

Data Management for Data Science

SQL Basics

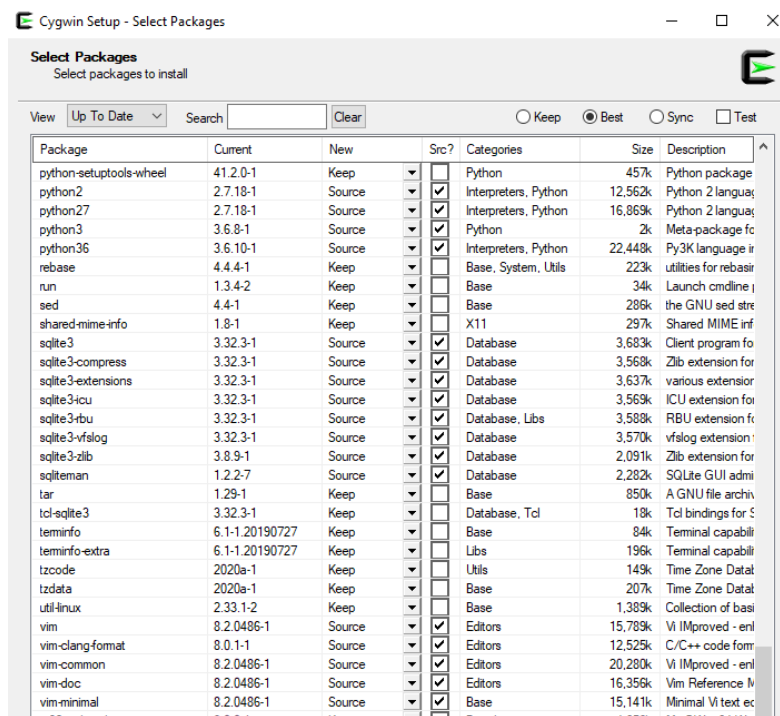
Paul G. Allen School of Computer Science and Engineering
University of Washington, Seattle

Announcements

- HW 1 released – due January 1/15 at **11:59 pm**
 - Collected via Gradescope
 - Try to do HW 1 setup today (should take ~5-10 minutes)

Announcements

- Note for Windows users:
 - Check boxes for sqlite3 in Cygwin



- [Windows Terminal](#) also recommended

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**

The diagram illustrates the components of a relational table. A blue arrow points from the label 'Table/Relation' to the entire table structure. Three blue arrows point from the label 'Columns/Attributes/Fields' to the column headers: 'UserID', 'Name', and 'Job'. A blue bracket on the left side of the table, labeled 'Rows/Tuples/Records', encompasses the four data rows.

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

Structured Query Language - SQL

- **Declarative** query language
 - Tell the computer what you want, not how to get it
- Languages like Java/Python are procedural
- Declarative query language allows **physical data independence**

Recap - Basic SQL query

Payroll

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



```
SELECT *  
FROM Payroll;
```

Recap - Basic SQL query

Payroll

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



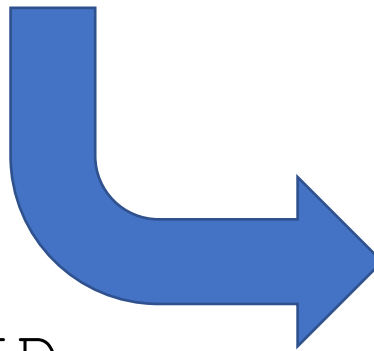
```
SELECT *  
FROM Payroll;
```

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

Recap - Basic SQL query

Payroll

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



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


Hello World

Payroll

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

SELECT
What kind of data
I want

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

Hello World

Payroll

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

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

SELECT
What kind of data
I want

FROM
Where the data
coming from

Hello World

Payroll

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

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

SELECT
What kind of data
I want

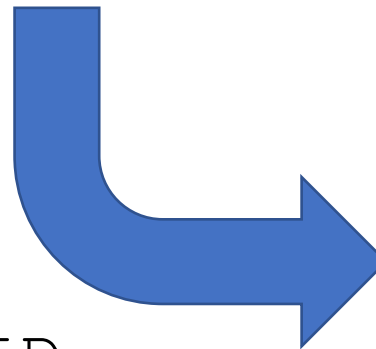
FROM
Where the data
coming from

WHERE
Filter the data

Recap - Basic SQL query

Payroll

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



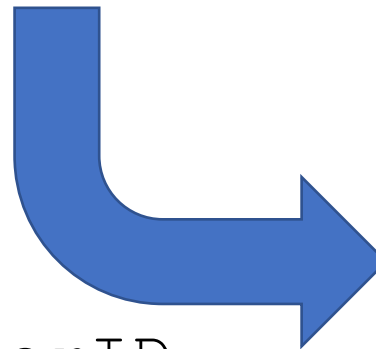
```
SELECT Name, UserID
FROM Payroll
WHERE Job = 'TA';
```

Name	UserID
Jack	123
Allison	345

Recap - Basic SQL query

Payroll

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

“Payroll AS P” makes P an alias.
This lets us specify that the
attributes come from Payroll

Wait!

