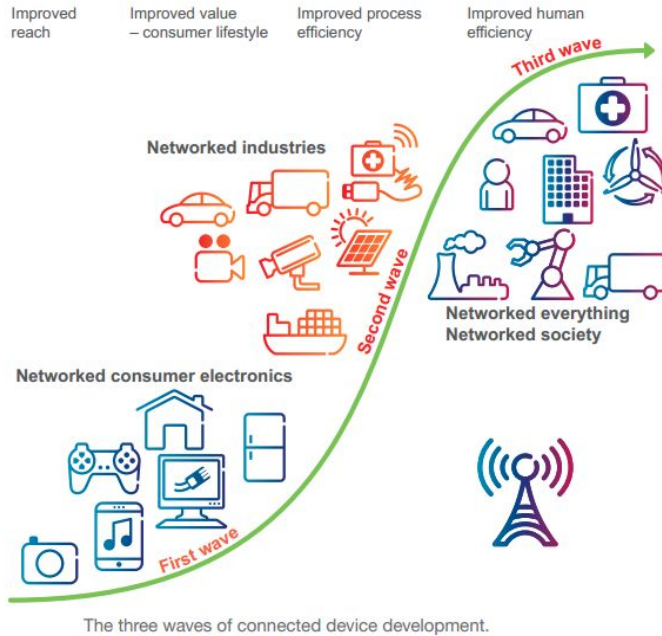ENERO 31, 2023

# FUNDAMENTOS DE
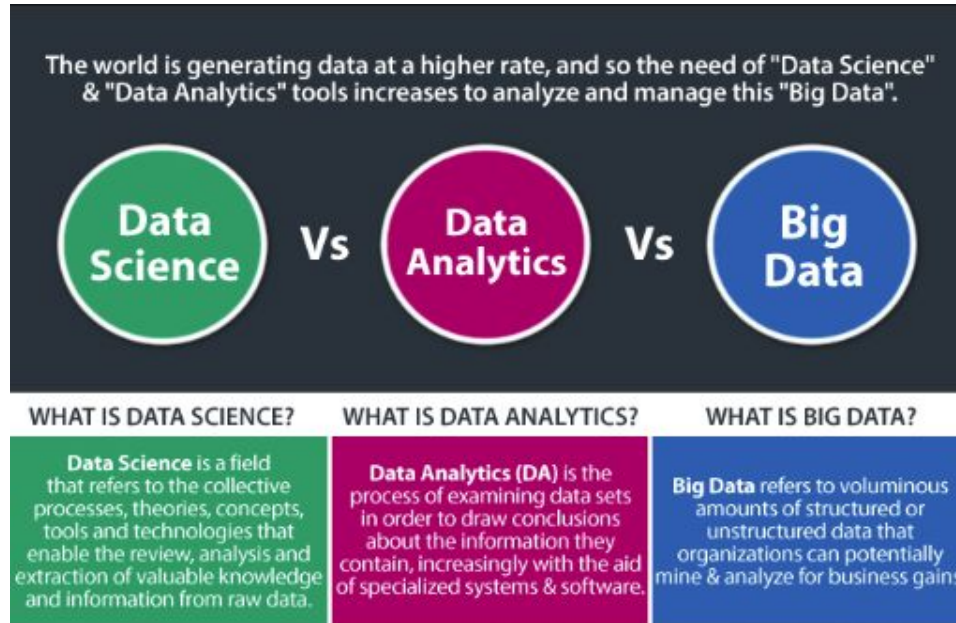# BIG DATA

UNIVERSIDAD
Panamericana

# Qué es BigData

Nuevo enfoque para entender los datos, dar valor y tomar decisiones, a partir de la descripción de datos (estructurados, no estructurados o semi estructurados) analizados desde un punto de vista no relacional (costo, tiempo, recursos), utilizando herramientas de supercómputo nativo o en la nube.

