# Big Data for Official Statistics

## The UN perspective

Karoly Kovacs

Data Innovation and Capacity Branch, United Nations Statistics Division

# Overview

- Big data – definition, data sources

- GWG on Big Data for Official Statistics

- Statistical Data Infrastructure

- Quality framework for Big Data

# Big data
## Definition – data sources -

# Big Data

Wikipedia:

The term has been in use since the 1990s, with some giving credit to John Mashey for popularizing the term. Big data usually includes data sets with sizes beyond the ability of commonly used software tools to capture, curate, manage, and process data within a tolerable elapsed time.

(https://en.wikipedia.org/wiki/Big_data)

United Nations Statistics Division

# Big Data

Wikipedia:

**"Big data"** is a field that treats ways to analyze, systematically extract information from, or otherwise deal with data sets that are too large or complex to be dealt with by traditional data-processing application software. Data with many cases (rows) offer greater statistical power, while data with higher complexity (more attributes or columns) may lead to a higher false discovery rate.
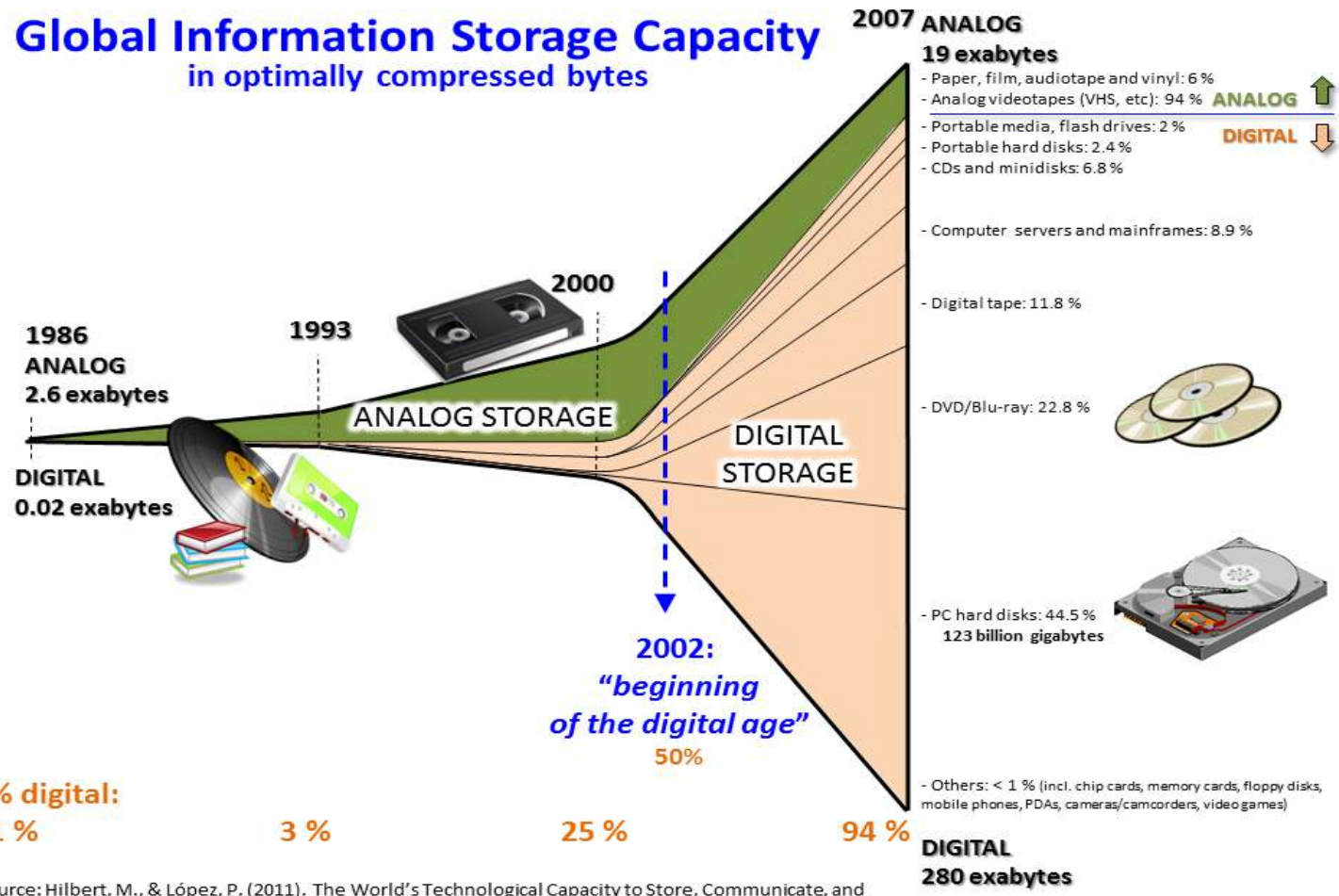
(https://en.wikipedia.org/wiki/Big_data)

