# Data Ingestion

**Week 2**

DeepLearning.AI

# Data Ingestion



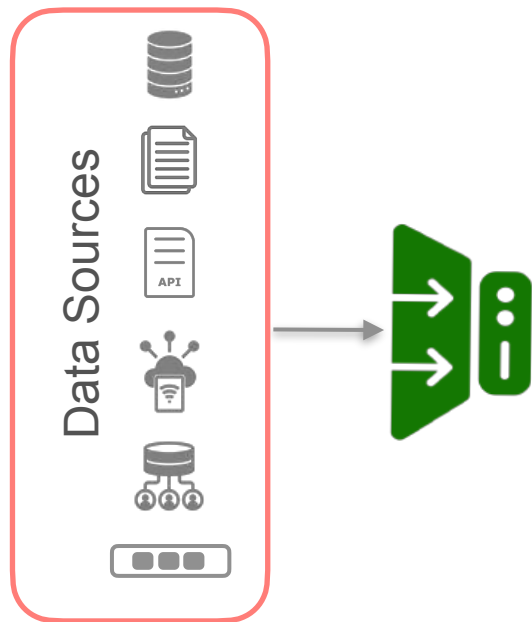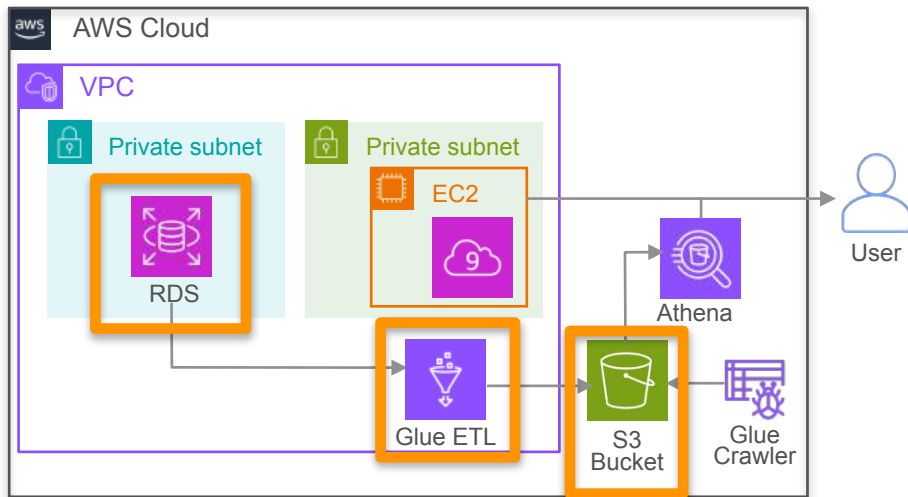Data Sources

**Get raw data** → **Turn it into something useful** → **Make it available for downstream use cases**

# Data Ingestion

## Course 1 Week 2 Lab

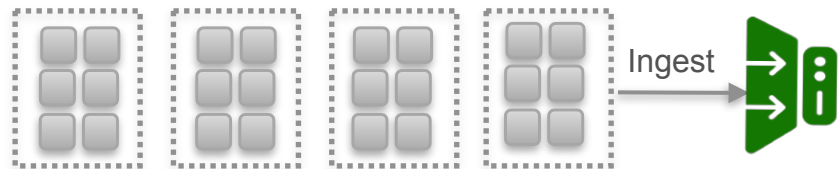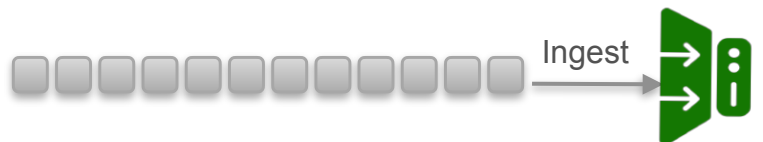

## Course 1 Week 4 Lab



## Course 2 Week 1 Lab



DeepLearning.AI

# Week 2 Overview



**Batch Ingestion** — Ingest — Processing data in chunks or batches

**Streaming Ingestion** — Ingest — Processing a continuous stream of events

# Week 2 Overview

**Lab: streaming ingestion**

**Lab: batch ingestion**

Investigate requirements for **streaming ingestion**:

