# Copyright Notice

These slides are distributed under the Creative Commons License.

# Data Engineering Lifecycle

# Course Plan

**Week 1**   **Common source systems**
- Databases, object storage, and streaming sources
- Working with source systems on AWS

**Week 2**   **Setting up ingestion from source systems**

**Week 3**   **DataOps undercurrent**
- Automating some of your pipeline tasks
- Monitoring data quality

**Week 4**   **Orchestration, monitoring, and automating data pipelines**
- Setting up directed acyclic graphs
- Working with infrastructure as code

Introduction to Source Systems

DeepLearning.AI

**Different Types of Source Systems**

**Structured Data**

Data organized as tables of rows and columns

| ID | Last | First | Card |
|----|------|-------|------|
| 14 | Barry | John | XXX878 |
| 25 | Goode | Cynthia | XXX980 |
| 14 | Barry | John | XXX978 |
| 25 | Goode | Cynthia | XXX990 |

→ **Rows**

↓ **Columns**

DeepLearning.AI

| | Avg. Profit per Item | ROI |
|---|---|---|
| | $36.91 | 419% |

**Listing Details** | Cost Details

| Category | Listing Type | Your Cost | | | | Total Price | Total Fees | Ship | Prep |
|---|---|---|---|---|---|---|---|---|---|
| | | Cost | Tax | S/H | Total Cost | | | | |
| | | | $0.00 | $0.00 | $0.00 | $24.00 | $6.50 | | |
| | | | $0.00 | $0.00 | $19.00 | $55.00 | $5.84 | | |
| | | $19.00 | $0.00 | $0.00 | $11.50 | $49.00 | $3.12 | | |
| | | $19.00 | $1.50 | $0.00 | $0.00 | $24.00 | $12.16 | | |
| | | $0.00 | $0.00 | $0.00 | $21.85 | $107.00 | $7.15 | | |
| | | $19.00 | $2.85 | $0.00 | $5.75 | $61.00 | $5.84 | | |
| | | $0.00 | $0.75 | $0.00 | $5.75 | $49.00 | $8.24 | | |
| | | $19.00 | $0.75 | $0.00 | $0.00 | $71.00 | $4.21 | | |
| | | $0.00 | $0.00 | $0.00 | $0.00 | $34.00 | $7.37 | | |
| | | | $0.00 | $0.00 | $21.85 | $63.00 | $8.46 | | |
| | | | $0.00 | $0.00 | $5.75 | $73.00 | $12.49 | | |
| | | | $2.85 | $0.00 | $21.85 | $110.00 | $4.86 | | |
| | | | $0.75 | $0.00 | $0.00 | $40.00 | $3.88 | | |
| | | | $2.85 | $0.00 | $0.00 | $31.00 | $4.97 | | |
| | | | | | | $41.00 | | | |

MyInventory   Count: 90   Sum: $1,427.05

**Customer**

- CustomerID
- FName
- LName
- Address
- ZipCode
- Gender
- BirthDate
- Email
- MobileNo

**Order**

- OrderID
- OrderDate
- CustomerID
- SalemanID
- TotalAmt
- OrderRemark
- Status

**OrderDetail**

- OrderID
- ProductID
- Qty
- Price
- Discount
- NetPrice
- Amt

**Product**

- ProductID
- ProductName
- Price
- Description
- Unit
- BrandID

**Brand**

- BrandID
- BrandName

**Saleman**

- SalemanID
- SalemanName
- Level

```python
import csv
with open('eggs.csv', newline='') as csvfile:
    spamreader = csv.reader(csvfile, delimiter=' ', quotechar='|')
    for row in spamreader:
        print(', '.join(row))
```
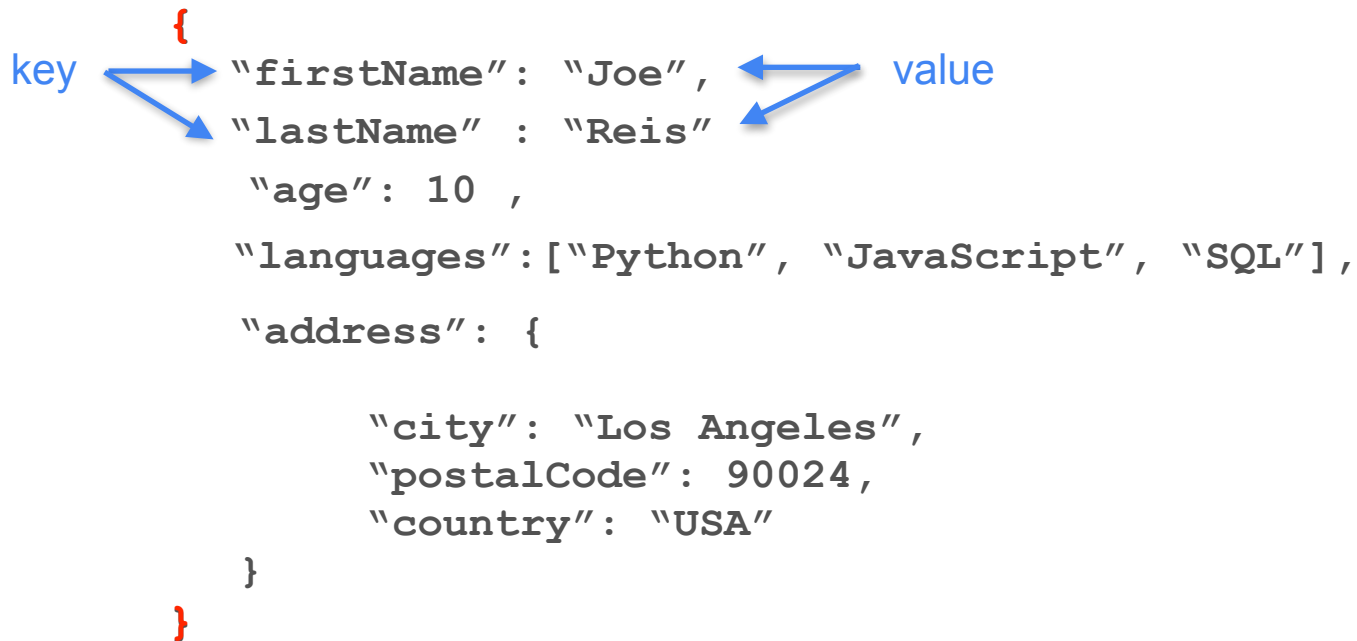
**Structured Data**  Data organized as tables of rows and columns

**Semi-Structured Data**  Data that is not in tabular form but still has some structure

**JavaScript Object Notation (JSON)**

A series of key-value pairs

key → value

```
{
    "firstName": "Joe",
    "lastName" : "Reis"
     "age": 10 ,
    "languages":["Python", "JavaScript", "SQL"],
    "address": {

        "city": "Los Angeles",
        "postalCode": 90024,
        "country": "USA"
    }
}
```

**Structured Data** — Data organized as tables of rows and columns

**Semi-Structured Data** — Data that is not in tabular form but still has some structure

**JavaScript Object Notation (JSON)**

A series of key-value pairs

key → "firstName": "Joe", ← value

"lastName" : "Reis"

"age": 10 ,

"languages": ["Python", "JavaScript", "SQL"],

"address": {

**Nested JSON format**

keys → "city": "Los Angeles", ← values
"postalCode": 90024,
"country": "USA"

}

}

DeepLearning.AI

**Structured Data** — Data organized as tables of rows and columns

**Semi-Structured Data** — Data that is not in tabular form but still has some structure

**Unstructured Data** — Data that does not have any predefined structure

**Text** — TXT

**Video** — MP4

**Audio** — MP3

**Images** — PNG
- dimensions
- pixel colors

# Databases

**Store data in an organized way**

Structured data
Semi-structured data

**C**reate
**R**ead
**U**pdate
**D**elete

Database
Storage

Database
Management
System
(DBMS)

Person/
Application

TXT PNG MP3 MP4 CSV

## Databases

**Store data in an organized way**

Structured data
Semi-structured data

**C**reate
**R**ead
**U**pdate
**D**elete

Database
Storage

Database
Management
System
(DBMS)

Person/
Application

### Relational databases

Tables with rows and columns

### Non-relational (NoSQL) databases

Key | Value | Document

k1 → value 1

k2 → value 2

k3 → value 3

```
{
    "firstName": "Joe",
    "lastName" : "Reis",
    "age": 10
}
```

Non-tabular data

DeepLearning.AI

**Files**

## Sequence of bytes representing information

TXT  PNG  MP3  MP4  CSV

|   | A | B | C |
|---|---|---|---|
| 1 |   |   |   |
| 2 |   |   |   |
| 3 |   |   |   |

```
{
 "firstName": "Joe",
 "lastName" : "Reis",
 "languages":["R", "SQL"],
}
```

Amazon S3

DeepLearning.AI

**Streaming Systems**

**Continuous flow of data**

Semi-structured data

Producer

Producer

Message queue/
Streaming platform

Consumer

DeepLearning.AI

# Streaming Systems

**Continuous flow of data**

Semi-structured data

## Source System

Producer

Smart Thermostat

Amazon Kinesis

kafka

**Your ingestion pipeline starts here**

TXT PNG MP3 MP4 CSV
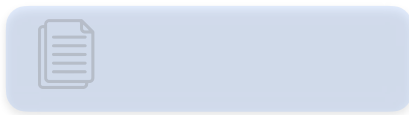
**Databases**

**Files**

**Streaming Systems**

**Store data in an organized way**

**Sequence of bytes representing information**

**Continuous flow of data**

**Ingest**

- Structured
- Semi-structured
- Unstructured

# Introduction to Source Systems

**Relational Databases**

# Relational Databases

Customers

| key | | |
|---|---|---|
| | | |
| | | |
| | | |
| | | |

Products

| key | | |
|---|---|---|
| | | |
| | | |
| | | |
| | | |

| key | | |
|---|---|---|
| | | |
| | | |
| | | |
| | | |

Orders

- Reduce redundancy
- Make data easier to manage

DeepLearning.AI

# Relational Databases

One big table for everything!

