



Statistics Data Warehouse

## Towards Streamlining The Data Management

By:

**Nur Hurriyatul Huda binti Abdullah Sani**

27-28 August 2015

INSTITUT LATIHAN STATISTIK MALAYSIA

**DEPARTMENT OF STATISTICS MALAYSIA**



# CONTENT



- 1. Introduction**
- 2. StatsDW Objective**
- 3. Issues & Challenges**
- 4. Way Forward**



# INTRODUCTION

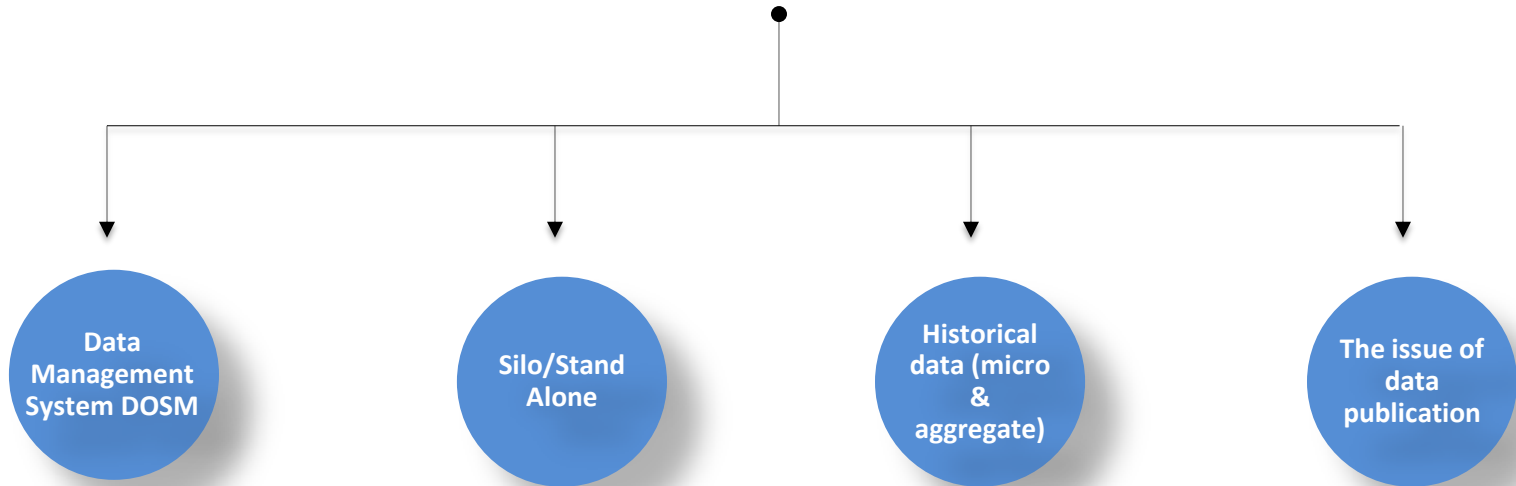


## Department of Statistics Malaysia as Producer of National Statistics Agency

### Function

- Under the "Statistics Act 1965 (Revised 1989)", the main functions of the department are: To collect and interpret statistics for the purpose of formulation or implementation of government policies in whatever fields as needed by the government or for fulfilling the requirements of trade, commerce, industry, agriculture or others
- To disseminate statistics which have been collected or interpretation based on statistics collected, not only to government agencies but also to authorities or persons where the information is useful to them

## SITUATION IN DOSM



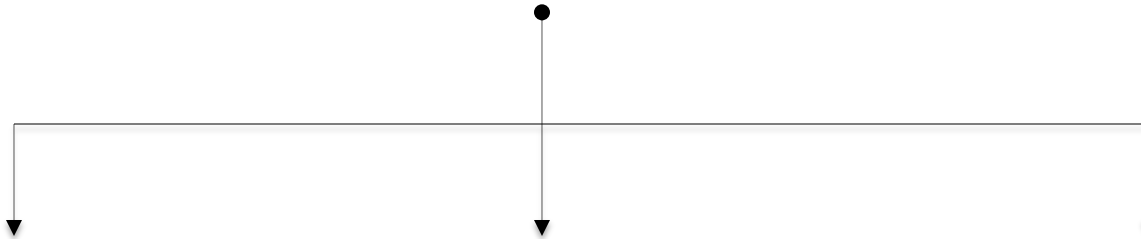
**Data Management System in DOSM built and saved by the application and Subject.**