# Big Data

Wikipedia:

Big data challenges include capturing data, data storage, data analysis, search, sharing, transfer, visualization, querying, updating, information privacy, and data source. Big data was originally associated with three key concepts: *volume*, *variety*, and *velocity*. Other concepts later attributed with big data are *veracity (i.e., how much noise is in the data)* and *value*.
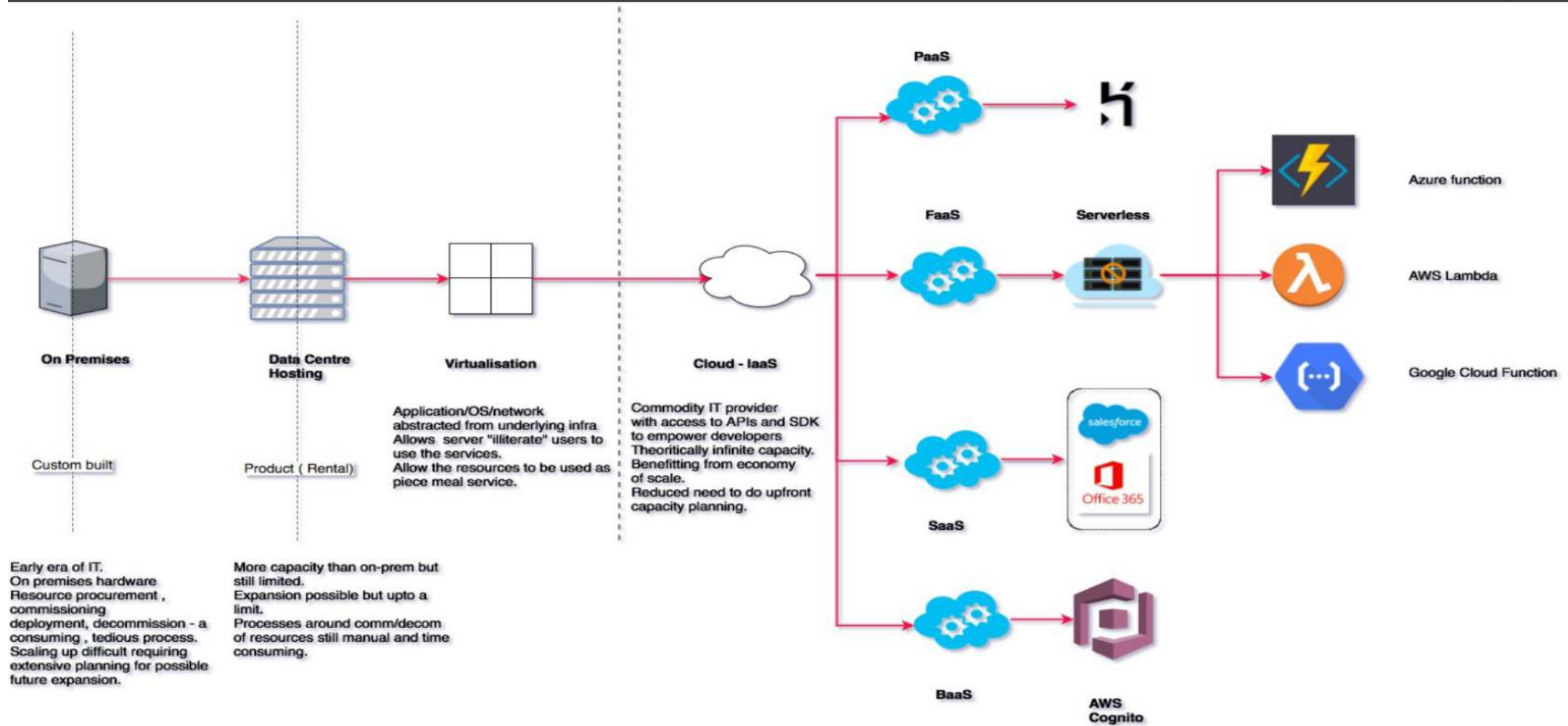
(https://en.wikipedia.org/wiki/Big_data)

