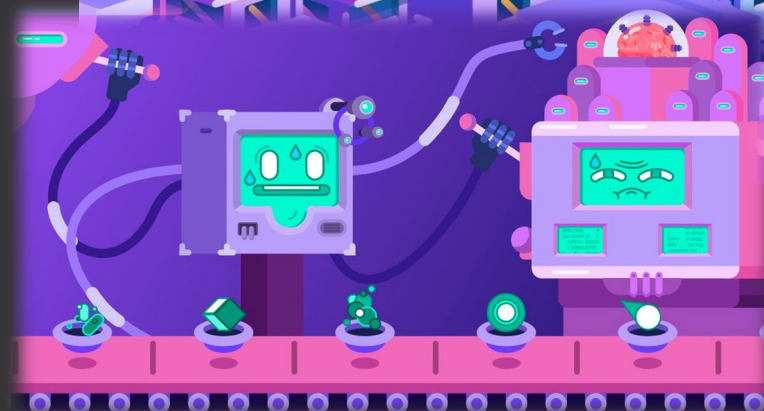




# Metody i narzędzia *Big Data*

Krajobraz metod i narzędzi *Big Data*

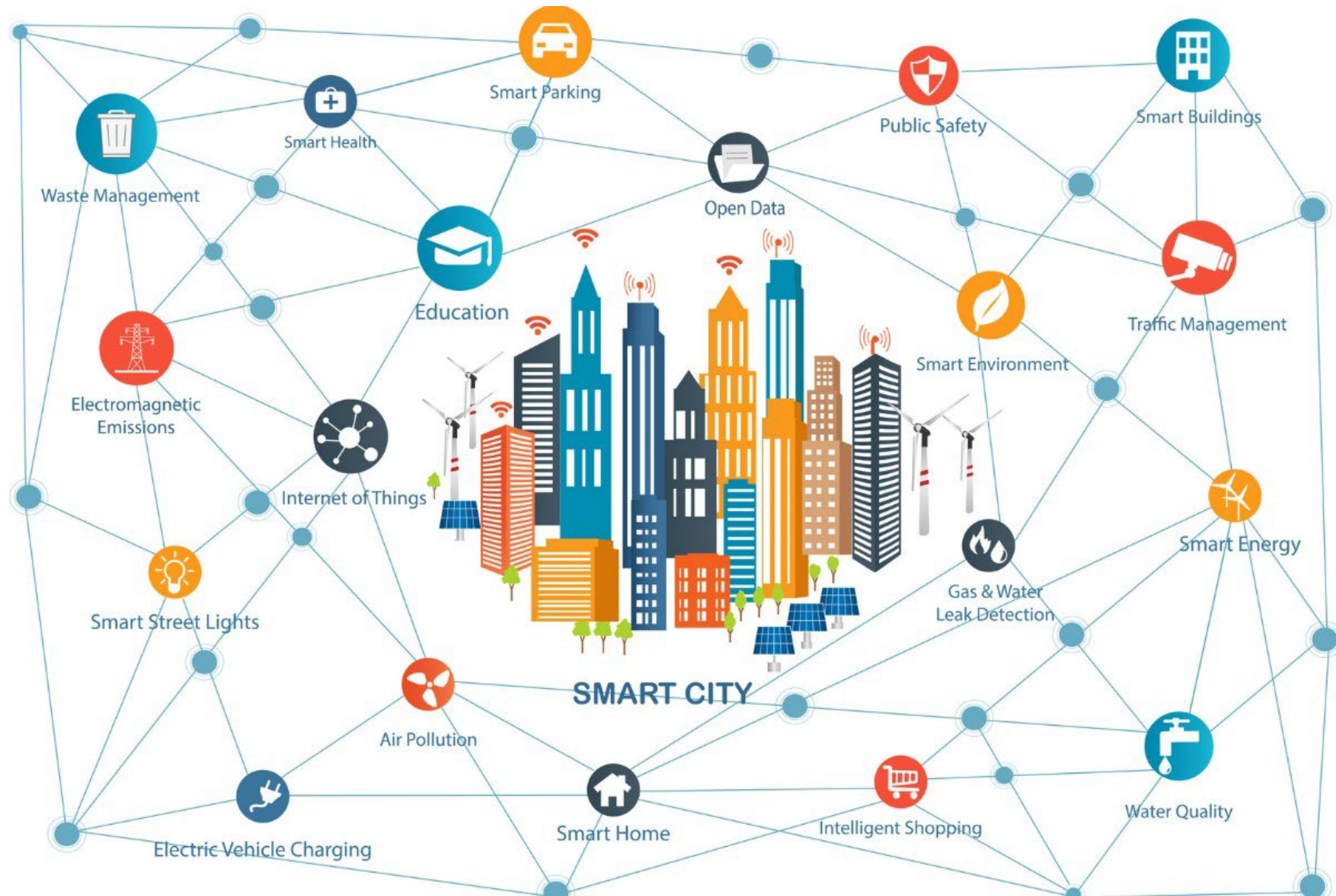




1956: transport dysku twardego o pojemności 5Mb



1963: gorąca linia nawigacyjna, powstała  
długo przed usługą Google Maps



# Problem

Dane są zbyt **duże**, aby dało się je **przechowywać** w klasycznych bazach danych i **przetwarzać** klasycznymi metodami



## Archiwa

Skany, oświadczenia, recepty, ...



## Dokumenty

PDF, DOC, XLSX, RTF, CSV, HTML, JSON, ...



## Aplikacje

E-recepty, Bankowość, gry, ...



## Media

Zdjęcia, filmy, audio, zdalne spotkania, ...



## Sieci społecznościowe

Facebook, Twitter, LinkedIn, Spotify ...



## Internet rzeczy

Klimatyzatory, TV, inteligentny dom, inteligentny magazyn, ...



## Dane sensoryczne, pomiary

Systemy monitoringu (kamery, czujniki), systemy czasu rzeczywistego, lokalizacja, GPS, Samoloty, Roboty przemysłowe,...



## Rejestry aplikacji

Serwery, logi, „kliknięcia” na witrynie,...



