

Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>

Data Modeling, Transformation and Serving



DeepLearning.AI

Serving Data



DeepLearning.AI

Serving Data for Analytics and Machine Learning

Week 4 Overview

Serving Data — Analytics Use Cases

Business Intelligence



Dashboards Reports



Dashboards Reports



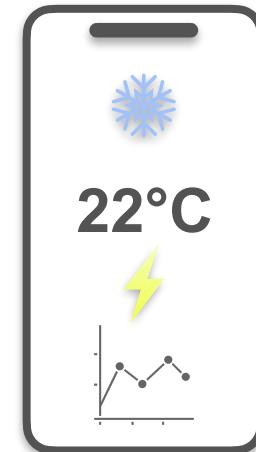
Operational Analytics

Monitor data to inform immediate action



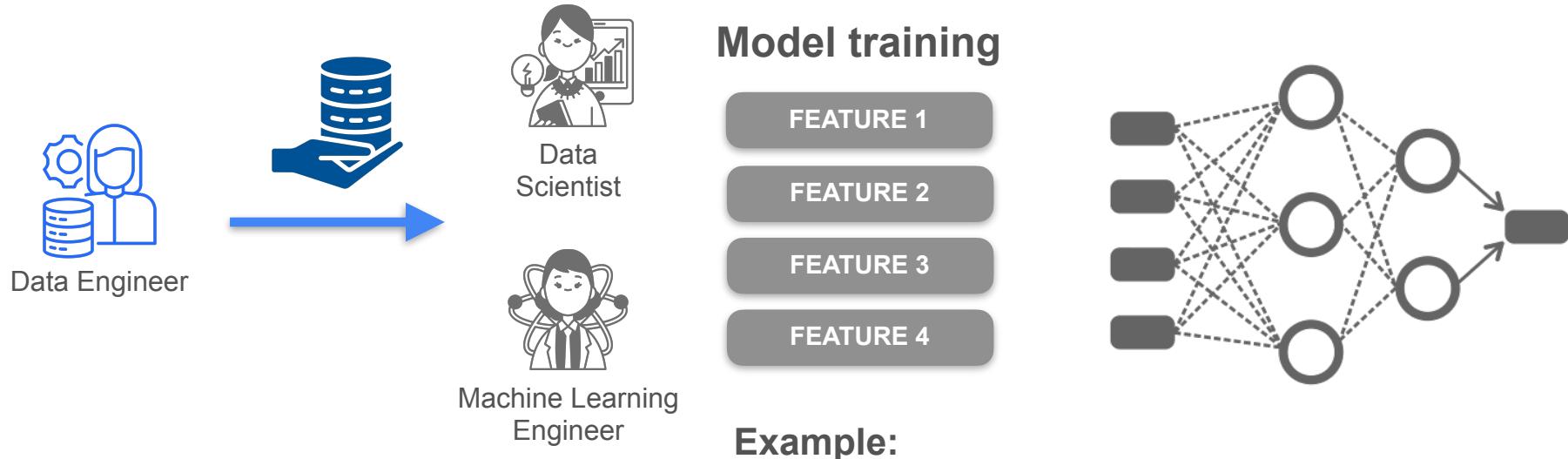
Serve data within the required latency

Embedded Analytics



Serving Data — Machine Learning Use Cases

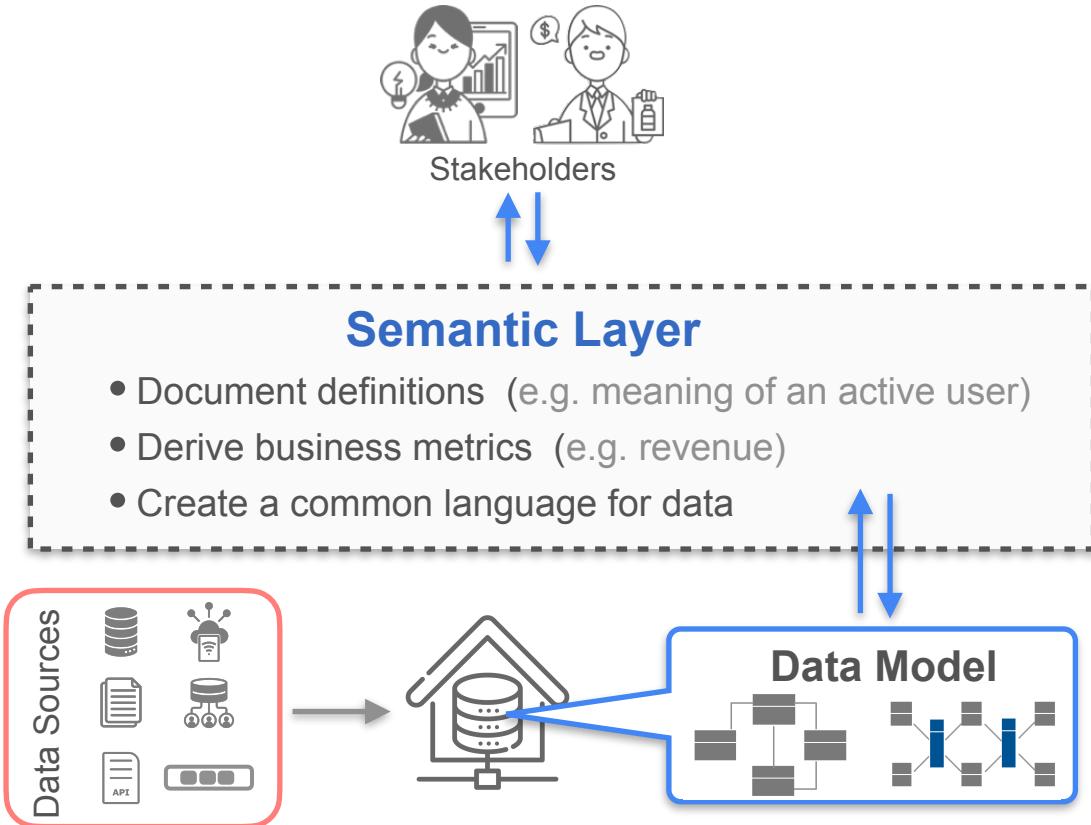
Becoming widely adopted in many companies



Example:

- Customer churn model
- Recommendation system

Serving Data



You can serve data as...