- the data payload and event rates
- how to configure the streaming pipeline

Identify requirements for **batch ingestion** from a **REST API**

Software Engineers

Data Engineer

Marketing Analyst

Data Ingestion

---

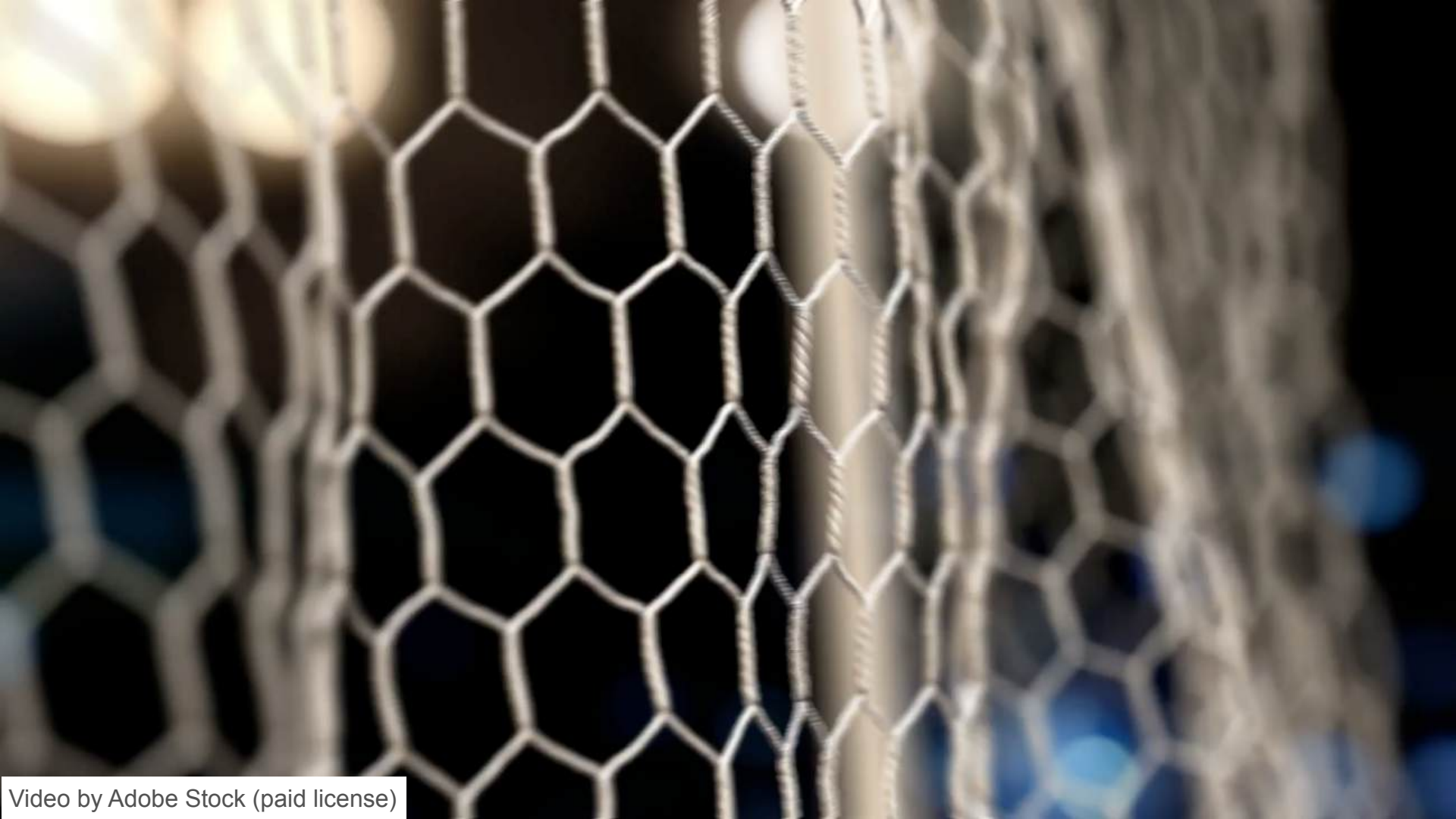**Data Ingestion on a Continuum**

DeepLearning.AI

Subtotal                              US $1.59

Total                              **US $1.59**

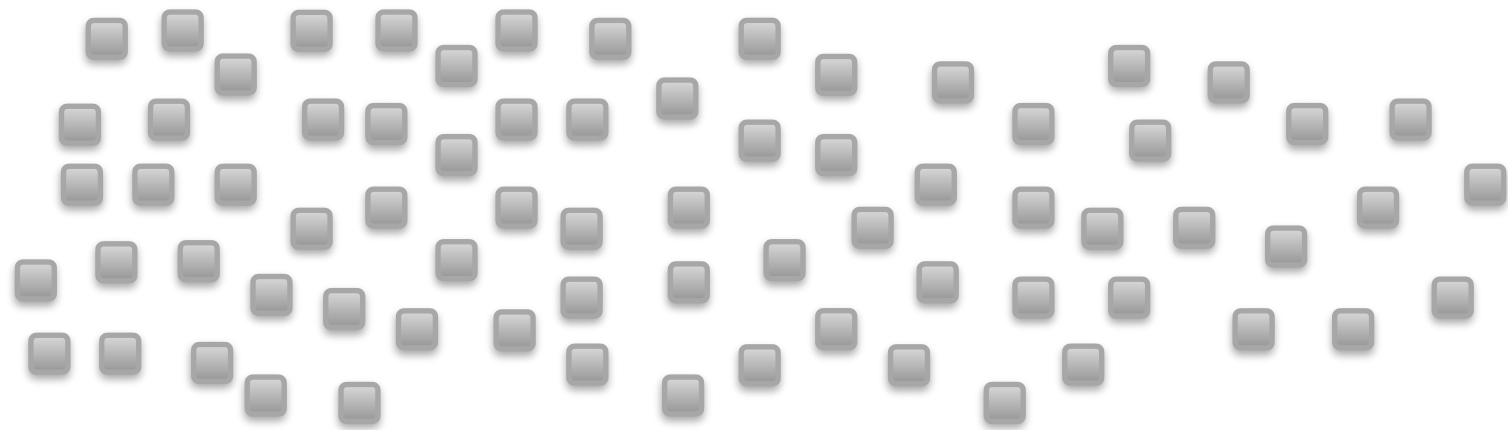(Approximately 44,20 грн.)

**BUY (2)**

Buy from this seller

# Data Ingestion

**Unbounded Data:** continuous stream of events

# Data Ingestion

**Stream Ingestion**



Ingest

# Data Ingestion



Batch Ingestion

# Data Ingestion

## Size-Based Batch Ingestion



10 GB

1000 events

10 GB

1000 events

10 GB

1000 events

10 GB

1000 events

Ingest

**Size-threshold**

**Total number of records**

# Data Ingestion

## Time-Based Batch Ingestion



Ingest

Weekly Chunk                    Weekly Chunk

DeepLearning.AI

# Data Ingestion



**Time-Based Batch Ingestion**

Ingest

**Daily Chunks**

# Data Ingestion



**Time-Based Batch Ingestion**

Ingest

**every hour?**   **every minute?**   **every second?**

DeepLearning.AI

# Data Ingestion

**Stream Ingestion**

Ingest

# Ingestion Frequencies



**Choice of ingestion frequency depends on:**

- the source systems you're working with

- the end use case

# Ways to Ingest Data from Databases



**Connectors**

**JDBC/ODBC API**

- Ingest at regular intervals
- Ingest when a certain amount of new data is recorded

**Ingestion Tool**

AWS Glue ETL

Ingest data on a regular basis

DeepLearning.AI

# Ways to Ingest Data from APIs

**Protocols**

API

Ingest

- How much can you ingest in one go?
- How frequently can you call the API?

Reading API documentation

Communicating with data owners

Writing custom API connection code

# Ways to Ingest Data from Files



Files

Ingest

**Manual File Download**

**Secure File Transfer**

**File Transfer Protocols**

SFTP: Secure File Transfer Protocol

SCP: Secure Copy Protocol

# Ways to Ingest Data from Streaming Systems

**Message Queue or Streaming Platform**

Event Producer → [ ] [ ] [ ] [ ] [ ] → Event Consumer

IoT Device

Batch Ingestion

ETL vs. ELT

# Goals of the Marketing Analyst

**Batch Ingestion**

Frequency of
request is limited

Marketing Analyst

Analysis of
Historical Trends

No need for real-
time analysis

# Batch Ingestion Patterns

**ETL**    **Extract - Transform - Load**

Tradeoffs

**ELT**    **Extract - Load - Transform**

# Batch Ingestion Patterns

**ETL**   **Extract** - **Transform** - **Load**

**ELT**   **Extract** - **Load** - **Transform**

# Batch Ingestion Patterns

**ETL**

Data Sources

**Extract**

raw data

**Transform**

Staging Area

**Load**

transformed data

Target Destination

**ELT**

**Extract** - **Load** - **Transform**

# Batch Ingestion Patterns

**ETL**

Data Sources
- (database)
- (document)
- API
- (IoT)
- (distributed)

**Extract**

raw data

**Transform**

Staging Area

**Load**

transformed data

Target Destination

**ELT**

**Extract** - **Load** - **Transform**

# Emergence of Cloud Storage Systems

**Early 2010s: Highly scalable cloud storage**

**Data Lake**

*built on top of object storage*

**Cloud
Data Warehouse**

Amazon Redshift

snowflake

- Store enormous amounts of data for relatively cheap
- Perform data transformations directly in the data warehouse

# Batch Ingestion Patterns

**ETL**

Data Sources



**Extract**

raw data

**Transform**

Staging Area

**Load**

transformed data

Target Destination

**ELT**

**Extract** - **Load** - **Transform**

DeepLearning.AI

# Batch Ingestion Patterns

# Advantages of Extract-Load-Transform

It is faster to implement.

It makes data available more quickly to end users.

Transformations can still be done efficiently.

You can decide later to adopt different transformations.

# Downsides of Extract-Load-Transform

HOW will you use the data?

**EL**

Pipeline

Data Sources

**Extract**

raw data

**Load**

raw data

Target Destination

**Transform**

Query data

**Data Swamp**

Data has become unorganized, unmanageable, and essentially useless.

# Downsides of Extract-Load-Transform

**Data Swamp**



Video by Adobe Stock (paid license)

# Conversation with the Marketing Analyst



Data Engineer

Marketing Analyst

API

Ingest

JSON   TXT   PNG

**Exploratory analysis of the data**

?

