Big Data
oo

Data warehouse
oo

Data Lake
oooooooooo

Methodology
ooooooooooo

# Slice 1
## Data Lake Architecture Overview

Le Nhu Chu Hiep

# Section 1

## Big Data

## What is big data

Data is too large for traditional storage to handle and retrieve.

Big Data
○○

Data warehouse
●○

Data Lake
○○○○○○○○○

Methodology
○○○○○○○○○○

# Section 2

## Data warehouse

Big Data
00

Data warehouse
○●

Data Lake
○○○○○○○○○

Methodology
○○○○○○○○○○

# What ?

- Centralize relation information for specific purpose.
- A RDBMS that has a schema to support BI tools and OLAP.

Big Data
oo

Data warehouse
oo

Data Lake
●oooooooo

Methodology
oooooooooo

# Section 3

## Data Lake

# What ?

- Centralized repository that allows you to store all your structured and unstructured data at any scale.
- Update version of data warehouse to handle large data with dirverse type.
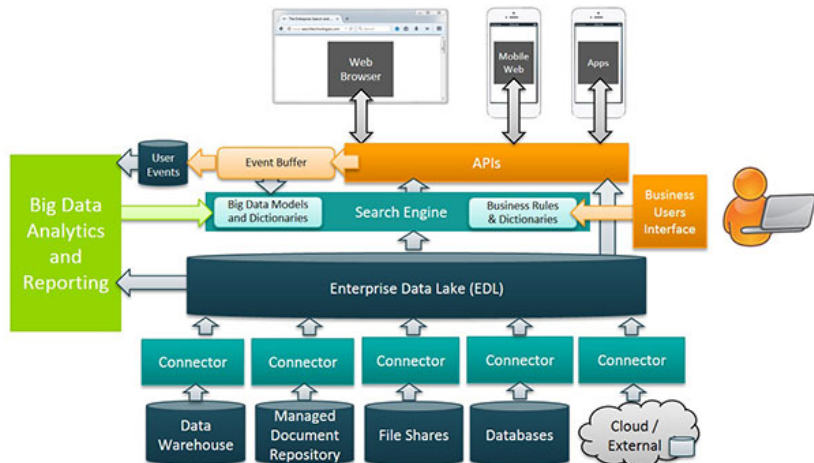- Big Data version of data warehouse

Big Data
oo

Data warehouse
oo

**Data Lake**
ooo●oooooo

Methodology
oooooooooo

# Why ?

- Data warehouse and its analytic tool can not handle big data.
- Requirement of storing and trieving unstructured data.

Big Data
○○

Data warehouse
○○

**Data Lake**
○○○●○○○○○

Methodology
○○○○○○○○○○

# Data Lake and Data Warehouse

| Characteristics | Data Warehouse | Data Lake |
|---|---|---|
| **Data** | Relational from transactional systems, operational databases, and line of business applications | Non-relational and relational from IoT devices, web sites, mobile apps, social media, and corporate applications |
| **Schema** | Designed prior to the DW implementation (schema-on-write) | Written at the time of analysis (schema-on-read) |
| **Price/Performance** | Fastest query results using higher cost storage | Query results getting faster using low-cost storage |
| **Data Quality** | Highly curated data that serves as the central version of the truth | Any data that may or may not be curated (ie. raw data) |
| **Users** | Business analysts | Data scientists, Data developers, and Business analysts (using curated data) |
| | | Machine Learning, Predictive |

Big Data
○○

Data warehouse
○○

Data Lake
○○○○○●○○○○

Methodology
○○○○○○○○○○

# Data Lake prototype

Big Data
○○

Data warehouse
○○

**Data Lake**
○○○○○●○○○

Methodology
○○○○○○○○○○○

# Data Lake Storage



Figure 3: Data Lake Stg

Big Data
oo

Data warehouse
oo

Data Lake
oooooooo●oo

Methodology
oooooooooo

# Stage of the Art

- Teradata
- Think Big

Big Data
○○

Data warehouse
○○

**Data Lake**
○○○○○○○●○

Methodology
○○○○○○○○○○○

## Stage of the Art

Big Data
○○

Data warehouse
○○

**Data Lake**
○○○○○○○○○●

Methodology
○○○○○○○○○○

# Stage of the Art



Figure 5: kylo_arch

Big Data
oo

Data warehouse
oo

Data Lake
oooooooooo

Methodology
●ooooooooo

# Section 4

## Methodology

Big Data
oo

Data warehouse
oo

Data Lake
ooooooooo

Methodology
o●oooooooooo

# General idea



Figure 6: DL_arch

Big Data
○○

Data warehouse
○○

Data Lake
○○○○○○○○○

Methodology
○○●○○○○○○○○

# Upload



Figure 7: Upload

Big Data
oo

Data warehouse
oo

Data Lake
ooooooooo

Methodology
oooooooooo

# Upload

```
+----------------+
| /User_1        |
|  /My_folder    |
|   .Catalog     |
|   /Data        |
|    /Patient_1  |
|     /CT_1      |
|       image_1  |
|       image_2  |
|     Cat.txt    |
|    /Video      |
|      Heo.mp4   |
+----------------+
```