UNIVERSIDAD
Panamericana

# La relación del Big Data con otras tecnologías



Improved reach
Improved value – consumer lifestyle
Improved process efficiency
Improved human efficiency

Third wave

Networked industries

Second wave

Networked everything
Networked society

Networked consumer electronics

First wave
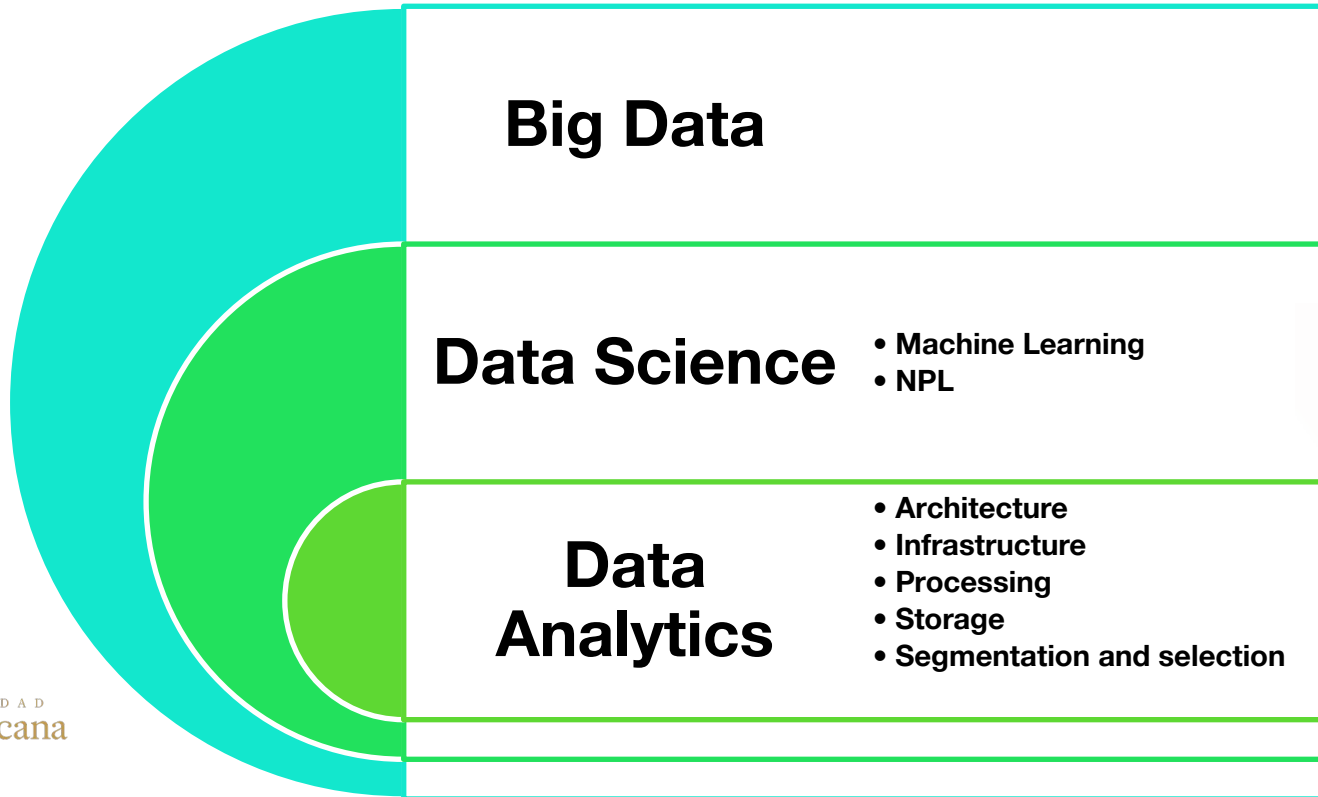
The three waves of connected device development.

- **Gran masa**
- **Diversidad**
- **Almacenamiento**
- **Seguridad**
- **Tomar decisiones y responder**