**ELT**

DeepLearning.AI

# Batch Ingestion

---

## REST API

# API Mandate



Amazon Web Services (AWS)

Application
Programming Interface

Stable & Predictable

# What is an API?

**API**  A set of rules and specifications that allows you to programmatically communicate and exchange data with an application.



Application

respond

request

API

respond

request

Code

Built into a wide range
of software applications

# What is an API?

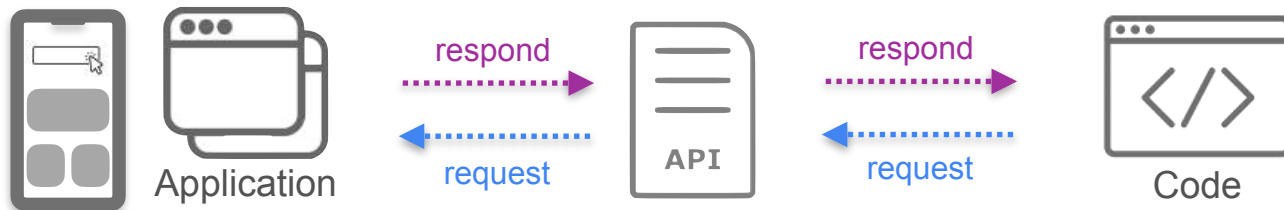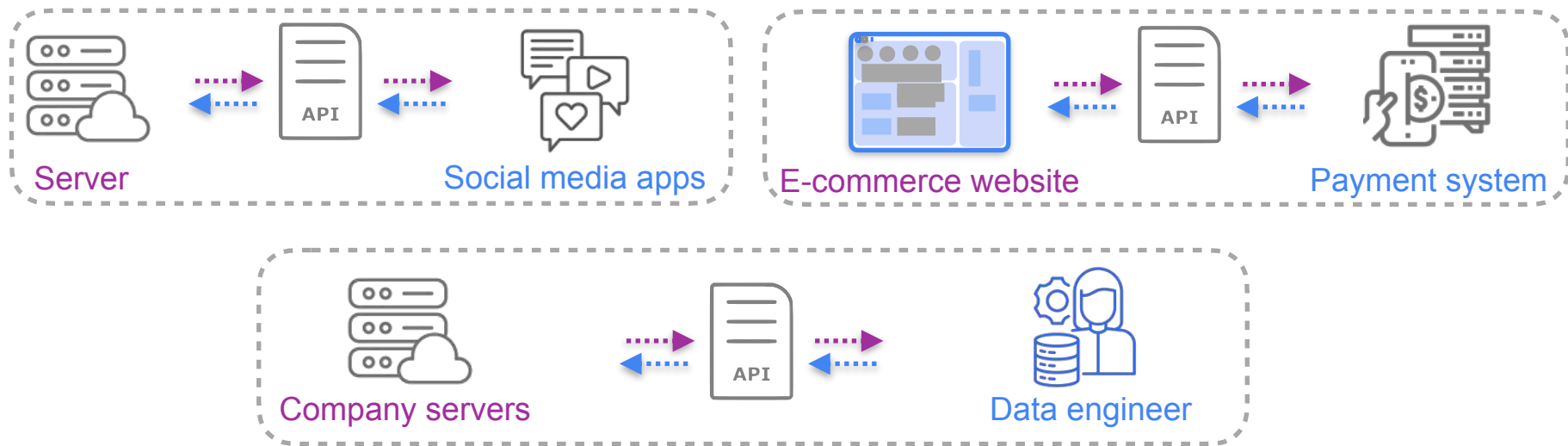API — A set of rules and specifications that allows you to programmatically communicate and exchange data with an application.



Server    Social media apps    E-commerce website    Payment system

Company servers    Data engineer

# What is an API?



**API Features**

- Metadata
- Documentation
- Authentication
- Error handling

# REST API

**REST API**

Representational State Transfer API

Use Hypertext Transfer Protocol (HTTP) as the basis for communication

**Respond with the requested resource**

Web server

**HTTP Request**

Link

**Respond to the HTTP request**

REST API

**HTTP Request**

Data engineer

# Batch Processing from an API

**Upcoming Lab**

Spotify®

**In this video**,

- Extract data from the Spotify API
- Explore what pagination means
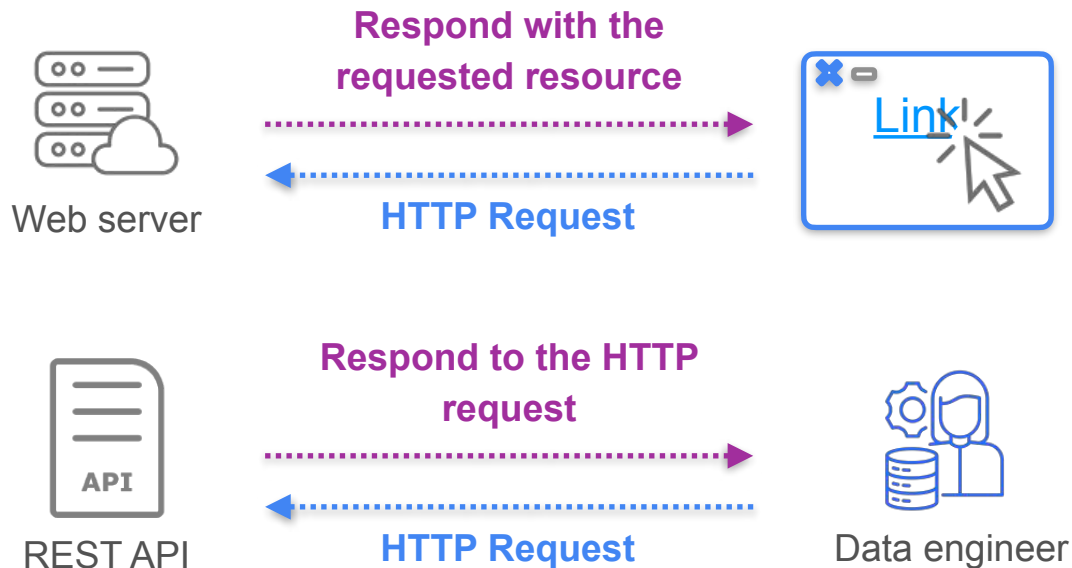- Send an API request that requires authorization

- Go through some API concepts
- Give you an overview of the lab tasks

**What you need**

- Spotify account   https://developer.spotify.com/

  When working with an API, it's very common that you'll have to sign up for an account

- Spotify Documentation   https://developer.spotify.com/documentation/web-api

# API Concepts

**Spotify Web API**

Restful API



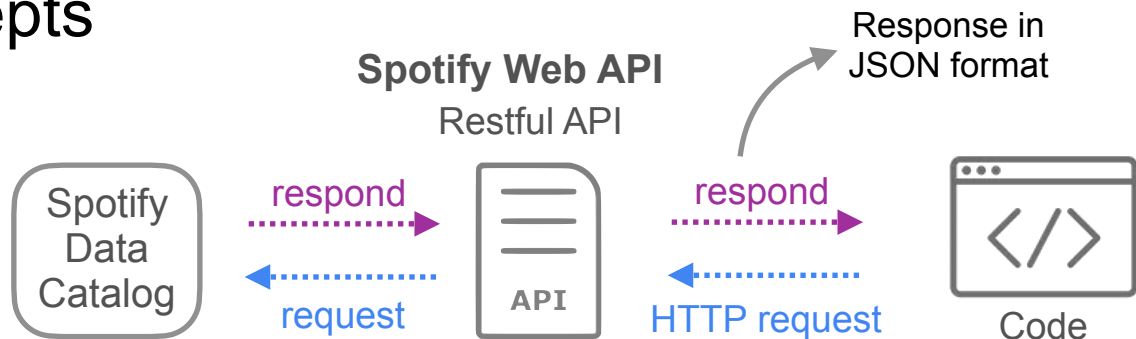Spotify Data Catalog → respond → API → respond → Code

← request ←        ← HTTP request ←

**Resource**
- Music
- Artists
- Albums
- Tracks
- Playlist

Each resource is represented by an **endpoint**.

| HTTP Request | Action |
|--------------|--------|
| GET | Retrieve a resource |
| POST | Create a resource |
| PUT | Change/Replace a resource |
| Delete | Delete a resource |

# API Concepts

**Spotify Web API**
Restful API

Response in JSON format



Spotify Data Catalog

respond →
← request

API

respond →
← HTTP request

Code

Each resource is represented by an **endpoint**.