What actually
happens when we
execute the SQL
query?

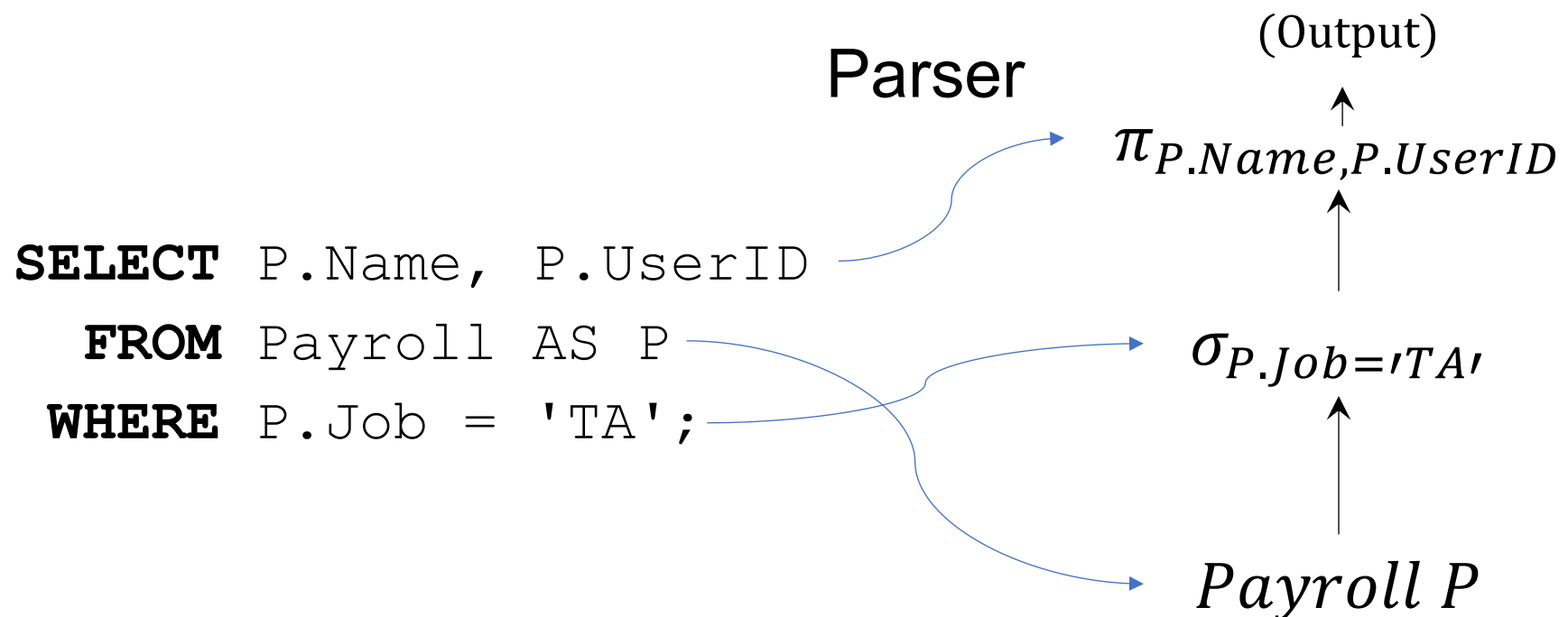
Database Internals

- Code has to boil down to instructions at some point
- Relational Database Management Systems (RDBMSs) use **Relational Algebra** (RA)