**Data is stored in a silo**

**Historical data (micro and aggregate) are stored separately in the Subject Matter Division respectively**

**The issue of data publication:**  
a. softcopy / hardcopy missing  
b. stored in various locations.

## IMPLICATION



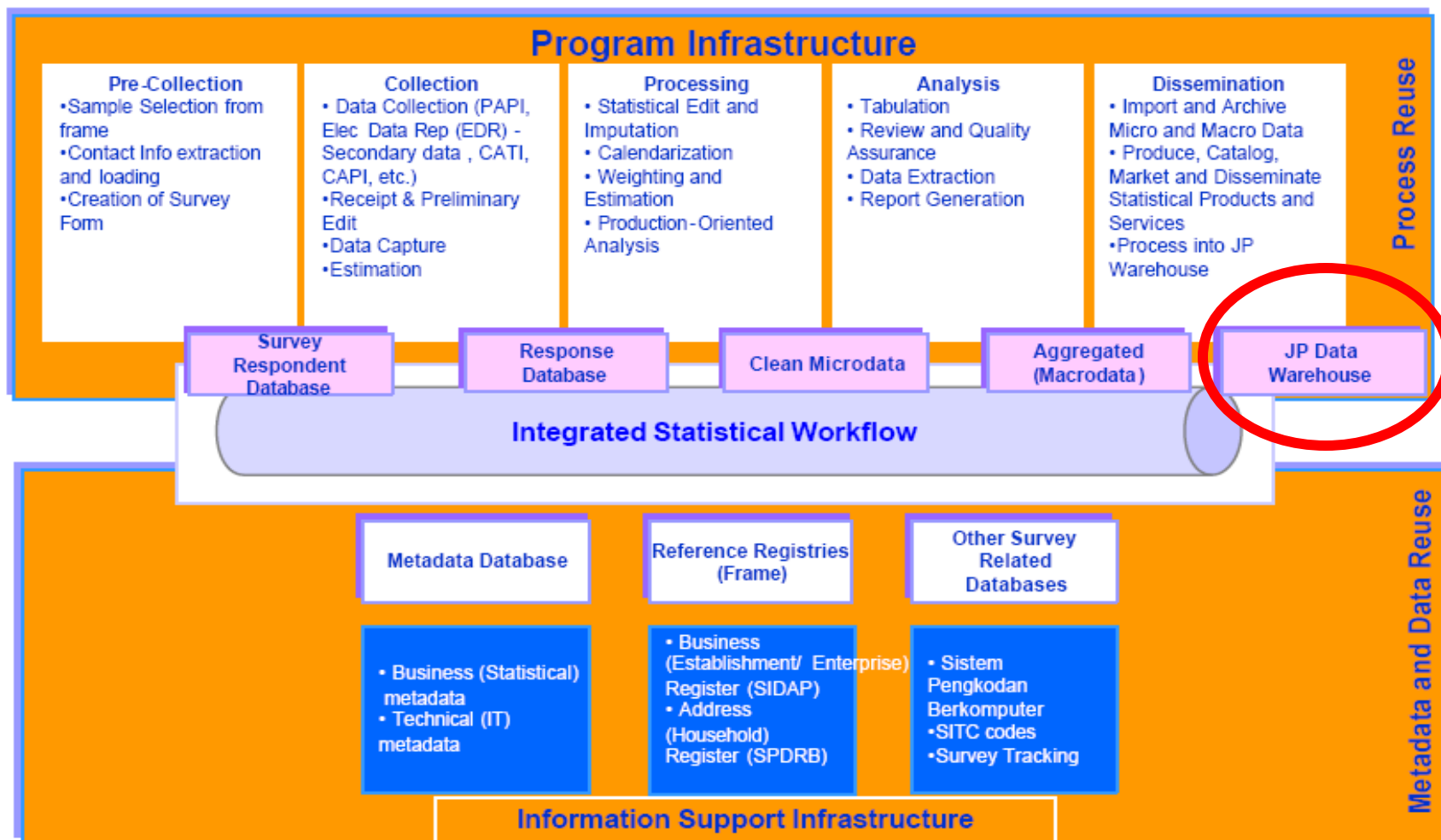
**Limited data mining / harvesting**



**Less efficient data management**

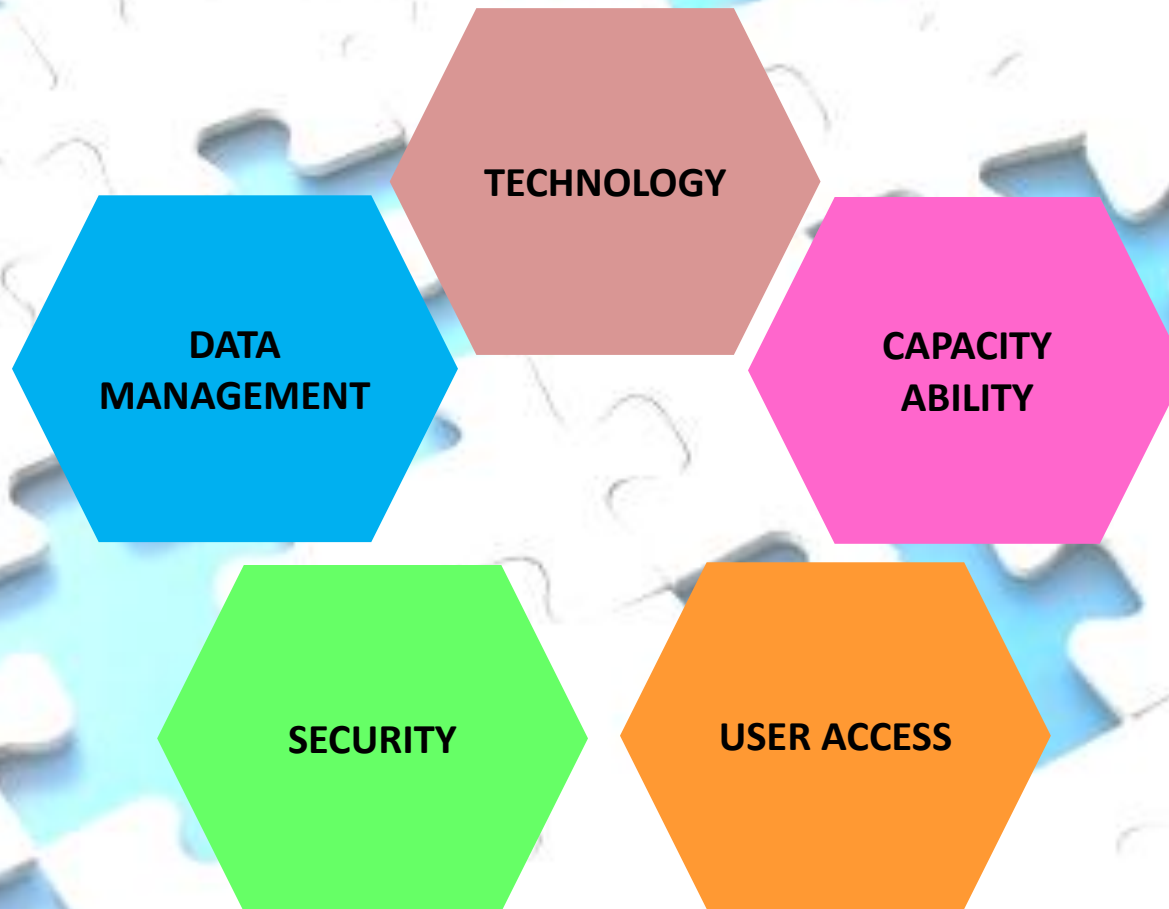


**Data archiving is not systematic**



- 1. To consolidate all Historical Micro Data and Aggregated Data in the Enterprise Data Warehouse**
- 2. To enhance data quality and consistency**
- 3. To enable fast and easy access to data stored**
- 4. To disseminate the data to wider users in a timely manner**









# ISSUE & CHALLENGES

## DATA MANAGEMENT

Data Source

Data Structure / Record Layout

System / Platform

Generation Issue

Where & Which the Final Data ?

## CAPACITY ABILITY

System Integration & Interaction

Data Preparation Duration

Data Access

Storage Capacity

## TECHNOLOGY

Data warehouse tools

Migration tools

Business Intelligence

User Acceptance

## USER ACCESS

Access Control Level

## SECURITY

ICT Security Policy

Statistic Act 1965 (Revised 1989)

Census Act 1960

Micro Data Dissemination Policy

# ISSUE & CHALLENGES

## DATA MANAGEMENT

Data Source

Data Structure / Record Layout

System / Platform

Generation Issue

Where & Which the Final Data ?