# Qué es BigData

# Qué es BigData



**Big Data**

**Data Science**
- Machine Learning
- NPL

**Data Analytics**
- Architecture
- Infrastructure
- Processing
- Storage
- Segmentation and selection

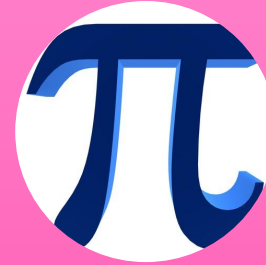# Qué es BigData, 10 V's

# Qué es Big Data

# Qué es necesario



**Infraestructura**

**Fuente(s) de datos**

**Modelo matemático**

# Big Data, generalidades

- Los motores de BD tradicionales no pueden tratar altas cantidades de datos, ni reaccionar inmediatamente a ellos. Almacena **datos no relacionales y relacionales**.
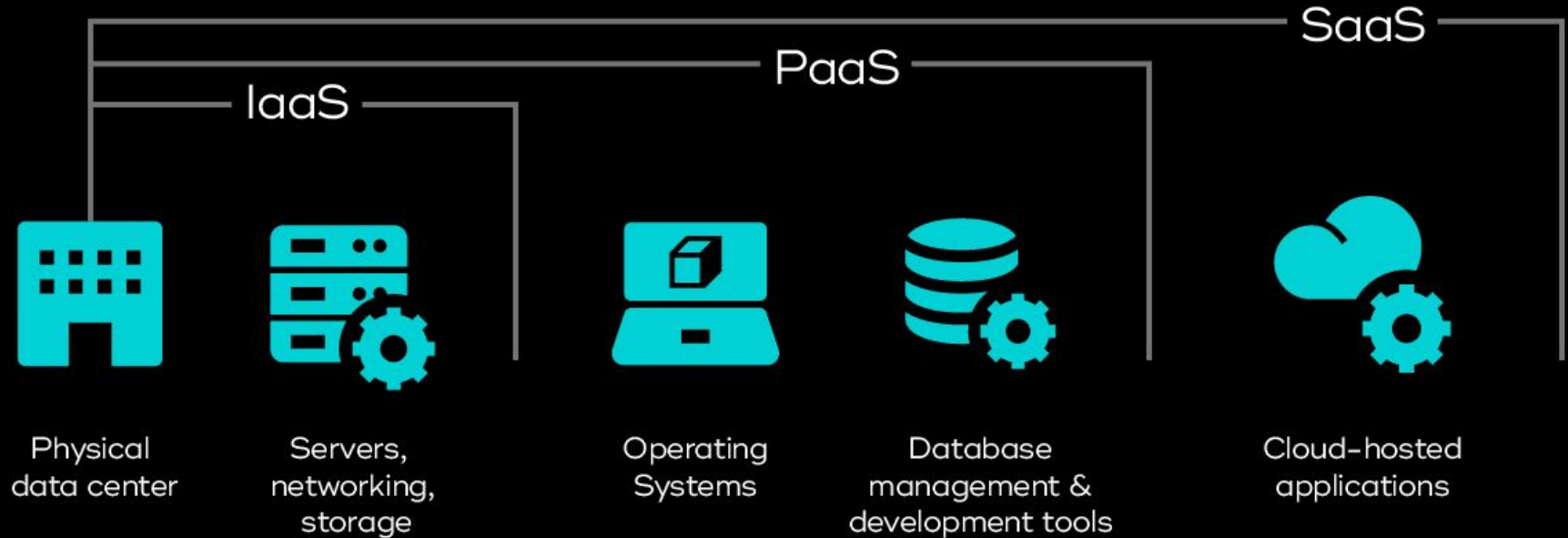
- Esquemas de consulta:

  1. **Key-value**

  2. **Columnal**

  3. **Documental**

UNIVERSIDAD
Panamericana

# Cloud Computing Overview



IaaS

PaaS

SaaS

Physical data center

Servers, networking, storage

Operating Systems

Database management & development tools

Cloud-hosted applications

IaaS, Infrastructure as a Service
PaaS, Platform as a Service
SaaS, Software as a Service

# Escenarios de cómputo en la nube

| On-site | IaaS | PaaS | SaaS |
|---------|------|------|------|
| Applications | Applications | Applications | Applications |
| Data | Data | Data | Data |
| Runtime | Runtime | Runtime | Runtime |
| Middleware | Middleware | Middleware | Middleware |
| O/S | O/S | O/S | O/S |
| Virtualization | Virtualization | Virtualization | Virtualization |
| Servers | Servers | Servers | Servers |
| Storage | Storage | Storage | Storage |
| Networking | Networking | Networking | Networking |

You manage

Service provider manages

UNIVERSIDAD Panamericana

# Tecnologías líderes Big Data engine

# Qué se puede hacer con BigData

**La nube consumía software, la nube ahora consume datos e integrando software desde donde se puede obtener información inteligente:**

**Análisis retrospectivo**
**Análisis en tiempo real**
**Análisis predictivo**
**App Inteligentes SaaS**

**Datos**
**Supervisados**
**No supervisados**

**IA**
**Machine Learning**
**Redes Neuronales**

# Queue adicional y backup

- Sharding



ORACLE SHARDING USE CASE

Oracle Sharding is implemented based on the Oracle Database partitioning feature.