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

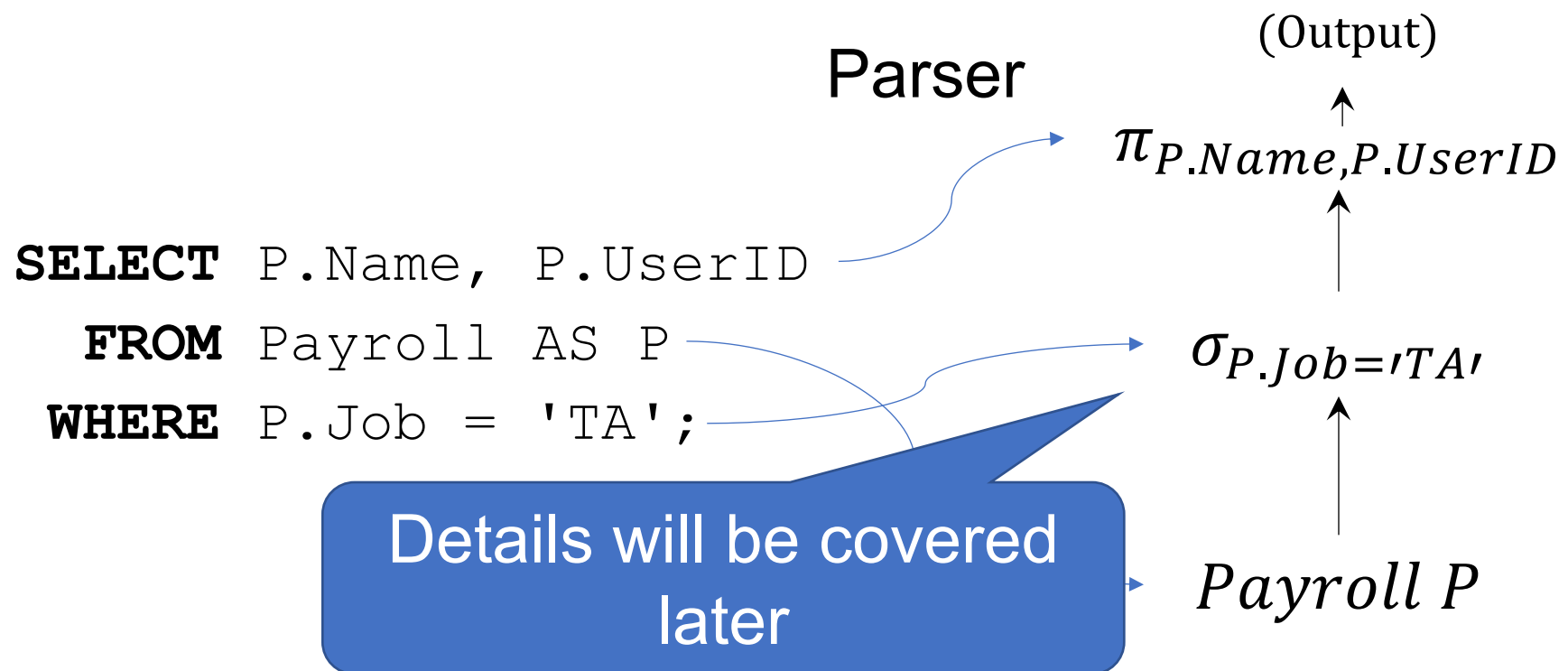
Database Internals

- Code has to boil down to instructions at some point
- Relational Database Management Systems (RDBMSs) use **Relational Algebra** (RA).



Database Internals

- Code has to boil down to instructions at some point
- Relational Database Management Systems (RDBMSs) use **Relational Algebra** (RA).

