# Intro to Data Engineering

Mils Burasakorn

# Burasakorn Sabyeying
## (Mils, มิลส์, มิล, มิว)

Mesodiar.com

Data Engineer @ CJ Express (TILDI team)

# Agenda

- **Why** do we need Data Engineer?

- **Who** is Data Engineer

- **What** does Data Engineer do?

- **Where** does Data Engineer stand

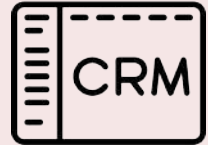- **How** do I become Data Engineer?

# Why do we need Data Engineer?
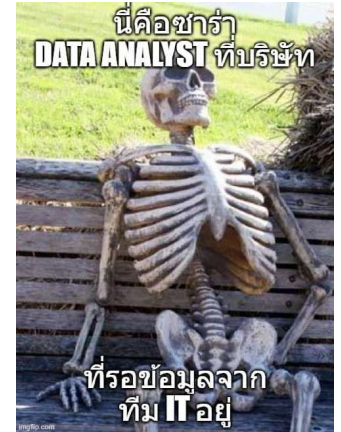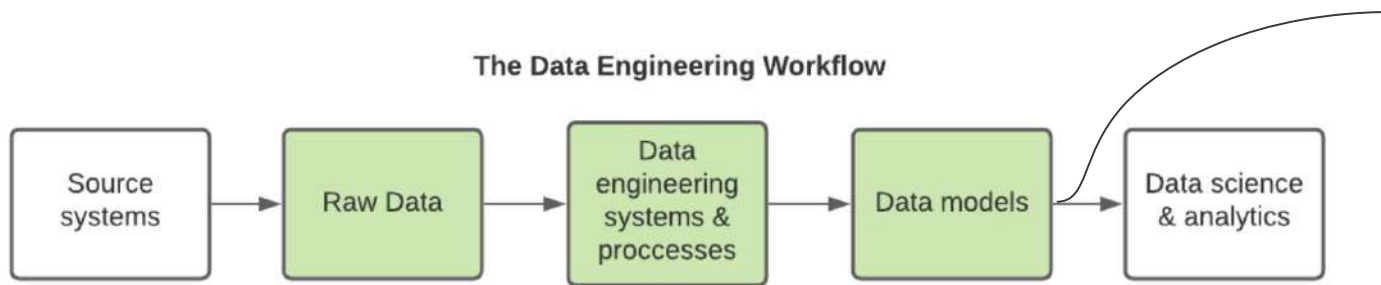
Database

Google Analytics
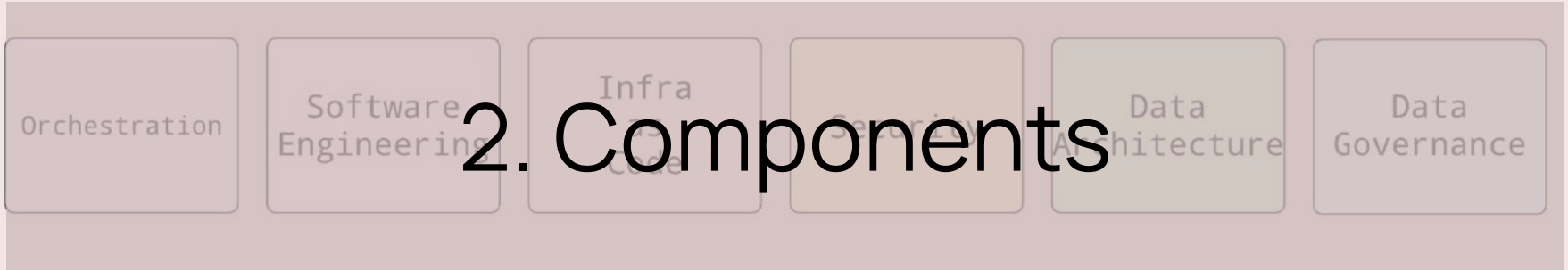
salesforce

SAP®

CRM

and more data sources..

