

IDS 706: Data Engineering Systems

Dr. Zhongyuan (Annie) Yu

Duke

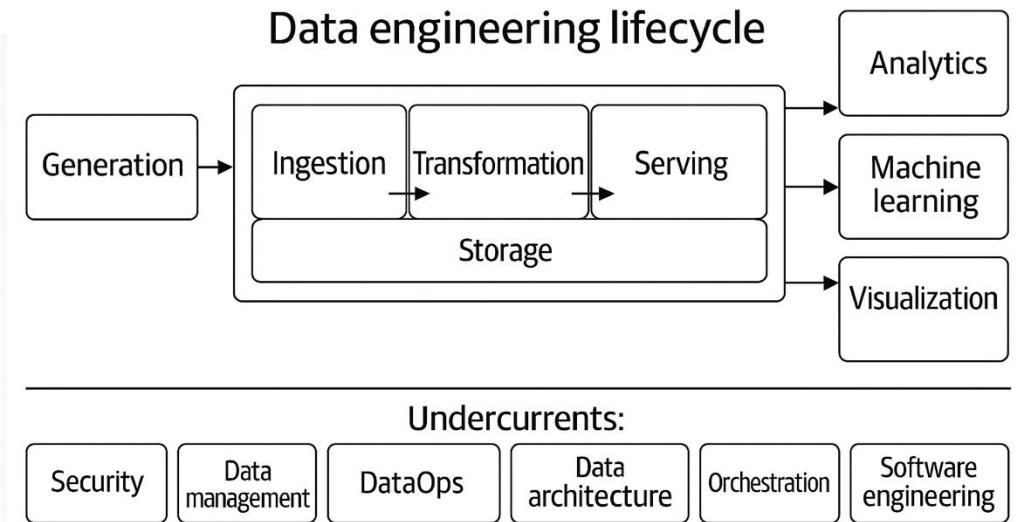
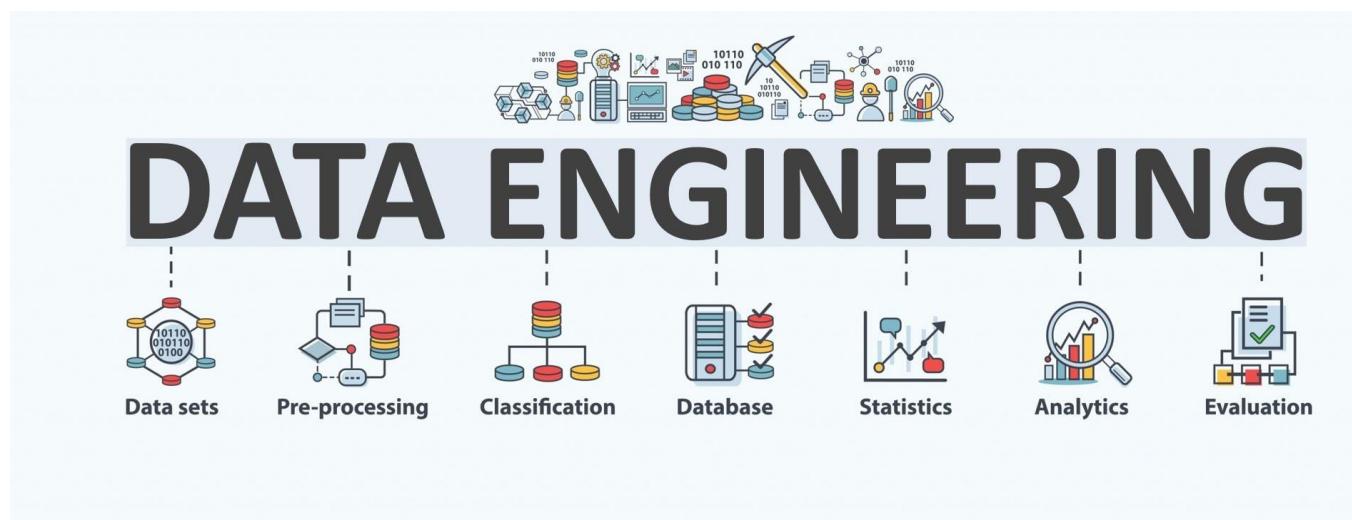
Course Basics

Instructor: Dr. Zhongyuan Yu

Contact Info: Zhongyuan.yu@duke.edu

Office Hours: Friday 11:45 am -12:45 pm Gross Hall 227 or Zoom

What is Data Engineering?



- Designing and managing the infrastructure that powers data generation, collection, and distribution.
- Building scalable pipelines to transform and move data across systems.
- Making data accessible, reliable, and analysis-ready for analysts and data scientists.

Core Engineering Mindset

- **Problem-solving focus:** Engineers approach challenges as opportunities to create effective solutions and break complex challenges into smaller, manageable pieces.
- **Systems thinking:** They understand how individual parts connect and influence the whole system, and understand dependencies and trade-offs when making decisions.
- **Iterative and continuous improvement:** Solutions evolve through prototyping, testing, and refining over time.
- **Creative but constraint-aware:** Innovation happens within practical limits like cost, time, and regulations.
- **Data-driven decision making:** Decisions are guided by metrics, evidence, and analysis, not assumptions.

What this course is about?

- Understand the Role of a Data Engineer
- Apply software engineering practices like version control, testing, and CI/CD
- Work with data using Python, SQL, and cloud platforms
- Build and manage scalable data pipelines using modern tools
- Understand big data processing and real-time streaming

Prerequisite(s): None. But programming experience is helpful.