| name | address | phone | date_time | amount | brand | SKU | description |
|------|---------|-------|-----------|--------|-------|-----|-------------|
| Jane Doe | 74th Street | 12345678 | 12/08/2024 | 700 | ABC | B32 | Blender |
| Jane Doe | 74th Street | 12345678 | 12/08/2024 | 99 | XYZ | i56 | Iron |
| Jane Doe | 74th Street | 12345678 | 12/08/2024 | 100 | GHJ | k70 | Kettle |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |

Jane Doe

# Relational Databases

One big table for everything!

| name | address | phone | date_time | amount | brand | SKU | description |
|------|---------|-------|-----------|--------|-------|-----|-------------|
| Jane Doe | 74th Street | 12345678 | 12/08/2024 | 700 | ABC | B32 | Blender |
| Jane Doe | 74th Street | 12345678 | 12/08/2024 | 99 | XYZ | i56 | Iron |
| Jane Doe | 74th Street | 12345678 | 12/08/2024 | 100 | GHJ | k70 | Kettle |
| Mary Ann | 19th Avenue | 98765432 | 13/08/2024 | 899 | STU | w40 | Washer |
| John Ken | 1st Link | 36891623 | 14/08/2024 | 899 | STU | w40 | Washer |
| Ivy Tan | 67th Street | 98639513 | 15/08/2024 | 899 | STU | w40 | Washer |

# Relational Databases

One big table for everything!

**Inconsistency**

| name | address | phone | date_time | amount | brand | SKU | description |
|---|---|---|---|---|---|---|---|
| Jane Doe | 11th Avenue | 12345678 | 12/08/2024 | 700 | ABC | B32 | Blender |
| Jane Doe | 11th Avenue | 12345678 | 12/08/2024 | 99 | XYZ | i56 | Iron |
| Jane Doe | 74th Street | 12345678 | 12/08/2024 | 100 | GHJ | k70 | Kettle |
| Mary Ann | 19th Avenue | 98765432 | 13/08/2024 | 899 | STU | w31 | Washer |
| John Ken | 1st Link | 36891623 | 14/08/2024 | 899 | STU | w31 | Washer |
| Ivy Tan | 67th Street | 98639513 | 15/08/2024 | 899 | STU | w40 | Washer |

Jane Doe

now lives on 11th Avenue

SKU
now w31

**Inconsistency**

DeepLearning.AI

# Relational Databases

### Customers

| id | first_name | last_name | age | address |
|----|-----------|-----------|-----|---------|
| 1 | Jane | Doe | 24 | 11th Ave. |
| 2 | Mary | Ann | 65 | 19th Ave. |
| 3 | John | Ken | 27 | 1st Link |
| 4 | Ivy | Tan | 18 | 67th St. |

**Single customer**

### Products

| id | brand | SKU | description |
|----|-------|-----|-------------|
| 1 | ABC | b32 | Blender |
| 2 | XYZ | i56 | Iron |
| 3 | GHJ | k70 | Kettle |
| 4 | STU | w31 | Washer |

**Single product**

### Orders

| id | customer_id | product_id | date_time | purchase_amount |
|----|-------------|------------|-----------|-----------------|
|  |  |  |  |  |
|  |  |  |  |  |
|  |  |  |  |  |
|  |  |  |  |  |
|  |  |  |  |  |
|  |  |  |  |  |

**Database schema**

# Relational Databases

**Keys**

Primary key:
uniquely
identifies each
row in a table

### Customers

| id | first_name | last_name | age | address |
|----|-----------|-----------|-----|---------|
| 1 | Jane | Doe | 24 | 11th Ave. |
| 2 | Mary | Ann | 65 | 19th Ave. |
| 3 | John | Ken | 27 | 1st Link |
| 4 | Ivy | Tan | 18 | 67th St. |

### Products

| id | brand | SKU | description |
|----|-------|-----|-------------|
| 1 | ABC | b32 | Blender |
| 2 | XYZ | i56 | Iron |
| 3 | GHJ | k70 | Kettle |
| 4 | STU | w31 | Washer |

Orders

| id | customer_id | product_id | date_time | purchase_amount |
|----|-------------|------------|-----------|-----------------|
| 1 | 1 | 1 | 12/08/2024 | 700 |
| 2 | 1 | 2 | 12/08/2024 | 99 |
| 3 | 1 | 3 | 12/08/2024 | 100 |
| 4 | 2 | 4 | 13/08/2024 | 899 |
| 5 | 3 | 4 | 14/08/2024 | 899 |
| | | | | |

**Database schema**

Foreign key:
references the primary key of the customer table

# Relational Databases

string

integer

## Customers

| id | first_name | last_name | age | address |
|----|-----------|-----------|-----|---------|
| 1 | Jane | Doe | 24 | 11th Ave. |
| 2 | Mary | Ann | 65 | 19th Ave. |
| 3 | John | Ken | 27 | 1st Link |
| 4 | Ivy | Tan | 18 | 67th St. |

## Products

| id | brand | SKU | description |
|----|-------|-----|-------------|
| 1 | ABC | b32 | Blender |
| 2 | XYZ | i56 | Iron |
| 3 | GHJ | k70 | Kettle |
| 4 | STU | w31 | Washer |

Orders

| id | customer_id | product_id | date_time | purchase_amount |
|----|------------|-----------|-----------|-----------------|
| 1 | 1 | 1 | 12/08/2024 | 700 |
| 2 | 1 | 2 | 12/08/2024 | 99 |
| 3 | 1 | 3 | 12/08/2024 | 100 |
| 4 | 2 | 4 | 13/08/2024 | 899 |
| 5 | 3 | 4 | 14/08/2024 | 899 |
| | | | | |

**Database schema**

Each row in a table has to follow the same column structure:
same sequence of columns and data types

# Relational Databases

## Customers

| id | first_name | last_name | age | address |
|----|-----------|-----------|-----|---------|
| 1 | Jane | Doe | 24 | 11th Ave. |
| 2 | Mary | Ann | 65 | 19th Ave. |
| 3 | John | Ken | 27 | 1st Link |
| 4 | Ivy | Tan | 18 | 67th St. |

## Products

| id | brand | SKU | description |
|----|-------|-----|-------------|
| 1 | ABC | b32 | Blender |
| 2 | XYZ | i56 | Iron |
| 3 | GHJ | k70 | Kettle |
| 4 | STU | w31 | Washer |

## Orders

| id | customer_id | product_id | date_time | purchase_amount |
|----|-------------|------------|-----------|-----------------|
| 1 | 1 | 1 | 12/08/2024 | 700 |
| 2 | 1 | 2 | 12/08/2024 | 99 |
| 3 | 1 | 3 | 12/08/2024 | 100 |
| 4 | 2 | 4 | 13/08/2024 | 899 |
| 5 | 3 | 4 | 14/08/2024 | 899 |
| 6 | 1 | 4 | 15/08/2024 | 899 |

DeepLearning.AI

# Relational Databases

One big table for everything!

| name | address | phone | date_time | amount | brand | SKU | description |
|------|---------|-------|-----------|--------|-------|-----|-------------|
| Jane Doe | 74th Street | 12345678 | 12/08/2024 | 700 | ABC | B32 | Blender |
| Jane Doe | 74th Street | 12345678 | 12/08/2024 | 99 | XYZ | i56 | Iron |
| Jane Doe | 74th Street | 12345678 | 12/08/2024 | 100 | GHJ | k70 | Kettle |
| Mary Ann | 19th Avenue | 98765432 | 13/08/2024 | 899 | STU | w40 | Washer |
| John Ken | 1st Link | 36891623 | 14/08/2024 | 899 | STU | w40 | Washer |
| Ivy Tan | 67th Street | 98639513 | 15/08/2024 | 899 | STU | w40 | Washer |

One Big Table (OBT) approach: use cases that need faster processing

# Relational Databases



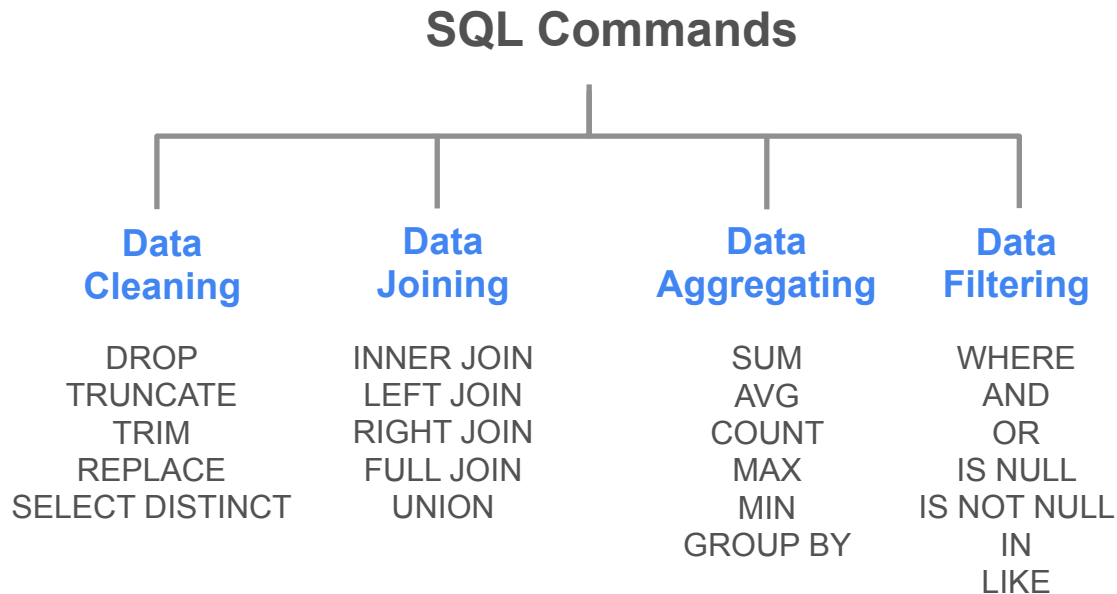**Relational Database Management System (RDBMS)**

Software layer that sits on top of a relational database to manage and interact with the data.

