#### CS 348 Lecture 2

#### **Relational Model Part 1**

#### Semih Salihoğlu

Jan 8<sup>th</sup>, 2025



#### More Info On \*\*Tentative\*\* Deadlines

> Note: Last lecture's slide had a different Milestone 1 deadline.

Week	M	W	<u>Deadline</u>
1	Jan 6	Jan 8	
2	Jan 13	Jan 15	A1 out (Fri Jan 17)
3	Jan 20	Jan 22	Project Milestone 0 due (Wed Jan 22)
4	Jan 27	Jan 29	A1 due / A2 out (Fri Jan 31)
5	Feb 3	Feb 5	Project Milestone 1 due (Fri Feb 7)
6	Feb 10	Feb 12	A2 due (Fri Feb 14)
			Reading Week
7	Feb 24	Feb 26	Midterm (Fri Feb 28 4:30-6:00pm)
8	Mar 3	Mar 5	A3 out (Fri Mar 7)
9	Mar 10	Mar 12	Milestone 2 due (Fri Mar 14)
10	Mar 17	Mar 19	
11	Mar 24	Mar 26	A3 due (Mar 28)
12	Mar 31	Apr 2	Milestone 3 due (TBD sometime this week)

#### More Info About Projects

- ➢ By Milestone 0:
  - Find teammates (use Piazza) and register to a "Learn Group"
- > Drop deadline: around Milestone 1 deadline, can drop out of a group
  - > might desolve a group, so need all teammates consent.
- After drop deadline you cannot drop out of Option 2
- Will provide more details over Piazza

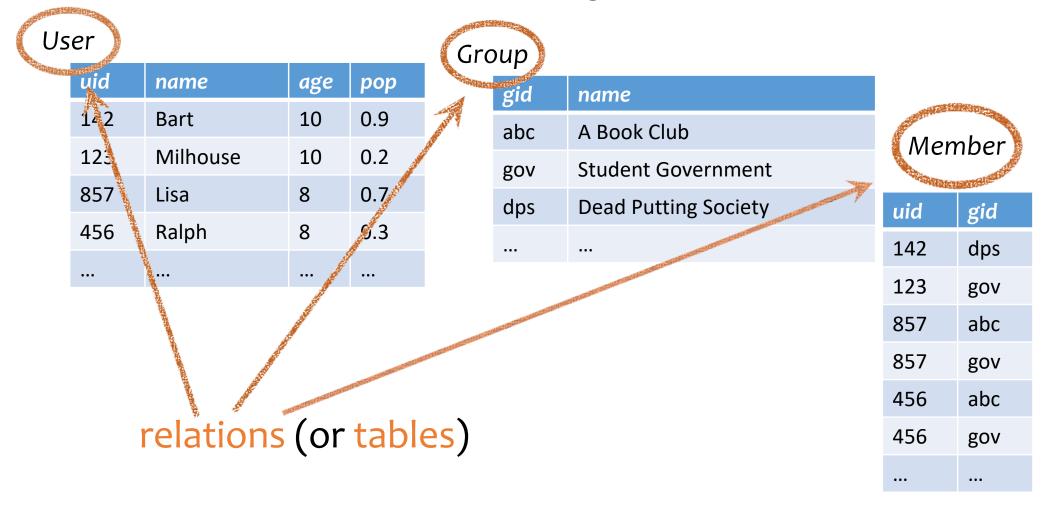
### Outline

• Part 1: Relational data model

• Part 2: Relational algebra

#### Relational data model

Modeling data as **relations** or **tables**, each storing logically related information together



#### Attributes

Us	er			
	uid	name	age	рор
	142	Bart	10	0.9
	123	Milhouse	10	0.2
	857	Lisa	8	0.7
	456	Ralph	8	0.3

#### Group

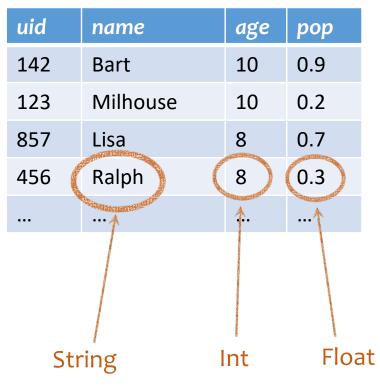
gid	name
abc	A Book Club
gov	Student Government
dps	Dead Putting Society