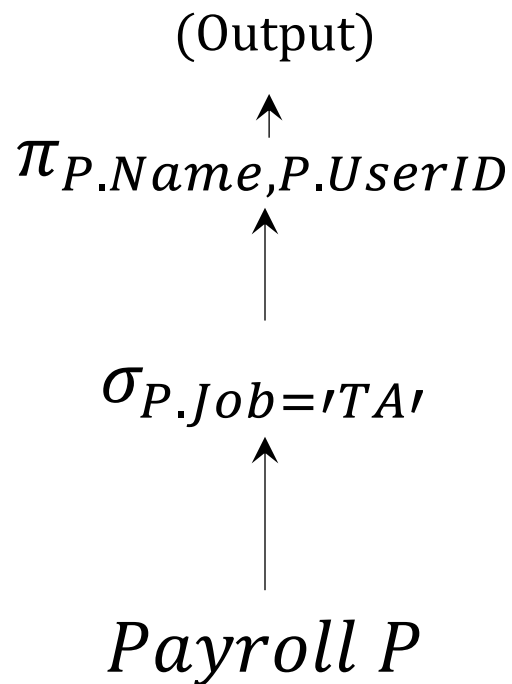


Database Internals

- It's important to define the semantics (meaning) of a query

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

For-each semantics



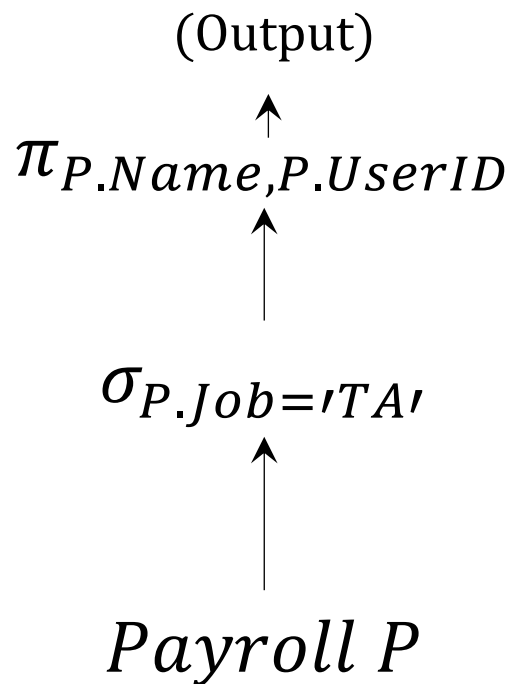
Database Internals

- It's important to define the semantics (meaning) of a query

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

For-each semantics

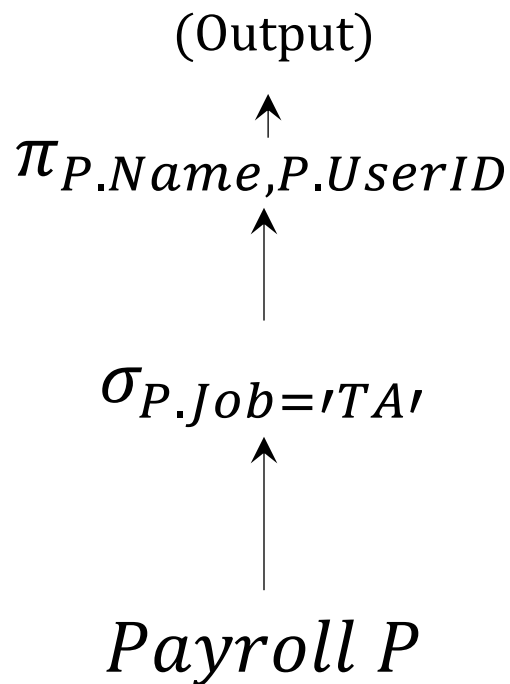
```
for each row in P:  
    if (row.Job == 'TA'):  
        output (row.Name, row.UserID)
```



Database Internals

- It's important to define the semantics (meaning) of a query

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



Tuples “flow” up the query plan, getting filtered and modified

For-each semantics

Payroll

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

```
for each row in P:  
    if (row.Job == 'TA'):  
        output (row.Name,  
                row.UserID)
```

Job == 'TA'?

Name	UserID
------	--------

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

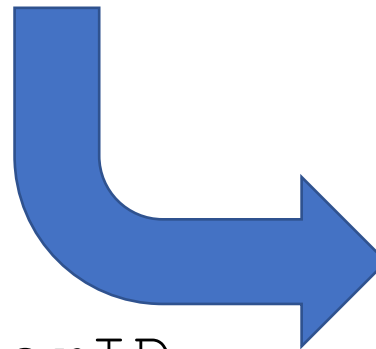
For-each semantics

Payroll

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

```
for each row in P:  
    if (row.Job == 'TA'):  
        output (row.Name,  
                row.UserID)
```

Job == 'TA'?



Name	UserID
Jack	123

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

For-each semantics

Payroll

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

```
for each row in P:  
    if (row.Job == 'TA'):  
        output (row.Name,  
                row.UserID)
```

Job == 'TA'?



Name	UserID
Jack	123

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

For-each semantics


Payroll

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

```
for each row in P:  
    if (row.Job == 'TA'):  
        output (row.Name,  
                row.UserID)
```

Job == 'TA'?

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



Name	UserID
Jack	123
Allison	345

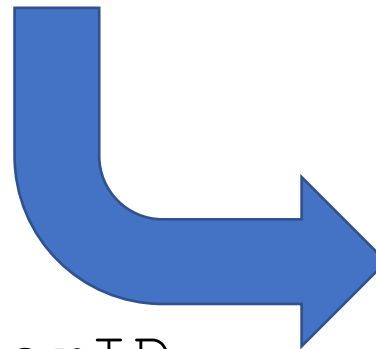
For-each semantics

Payroll

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

```
for each row in P:  
    if (row.Job == 'TA'):  
        output (row.Name,  
                row.UserID)
```

Job == 'TA'?



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

Name	UserID
Jack	123
Allison	345

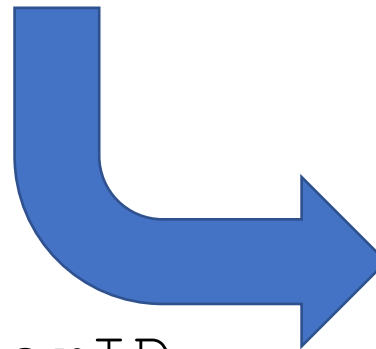
For-each semantics

Payroll

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

```
for each row in P:  
    if (row.Job == 'TA'):  
        output (row.Name,  
                row.UserID)
```

Job == 'TA'?



