

# Introduction to Data Management

## SQL, Keys, and Joins







Alyssa Pittman

Based on slides by Jonathan Leang, Dan Suciu, et al

Paul G. Allen School of Computer Science and Engineering  
University of Washington, Seattle

# Announcements

- TA OH start next tomorrow

	<p>Andrew Guterman</p> <p>anguterm@cs</p> <p>OH: TBD</p>		<p>Jack Khuu</p> <p>jackkhuu@cs</p> <p>OH: Mon 5:30 - 6:30, CSE2 121</p>
	<p>Matthew Liu</p> <p>liux44@cs</p> <p>OH: TBD</p>		<p>Phong Ly</p> <p>pdly@cs</p> <p>OH: Tue 12:00-1:00pm, CSE2 151</p>
	<p>Samuel Oh</p> <p>mnkloip@cs</p> <p>OH: TBD</p>		<p>Khang Phan</p> <p>kphan000@cs</p> <p>OH: Mon 2:30-3:30 CSE2 151</p>

- HW1 released on the website
  - Let us know if you can't access your git repository!
  - Ask questions on Piazza

# Recap – The Relational Model

- Flat tables, static and typed attributes, etc.
  - “It’s a spreadsheet with rules”

**Table/  
Relation**

**Columns/Attributes/Fields**

**Rows/  
Tuples/  
Records**

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

# Recap – The Relational Model

- **Set semantics**

- No duplicate tuples

- Attributes are **typed** and **static**

- INTEGER, FLOAT, VARCHAR(n), DATETIME, ...

- Tables are **flat**

# Recap – SQL and RA

## ■ SQL

(Next several lectures)

- “What data do I want”

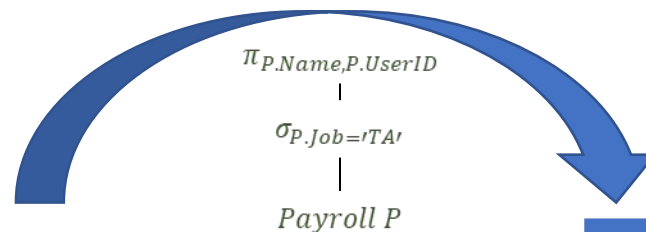
## ■ RA

(After SQL)

- “How does the computer get the data”

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

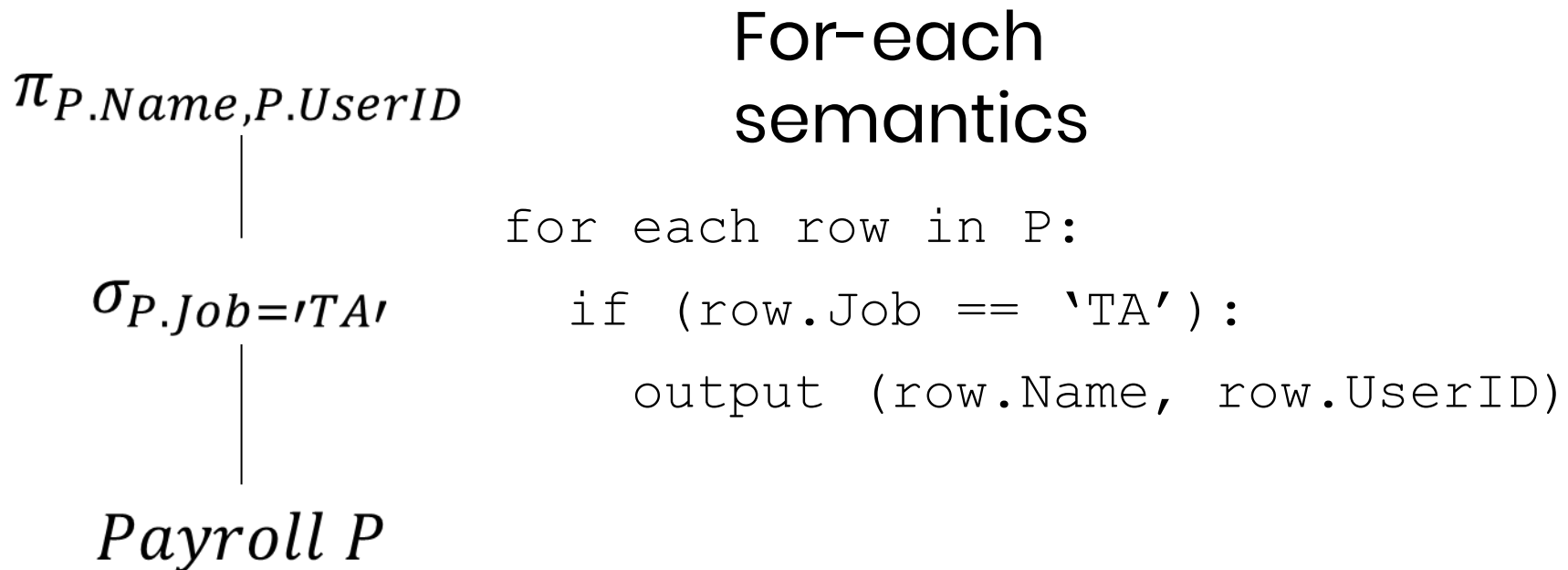
```
SELECT P.Name, P.UserID  
FROM Payroll AS P  
WHERE P.Job = 'TA';
```