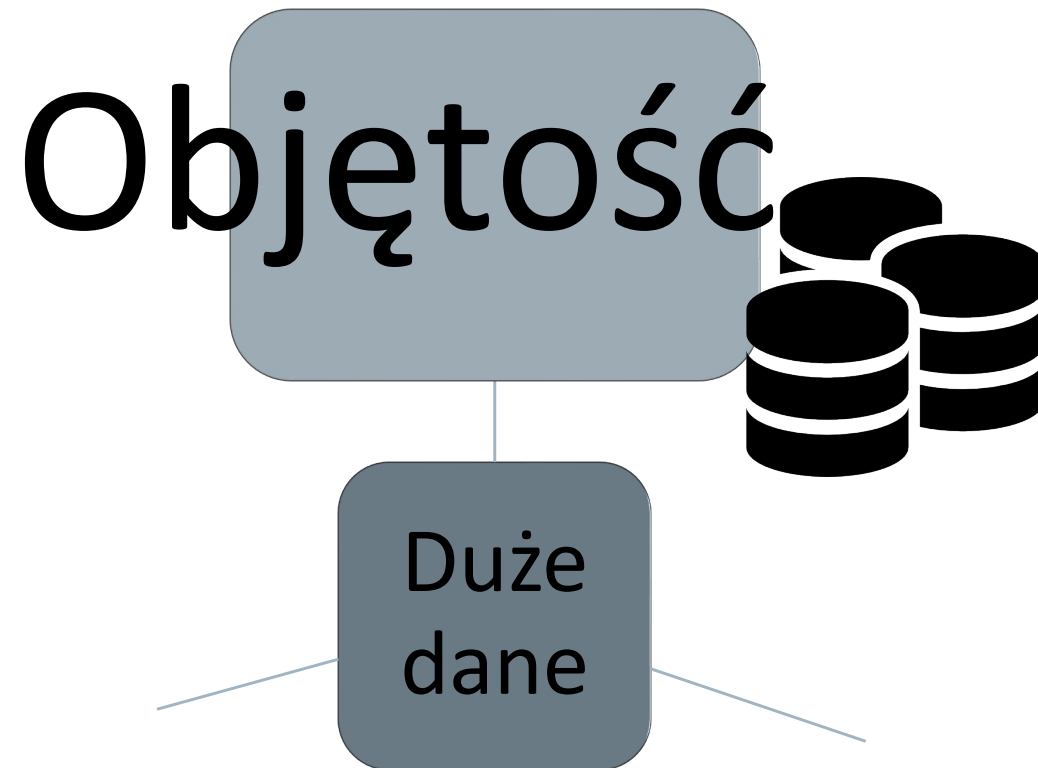
## Repozytoria

NoSQL, Hadoop, RDBMS ,...

**INTERNET RZECZY**

Co to znaczy, że dane są  
**duże?**

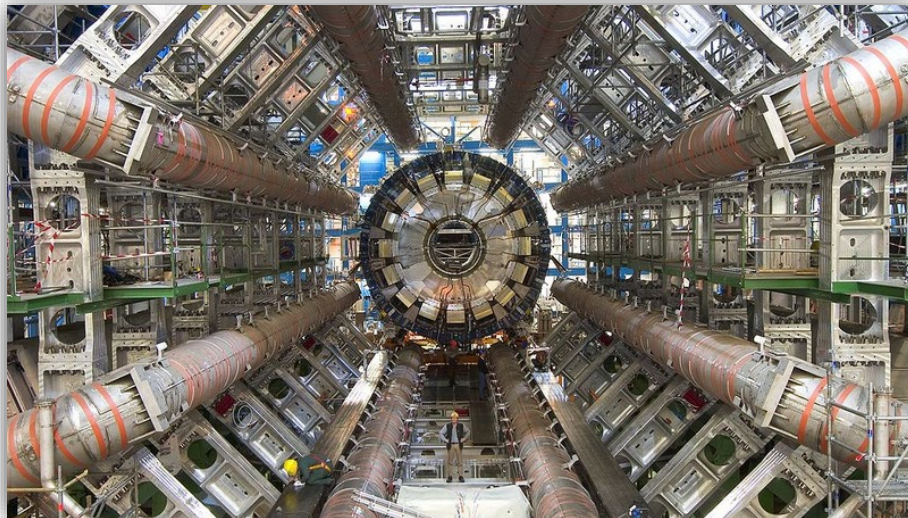
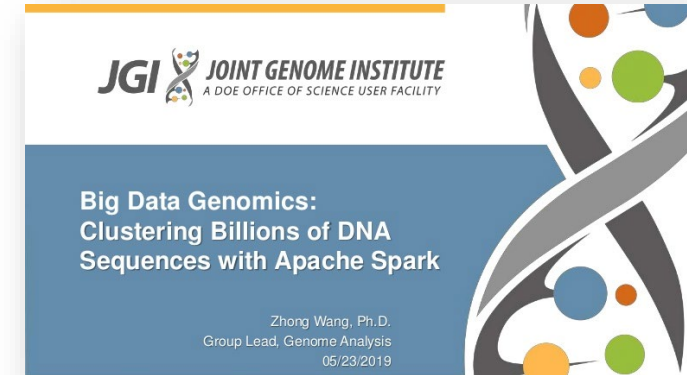
# Co to znaczy, że dane są **duże**?



# Duża objętość

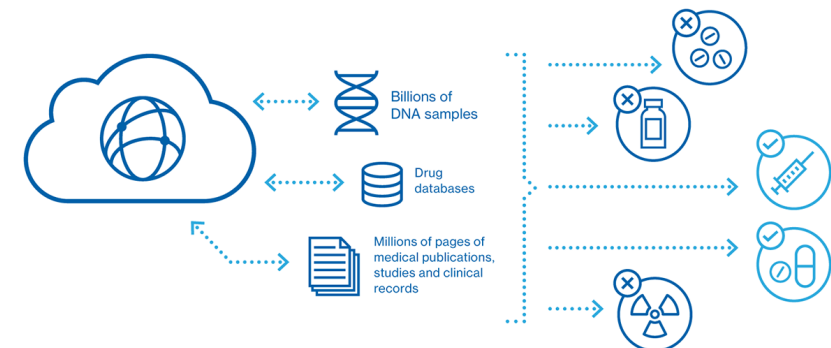
## Przykład

- Objętość liczona w petabajtach
- Wyzwania: przechowywanie i przetwarzanie
  - nie da się ich załadować na pojedynczą maszynę i przetwarzać



### In the cloud – Better therapies for patients

New cloud-based cognitive systems compare the patient's data with big data sets. This helps in the search for suitable treatment options.