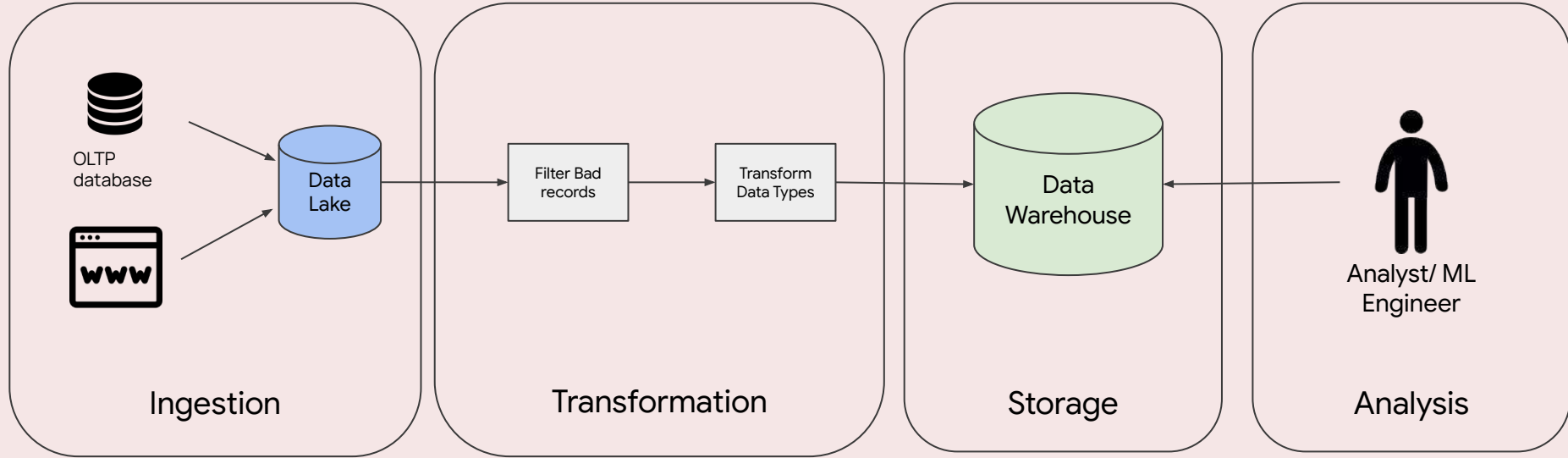Cr. Data TH.com – Data Science ชิลชิล

# What is Data Engineer?

data engineers set up and operate the organization's data infrastructure preparing it for further analysis by data analysts and scientists
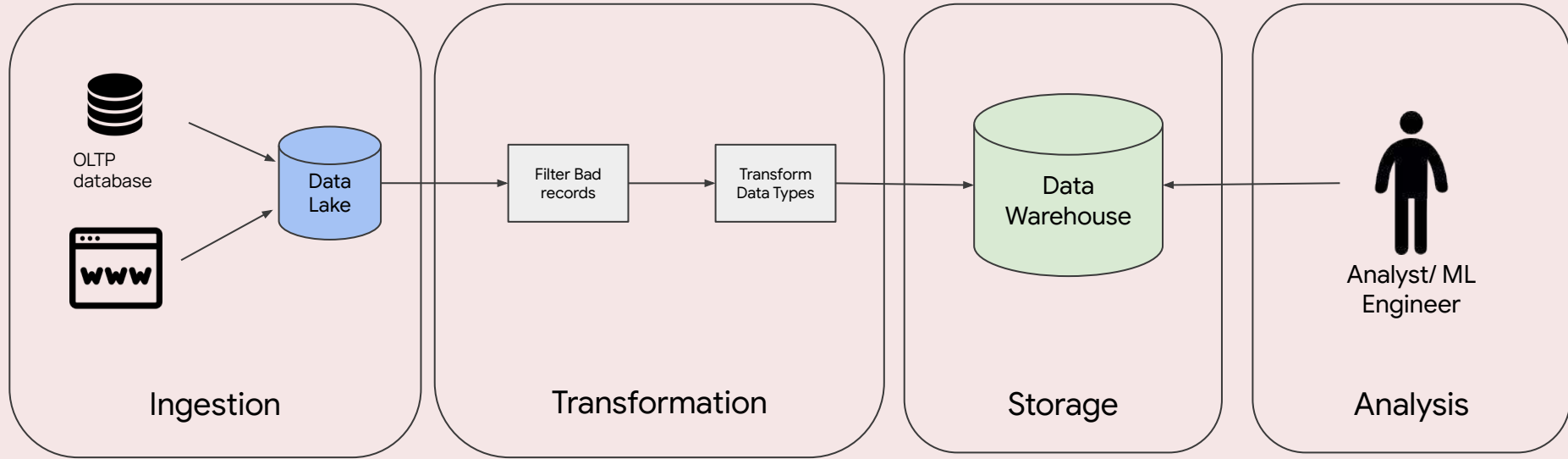
## The Data Engineering Workflow