A table

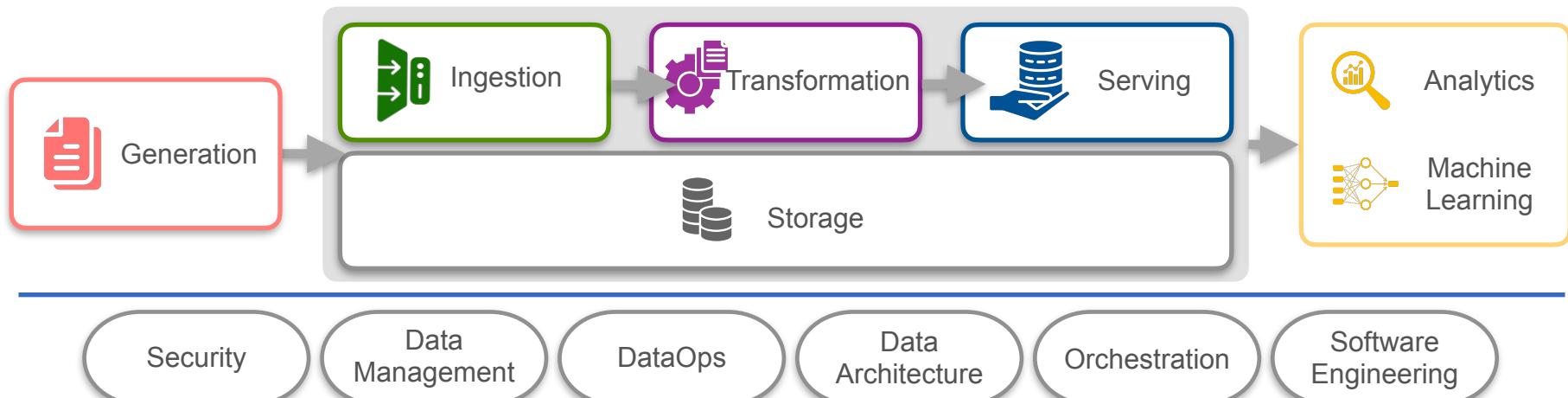
View

Materialized View

Thinking Like a Data Engineer



Data Engineering Lifecycle and Undercurrent



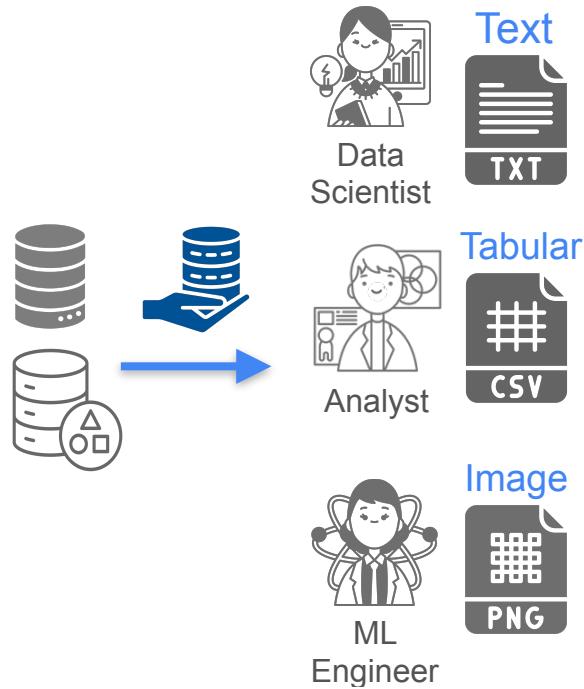


DeepLearning.AI

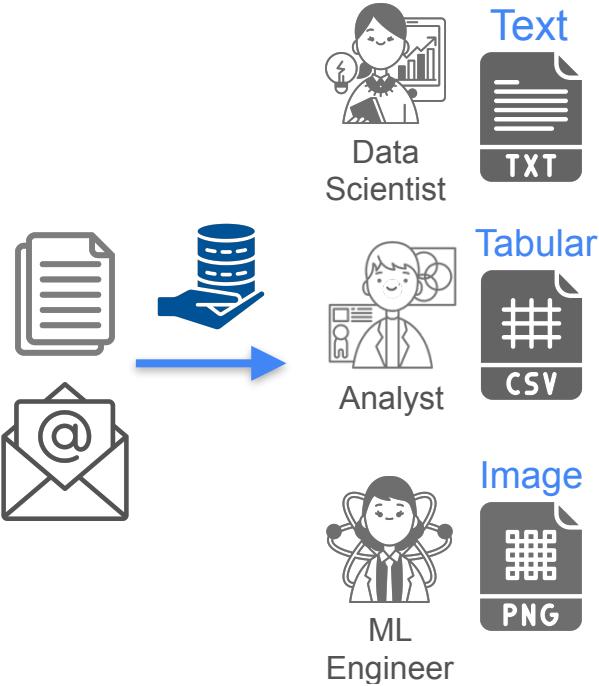
Serving Data for Analytics and Machine Learning

Serving Data

As Files

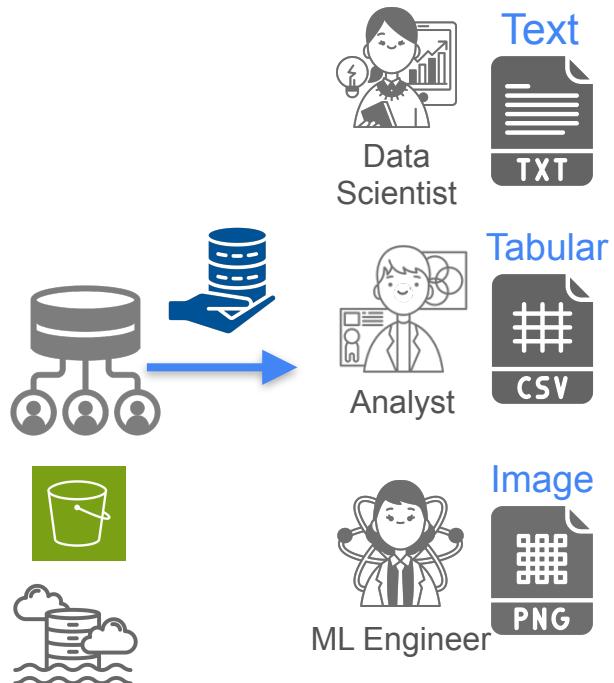


As Files

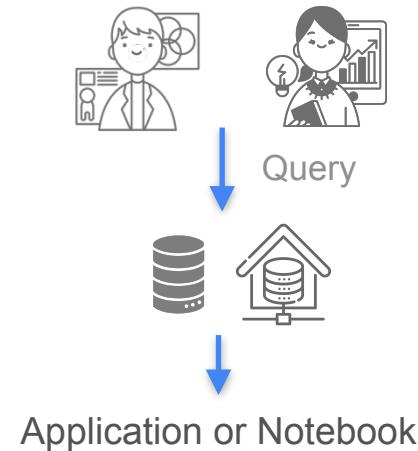


As Files

For ad hoc requests



From Databases and Data Warehouses



Benefits:

- Imposes order & structure through schema
- Gives you fine-grained permissions controls
- Offers high performance for queries

From Streaming Systems

Streaming Systems



real-time data

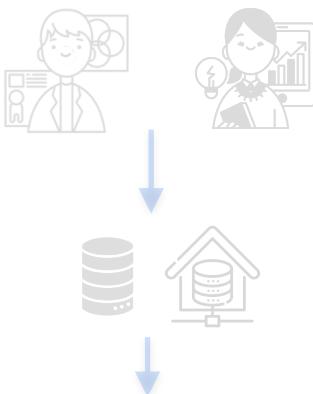


Example

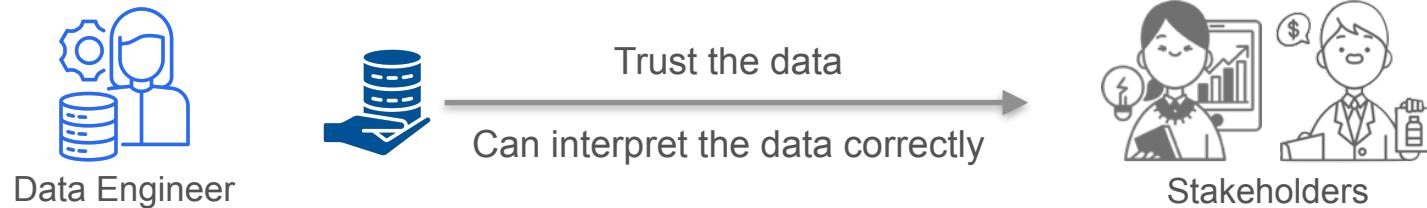


Operational
analytics database

- Enables low latency analytical queries across historical and current data
- Effectively combining the features of an OLAP database with a stream-processing system



Data Management



Ensuring data
correctness,
consistency,
trustworthiness

Data Definition

Meaning of data as it is understood throughout the organization
e.g. document and make available the definition of "Customer"

Data Logic

Consists of formulas for deriving metrics from data
e.g. gross sales or customer lifetime value

Customer
Churn Metric

- Understand the definition of "Customer"
- Write SQL query
- Define this metric so it can be reused

Semantic Layer

- Translate the underlying data elements and structures into more intuitive business terms
- Ensures consistent definitions for business terms
- Helps end users more easily navigate the data

Created in a BI tool

OR

Created using software

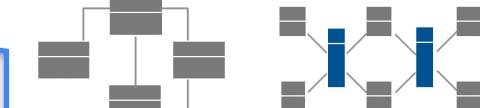


(using YAML files & SQL queries)

Data Sources

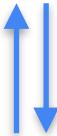


Data Model





Stakeholders



Use BI platform to visualize and analyze data



Amazon Quicksight



Looker

Use a notebook to explore data, engineer features, or train a model



Amazon SageMaker



Google Cloud Vertex AI



Azure Machine Learning

Semantic Layer

- Translate the underlying data elements and structures into more intuitive business terms
- Ensures consistent definitions for business terms
- Helps end users more easily navigate the data

Created in a BI tool

OR

Created using software

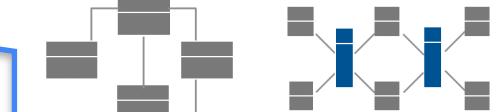


(using YAML files & SQL queries)

Data Sources



Data Model





DeepLearning.AI

Serving Data for Analytics and Machine Learning

Views and Materialized Views

View

- Stored in the database to provide easier access to common queries
- Represents a virtual table, not a physical one
- When selecting from a view:
 - The database creates a new query
 - The query optimizer optimizes and runs the query

```
CREATE VIEW customer_info AS
SELECT
    first_name,
    last_name,
    email,
    phone,
    city,
    postal_code,
    country
FROM
    customer
JOIN address ON address.id = customer.address_id
JOIN city ON city.id = address.city_id
JOIN country ON country.id = address.country_id
```