Cloud-based Containerized Python Template

The screenshot shows a GitHub repository named "IDS-706-week-1-template" owned by "zhongyuan-duke". The repository is public and contains 1 branch and 0 tags. The code tab is selected. The repository has 5 commits from "zhongyuan-duke" pushing changes to the main branch. The commits are:

Commit	Message	Time
modify readme	add devcontainer	44 minutes ago
.devcontainer	add devcontainer	44 minutes ago
.github	Initial commit	yesterday
.gitignore	Initial commit with Python template setup	2 hours ago
Makefile	modify readme	1 minute ago
README.md	Initial commit with Python template setup	2 hours ago
hello.py	Initial commit with Python template setup	2 hours ago
requirements.txt	Initial commit with Python template setup	2 hours ago
test_hello.py	Initial commit with Python template setup	2 hours ago

The repository also contains a README file.

On the left, there is a sidebar for the repository showing workflow runs and actions. The workflow runs section lists "modify readme", "devcontainers in ./ - Update #1079070076", "add devcontainer", and "Create main.yml". The actions section lists "Help us improve GitHub Actions" and "4 workflow runs".

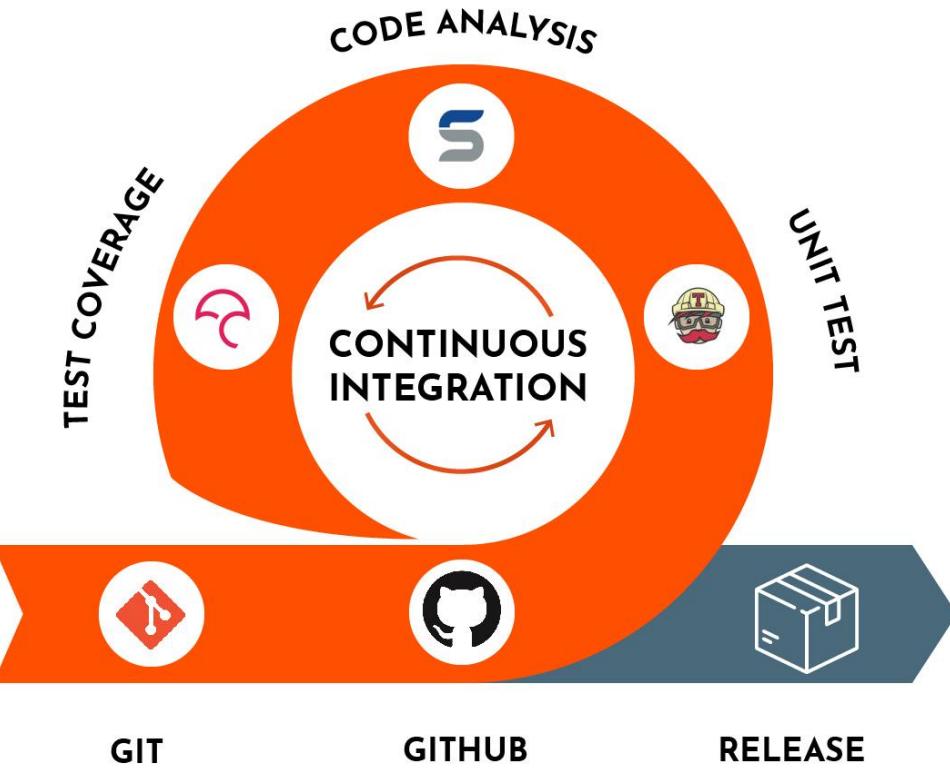
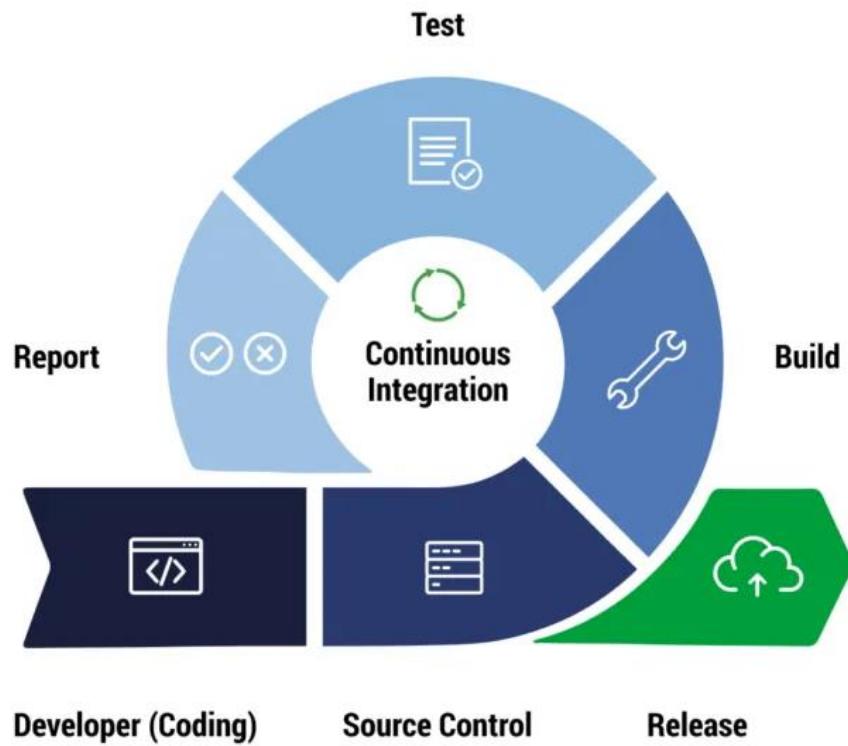
At the bottom, there is a call-to-action: "Create your first project for Data Engineering Systems (IDS 706)".