**Structured Query Language (SQL)**

# Relational Databases



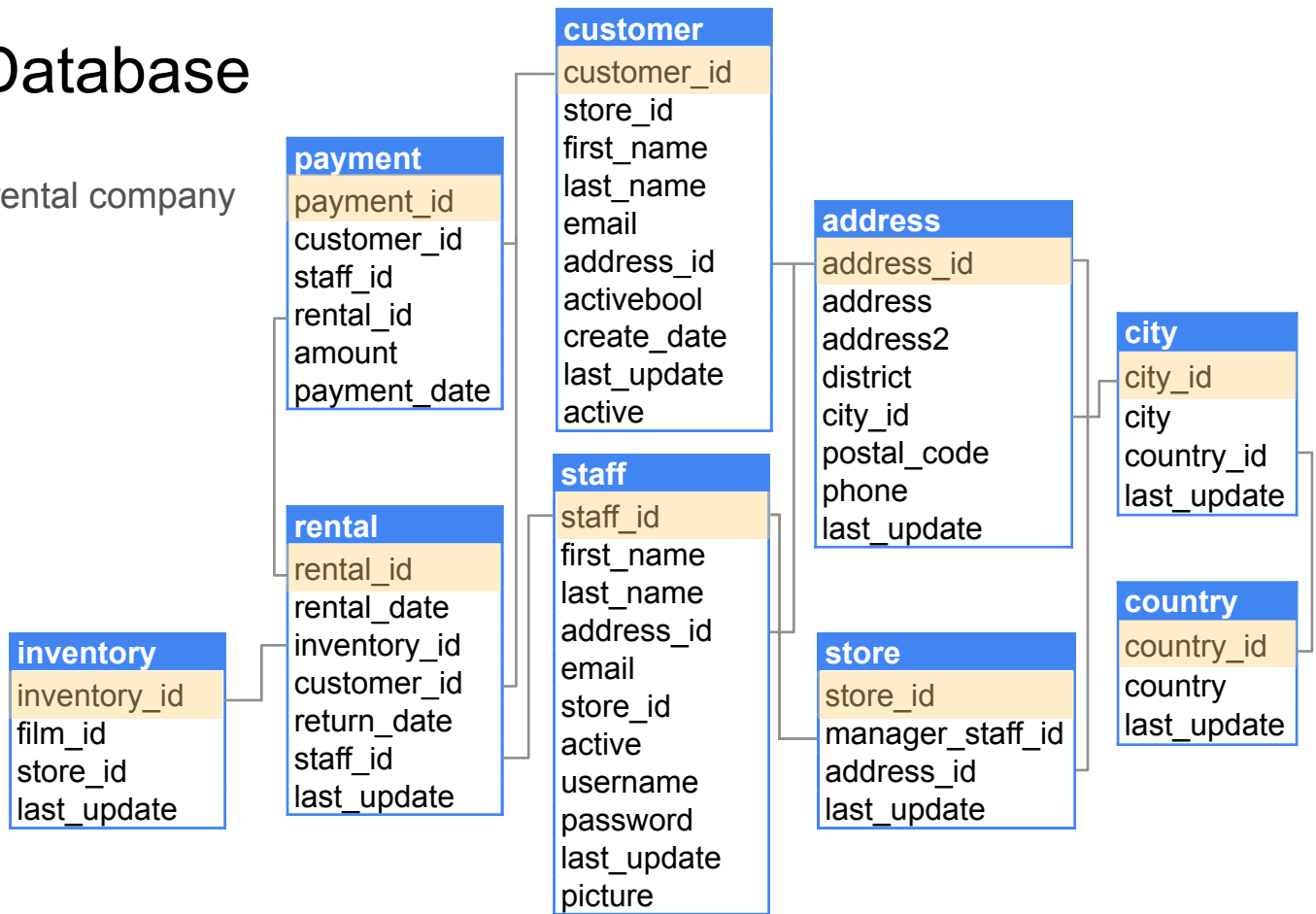**SQL Commands**

**Data Cleaning**

DROP
TRUNCATE
TRIM
REPLACE
SELECT DISTINCT

**Data Joining**

INNER JOIN
LEFT JOIN
RIGHT JOIN
FULL JOIN
UNION

**Data Aggregating**

SUM
AVG
COUNT
MAX
MIN
GROUP BY

**Data Filtering**

WHERE
AND
OR
IS NULL
IS NOT NULL
IN
LIKE

# The Relational Database

- Database for a fictitious DVD rental company called **Rentio**

- Database schema

SQL queries

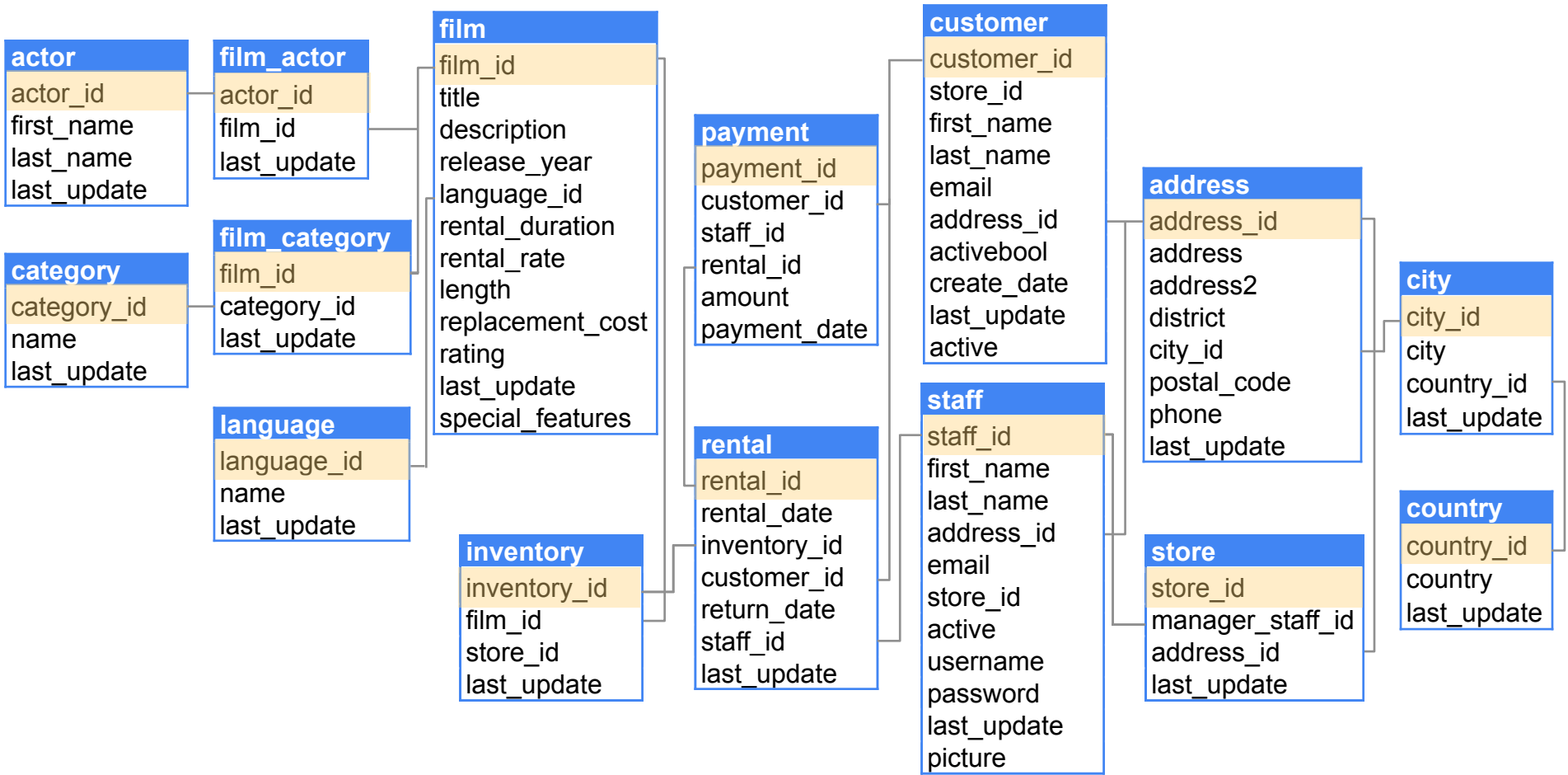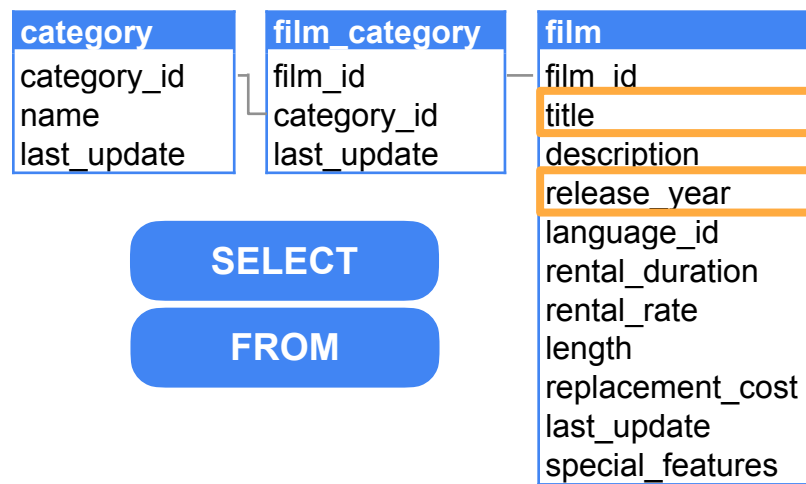Answer business questions