View

- Stored in the database to provide easier access to common queries
- Represents a virtual table, not a physical one
- When selecting from a view:
 - The database creates a new query
 - The query optimizer optimizes and runs the query

```
CREATE VIEW customer_info AS
```

```
SELECT
```

```
    first_name,  
    last_name,  
    email,  
    Phone,  
    city,  
    postal_code,  
    country
```

```
FROM
```

```
    customer  
    JOIN address ON address.id = customer.address_id  
    JOIN city ON city.id = address.city_id  
    JOIN country ON country.id = address.country_id
```



Effectively restricting data access to
only the data needed



Marketing Analyst

query

Wide Table:
customer, address, city, country

View

- Stored in the database to provide easier access to common queries
- Represents a virtual table, not a physical one
- When selecting from a view:
 - The database creates a new query
 - The query optimizer optimizes and runs the query
- Database object
- Stored and persists on disk

Common Table Expressions (CTE)

- Represents temporary results that you can reference
- Only exists within scope of the main query where they are referenced

```
WITH name of the CTE AS (
Query),
Subsequent query
```

View

- Stored in the database to provide easier access to common queries
- Represents a virtual table, not a physical one
- When selecting from a view:
 - The database creates a new query
 - The query optimizer optimizes and runs the query
- Database object
- Stored and persists on disk



Can be expensive
for frequently-run
complex queries

Common Table Expressions (CTE)

- Represents temporary results that you can reference
- Only exists within scope of the main query where they are referenced

```
CREATE MATERIALIZED VIEW rental_by_category AS
SELECT category.name AS category,
       sum(payment.amount) AS total_sales
FROM payment
JOIN ... rental, inventory, film, film_category, category ...
GROUP BY category.name
```

Materialized Views

- Does some or all of the view computations in advance
- Caches the query results and allows you to refresh the data
- Some latency between refreshes

query





DeepLearning.AI

Data Engineering

Summary of the Program Concepts

Thinking Like a Data Engineer



1

Identify business goals & stakeholder needs

1. Identify business goals & stakeholders you will serve
2. Explore existing systems and stakeholder needs
3. Ask stakeholders what actions they will take with the data product



2

Define system requirements

1. Translate stakeholder needs to functional requirements
2. Define non-functional requirements
3. Document and confirm requirements with stakeholders



3

Choose tools & technologies

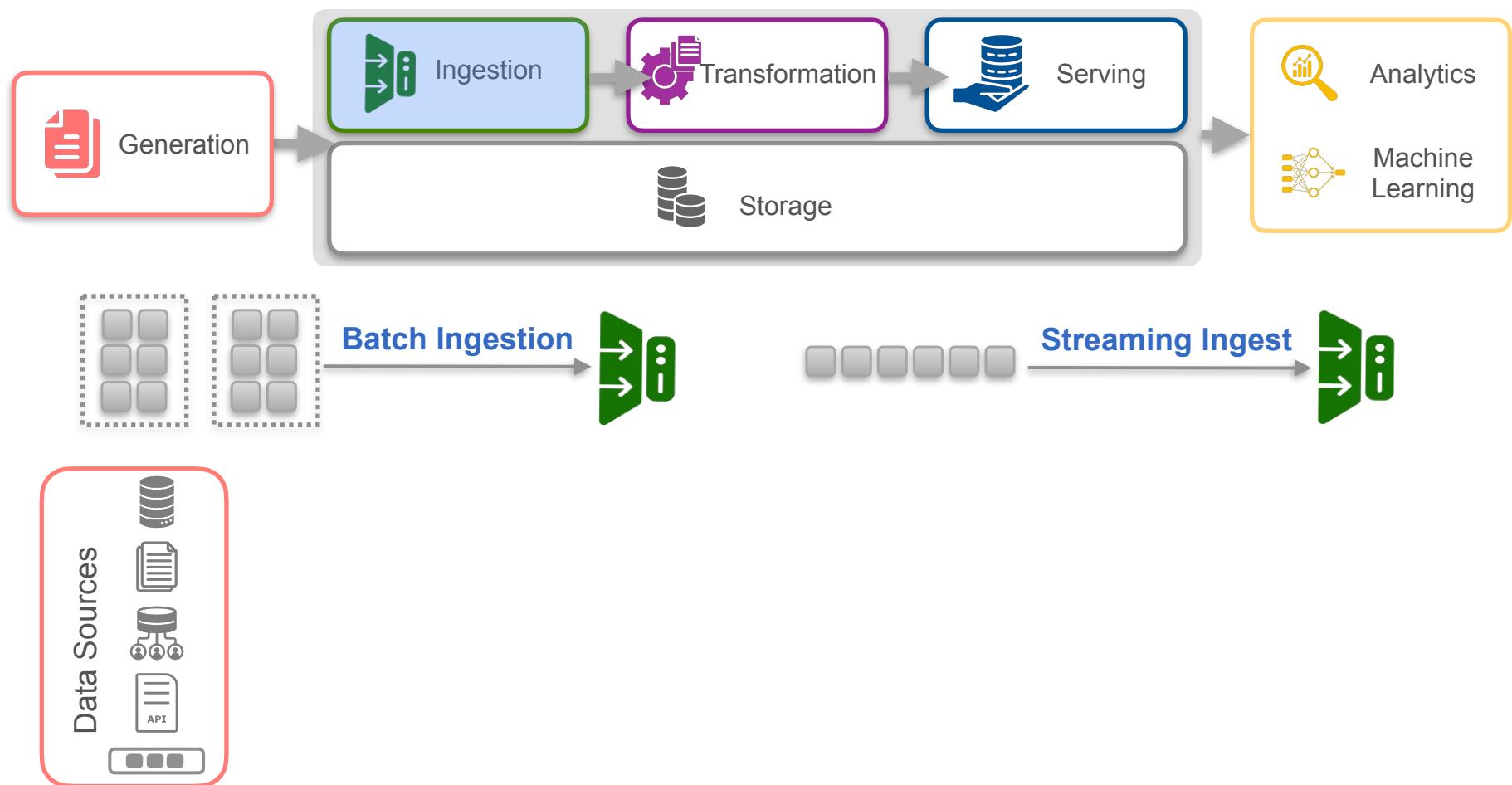
1. Identify tools & tech to meet non-functional requirements
2. Perform cost / benefit analysis and choose between comparable tools & tech
3. Prototype and test your system, align with stakeholder needs