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

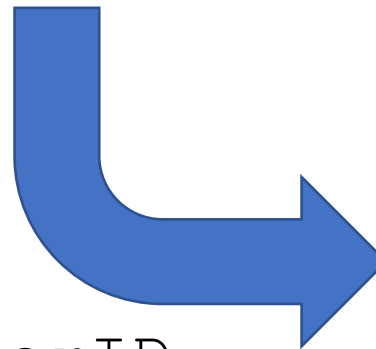
Name	UserID
Jack	123
Allison	345

For-each semantics

Payroll

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

```
for each row in P:  
    if (row.Job == 'TA'):  
        output (row.Name,  
                row.UserID)
```



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

Name	UserID
Jack	123
Allison	345

Recap – SQL and RA

■ SQL

(Next few lectures)

- “What data do I want”

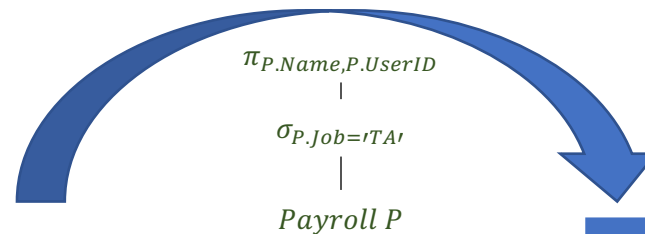
■ RA

(After SQL)

- “How do I 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

What's Next?

- Creating tables
- Keys → Identification
- Foreign Keys → Relationships
- Joins in SQL and RA
 - Inner joins
 - Outer joins
 - Self joins

Create Table Statement

Payroll(UserId, Name, Job, Salary)



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

Data types

- Each attribute has a type
 - Examples types:
 - Strings: CHAR(20), VARCHAR(50), TEXT
 - Numbers: INT, SMALLINT, FLOAT
 - MONEY, DATETIME, ...
 - Few more that are DBMS specific
 - Statically and strictly enforced

Data types

- Generally you will use:
 - **VARCHAR(N)** for strings where **N** is the maximum character length
 - Generally set this to as large as you need, like 256 or 1000.
 - **INT**, **FLOAT** for numbers (**INTEGER** works in SQLite)
 - **DATETIME** for dates
 - Can use VARCHAR(N) in SQLite

Create Table Statement

Payroll(**UserId**, Name, Job, Salary)



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

Create Table Statement

Payroll(**UserId**, Name, Job, Salary)



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

Everything is case-insensitive, but having your own guidelines is useful for readability

Keys

Key

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

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

Keys

```
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

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

Can also define the
PK on a new line

Payroll(UserId, Name, Job, Salary)

Keys of more than one attribute

Sometimes no single attribute is unique, but combinations of attributes are a unique key for the table.

Must use the PK definition on a new line for multi-attribute keys

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

Keys of more than one attribute

Sometimes no single attribute is unique, but combinations of attributes are a unique key for the table.

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

Must use the PK definition on a new line for multi-attribute keys

Here the combination of Name and Job are unique
e.g. only one "Ryan, Professor"
but some "Ryan, CEO" or "Mary, Professor" also exist

Payroll(Name, Job, Salary)

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;
```

(inverse **ORDER BY** P.Salary DESC)

- DISTINCT – Deduplicates result tuples

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


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?

