# Data Engineering

Powered by





# Presenter



in /arifmazumder





Mohammed Arif has more than nineteen (19+) years of working experience in Information Communication and Technology (ICT) industry. The highlights of his career are more than seven (7) years of holding various senior management and/or C-Level and had five (5) years of international ICT consultancy exposure in various countries (APAC and Australia), specially on Big Data, Data Engineering, Machine Learning and AI arena.

He is also Certified Trainer for Microsoft & Cloudera.



# **Related Resources**

https://arif.works/de



# Disclaimer

- $\circ$  Data Engineering is a vast domain
- Emphasize on Practical / Industry Implementation
- Take it as Experience Sharing Session



## What is data engineering?

SQL	Python	ETL	Data warehouse
Data ingestion	ML/AI	DataOps	Cloud infra
Data modeling	Data quality	Data governance	Security
Data observability	Data versioning	Data cleaning	Optimizations



## What is data engineering? Data → "Flow"



#### Many producers and consumers



### Example - Business analytics (or BI)



## Example - Real-time



## Why is it important?



## Why is it important?



#### End-to-end data pipeline





#### End-to-end data pipeline





### End-to-end data pipeline





#### Compute vs Storage





#### Undercurrents





[Historical Context: Data Engineering]

#### 1980-2000: Data warehouses



BI engineer ETL developer Data warehouse engineer

On-prem (vs the Cloud) Monolithic data stores

Farif.works

### 2000s: The beginning of "Big Data"

Internet + Google, Amazon, etc.

The birth of the Cloud (built on cheap, commodity hardware)

AWS starting with S3 and EC2

2003 Google publishes a paper on MapReduce - Hadoop

## 2000s-2010s: Big Data engineering

Big Data = Hadoop

Hadoop Ecosystem: Hadoop, YARN, HDFS

- A lot of work to manage
- Requires a specialized engineering team



### 2020s: Modern Data Stack



	Start		-> Lead
	Start	Scale	Lead
Typical use case	Data analytics	+ Data products	+ Real-time, ML/AI
Data volume	Small	Medium to large	Large to huge
Team size	Small team	Bigger team	+ Data science team
Tech maturity	Very basic (e.g. monolithic DB)	Proper data practices (e.g. automation,	ML/AI features Enterprise features



Data engineer's place within the data team





#### **Responsibilities - business & technical**

#### **Business responsibilities**

Communicate with both technical and non-technical people

Understand how to scope and manage a data project

Minimize cost and work within budget

#### **Technical responsibilities**

Create a good data architecture

Build and manage a reliable and performant data pipelines

Security, data governance, automation, observability, etc.



#### Data Engineers:

#### Required technical skills









#### Generation [Source Systems]





#### Structured vs unstructured data



### Types of databases

Relational database (RDBMS) Non-relational database (NoSQL)

Key-value store

#### Relational database (RDBMS)



order_id	amount	
1	\$20	
2	\$25	
3	\$10	



price

supplier

float

varchar

#### NoSQL database

Relational databases	
Tabular structure Strict schema Normalized Single-machine	

customer_id	name	
1	John Doe	
2	Jane Doe	
3	Alan Smith	

NoSQL databases

JSON structure Flexible schema De-normalized Distributed



#### Key-value store



#### Third-party systems via API



#### **Event streams**

#### Event producers



#### **Event streams**



#### **Event streams**



## Messages arriving in the wrong order


#### Exactly-once vs at-least-once delivery



# Idempotency [An operations is "idempotent" when the same result comes out no matter how many times you run it.]



J Farif.works

#### Why is idempotency so important?





#### Popular event streaming platforms





#### Caveats

Do NOT consider the source systems as someone else's problem









# Storage



## Storage hardware

HDD (Hard Disk Drive)



#### Memory (RAM or Random Access Memory)



SSD (Solid State Drive)



## **Storage formats - serialization**

: Turning data into byte streams to easily save or transport it



We serialize the data into standard format which is sent around and deserialized on the receiving end

#### **Storage formats - serialization**

# Row-based serialization file formats

- XML
- JSON
- CSV
- Fast lookup of individual rows

# Column-based serialization file formats

- Parquet
- ORC
- Avro
- Arrow
- Aggregation by Column

Row-based databases

#### Columnar databases

## **Compression and caching**

#### Compression

- · Reduce the amount of data stored
- Faster query (less to search)
- Faster transport (less data to move)

#### Caching

- Some storages are faster to access
- Memory is 1,000x faster than SSDs and CPU caches are orders of magnitude faster than memory
- Faster storage is more limited and expensive

## Single machine vs distributed storage



## Strong vs eventual consistency

Strong consistency





#### **Eventual consistency**



## ACID vs BASE

Single-machine transactional database

#### ACID

- Atomicity
- Consistency
- Isolation
- Durability



#### BASE

- Basically Available
- Soft-state
- Eventual consistency

## Type of storage systems



persisted

## File storage vs object storage



## OLTP vs OLAP (or Row vs columnar)

order_id	customer_id	order_date	amount
1	а	Jan 29	1000
2	f	Jan 30	1250
100	d	Mar 10	3500