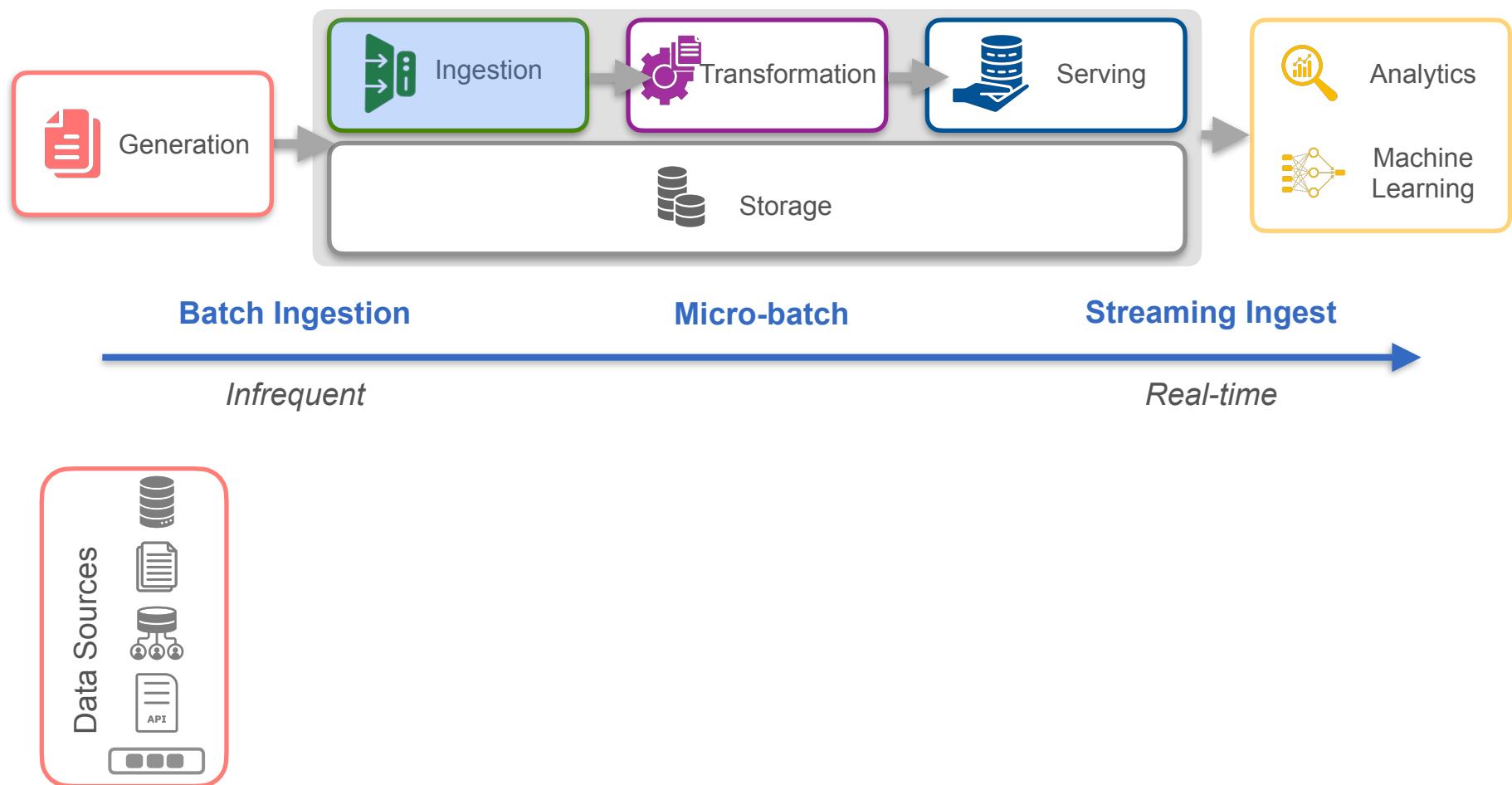


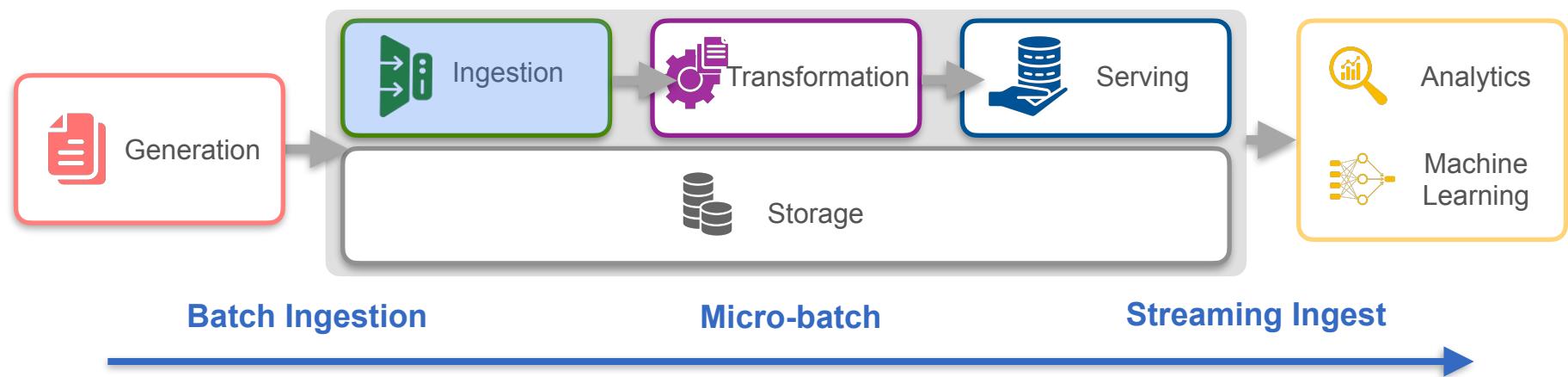
4

Build, evaluate, iterate & evolve

1. Build & deploy your production data system
2. Monitor, evaluate, and iterate on your system to improve it
3. Evolve your system based on stakeholder needs



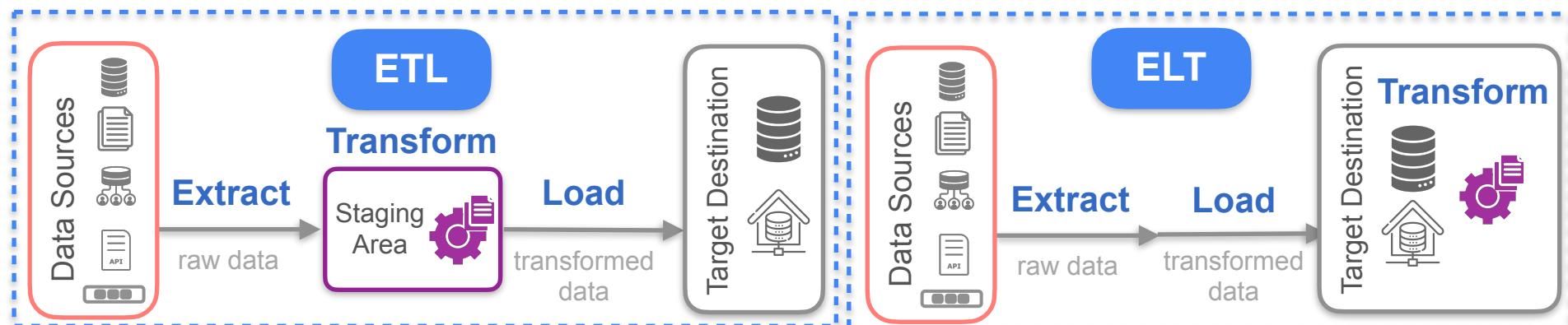




- Type of transformation

- Hardware specification

- Size of the data





Transformation

- Cleaning and combining data from multiple sources
- Convert data into a target schema based on a data model

- Apply transformations with SQL queries or Python
- Use a non-distributed or distributed processing tool



Analytics



Process large amounts of data inside a cloud data warehouse to leverage MPP



Machine Learning

Transformations depend on end use case:

- Data exploration
- Training a machine learning model
- Making predictions



Data Warehouse

- Expects structured datasets with a well-defined schema
- Query optimizer returns results based on best execution plan
- Suitable for analytical workloads because it's based on columnar storage



Data Lake

- Built on top of low-cost object storage
- Supports ML & big data processing
- To prevent creating a data swamp, you can create a data catalog to track and manage data



Data Lakehouse



+



Low-cost &
scalable
storage

Structured query
& data
management

- Serves low-latency analytics and ML

Security



IAM



VPC



Network ACL



Route table

Security group

Data Management

- Model your data



AWS Glue
Data Catalog

DataOps

- Infrastructure as code



Data Architecture

Orchestration



Software Engineering

- Test & monitor data quality



AWS Glue
Data Quality



DeepLearning.AI

Lab Walkthrough

Capstone Project (Part 1)

Capstone Lab

Solution to a business use case



Part 1 Create and configure the resources for the pipeline

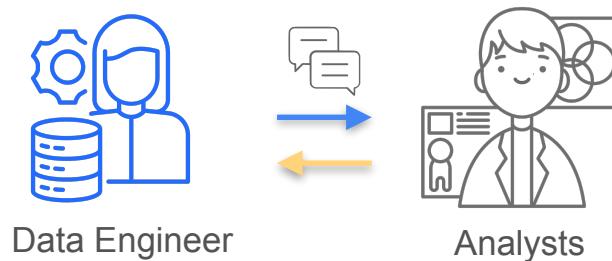
Part 2 Integrate data quality checks and orchestration,
& create data visualizations

Capstone Lab



- Offers subscription-based music streaming services
- Added a new feature allowing its clients to purchase and download music

GET	/users	Get Users
GET	/sessions	Get Sessions



Building a data pipeline to serve the purchase data

Analyzing the purchase data

- with respect to the song and artist features
- how the purchases are evolving with time
- how the purchases are affected by the users' personal information

Data Sources



GET

/users Get Users

GET

/sessions Get Sessions

Response of “Get Users” Request

```
{  
  "user_id": "a3141825-3a8c-4968-a3af-5362011ef7d5",  
  "user_name": "Elizabeth",  
  "user_lastname": "Carey",  
  "user_location": [  
    "46.32313",  
    "-0.45877",  
    "Niort",  
    "FR",  
    "Europe/Paris"  
],  
  "user_since": "2020-12-22T14:15:35.936090"  
}
```