Payroll

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

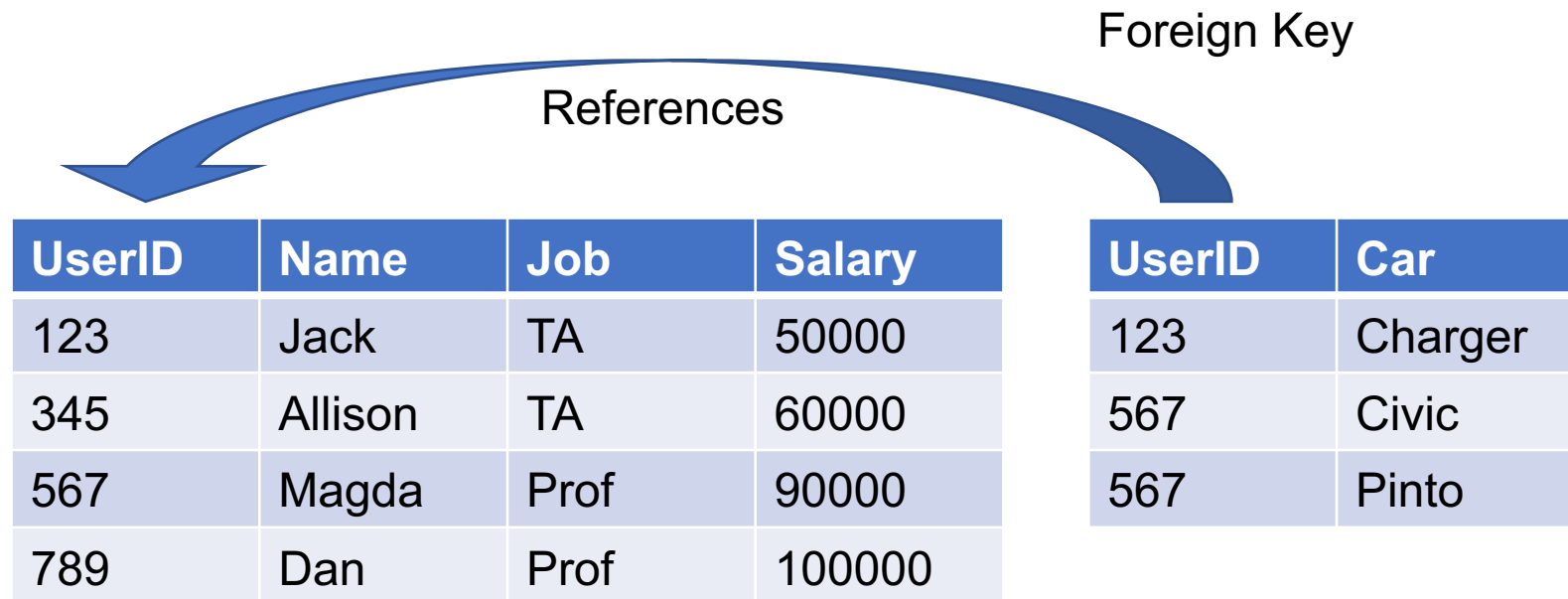
Foreign Key
UserID

Regist

UserID	Car
123	Charger
567	Civic
567	Pinto

Foreign Keys

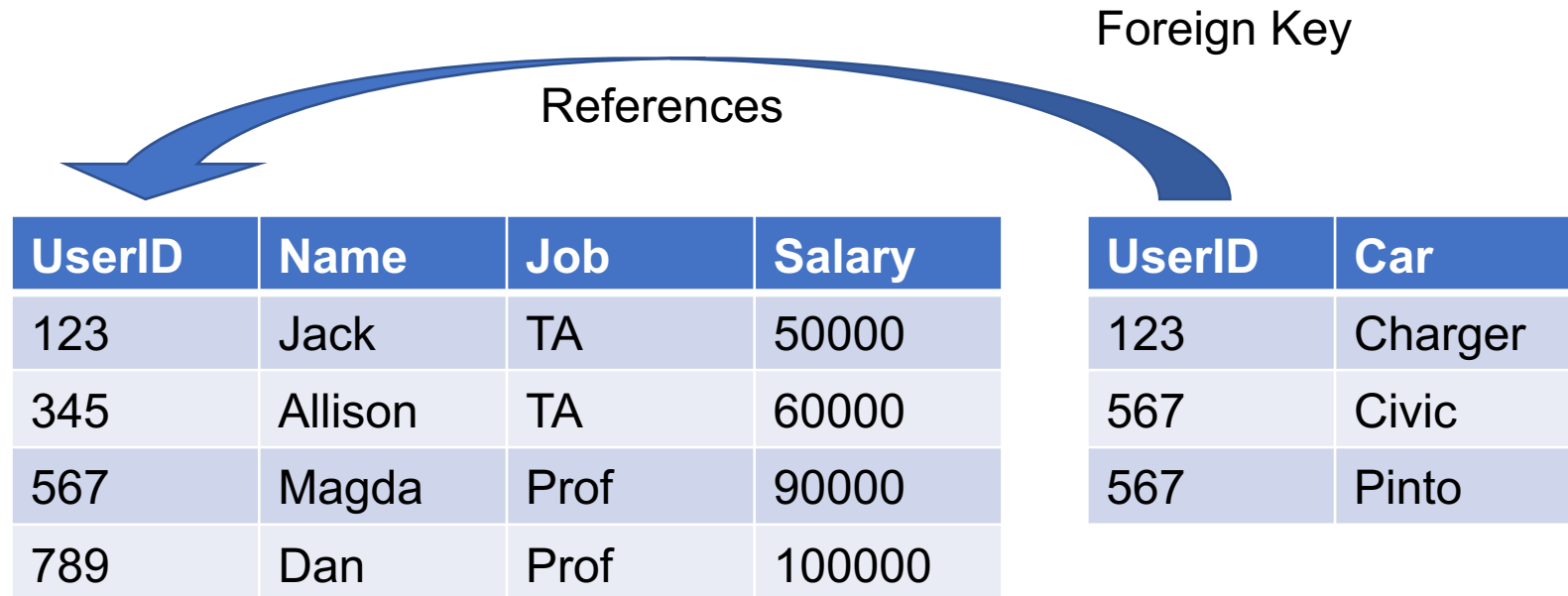
- Databases can hold multiple tables
- How do we capture relationships *between* tables?