OLTP (row-based database) in memory:

✓ Specific row, Single transactions

 Online Transaction Processing (OLTP) : row-based database (e.g. PostgreSQL, MongoDB)

1	а	Jan 29	1000	2	f	Jan 30	1250		100	d	Mar 10	3500	$\rightarrow$
---	---	--------	------	---	---	--------	------	--	-----	---	--------	------	---------------

OLAP (column-based database) in memory:

Analytics use case ( aggregating total Orders, total customers etc. ) Online Analytical Processing (OLAP) : column-based database

 $\checkmark$ 

 $\checkmark$ 

## OLTP vs OLAP (or Row vs columnar)

HDD (Hard Disk Drive)



OLAP (columnar) databases

Volumes are important

SSD (Solid State Drive)



OLTP (row-based) databases

#### Speed is important

## Data warehouse, lake, lakehouse



- - Good integration with external query engines

# Separation of compute from storage





## Data storage lifecycle - hot, warm, cold









# Ingestion







# Frequency of ingestion: batch, micro-batch, streaming

Batch	Micro-batch	Streaming
<ul><li>Daily or hourly</li><li>Ex. Tax reporting</li></ul>	<ul> <li>Minutes to seconds</li> <li>Ex. Near real-time analytics</li> </ul>	<ul><li>Seconds to sub-seconds</li><li>Ex. Uber</li></ul>

# ETL vs ELT





ELT



## Scalability and schema changes

#### Scalability

- Make sure your system can scale up and down automatically based on workload
- For streaming data, add a sufficient buffer for fluctuating event volumes

#### Schema changes

- Upstream changes to fixed schema will break the ingestion pipeline
- Schema migration features help, but not perfect - frequently communicate with upstream teams

# Ways to ingest data





## Ways to ingest data - CDC

#### Customers table





arif.works

#### Customers table

customer_id	Address		customer_id	Address
1	123 Disney St	$\rightarrow$	1	456 Pixar Way

#### Customers table log

{

}

op: "update", field: "address", before: "123 Disney St", after: "456 Pixar Way"



## Ways to ingest data





# 3.2 Data Transformation

# Transformations



# Queries

order_id	amount	
1	1250	· • • •
		····
98292	2000	

SELECT \* FROM orders WHERE order\_id=1

ALTER TABLE orders ADD shipping\_amount INT

# Queries - DDL, DML, DCL

DDL (Data definition language)

CREATE table customers

DML (Data manipulation language)

SELECT \* FROM orders WHERE order\_id=1

DCL (Data control language)

GRANT SELECT ON main\_db TO user\_name Dave
# What happens when you run a query?





# Improving the query performance

- Reduce the search space (prune, filter, etc)
- Use "EXPLAIN" SQL statement
- Pre-join frequently joined tables
- Vacuum dead records
- Leverage caching

## Querying streaming data





#### Data modeling





### Why do we need data modeling?

Easier to understand how different business data are connected

Allows you to have consistent definitions



## Ingredients of successful data modeling

Communicate with the data consumers from the beginning

Model the data at the lowest granularity possible



### Normalization

Data modeling with no duplicates

DRY Principle - "Don't Repeat Yourself"



### Normalization example

#### Denormalized "orders" table

id	order_date	order_items	customer_email	customer_name	amount	weight
1	Feb 12, 2024	MacBook, iPhone	seungchan@test.com	Seungchan Lee	350000	6
2	Feb 17, 2024	iPhone, iPad	nami@xyz.com	Nami Kim	250000	2.5
3	Mar 4, 2024	MacBook	nami@xyz.com	Nami Kim	250000	5
4	Apr 20, 2024	Vision Pro	nami@test.com	Nami Kim	300000	2

"orders" table

id	order_date	order_items	customer_id	amount
1	Feb 12, 2024	1, 3	1	350000
2	Feb 17, 2024	3, 4	2	250000
3	Mar 4, 2024	1	2	250000
4	Apr 20, 2024	2	2	300000

### Normalized tables example

#### "orders" table

id	order_date	order_items	customer_id	amount
1	Feb 12, 2024	1, 3	1	350000
2	Feb 17, 2024	3, 4	2	250000
3	Mar 4, 2024	1	2	250000
4	Apr 20, 2024	2	2	300000

#### "products" table

id	name	price	weight
1	MacBook	250000	5
2	Vision Pro	300000	2
3	iPhone	100000	1
4	iPad	150000	1.5

#### "customers" table

id	email	full_name
1	seungchan@test.com	Seungchan Lee
2	nami@xyz.com	Nami Kim



### Normalized tables - ER diagram





### Different modeling techniques





### Wide tables

#### "orders" table

id	order_date	order_items	customer_id	amount
1	Feb 12, 2024	1, 3	1	350000
2	Feb 17, 2024	3, 4	2	250000
3	Mar 4, 2024	1	2	250000
4	Apr 20, 2024	2	2	300000

#### "products" table

id	name	price	weight
1	MacBook	250000	5
2	Vision Pro	300000	2
3	iPhone	100000	1
4	iPad	150000	1.5