#### attributes (or columns)

Member	uid	gid
	142	dps
	123	gov
	857	abc
	857	gov
	456	abc
	456	gov

## Domain

#### User



domain (or type)

gid	name
abc	A Book Club
gov	Student Government
dps	Dead Putting Society

Member	uid	gid
	142	dps
	123	gov
	857	abc
	857	gov
	456	abc
	456	gov

# Tuples

#### User

	uid	name	age	рор	
	142	Bart	10	0.9	
	123	Milhouse	10	0.2	
<	857	Lisa	8	0.7	>
<	456	Ralph	8	0.3	>
				•••	
tuples (or rows)					

Duplicates (all attr. have same val) are not allowed

Ordering of rows doesn't matter (even though output can be ordered)

gid	name
abc	A Book Club
gov	Student Government
dps	Dead Putting Society

Member	uid	gid
	142	dps
	123	gov
	857	abc
ed	857	gov
	456	abc
	456	gov

## Set representation of tuples

#### User

uid	name	age	рор
142	Bart	10	0.9
123	Milhouse	10	0.2
857	Lisa	8	0.7
456	Ralph	8	0.3

gid	name
abc	A Book Club
gov	Student Government
edu	Dead Putting Society

	Member	uid	gid
		142	dps
User: {(142, Bart, 10, 0.9),		123	gov
(857, Milhouse, 10, 0.2),}		857	abc
<pre>Group: {(abc, A Book Club),     (gov, Student Government),}</pre>		857	gov
Member: { $(142, dps)$ , $(123, gov)$ ,}		456	abc
		456	gov

#### Relational data model

- A database is a collection of relations (or tables)
- Each relation has a set of attributes (or columns)
- Each attribute has a unique name and a domain (or type)
  - The domains are required to be **atomic**

Single, indivisible piece of information

- Each relation contains a set of tuples (or rows)
  - Each tuple has a value for each attribute of the relation
  - Duplicate tuples are not allowed
    - Two tuples are duplicates if they agree on all attributes
- Simplicity is a virtue!

### Schema vs. instance

- Schema (metadata)
  - Specifies the logical structure of data
  - Is defined at setup time, rarely changes

```
User (uid int, name string, age int, pop float)
Group (gid string, name string)
Member (uid int, gid string)
```

#### Instance

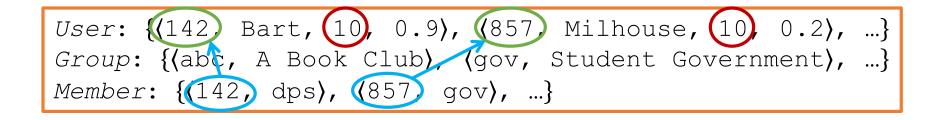
- Represents the data content
- Changes rapidly, but always conforms to the schema
- Typically has additional rules

User: {(142, Bart, 10, 0.9), (857, Milhouse, 10, 0.2), ...}
Group: {(abc, A Book Club), (gov, Student Government), ...}
Member: {(142, dps), (123, gov), ...}

## Integrity constraints

- A set of rules that database instances should follow
- Example:
  - age cannot be negative
  - *uid* should be unique in the User relation
  - uid in Member must refer to a row in User

```
User (uid int, name string, age int, pop float)
Group (gid string, name string)
Member (uid int, gid string)
```



## Integrity constraints

- An instance is only valid if it follows the schema and satisfies all the integrity constraints.
- Reasons to use constraints:
  - Address consistency challenges (last class: duplicate entry for Bob)
  - Ensure data entry/modification respects to database design
  - Protect data from bugs in applications

# Types of integrity constraints

- Tuple-level
  - Domain restrictions, attribute comparisons, etc.
    - E.g. *age* cannot be negative
    - E.g. for flights table, arrival time > take off time
- Relation-level
  - Key constraints (focus in this lecture)
    - E.g. uid should be unique in the User relation
  - Functional dependencies (week 5/6)
- Database-level
  - Referential integrity foreign key (focus in this lecture)
    - uid in Member must refer to a row in User with the same uid

# Key (Candidate Key)

Def: A set of attributes *K* for a relation *R* if

- Condition 1: In no instance of *R* will two different tuples agree on all attributes of *K* 
  - That is, *K* can serve as a "tuple identifier"
- Condition 2: No proper subset of *K* satisfies the above condition
  - That is, *K* is minimal
- Example: User (uid, name, age, pop)
  - uid is a key of User
  - age is not a key (not an identifier)
  - {uid, name} is not a key (not minimal), but a superkey