Oracle Sharding is "Distribute Partitioning".

# Algo de historia

- Big Data fue precedido por Google (mapreduce), Amazon (Dynamo) y Software Open Source (Hadoop, MongoDB, Cassandra, RabbitMQ, etc)

- RDBMS tienen funciones específicas donde las relaciones entre los datos toman valor único

# Propiedades deseadas de un BD

1. Robustes y tolerancia a fallas

2. Baja latencia y actualizaciones

3. Escalabilidad

4. Generalización

5. Extensabilidad

6. Queries ad-hoc

7. Mantenimiento mínimo

8. Debugabilidad

9. Problemas de Arq. Incremental (nuevas funciones)

Universidad
Panamericana

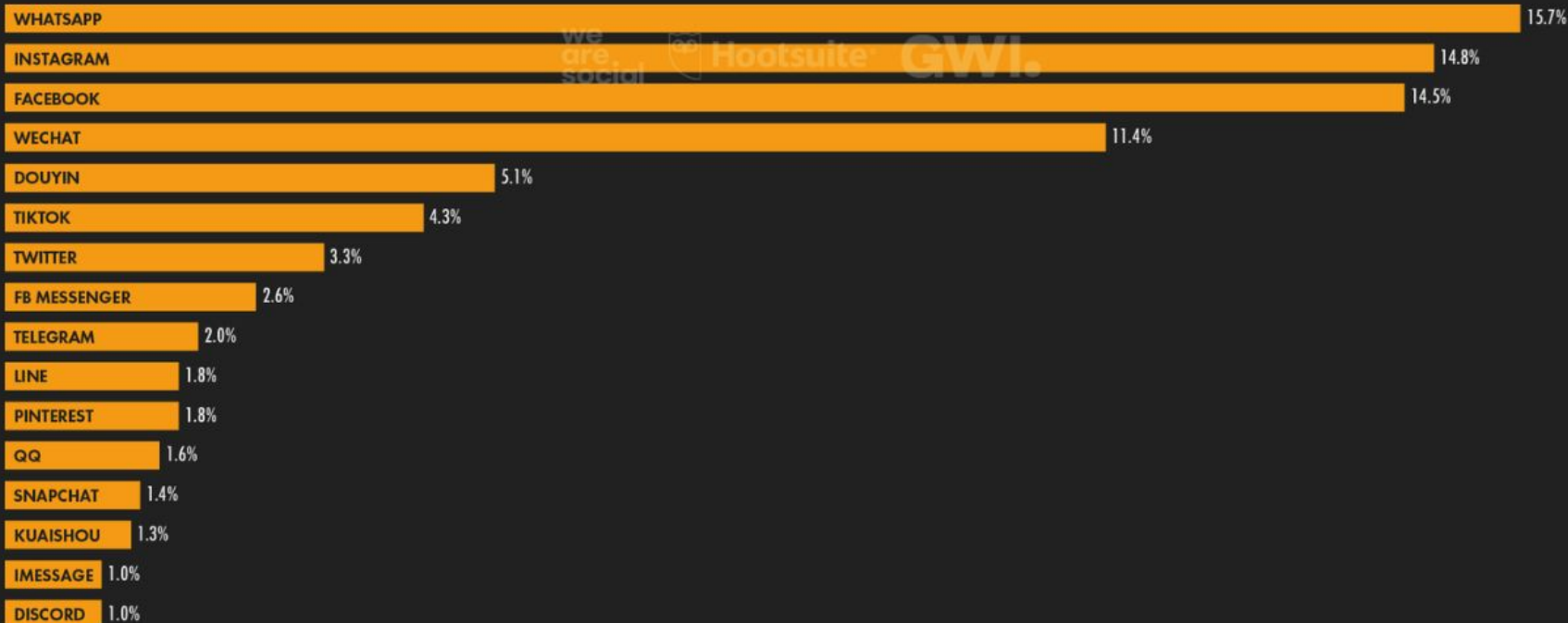Panamericana

JAN 2022

# FAVOURITE SOCIAL MEDIA PLATFORMS
PERCENTAGE OF INTERNET USERS AGED 16 TO 64 WHO SAY THAT EACH OPTION IS THEIR "FAVOURITE" SOCIAL MEDIA PLATFORM