Foreign Keys

Foreign Key

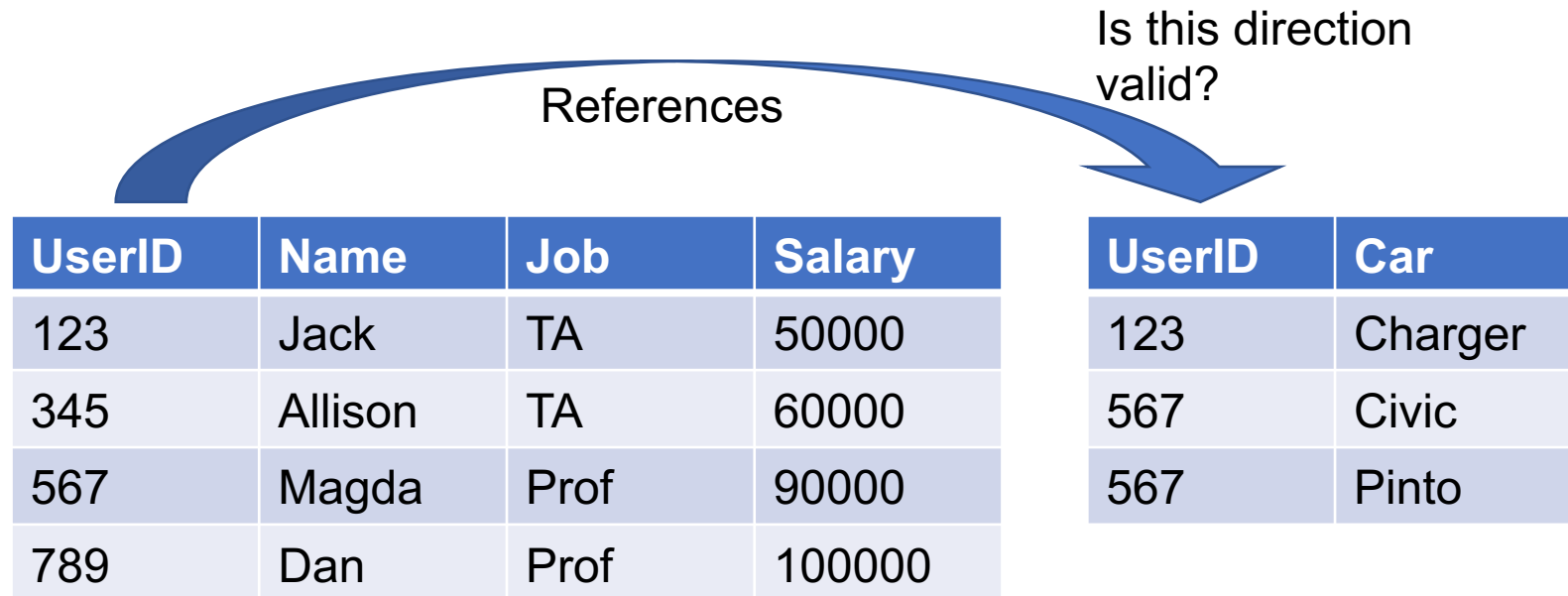
A **Foreign Key** is one or more attributes that uniquely identify a row in *another table*.



Foreign Keys

Foreign Key

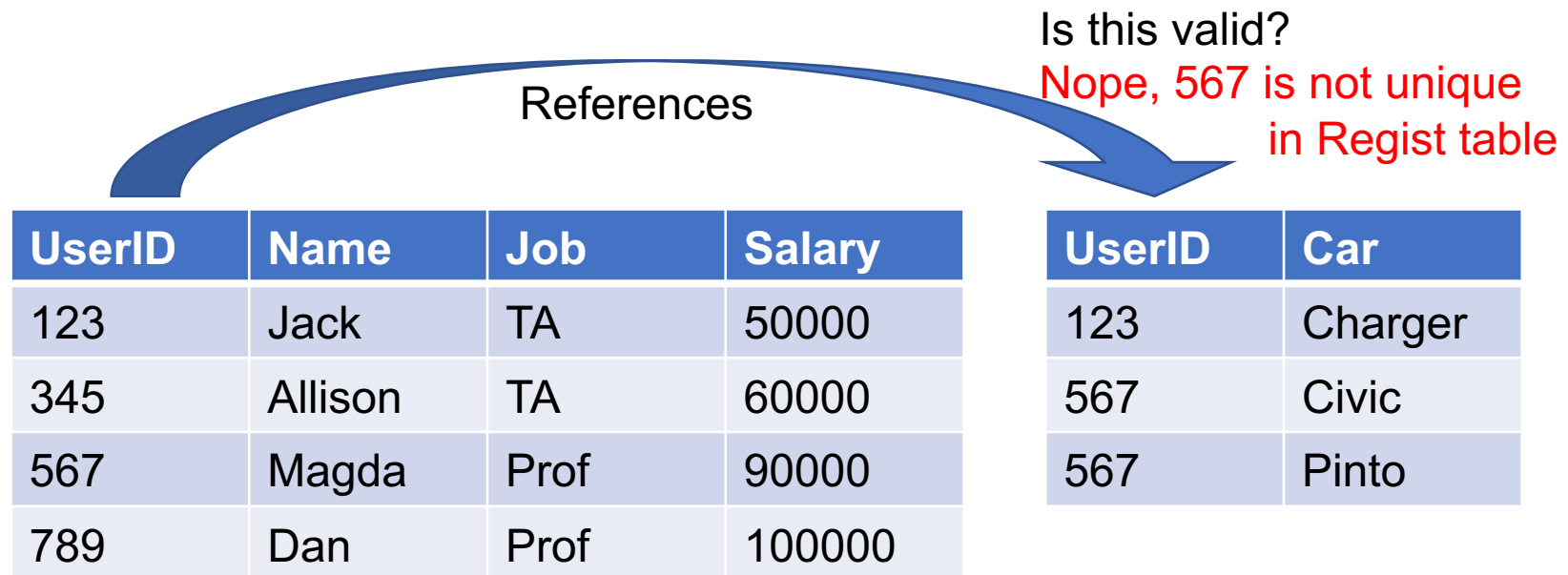
A **Foreign Key** is one or more attributes that uniquely identify a row in *another table*.