15

Satisfies only

Condition 1

# Key (Candidate key)

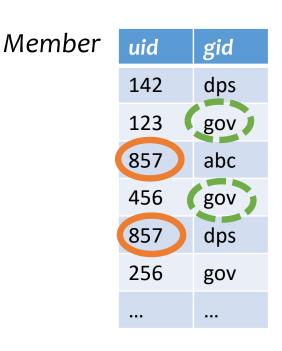
uid	name	age	рор
142	Bart	10	0.9
123	Milhouse	10	0.2
857	Lisa	8	0.7
456	Ralph	8	0.3

- Is name a key of User?
  - Yes? Seems reasonable for this instance
  - No! User names are not unique in general
- Key declarations are part of the schema

### More examples of keys

- Member (uid, gid)
- Only uid?
  - No, because of repeated entries
- Only gid?
  - No, again due to repeated entries
- Use both!
  - {uid, gid}

#### A key can contain multiple attributes



#### More examples of keys

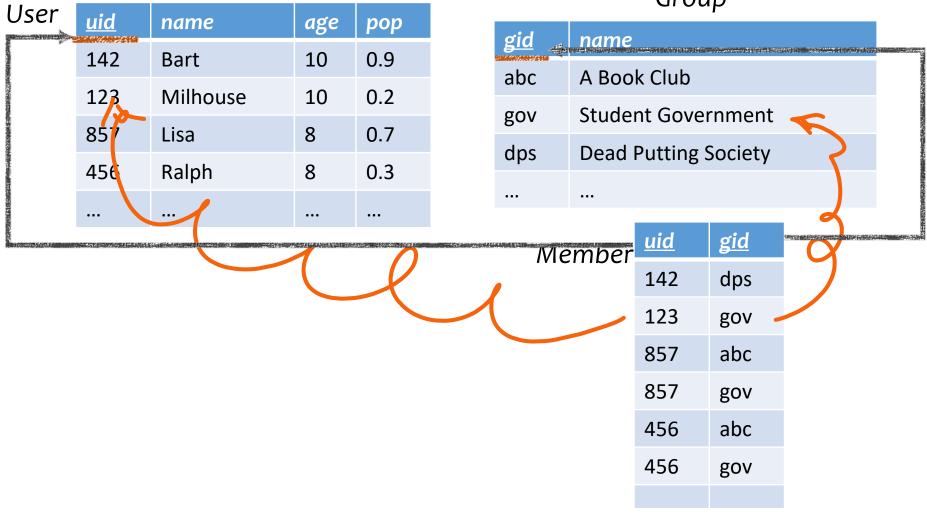
- Address (street\_address, city, province, zip)
  - Key 1: {street\_address, city, province}
- Primary key: a designated candidate key in the schema declaration
  - <u>Underline</u> all its attributes, e.g., Address (<u>street\_address</u>, city, province, <u>zip</u>)

## Use of keys

- More constraints on data, fewer mistakes
- Look up a row by its key value
  - Many selection conditions are "key = value"
- "Pointers" to other rows (often across tables)

### "Pointers" to other rows

• Foreign key: primary key of one relation appearing as attribute of another relation



#### "Pointers" to other rows

 Referential integrity: A tuple with a non-null value for a foreign key must match the primary key value of a tuple in the referenced relation

		Member	<u>uid</u>	<u>gid</u>
	Crown	Member	142	dps
	Group		123	gov
<u>gid</u>	name		857	ON 🜔
abc	A Book Club		857	gov
gov	Student Government		456	abc
dps	Dead Putting Society		456	gov

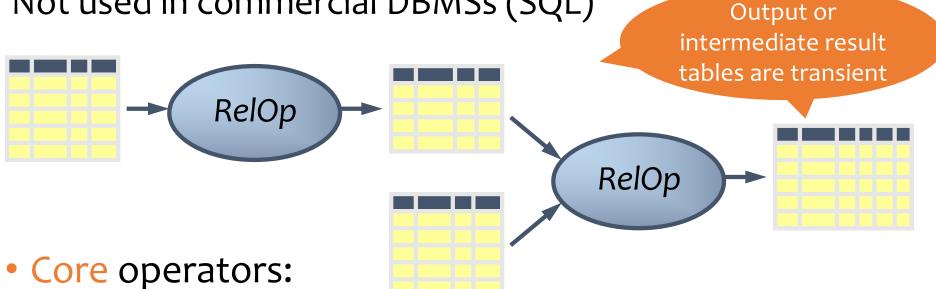
Referential integrity violation!