**payment**
- payment_id
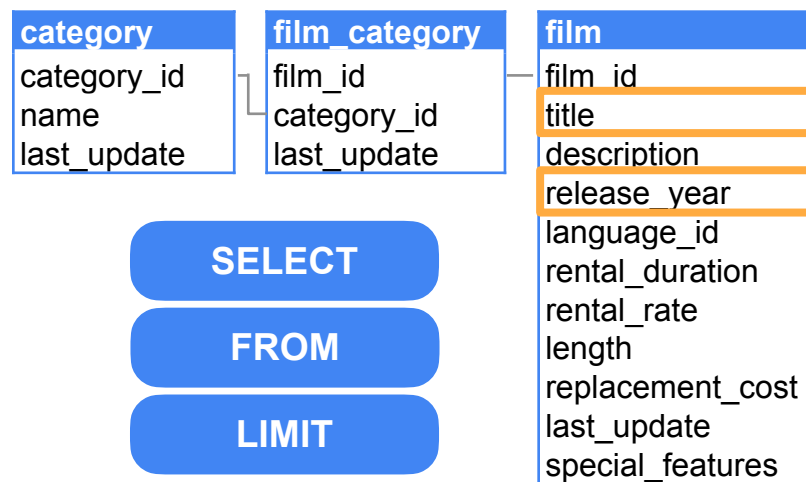- customer_id
- staff_id
- rental_id
- amount
- payment_date

**customer**
- customer_id
- store_id
- first_name
- last_name
- email
- address_id
- activebool
- create_date
- last_update
- active

**address**
- address_id
- address
- address2
- district
- city_id
- postal_code
- phone
- last_update

**city**
- city_id
- city
- country_id
- last_update

**rental**
- rental_id
- rental_date
- inventory_id
- customer_id
- return_date
- staff_id
- last_update

**staff**
- staff_id
- first_name
- last_name
- address_id
- email
- store_id
- active
- username
- password
- last_update
- picture

**store**
- store_id
- manager_staff_id
- address_id
- last_update

**country**
- country_id
- country
- last_update

**inventory**
- inventory_id
- film_id
- store_id
- last_update

# Entity Relationship

**actor**
- actor_id
- first_name
- last_name
- last_update

**film_actor**
- actor_id
- film_id
- last_update

**category**
- category_id
- name
- last_update

**film_category**
- film_id
- category_id
- last_update

**language**
- language_id
- name
- last_update

**film**
- film_id
- title
- description
- release_year
- language_id
- rental_duration
- rental_rate
- length
- replacement_cost
- last_update
- special_features

**inventory**
- inventory_id
- film_id
- store_id
- last_update

**payment**
- payment_id
- customer_id
- staff_id
- rental_id
- amount
- payment_date

**rental**
- rental_id
- rental_date
- inventory_id
- customer_id
- return_date
- staff_id
- last_update

**customer**
- customer_id
- store_id
- first_name
- last_name
- email
- address_id
- activebool
- create_date
- last_update
- active

**staff**
- staff_id
- first_name
- last_name
- address_id
- email
- store_id
- active
- username
- password
- last_update
- picture

**store**
- store_id
- manager_staff_id
- address_id
- last_update

**address**
- address_id
- address
- address2
- district
- city_id
- postal_code
- phone
- last_update

**city**
- city_id
- city
- country_id
- last_update

**country**
- country_id
- country
- last_update

DeepLearning.AI

# Entity Relationship

**actor**
- actor_id
- first_name
- last_name
- last_update

**film_actor**
- actor_id
- film_id
- last_update

**film_category**
- film_id
- category_id
- last_update

**category**
- category_id
- name
- last_update

**language**
- language_id
- name
- last_update

**film**
- film_id
- title
- description
- release_year
- language_id
- rental_duration
- rental_rate
- length
- replacement_cost
- last_update
- special_features

**inventory**
- inventory_id
- film_id
- store_id
- last_update

**payment**
- payment_id
- customer_id
- staff_id
- rental_id
- amount
- payment_date

**rental**
- rental_id
- rental_date
- inventory_id
- customer_id
- return_date
- staff_id
- last_update

**customer**
- customer_id
- store_id
- first_name
- last_name
- email
- address_id
- activebool
- create_date
- last_update
- active

**staff**
- staff_id
- first_name
- last_name
- address_id
- email
- store_id
- active
- username
- password
- last_update
- picture

**store**
- store_id
- manager_staff_id
- address_id
- last_update

**address**
- address_id
- address
- address2
- district
- city_id
- postal_code
- phone
- last_update

**city**
- city_id
- city
- country_id
- last_update

**country**
- country_id
- country
- last_update

## category
category_id
name
last_update

## film_category
film_id
category_id
last_update

## film
film_id
title
description
release_year
language_id
rental_duration
rental_rate
length
replacement_cost
last_update
special_features

SELECT

FROM

```
In [1]:  ▶  %load_ext sql
            %sql mysql+pymysql://root:adminpwrd@localhost:3306/sakila

In [ ]:  ▶  %%sql

In [ ]:  ▶
```

DeepLearning.AI

## category
category_id
name
last_update

## film_category
film_id
category_id
last_update

## film
film_id
title
description
release_year
language_id
rental_duration
rental_rate
length
replacement_cost
last_update
special  features

SELECT

FROM

LIMIT

| | |
|---|---|
| ADAPTATION HOLES | 2006 |
| AFFAIR PREJUDICE | 2006 |
| AFRICAN EGG | 2006 |
| AGENT TRUMAN | 2006 |
| AIRPLANE SIERRA | 2006 |
| AIRPORT POLLOCK | 2006 |
| ALABAMA DEVIL | 2006 |
| ALADDIN CALENDAR | 2006 |

In [ ]: ▶ %%sql

## category
category_id
name
last_update

## film_category
film_id
category_id
last_update

## film
film_id
title
description
release_year
language_id
rental_duration
rental_rate
length
replacement_cost
last_update
special_features

SELECT

FROM

WHERE

Exploring the films that are less than 60 minutes long.



| | 3 | 4.99 | 48 | | 12.99 | G | Trailers,Deleted Scer |

```
In [ ]:    ▶    %%sql
```

## category

category_id
name
last_update

## film_category

film_id
category_id
last_update

## film

film_id
title
description
release_year
language_id
rental_duration
rental_rate
length
replacement_cost
last_update
special_features

**SELECT**

**FROM**

**WHERE**

**ORDER BY**

**LIMIT**

| 3 | 4.99 | 48 | 12.99 | G | Trailers,Deleted Scer |

```
In [5]:  %%sql
         SELECT *
         FROM film
         WHERE length<60
         ORDER BY length
```

| original_language_id | rental_duration | rental_rate | length | replacement_cost | rating |
| --- | --- | --- | --- | --- | --- |

**category**
- category_id
- name
- last_update

**film_category**
- film_id
- category_id
- last_update

**film**
- film_id
- title
- description
- release_year
- language_id
- rental_duration
- rental_rate
- length
- replacement_cost
- last_update
- special_features

SELECT

FROM

JOIN

WHERE

ORDER BY

LIMIT

Monkey in
Ancient India

```sql
%%sql
SELECT *
FROM film
JOIN film_category
ON film.film_id = film_category.film_id
JOIN category
ON film_category.category_id = category.category_id
WHERE length<60
```

In [6]:

| date | film_id_1 | category_id | last_update_1 | category_id_1 | name | last_update_2 |
|------|-----------|-------------|---------------|---------------|------|---------------|
| 2-15 | 2 | 11 | 2006-02-15 | 11 | Horror | 2006-02-15 |

DeepLearning.AI

## category

category_id
name
last_update

## film_category

film_id
category_id
last_update

## film

film_id
title
description
release_year
language_id
rental_duration
rental_rate
length
replacement_cost
last_update
special_features

**SELECT**

**FROM**

**JOIN**

**WHERE**

**ORDER BY**

**LIMIT**

## INNER JOIN

JOIN: combine the records from both tables that have a matching column value specified in the ON statement.

film has a row with film_id = 123
film_category does not have a row with film_id= 123

The row with film_id = 123 will not be in the join results

## category
category_id
name
last_update

## film_category
film_id
category_id
last_update

## film
film_id
title
description
release_year
language_id
rental_duration
rental_rate
length
replacement_cost
last_update
special_features

SELECT

FROM

JOIN

WHERE

GROUP BY

LIMIT

Monkey in
Ancient India

In [7]:

```sql
%%sql
SELECT film.title, category.name
FROM film
JOIN film_category
ON film.film_id = film_category.film_id
JOIN category
ON film_category.category_id = category.category_id
WHERE length<60
```

Out[7]:

| title | name |
|---|---|
| ACE GOLDFINGER | Horror |
| ADAPTATION HOLES | Documentary |
| AIRPORT POLLOCK | Horror |
| ALIEN CENTER | Foreign |
| ALTER VICTORY | Animation |

## category

category_id
name
last_update

## film_category

film_id
category_id
last_update

## film

film_id
title
description
release_year
language_id
rental_duration
rental_rate
length
replacement_cost
last_update
special_features

**SELECT**

**COUNT**

**FROM**

**JOIN**

**WHERE**

**GROUP BY**

**ORDER BY**

**LIMIT**

Monkey in
Ancient India

In [7]:

```sql
%%sql
SELECT film.title, category.name
FROM film
JOIN film_category
ON film.film_id = film_category.film_id
JOIN category
ON film_category.category_id = category.category_id
WHERE length<60
```

| | |
|---|---|
| DAWN POND | Games |
| DEEP CRUSADE | Documentary |
| DESTINY SATURDAY | New |
| DIVORCE SHINING | Sports |
| DOCTOR GRAIL | Children |
| DOORS PRESIDENT | Animation |

# Common SQL Commands

- SELECT
- COUNT
- FROM
- JOIN
- WHERE
- GROUP BY
- ORDER BY
- LIMIT

# Data Manipulation Operations

- CREATE
- INSERT INTO
- UPDATE
- DELETE

# Introduction to Source Systems

DeepLearning.AI

## NoSQL Databases

# NoSQL Databases

**NoSQL**

# NoSQL Databases

**No    SQL**

# NoSQL Databases

## Not Only SQL

### Non-Relational Databases

It can still support SQL or SQL-like query languages.