## CAPACITY ABILITY

System Integration & Interaction

Data Preparation Duration

Data Access

Storage Capacity

## TECHNOLOGY

Data warehouse tools

Migration tools

Business Intelligence

User Acceptance

## SECURITY

ICT Security Policy

Statistic Act 1965 (Revised 1989)

Census Act 1960

Micro Data Dissemination Policy

## ACCESS

Access Control

## Data Sources

### Structured Data



Micro  
Data



Aggregated  
Data



Historical  
Data

### Unstructured Data



## Data Governance, Data Integration (ETL)



Data Profiling &  
Data Quality



Data Migration,  
Data Integration  
(ETL)

## Enterprise Data Warehouse (EDW)



Data Warehouse  
Appliance



Operational Data Store

## Module

Visualisation,  
eData Bank,  
Mobile

## End Users

Internal User

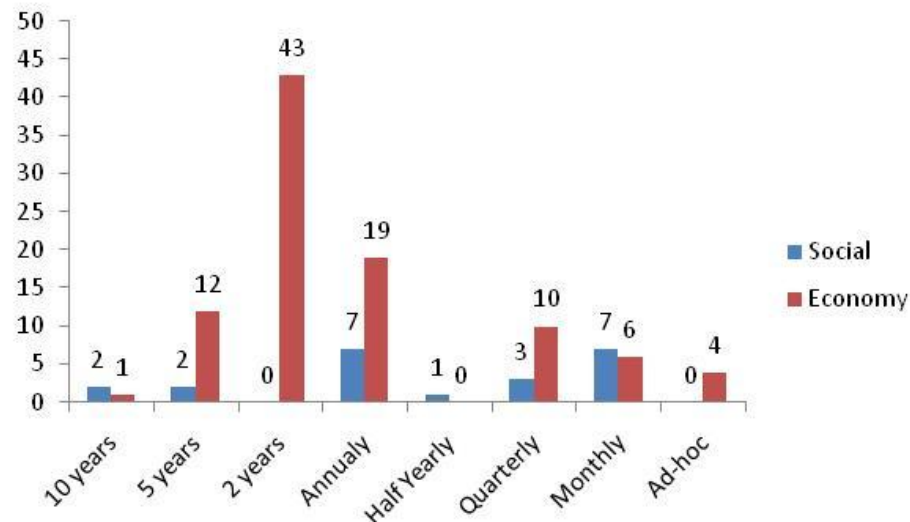
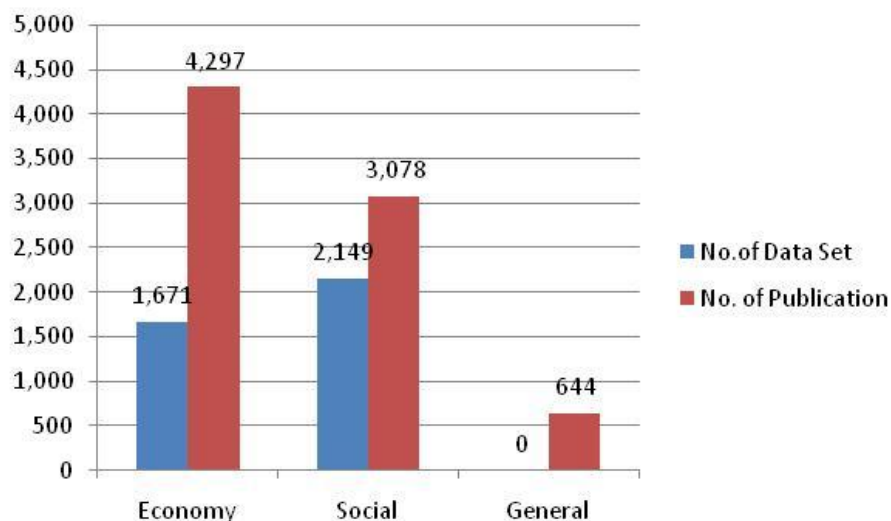
External User

## DATA MANAGEMENT

### Data Source

Source diversification and data form :

- Type of census/ survey ;
- Frequency of census/ survey ;
- Structured and unstructured data.