## Outline

- Part 1: Relational data model
  - Data model
  - Database schema
  - Integrity constraints (keys)
  - Languages
    - Relational algebra (focus in this lecture)
    - SQL (next lecture)
    - Relational calculus (textbook, Ch. 27)
- Part 2: Relational algebra

# **Relational algebra**

- A language for querying relational data based on "operators"
- Not used in commercial DBMSs (SQL)

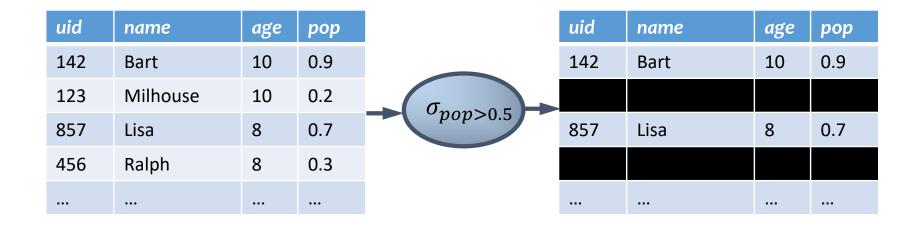


- Selection, projection, cross product, union, difference, and renaming
- Additional, derived operators:
  - Join, natural join, intersection, etc.
- Compose operators to make complex queries

#### Core operator 1: Selection $\sigma$

• Example query: Users with popularity higher than 0.5

 $\sigma_{pop>0.5}$ User



#### Core operator 1: Selection

- Input: a table *R*
- Notation:  $\sigma_p R$ 
  - *p* is called a selection condition (or predicate)
- Purpose: filter rows according to some criteria
- Output: same columns as *R*, but only rows of *R* that satisfy *p*

#### More on selection

- Selection condition can include any column of R, constants, comparison (=, ≤, etc.) and Boolean connectives (∧: and, ∨: or, ¬: not)
  - Example: users with popularity at least 0.9 and age under 10 or above 12

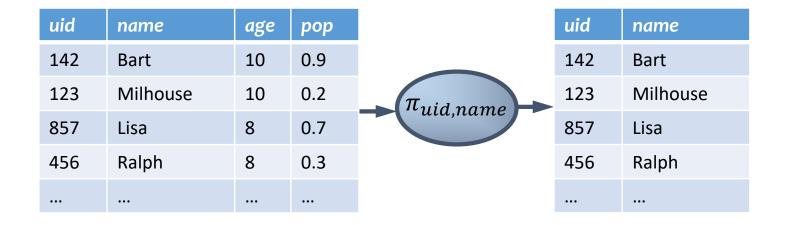
 $\sigma_{pop \ge 0.9 \land (age < 10 \lor age > 12)} User$ 

- You must be able to evaluate the condition over each single row of the input table!
  - Example: the most popular user  $\sigma_{pop \ge every \ pop \ in \ User} \ User \ WRONG!$

#### Core operator 2: Projection $\pi$

• Example: IDs and names of all users

 $\pi_{uid,name}$  User



#### Core operator 2: Projection

- Input: a table *R*
- Notation:  $\pi_L R$ 
  - L is a list of columns in R
- Purpose: output chosen columns
- Output: "same" rows, but only the columns in L

#### More on projection

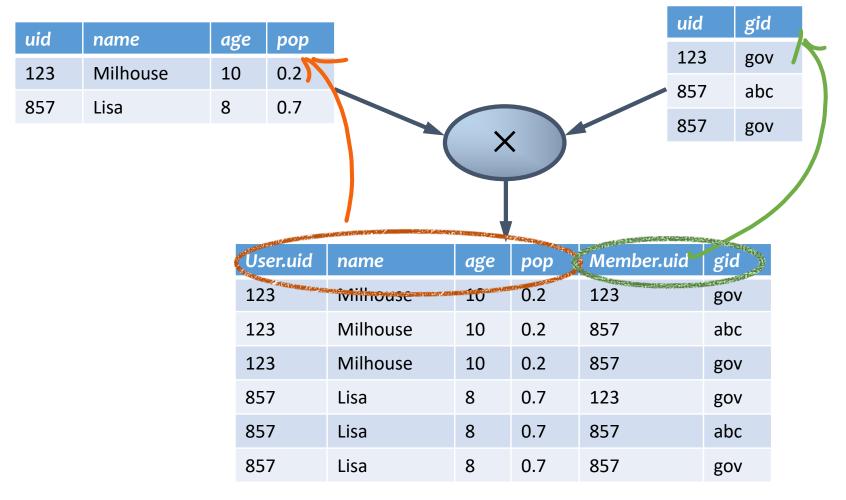
- Duplicate output rows are removed (by definition)
  - Example: user ages