Response of “Get Sessions” Request

```
{  
  "user_id": "6b287203-7cab-4f1a-b1a4-2b5076294682",  
  "session_id": "04a5e8ac-1acd-48dc-88b9-651c4ddf489c",  
  "session_items": [  
    {  
      "song_id": "TRXKAGX128F9342DD7",  
      "song_name": "3 Cards",  
      "artist_id": "AR475MP1187B9A5449",  
      "artist_name": "The Balancing Act",  
      "price": 1.03,  
      "currency": "USD",  
      "liked": true,  
      "liked_since": "2023-01-27T08:29:54.970697"  
    },  
    {  
      "song_id": "TRUKGBT128F4292C9B",  
      "song_name": "Parisian Walls (gband Version_ Barcelona)",  
      "artist_id": "ARP9HJX1187FB4E5DA",  
      "artist_name": "Apostle Of Hustle",  
      "price": 1.31,  
      "currency": "USD",  
      "liked": true,  
      "liked_since": "2023-06-14T00:27:55.876873"  
    }  
  ]  
}
```

Data Sources

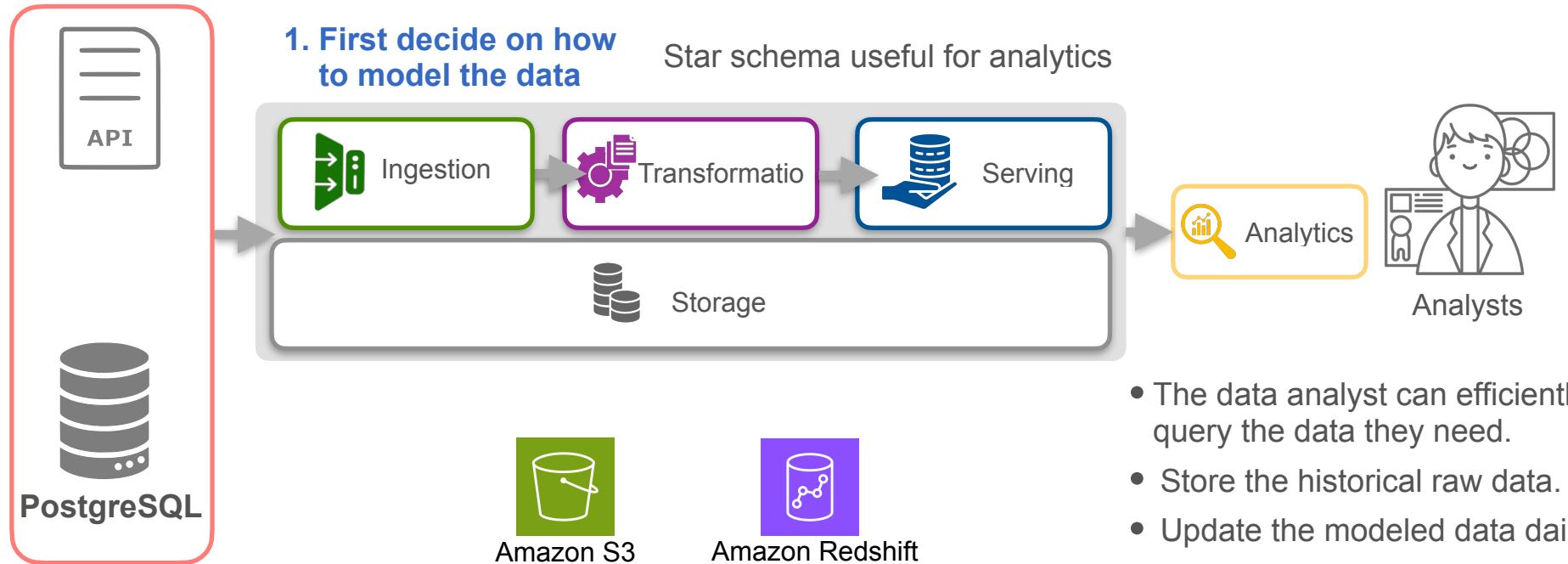


PostgreSQL

songs	
track_id	
title	
song_id	
release	
artist_id	
artist_mbid	
artist_name	
duration	
artist_familiarity	
artist_hottness	
year	
track_7digitalid	

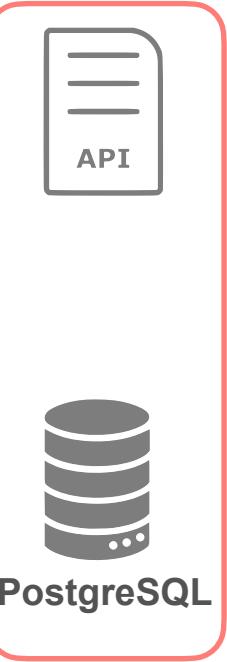
schema:
deftunes

Requirements



Data Model

1. First decide on how to model the data



Business process: Song purchases

Atomic Grain: Individual song item purchased within a session

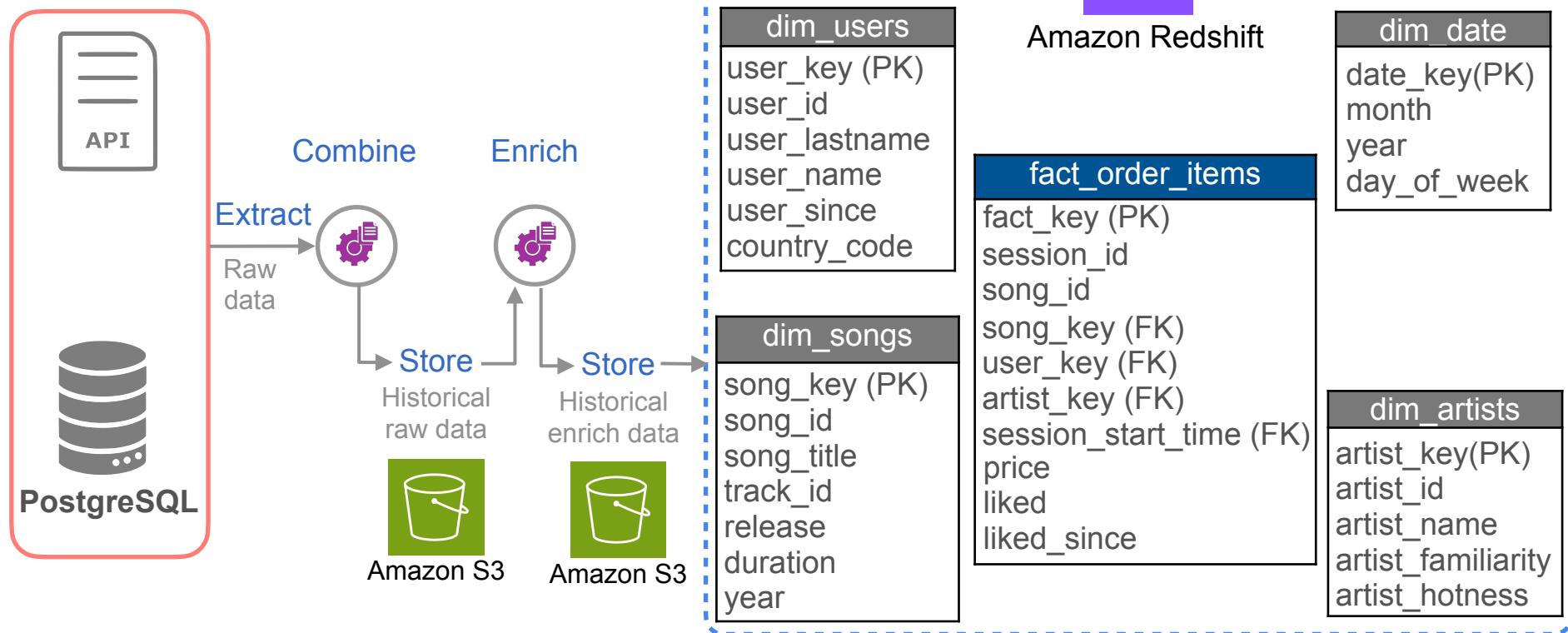
dim_users	dim_date
user_key (PK)	date_key(PK)
user_id	month
user_lastname	year
user_name	day_of_week
user_since	
country_code	

fact_order_items
fact_key (PK)
session_id
song_id
song_key (FK)
user_key (FK)
artist_key (FK)
session_start_time (FK)
price
liked
liked_since

dim_songs	dim_artists
song_key (PK)	artist_key(PK)
song_id	artist_id
song_title	artist_name
track_id	artist_familiarity
release	artist_hotness
duration	
year	