| Platform | Percentage |
|---|---|
| WHATSAPP | 15.7% |
| INSTAGRAM | 14.8% |
| FACEBOOK | 14.5% |
| WECHAT | 11.4% |
| DOUYIN | 5.1% |
| TIKTOK | 4.3% |
| TWITTER | 3.3% |
| FB MESSENGER | 2.6% |
| TELEGRAM | 2.0% |
| LINE | 1.8% |
| PINTEREST | 1.8% |
| QQ | 1.6% |
| SNAPCHAT | 1.4% |
| KUAISHOU | 1.3% |
| IMESSAGE | 1.0% |
| DISCORD | 1.0% |

103

we are social

Hootsuite®

Panamericana

# Necesidades Globales de datos:

1. Crecimiento de logs exponencial

2. Crecimiento de demanda

3. Crecimiento de popularidad

4. Mayor inteligencia

¿Existe una correlación entre el nivel de vida de un país medido por el ingreso conoómico de su población, edad y el acceso a internet?

https://global-internet-map-2022.telegeography.com/
https://datareportal.com/reports/digital-2022-global-overview-report

# Data sources types in Big Data

1. **Structured**
2. **None structured**
3. **Semi-structured**

# Structured Data



- **Big mass**
- **Diversity**
- **Storage**
- **Security**
- **Decision Support Systems and Answers Recovery**

# None Structured Data, Organizations

- **80 -90%** world wide data

- **Images, videos, email, searches, text, pdfs, etc.**

- **V = $\frac{\Delta X}{\Delta t}$**

# Data schemas, samples

# Structured Data, **Organizations**

**Past stored in SILOS**

    **-** Unconnected stored islands

    - Hindered stored and connect silos for pattern recognition

    - Outdated and unsynchronized

**RDMS (Highly structured data) – SILOS ▯ Value**

UNIVERSIDAD
Panamericana

# Structured Data, Organizations

Using pattern recognition a organizations can use it for:

1.- Detect correlated products

2.- Estimated demand

3.- Dapture fraudulent actions

Commerce+Open Data+Analytics☐Better predictions (Business Intelligence)

# Semi- Structured Data

16 PB data per year
1 mile per driver route
    optimized, savings 50 mdd

250 millions clients, 10000
    stores
2.5 PB per hour
New products, customize
    recommendations,
    predictive support

# None-Structured Data, work path

# For all data workload

**ETL data
(Extract, Transform and Load)**

To integrate different sources and loaded it into:
Data Warehouse-Data Lake-Data Mart

# 4 key properties of a data transaction

**A**tomicity

All changes to data are performed as if they are a single operation

**C**onsistency

Data is in a consistent state when a transaction starts and when it ends.

**I**solation

The intermediate state of a transaction is invisible to other transactions. As a result, transactions that run concurrently appear to be serialized.

**D**urability

After a transaction successfully completes, changes to data persist and are not undone, even in the event of a system failure.

# None Structured Data, NoSQL

1. Graph Data Base.– used to find connections between data sets (Neo4j)
2. Key Value Pairs.- access and process data with key value pairs (Cassandra)
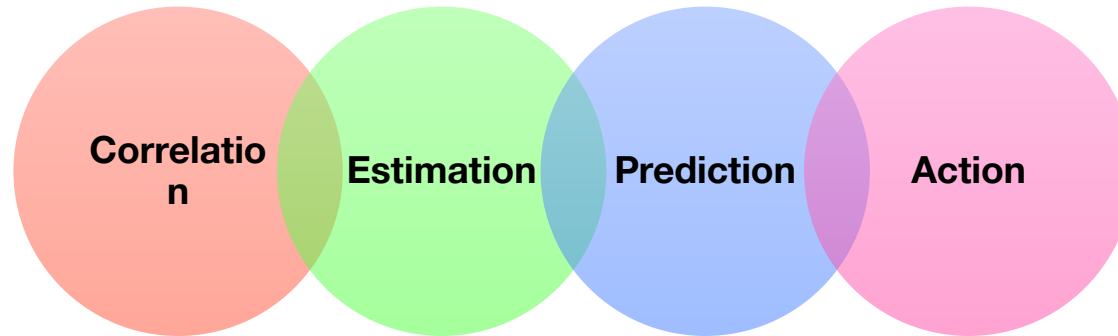
Layers for Value Big Data