SOFTWARE TESTING



AI Generated

Continuous Integration



<https://wwwtierpoint.com/blog/benefits-of-continuous-integration/>
<https://devonblog.com/continuous-delivery/continuous-integration-best-practices/>

REFACTORING OPTIMIZED



AI Generated

SQL Cheat Sheet

Created by @JessRamosData | Big Data Energy

The Basics

```
select group by = <= not
from having != >= in
where order by < True between
      limit > False
```

String Manipulation

```
concat() upper() wildcards (%)
replace() lower() like / ilike
reverse() len()
trim() str()
```

Joins

inner join	anti join
left join	joining on multiple keys
outer join	joining on a condition
self join	

Window Functions

```
sum() row_number()
count() rank() over(
avg() denserank() partition by ...
max() lead() order by ...)
min() lag()
```

Follow Jess Ramos for Data, SQL, & Tech Career tips

Aggregate Functions

```
min() count()
max() median()
avg() mode()
stddev()
```

Date Manipulation

```
day() date_a
month() datediff
year() date_t
getdate() date_f
```

Cleaning & Transforming

```
cast() iif()
coalesce() listagg()
case when except
ifnull()
```

Advanced

CTEs
Subqueries (correlated vs non-correlated)
UDFs
Data Modeling

Try out the interactive SQL editor!



Big Data Energy

SQL Basics Cheat Sheet

SQL

SQL or Structured Query Language, is a language to talk to databases. It allows you to select specific data and to build complex reports. Today, SQL is a universal language of data. It's used in practically all technologies that process data.

SAMPLE DATA

COUNTRY	id	name	population	area
1	France	666000000	640680	
2	Germany	807000000	357900	
...

CITY	id	name	country_id	population	rating
1	Paris	1	22430000	4	
2	Berlin	2	34660000	3	
...

FILTERING THE OUTPUT COMPARISON OPERATORS

Fetch names of cities that have a rating above 3:

```
SELECT name
FROM city
WHERE rating > 3;
```

Fetch names of cities that are neither Berlin nor Madrid:

```
SELECT name
FROM city
WHERE name != 'Berlin'
AND name != 'Madrid';
```

QUERYING MULTIPLE TABLES INNER JOIN

JOIN (or explicitly INNER JOIN) returns rows that have matching values in both tables.

```
SELECT city.name, country.name
FROM city
INNER JOIN country
ON city.country_id = country.id;
```

CITY	id	name	country_id	sd	name
1	Paris	1	1	1	France
2	Berlin	2	2	2	Germany
3	Warsaw	3	4	NULL	NULL
4	London	4	5	5	United Kingdom

FULL JOIN

FULL JOIN (or explicitly FULL OUTER JOIN) returns all rows from both tables - if there's no matching row in the second table, NULLs are returned.

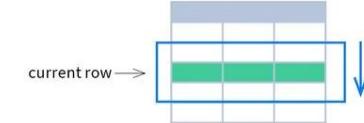
```
SELECT city.name, country.name
FROM city
FULL OUTER JOIN country
ON city.country_id = country.id;
```

CITY	id	name	country_id	sd	name
1	Paris	1	1	1	France
2	Berlin	2	2	2	Germany
3	Warsaw	3	4	NULL	NULL
4	London	4	5	5	United Kingdom
5	Rome	5	6	6	Italy

SQL Window Functions Cheat Sheet

WINDOW FUNCTIONS

compute their result based on a sliding window frame, a set of rows that are somehow related to the current row.



SYNTAX

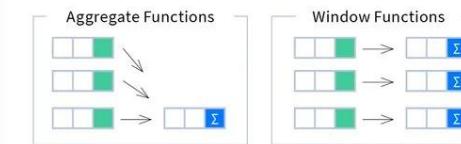
```
SELECT city, month,
       sum(sold) OVER (
        PARTITION BY city
        ORDER BY month
        RANGE UNBOUNDED PRECEDING) total
FROM sales;
```

Named Window Definition

```
SELECT country, city,
       rank() OVER country_sold_avg
FROM sales
WHERE month BETWEEN 1 AND 6
GROUP BY country, city
HAVING sum(sold) > 10000
WINDOW country_sold_avg AS (
    PARTITION BY country
    ORDER BY avg(sold) DESC)
```

AGGREGATE FUNCTIONS VS. WINDOW FUNCTIONS

unlike aggregate functions, window functions do not collapse rows.



```
SELECT <column_1>, <column_2>,
       <>window_function>() OVER (
        PARTITION BY <...>
        ORDER BY <...>
        <>window_frame>) <window_column_alias>
FROM <table_name>;
```

```
SELECT <column_1>, <column_2>,
       <>window_function>() OVER <window_name>
FROM <table_name>
WHERE <...>
GROUP BY <...>
HAVING <...>
WINDOW <window_name> AS (
    PARTITION BY <...>
    ORDER BY <...>)
```

PARTITION BY

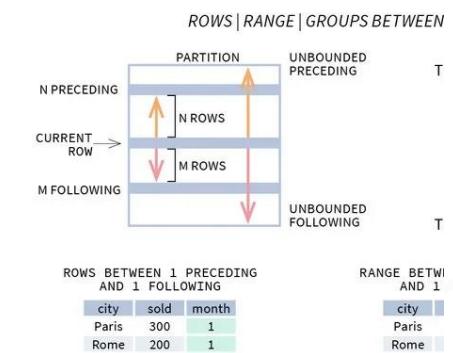
divides rows into multiple groups, called partitions, to which the window function is applied.

month	city	sold	sum
1	Rome	200	800
2	Paris	500	800
1	London	100	900
1	Paris	300	900
2	Rome	300	900
2	London	400	500
3	Rome	400	500

