the empty set

# Natural Join Example

**R**

A	B	C	D
111	Buffy	Summers	1981
222	Xander	Harris	1981
333	Cordelia	Chase	1980
444	Rupert	Giles	1955

**S**

A	B	C	E
208	Clark	Kent	80000.55
111	Buffy	Summers	22000.78
396	Dawn	Allen	41000.21

$$R \bowtie S = \pi_{A,B,C,D,E}(\sigma_{R.A=S.A \wedge R.B=S.B \wedge R.C=S.C}(R \times S))$$

A	B	C	D	E
111	Buffy	Summers	1981	22000.78

# Theta Join

$$R \bowtie_{\theta} S = \sigma_{\theta}(R \times S)$$

- Most general form
  - $\theta$  can be any condition
- No projection in this case!
  - Result schema same as cross product



# Theta Join Example

**Customer**

sin	firstName	lastName	birth
111	Buffy	Summers	1981
222	Xander	Harris	1981
333	Cordelia	Chase	1980
444	Rupert	Giles	1955
555	Dawn	Summers	1984

**Employee**

sin	firstName	lastName	salary
208	Clark	Kent	80000.55
111	Buffy	Summers	22000.78
412	Carol	Danvers	64000.00

**Customer** ⋈<sub>Customer.sin < Employee.sin</sub> **Employee**

1	2	3	birth	5	6	7	salary
111	Buffy	Summers	1981	208	Clark	Kent	80000.55
111	Buffy	Summers	1981	412	Carol	Danvers	64000.00
222	Xander	Harris	1981	412	Carol	Danvers	64000.00
333	Cordelia	Chase	1980	412	Carol	Danvers	64000.00

# Equi-Joins

$$R \bowtie_{\theta} S = \sigma_{\theta}(R \times S)$$

- A theta join where  $\theta$  is an equality predicate

**Customer**

sin	firstName	lastName	birth
111	Buffy	Summers	1981
222	Xander	Harris	1981
333	Cordelia	Chase	1980
444	Rupert	Giles	1955

**Employee**

sin	firstName	lastName	salary
208	Clark	Kent	80000.55
111	Buffy	Summers	22000.78
396	Dawn	Allen	41000.21

**Customer**  $\bowtie_{\text{Customer.sin} = \text{Employee.sin}}$  **Employee**

1	2	3	birth	5	6	7	salary
111	Buffy	Summers	1981	111	Buffy	Summers	22000.78

# (Inner) Joins Summary

- **Natural Join:**  $R \bowtie S = \pi_A (\sigma_\theta(R \times S))$ 
  - Equality on all fields with same name in R and in S
  - Projection  $\pi_A$  drops all redundant attributes
- **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 selection  $\theta$
  - No projection
- **Equijoin:**  $R \bowtie_\theta S = \sigma_\theta (R \times S)$ 
  - Join condition  $\theta$  consists only of equalities
  - No projection

# Relational Algebra Exercises

- **Student** (sID, lastName, firstName, cgpa)
  - 101, Jordan, Michael, 3.8
- **Course** (dept, cNum, name, breadth)
  - CMPT, 354, DB, True
- **Offering** (oID, dept, cNum, term, instructor)
  - abc, CMPT, 354, Fall 2018, Jiannan
- **Took** (sID, oID, grade)
  - 101, abc, 95

The names of all students who have passed a breadth course (grade  $\geq 60$  and breadth = True) with Martin

$$\pi_{lastName, firstName} (\sigma_{breadth = True \wedge grade > 60 \wedge instructor = 'Martin'} (Student \bowtie Took \bowtie Offering \bowtie Course))$$

# Different Plans, Same Results

- Semantic equivalence: results are *always* the same

$$\pi_{\text{name}}(\sigma_{\text{cNum}=354} (R \bowtie S))$$

$$\pi_{\text{name}}(\sigma_{\text{cNum}=354} (R) \bowtie S)$$

- Are they equivalent?
- Which one is more efficient?
- Can you make it even more efficient?

# Other Operators

- There are additional relational algebra operators
  - Usually used in the context of query optimization
- Duplicate elimination –  $\delta$ 
  - Used to turn a bag into a set
- Aggregation operators
  - e.g. sum, average
- Grouping –  $\gamma$ 
  - Used to partition tuples into groups
    - Typically used with aggregation

# Summary

- Relational Algebra (RA) operators
  - Five core operators: selection, projection, cross-product, union and set difference
  - Additional operators are defined in terms of the core operators: rename, intersection, join
- Theorem: SQL and RA can express exactly the same class of queries
- Multiple RA queries can be equivalent
  - Same semantics but difference performance
  - Form basis for optimizations

RDBMS translate SQL  $\rightarrow$  RA, then optimize RA

# Acknowledge

- Some lecture slides were copied from or inspired by the following course materials
  - “W4111: Introduction to databases” by Eugene Wu at Columbia University
  - “CSE344: Introduction to Data Management” by Dan Suciu at University of Washington
  - “CMPT354: Database System I” by John Edgar at Simon Fraser University
  - “CS186: Introduction to Database Systems” by Joe Hellerstein at UC Berkeley
  - “CS145: Introduction to Databases” by Peter Bailis at Stanford
  - “CS 348: Introduction to Database Management” by Grant Weddell at University of Waterloo