

SQL - Basics

Contents

- What is Data?
- Relational Databases
- SQL - A Brief Intro
- Advantages and Disadvantages of SQL
- Entity Relationship Diagram
- SQL Statements
- Read Tables: The SELECT Query
- Set rules to our SELECT Queries

What is Data

First of all, what is data? This course covers Data Science after all, we need to know a brief definition of data

- **Data** can be thought of as recorded measurements of something in the real world. This *something* is a unit of observation.
 - *For example, a list of people's height is data.*
 - *A person would be the unit of observation (**sample**)*
- Data can describe a vast amount of different things.
 - *For example, there is a lot of data we can use to describe a person.*
 - *Each measurement is called a variable (**feature**)*
 - *Each observation (the height of a person) is a **data point***
- Several data points together create a **Dataset**

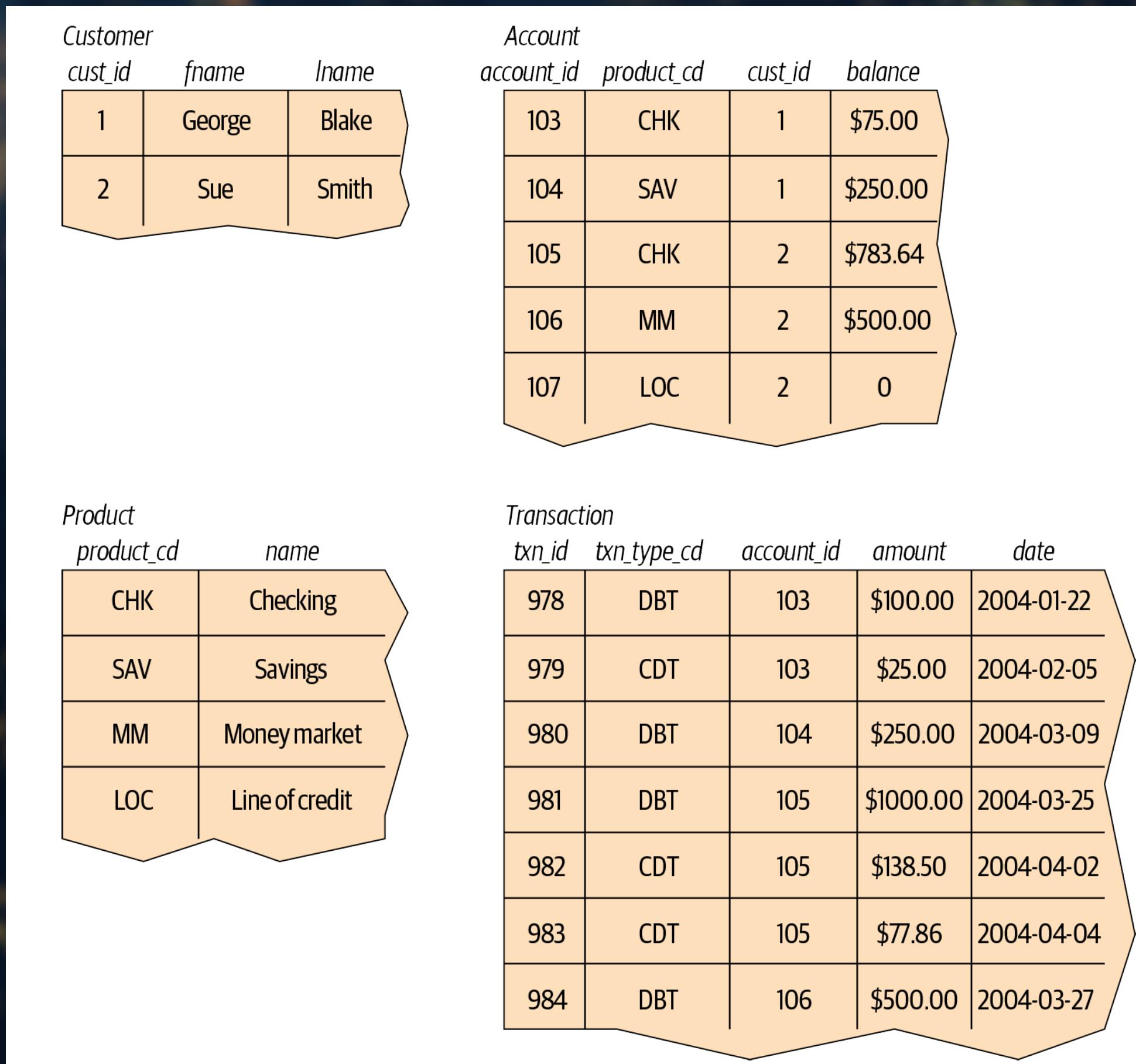
Relational Database

Data is useful to obtain valuable information. We could analyze now that data by hand, but of course, computers can analyze much more information than humans. There are two tools for storing, organizing, and processing data in a computer. The first one is:

- Relational Database:
 - Utilizes the relation model of data: **Data corresponding to the same ID is used to link records in different tables**
 - Data is organized as relations, containing a **relation key (ID)**
 - Relations are usually implemented as **Tables**
 - Tables are assimilated in common collections in databases called **Schemas**
 - For example: One Table contains (ID, Name, Last Name, Age)
 - Another Table contains (ID, Height, Smoker?)
 - And another Table contains (Smoker?, Cancer Development)
 - The software used to manage relational databases is referred to as a **Relational DataBase Management System (RDBMS)**

Relational Database

- Data corresponding to the same ID is used to link records in different tables***



Term	Definition
Entity	Something of interest to the database user community. Examples include customers, parts, geographic locations, etc.