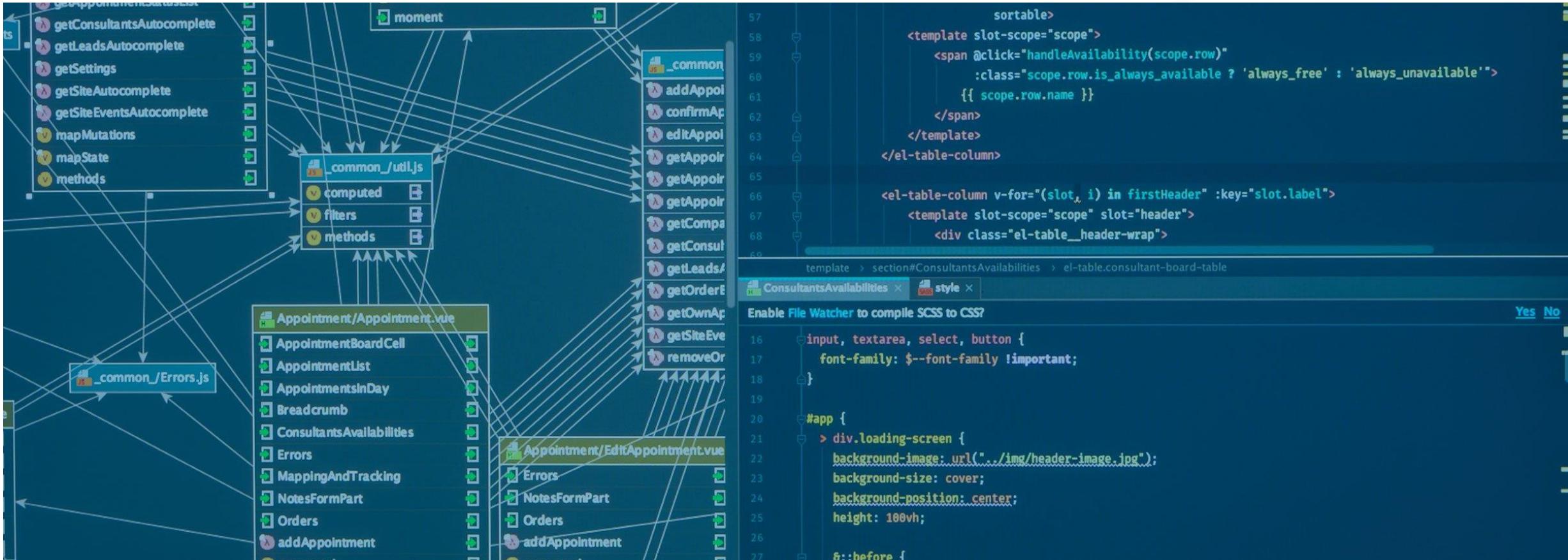
Default Partition: with no PARTITION BY clause, the entire result set is the partition.

WINDOW FRAME

is a set of rows that are somehow related to the current row. This

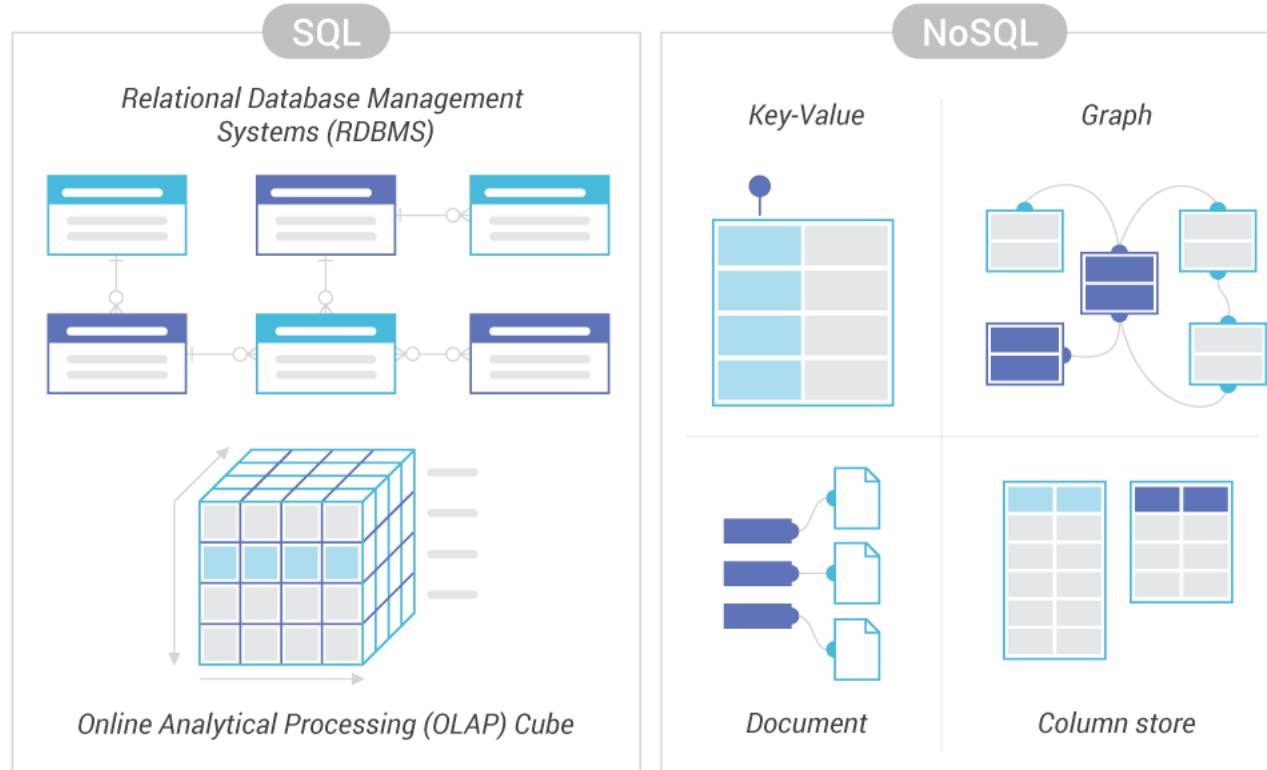


Relational Database Management System Design



<https://www.ucsc-extension.edu/courses/relational-database-design-and-sql-programming/>