Data Classification	No. of Data Set	No. of Publication
Economy	1,671	4,297
Social	2,149	3,078
General	0	644
<b>Jumlah</b>	<b>3,820</b>	<b>8,019</b>

Total number of data set and publication release by DOSM as at Jun 2015



# ISSUE & CHALLENGES

## DATA MANAGEMENT

### Data Structure / Record Layout

Changes found in;

- record layout and data structure;
- variables;
- questionnaire / survey form;

Need to have data mapping for standard layout and structure.

The image displays two overlapping Excel spreadsheets. The background spreadsheet, titled '2. Mapping Variable dan Jadual Pembuatan 2000-2013\_Rev1.xlsx', shows a data table for 'Pembuatan coklat dan produk coklat' with columns for years (2011, 2010, 2009, 2008, 2007) and rows for various variables like 'Gaji & upah yang diterima' and 'Nilai harta tetap'. The foreground spreadsheet, titled '2. Mapping Variable dan Jadual Pembuatan 2000-2013\_Rev1.xlsx [Compatibility Mode]', shows a 'RECORD LAYOUT CONTENTS' table. This table maps specific content codes (e.g., 31, 32, 33) to their respective content descriptions (e.g., 'Taraf Pendidikan', 'Sijil Tertinggi') and then maps these to specific file sizes and positions across different years (2007, 2005-2006, 2004, 2001-2003, 2000 (rebase), 3-2000 (before rebase)).

NO.	CONTENT	2007	2005-2006	2004	2001-2003	2000 (rebase)	3-2000 (before rebase)
31	Taraf Pendidikan	✓	✓	✓	✓	✓	✓
32	Sijil Tertinggi	✓	✓	✓	✓	✓	✓
33	Bidang Pengajian	✓	✓	✓	✓	✓	✓
34	Institut Pengajian	✓	✓	✓	✓	✓	✓
35	No. Ahli 15 tahun ke atas	✓	✓	✓	✓	✓	✓
36	Tempat Kerja	✓	✓	✓	✓	✓	✓
37	Q1- Bekerja sekurang-kurangnya sejam dalam minggu rujukan	✓	✓	✓	✓	✓	✓
38	Q2- Ada pekerjaan untuk dikembalikan ?	✓	✓	✓	✓	✓	✓
39	Q3- Bilangan jam bekerja?	✓	✓	✓	✓	✓	✓
40	Q4- Mengapa bekerja kurang dari 30 jam?	✓	✓	✓	✓	✓	✓
41	Q5- sanggup dan boleh menerima tambahan jam bekerja?	✓	✓	✓	✓	✓	✓
42	Q6- mengapa tidak bekerja dalam minggu rujukan?	✓	✓	✓	✓	✓	✓
43	Q7- mencari kerja minggu lalu?	✓	✓	✓	✓	✓	✓
44	Q8- kenapa tidak mencari pekerjaan pada minggu lalu?	✓	✓	✓	✓	✓	✓
45	Q9- Apakah anda buat minggu lalu?	✓	✓	✓	✓	✓	✓
46	Q10- Apakah langkah-langkah yang telah	✓	✓	✓	✓	✓	✓



# ISSUE & CHALLENGES



## DATA MANAGEMENT

### System / Platform

**Different versioning found in:**

- **Data format used.**
- **Different application.**
- **Different system.**

**Need data conversion before data can be migrated into DW.**



# ISSUE & CHALLENGES

## DATA MANAGEMENT

### Generation Issue

Generation issues found in :

- Questionnaire ;
- Base year ;
- Code and Classification.

### Area

Country, State, District, *Mukim*, EB

### Social

- Ethnic
- Religion
- Nationality
- Type of houses
- Relationship to the head of HH
- Occupation status
- Field of study
- Classification of Disease & health care
- etc.

### Economy

- Malaysia Standard Classification of Occupations (MASCO)
- Malaysia Standard Industrial Classifications (MSIC).
- Malaysia Classification of Products by Activity (MCPA)
- Classifications of Functions of Government (COFOG).
- Standard International Trade Classification (SITC).
- Harmonized Commodity Description & Coding System (HS) .
- ASEAN Harmonized Tariff Nomenclature (AHTN).
- Classification of Individual Consumption by Purpose (COICOP).
- National Account Industrial Classification
- etc.