Column	An individual piece of data stored in a table.
Row	A set of columns that together completely describe an entity or some action on an entity. Also called a record.
Table	A set of rows, held either in memory (nonpersistent) or on permanent storage (persistent).
Result set	Another name for a nonpersistent table, generally the result of an SQL query.
Primary key	One or more columns that can be used as a unique identifier for each row in a table.
Foreign key	One or more columns that can be used together to identify a single row in another table.

SQL - A brief intro

*"Along with Codd's definition of the relational model, he proposed a language called **DSL/Alpha** for manipulating the data in relational tables. Shortly after Codd's paper was released, IBM commissioned a group to build a prototype based on Codd's ideas. This group created a **simplified version of DSL/Alpha** that they called **SQUARE**. Refinements to SQUARE led to a language called SEQUEL, which was, finally, **shortened to SQL**. While SQL began as a language used to manipulate data in relational databases, it has evolved [...] to be a **language for manipulating data across various database technologies** [...]*

One final note: SQL is not an acronym for anything (although many people will insist it stands for "Structured Query Language"). When referring to the language, it is equally acceptable to say the letters individually (i.e., S. Q. L.) or to use the word sequel."

Extracted from 'Learning SQL' by Alan Beaulieu

SQL - Pros and Cons

SQL Databases provide a ton of **Advantages** that may it the de facto choice for many applications

- **Intuitive:** Relations that almost anyone can understand
- **Efficient:** They use normalization so it doesn't repeat its representation (less space)
- **Declarative:** You tell the data that you want, and the system takes care of how to execute the query
- **Robust:** Most databases have the ACID compliance (Atomicity, Consistency, Isolation, and Durability), guaranteeing the validity of data even if the hardware fails