SQL vs NoSQL, Data Lakes & Warehouse

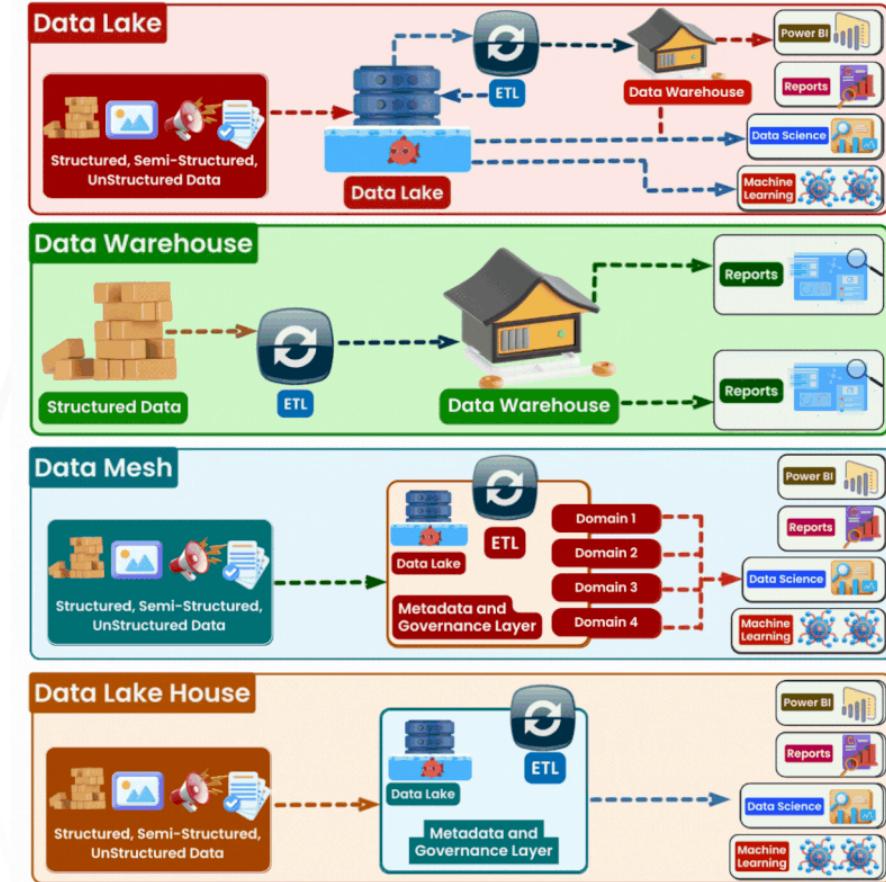


<https://www.scylladb.com/learn/nosql/nosql-vs-sql/>
[https://www.linkedin.com/posts/rocky-bhatia-a4801010_%F0%9D%90%83%F0%9D%90%9A%F0%9D%90%A%D%F0%9D%90%9A-%F0%9D%90%8B%F0%9D%90%9A%F0%9D%90%9A%D%F0%9D%90%9E-%F0%9D%90%9A%D%F0%9D%90%9A%D%F0%9D%90%9E-activity-7274385298258575360-m_p6/](https://www.linkedin.com/posts/rocky-bhatia-a4801010_%F0%9D%90%83%F0%9D%90%9A%F0%9D%90%A%D%F0%9D%90%9A-%F0%9D%90%8B%F0%9D%90%9A%F0%9D%90%9A%D%F0%9D%90%9E-%F0%9D%90%9A%D%F0%9D%90%9A%C-%F0%9D%90%83%F0%9D%90%9A%D%F0%9D%90%9A%D%F0%9D%90%9A%D%F0%9D%90%9E-activity-7274385298258575360-m_p6/)



Save for later

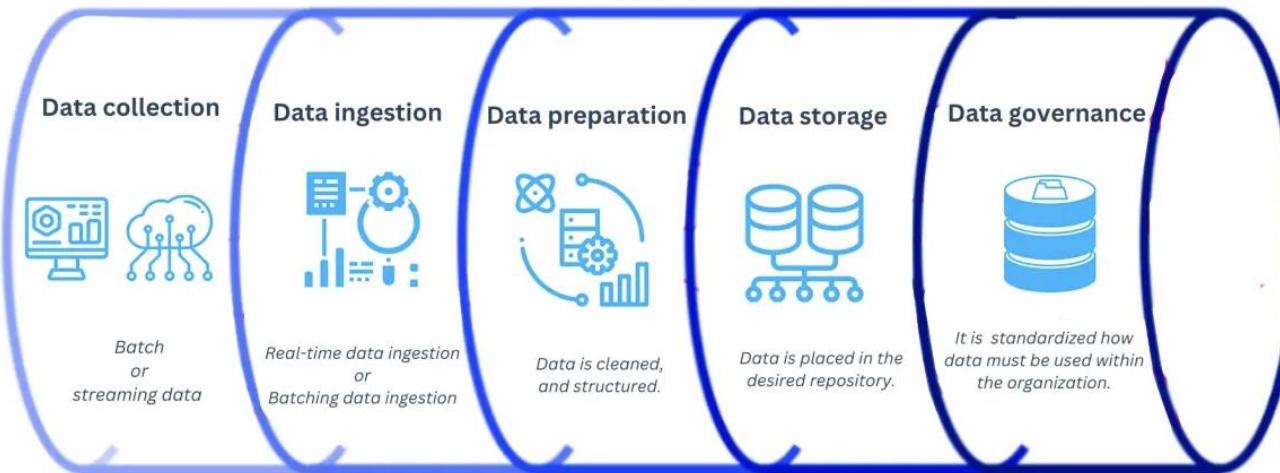
Data Lake Vs. Data Warehouse Vs. Data Mesh Vs Data Lake House