**Resource**
- Music
- Artists
- Albums
- Tracks
- Playlist

| HTTP Request | Action |
|---|---|
| GET | Retrieve a resource |
| POST | Create a resource |
| PUT | Change/Replace a resource |
| Delete | Delete a resource |

# API Concepts

**Spotify Web API**
Restful API

Error object containing a status code:
- 400: bad request
- 404: requested resource not found

Spotify Data Catalog

respond →

← request

API

respond →

← HTTP request

Code
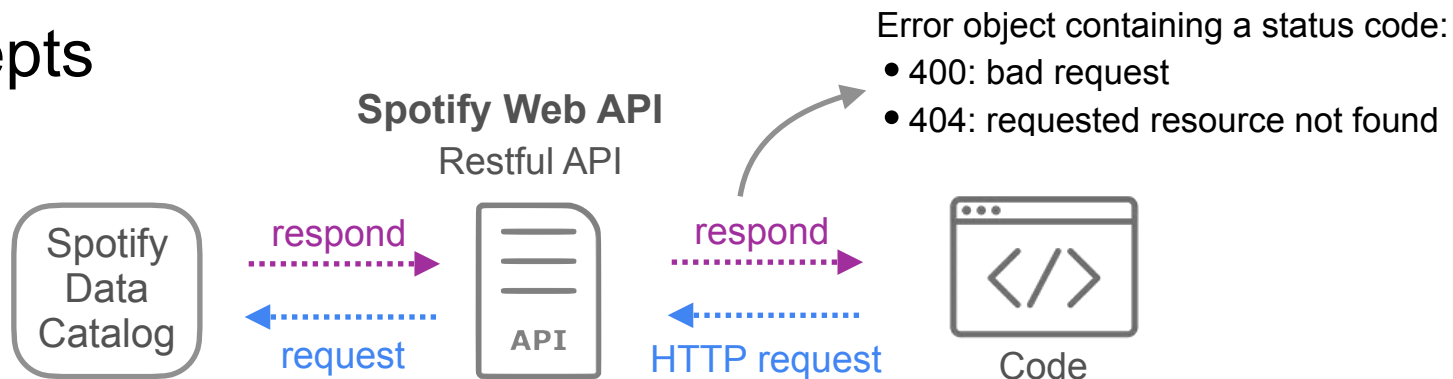
Each resource is represented by an **endpoint**.

**Resource**
- Music
- Artists
- Albums
- Tracks
- Playlist

| HTTP Request | Action |
|---|---|
| GET | Retrieve a resource |
| POST | Create a resource |
| PUT | Change/Replace a resource |
| Delete | Delete a resource |

DeepLearning.AI

# API Concepts

**Spotify Web API**
Restful API

Spotify
Data
Catalog

respond →

← request

API

respond →

← HTTP request

Code

: Endpoint + Access token

**Resource**
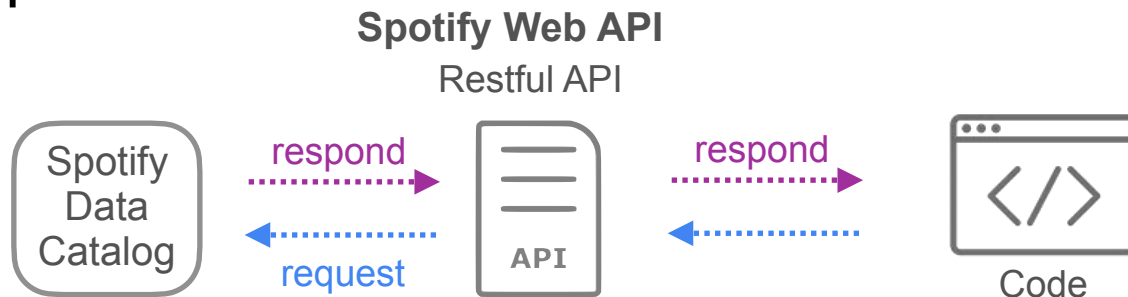- Music
- Artists
- Albums
- Tracks
- Playlist

Each resource is
represented by
an **endpoint**.

**Access token**: string that contains the permissions to access
a given resource. *(valid for 1 hour)*
- create a Spotify account
- get a client ID and a client secret and use them to generate
  the access token *(provided with the code)*

https://developer.spotify.com/documentation/web-api/concepts/authorization

# API Concepts

**Spotify Web API**
Restful API

Spotify Data Catalog → respond → API → respond → Code

request

**https://developer.spotify.com/documentation/web-api** s token

**Resource**
- Music
- Artists
- Albums

Each resource is represented by an **endpoint**.

ontains the permissions to access *hour)*

**Get Playlist** 🔒 OAuth 2.0

Get a playlist owned by a Spotify user.

secret and use them to generate *de)*

**Get Featured Playlists** 🔒 OAuth 2.0

Get a list of Spotify featured playlists (shown, for example, on a Spotify player's 'Browse' tab).

https authorization

**Pagination**    Extract the items chunk by chunk.

- Using offset and limit

```
https://api.spotify.com/v1/me/shows?offset=0&limit=20
```

```
https://api.spotify.com/v1/me/shows?offset=20&limit=20
```
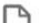
```
https://api.spotify.com/v1/me/shows?offset=40&limit=20
```

```
https://api.spotify.com/v1/me/shows?offset=60&limit=20
```

```
https://api.spotify.com/v1/me/shows?offset=80&limit=20
```

- Using the next field

```
response.get('playlists').get('next')
```

/ src /

| Name | Last Modified |
|------|---------------|
| authenticat... | 4 months ago |
| endpoint.py | 4 months ago |
| env | 56 minutes ago |
| main.py | 3 months ago |

Contains the scripts of the get_token function

1. Paginated call to the endpoint "Get featured playlists"
2. Paginated call to the endpoint "Get playlist"
3. Automatically generate a new token

1. Get the ids of the featured playlists
2. Extract the track information for each playlist id

# Streaming Systems



Message Queues or
Event streaming platforms

Event
Producer

Event
Consumer

Source System

Data engineer

Your ingestion
pipeline starts
here

DeepLearning.AI

# Streaming Systems



Message Queues or
Event streaming platforms

Event
Producer

Event
Consumer

**Source System**

Data engineer

**Your ingestion
pipeline starts
here**

**Message Queue**

**Event Streaming Platform**

# Event Streaming Platforms

**In this week's lab:**



Amazon Kinesis
Data Streams

**In this video:**



Apache Kafka

# kafka

Open-source event streaming platform

Event Producer → **push messages** → **Kafka Cluster** (contains servers / brokers)