Source systems → Raw Data → Data engineering systems & proccesses → Data models → Data science & analytics

นี่คือซาร่าๆ
DATA ANALYST ที่บริษัท
ที่รอข้อมูลจาก
ทีม IT อยู่

# Data Engineering Life Cycle



1. Pipeline

2. Components

# 4 stages in Data Pipeline



OLTP database

www

Data Lake

Filter Bad records

Transform Data Types

Data Warehouse

Analyst/ ML Engineer

Ingestion

Transformation

Storage

Analysis

# Data Lake vs Data Warehouse vs Data Mart



Ingestion | Transformation | Storage | Analysis

OLTP database → Data Lake → Filter Bad records → Transform Data Types → Data Warehouse ← Analyst/ ML Engineer

Data Lake → transform → Data Warehouse → Data Mart

# ETL vs ELT

**E**xtract - pull data from all your data sources
**T**ransform - clean and process data
**L**oad - load into storage destination

ML

| Data Sources | Ingestion | Transformation | Storage | Analysis/Visualization |

OLTP, Excel, Google Analytics — Data Sources

python, kafka, Airbyte — Ingestion (Storage)

Apache Spark, pandas, python — Transformation

Google BigQuery, Cloud Storage, amazon S3 — Storage

Metabase, tableau, Google Data Studio — Analysis/Visualization

Data Scientist

Data Analyst

Data Engineer

Orchestration — Apache Airflow

Data Lineage, Data Monitoring

# Structure

**Unstructured data**

The university has 5600 students.
John's ID is number 1, he is 18 years old and already holds a B.Sc. degree.
David's ID is number 2, he is 31 years old and holds a Ph.D. degree. Robert's ID is number 3, he is 51 years old and also holds the same degree as David, a Ph.D. degree.

**Semi-structured data**

```
<University>
 <Student ID="1">
  <Name>John</Name>
  <Age>18</Age>
  <Degree>B.Sc.</Degree>
 </Student>
 <Student ID="2">
  <Name>David</Name>
  <Age>31</Age>
  <Degree>Ph.D. </Degree>
 </Student>
 ....
</University>
```

**Structured data**

| ID | Name | Age | Degree |
|----|---------|-----|--------|
| 1 | John | 18 | B.Sc. |
| 2 | David | 31 | Ph.D. |
| 3 | Robert | 51 | Ph.D. |
| 4 | Rick | 26 | M.Sc. |
| 5 | Michael | 19 | B.Sc. |

Text, Audio, Video, PDF, Internet of Things (IoT) sensor data

XML, CSV, JSON, Web pages

PostgreSQL, MySQL

# Q: Facebook post เป็นข้อมูลแบบไหน?

1. Structured data
2. Semi-structured data
3. Unstructured data

🍌 Google Sheets Dashboard !! ปั่นแดชบอร์ดสนุกเลย ง่ายเหลือเชื่อ 555+

.

เมื่อเข้าสอน Free Fire vs. PubG Facebook Post Analysis สนุกมาก เรียนกันแบบเน้นๆ สอนสดแบบสดจริงๆ 555+

.

สรุป Steps ที่เราสอนในคลาส... **See more**

# Databases

Relational(SQL)

**Traditional database/ DBMS**

Non-relational (NoSQL)

# Databases

## Relational (SQL)

### Row - oriented

| 1 | red | square | .. |
|---|-----|--------|-----|
| 2 | blue | square | .. |
| 3 | yellow | square | .. |
| 4 | none | triangle | .. |

| square | 4 |
|--------|---|
| triangle | 3 |

## Non-relational (NoSQL)

### Key-Value

Key → Value

Key → Value

Key → Value

### Columnar

**HBase table**

| | color | shape |
|---|-------|-------|
| "first" | "red": "#F00", "blue": "#00F", "yellow": "#FF0" | "square": "4" |
| "second" | | "square": "4" "triangle": "3" |

### Document

```
{
   "_id": "ObjectId("d7caskf00010dsa")";
   "firstName": "Burasakorn",
   "lastName": "Sabyeying",
   "nickName": "Mils",
   "role": "Data Engineer"
}
```

### Graph

# Databases

## Relational (SQL)

### Traditional database/ DBMS



## Non-relational (NoSQL)

### Key-Value



### Columnar



### Document



### Graph

# OLTP vs OLAP

## Online **Transaction** Processing

- captures, **stores**, and processes data from transactions **in real time**

- banking and credit card activity or retail checkout scanning.