Data Pipelines and Cloud

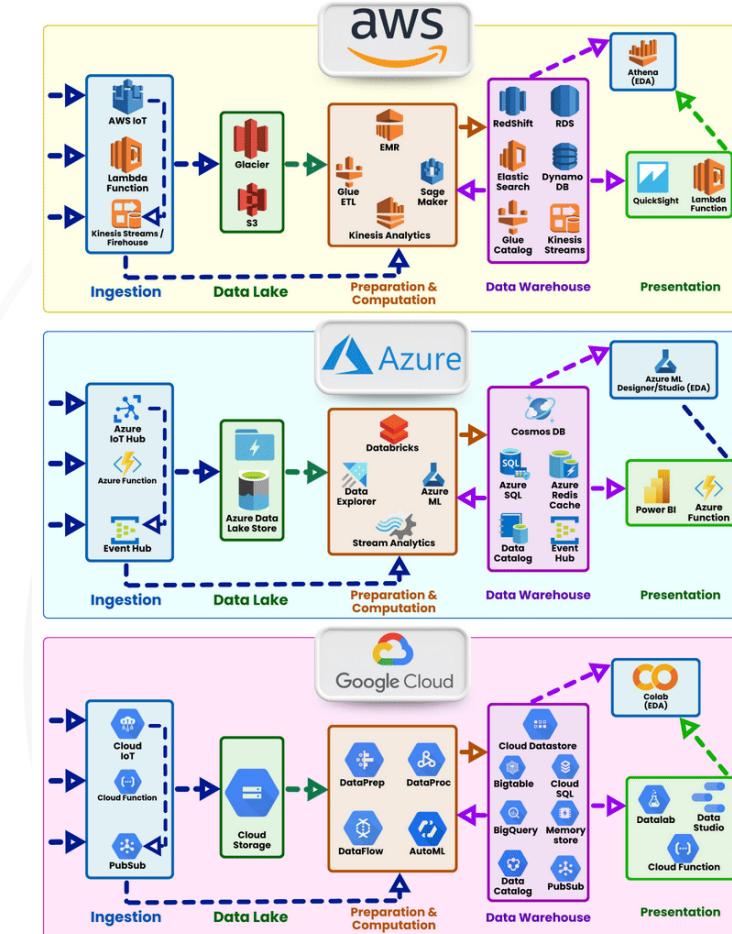
Created by:
Rocky Bhatia [in](#)

 learnwithrockybhatia.com
Your Gateway for Tech Learning

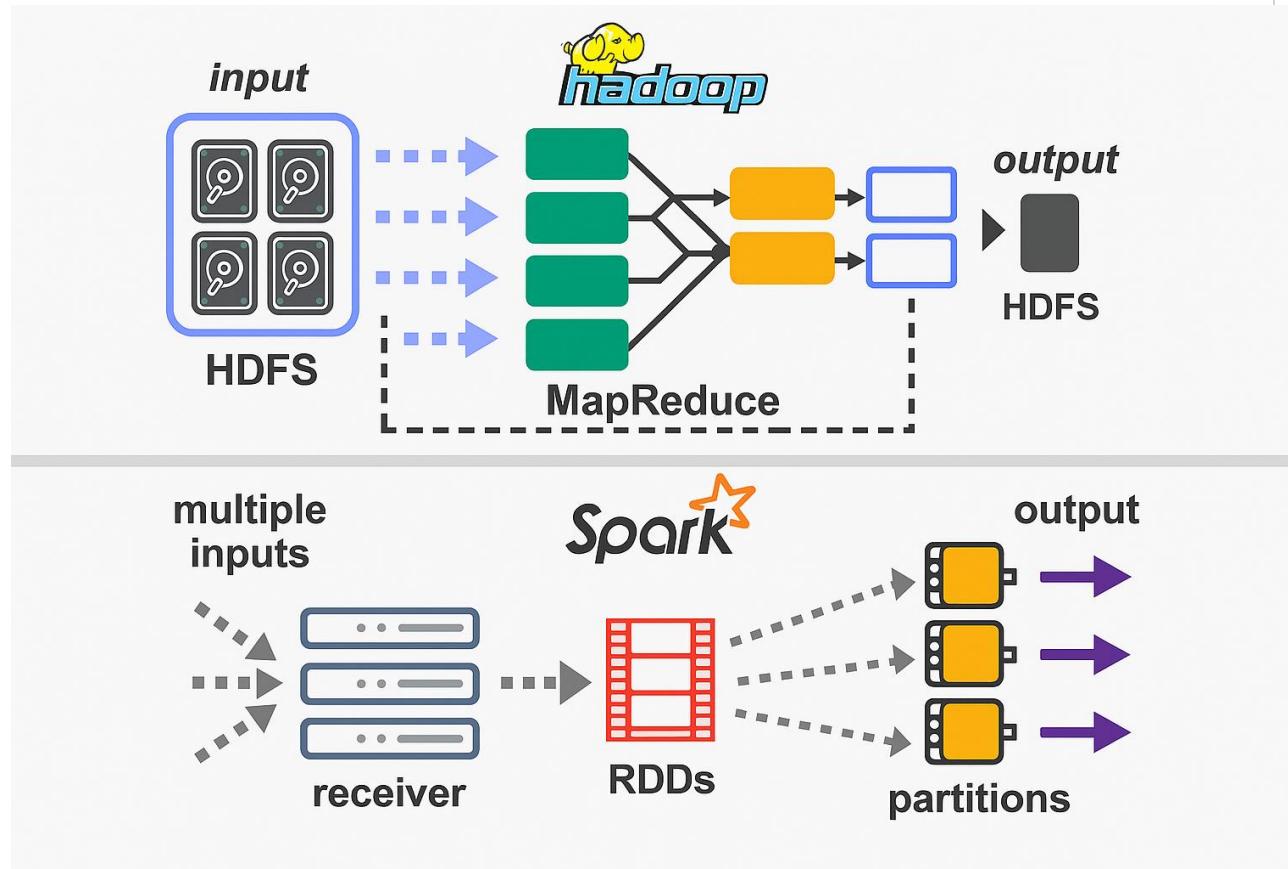


<https://www.secoda.co/blog/10-best-practices-to-build-data-pipelines>
<https://rockybhatia.substack.com/p/data-pipelines-in-the-cloud-azure>