Name	UserID
Jack	123
Allison	345

# Recap: For-each semantics

- Don't care about physical data layout or query plan
- But need to know the meaning of a query



# Goals for Today

- Last time we talked about the barebone building blocks of an RDBMS
  - Individual tables with no special properties
  - SQL and RA that work over individual tables
- Today is about semantics and relationships in the relational data model

# Outline

- Keys □ Identification
- Foreign Keys □ Relationships
- Joins in SQL
  - Inner joins
  - Outer joins
  - Self joins



# Keys

## Key

A **Key** is one or more attributes that uniquely identify a row.

## Payroll

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

# Keys

## Key

A **Key** is one or more attributes that uniquely identify a row.

Definitely not a key

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

# Keys

## Key

A **Key** is one or more attributes that uniquely identify a row.

Good candidate  
for a key

Definitely not a key

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

# Keys

## Key

A **Key** is one or more attributes that uniquely identify a row.

Is this a good candidate for a key?

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

# Keys

## Key

A **Key** is one or more attributes that uniquely identify a row.



Is this a good candidate for a key?

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

# Keys

## Key

A **Key** is one or more attributes that uniquely identify a row.



Is this a good candidate for a key?

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000
913	Peter	TA	60000

# Keys

## Key

A **Key** is one or more attributes that uniquely identify a row.

Data comes from the real world so models ought to reflect that

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000
913	Peter	TA	60000

```
CREATE TABLE Payroll (  
    UserID INT,  
    Name VARCHAR(100),  
    Job VARCHAR(100),  
    Salary INT);
```

Payroll(UserId, Name, Job, Salary)



# Keys

Unique Identifier

```
CREATE TABLE Payroll (  
  UserID INT,  
  Name VARCHAR(100),  
  Job VARCHAR(100),  
  Salary INT);
```

Payroll(UserId, Name, Job, Salary)

# Keys

Unique Identifier

```
CREATE TABLE Payroll (  
  UserID INT PRIMARY KEY,  
  Name VARCHAR(100),  
  Job VARCHAR(100),  
  Salary INT);
```

Payroll(UserId, Name, Job, Salary)

# Keys

Unique Identifier

```
CREATE TABLE Payroll (  
  UserID INT,  
  Name VARCHAR(100),  
  Job VARCHAR(100),  
  Salary INT);
```

Payroll(UserId, Name, Job, Salary)

# Keys

Unique Identifier

```
CREATE TABLE Payroll (  
  UserID INT,  
  Name VARCHAR(100),  
  Job VARCHAR(100),  
  Salary INT,  
  PRIMARY KEY (UserId, Name));
```

Payroll(UserId, Name, Job, Salary)

# Foreign Keys

- Databases can hold multiple tables
- How do we capture relationships *between* tables?

**Payroll**

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

**Regist**

UserID	Car
123	Charger
567	Civic
567	Pinto

# Foreign Keys

- Databases can hold multiple tables
- How do we capture relationships *between* tables?