**Relational Databases**

# NoSQL Databases

**Non-tabular structures**

## Key-Value

| k1 | → | value 1 |
| k2 | → | value 2 |
| k3 | → | value 3 |

## Document

```
{
  "firstName": "Joe",
  "lastName" : "Reis",
  "age": 10,
  "address": {
      "city": "Los Angeles",
      "postalCode": 90024,
      "country": "USA"
  }
}
```

## Wide-Column

| Row A | Column 1 | Column 2 | Column 3 |
|-------|----------|----------|----------|
|       | Value 1  | Value 2  | Value 3  |

| Row B | Column 1 | Column 2 | Column 3 |
|-------|----------|----------|----------|
|       | Value 1  | Value 2  | Value 3  |

## Graph

- No predefined schemas
- More flexibility when storing your data

DeepLearning.AI

# Horizontal Scaling



write

user 1

NoSQL
Database

Not updated Updated

read

Data received is
is not up-to-date

user 2

Updated

**Eventual consistency**

DeepLearning.AI

# Consistency

| NoSQL Databases | Relational Databases |
|---|---|
| Eventual Consistency | Strong Consistency |
| • Speed is prioritized<br><br>• System availability and scalability are important | • Read data only when all nodes have been updated |

# NoSQL Databases

Not all NoSQL databases guarantee:

## **ACID** compliance

**A**tomicity

**C**onsistency

**I**solation

**D**urability

# Specialized Query Language

**Example of NoSQL Data**

```
{
  "id": 1,
  "key": "Blender",
  "qty": 6,
  "sku": "b32"
}
```

**Query**

```
db.products.find({qty: {$gt: 4}})
```

*Ref: AWS docs*

DeepLearning.AI

# Types of No-SQL Databases

## Key-Value

key                    value

84620 → {"name": "blender", "sku": "b32", "quantity":6}

64820 → {"name": "iron", "sku": "i56", "quantity":5}

46173 → {"name": "washer", "sku": "w40", "quantity":6}

## Document

```
{
  "firstName": "Joe",
  "lastName" : "Reis",
  "age": 10,
  "address": {
     "city": "Los Angeles",
     "postalCode": 90024,
     "country": "USA"
  }
}
```

DeepLearning.AI

# Key-Value Database

**Key-Value**

key                                    value



Unique Identifier

| 84620 | → | {"name": "blender", "sku": "b32", "quantity":6} |
| 64820 | → | {"name": "iron", "sku": "i56", "quantity":5} |
| 46173 | → | {"name": "washer", "sku": "w40", "quantity":6} |

Fast lookup: such as caching user session data

- viewing different products
- adding items to the shopping cart
- checking out

key                        value

| user_session_id | → | values |

# Document Database

**Collection** *(Like a table)*

```
{
    "users" : [
        {
            "id": 1234,
            "name": {
                    "first":"Joe",
                    "last":"Reis"
                    },
            "favorite_bands" : ["AC/DC", "Slayer", "WuTang Clan", "Action Bronson" ]
        },
        {
            "id":1235,
            "name": {
                    "first": "Matt",
                    "last":"Housley"
                    },
            "favorite_bands" : ["Dave Matthews Band", "Creed", "Nickelback"]
        }
    ]
}
```

keys

**Single users**
**Documents**
*(Like a row)*

# Document Database

```
{
 "users" : [
        {
         "id": 1234,
         "name": {
                "first":"Joe",
                "last":"Reis"
                },
         "favorite_bands" : ["AC/DC", "Slayer", "WuTang Clan", "Action Bronson" ]
        },

        {
         "id":1235,
         "name": {
                "first": "Matt",
                "last":"Housley"
                },
         "favorite_bands" : ["Dave Matthews Band", "Creed", "Nickelback"]
        }
    ]
}
```

| user_id | band_id |
|---------|---------|
| 1234 | 1 |
| 1234 | 2 |
| 1234 | 5 |
| 1234 | 6 |
| 1235 | 7 |
| 1235 | 3 |
| 1235 | 4 |

| band_id | band_name |
|---------|-----------|
| 1 | AC/DC |
| 2 | Slayer |
| 3 | Creed |
| 4 | Nickelback |
| 5 | Wutan Clan |
| 6 | Action Bronson |
| 7 | Dave Matthews Band |

| user_id | first_name | last_name |
|---------|-----------|-----------|
| 1234 | Joe | Reis |
| 1235 | Matt | Housely |

- Easy to retrieve all the information about a user (locality)

- Document stores don't support joins

- Flexible schema

Fixed schema

# Document Database

## Use cases

- Content management

- Catalogs

- Sensor readings

```
{
 "iot" : [
        {
        "id": 24,
        "interaction": "some_interaction",
        "device": "my_device",
        "sensor_reading": 34
        }
        ]
}
```

Flexible Schema

# Document Database

Document databases become absolute nightmares to manage and query.



Source system owner → Change in data → NoSQL Document → Ingest → ✗ → Downstream use

# Introduction to Source Systems

**Database ACID Compliance**

# OLTP Systems

**OLTP**



**Online Transaction Processing**

Support very high transaction rates  (bank account balances, online orders)

# ACID Compliance

| Relational Databases | NoSQL Databases |
| --- | --- |
| ACID compliant | Not ACID compliant by default |
| **A**tomicity | |
| **C**onsistency | |
| **I**solation | |
| **D**urability | |

They help ensure transactions are processed reliably and accurately in an OLTP system.

# ACID Compliance

**Atomicity**

Atomicity ensures that transactions are **atomic**, treated as a single, indivisible unit.

**A transaction: placing an order**

Deducting the total cost from the customer's account

Updating the inventory to reflect the purchased item

Both operations must happen as a single transaction

## Atomicity

Atomicity ensures that transactions are **atomic**, treated as a single, indivisible unit.

## Consistency

Any changes to the data made within a transaction follow the set of rules or constraints defined by the database schema.

| id | product_name | quantity |
|----|--------------|----------|
| 1  | blender      | 1        |
|    |              |          |
|    |              |          |
|    |              |          |

Buy 2 blenders

**Transaction**

| id | product_name | quantity |
|----|--------------|----------|
| 1  | blender      | -1       |
|    |              |          |
|    |              |          |
|    |              |          |

Rule: stock level ≥ 0

DeepLearning.AI

**Atomicity**  Atomicity ensures that transactions are **atomic**, treated as a single, indivisible unit.

**Consistency**  Any changes to the data made within a transaction follow the set of rules or constraints defined by the database schema.

**ACID** compliance

**A**tomicity

**C**onsistency

**I**solation

**D**urability

**Strong Consistency**

All nodes provide the same up-to-date

Up-to-date

Up-to-date

**Atomicity**
Atomicity ensures that transactions are **atomic**, treated as a single, indivisible unit.

**Consistency**
Any changes to the data made within a transaction follow the set of rules or constraints defined by the database schema.

**Isolation**
Each transaction is executed independently in sequential order.

| id | product_name | quantity |
|----|--------------|----------|
| 1  | blender      | 5        |
|    |              |          |
|    |              |          |
|    |              |          |

**Transaction**

Buy 5 blenders

**Transaction**

Buy 10 blenders

DeepLearning.AI

**Atomicity** — Atomicity ensures that transactions are **atomic**, treated as a single, indivisible unit.

**Consistency** — Any changes to the data made within a transaction follow the set of rules or constraints defined by the database schema.

**Isolation** — Each transaction is executed independently in sequential order.

**Durability** — Once a transaction is completed, its effects are permanent and will survive any subsequent system failures.

*Essential for maintaining the reliability of the database*

# ACID Compliance

The ACID principles guarantee that a database will maintain a consistent picture of the world.



