# Relational Algebra (continued)
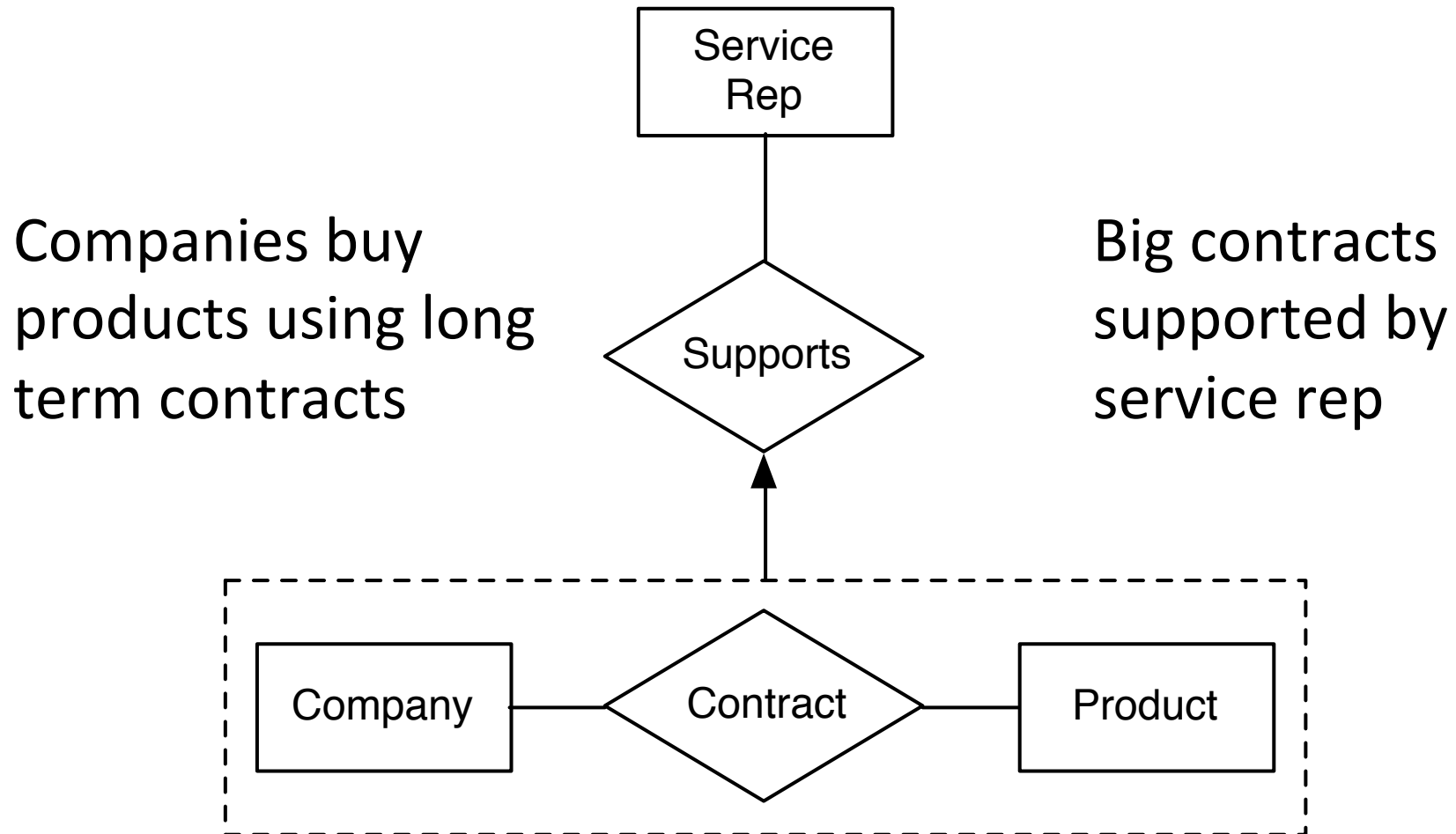
# Annoucements

HW1 Due

Project 1 Part 1:

    Passwords for part 2 on front for each user

    All capital letters
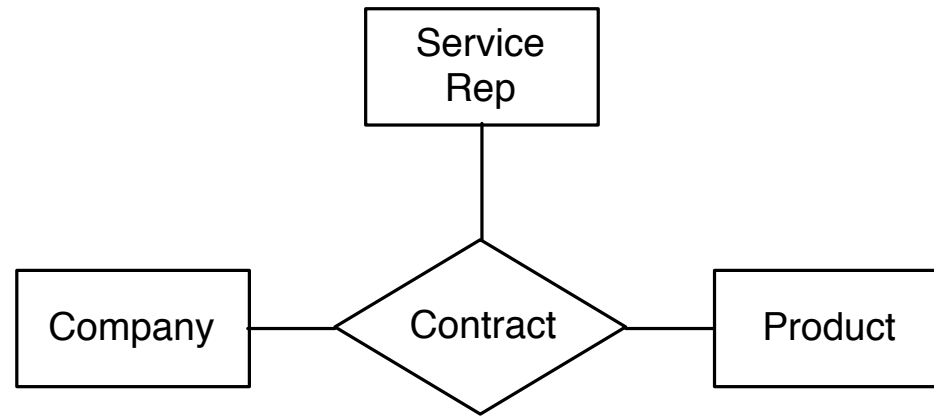
    Will provide access to list of Azure codes

Part 2: Available now

# Aggregate example: Why not ternary?

Service
Rep

Companies buy
products using long
term contracts

Supports

Big contracts
supported by
service rep

Company — Contract — Product

# Aggregate example: Why not ternary?



Companies buy products with a contract;
   **all** contracts have service reps
Relationship sets: Connects N entities
   All entities are required

# Cross-Product

S1

| sid | name | gpa | age |
|-----|------|-----|-----|
| 1 | eugene | 4 | 20 |
| 2 | barak | 3 | 21 |
| 3 | trump | 2 | 88 |

R1

| sid | rid | day |
|-----|-----|-----|
| 1 | 101 | 10/10 |
| 2 | 102 | 11/11 |

S1 x R1 =

| (sid) | name | gpa | age | (sid) | rid | day |
|-------|------|-----|-----|-------|-----|-----|
| 1 | eugene | 4 | 20 | 1 | 101 | 10/10 |
| 2 | barak | 3 | 21 | 1 | 101 | 10/10 |
| 3 | trump | 2 | 88 | 1 | 101 | 10/10 |
| 1 | eugene | 4 | 20 | 2 | 102 | 11/11 |
| 2 | barak | 3 | 21 | 2 | 102 | 11/11 |
| 3 | trump | 2 | 88 | 2 | 102 | 11/11 |

# Rename

ρ(<new_name>(<mappings>), Q)

Explicitly defines/changes field names of schema

ρ(C(1 → sid1, 5 → sid2), S1 x R1)

C =

| sid1 | name | gpa | age | sid2 | rid | day |
|------|------|-----|-----|------|-----|-----|
| 1 | eugene | 4 | 20 | 1 | 101 | 10/10 |
| 2 | barak | 3 | 21 | 1 | 101 | 10/10 |
| 3 | trump | 2 | 88 | 1 | 101 | 10/10 |
| 1 | eugene | 4 | 20 | 2 | 102 | 11/11 |
| 2 | barak | 3 | 21 | 2 | 102 | 11/11 |
| 3 | trump | 2 | 88 | 2 | 102 | 11/11 |

# How do I know column # in large DB?

You count

Yes, that seems silly, but need to assign identifiers automatically; all solutions are arbitrary