SQL - Pros and Cons

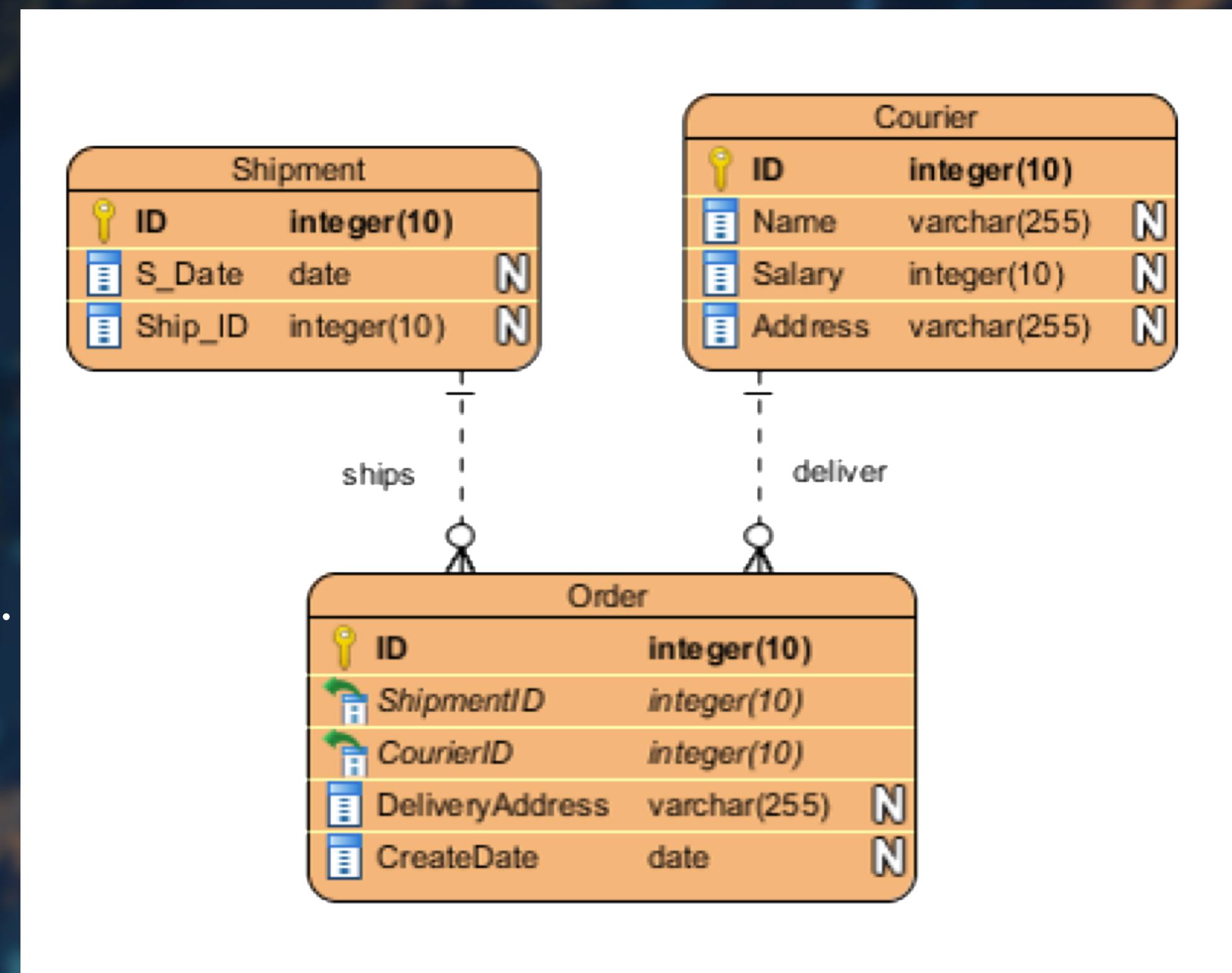
However, we might find some **Downsides** when working with SQL databases:

- **Lower specificity**: Sometimes, SQL's functionality is limited to what it has been programmed to do. Not common though, since it management systems get updated
- **Limited Scalability**: Due to the robustness might be an impediment for scaling data.
- **Object-relation mismatch impedance**: Sometimes objects have attributes with many-to-many relationships. For example, a costumer may own multiple products, but each product may have multiple objects