**Kafka Cluster**
contains servers / brokers

| Topic 1 | Topic 2 | Topic 3 |
|---|---|---|
| **Fraud alerts** | **Customer orders** | **Temperature readings from IoT** |

**push messages**

**pull messages**

Event Consumer

Event Consumer

Event Consumer

**Topics**:
Categories to hold related events

**kafka** Open-source event streaming platform

Event Producer

Event Producer

Event Producer

**push messages**

**Kafka Cluster**
contains servers / brokers

Topic 1

Fraud alerts

Topic 2

Customer orders

Topic 3

Temperature readings from IoT

**pull messages**

Event Consumer

Event Consumer

Event Consumer

**Anatomy of a Kafka Topic**

**Partition 0**

**Partition 1**

**Partition 2**

**Topics**:
Categories to hold related events

**Partitions** (logs):
Ordered immutable sequences of messages

DeepLearning.AI

**kafka** Open-source event streaming platform

Event Producer

Event Producer

Event Producer

**push messages**

**Kafka Cluster**
contains servers / brokers

Topic 1

Fraud alerts

Topic 2

Customer orders

Topic 3

Temperature readings from IoT

**pull messages**

Event Consumer

Event Consumer

Event Consumer

**Topics**:
Categories to hold related events

**Partitions** (logs):
Ordered immutable sequences of messages

**Anatomy of a Kafka Topic**
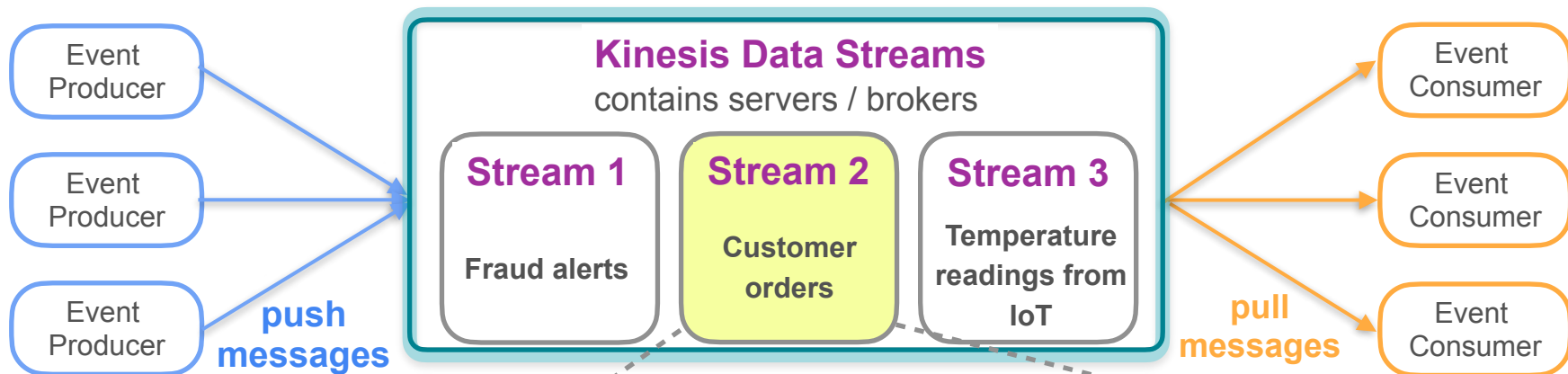
Partition 0

Partition 1

Partition 2

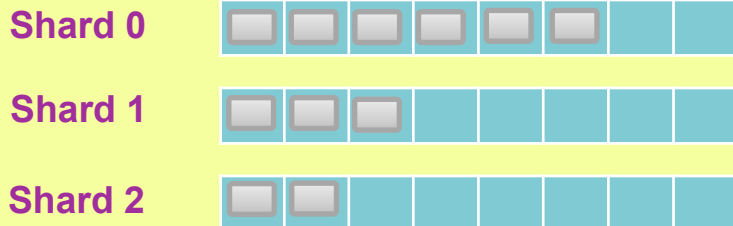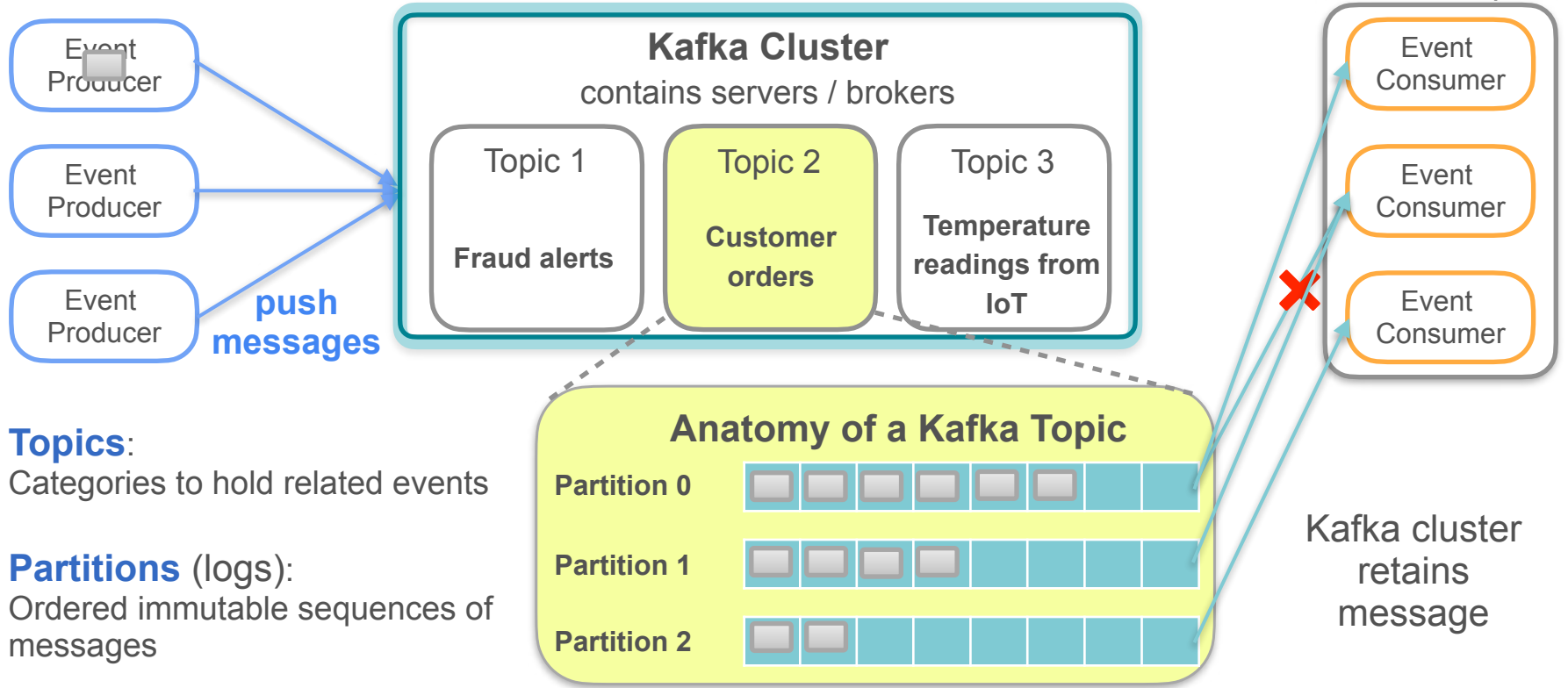Partition decision:

