# Big Data Fundamental
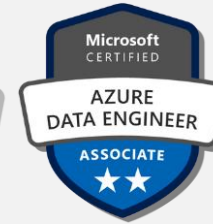
# Presenter

## Mohammed Arif, PhD
### Lead Data Scientist
**Big Data | Machine Learning | AI**

Mohammed Arif has more than fifteen (15) years of working experience in Information Communication and Technology (ICT) industry. The highlights of his career are more than six (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.

in /arifmazumder

# Agenda

- Data vs Big Data
- Big Data Characteristics
- Big Data Reference Architecture
- Big Data Ecosystem Components
- Analysis vs Analytics
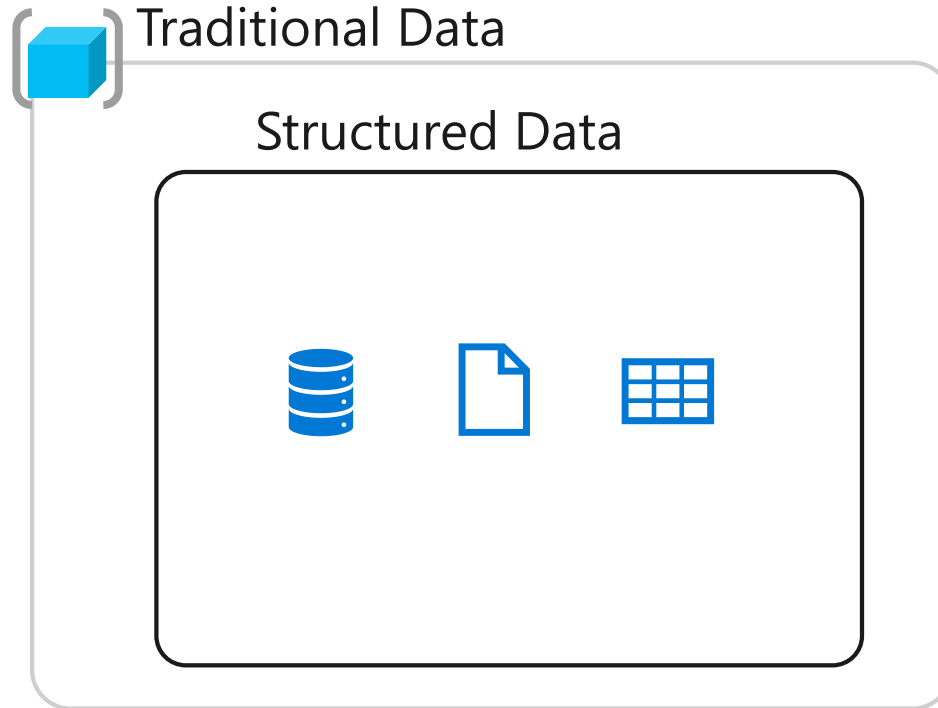- Data Analysis (Hands-on)
- Big Data – Career Path

SIMPLIFIED
ACADEMY
learn • earn • grow

Resource Link

http://arif.works/bdf/

# Data vs Big Data

# Data

Big Data

Both Structured and Unstructured Data

Traditional Data

Structured Data

Volume

Variety

Velocity

# Big Data Characteristics

**VOLUME**
- Amount of data generated
- Online & offline transactions
- In kilobytes or terabytes
- Saved in records, tables, files

**VELOCITY**
- Speed of generating data
- Generated in real-time
- Online and offline data
- In Streams, batch or bits

**VARIETY**
- Structured & unstructured
- Online images & videos
- Human generated - texts
- Machine generated - readings

SIMPLIFIED ACADEMY
learn • earn • grow

# Big Data Reference Architecture

In summary,

Generally Big Data Architecture Data Pipeline has five stages:

- Collection

- Ingestion

- Preparation

- Computation

- Presentation

| Collection | Ingestion | Preparation | Computation | Presentation |

| Data Sources | Integration | Data Storages | Analytics | Presentation |
|---|---|---|---|---|
| Structured | ETL | NoSQL Databases | Query & Reporting | Web Browsers |
| Semi-Structured | Messaging | Distributed File Systems | Map Reduce | Native Desktop |
| Unstructured | API | | Search Engines | Mobile Devices |
| | | | Advanced Analytics | Web Services |

SIMPLIFIED ACADEMY
learn • earn • grow

# ETL (Extract, Transform & Load)

# Data Lake

Data lake is one place to put all the data enterprises may want to use, including structured and unstructured data.

# Data Lake

# Data Warehouse

Single Source of Truth.