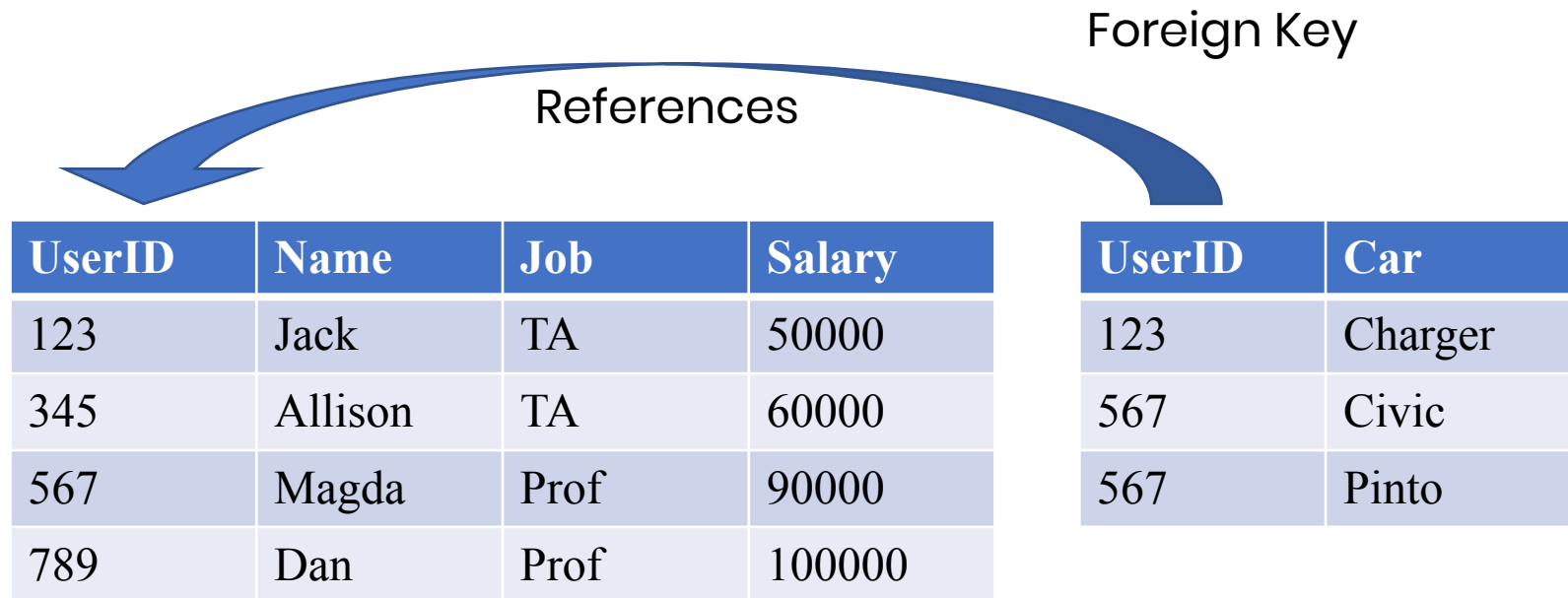
Foreign Key

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

# Foreign Keys

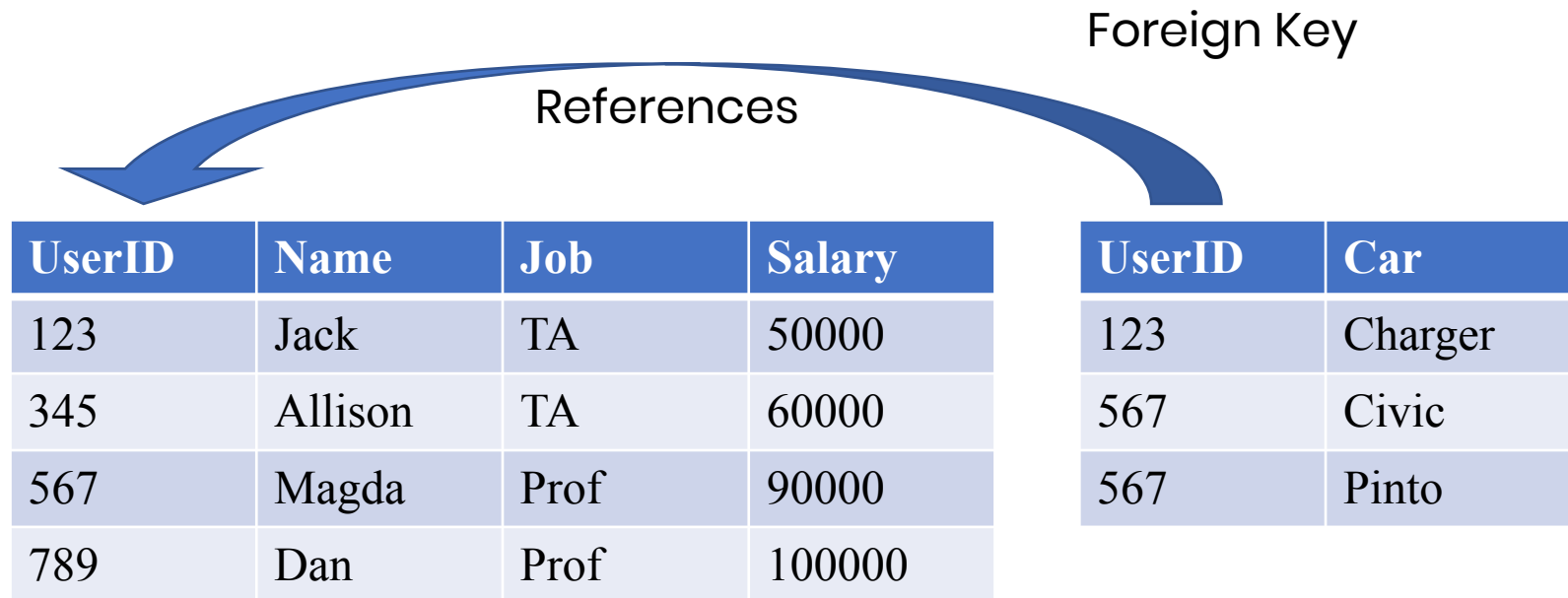
- Databases can hold multiple tables
- How do we capture relationships *between* tables?



# Foreign Keys

## Foreign Key

A **Foreign Key** is one or more attributes that uniquely identify a row in *another table*.





# Foreign Keys

## Foreign Key

A **Foreign Key** is one or more attributes that uniquely identify a row in *another table*.

Is this valid?