- Round-robin strategy
- Message key

DeepLearning.AI

# Amazon Kinesis Data Stream

**Event Producer**

**Event Producer**

**Event Producer**

**push messages**

**Kinesis Data Streams**
contains servers / brokers

**Stream 1**

Fraud alerts

**Stream 2**

Customer orders

**Stream 3**

Temperature readings from IoT

**pull messages**

**Event Consumer**

**Event Consumer**

**Event Consumer**

**Stream:**
Categories to hold related events

**Shard** (logs):
Ordered immutable sequences of messages

**Anatomy of a Data Stream**

**Shard 0**

**Shard 1**

**Shard 2**

DeepLearning.AI

# Conversation with the Software Engineer



Software Engineers

Data Engineer

User Id: 7945
IP address: 127.168.10.32
Action: User added a product x to their cart
Status: Success
Time Stamp: 01-01-2025:10.30

Web-Server Log

Event
Producer

Amazon Kinesis
Data Streams

kafka

**Ingestion pipeline starts here**

DeepLearning.AI

# Conversation with the Software Engineer

Software Engineers

Data Engineer

**Ingestion pipeline starts here**

User Id: 7945
IP address: 127.168.10.32
Action: User added a product x to their cart
Status: Success
Time Stamp: 01-01-2025:10.30

Web-Server Log
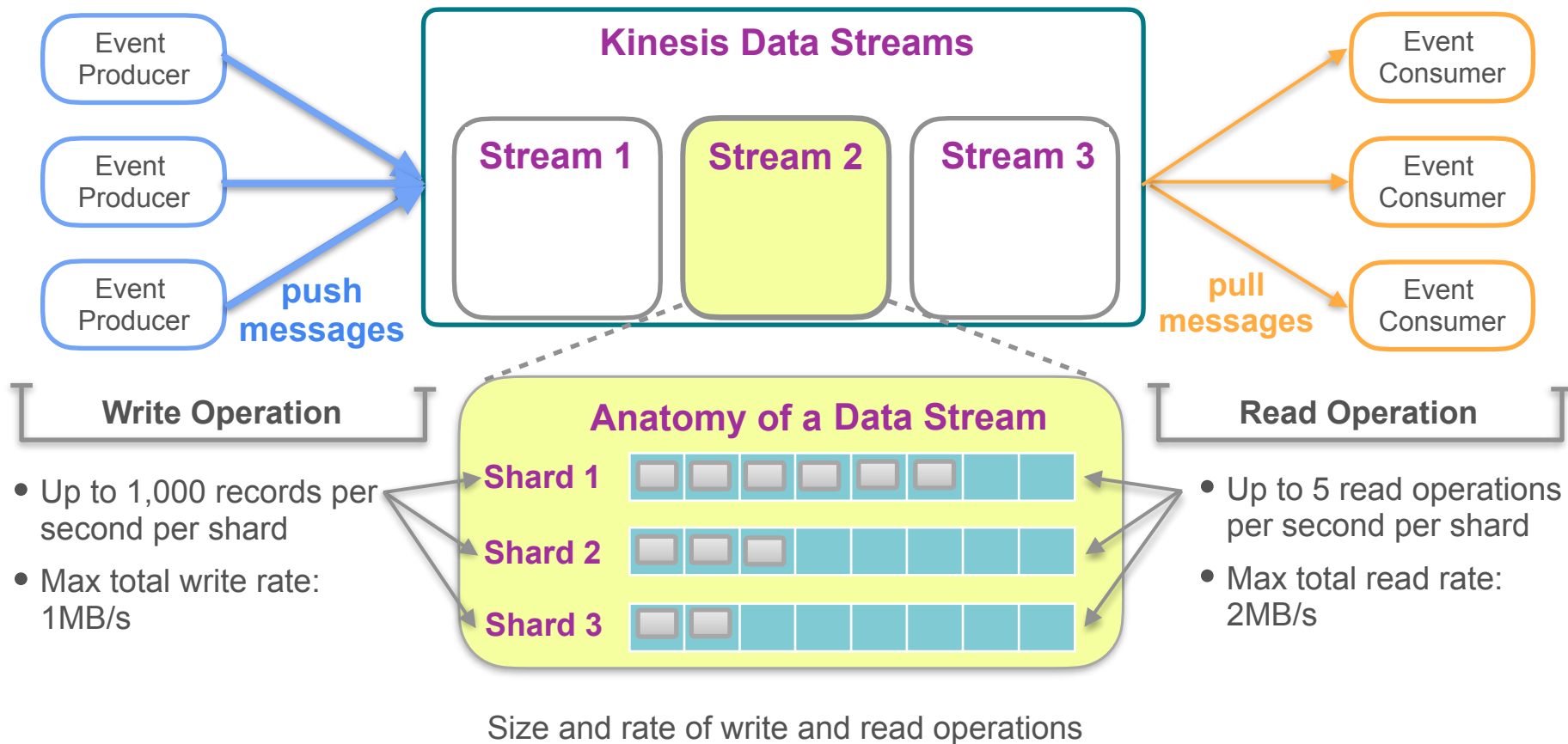
Event Producer

Ingest

Amazon Kinesis Data Streams

kafka

Size and rate of write and read operations

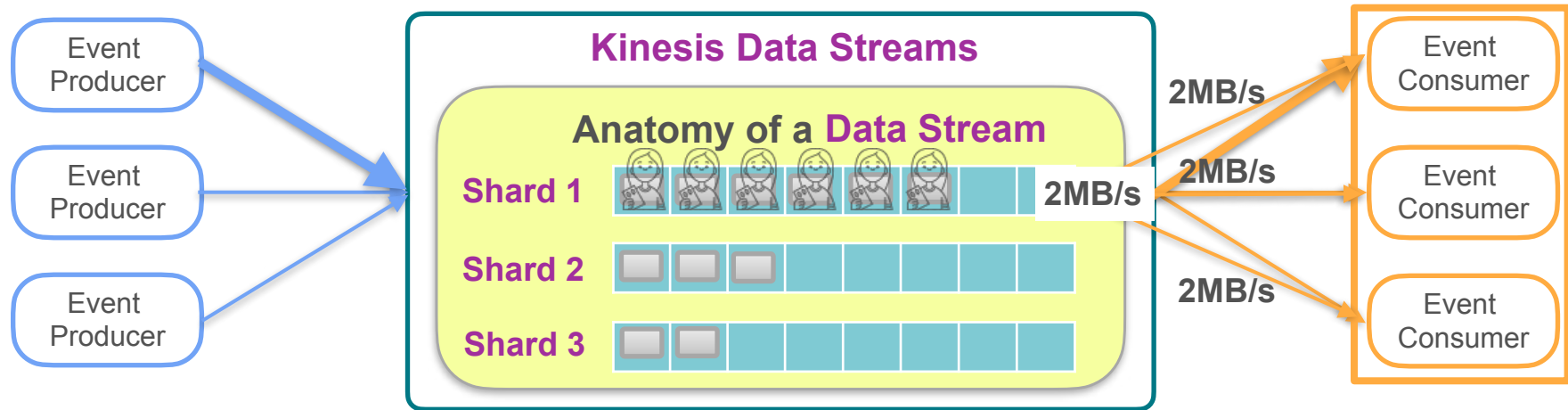| Kinesis in "on-demand" Mode | Kinesis in "Provisioned" Mode |
| --- | --- |
| • Automatically manage the scaling of the shards up or down as needed<br><br>• Only charged for what you use<br><br>• More convenient from an operational perspective | • Specify the number of shards necessary for your application based on the expected write and read request rate<br><br>• Manually add more shards or re-shard when needed<br><br>• A good fit if…<br>   • you have predictable application traffic<br>   • you are able to control your costs more carefully |