**Assuming all prices are in USD.
**If the prices were in different currencies, convert all prices to the same currency.

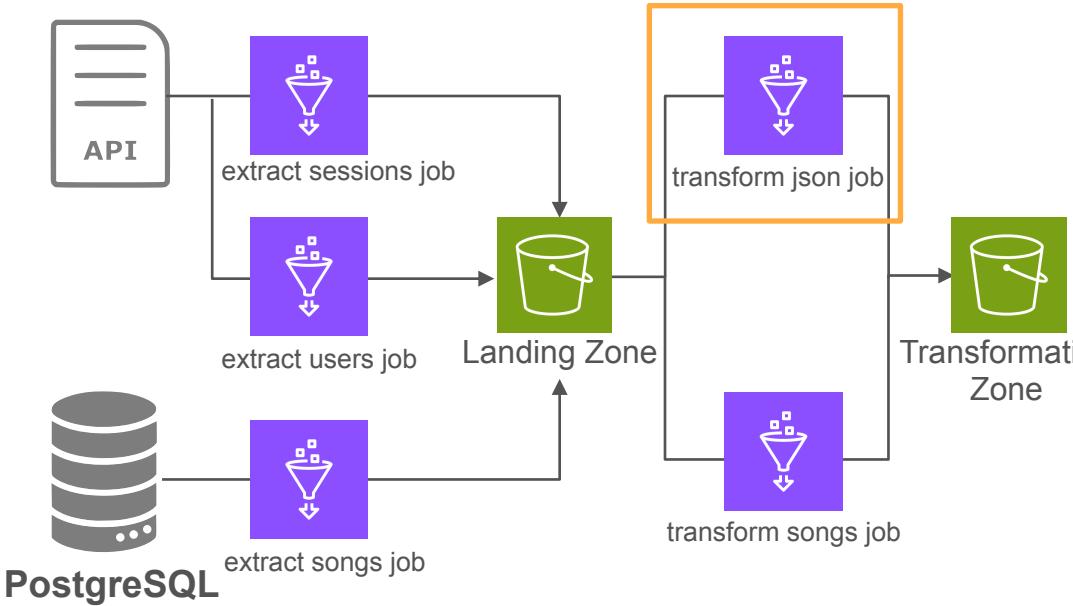
Data Model



Solution Architecture



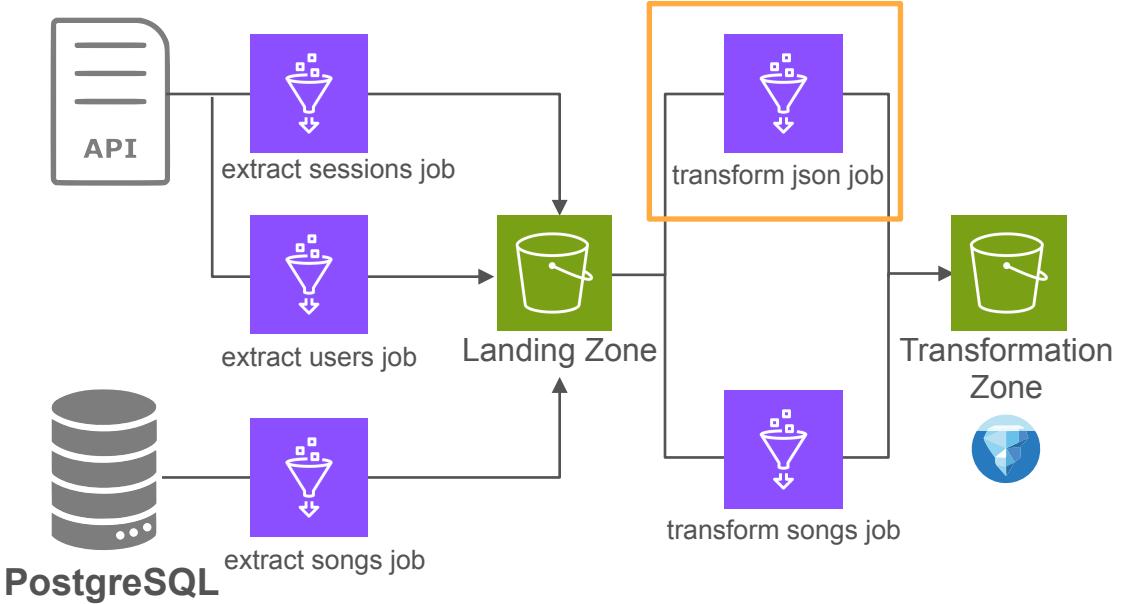
Solution Architecture



```
{  
  "user_id": "a3141825-3a8c-4968-a3af-5362011ef7d5",  
  "user_name": "Elizabeth",  
  "user_lastname": "Carey",  
  "user_location": [  
    "46.32313",  
    "-0.45877",  
    "Niort",  
    "FR",  
    "Europe/Paris"  
  ],  
  "user_since": "2020-12-22T14:15:35.936090"  
}
```

- Transformation Zone
1. Extract the users data into a DataFrame
 2. Unnest “user_location” :
 - latitude
 - longitude
 - place name
 - country code
 - timezone
 3. Add a column indicating the ingestion time
 4. Add a column indicating the processing time

Solution Architecture

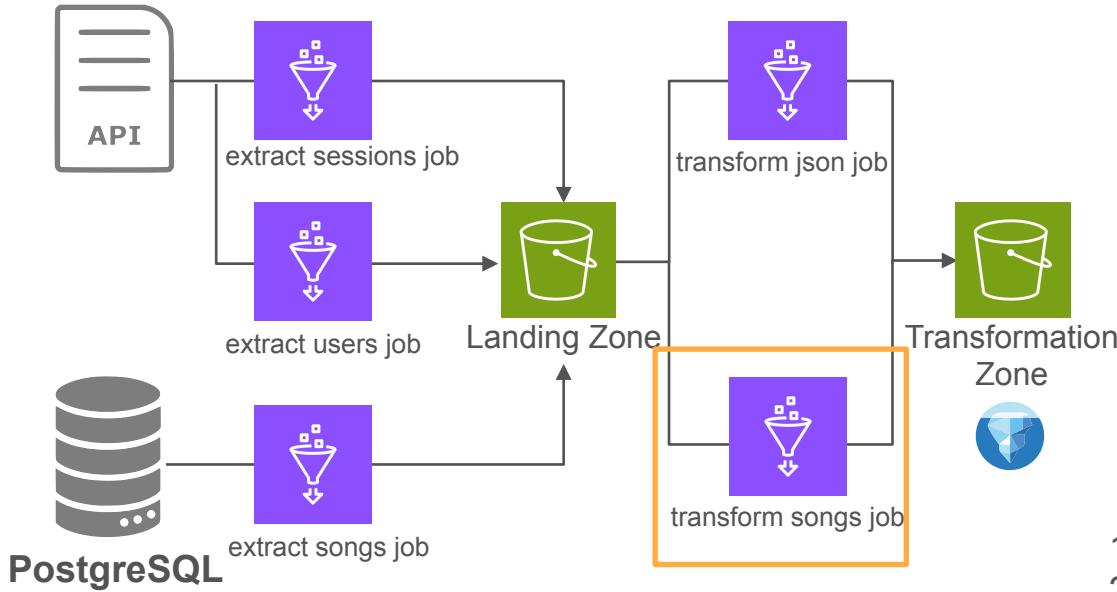


 s3://{{data_lake_bucket_name}}/transform_zone/users
s3://{{data_lake_bucket_name}}/transform_zone/sessions

```
{  
  "user_id": "6b287203-7cab-4f1a-b1a4-2b5076294682",  
  "session_id": "04a5e8ac-1acd-48dc-88b9-651c4ddf489c",  
  "session_items": [  
    {  
      "song_id": "TRXKAGX128F9342DD7",  
      "song_name": "3 Cards",  
      "artist_id": "AR475MP1187B9A5449",  
      "artist_name": "The Balancing Act",  
      "price": 1.03,  
      "currency": "USD",  
      "liked": true,  
      "liked_since": "2023-01-27T08:29:54.970697"  
    },  
    {  
      "song_id": "TRUKGBT128F4292C9B",  
      "song_name": "Parisian Walls (gband Version_1)",  
      "artist_id": "ARP9HJX1187FB4E5DA",  
      "artist_name": "Apostle Of Hustle",  
      "price": 1.31,  
      "currency": "USD",  
      "liked": true,  
      "liked_since": "2023-06-14T00:27:55.876873"  
    },  
  ]  
},  
}
```