#### "customers" table

id	email	full_name
1	seungchan@deepintuitions.com	Seungchan Lee
2	nami@deepintuitions.com	Nami Kim



### Wide tables - denormalized order items

id	order_id	order_date	order_item_id	order_item	customer_email	customer_name	amount	weight
1	1	Feb 12, 2024	1	MacBook	seungchan@test.com	Seungchan Lee	250000	5
2	1	Feb 12, 2024	3	iPhone	seungchan@test.com	Seungchan Lee	100000	1
3	2	Feb 17, 2024	3	iPhone	nami@test.com	Nami Kim	100000	1
4	2	Feb 17, 2024	4	iPad	nami@test.com	Nami Kim	150000	1.5
5	3	Mar 4, 2024	1	MacBook	nami@test.com	Nami Kim	250000	5
6	4	Apr 20, 2024	2	Vision Pro	nami@test.com	Nami Kim	300000	2

- Faster Analytical Queries



### Transformations







### Views and materializations

#### Tables

Stored in physical storage Retrieved from storage Requires less compute Data could be stale Views (~virtual tables)

Not stored in physical storage Re-computed every time Requires more compute Data is always fresh

Materialization is a process of persisting the query in physical storage



### Events vs states





### Streaming vs batch









# Serving



## Analytics (BI or Business Intelligence)





Think backward : Why this metrics are useful? Trust : Trust is very hard to gain but easy to lose

## ML





## Deep Learning

ML

Structured (tabular) data Limited # of features (100-1,000s) Smaller datasets Relatively small compute (CPUs)

#### Deep Learning

Unstructured data - e.g. images Large # of features (millions) Much larger datasets Requires lots of compute (GPUs)

## Reverse ETL

Ingestion (or ETL/ELT)

**Reverse ETL** 







## DataOps

DevOps

Build/manage cloud infra Observability of cloud infra Build automated CI/CD pipeline DataOps (~DevOps for data)

Build/manage cloud infra for data tools Observability of data systems Automation (including CI/CD)



Observability and monitoring

Data fails silently!



### Automation

Cron jobs



#### Orchestration



## Automation - CI/CD pipeline

Most small teams overlook DataOps

Learn from software engineering

#### Example CI/CD pipeline



## Orchestration

### **Orchestration layer**



## **Orchestration - DAG**



## Why is orchestration so important?



## Why is orchestration so important?

### End-to-end visibility



# Why is orchestration so important?

Data lineage



## **Orchestration tools**





# Other undercurrents



- Authorization or "Access Control"
- Principle of least privilege

- Coming into focus
- High profile breaches
- Personally identifiable information (PII)


# Data quality

More nuanced than it seems at first

It's all about TRUST! Data fails silently

Validation Testing Observability



# Data quality

Software problems only happen when you deploy software, but data problems can happen independently from code

Data testing is different from software testing - it's trickier



# **Development workflow**



#### Benefits

- Fast iteration
- Isolation (safe to tinker without affecting prod env)
- Cost reduction (use local laptop for dev)

#### Challenges

- Sampling a large dataset is not always easy
- Dev/prod environment parity not always possible
- Cloud tools not always available locally







# What is a good data architecture?

Performance	Scalability	Reliability	Expect failure
Security	Modularity	Cost-effectiveness	Flexibility



# Lambda & Kappa architectures



# Modern data stack





# What factors should you consider when designing your data architecture?





#### **BI Stack**

#### Data sources



#### BI stack v1 - Data warehouse



#### Central storage





#### **Orchestrator - Airflow vs Dagster**



No local env setup

Dependency management issues

Complex to deploy and manage

Local dev env

🌀 dagster

Cloud-native, easy deployment, CI/CD

Declarative (rather than imperative)

#### Imperative vs declarative

Imperative

The "how" Spell out the exact steps Declarative

The "what" Specify the desired outcome



🛞 kubernetes



#### Ingestion 🌖 dagster X dbt Core. Data S meltano **þ** AWS S3 + 🔆 snowflake $\rightarrow$ **CO** Superset sources Alternatives: **Fivetran** nirbyte Production Development DuckDB

### Ingestion tool comparison

	<b>K</b> Fivetran	nirbyte	S meltano	dltHub
Open-source?	×	$\checkmark$	$\checkmark$	$\checkmark$
UI or code-first	UI	UI	Code	Code
# of connectors	100s	100s	100s	<20
Embedded?	×	×	×	$\checkmark$

#### Transform



#### Visualization



#### BI stack v2 - Data lakehouse



#### Apache Iceberg



#### ML stack





#### Model training





#### Model serving





#### Labeling

#### For training the AI model



	Al model		Output		Labe
>	?	$\rightarrow$	?	$\leftrightarrow$	T/F

#### Dataset (human labeled)

id	img_url	label
1	s3://john.jpg	True
2	s3://lisa.jpg	False
	•••	
10000	s3://blah.jpg	True



#### Deep Learning stack





#### LLMs?





#### Other considerations



# Great ! We've just finished the overview of **Data Engineering**