uid name age age pop 0.9 142 Bart 10 10 Milhouse 0.2 123 10  $\pi_{age}$ 0.7 857 Lisa 8 8 456 Ralph 8 0.3 ••• ••• ••• ••• •••

 $\pi_{age}$  User

### Core operator 3: Cross product ×

#### *User*×*Member*



### Core operator 3: Cross product

- Input: two tables R and S
- Notation:  $R \times S$
- Purpose: pairs rows from two tables
- Output: for each row r in R and each s in S, output a row rs (concatenation of r and s)

### A note on column ordering

 Ordering of columns is unimportant as far as contents are concerned

U.uid	name	age	рор	M.uid	gid		M.uid	gid	U.uid	name	age	рор
123	Milhouse	10	0.2	123	gov		123	gov	123	Milhouse	10	0.2
123	Milhouse	10	0.2	857	abc		857	abc	123	Milhouse	10	0.2
123	Milhouse	10	0.2	857	gov		857	gov	123	Milhouse	10	0.2
857	Lisa	8	0.7	123	gov	=	123	gov	857	Lisa	8	0.7
857	Lisa	8	0.7	857	abc		857	abc	857	Lisa	8	0.7
857	Lisa	8	0.7	857	gov		857	gov	857	Lisa	8	0.7

• So cross product is commutative, i.e., for any R and  $S, R \times S = S \times R$  (up to the ordering of columns)

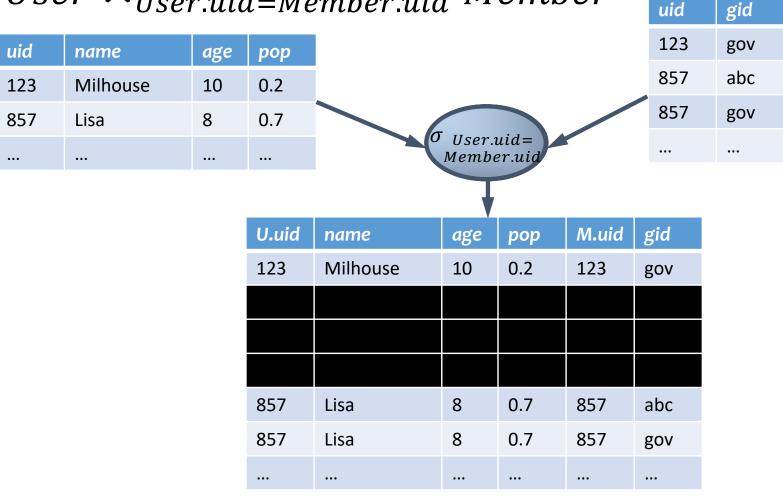
#### Derived operator 1: Join ⋈

 Info about users, plus IDs of their groups User ⋈<sub>User.uid=Member.uid</sub> Member

									gia
uid	name	age	рор					123	gov
123	Milhouse	10	0.2					857	abc
857	Lisa	8	0.7					857	gov
			U.uid	name	age	рор	M.uid	gid	
			123	Milhouse	10	0.2	123	gov	
			123	Milhouse	10	0.2	857	abc	
			123	Milhouse	10	0.2	857	gov	
			857	Lisa	8	0.7	123	gov	
			857	Lisa	8	0.7	857	abc	
			857	Lisa	8	0.7	857	gov	

#### Derived operator 1: Join ⋈

 Info about users, plus IDs of their groups User ⋈<sub>User.uid=Member.uid</sub> Member

