

# EE3801 Cheatsheet

Intro to Data Engineering

github.com/securespider

## Architecture

### Batch-based DP

- $N$  independent tasks to process with  $k$  stages
- Each stage takes max of  $T$  time process input
- Diff stage can operate concurrently
- $t(N, k) = T \times (N + k - 1)$

### Streaming-based DP

- 

## 01.1 Intro

### Data science vs engineering

- **Science** - Learn, optimise, analytics, aggregate and labelling
- **Engineering** - Cleaning, data storage, logging, sensors, pipelines

### Data structure

#### Unstructured data

- Chaotic no order to data

#### Structured data

- Data stored access in the same format

#### Semi structured data

- Can contain both forms of data
- Some structure but not all data points follow same format

### Big data

#### Volume, Variety, Variability

#### Velocity

High rate of data generation

- Must create a robust and scalable pipeline

### Raw Data

- Tend to have gaps

### Data wrangling

Used to understand raw data

**Discovery** Understand what is in your data

#### Structure

**Cleaning** Dealing with gaps (nulls), outliers, formatting bugs

**Enrichment** Derive other data from other information/ additional data augmentation (feature selection)

**Validation** Verify data quality, sources

**Publishing** Give data scientist

### Process

**Extraction** Retrieve raw data from unstructured pool and migrate to temp repo

**Transformation** Structure enrich and convert raw data

**Loading** Loading structured data into data warehouse

### Data warehouse

Decision support system storing historical data from organisations

### Data Pipeline

- Processing underlying raw data in ordered sequence of steps

## 01.2 Data Pipelines

### Considerations

#### Big data

**Velocity** Streaming, captured and processed in real time

**Volume** Scalable wrt time

**Variety** Recognise and process diff formats

#### Business

- Handling streaming data?
- How much data to expect (Time horizon/how much storage consumed)
- What type/how much processing in DP
- Where is data source? Need micro-services?