Figure 8: Upload

Big Data
○○

Data warehouse
○○

Data Lake
○○○○○○○○○

Methodology
○○○○●○○○○○

# Raw



```
+--------------------+
|                    |
| /Log               |
|  Flow_log.txt      |
| /Stg               |
|  /Group_1          |
|   /User_1          |
|    /My_folder_0    |
|    /My_folder_1    |
|    /My_folder_2    |
|     /Japan_film_0  |
|                    |
|                    |
+--------------------+
```

Figure 9: Raw

Big Data
00

Data warehouse
00

Data Lake
000000000

Methodology
0000000●0000

# Transformed

```
+------------------+
| /Transformed     |
|  /Group_0        |
|   /User_0        |
|    /My_folder_0  |
|     /Ver_0       |
|      /Patient_1  |
|       /CT_1      |
|        image_1.h |
|        image_2.h |
|     /Ver_1       |
|      /Patient_1  |
|       /CT_1      |
|        image_1.k |
|        image_2.k |
|                  |
+------------------+
```
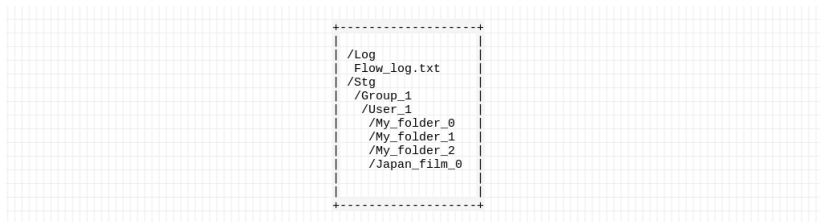
Figure 10: Transformed

Big Data
oo

Data warehouse
oo

Data Lake
ooooooooo

Methodology
ooooooo●ooo

## Trusted

```
+-------------------+
| /Transformed      |
|  /Group_0         |
|   /User_0         |
|    /My_folder_0   |
|     /Ver_0        |
|      /Patient_1   |
|       /CT_1       |
|         image_1.h |
|         image_2.h |
|         gen_img.h |
|     /Ver_1        |
|      /Patient_1   |
|       /CT_1       |
|                   |
+-------------------+
```
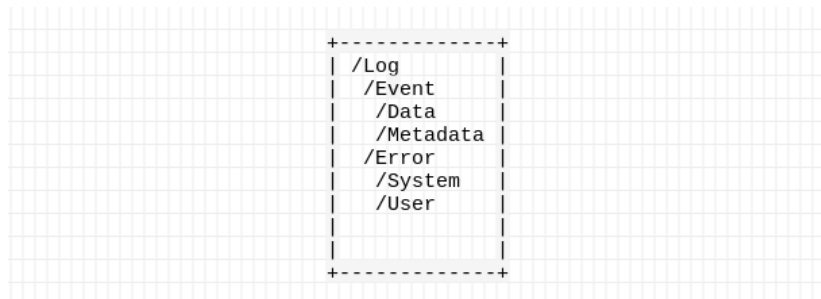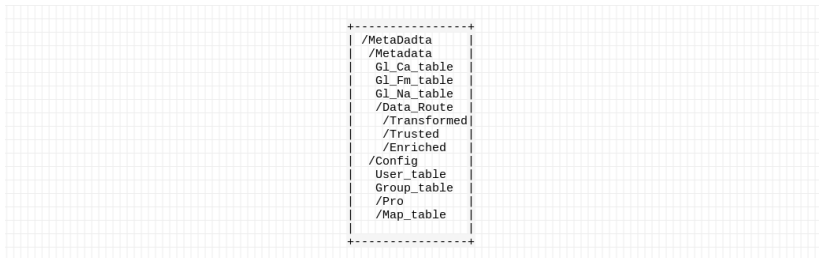
Figure 11: Trusted

# Log_struct



Figure 12: Log

Big Data
○○

Data warehouse
○○

Data Lake
○○○○○○○○○

Methodology
○○○○○○○○○●○

## Metadata

```
                    +----------------+
                    | /MetaDadta     |
                    |  /Metadata     |
                    |   Gl_Ca_table  |
                    |   Gl_Fm_table  |
                    |   Gl_Na_table  |
                    |   /Data_Route  |
                    |    /Transformed|
                    |    /Trusted    |
                    |    /Enriched   |
                    |  /Config       |
                    |   User_table   |
                    |   Group_table  |
                    |   /Pro         |
                    |    /Map_table  |
                    |                |
                    +----------------+
```

Figure 13: Meta Data

Big Data
○○

Data warehouse
○○

Data Lake
○○○○○○○○○

Methodology
○○○○○○○○○●

# Q/A ?