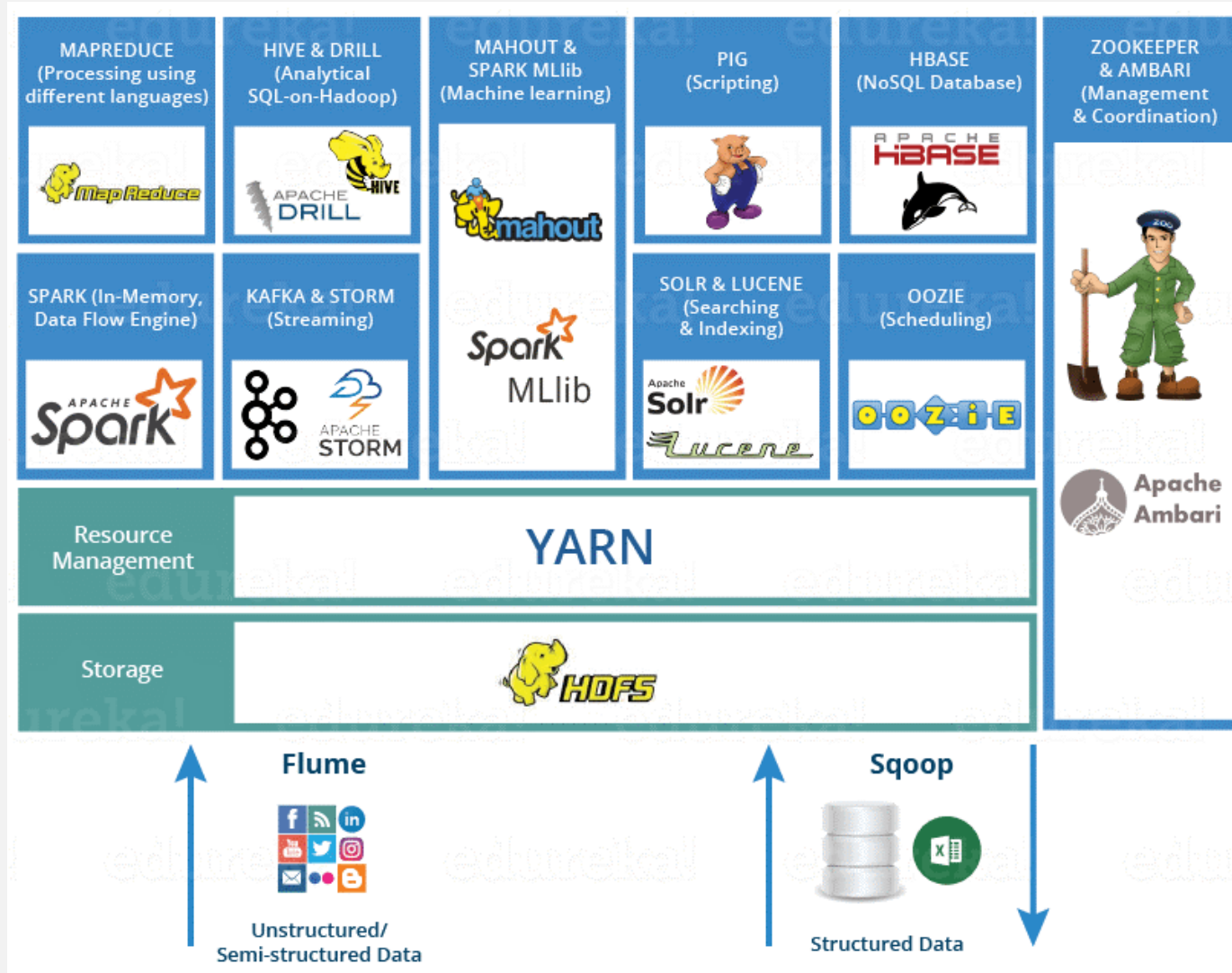Structuring all the Best Quality data in one place.

# Data Warehouse

# Hadoop

Hadoop is a collection of open source programs/procedures/platform relating to Big Data analysis. Being open source, it is freely available for use, reuse and modification (with some restrictions) for anyone who is interested in it. Big Data scientists call Hadoop the 'backbone' of their operations.

# Big Data/Hadoop Eco System Component

# Case Study (Uber)

Transformation Journey towards Big Data Platform.

Please read this article to get more info on how Big Data & its Architecture



https://eng.uber.com/uber-big-data-platform/

# Business Analytics



**Descriptive**
Explains what happened.

**Diagnostic**
Explains why it happened.

**Predictive**
Forecasts what might happen.

**Prescriptive**
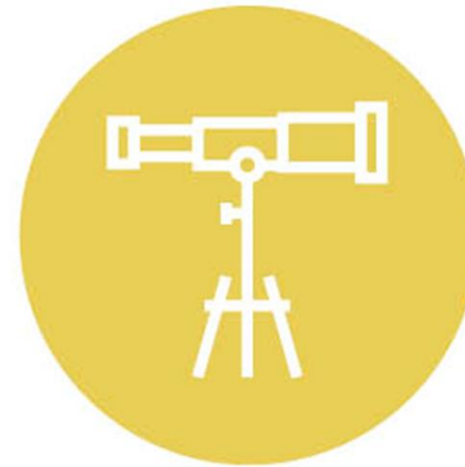Recommends an action based on the forecast.

# Business Analytics



**Descriptive**
What has happened?

**Diagnostic**
Why did it happen?

**Predictive**
What will happen next?

**Prescriptive**
What should I do?

Looking back | Looking forward

SIMPLIFIED ACADEMY
learn • earn • grow

# Hands-on

We will do some Data Analysis using BI Tools (Tableau)

Download Tableau Public

https://public.tableau.com/en-us/s/download

+⁺⁺ + t a b l e a u

# Data Analysis

❑ Data in Raw format might not help.

❑ Data transformation through calculation help to do :

    ❑ Draw better insights

    ❑ Generate Report

    ❑ Data Driven Decision

    ❑ Self Serving Data

**Get Dataset (Super Store Sales Data)**
https://drive.google.com/file/d/131VI-hVyeLFRwkFNqaa01N_Rk8NHlZvs/view

SIMPLIFIED
ACADEMY
learn • earn • grow

# Business Questions

❑ What is the growth of various Sub-Categories over 4 years?

❑ Which category in each segment is yielding more profit?

❑ Monthly fluctuations in sales in various years?

❑ Running total of sales in each year?

❑ Rank sub-categories based on Quantity sold and compare their profits.

❑ Find the average discount to Sales ratio for each sub-category in different regions.

❑ What is the average order to ship time for various sub-categories?

# Big Data – Career Path