References

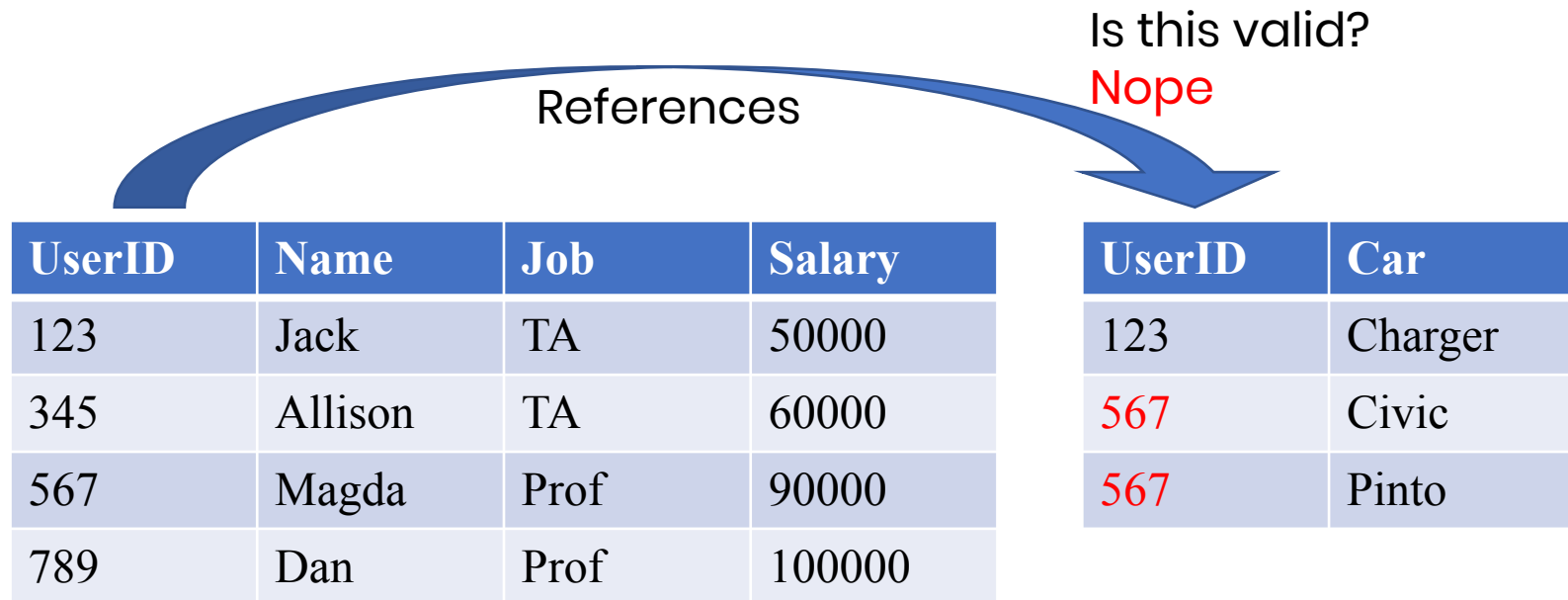
UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

# Foreign Keys

## Foreign Key

A **Foreign Key** is one or more attributes that uniquely identify a row in *another table*.



# Foreign Keys

```
CREATE TABLE Payroll (  
    UserID INT PRIMARY KEY,  
    Name VARCHAR(100),  
    Job VARCHAR(100),  
    Salary INT);
```

Payroll(UserId, Name, Job, Salary)

```
CREATE TABLE Regist (  
    UserID INT,  
    Car VARCHAR(100));
```

Regist(UserId, Car)

# Foreign Keys

```
CREATE TABLE Payroll (  
    UserID INT PRIMARY KEY,  
    Name VARCHAR(100),  
    Job VARCHAR(100),  
    Salary INT);
```

Payroll(UserId, Name, Job, Salary)

```
CREATE TABLE Regist (  
    UserID INT  
        REFERENCES Payroll(UserID),  
    Car VARCHAR(100));
```

Regist(UserId, Car)

# The Relational Model Revisited

- More complete overview of the Relational Model:
  - Database □ collection of tables
  - All tables are flat
  - Keys uniquely ID rows
  - Foreign keys act as a “semantic pointer”
  - **Physical data independence**

# Joins

- Foreign keys are able to *describe* a relationship between tables
- Joins are able to *realize* combinations of data

# Inner Joins

- Bread and butter of SQL queries
  - “Inner join” is often interchangeable with just “join”

# Nested-Loop Semantics

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

```
SELECT P.Name, R.Car
FROM Payroll AS P JOIN Regist AS R
ON P.UserID = R.UserID;
```

How do we  
algorithmically  
get our results?

Name	Car
Jack	Charger
Magda	Civic
Magda	Pinto



# Nested-Loop Semantics


UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

```
SELECT P.Name, R.Car
FROM Payroll AS P JOIN Regist AS R
ON P.UserID = R.UserID;
```


```
for each row1 in Payroll:
    for each row2 in Regist:
        if (row1.UserID = row2.UserID):
            output (row1.Name, row2.Car)
```

# Nested-Loop Semantics



UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000


UserID	Car
123	Charger
567	Civic
567	Pinto



Name	Car
------	-----


```
for each row1 in Payroll:
    for each row2 in Regist:
        if (row1.UserID = row2.UserID):
            output (row1.Name, row2.Car)
```

# Nested-Loop Semantics



UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000


UserID	Car
123	Charger
567	Civic
567	Pinto



Name	Car
Jack	Charger


```
for each row1 in Payroll:
    for each row2 in Regist:
        if (row1.UserID = row2.UserID):
            output (row1.Name, row2.Car)
```