Retrival and Storage

Pre-processing

Analysis

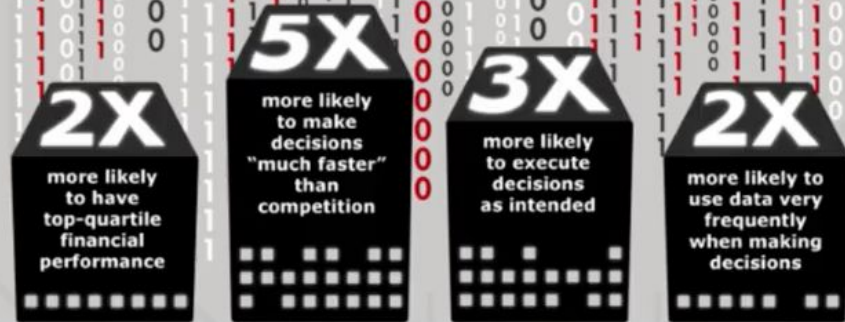12 TB/day representative data to sentiment analysis for a product or service (crisis mappers)

# None Structured Data, NoSQL
# Flow decision making

**Correlation**  **Estimation**  **Prediction**  **Action**

# Strategical Big Data growth, new economical system approach



THE COMPANIES THAT USE ANALYTICS BEST ARE...

2X more likely to have top-quartile financial performance

5X more likely to make decisions "much faster" than competition

3X more likely to execute decisions as intended

2X more likely to use data very frequently when making decisions

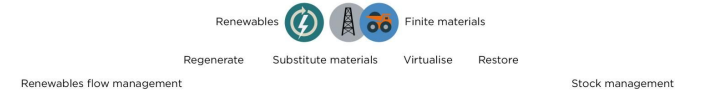UNIVERSIDAD Panamericana



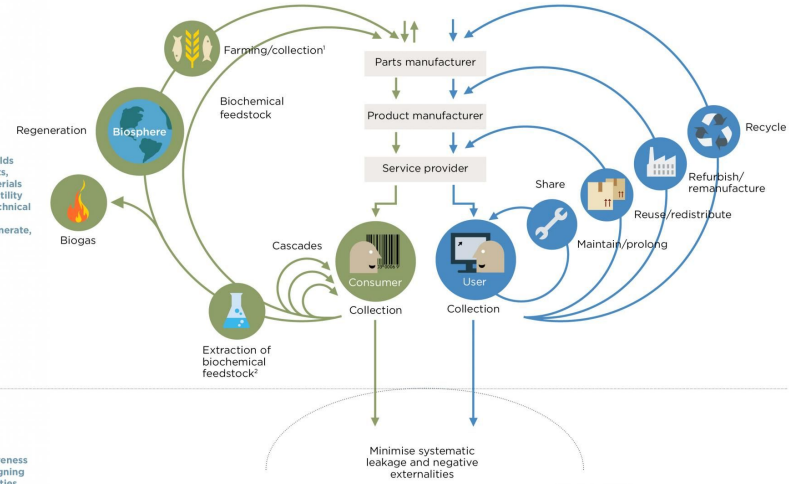OUTLINE OF A CIRCULAR ECONOMY

PRINCIPLE 1
Preserve and enhance natural capital by controlling finite stocks and balancing renewable resource flows ReSOLVE levers: regenerate, virtualise, exchange

Renewables        Finite materials

Regenerate    Substitute materials    Virtualise    Restore

Renewables flow management                                    Stock management

Farming/collection¹

Parts manufacturer

Biochemical feedstock

PRINCIPLE 2
Optimise resource yields by circulating products, components and materials in use at the highest utility at all times in both technical and biological cycles ReSOLVE levers: regenerate, share, optimise, loop