# Big Data

Wikipedia:



**Global Information Storage Capacity**
in optimally compressed bytes

2007 **ANALOG**
**19 exabytes**
- Paper, film, audiotape and vinyl: 6 %
- Analog videotapes (VHS, etc): 94 % **ANALOG** ⬆
- Portable media, flash drives: 2 % **DIGITAL** ⬇
- Portable hard disks: 2.4 %
- CDs and minidisks: 6.8 %

- Computer servers and mainframes: 8.9 %

- Digital tape: 11.8 %

**2000**

**1986**
**ANALOG**
**2.6 exabytes**

**1993**

ANALOG STORAGE

**DIGITAL**
**0.02 exabytes**

- DVD/Blu-ray: 22.8 %

DIGITAL
STORAGE

- PC hard disks: 44.5 %
  123 billion gigabytes

**2002:**
*"beginning
of the digital age"*
**50%**

- Others: < 1 % (incl. chip cards, memory cards, floppy disks,
mobile phones, PDAs, cameras/camcorders, video games)

**% digital:**

**1 %**          **3 %**          **25 %**          **94 %**    **DIGITAL**
**280 exabytes**

Source: Hilbert, M., & López, P. (2011). The World's Technological Capacity to Store, Communicate, and
Compute Information. *Science*, 332(6025), 60 –65. http://www.martinhilbert.net/WorldInfoCapacity.html

PaaS

FaaS    Serverless

Azure function

AWS Lambda

Google Cloud Function

**On Premises**

Custom built

**Data Centre Hosting**

Product ( Rental)

**Virtualisation**

Application/OS/network abstracted from underlying infra
Allows server "illiterate" users to use the services.
Allow the resources to be used as piece meal service.

**Cloud - IaaS**

Commodity IT provider with access to APIs and SDK to empower developers
Theoritically infinite capacity.
Benefitting from economy of scale.
Reduced need to do upfront capacity planning.

SaaS

BaaS

AWS Cognito

Early era of IT.
On premises hardware
Resource procurement , commissioning deployment, decommission - a consuming , tedious process.
Scaling up difficult requiring extensive planning for possible future expansion.

More capacity than on-prem but still limited.
Expansion possible but upto a limit.
Processes around comm/decom of resources still manual and time consuming.

Source: https://medium.freecodecamp.org/a-brief-history-of-serverless-or-how-i-learned-to-stop-worrying-and-start-loving-the-cloud-7e2fc633310d

# What are common sources of Big Data?

o **Automatically generated data** in electronic format, **such as mobile phone data, social media data, electronic commercial transactions, sensor networks, smart meters, GPS tracking device, or satellite images**

o **High frequency, and/or fine granularity, and/or wide coverage**

Datafication

Digital footprint

Sensors

United Nations Statistics Division

As a "special case" of human mobility, tourism is a human activity that leaves multiple traces, as a digital footprint or captured by sensors