**Strong Consistency**

- Data is consistent across the entire network

- Key feature of relational databases that ensures ACID

Data is partitioned or replicated

# Interacting with Amazon DynamoDB



Amazon DynamoDB

Apply some **C**reate, **R**ead, **U**pdate and **D**elete (**CRUD**) operations

In this video,

- Overview of DynamoDB features

- Data you will work on

- DynamoDB methods that you will use to apply **CRUD** operations

**Amazon DynamoDB**
Key-value Database

**Key-value Items**

⬇

**Table**

- Row: attributes of one item
- Uniquely identified by the item's key.
- **Simple Primary Key:** partition key
- **Composite Primary Key:**

| PersonID | Attributes |
|---|---|

**Person** | **key** | **Value** |

| 101 | → | {"**FirstName**": "Joe", "**LastName**": "Reis", "**Phone**": "111-222", "**Country**": "USA", "**FavoriteBands**": {"Action Bronson", "Slayer", "WuTang Clan"}} |

| 102 | → | {"**FirstName**": "Matt", "**LastName**": "Housley", "**Phone**": "222-333", "**Country**": "USA"} |

Simple primary key ↘

| Key: PersonID | Attributes | | | | |
|---|---|---|---|---|---|
| **101** | FirstName | LastName | Phone | Country | FavoriteBands |
| | Joe | Reis | 111-222 | USA | {"Action Bronson", "Slayer", "WuTang Clan"} |
| **102** | FirstName | LastName | Phone | Country | |
| | Matt | Housley | 222-333 | USA | |

**Amazon DynamoDB**
Key-value Database

**Key-value Items**

⬇

**Table**

- Row: attributes of one item
- Uniquely identified by the item's key.
- **Simple Primary Key:** partition key
- **Composite Primary Key:** partition key & sort key

| Composite Primary Key | | Attributes | | | | |
|---|---|---|---|---|---|---|
| **Partition Key** | **Sort Key** | | | | | |
| OrderID | ItemNum | | | | | |
| 1234 | Item1 | Price | Quantity | ProductType | ISBN | Title |
| | | 10 | 1 | Book | 45679 | Data |
| 1234 | Item2 | Price | Quantity | ProductType | Brand | Color |
| | | 50 | 1 | Bike | AZY | Black |
| 1235 | Item1 | Price | Quantity | ProductCode | | |
| | | 23 | 4 | 23697 | | |
| 1235 | Item2 | Price | Quantity | ProductType | Brand | |
| | | 1200 | 2 | Laptop | XYZ | |

**Schema-less:** Each item can have its own distinct attributes.

Simple primary key →

| Key: PersonID | Attributes | | | | |
|---|---|---|---|---|---|
| 101 | FirstName | LastName | Phone | Country | FavoriteBands |
| | Joe | Reis | 111-222 | USA | {"Action Bronson", "Slayer", "WuTang Clan"} |
| 102 | FirstName | LastName | Phone | Country | |
| | Matt | Housley | 222-333 | USA | |

DeepLearning.AI

# Interacting with Amazon DynamoDB

**Table**

**Table**

**Table**

**Table**

**Interact with the tables using Python**

**Boto3**

AWS Software Development Kit  (SDK) for Python
Allows you to create and configure AWS services using Python

# Boto3 documentation

You use the AWS SDK for Python (Boto3) to create, configure, and manage AWS services, such as Amazon Elastic Compute Cloud (Amazon EC2) and Amazon Simple Storage Service (Amazon S3). The SDK provides an object-oriented API as well as low-level access to AWS services.

> ✏️ **Note**
>
> Documentation and developers tend to refer to the AWS SDK for Python as "Boto3," and this documentation often does so as well.

# Quickstart

- Quickstart
  - Installation
  - Configuration
  - Using Boto3
- A Sample Tutorial
  - SQS
  - Creating a queue
  - Using an existing queue
  - Sending messages
  - Processing messages
- Code Examples
  - Amazon CloudWatch examples
  - Amazon DynamoDB

# Interacting with Amazon DynamoDB


Table


Table


Table


Table

**Interact with the tables using Python**

**Boto3**   AWS Software Development Kit  (SDK) for Python
Allows you to create and configure AWS services using Python

```python
import boto3

client = boto3.client('dynamodb')
```

| **C**reate | create_table |
|---|---|
| **R**ead | scan<br>get_item<br>query |
| **U**pdate | put_item<br>write_batch_items<br>update_item |
| **D**elete | delete_item |

# Data

**Load**

**Product Catalog** (JSON)

- Information about some products sold on Amazon
- *ID*: simple primary key

Table

**Forum** (JSON)

- Information about some AWS forums
- Each forum: total number of threads, messages and views
- *Name*: simple primary key

Table

**Thread** (JSON)

- Information about each forum thread
- Each thread: subject, message, total number of views and replies
- *ForumName* & *Subject*: composite primary key

Table

**Reply** (JSON)

- Information about each thread reply
- Each reply: time, reply message, user, ID (Forum and thread subject)
- *ID* & *ReplyDateTime:* composite primary key

Table

# Data

Table



**Load**

**Product Catalog** JSON → Table

**Forum** JSON → Table

**Thread** JSON → Table

**Reply** JSON → Table

```
{
    "Forum": [
        {
            "PutRequest": {
                "Item": {
                    "Name": {"S":"Amazon DynamoDB"},
                    "Category": {"S":"Amazon Web Services"},
                    "Threads": {"N":"2"},
                    "Messages": {"N":"4"},
                    "Views": {"N":"1000"}    N: Number
                }
            }
        },
        {
            "PutRequest": {
                "Item": {
                    "Name": {"S":"Amazon S3"},   S: String
                    "Category": {"S":"Amazon Web Services"}
                }
            }
        }
    ]
}
```

Table

# 1 - Import Packages

First, let's import some packages. Among these packages, you can find `boto3`, which is the AWS Software Development Kit (SDK) for Python that allows you to interact with various AWS services using Python code. With `boto3`, you can programmatically access AWS resources such as EC2 instances, S3 buckets, Amazon DynamoDB tables, and more. It provides a simple and intuitive interface for managing AWS services, enabling developers to automate tasks, build applications, and integrate AWS services into their Python applications efficiently.

For more information on each of the methods that you will use throughout this lab, you can check out boto3 documentation.

```python
import decimal
import json
import logging
from typing import Any, Dict, List

import boto3
from botocore.exceptions import ClientError
```

# Introduction to Source Systems

**Object Storage**

# Object Storage

**Object Storage**

**Traditional File System Hierarchy**

files

No hierarchy!

# Object Storage



Amazon S3

# Object Storage

# Object Storage

**Object Storage**



files

- Storing semi-structured and unstructured data

- Serving data for training machine learning models

# Object Storage

**Object Storage**



*Objects are immutable*

UUID

Programs

write

For each object,

- Universal Unique Identifier or UUID (key)

- Metadata: creation date, file type, owner

# Object Storage

**Object Storage**



write

Programs

*Enable Versioning*

UUID

For each object,

- Universal Unique Identifier or UUID (key)
- Metadata: creation date, file type, owner, version

# Why Use Object Storage?

- Store files of various data formats without a specific file system structure

- Easily scale out to provide virtually limitless storage space

- Replicate data across several availability zones



Amazon S3

99.999999999% : data durability

- Cheaper than other storage options

# Introduction to Source Systems

## Logs

# Logs

```
01-01-2025:10.30    67945    success     user added a product x to their cart
01-01-2025:10.32    38910    fail        invalid values typed for product quantity
01-01-2025:10.38    17462    fail        customer table corrupted
```

Software
Application

Exhaust / Byproduct

## Monitoring or Debugging a system

- User activity:
  - Signing in
  - navigating to a particular page
- An update to a database
- An error from a procedure

DeepLearning.AI

**Log**

An append-only sequence of records ordered by time, capturing information about events that occur in systems.

| timestamp | user id | status | action |
|---|---|---|---|
| 01-01-2025:10.30 | 67945 | success | user added a product x to their cart |
| 01-01-2025:10.32 | 38910 | fail | invalid values typed for product quantity |
| 01-01-2025:10.38 | 17462 | fail | customer table corrupted |

Event timestamp

Person, system, or account associated with the event

Event & event metadata

DeepLearning.AI

**Log**

An append-only sequence of records ordered by time, capturing information about events that occur in systems.

| timestamp | user id | status | action |
|---|---|---|---|
| 01-01-2025:10.30 | 67945 | success | user added a product x to their cart |
| 01-01-2025:10.32 | 38910 | fail | invalid values typed for product quantity |
| 01-01-2025:10.38 | 17462 | fail | customer table corrupted |

```
{
  "user id": 67945,
  "action": "user added a product x to their cart",
  "status": "success",
  "time-stamp": 01-01-2025:10.30
}
```

[00101011 11000101 11001001 11000101 110001001]

| user id | action | status | timestamp |
|---|---|---|---|
| 67945 | user added a product x to their cart | success | 01-01-2025:10.30 |
| 38910 | invalid values typed for product quantity | fail | 01-01-2025:10.32 |
| 17462 | customer table corrupted | fail | 01-01-2025:10.38 |

DeepLearning.AI

# Log Levels

**A tag to categorize the event (log level)**

- "debug"

- "info"

- "warn"

- "error"

- "fatal"

| user id | action | status | timestamp | level |
|---------|--------|--------|-----------|-------|
| 67945 | user added a product x to their cart | success | 01-01-2025:10.30 | Info |
| 38910 | invalid values typed for product quantity | fail | 01-01-2025:10.32 | error |
| 17462 | customer table corrupted | fail | 01-01-2025:10.38 | fatal |

# Introduction to Source Systems

## Streaming Systems

DeepLearning.AI

# Terminology

**Event**

Something that happened in the world or a change to the state of a system.

User clicking on a link

Sensor measuring a temperature change

P**Event**er

Data: record of events

DeepLearning.AI

# Terminology

**Message**

A record of information about an event.

**Message**
Event Details
Event Metadata
Event Timestamp

Producer

Event    Event    Event    Event

# Terminology

**Stream**

A sequence of messages.