SQL - Entity Relationship Diagram

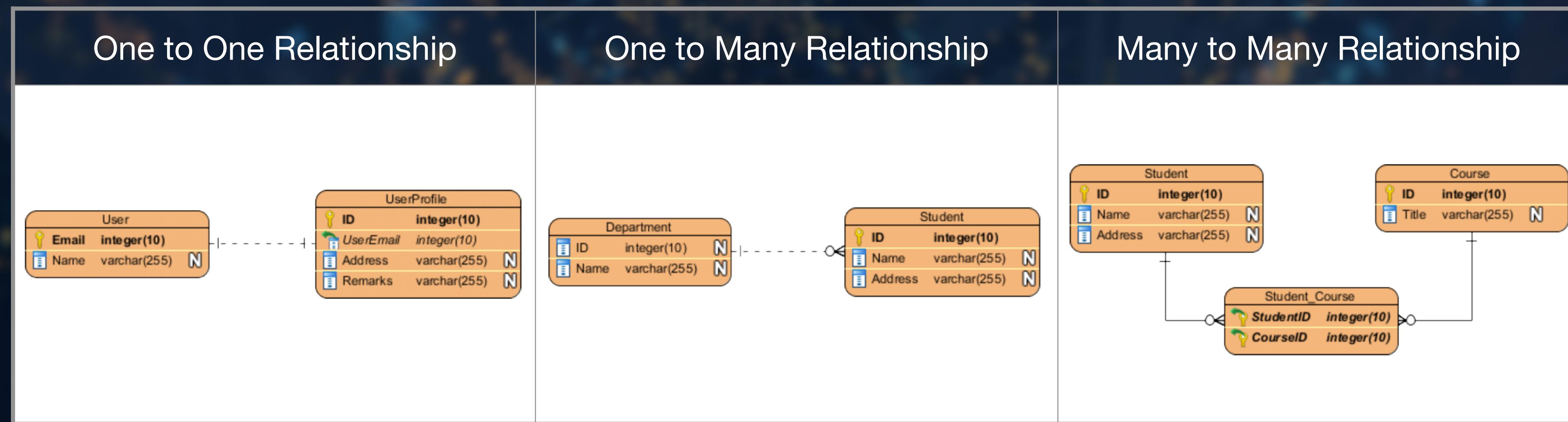
Entity Relationship Diagram (ERD) is a type of structural diagram for use in database design. An ERD contains different symbols and connectors that visualize two important information: **The major entities within the system scope**, and the **inter-relationships among these entities**.

- **Primary Keys:** uniquely defines a record in a database table.
 - There must not be two (or more) records that share the same value for the primary key attribute.
 - For example (*ID*)
- **Foreign Keys:** reference to a primary key in another table. It is used to identify the relationships between entities.
 - They don't need to be unique
 - For example (*ShipmentID*, *CourierID*)



SQL - Entity Relationship Diagram

- **Cardinality:** possible number of occurrences in one entity which is associated with the number of occurrences in another.
 - *For example, ONE team has MANY players. Then Team and Player are inter-connected with a one-to-many relationship.*



SQL - Statements

- Databases have four main statements, following the *CRUD (Create, Read, Update, Delete)* paradigm. Statements can essentially be thought of as functions we want to perform over the database. Within a given database, relevant statements could include:
 - **CREATE TABLE:** A statement which creates a new table in the database
 - **SELECT:** Allows you to read and *query* the database to extract records you want
 - **UPDATE:** Allows you to update records within a table
 - **DROP TABLE:** Gives you the power to remove a table from a database
- SQL Statements and Keywords are, by convention, UPPERCASE

SQL - Use Case

It's time to get our hands dirty. By now you should already have PostgreSQL installed. If not, open the SQL_Setup.md in the repo and follow the steps included in the file.

- We're going to be working with a database known as Pagila
 - PostgreSQL port of an open-source sample database known as [Sakila](#) (A sample database for learning purposes)
- Pagila is a database which models a DVD rental store. It features films, actors, film-actor relationships, and a central inventory table that connects films, stores, and rentals.

SQL - Use Case



SQL - SELECT Statement

The most common operation in a database is reading data from it. We can do it using the SELECT keyword.

Let's see the basic anatomy of a SELECT keyword:

- **Operation:** What is going to be done. SELECT followed by the names of columns (* indicates all columns)
- **Data:** FROM where we are getting the data (SELECT * FROM payments: select all columns from the payments entity). Data can also be obtained from the combination of existing columns using arithmetic operations
- **Post-Processing:** It takes the results of a query and sort them or limit them by using ORDER BY and LIMIT
- **Conditional:** Acts as a filter. Usually indicated by WHERE
- **Grouping:** Assemble the rows of a data source using a key created by a GROUP BY clause.