Project     $\pi($  $) =$ 

Select     $\sigma($  $) =$ 

Cross product      $\times$  $=$ 

Difference      $-$  $=$ 

Union      $\cup$  $=$ 

Intersect      $\cap$  $=$ 

# Basic Single Table SELECT

```
SELECT * FROM Students
SELECT name FROM Students
SELECT name FROM Students WHERE age < 21
SELECT name, login FROM Students WHERE gpa >= 3
```

| sid | name | login | age | gpa |
|-----|------|-------|-----|-----|
| 1 | eugene | ewu | 20 | 2.5 |
| 2 | luis | gravano | 25 | 3.5 |
| 3 | martha | martha | 32 | 3.9 |

Students

$\pi_{name}(Students)$

$\pi_{name}(\sigma_{age<21}(Students))$

$\pi_{name,login}(\sigma_{gpa\geq3}(Students))$

# Ambiguous names

E.g.   Students: (sid, name, ...)
        Enrolled: (sid, cid, grade)

*Qualified* names: Use table name: Students.age

Rename: AS (optional): shortcuts, ambiguity, clarity


SELECT Students.sid, Students.name
    FROM Students

SELECT S.sid, S.name FROM Students AS S

SELECT S.sid, S.name FROM Students S

# Related data: Multiple tables

What does this return?

```
SELECT S.name, E.cid
FROM Students S, Enrolled E
WHERE S.sid = E.sid AND
        E.grade = 'A'
```

Enrolled

| sid | cid | grade |
|-----|-----|-------|
| 1   | 2   | A     |
| 1   | 3   | B     |
| 2   | 2   | A+    |

```
SELECT S.name, E.cid
FROM Students S, Enrolled E
WHERE S.sid = E.sid AND
        E.grade = 'A'
```

Students

| sid | name   |
|-----|--------|
| 1   | eugene |
| 2   | luis   |

Result

| name   | cid |
|--------|-----|
| eugene | 2   |

# Multi-Table Semantics

- Modify the FROM clause evaluation
  - 1. FROM clause: compute *cross-product* of Students and Enrolled

Enrolled

| sid | cid | grade |
|-----|-----|-------|
| 1 | 2 | A |
| 1 | 3 | B |
| 2 | 2 | A+ |

Students

| sid | name |
|-----|------|
| 1 | eugene |
| 2 | luis |

Cross-product

| sid | cid | grade | sid | name |
|-----|-----|-------|-----|------|
| 1 | 2 | A | 1 | eugene |
| 1 | 3 | B | 1 | eugene |
| 2 | 2 | A+ | 1 | eugene |
| 1 | 2 | A | 2 | luis |
| 1 | 3 | B | 2 | luis |
| 2 | 2 | A+ | 2 | luis |

# Multi-Table Semantics

Modify the FROM clause evaluation

1. FROM clause: compute *cross-product* of Students, Enrolled

2. WHERE clause: Check conditions, discard tuples that fail

3. SELECT clause: Delete unwanted fields

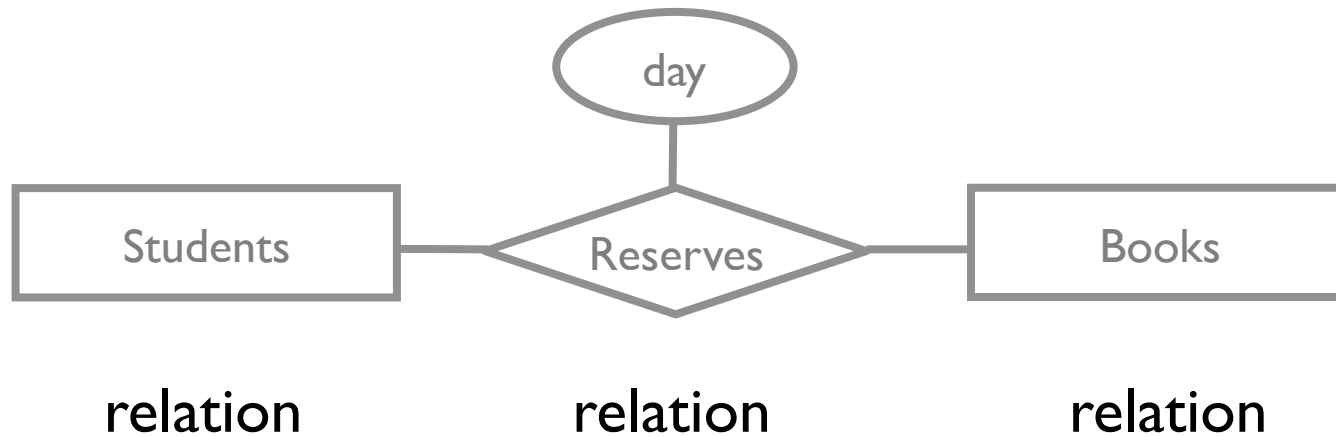# Joins (high level)