# ISSUE & CHALLENGES



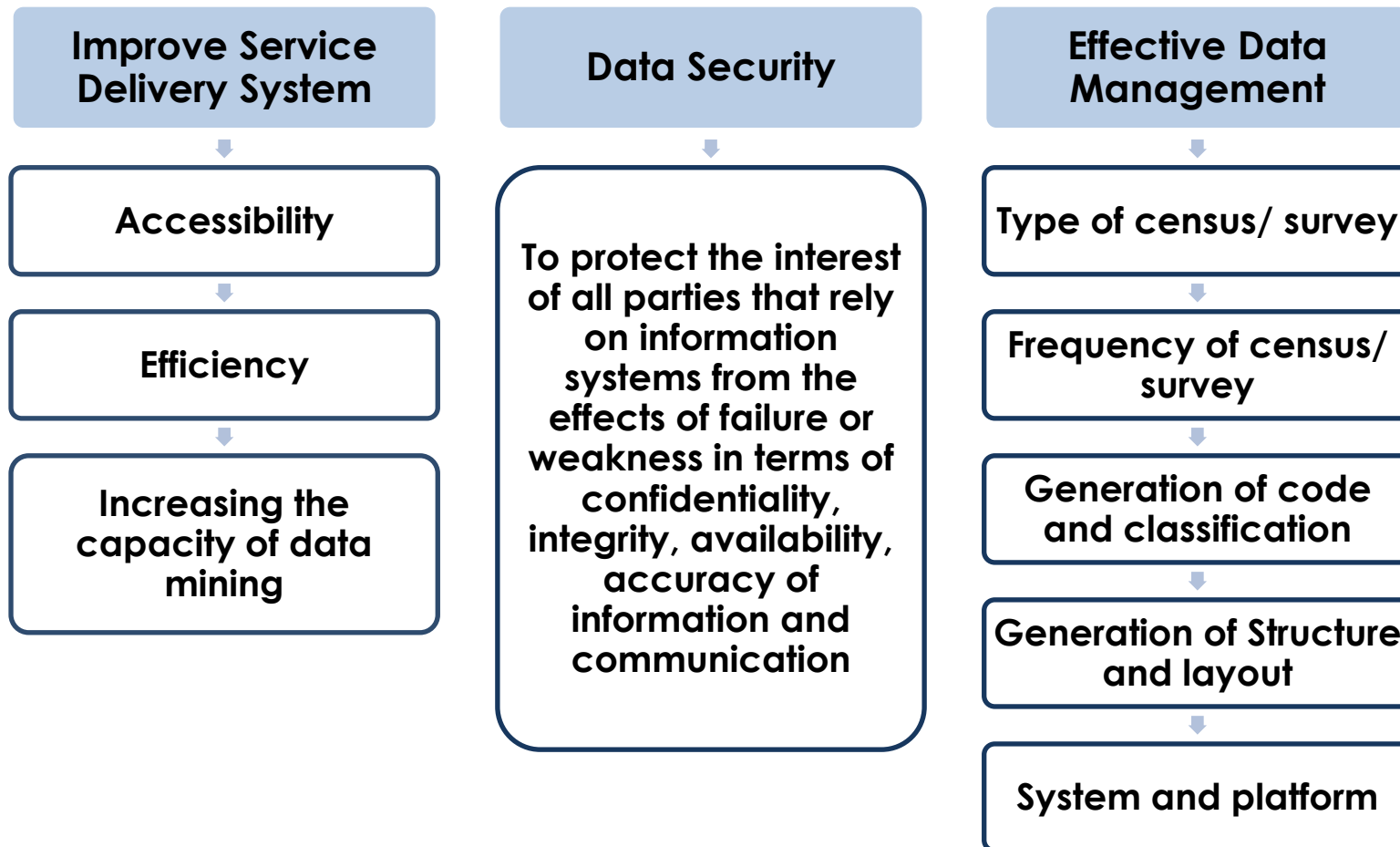
## DATA MANAGEMENT

**Where & Which the Final Data?**

- Which data final at what level?
- Who keep the final data?

To have Standard Operation Procedure :  
Where should we keep the **real** final data.







# WAY FORWARD

**Standard Operating Procedure**

**User Acceptance**

**Complete and comprehensive Data Management method and workflow**

# Thank You



**For Better Data. Better Decisions.**