# Nested-Loop Semantics



UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000


UserID	Car
123	Charger
567	Civic
567	Pinto



Name	Car
Jack	Charger


```
for each row1 in Payroll:
    for each row2 in Regist:
        if (row1.UserID = row2.UserID):
            output (row1.Name, row2.Car)
```

# Nested-Loop Semantics



UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto



Name	Car
Jack	Charger

```
for each row1 in Payroll:
    for each row2 in Regist:
        if (row1.UserID = row2.UserID):
            output (row1.Name, row2.Car)
```

# Nested-Loop Semantics


UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

Name	Car
Jack	Charger


```
for each row1 in Payroll:
    for each row2 in Regist:
        if (row1.UserID = row2.UserID):
            output (row1.Name, row2.Car)
```

# Nested-Loop Semantics



UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto



Name	Car
Jack	Charger

```
for each row1 in Payroll:
    for each row2 in Regist:
        if (row1.UserID = row2.UserID):
            output (row1.Name, row2.Car)
```

# Nested-Loop Semantics

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

Name	Car
Jack	Charger

```
for each row1 in Payroll:
  for each row2 in Regist:
    if (row1.UserID = row2.UserID):
      output (row1.Name, row2.Car)
```



# Nested-Loop Semantics

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

Name	Car
Jack	Charger

```
for each row1 in Payroll:
    for each row2 in Regist:
        if (row1.UserID = row2.UserID):
            output (row1.Name, row2.Car)
```

# Nested-Loop Semantics

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

Name	Car
Jack	Charger

```
for each row1 in Payroll:
  for each row2 in Regist:
    if (row1.UserID = row2.UserID):
      output (row1.Name, row2.Car)
```

# Nested-Loop Semantics

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000


UserID	Car
123	Charger
567	Civic
567	Pinto

Name	Car
Jack	Charger
Magda	Civic


```
for each row1 in Payroll:
  for each row2 in Regist:
    if (row1.UserID = row2.UserID):
      output (row1.Name, row2.Car)
```

# Nested-Loop Semantics

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000



UserID	Car
123	Charger
567	Civic
567	Pinto



Name	Car
Jack	Charger
Magda	Civic

```
for each row1 in Payroll:
    for each row2 in Regist:
        if (row1.UserID = row2.UserID):
            output (row1.Name, row2.Car)
```

# Nested-Loop Semantics

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

Name	Car
Jack	Charger
Magda	Civic
Magda	Pinto

```
for each row1 in Payroll:
    for each row2 in Regist:
        if (row1.UserID = row2.UserID):
            output (row1.Name, row2.Car)
```

# Nested-Loop Semantics

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

Name	Car
Jack	Charger
Magda	Civic
Magda	Pinto

```
for each row1 in Payroll:
    for each row2 in Regist:
        if (row1.UserID = row2.UserID):
            output (row1.Name, row2.Car)
```

# Nested-Loop Semantics

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000


UserID	Car
123	Charger
567	Civic
567	Pinto

Name	Car
Jack	Charger
Magda	Civic
Magda	Pinto


```
for each row1 in Payroll:
  for each row2 in Regist:
    if (row1.UserID = row2.UserID):
      output (row1.Name, row2.Car)
```

# Nested-Loop Semantics

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000



UserID	Car
123	Charger
567	Civic
567	Pinto



Name	Car
Jack	Charger
Magda	Civic
Magda	Pinto

```
for each row1 in Payroll:
    for each row2 in Regist:
        if (row1.UserID = row2.UserID):
            output (row1.Name, row2.Car)
```



# Nested-Loop Semantics

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

Name	Car
Jack	Charger
Magda	Civic
Magda	Pinto

```
for each row1 in Payroll:
    for each row2 in Regist:
        if (row1.UserID = row2.UserID):
            output (row1.Name, row2.Car)
```

# Inner Joins

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

Explicit

```
SELECT P.Name, R.Car
FROM Payroll AS P JOIN Regist AS R
ON P.UserID = R.UserID;
```

Implicit

```
SELECT P.Name, R.Car
FROM Payroll AS P, Regist AS R
WHERE P.UserID = R.UserID;
```

# Outer Joins

Now I want to include everyone, even if they don't drive.