relation       relation       relation

What if you want to query across all three tables?
e.g., names of all students that reserved "The Purple Crayon"

Need to combine these tables
Cross product?  But that ignores foreign key references

# Joins (high level)



day

Students — Reserves — Books

relation          relation          relation

## S1

| sid | name | gpa | age |
|-----|------|-----|-----|
| 1 | eugene | 4 | 20 |
| 2 | barak | 3 | 21 |
| 3 | trump | 2 | 88 |

## R1

| sid | rid | day |
|-----|-----|-----|
| 1 | 101 | 10/10 |
| 2 | 102 | 11/11 |

## B1

| rid | name |
|-----|------|
| 101 | The Purple Crayon |
| 102 | 1984 |

# Joins (high level)



**day**

| Students | Reserves | Books |
|:---:|:---:|:---:|

relation              relation              relation

### S1

| sid | name | gpa | age |
|:---:|:---:|:---:|:---:|
| 1 | eugene | 4 | 20 |
| 2 | barak | 3 | 21 |
| 3 | trump | 2 | 88 |

### R1

| sid | rid | day |
|:---:|:---:|:---:|
| 1 | 101 | 10/10 |
| 2 | 102 | 11/11 |

### B1

| rid | name |
|:---:|:---:|
| 101 | The Purple Crayon |
| 102 | 1984 |

# Joins (high level)



day

| Students | Reserves | Books |
| --- | --- | --- |
| relation | relation | relation |

## S1

| sid | name | gpa | age |
| --- | --- | --- | --- |
| 1 | eugene | 4 | 20 |
| 2 | barak | 3 | 21 |
| 3 | trump | 2 | 88 |

## RB1

| sid | (rid) | day | (rid) | name |
| --- | --- | --- | --- | --- |
| 1 | 101 | 10/10 | 101 | The Purple Crayon |
| 2 | 102 | 11/11 | 102 | 1984 |

# Joins (high level)



relation                    relation                    relation

## S1

| sid | name | gpa | age |
|-----|------|-----|-----|
| 1 | eugene | 4 | 20 |
| 2 | barak | 3 | 21 |
| 3 | trump | 2 | 88 |

## RB1

| sid | (rid) | day | (rid) | name |
|-----|-------|-----|-------|------|
| 1 | 101 | 10/10 | 101 | The Purple Crayon |
| 2 | 102 | 11/11 | 102 | 1984 |

# Joins (high level)



relation      relation      relation

### SRB1

| (sid) | (name) | gpa | age | (sid) | (rid) | day | (rid) | (name) |
|-------|--------|-----|-----|-------|-------|-------|-------|--------|
| 1 | eugene | 4 | 20 | 1 | 101 | 10/10 | 101 | The Purple Crayon |
| 2 | barak | 3 | 21 | 2 | 102 | 11/11 | 102 | 1984 |

# theta ($\theta$) Join

$$A \bowtie_c B = \sigma_c(A \times B)$$

Most general form

Result schema same as cross product

Often *far* more efficient to compute than cross product

Commutative

$$(A \bowtie_c B) \bowtie_c C = A \bowtie_c (B \bowtie_c C)$$

# theta (θ) Join

### S1

| sid | name | gpa | age |
|-----|------|-----|-----|
| 1 | eugene | 4 | 20 |
| 2 | barak | 3 | 21 |
| 3 | trump | 2 | 88 |