**Message**
Event Details
Event Metadata
Event Timestamp

Producer

Event  Event  Event  Event ← **Stream**

# Stream Processing



Producer → Batch Processing

Streaming system

Producer → Process message as it is received

# Streaming System

Message
Producer

→

Message Router /
Streaming Broker

→

Message
Consumer

# Streaming System



Event Producer

**Generates the messages**

IoT Device    Mobile App

API    Website

Event Router / Streaming Broker

Event Consumer

**Processes the messages**

Event Consumer

DeepLearning.AI

# Streaming System

Event
Producer

**Order system**

Event Router /
Streaming Broker

Event
Consumer

**Payment service**

Event
Consumer

**Inventory service**

DeepLearning.AI

# Streaming System



Event Producer

*Producer can send messages anytime*

Event Router / Streaming Broker

- Acts as a buffer to filter and distribute the messages

- Decouples producer from consumer

- Prevents message from being lost

Event Consumer

Event Consumer

DeepLearning.AI

# Your Data System



Source System

Event Producer

Event Router / Streaming Broker

Event Consumer

Event Consumer

The system you build

DeepLearning.AI

# Your Data System



Event Producer

Event Router / Streaming Broker

Event Consumer

Event Consumer

**Source System**

**Your ingestion pipeline starts here**

DeepLearning.AI

# Streaming System

**Message Queue**

A queue/buffer that accumulates messages

**Message Queue**

Event
Producer

`1`

**Message Queue**

A queue/buffer that accumulates messages
and delivers those messages to consumers asynchronously.

**Message Queue**

Producer

| 5 |

| 4 | | 3 | | 2 |

Event
Consumer

First-in first-out (FIFO) basis

DeepLearning.AI

**Message Queue**

A queue/buffer that accumulates messages and delivers those messages to consumers asynchronously.

**Message Queue**

Event Producer → 

| 5 | 4 | 3 | 2 |

→ Event Consumer

First-in first-out (FIFO) basis

*Temporary storage*

Amazon Simple Queue Service (Amazon SQS)

# Connecting to Source Systems



Data Sources

Data Engineer

- Improper identity and access managements (IAM) definitions
- Broken networking configurations
- Wrong set of access credentials

DeepLearning.AI

# Lesson's Plan

**1**    Ways of connecting to source systems

**2**    IAM roles and permissions

*Key to controlling and managing access to cloud-based data sources*

Role

Permissions

# Lesson's Plan

**1**    Ways of connecting to source systems

**2**    IAM roles and permissions

*Key to controlling and managing access to cloud-based data sources*

**3**    Basics of networking

*VPCs and Subnets, Gateways, Routing, Security groups*

aws

Role

Permissions

# Lesson's Plan

**1**    Ways of connecting to source systems

**2**    IAM roles and permissions

*Key to controlling and managing access to cloud-based data sources*

Role

Permissions

**3**    Basics of networking

*VPCs and Subnets, Gateways, Routing, Security groups*

**Real world scenario**

**4**    Lab exercise: put your skills to the test

Your job: troubleshoot and figure out the cause of the problem

# Connecting to Source Systems

```python
def create_client():
        dynamodb_client = boto3.client("dynamodb")

        return dynamodb_client
```

boto3: AWS Software Development Kit (SDK) for Python

# Connecting to Source Systems

Running this command in Cloud9 IDE

```
mysql --host=<MySQLEndpoint> --user=<DatabaseUserName> --password=<Password> --port=3306
```

# Programmatic Way

## SDK (boto3)



### IDE (Cloud 9)

```
connect_to_source.py ×        ⊕

1   import pymysql
2   import boto3
3
4   ENDPOINT="............"
5   PORT="3306"
6   USER="jane_doe"
7   REGION="us-east-1"
8   DBNAME="mydb"
9
10  #gets the credentials from .aws/credentials
11  session = boto3.Session(profile_name='default')
12  client = session.client('rds')
```

### Jupyter Notebook

```
import pymysql
import boto3

ENDPOINT="............"
PORT="3306"
USER="jane_doe"
REGION="us-east-1"
DBNAME="mydb"

#gets the credentials from .aws/credentials
session = boto3.Session(profile_name='default')
client = session.client('rds')
```

# API Connector

Interacting with Source Systems

Basics of IAM and Permissions

# Security on the Cloud



Encryption Methods

Identity and Access Management (IAM)

Networking Protocols

**We're only human:**
The #1 root cause of cloud data breaches is human error

- Insecure storage of passwords

- IAM misconfigurations

DeepLearning.AI

# Mistakes



Confidential data

**Public** S3 Bucket

**GitHub**

Access Credentials

Admin access

# IAM

**IAM is a framework for managing permissions.**

Permissions define which actions an identity
can perform on a specific set of resources



Person

Application

# Principle of Least Privilege

# Principle of Least Privilege

# Principle of Least Privilege



Read from specific tables

Ingestion System

# AWS IAM



IAM
services

AWS Identity and Access
Management (IAM)

# AWS IAM

**Root User** — Has unrestricted access to all resources

**IAM User** — Has specific permissions to certain resources
- Username & password
- Access key

**IAM Group** — A collection of users that inherit the same permission from the group policy

**IAM Role** — A user, application, or service that's been granted temporary permissions



AWS account

Root User

**Identities**

IAM User

IAM Group

IAM Role (User/ application/ service)

Policies

**Resources**

Amazon S3

Amazon RDS

Amazon EC2

DeepLearning.AI

# AWS IAM

Role



Amazon EC2

**Not allowed to read from or write to**

Amazon S3

# AWS IAM

Role



Amazon EC2

**Allowed to read from or write to**

Amazon S3

More secure than storing long-term user credentials within the EC2 configurations

**ACCESS DENIED**

Check if credentials have expired!

# IAM Policies

**permission to access the specified S3 buckets**

**permission to access the AWS Glue job**

```
{
  "Version":"2012-10-17",
  "Statement":[
    {
      "Sid": "S3AccessDLAIBucket",
      "Action": [
              "s3:List*",
              "s3:Get*"
      ],
      "Effect":"Allow",
      "Resource": [
              "arn:aws:s3:::dlai-data-engineering",
              "arn:aws:s3:::dlai-data-engineering/*"
      ]
    },
    {
      "Sid": "GlueMgmt",
      "Action": [
              "glue:*"
      ],
      "Effect":"Allow",
      "Resource": [
              "arn:aws:glue:*:*:catalog",
              "arn:aws:glue:*:*:*/de-c1w2*"
      ]
    }
  ]
}
```

Interacting with Source Systems

**Basics of Networking**

DeepLearning.AI

# AWS Cloud

**What does cloud in "cloud computing" mean?**

The "cloud" is made up of very real physical data centers that are spread out around the world.



Each dot represents a region

Screenshot from AWS Global Infrastructure (2023)

DeepLearning.AI

# AWS Cloud



Resources are replicated across availability zones to ensure that your systems keep working even if a data center goes down.

# AWS Cloud

**Region considerations:**

- Legal compliance

- Latency *: the closer your end users are to the region, the lower the latency*

- Availability *: the more availability zones, the better you will be able to recover from a disaster*