# Taxonomy of big data sources (Eurostat 2017)



| Communication systems | World Wide Web | Business process generated data | Sensors | Crowd sourcing |
|---|---|---|---|---|
| Mobile network operator data | Web activity | Flight booking systems | Traffic loops | Volunteered geographic information (OpenStreetMap) |
| Smart mobile devices data | Dynamic websites | Stores cashier data | Smart energy meters | Wikipedia contents |
| Social media posts | Static websites | Financial transactions | Vessel radio identification | Picture collections |
| | | | Satellite images | |

Potentially relevant for tourism statistics

No direct relevance for tourism statistics

# Why are Big Data important?

✓ **Big Data can keep official statistics relevant** – private sector moves fast

✓ **Big Data are part of modernization of statistical systems** – new production processes and partnerships

✓ **Big Data can help core national statistics** – for integrated economic, social and environmental policies

✓ **Big Data can help meeting the data demand of the 2030 agenda** – monitoring policies – "leave no one behind"

✓ **Big Data are needed for agile statistics** – for emergency issues

United Nations Statistics Division

# UN Global Working Group
# on Big Data for Official Statistics

United Nations Statistics Division

# Big Data for Official Statistics

**Drivers:**

o **Availability of automatically generated data** in electronic format, such as mobile phone, social media, electronic commercial transactions, sensor networks, smart meters, GPS tracking device, or satellite images

o **Higher frequency, more granularity, wider coverage, lower cost for data collection**

o **Modernisation of statistical production and services & the 2030 Agenda for sustainable development**

# United Nations Global Working Group on Big Data for Official Statistics

o Created in March 2014

o 44 members (28 countries and 16 international agencies)

o 8 Task Teams

o Coordination of the work of the TTs

o Preparation of meetings, including international conferences

o UN Global Platform

o Reporting to the United Nations Statistical Commission

# Composition of the GWG

**Countries:**
o   Australia, Bangladesh, Brazil, Cameroon, Canada, China, Colombia, Denmark, Egypt, Georgia, Germany, Indonesia, Ireland, Italy, Mexico, Morocco, Netherlands, Oman, Pakistan, Philippines, Poland, Republic of Korea, Saudi Arabia, Switzerland, UAE, UK, Tanzania, US

**Organizations:**
o   AfDB, CARICOM, Eurostat, FAO, IMF, OECD, GCC-Stat, ITU, UN GP, UNECA, UNECE, UNESCAP, UN SIAP, UNSD, UPU, WB

United Nations Statistics Division

## 👥 TASK TEAMS

Access and Partnerships

Big Data and the Sustainable Development Goals

**Mobile Phone Data**

Satellite Imagery and Geo-Spatial Data

Scanner Data

Social Media Data

Training, Skills and Capacity-building

Committee on Global Platform for Data, Services and Applications

### Mobile Phone Data

Mobile Phone Data has surfaced in recent years as one of the Big Data sources with a lot of promise. It is expected that Mobile Phone data could fill data gaps especially for developing countries given their high penetration rates. In its 2014 'Measuring the Information Society Report', ITU shows that the average mobile subscription rate is 96.4 per 100 inhabitants world-wide, with some lower averages in Asia (89.2) and Africa (69.3). Nevertheless, these numbers show how pervasive mobile phone use is. ITU elaborates that rural areas are still lacking behind urban areas, and this should be considered in studies using Mobile Phone data, but it is clear that the coverage of these data is global. Almost every person in the world lives within reach of a mobile-cellular signal.

# Handbook on the use of mobile phone data for official statistics – draft version is available at:

## Table of Contents

# UNSD project on measuring human mobility with using mobile phone data

https://unstats.un.org/bigdata/events/2019/tbilisi/default.asp



International Meeting on Measuring Human Mobility

Hosted by the National Statistics Office of Georgia (GeoStat)