# Kinesis Data Streams

## Anatomy of a Data Stream

Shard 1
Shard 2
Shard 3

2MB/s
2MB/s
2MB/s
2MB/s

Event Producer
Event Producer
Event Producer

Event Consumer
Event Consumer
Event Consumer

## Data Record

Partition Key — customerID *12567910*

Sequence Number

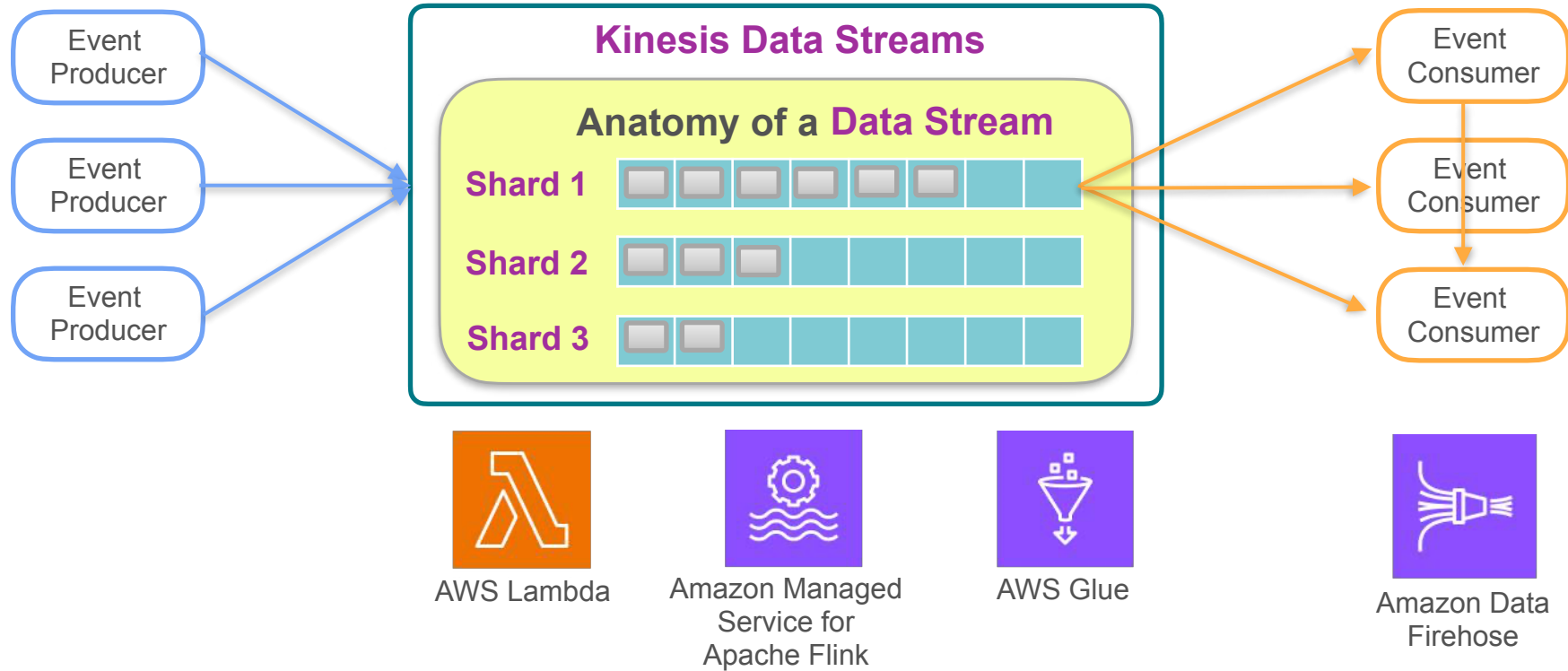Binary Large Object (BLOB)

Used to determine which shard the data record is placed into

**Shared Fan-Out** — When consumers share a shard's read capacity

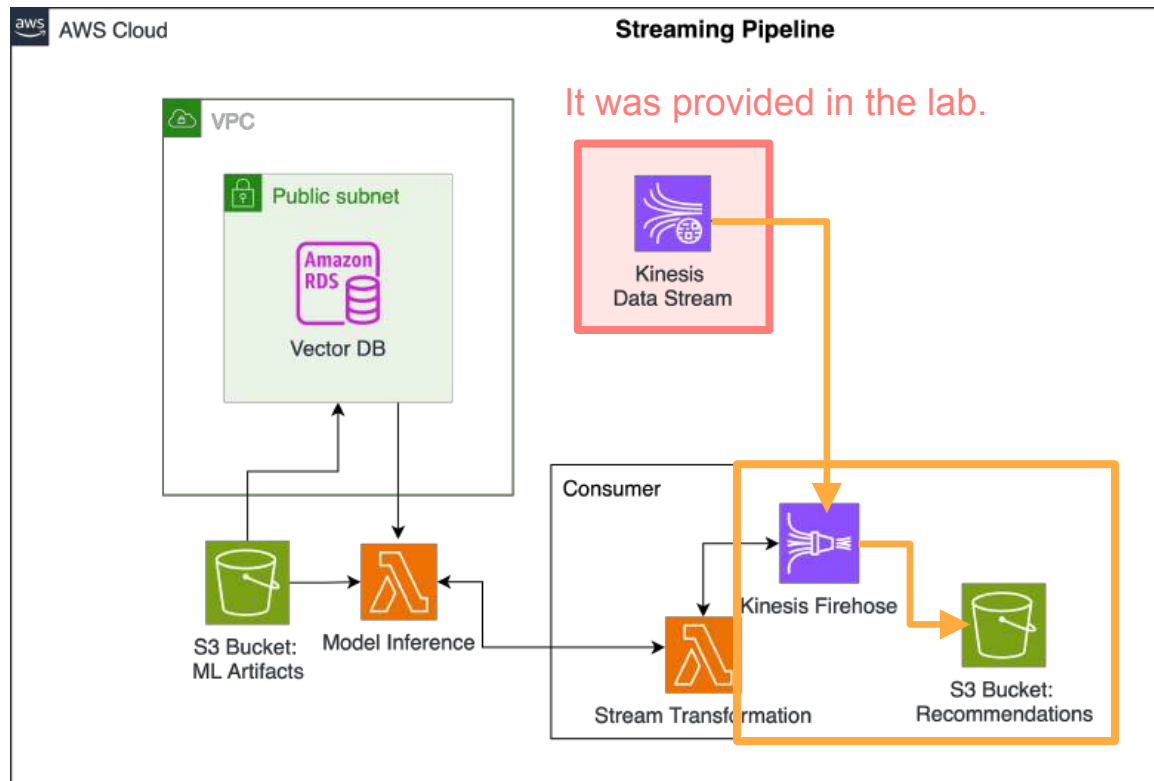**Enhanced Fan-Out** — When consumers are able to read at the full read capacity of the shard

DeepLearning.AI

# Course 1 Lab

# Part 1

Producer → Amazon Kinesis Data Stream → Consumer

Amazon Kinesis
Data Stream

## Streaming Ingestion

In this lab, you will interact with Amazon Kinesis Data Streams and gain a better understanding of how the streaming ingestion process is performed. The lab consists of two parts:

1. You will manually generate data and write it to a Kinesis Data Stream; after that, you will consume the generated data from that stream.
2. You will perform a streaming ETL process: you will consume data from a Kinesis Data Stream that is fed by a producer. You will apply some simple transformations to this data, and then put the transformed data into one of two other data streams. From each of these two new data streams, data will be taken by a Kinesis Firehose and delivered to their respective S3 bucket.