- Cost

# AWS Cloud

**Region considerations:**

- Legal compliance

- Latency

- Availability

- Cost

# Virtual Private Cloud



**Virtual Private Cloud (VPC)**

Smaller networks that span multiple availability zones

# Virtual Private Cloud

# Virtual Private Cloud



For internet-facing resources

For internal resources

VPC

Availability Zone

Public subnet

Private subnet

ACL

ACL

172.16.0.0
172.16.1.0
172.16.2.0

172.16.0.0
172.16.1.0
172.16.2.0

Internet Gateway

Internet

DeepLearning.AI

# Virtual Private Cloud

# Virtual Private Cloud

Interacting with Source Systems

**AWS Networking - VPCs & Subnets**

# Example Scenario



DeepLearning.AI

# AWS Networking - VPCs & Subnets

AWS Cloud

VPC

Availability Zone 1 | Availability Zone 2

VPC

Needs to be configured

DeepLearning.AI

# AWS Networking - VPCs & Subnets

Define the network    Host addresses

| 10 | 0 | 0 | 0 | /24 |

8 bit    8 bit    8 bit    8 bit

32 bits

16 bits

prefix length

*How many bits used for the network part of the address*

Amazon RDS

EC2 Instance

10.0._._    10.0._._

0 to 255

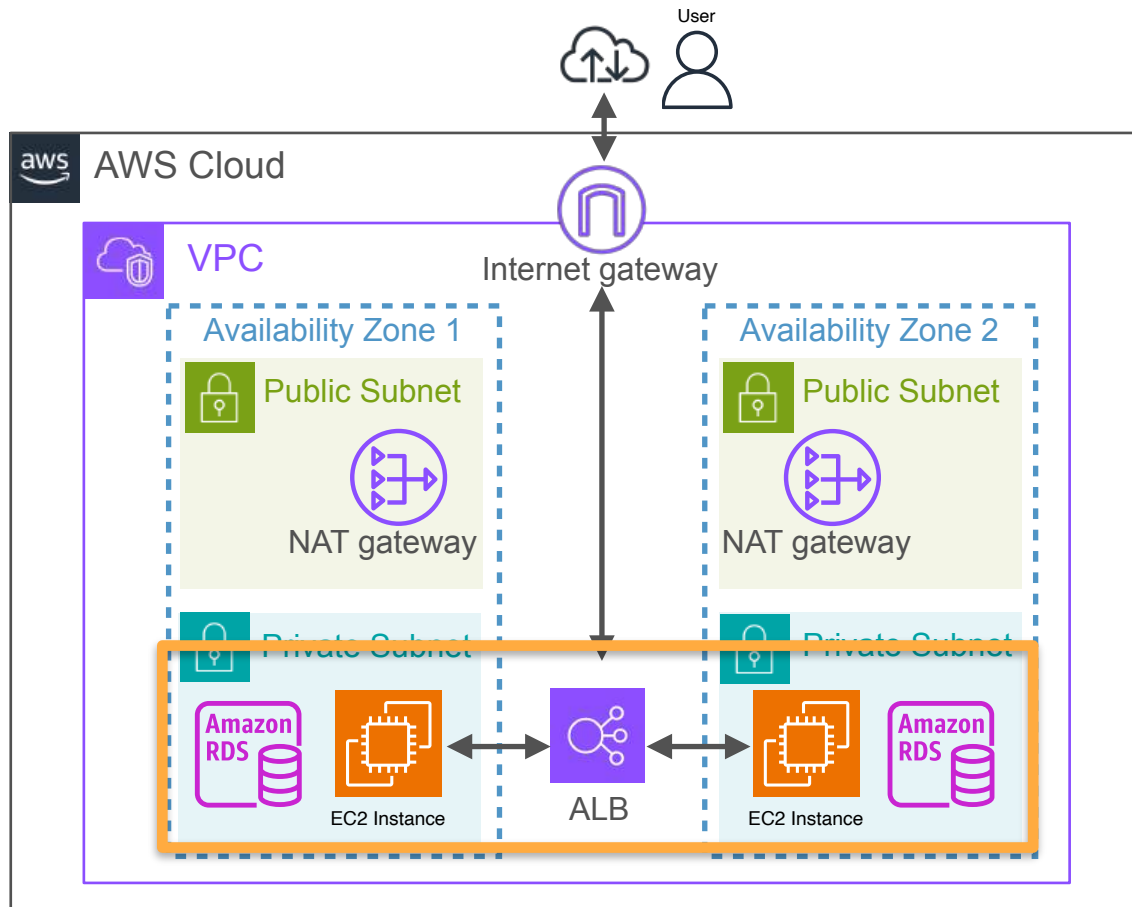# AWS Networking - VPCs & Subnets

# Example Scenario



## Considerations

1. Applications running on EC2 need to occasionally download updates from resources on the internet
   - Upgrades, patching

2. Need a way to submit requests to the application running on the EC2 instance

## Considerations

1. Applications running on EC2 need to occasionally download updates from resources on the internet
   - Upgrades, patching

2. Need a way to submit requests to the application running on the EC2 instance

# Considerations

1. Applications running on EC2 need to occasionally download updates from resources on the internet
   - Upgrades, patching

2. Need a way to submit requests to the application running on the EC2 instance

# Considerations

1. **Applications running on EC2 need to occasionally download updates from resources on the internet**
   - **Upgrades, patching**

2. Need a way to submit requests to the application running on the EC2 instance

**NAT Gateway** — **Network Address Translation Gateway**

- Allows resources in a private subnet to connect to the internet or other AWS services
- Prevents the internet from initiating connections with those resources

DeepLearning.AI

# Considerations

1. Applications running on EC2 need to occasionally download updates from resources on the internet
   - Upgrades, patching

2. Need a way to submit requests to the application running on the EC2 instance

**ALB:**

- Distributes incoming application traffic across multiple backend targets

- Handles the load and ensures the application remains responsive and available

- Keeps those EC2 instances private

DeepLearning.AI

Route Tables

- Essential for directing network traffic within your VPC

- Default route table allows internal communication within the VPC

DeepLearning.AI

# Security Groups

Instance level virtual firewalls, controlling both inbound and outbound traffic
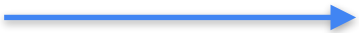
deny

inbound traffic

**Security Groups**

allow

outbound traffic

## Inbound Rules

- Determine what types of traffic you want to allow
- Where you want to allow that traffic to come from

**Port 80**

inbound HTTP requests
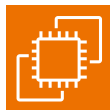
outbound HTTP responses

## Security Groups are stateful

- Allow inbound traffic to an instance automatically allows the return traffic

DeepLearning.AI

**Security Groups** — Instance level virtual firewalls, controlling both inbound and outbound traffic

ALB

EC2 Instance

Amazon RDS

**Security group ID: sg-123**

| Source | Protocol | Port |
|---|---|---|
| 0.0.0.0/0 (internet) | HTTP | 80 |
| 0.0.0.0/0 (internet) | HTTPS | 443 |

**Security group ID: sg-456**

| Source | Protocol | Port |
|---|---|---|
| sg-123 (ALB) | HTTP | 80 |
| sg-123 (ALB) | HTTPS | 443 |

**Security group ID: sg-789**

| Source | Protocol | Port |
|---|---|---|
| sg-456 (EC2) | TCP | 3306 |

**Security Group Chaining**

DeepLearning.AI

**Network Access Control Lists (ACL)**

- They provide an additional layer of security at the subnet level

- Network ACLs are stateless

- You need to define both inbound and outbound rules explicitly

- Useful for implementing security policies at the subnet level

**If you encounter connectivity issues:**

1. Verify that your VPC has an internet gateway properly attached

2. Verify that the route tables have appropriate rules to direct traffic correctly

3. Verify that the route table associations with the subnets are configured correctly

4. Check security groups to make sure they have the needed rules in place

5. Review network ACLs to confirm they allow the necessary traffic

6. Double-check instance configurations to ensure they are associated with the correct security groups and subnets
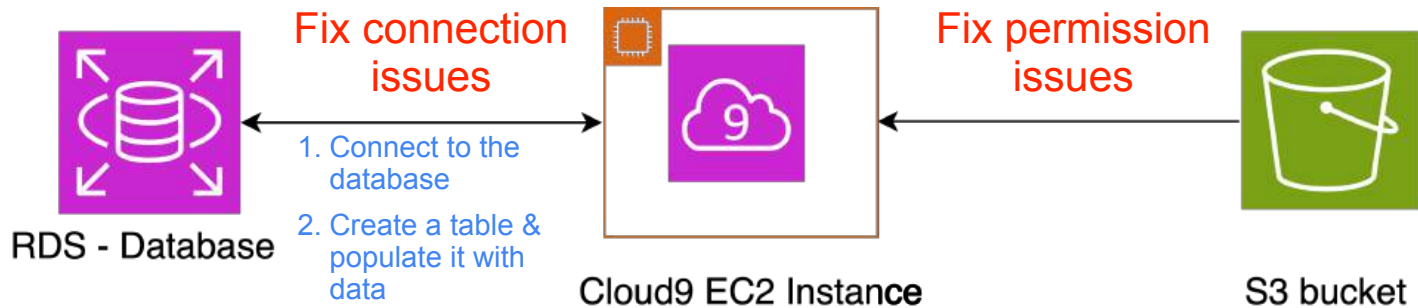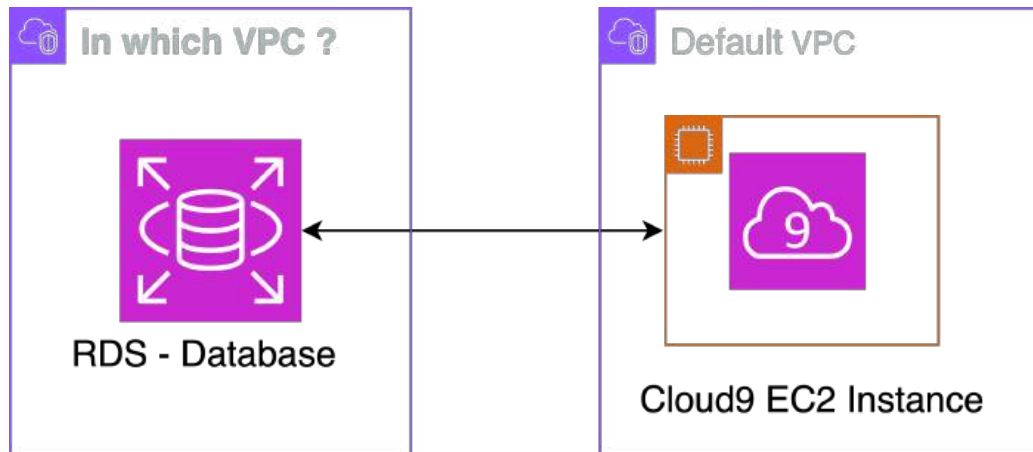
# Database Connectivity and Troubleshooting on AWS



Fix connection issues

1. Connect to the database
2. Create a table & populate it with data

RDS - Database

Cloud9 EC2 Instance

Fix permission issues

S3 bucket

- Skip this video, jump straight into the lab and go for it
  - The lab instructions contain hints
- Or, start the lab and follow along with me as you go through this video.
  - When an issue occurs, I'll be inviting you to pause the video
  - After that, I'll show you how to fix it.
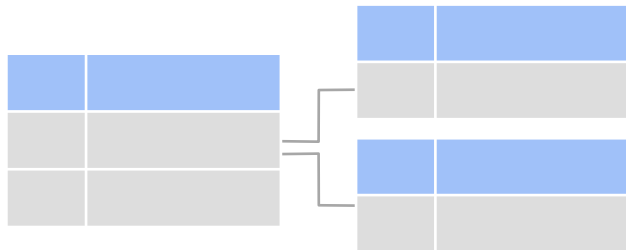
# Database Connectivity and Troubleshooting on AWS

# Week 1 Summary

## Understand how source systems work

### Relational databases



### NoSQL databases

| Key | Value | Document |
|-----|-------|----------|
| k1 | value 1 | |
| k2 | value 2 | |
| k3 | value 3 | |

```
{
    "firstName": "Joe",
    "lastName" : "Reis",
    "age": 10,
    "address": {
        "city": "Los Angeles",
        "postalCode": 90024,
        "country": "USA"
    }
}
```

### Object Storage



### Logs

| user id | action | status | timestamp |
|---------|--------|--------|-----------|
| 67945 | user added a product x to their cart | success | 01-01-2025:10.30 |
| 38910 | invalid values typed for product quantity | fail | 01-01-2025:10.32 |

### Streaming Systems

Producer → Consumer

Event Router / Streaming Broker

DeepLearning.AI

# Week 1 Summary

- How to connect to data sources

- Basics of networking

- Importance of IAM in ensuring security in source systems

# Week 1 Summary



Data Sources

- Improper identity and access managements (IAM) definitions
- Broken networking configurations
- Wrong set of access credentials

Data Engineer