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

# Outer Joins

Now I want to include everyone, even if they don't drive.

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

```
SELECT P.Name, R.Car
FROM Payroll AS P LEFT OUTER JOIN Regist AS R
ON P.UserID = R.UserID;
```

# Outer Joins

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

Name	Car
Jack	Charger
Allison	NULL
Magda	Civic
Magda	Pinto
Dan	NULL

# Outer Joins

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

Name	Car
Jack	Charger
Allison	NULL
Magda	Civic
Magda	Pinto
Dan	NULL

NULL is a value placeholder. Depending on context, it may mean unknown, not applicable, etc.

# Outer Joins

- LEFT OUTER JOIN
  - All rows in left table are preserved
- RIGHT OUTER JOIN
  - All rows in right table are preserved
- FULL OUTER JOIN
  - All rows are preserved

# Self Joins

Find all people who drive a Civic and Pinto

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

```
SELECT P.Name, R.Car
FROM Payroll AS P, Regist AS R
WHERE P.UserID = R.UserID AND
        R.Car = 'Civic';
```



# Self Joins

Find all people who drive a Civic and Pinto

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

```
SELECT P.Name, R.Car
FROM Payroll AS P, Regist AS R
WHERE P.UserID = R.UserID AND
      R.Car = 'Civic' AND
      R.Car = 'Pinto';
```

Will this work?

# Self Joins

Find all people who drive a Civic and Pinto

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

```
SELECT P.Name, R.Car
FROM Payroll AS P, Regist AS R
WHERE P.UserID = R.UserID AND
      R.Car = 'Civic' AND
      R.Car = 'Pinto';
```

Will this work?  
Nope, empty set  
is returned

# Self Joins

Find all people who drive a Civic and Pinto

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

```
SELECT P.Name, R.Car
FROM Payroll AS P, Regist AS R
WHERE P.UserID = R.UserID AND
      R.Car = 'Civic' OR
      R.Car = 'Pinto';
```

Will this work?

# Self Joins

Find all people who drive a Civic and Pinto

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto
789	Civic

```
SELECT P.Name, R.Car
FROM Payroll AS P, Regist AS R
WHERE P.UserID = R.UserID AND
      R.Car = 'Civic' OR
      R.Car = 'Pinto';
```

Will this work?

Nope, returns  
people who had  
just one type of  
car

# Self Joins

Find all people who drive a Civic and Pinto

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

```
SELECT P.Name, R1.Car
FROM Payroll AS P, Regist AS R1, Regist AS R2
WHERE P.UserID = R1.UserID AND
      P.UserID = R2.UserID AND
      R1.Car = 'Civic' AND
      R2.Car = 'Pinto';
```

# Self Joins

Find all people who drive a Civic and Pinto

UserID	Name	Job	Salary
123	Jack	TA	50000
345	Allison	TA	60000
567	Magda	Prof	90000
789	Dan	Prof	100000

UserID	Car
123	Charger
567	Civic
567	Pinto

All pairs of cars a person can drive

```
SELECT P.Name, R1.Car
FROM Payroll AS P, Regist AS R1, Regist AS R2
WHERE P.UserID = R1.UserID AND
      P.UserID = R2.UserID AND
      R1.Car = 'Civic' AND
      R2.Car = 'Pinto';
```

# A little extra SQL

- ORDER BY – Orders result tuples by specified attributes (default ascending)

```
SELECT P.Name, P.UserID  
      FROM Payroll AS P  
      WHERE P.Job = 'TA'  
      ORDER BY P.Salary, P.Name;
```

- DISTINCT – Deduplicates result tuples

```
SELECT DISTINCT P.Job  
      FROM Payroll AS P  
      WHERE P.Salary > 70000;
```

# Takeaways

- We can describe relationships between tables with keys and foreign keys
- Different joining techniques can be used to achieve particular goals
- Our SQL toolbox is growing!
  - Not just reading and filtering data anymore
  - Starting to answer complex questions