# RoadMap of Data Engineer

## Role of Data Engineer:

| ① | ② | ③ | ④ |
|---|---|---|---|
| Data Source | Data Processing | Data Storage | Data View |

**Data Source:**
- file
- Cloud storage
- Database
- IOT
- Applications

**Data Processing:**
- ETL Pipeline
- ELT pipeline

entract, load, Transform

**Data Storage**

**Data View:**
- Stakeholder/ Business Users
- Analytical Databases
- Data Professionals
- Reporting
- Visualization

## Pre-requisities:

* SQ2 → bread & butter of data fields
* Python → It is crucial for advanced data engineering tasks

### Cloud Platforms

✓ * GCP      $400 → 3 Mnths (data proc)
  * Azure    (hdinsight)
  * AWS      (EMR)

## Technologies Covered:

(1, **Hadoop :** (Foundational)
  HDFS (Storage) → MapReduce (Processing) → YARN (resource management)

  eng. Amazon

  HDFS :   (transaction log) → distributed Servers

  → MapReduce : count product purchases         $\frac{100 → Items}{}$
                                                ↑ mem Count ()
  YARN :   Efficiently processing across distributed system  ①,②,③④

② **Apache Spark:**

Components:
- Row-level API's, High-level API's  } ETL
- Spark SQL
- Caching
- join
- Optimization

<u>examples</u>: Recommendation System → Netflix

1, History of user views ⟶ Spark

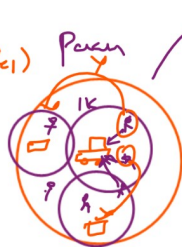2, Spark SQL ⟶ user preferences

3, ML Algorithm ⟶ (Collaborative filtering) → Spark MLlib → recommend movies

③ **Hive:**

simplifies

MapReduce ⟶ HiveQL

④ **Katka:** (Time-series) Producer ⟶ Centralized

synchronize

user ⟶ Cab



⑤ Apache Airflow

⑥ Databricks

⑦ Azure Cloud:

Tools: Azure data factory

ADLS Gen2

Databricks on Azure

MongoDB

Azure Synapse

## Promise

End-to-end industry level project Integrating multiple tools

## Goal:

: * equip skills to clear interviews

* excel in data engineering roles