If you get stuck on any of the lab steps, you can check the solution notebook *C2_W2_Lab_1_Streaming_Ingestion_Solution.ipynb* and script *src/etl/consumer_Solution.py* that you can download by running the

**Part 1**

**produ**
- write
- uses
- can

pytho

```python
def main():
    logging.info("Starting PutRecord Producer")
    args = parser.parse_args()


    kinesis_stream_name = args.stream
    data_record = json.loads(args.json_string)


    kinesis = boto3.client("kinesis")


    try:
        # execute single PutRecord request
        response = kinesis.put_record(
            StreamName=kinesis_stream_name,
            Data=json.dumps(data_record).encode("utf-8"),
            PartitionKey=data_record["session_id"],
        )
        logging.info(
            f"Produced record {response['SequenceNumber']} to Shard {response['ShardId']}"
        )
```

⟲ DeepLearning.AI

# Part 1

**producer_**

- writes a s
- uses bot
- can be ru

```
python pr
        --
        --
```

>eam>

```python
def poll_shards(kinesis, shard_iterators):
    """This function continuously polls the shards for data. It iterates
    over the list of shard iterators, fetching records from each shard using
    the respective iterator. For each record retrieved, it logs the order
    data along with the shard ID and sequence number. It updates the shard
    iterator to the next iterator if available.

    Args:
        kinesis (boto3 client): Boto3 client for kinesis resources
        shard_iterators (List): Pair of ShardId and corresponding Iterator
    """
    while True:
        for shard_itr in shard_iterators:
            try:
                records_response = kinesis.get_records(
                    ShardIterator=shard_itr.iterator, Limit=200
                )
                for record in records_response["Records"]:
                    order = json.loads(record["Data"].decode("utf-8"))
                    logging.info(
                        f"Read Order {order} from Shard {shard_itr.shard_id} at position {record['SequenceNumber']}"
                    )

                if records_response["NextShardIterator"]:
                    shard_itr.iterator = records_response["NextShardIterator"]
            except Exception as e:
                logging.error(
                    {"message": "Failed fetching records", "error": str(e)}
                )

        time.sleep(1)
```

record

**Part 1**



Producer → Amazon Kinesis Data Stream → Consumer
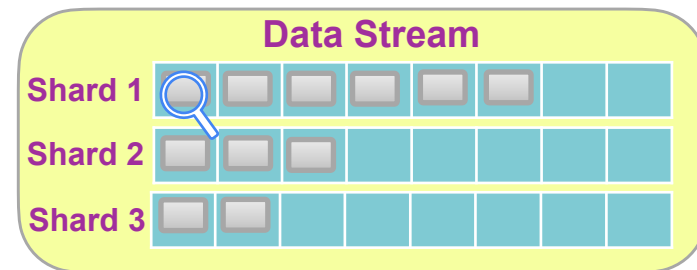
**producer_from_cli.py**

- writes a single data record into the data stream
- uses boto3 to interact with Kinesis
- can be run from the terminal:

```
python producer_from_cli.py
      --stream <name of the data stream>
      --json_string <record as json string>
```
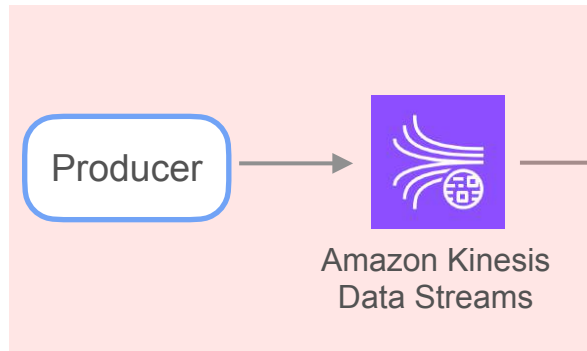
**consumer_from_cli.py**

- simple consumer application
- uses boto3 to interact with Kinesis
- can be run from the terminal:

```
python consumer_from_cli.py
      --stream <name of the data stream>
```
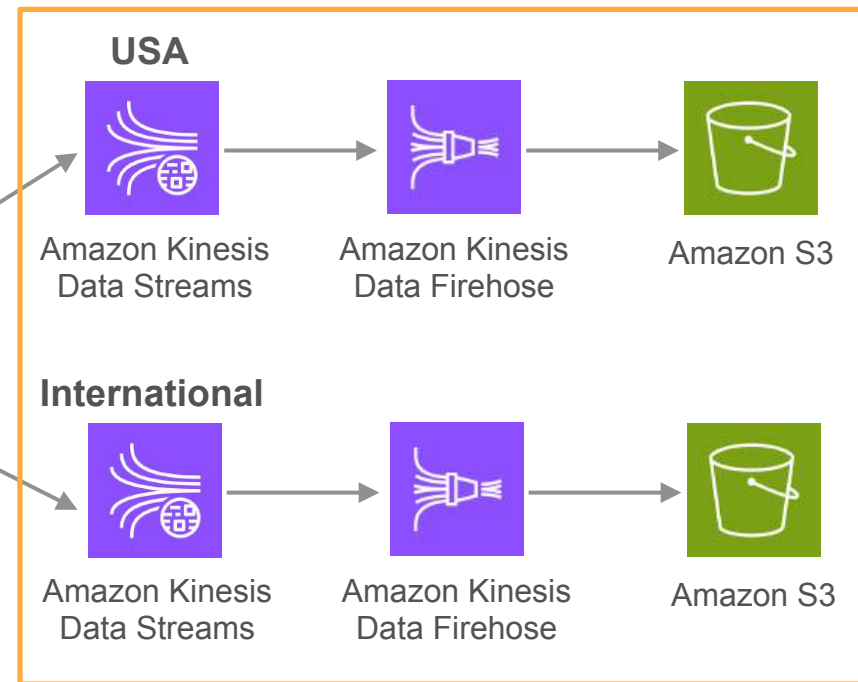
**Data Stream**

Shard 1

Shard 2

Shard 3

print information in the terminal about each record

# Part 2

**You will be provided with:**

Producer → Amazon Kinesis Data Streams → Consumer

Amazon Kinesis Data Streams → Amazon Data Firehose → Amazon S3

Amazon Kinesis Data Streams → Amazon Data Firehose → Amazon S3

# Batch and Streaming Ingestion



You determine your approach based on the stakeholder needs.

# ETL and ELT

ETL

Data Sources

**Extract**

raw data

**Transform**

Staging Area

**Load**

transformed data

Target Destination

ELT

Data Sources

**Extract**

raw data

**Load**

raw data

Target Destination

**Transform**

Query data

DeepLearning.AI

# Week 2 Labs



**Lab 1**

Ingest

API

- Connection to API
- Authentication
- Pagination

**Lab 2**

Ingest

Amazon Kinesis
Data Streams

Consumer

A recommender system

Stream 1

Stream 2

You might also like…