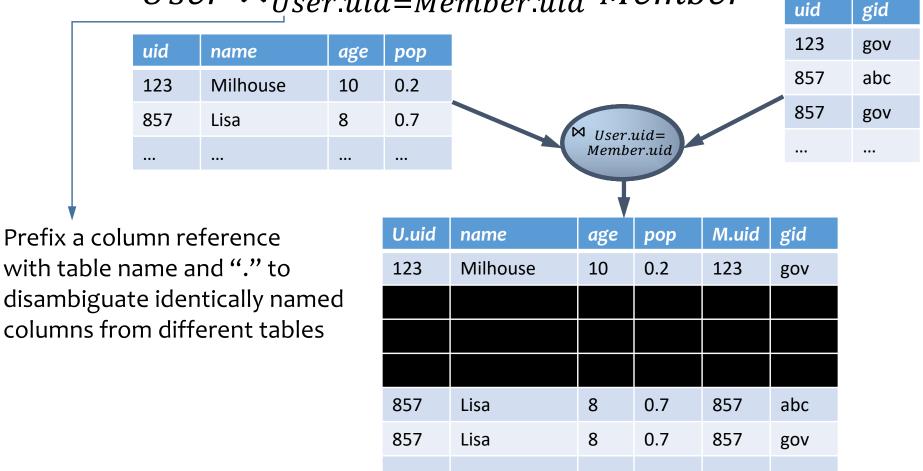


#### Derived operator 1: Join ⋈

 Info about users, plus IDs of their groups User ⋈<sub>User.uid=Member.uid</sub> Member

• • •

• • •



...

...

• • •

• • •

#### Derived operator 1: Join

- Input: two tables R and S
- Notation:  $R \bowtie_p S$ 
  - *p* is called a join condition (or predicate)
- Purpose: relate rows from two tables according to some criteria
- Output: for each row r in R and each row s in S, output a row rs if r and s satisfy p
- Shorthand for  $\sigma_p(R \times S)$
- (A.k.a. "theta-join")

# Derived operator 2: Natural join

### $User \bowtie Member$

 $= \pi_{uid,name,age,pop,gid} \begin{pmatrix} User \bowtie User.uid = Member \\ Member.uid \end{pmatrix}$ uid gid uid age pop name 123 gov Milhouse 0.2 123 10 857 abc 0.7 857 Lisa 8  $\bowtie$ 857 gov ... ... ... ... • • • • • • uid gid age pop name 123 Milhouse 10 0.2 gov 857 0.7 8 Lisa abc 857 0.7 Lisa 8 gov ... ... • • • ... ...

# Derived operator 2: Natural join

- Input: two tables R and S
- Notation:  $R \bowtie S$
- Purpose: relate rows from two tables, and
  - Enforce equality between identically named columns
  - Eliminate one copy of identically named columns
- Shorthand for  $\pi_L(R \bowtie_p S)$ , where
  - *p* equates each pair of attributes common to *R* and *S*
  - *L* is the union of attributes from *R* and *S* (with common attributes removed)

# Core operator 4: Union

- Input: two tables R and S
- Notation:  $R \cup S$ 
  - *R* and *S* must have identical schema
- Output:
  - Has the same schema as *R* and *S*
  - Contains all rows in *R* and all rows in *S* (with duplicate rows removed)

uid	gid		uid	gid		uid	gid
123	gov	U	123	gov	=	123	gov
857	abc		901	edf		857	abc
						901	edf

# Core operator 5: Difference

- Input: two tables R and S
- Notation: R S
  - R and S must have identical schema
- Output:
  - Has the same schema as *R* and *S*
  - Contains all rows in R that are not in S

uid	gid		uid	gid		uid	gid
123	gov	—	123	gov	=	857	abc
857	abc		901	edf			

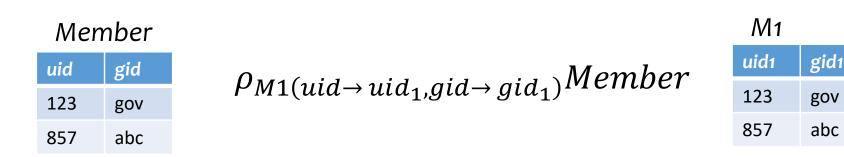
# Derived operator 3: Intersection

- Input: two tables R and S
- Notation:  $R \cap S$ 
  - R and S must have identical schema
- Output:
  - Has the same schema as *R* and *S*
  - Contains all rows that are in both R and S
- Shorthand for R (R S)
- Also equivalent to S (S R)
- And to  $R \bowtie S$

 Find tuples in R not in S
 Remove those tuples from R

# Core operator 6: Renaming

- Input: a table (or an expression) R
- Notation:  $\rho_S R$ ,  $\rho_{(A_1 \to A'_1,...)}R$ , or  $\rho_{S(A_1 \to A'_1,...)}R$
- Purpose: "rename" a table and/or its columns
- Output: a table with the same rows as *R*, but called differently