Data Pipelines on AWS, Microsoft Azure and GCP



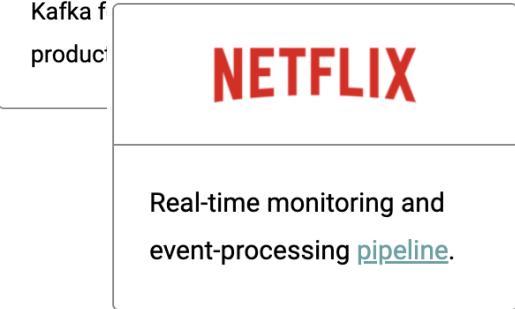
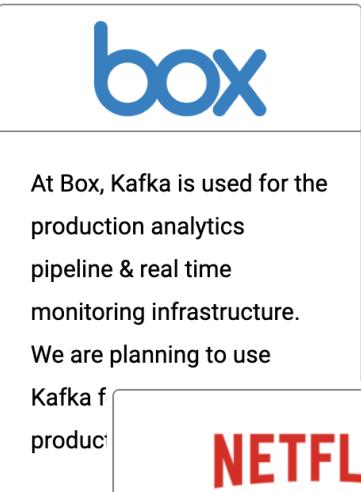
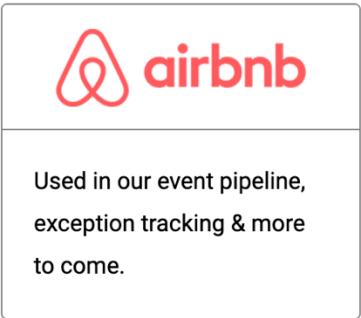
Big Data Processing: Spark



<https://phoenixnap.com/kb/hadoop-vs-spark>

```
+ comparison_BMw.py > ...
92
93     # Function to benchmark PySpark
94     def benchmark_pyspark(file_path):
95         results = {}
96         spark = SparkSession.builder.appName("BMWAnalysis").getOrCreate()
97
98         start = time.time()
99         df = spark.read.csv(file_path, header=True, inferSchema=True)
100        results['load'] = time.time() - start
101
102        start = time.time()
103        filtered = df.filter(col('Sales_Classification') == 'High')
104        filtered.count()
105        results['filter'] = time.time() - start
106
107        start = time.time()
108        agg = df.select(mean('Price_USD')).collect()
109        results['agg'] = time.time() - start
110
111        start = time.time()
112        group = df.groupBy('Region').agg(_sum('Sales_Volume')).collect()
113        results['group'] = time.time() - start
114
115        start = time.time()
116        df.write.mode("overwrite").parquet("spark_output.parquet")
117        results['write'] = time.time() - start
118
119        spark.stop()
```

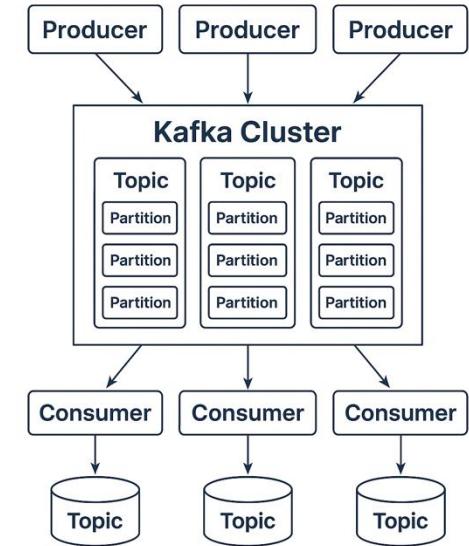
Real-time Streaming



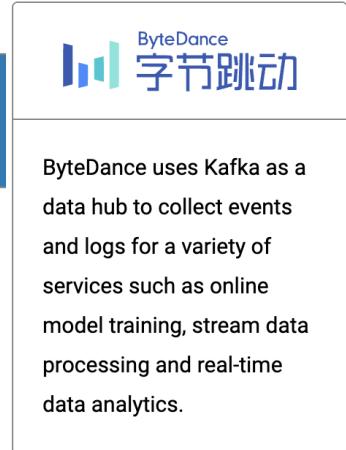
APACHE KAFKA

More than 80% of all Fortune 100 companies trust, and use Kafka.

Apache Kafka is an open-source distributed event streaming platform used by thousands of companies for high-performance data pipelines, streaming analytics, data integration, and mission-critical applications.



Apache Kafka is used at LinkedIn for activity stream data and operational metrics. This powers various products like LinkedIn Newsfeed, LinkedIn Today in addition to our offline analytics systems like Hadoop.



Data Orchestration with Apache Airflow

The screenshot displays the Apache Airflow web interface with two main panels. The left panel shows a DAG named 'demo' with a single run from April 14, 2025, at 17:54:11. It contains two tasks: 'hello' (BashOperator) and 'airflow' (@task). Both tasks are marked as successful. The right panel shows a DAG named 'toy_chain_linear_vs_chain_complex' with a run from April 14, 2025, at 17:59:04. This DAG has many tasks, mostly labeled 'chain_t...' or 'empty_t...', which are shown in a grid view. A histogram indicates the distribution of task execution times. Below the histogram, sections show failed tasks and failed runs, along with recent task logs.