### R1

| sid | rid | day |
|-----|-----|-----|
| 1 | 101 | 10/10 |
| 2 | 102 | 11/11 |

$$S1 \bowtie_{S1.sid \leq R1.sid} R1 = $$

| (sid) | name | gpa | age | (sid) | rid | day |
|-------|------|-----|-----|-------|-----|-----|
| 1 | eugene | 4 | 20 | 1 | 101 | 10/10 |
| 1 | eugene | 4 | 20 | 2 | 102 | 11/11 |
| 2 | barak | 3 | 21 | 2 | 102 | 11/11 |

# Equi-Join

Common case where the condition is attribute equality

$$A \bowtie_{attr} B = \pi_{all\ attrs\ except\ B.attr}(A \bowtie_{A.attr\ =\ B.attr} B)$$

Result schema only keeps *one copy* of equality fields

Natural Join (A⋈B):

    Equijoin on *all* shared fields (fields w/ same name)

# Equi-Join

### S1

| sid | name | gpa | age |
|-----|------|-----|-----|
| 1 | eugene | 4 | 20 |
| 2 | barak | 3 | 21 |
| 3 | trump | 2 | 88 |

### R1

| sid | rid | day |
|-----|-----|-----|
| 1 | 101 | 10/10 |
| 2 | 102 | 11/11 |

$$S1 \bowtie_{sid} R1 =$$

| sid | name | gpa | age | rid | day |
|-----|------|-----|-----|-----|-----|
| 1 | eugene | 4 | 20 | 101 | 10/10 |
| 2 | barak | 3 | 21 | 102 | 11/11 |

Names of students that reserved book 2

$$\pi_{name}(\sigma_{rid=2} (R1) \bowtie S1)$$

# Equivalent Queries

$p(tmp1, \sigma_{rid=2} (R1))$
$p(tmp2, tmp1 \bowtie S1)$
$\pi_{name}(tmp2)$

$$\pi_{name}(\sigma_{rid=2}(R1 \bowtie S1))$$

# Names of students that reserved db books

Book(rid, type)     Reserve(sid, rid)     Student(sid)

Need to join DB books with reserve and students

$\sigma_{type='db'}$ (Book)

# Names of students that reserved db books

Book(rid, type)    Reserve(sid, rid)    Student(sid)

Need to join DB books with reserve and students

$$\sigma_{type=\text{'db'}} (\text{Book}) \bowtie \text{Reserve}$$

# Names of students that reserved db books

Book(rid, type)    Reserve(sid, rid)    Student(sid)

Need to join DB books with reserve and students

$\sigma_{type='db'}$ (Book) ⋈ Reserve ⋈ Student

# Names of students that reserved db books

Book(rid, type)    Reserve(sid, rid)    Student(sid)

Need to join DB books with reserve and students

$\pi_{name}(\sigma_{type=\text{'db'}} (\text{Book}) \bowtie \text{Reserve} \bowtie \text{Student})$

# Names of students that reserved db books

Book(rid, type)     Reserve(sid, rid)     Student(sid)

Need to join DB books with reserve and students

$\pi_{name}(\sigma_{type='db'} (Book) \bowtie Reserve \bowtie Student)$

More efficient query

$\pi_{name}(\pi_{sid}(( \pi_{rid} \sigma_{type='db'} (Book)) \bowtie Reserve) \bowtie Student)$

Query optimizer can find the more efficient query!

# Names of students that reserved db books

Book(rid, type)     Reserve(sid, rid)     Student(sid)

Need to join DB books with reserve and students

$\pi_{name}(\sigma_{type=\text{'db'}} (Book) \bowtie Reserve \bowtie Student)$

More efficient query

$\pi_{name}(\pi_{sid}(( \pi_{rid} \sigma_{type=\text{'db'}} (Book)) \bowtie Reserve) \bowtie Student)$

Query optimizer can find the more efficient query!