# Duża objętość

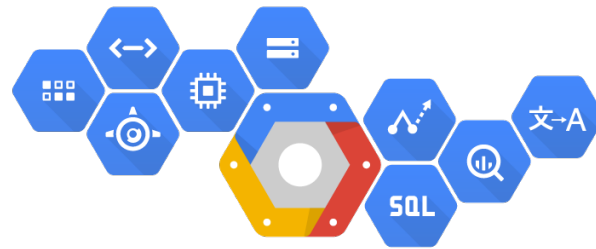
*Przykład*

Czy mieszczą się w pamięci?	Czy mieszczą się na dysku?	Jakie to dane?
tak	tak	małe
nie	tak	średnie
nie	nie	<b><u>duże</u></b>



# Duża objętość

## Technologie



Google Cloud Platform



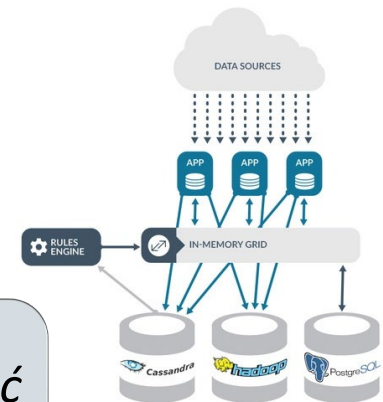
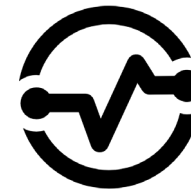
# Co to znaczy, że dane są **duże**?

# Objętość



Duże dane

Szybkość i zmienność  
(strumienie)