Regeneration    Biosphere

Product manufacturer

Service provider

Share

Recycle

Refurbish/ remanufacture

Reuse/redistribute

Biogas

Cascades

Maintain/prolong

Consumer       User

Collection       Collection

Extraction of biochemical feedstock²

PRINCIPLE 3
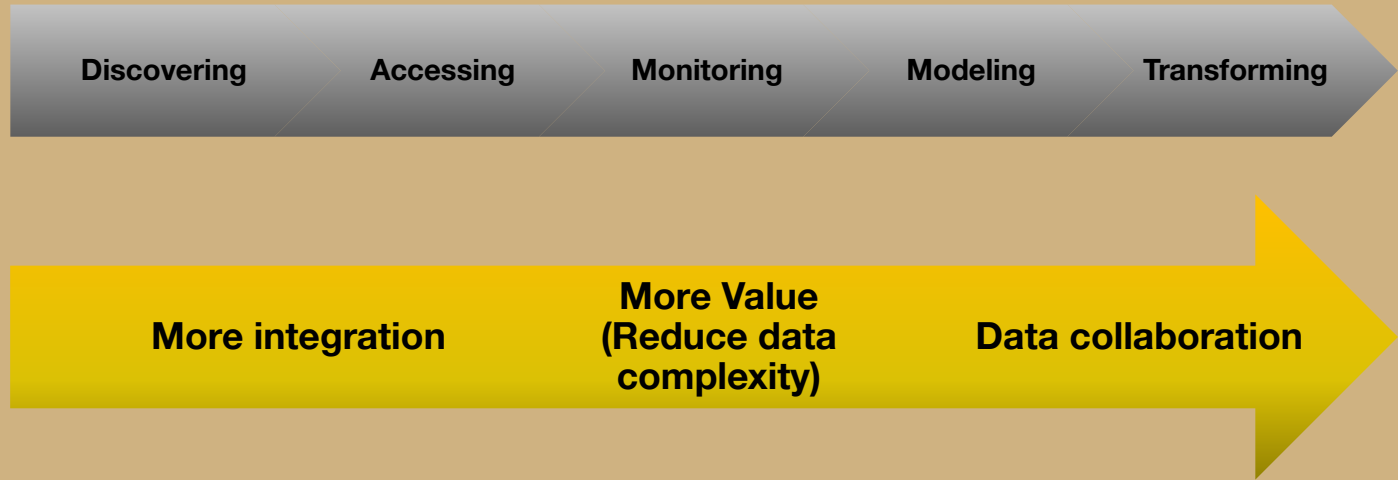Foster system effectiveness by revealing and designing out negative externalities All ReSOLVE levers

Minimise systematic leakage and negative externalities

1. Hunting and fishing
2. Can take both post-harvest and post-consumer waste as an input

Source: Ellen MacArthur Foundation, SUN, and McKinsey Center for Business and Environment; Drawing from Braungart & McDonough, Cradle to Cradle (C2C).

# Data integration Process

# Activity 2:

1. Read the appropiate article according the assignated team.
2. Answer the questiones asignated to the article.
3. Every team must prepare and present a 5 minute resume

**Team 1**
- [Link](Link)

**Team 2**
- [Link](Link)

**Team 3**
- [Link](Link)

**Team 4**
- [Link](Link)

**Team 5**
- [Link](Link)

# Questions Team 1:

1. Does the Reels are changing the consumer experiences and it reflect a new data Vs value?
2. Wich one's and how?
3. What trends (3) do you identify after Covid-19 over data user production?
4. Are those data commonly Structured, NoStructured or SemiStructured? Explain, why?

# Questions Team 2:

1. What are the relations between Data Analytics and Business Intelligence (3)?
2. What are the relations between Data Science and Business Intelligence (3)?
3. How do you explain this for a new user or Company who wants to be involved in BI?
4. How could they start with it?

# Questions Team 3:

1.  For the 12 data keypoint, how many of them are affeted by AI?
2.  In conclusion, the D&A should be focused on people or metaverse? Explain, why?
3.  How do you resume this article in one paragraph?
4.  Does geopolitics are changing data production and consumtion?

# Questions Team 4:

1.  How do you represent graphically the relation between Data Science-ML-AI?
2.  How do you explain the role of Big Data in Data Science, ML and AI?
3.  Describe three agrees, and three disagrees about the article's content
4.  According the study case: "Self-driving car". Do you agree with the ML, AI and Data Science descriptions? Explain, What is the Big Data role in this case according 10 V's?

# Questions Team 5:

1. Does Data Scientist and Data Engineer would do the same job? Explain, why?
2. Define five key skills that make difference between Data Scient and Data Engineer
3. Does DASCA is leading changing the employment bases for Data Scientist and Data Engineers?
4. How has been involved the Big Data in this new industrial revolution called 4.0?