Tbilisi, Georgia    27 – 29 March 2019

# Within the next 18 months, the Task Team on the use of mobile phone data would like to achieve the following:

- **Develop handbook, training materials, e-learning course and update guidelines on using mobile phone data for official statistics**
- **Document and further develop methodologies and algorithms on using mobile phone data for statistical applications** (Tourism statistics, Migration statistics, Population density statistics)
- **Develop methodologies on using mobile phone data for quality checks and getting complementary information on SDG indicators**
- **Organize project meeting on the use of mobile phone data to measure human mobility, Tbilisi, Georgia, March 2018**
- **Organize regional workshop in Indonesia, June 2019**

# Strong "outside" participation

- Positium,
- Telenor,
- IBM
- Google,
- Data Pop,
- World Pop,
- Flowminder,
- Orange,
- UNU-EHS,
- World Economic Forum,
- NASA,
- Harvard

# Statistical Data Infrastructure

# Statistical Data Infrastructure

Population and
housing census

Economic census

Agriculture census

Household surveys

Business surveys

Country

District

Community

Person

Business

## 1. Statistical data

# Statistical Data Infrastructure

Civil Register

Business Register

Road and Waterway Register

Land Register

Building Register

Country

District

Community

Person

Business

## 2. Register Data

# Statistical Data Infrastructure

**Country**

**District**

**Community**

**Person**

**Business**

Satellite imagery and aerial data

Mobile phone data

Social media or web-scraping data

Smart meter or sensor data

Credit card or scanner data

## 4. (Privately held) Big Data

# Statistical Data Infrastructure

**BigData**
UN Global Working Group

1. Statistical data

2. Register Data

3. Other Administrative Data

4. Big Data

Country

District

Community

Person

Business

**Geo-Spatial integration**

# Statistical Data Infrastructure

1. Statistical data
2. Register Data
3. Other Administrative Data
4. Big Data

Country

District

Community

Person

Business

Statistical integration: SNA and SEEA

# Integrated statistics approach

Statistical operations

Outputs / Dissemination

Inputs

**Macroeconomic accounts**

| Household and demographic statistics | Economic & environmental statistics |
|---|---|

Data integration

Data processing

Data collection

| Registers and frames | Surveys |
|---|---|

Statistical infrastructure

Standards and methods

Institutional setting

Information, Communication Technology (ICT)

Management and internal policy

**Institutional arrangements**

# Quality framework for big data

# Framework for NSO to assess quality of big data

## General approach

Quality: to be evaluated in the of intended use ('fitness for use')

Generic statistical business process model:

input                    throughput                    output

acquisition              transformation                reporting

Framework: For each phase define appropriate quality dimensions and quality indicators

# Framework for NSO to assess quality of big data

## Hyperdimensions

The concept of hyperdimension was taken from the Netherland administrative data quality framework.

- Source: Related to the type of data, the entity from which the data is obtained, and how it is administered and regulated.
- Metadata: Description of concepts, file contents, and processes.
- Data: Related to quality of the data itself.

# Framework for NSO to assess quality of big data

## Quality dimensions

- Institutional/business environment
- Privacy and security, complexity
- Completeness, usability, time factor
- Accuracy
  - selectivity
- Coherence
  - linkability
- Validity
- Accessibility, clarity, relevance

| Hyperdi-mension | Quality Dimension | Factors to consider |
|---|---|---|
| Source | Institutional Environment | Sustainability of the entity-data provider <br> Reliability status, transparency, interpretability |
| | Privacy and Security | Legislation, Data Keeper vs. Data provider <br> Restrictions, Perception |
| Metadata | Complexity | Technical constraints, Sructured or Unstructured <br> Readability, Presence of hierarchies and nesting |
| | Completeness | Metadata is available, interpretable and complete |
| | Usability | Resources required to import and analyse <br> Risk analysis |
| | Time-related | Timeliness, Periodicity, Changes through time |
| | Linkability | Presence and quality of linking variables |
| | Coherence | Use of standards |
| | Validity | Transparency of methods and processes <br> Soundness of methods and processes |

*Thank you!*
*Murakoze!*

Karoly Kovacs

United Nations Statistics Division | Department of Economic and

Social Affairs

Email: bigdata@un.org

http://unstats.un.org/unsd/bigdata