Two small notes:

- *Statements terminate with a semi-colon*
- *SQL ignores whitespace*

SQL - Post Processing

LIMIT

- It allows us to limit the number of records
- It's the last part part of the query

```
SELECT *  
FROM payment  
LIMIT 10;
```

ORDER BY

- It allows us to specify which columns we want our returned data ordered by
- Ascending by default. We can specify DESC
- We can also sort multiple columns, and the sorting occurs from left to right
- Placed before LIMIT

```
SELECT *  
FROM film  
ORDER BY rental_rate DESC, length DESC;  
LIMIT 10
```

SQL - Practicals (Part I)

Take a look at the ERD in the Pagila database and attempt the first practical ("Use the SELECT query to read the following data from the Pagila database")

SQL - Conditional

WHERE :

- It allows us to specify which columns we want our returned data ordered by
- Common condition symbols are:
 - >, >= Greater than, greater than or equal to
 - <, <= Less than, less than or equal to
 - = Equal to
 - != Not equal to

```
SELECT *
FROM payments
WHERE customer_id = 72;
```

SQL - Arithmetic Operations

- It's possible to return a new column from a combination of existing columns
- We derive/compute this new column using mathematical operators: + , - , * , /
- Usually we **alias** the new column using the AS keyword

```
SELECT title, (rental_rate/rental_duration) AS rental_rate_per_day  
FROM film;
```

SQL - Practicals (Part II)

Take a look at the ERD in the Pagila database and attempt the second practical ("Using Simple Conditional Statements to filter the data")

SQL - More on Conditional

SQL allows us access to operations that improves the conditional statement WHERE

- **LIKE:** Similar to WHERE =, but in cases you're not entirely sure what you're searching for
- **IN:** Similar to WHERE =, but for more than one condition
- **NOT:** Used to return opposite results from LIKE or IN (e.g. NOT LIKE or NOT IN)
- **AND & BETWEEN:** Allows us to combine operations
- **OR:** Allows us to find data where one of the conditions we provide to the query is true

SQL - LIKE

- LIKE is typically performed on string based columns
- Allows us to match strings using wildcards:
 - One wildcard character in SQL is %
 - Represents zero or more characters
 - Case sensitive
 - Remember to use single quotes in the query

```
SELECT *
FROM film
WHERE description LIKE '%Drama%'
```

SQL - IN

- IN can be used with both string and numeric data types
- Essentially allows us to use an =, but over more than one condition

```
SELECT *
FROM language
WHERE name IN ('English', 'Italian')
```

```
SELECT *
FROM film
WHERE language_id IN (1, 2)
```

SQL - NOT

- NOT allows us to inverse the results of the previous queries

```
SELECT *
FROM category
WHERE name NOT IN ('Action', 'Animation')
```

```
SELECT *
FROM actor
WHERE last_name NOT LIKE '%SON'
```

SQL - AND

- We use AND when we want to run/check against multiple conditions
- We can link as many expressions as we want
 - Of different types too... LIKE, IN, and NOT can all be linked by AND

```
SELECT *
FROM payment
WHERE payment_date >= '2017-01-25'
AND payment_date <= '2017-01-29'
```

```
SELECT *
FROM payment
WHERE payment_date < '2017-01-25' AND staff_id = 2
```

SQL - OR

- Similar to AND, OR combines multiple statements. However, with OR, only one of the conditions we specify needs to be true (instead of all cases as is the case with AND)

```
SELECT *
FROM film
WHERE rental_duration < 5 OR length > 120
```

SQL - BETWEEN

- When we want to specify a range between data on the same column (like we did in the previous slide), it is easier and cleaner to use BETWEEN
 - Works with dates, strings and numbers

```
SELECT *
FROM payment
WHERE payment_date BETWEEN '2017-01-25' AND '2017-01-29';
```

```
SELECT customer_id, payment_date, amount
FROM payment
WHERE amount BETWEEN 10.0 AND 11.99;
```

SQL - Practicals (Part III)

Take a look at the ERD in the Pagila database and attempt the third practical ("Applying filters to SQL queries")