- **traditional DBMS**

- Based on INSERT, UPDATE, DELETE commands

## Online **Analytical** Processing

- analyze **aggregated historical data** from OLTP systems.

- designed for use by data scientists, business analysts

- For **data warehouse** and data mart applications

- Based on SELECT commands to aggregate data for reporting

# Scenario

OLTP

OLAP

```
SELECT *
FROM table
WHERE date=today
```

```
SELECT
        country.country_name_eng,
        SUM(CASE WHEN call.id IS NOT NULL THEN 1 ELSE
0 END) AS calls,
        AVG(ISNULL(DATEDIFF(SECOND, call.start_time,
call.end_time),0)) AS avg_difference
FROM country
LEFT JOIN city ON city.country_id = country.id
LEFT JOIN customer ON city.id = customer.city_id
LEFT JOIN call ON call.customer_id = customer.id
GROUP BY
        country.id,
        country.country_name_eng
HAVING AVG(ISNULL(DATEDIFF(SECOND, call.start_time,
call.end_time),0)) > (SELECT AVG(DATEDIFF(SECOND,
call.start_time, call.end_time)) FROM call)
ORDER BY calls DESC, country.id ASC;
```

# Data Lake, Data Warehouse, Data Lakehouse



Data Lake

Unstructured data

+

Data Warehouse

Structured data

=

Data Lakehouse

Unstructured data

# Data Lakehouse



(a) First-generation platforms.

(b) Current two-tier architectures.

(c) Lakehouse platforms.

# Centralized Data vs Decentralized Data

# Data Engineering Life Cycle

# Orchestration

# Orchestration

process of coordinating many jobs to run efficiently



Managed service

# Software Engineer



- Web Scraping

- Get data through API

1 page (1-100 rows)
Total: 300 pages
**300 API calls**

- Process data

  Pyspark, Pandas

- Pipeline as code

# Which language should I know?

# Which language should I know?

```
SELECT * FROM Customers;
```

SQL

Python

JVM languages
(Java, Scala)

Bash

# Infrastructure as a Code

## Containers

**important



## Provisioning



## Version control

**important

# Security

In both Data and System



**Who should see the data?**
Owner
Editor
Viewer

**Key management**



## Options for encryption at rest

Fully-automated management → Finer-grained customer control

**Default Google encryption**

World-class encryption without further need for configurations.

By default

**Customer-managed encryption keys (CMEK) using Cloud KMS**

Keep keys in the cloud for direct use by cloud services.

Available for BigQuery, Cloud Storage (Beta), Compute Engine (Beta), and Dataproc (Beta).

**Customer-supplied encryption keys (CSEK)**

Keep keys on premise, and use them to encrypt your cloud services.

Available for Cloud Storage and Compute Engine

Google Cloud

**Encryption**

# Data Architecture



Architecture != Tools

ต้องเข้าใจ
- **business requirement**
- Nature การเกิดของข้อมูล
- Nature การใช้ของข้อมูล

แล้วจะปรับสิ่งพวกนี้มา design ในการ serve data ยังไง

# On premises vs Cloud

**On premises/ On-prem**

**= purchase hardware/ data centers they own**

- Still default for established companies

- install/ maintain/ upgrade by their own

- Direct control over configuration, management, security

Cloud

**= Cloud provider** (AWS, Azure, Google Cloud, etc)

- Infrastructure as a Service(IaaS)

- **Serverless products** and managed service

- Billed on **pay-as-you-go**

- Unpredictable scale requirements

# Serverless

cloud-native development model that allows developers to build and run applications **without having to manage servers**

**Serverless = Still have server**

Google Cloud Platform (GCP) serverless products



| 2008 | 2009 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 |

# Hybrid and multi-cloud

Hybrid

Multi-cloud



Cr. cloudflare.com

Cr. xenonstack.com

# Best Data Pipeline ?



cr.  Richard Burlton on Unsplash

# Data Governance

*Data Governance: The Definitive Guide -*
"Data governance is, first and foremost, a **data management function** to ensure the quality, integrity, security, and usability of the data collected **by an organization**."