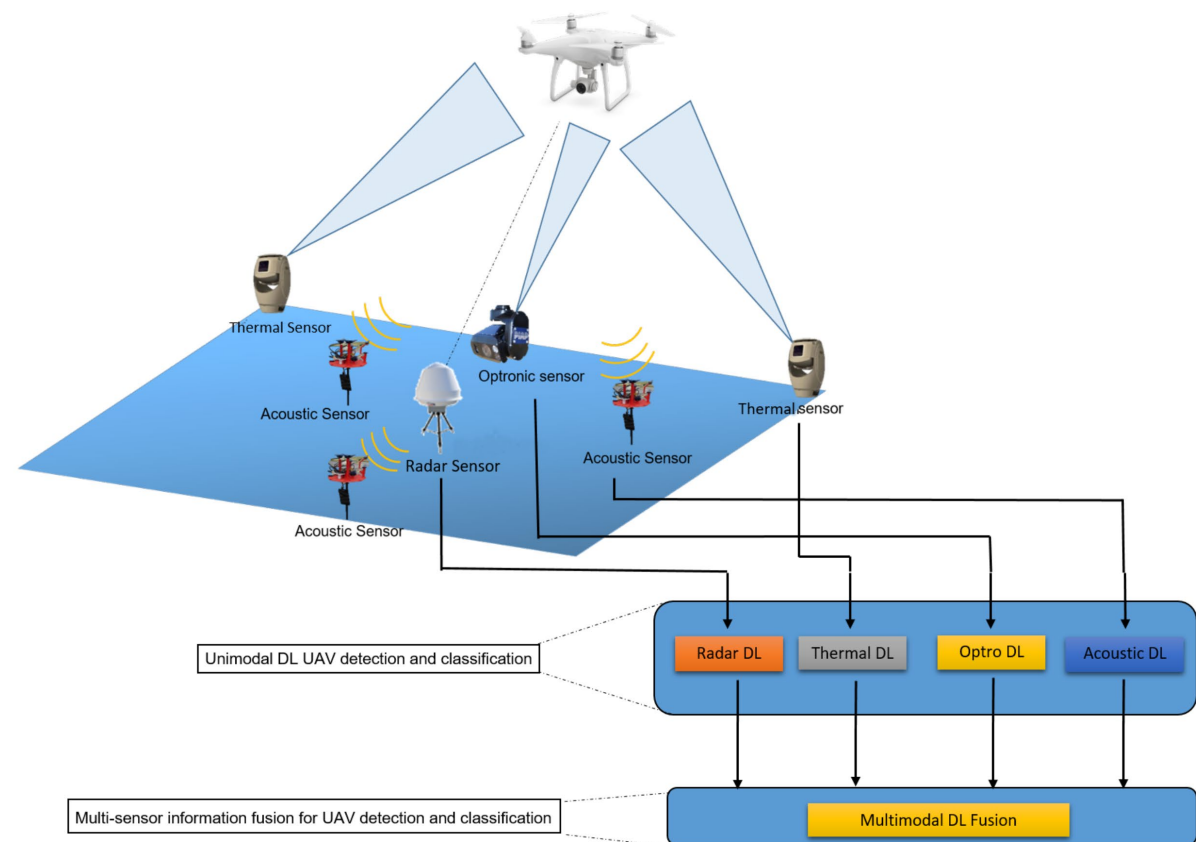


Szybko  
narastają

# Duża szybkość i zmienność (dane strumieniowe)

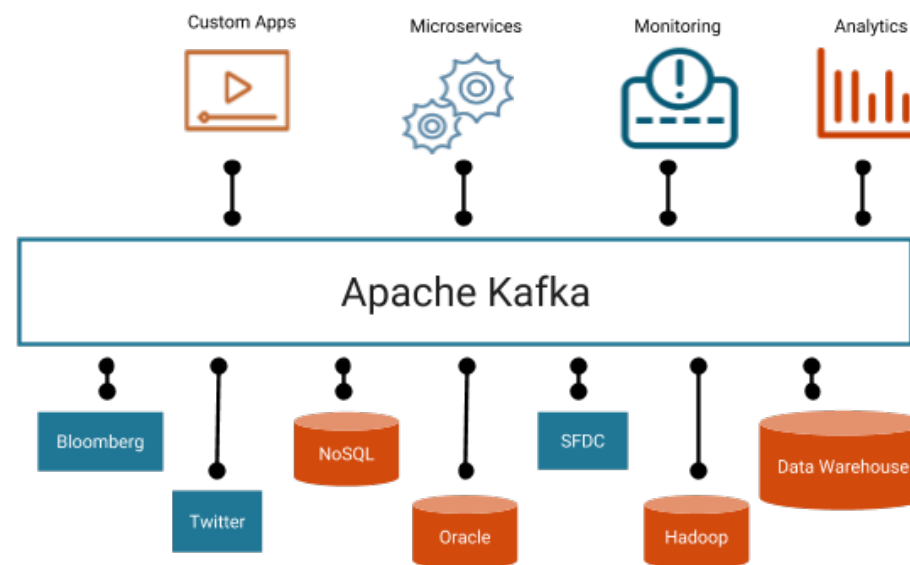
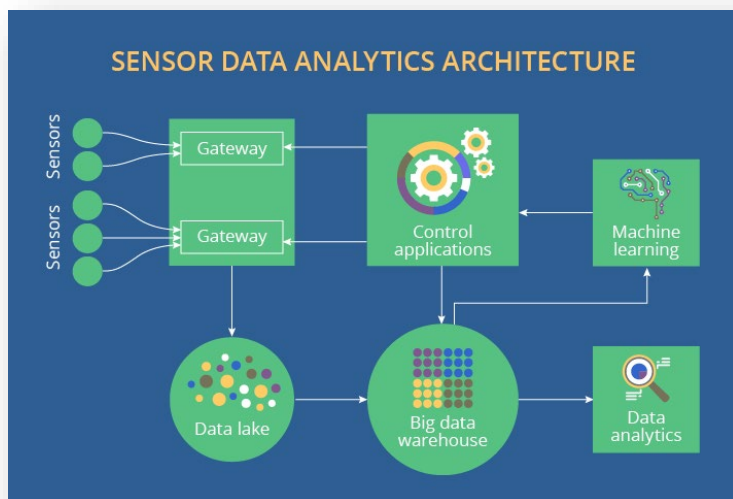
## Przykład