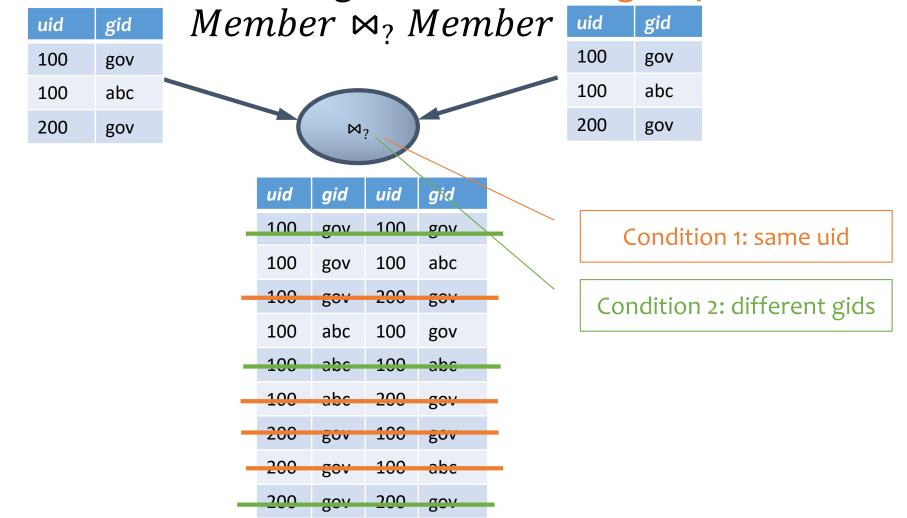


# 9. Core operator: Renaming

- As with all other relational operators, it doesn't modify the database
  - Think of the renamed table as a copy of the original
- Used to: Avoid confusion caused by identical column names

# 9. Core operator: Renaming

#### IDs of users who belong to at least two groups



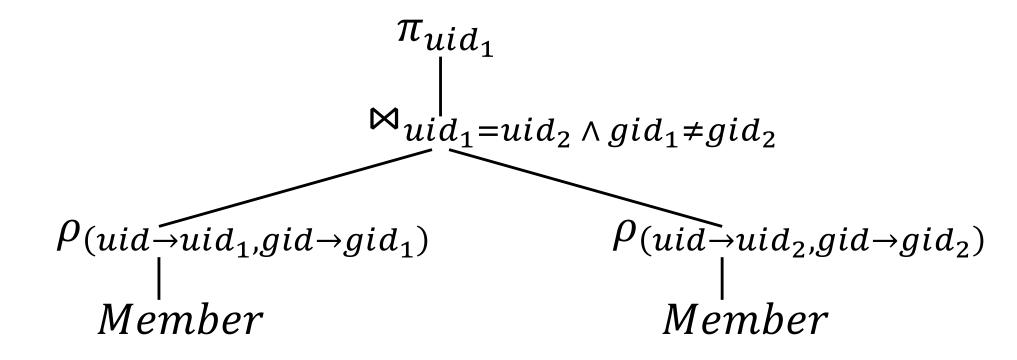
# Renaming example

 IDs of users who belong to at least two groups *Member* ⋈<sub>?</sub> *Member*

 $\pi_{uid} \left( \begin{matrix} Member \bowtie_{Member.uid=Member.uid \land Member} \\ Member.gid \neq Member.gid \end{matrix} \right)$ 

$$\pi_{uid_{1}} \begin{pmatrix} \rho_{(uid \to uid_{1}, gid \to gid_{1})} Member \\ \bowtie_{uid_{1} = uid_{2} \land gid_{1} \neq gid_{2}} \\ \rho_{(uid \to uid_{2}, gid \to gid_{2})} Member \end{pmatrix}$$

### Expression tree notation



## Take-home Exercises

- Exercise 1: IDs of groups who have at least 2 users?
- Exercise 2: IDs of users who belong to at least three groups?