1. Data Monitoring

2. Data Discovery & Data Catalog

3. Data Lineage

4. Data Quality

# Data Monitoring

# Data Monitoring

# Data Discovery & Data Catalog

# Data Lineage

allows you to know where that data is stored and its dependencies.

# Data Quality* (!important)

According to *Data Governance: The Definitive Guide*, data quality is defined by **three main characteristics**

**Accuracy**
Is the collected data factually correct? Are there duplicate values? Are the numeric values accurate?

**Completeness**
Are the records complete? Do all required fields contain valid values?

**Timeliness**
Are records available in a timely fashion?

---

great_expectations  Home / Validations / taxi / 20211214-205807-my-run-name-template / 2021-12-14T20:58:07Z

passenger_count

**Actions**

Validation Filter:

Show All | Failed Only

How to Edit This Suite

Show Walkthrough

**Table of Contents**

Overview

Table-Level Expectations

passenger_count

| Status | Expectation | Observed Value |
|---|---|---|
| ✕ | minimum value must be greater than or equal to `1` and less than or equal to `1`. | 0 |
| ✓ | maximum value must be greater than or equal to `6` and less than or equal to `6`. | 6 |
| ✕ | mean must be greater than or equal to `1.5716` and less than or equal to `1.5716`. | 1.3577 |
| ✓ | median must be greater than or equal to `1.0` and less than or equal to `1.0`. | 1 |

quantiles must be within the following value ranges.

| Quantile | Min Value | Max Value |
|---|---|---|
| 0.05 | 1 | 1 |
| Q1 | 1 | 1 |
| Median | 1 | 1 |
| Q3 | 2 | 2 |
| 0.95 | 5 | 5 |

| Quantile | Value |
|---|---|
| 0.05 | 0 |
| Q1 | 1 |
| Median | 1 |
| Q3 | 1 |
| 0.95 | 5 |

✕ values must belong to this set: `1` `2` `3` `4` `5` `6`.

**1579 unexpected values found. ≈15.79% of 10000 total rows.**

Sampled Unexpected Values
0

≈15.79% unexpected

# สรุป



## Data Engineering Life Cycle

Ingestion → Transformation → Storage → Machine Learning / Analytics / Reporting

Orchestration | Software Engineering | Infra as Code | Security | Data Architecture | Data Governance
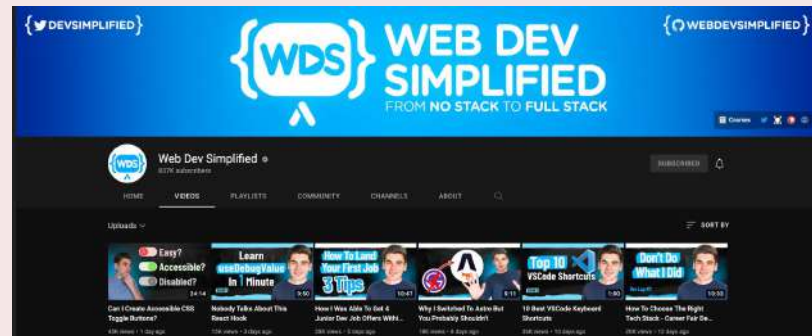
# How do I become Data Engineer?

# ~~How to become Data Engineer ?~~

# how I learn

- Data Engineer Roadmap:
  https://github.com/datastacktv/data-engineer-roadmap

- Compare tools
- Compare services https://comparecloud.in/

# Listen

**random\*\***

# READ !

medium.com

learning.oreilly.com

# Write !

Note in Notion

Mesodiar.com (blog on medium)

โน้ตว่าวันนี้เราเรียนรู้อะไร จะได้ไม่ท้อ !

# Facebook Group



เราอยากเป็นแบบไหน จงเอาตัวเองเข้าไปอยู่วงนั้น

**FB**: Mesodiar.com
**Medium**: www.mesodiar.com