- Wyzwania:
  - analiza na bieżąco: **przetwarzanie w czasie rzeczywistym**, niestacjonarność
  - integracja danych z różnych źródeł
  - integracja z historycznymi danymi
  - niekompletność, zdublowane dane

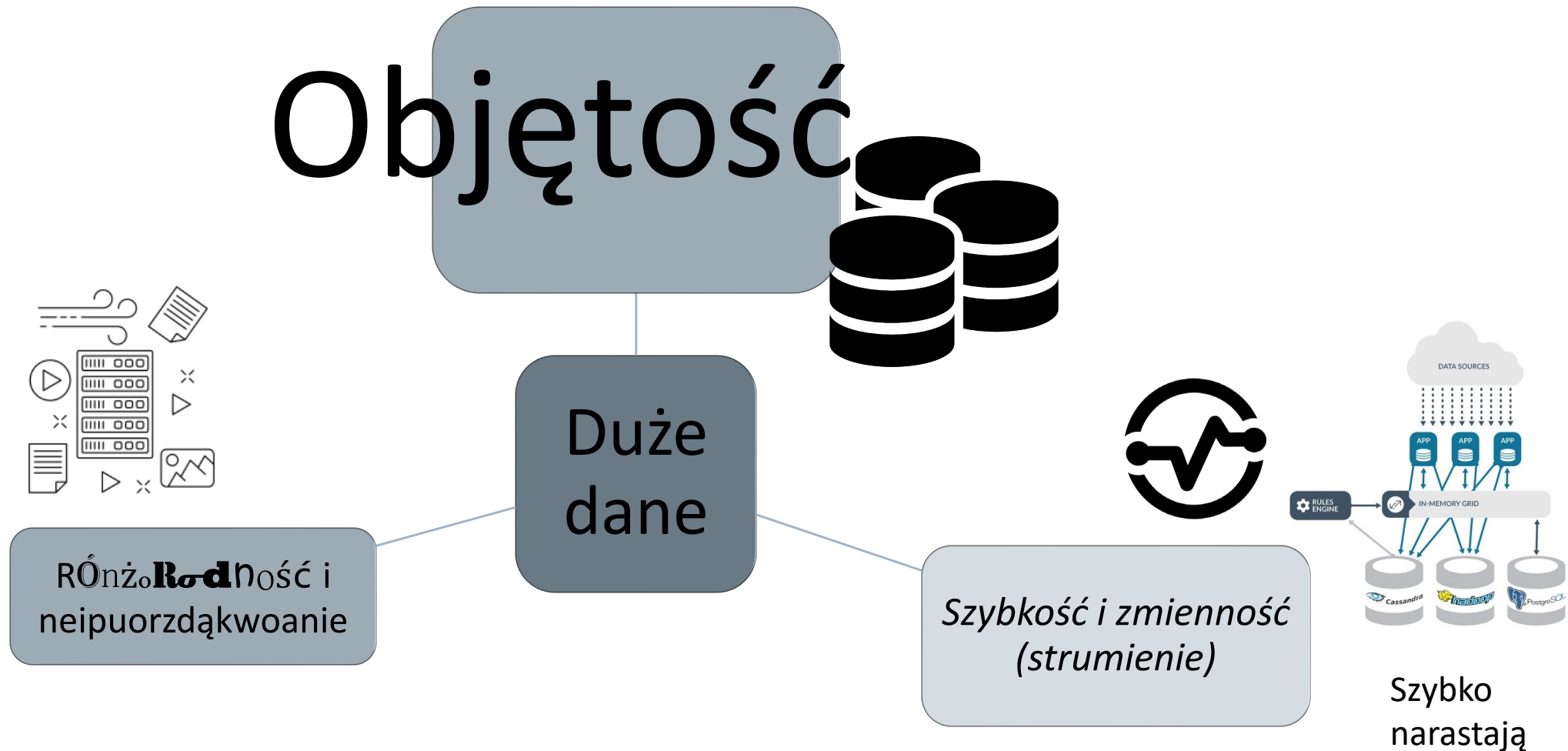


# Duża szybkość i zmienność (dane strumieniowe)

## Technologie



# Co to znaczy, że dane są **duże**?