# Names of students that reserved db books

Book(rid, type)　　Reserve(sid, rid)　　Student(sid)

Need to join DB books with reserve and students

$$\pi_{name}(\sigma_{type='db'} (Book) \bowtie Reserve \bowtie Student)$$

More efficient query

$$\pi_{name}(\pi_{sid}(( \pi_{rid} \sigma_{type='db'} (Book)) \bowtie Reserve) \bowtie Student)$$

Query optimizer can find the more efficient query!

# Names of students that reserved db books

Book(rid, type)    Reserve(sid, rid)    Student(sid)

Need to join DB books with reserve and students

$\pi_{name}(\sigma_{type=\text{‘db’}} (\text{Book}) \bowtie \text{Reserve} \bowtie \text{Student})$

More efficient query

$\pi_{name}(\pi_{sid}(( \pi_{rid} \sigma_{type=\text{‘db’}} (\text{Book})) \bowtie \text{Reserve}) \bowtie \text{Student})$

Query optimizer can find the more efficient query!

# Students that reserved DB or HCI book

1. Find all DB or HCI books
2. Find students that reserved one of those books
   - $p(tmp, (\sigma_{type='DB' \lor type='HCI'} (Book))$
   - $\pi_{name}(tmp \bowtie Reserve \bowtie Student)$

Alternatives

define tmp using UNION (how?)

what if we replaced v with ^ in the query?

# Students that reserved a DB and HCI book

Does previous approach work?

$$\rho(tmp, (\sigma_{type='DB' \wedge type='HCI'} (Book))$$
$$\pi_{name}(tmp \bowtie Reserve \bowtie Student)$$

## NO

# Students that reserved a DB and HCI book

Does previous approach work?
1. Find students that reserved DB books
2. Find students that reversed HCI books
3. Intersection

$$p(tmpDB, \pi_{sid}(\sigma_{type='DB'} \; Book \bowtie Reserve))$$
$$p(tmpHCI, \pi_{sid}(\sigma_{type='HCI'} \; Book \bowtie Reserve))$$
$$\pi_{name}((tmpDB \cap tmpHCI) \bowtie Student)$$

# Let's step back

Relational algebra is expressiveness benchmark

> A language equal in expressiveness as relational algebra is relationally complete

But has limitations

> nulls
>
> aggregation
>
> recursion
>
> duplicates

# Equi-Joins are a way of life

Matching of two sets based on shared attributes

Yelp:        Join between your location and restaurants

Market:      Join between consumers and suppliers

High five:   Join between two hands on time and space

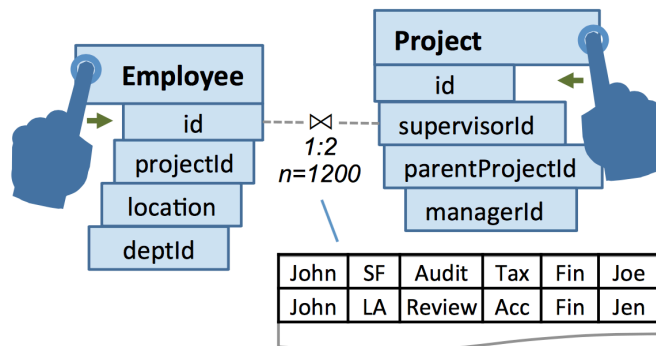Comm.:       Join between minds on ideas/concepts

# What can we do with RA?

Query optimization

Query by example

    Here is my data and examples, *generate the query*

Novel relationally complete interfaces



GestureDB.  Nandi et al.

# Summary

Relational Algebra (RA) operators

Operators are closed

    inputs & outputs are relations

Multiple Relational Algebra queries can be equivalent

    It is operational

    Same semantics but different performance

    Forms basis for optimizations

# Next Time

~~Relational Calculus~~

SQL