1. Extract the sessions data into a DataFrame
2. Unnest “session items”
3. Add a column indicating the ingestion time
4. Add a column indicating the processing time

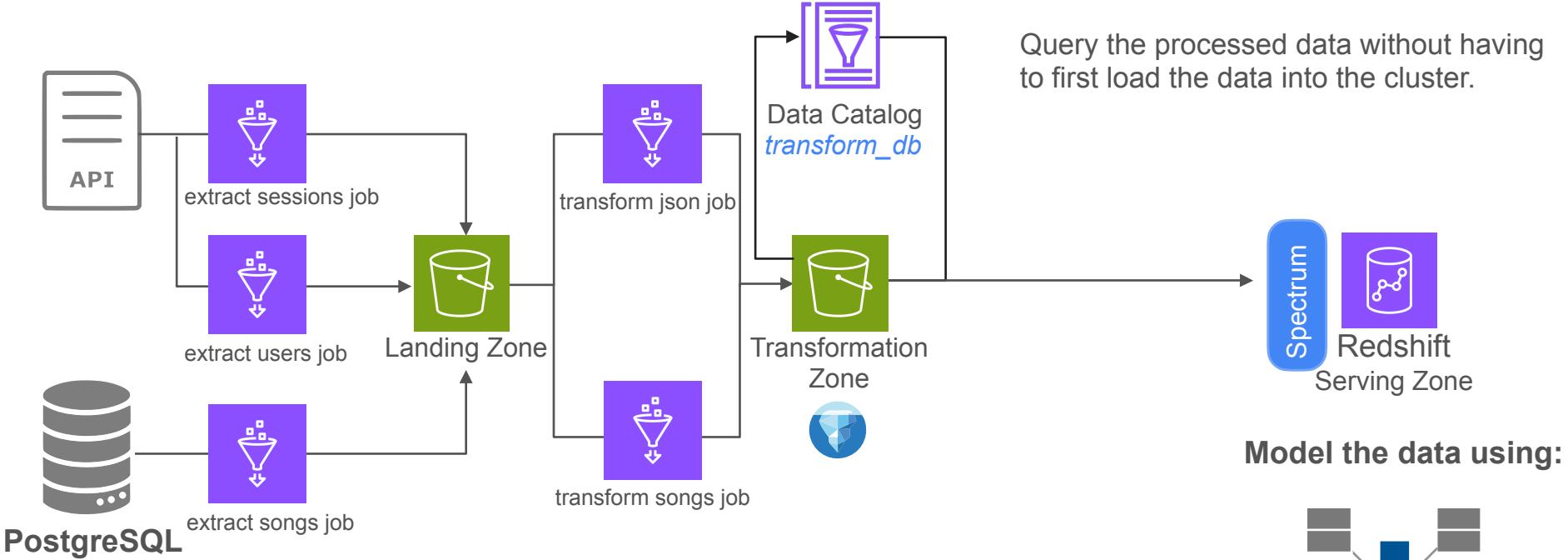
Solution Architecture



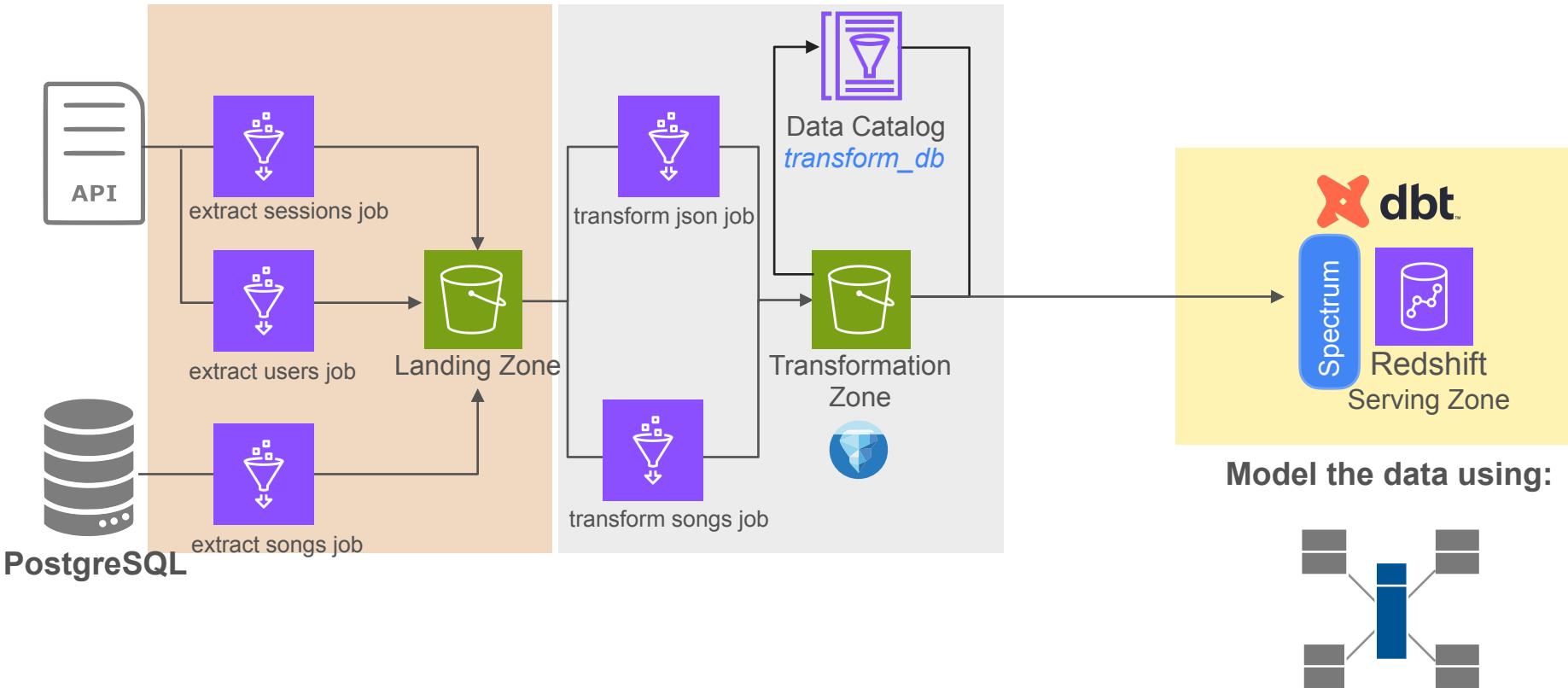
songs
track_id
title
song_id
release
artist_id
artist_mbid
artist_name
duration
artist_familiarity
artist_hottness
year
track_7digitalid

1. Extract the csv data into a DataFrame
2. Perform schema enforcement
3. Add a column indicating the ingestion time
4. Add a column indicating the name of the database source