Foreign Keys

Foreign Key

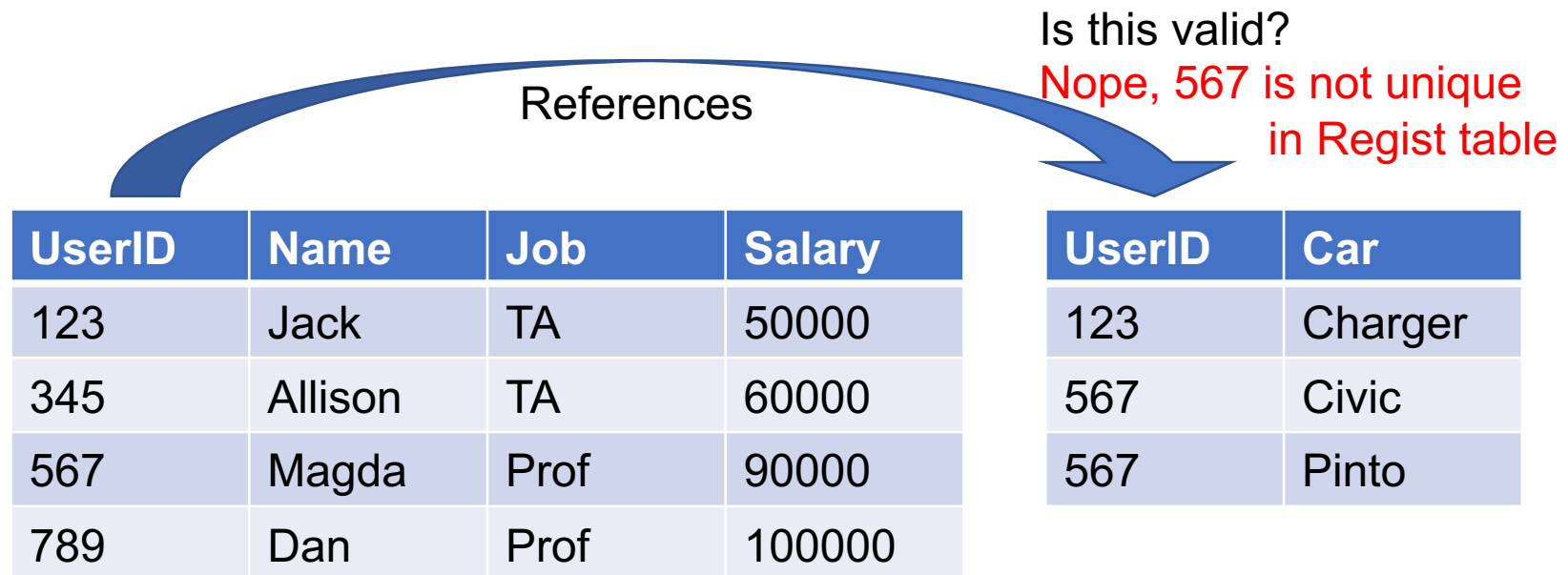
A **Foreign Key** is one or more attributes that uniquely identify a row in *another table*.



Foreign Keys

Foreign Key

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



Foreign keys must reference (point to) a unique attribute, almost always a primary key

Foreign Keys

We add foreign key declaration in the same way as a primary key.

```
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

We add foreign key declaration in the same way as a primary key.

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

Payroll(UserId, Name, Job, Salary)

Regist(UserId, Car)

Foreign Keys

We add foreign key declaration in the same way as a primary key.

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

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

or, when attribute name is the same:

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

Payroll(UserId, Name, Job, Salary)

Regist(UserId, Car)

Foreign Keys

Can also put foreign key declaration on a new line, need to do this for multiple attributes

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

Payroll(UserID, Name, Job, Salary)

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

Regist(UserID, Name, 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

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