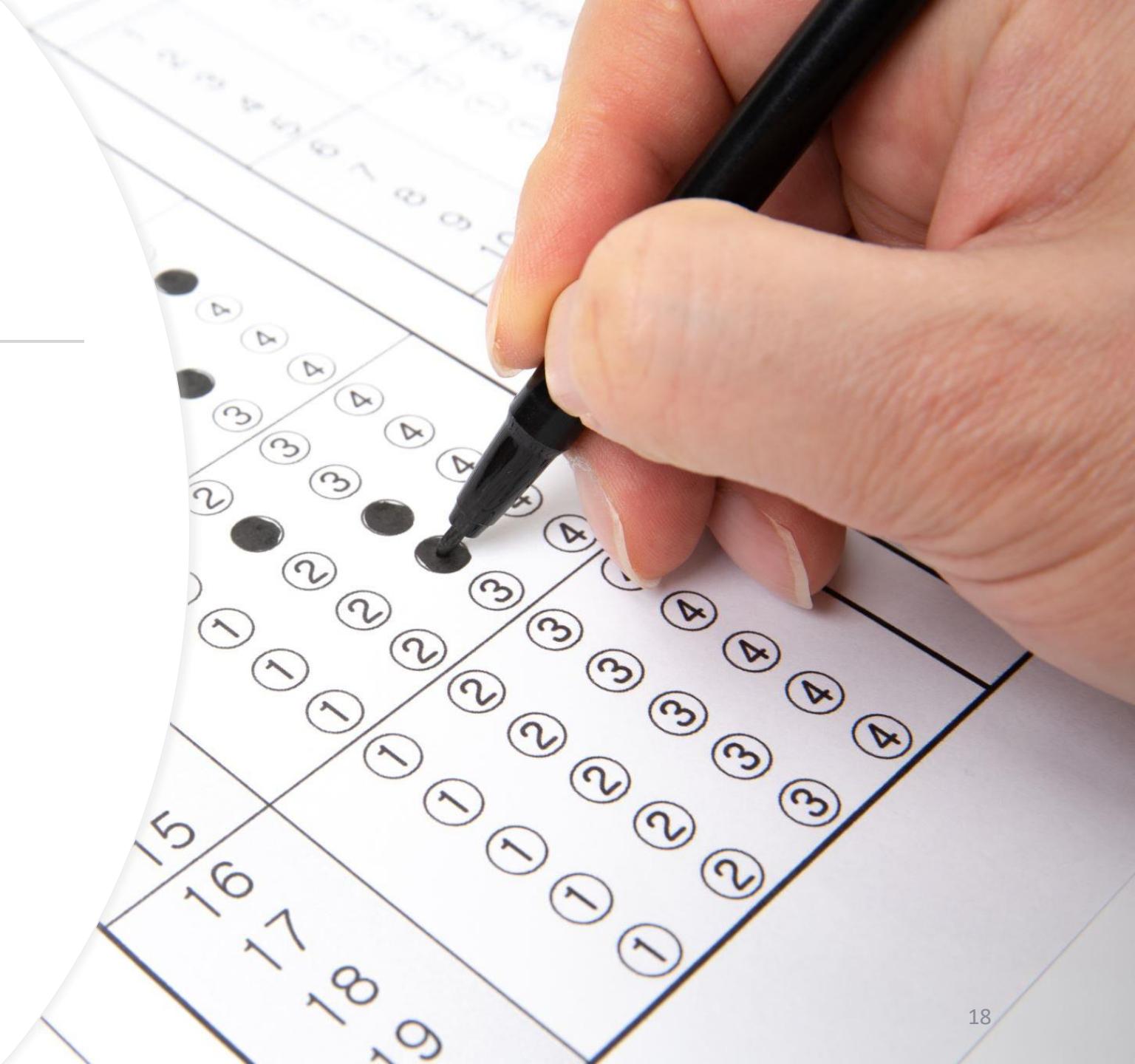
<https://airflow.apache.org/docs/apache-airflow/stable/index.html>



Apache
Airflow

Grading

- Grades will be based on:
 - **Mini Assignments (45%)**
 - **Major Assignments (30%)**
 - **Final Project (25%)**
 - **Participation and Bonus points (+)**
- Grades will be posted on Canvas



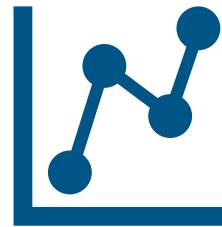
Late Policy and Assignment Policy

- All mini-assignments are due at the same time each week if any, so that you can plan your schedule
- Late assignments
 - Less than 24 hours: 5% deduction
 - 1-3 days: 10% deduction
 - 4-7 days, 30% deduction
 - Greater than a week, you will receive 0 for the corresponding assignment if without a written explanation **in advance** regarding a situation of unavoidable emergency
- The **lowest mini-assignment grade will be dropped** for the final homework calculation.

Data Engineering vs Data Science vs Data Analytics



Data Engineers: Build the highways for data to travel—pipelines, storage systems, and infrastructure.



Data Scientists: Use advanced analytics and ML to predict trends and answer complex questions.



Data Analysts: Dig into the datasets provided by engineers to extract actionable insights.

Some Trends in Modern Data Engineering

- **Pay-as-you-go cloud infrastructure**

Cloud services like Snowflake and BigQuery enable scalable, low-cost data analysis without upfront investment.

- **Shift from ETL to ELT?**

Modern warehouses handle transformations post-load, leveraging SQL and cloud compute power.

- **Templated SQL & YAML pipelines**

Modular, reusable SQL with YAML configs for scalable pipeline development.

- **Data observability (DataOps)**

Monitoring + alerting tools for data similar to CI/CD in software.

- **Decentralized data ownership**

Teams manage their own data products, reducing reliance on centralized data teams.



Questions?



Thank you