# Summary of operators

**Core Operators** 

- 1. Selection:  $\sigma_p R$
- 2. Projection:  $\pi_L R$
- 3. Cross product:  $R \times S$
- 4. Union:  $R \cup S$
- 5. Difference: R S
- 6. Renaming:  $\rho_{S(A_1 \rightarrow A'_1, A_2 \rightarrow A'_2, \dots)} R$

**Derived Operators** 

- 1. Join:  $R \bowtie_p S$
- 2. Natural join:  $R \bowtie S$
- 3. Intersection:  $R \cap S$

Note: Only use these operators for assignments & exams

User (<u>uid</u> int, name string, age int, pop float) Group (<u>gid</u> string, name string) Member (<u>uid</u> int, <u>gid</u> string)

• All groups (ids) that Lisa belongs to

User (<u>uid</u> int, name string, age int, pop float) Group (<u>gid</u> string, name string) Member (<u>uid</u> int, <u>gid</u> string)

• All groups (ids) that Lisa belongs to

Writing a query bottom-up:

	uid	name	age	рор	
	857	Lisa	8	0.7	
Who's Lisa?	$\sigma_{name="Lisa"}$				
uid	nan	ne	age	рор	

uid	name	age	рор
123	Milhouse	10	0.2
857	Lisa	8	0.7

### Member

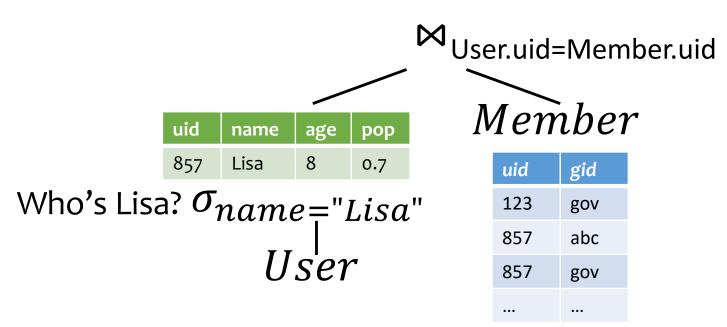
uid	gid
123	gov
857	abc
857	gov

User (<u>uid</u> int, name string, age int, pop float) Group (<u>gid</u> string, name string) Member (<u>uid</u> int, <u>gid</u> string)

• All groups (ids) that Lisa belongs to

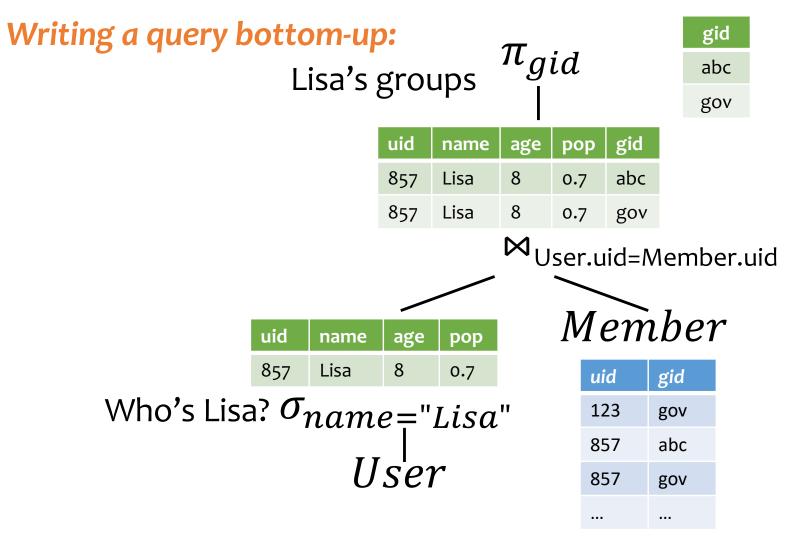
Writing a query bottom-up:

uid	name	age	рор	gid
857	Lisa	8	0.7	abc
857	Lisa	8	0.7	gov



User (<u>uid</u> int, name string, age int, pop float) Group (<u>gid</u> string, name string) Member (<u>uid</u> int, <u>gid</u> string)

• All groups (ids) that Lisa belongs to



# Take home ex.

User (<u>uid</u> int, name string, age int, pop float) Group (<u>gid</u> string, name string) Member (<u>uid</u> int, <u>gid</u> string)

 All groups (ids) that Lisa belongs to names?

# Summary

- Part 1: Relational data model
  - Data model
  - Database schema
  - Integrity constraints (keys)
  - Languages (relational algebra, relational calculus, SQL)
- Part 2: Relational algebra basic language
  - Core operators & derived operators (how to write a query)
- What's next?
  - More examples in RA
  - Relational calculus
  - SQL