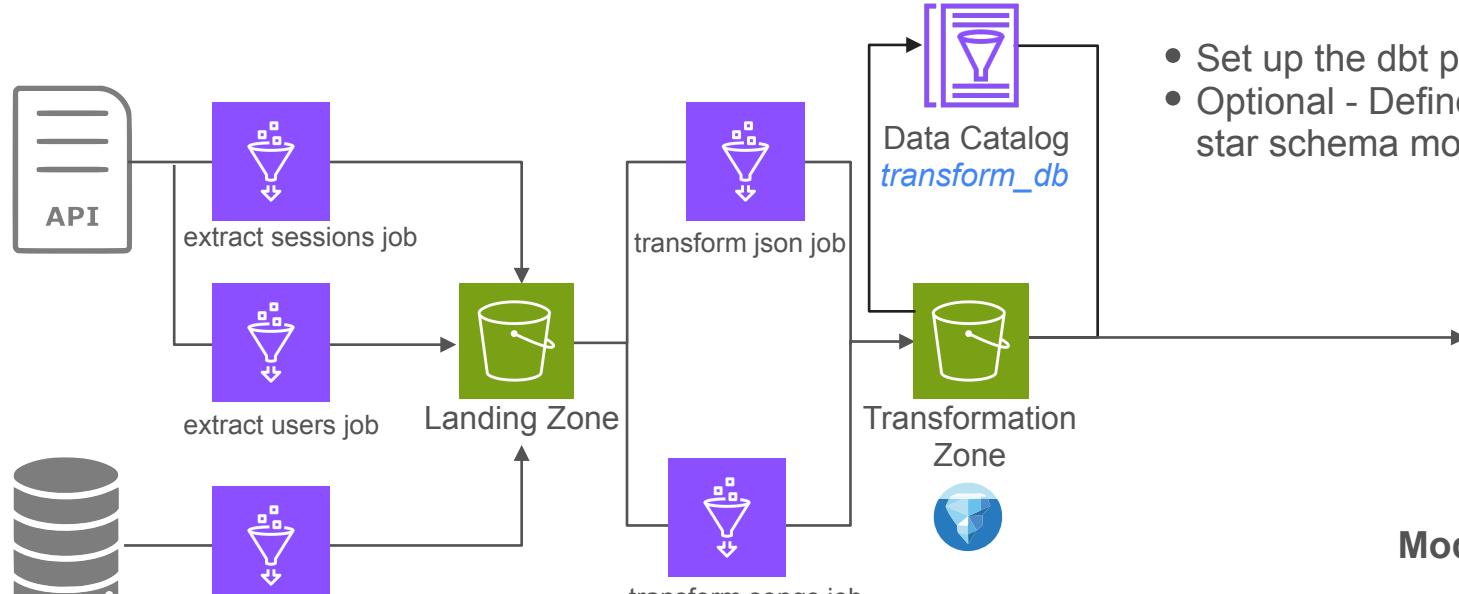
Solution Architecture



Solution Architecture



Capstone Lab - Part 1



- Set up the dbt project folder
- Optional - Define the tables of the star schema model

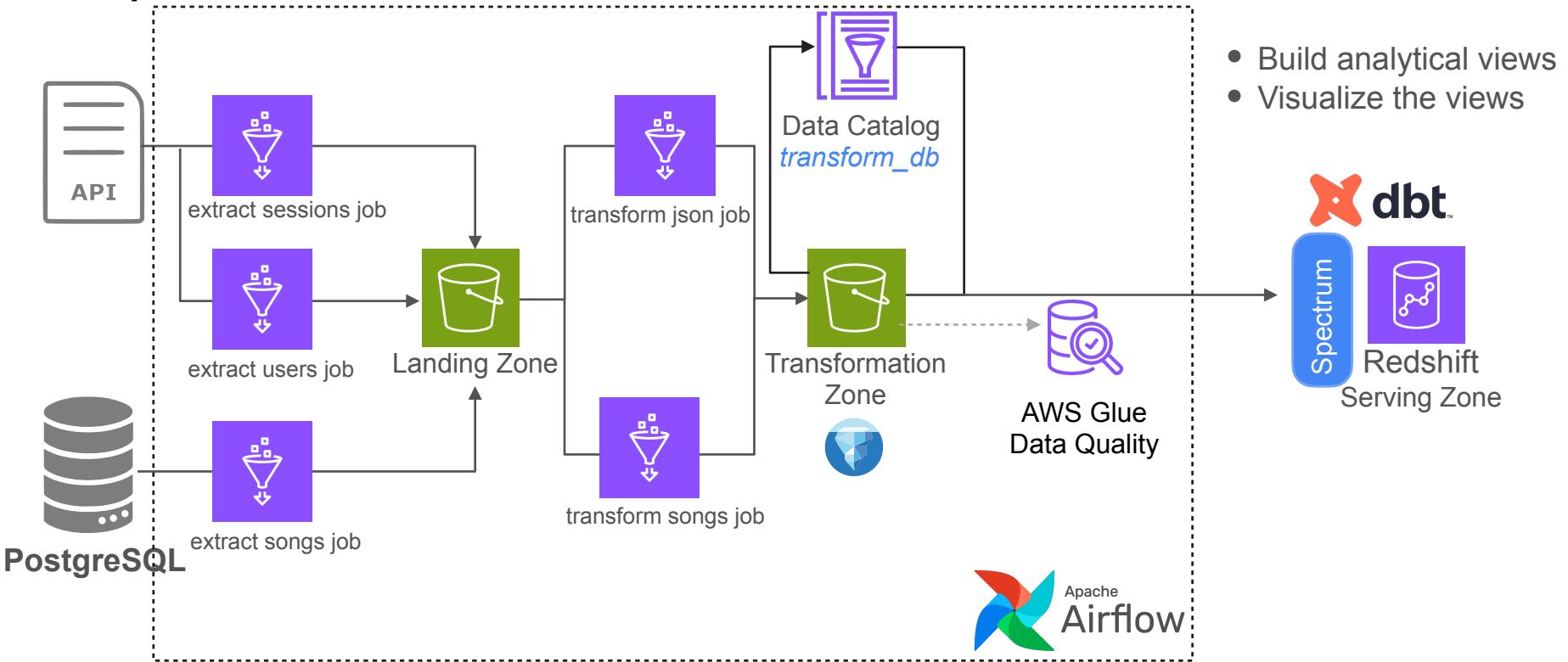


Model the data using:



Use Terraform to create the resources of the pipeline.

Capstone Lab - Part 2





DeepLearning.AI

Lab Walkthrough

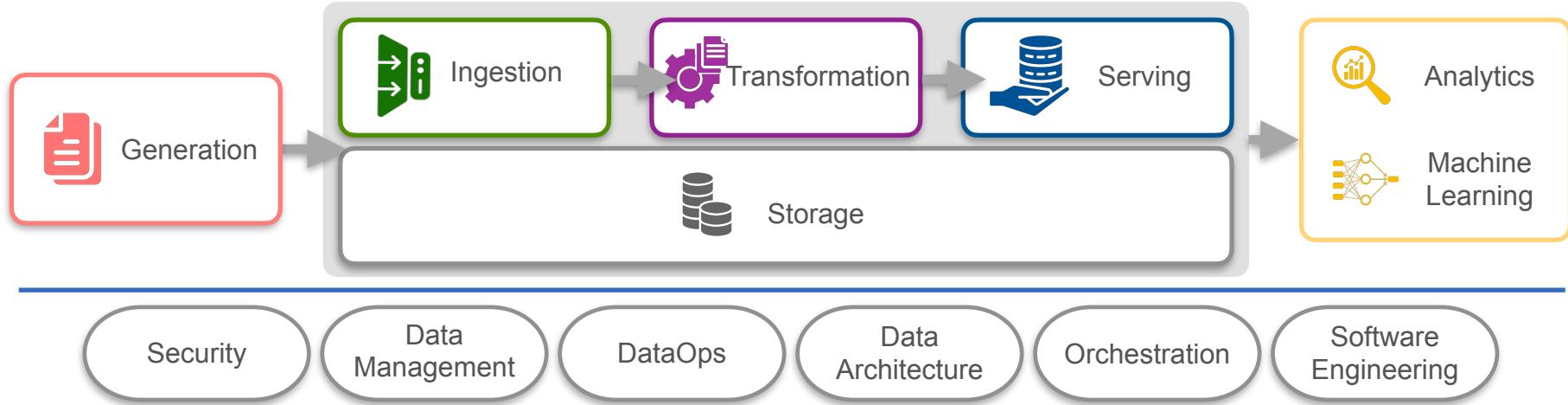
Capstone Project (Part 2)



DeepLearning.AI

Data Engineering

Conclusion & Thank You



The Role of the Data Engineer

Data Team



New role in the realm of
data and algorithms



Data Engineer



Software Engineer



Data Scientist



ML Engineer

The Role of the Data Engineer

Data Team



New role in the realm of
data and algorithms



Data Engineer



Software Engineer



Data Scientist ML Engineer

Operational discipline

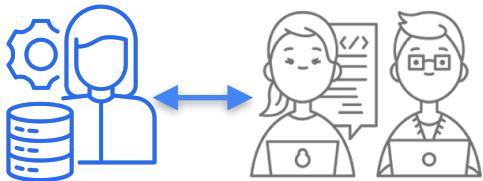
- Create or use systems that automate the development and operationalization of well-understood ML models
- Monitor data pipelines and quality

The Role of the Data Engineer

Data Team

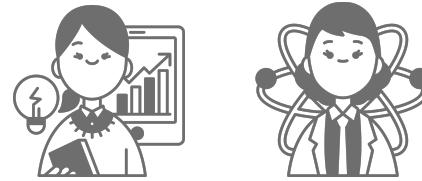


New role in the realm of data and algorithms



Data Engineer Software Engineer

- Deep understanding in software and DE
- Expertise in streaming, data pipelines, data modeling, and data quality
- Boundaries between application backend systems and DE tools will be lowered



Data Scientist ML Engineer

Operational discipline

- Create or use systems that automate the development and operationalization of well-understood ML models
- Monitor data pipelines and quality