# Duży bałagan (ang. *complexity*)

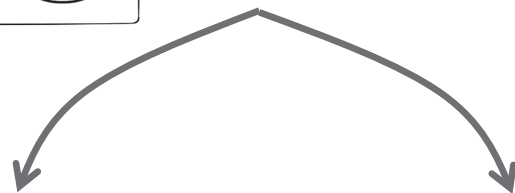
## Przykłady



- **Nieustrukturyzowane**, słabo udokumentowane, wielopoziomowe, pochodzące z **różnych źródeł**
- Wyzwania: przekształcanie danych do określonego formatu

# Duży bałagan (ang. *complexity*)

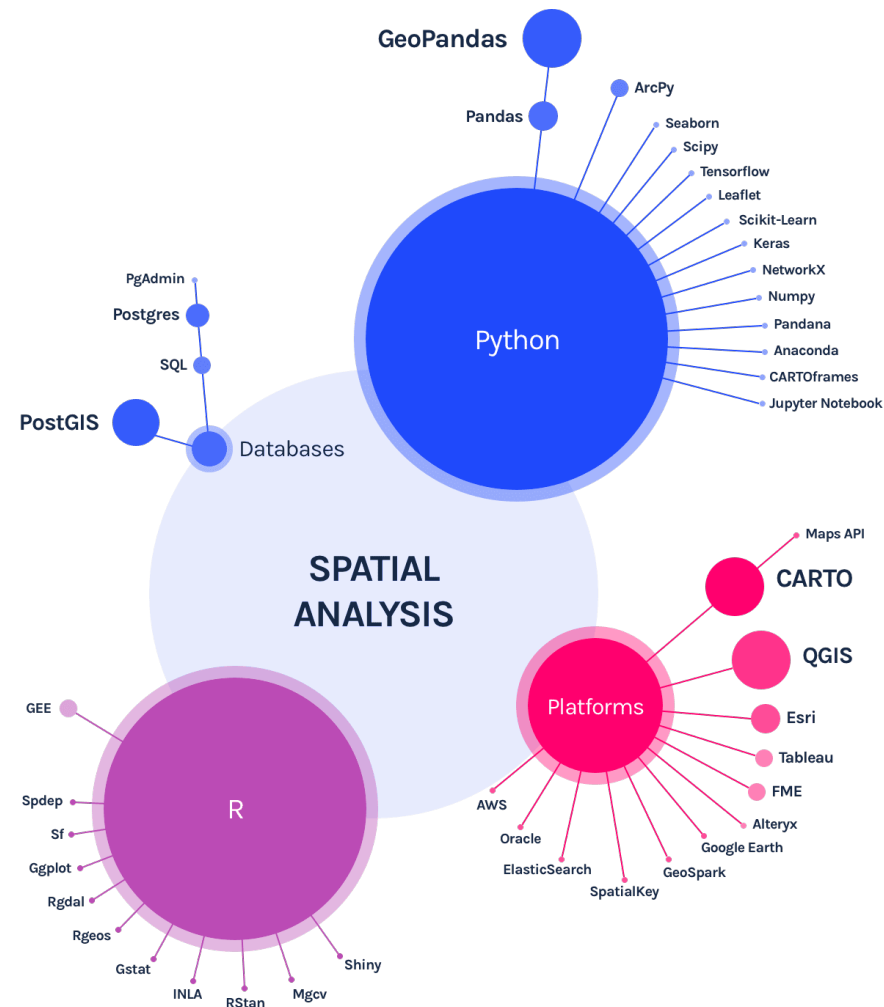
## Technologie



**DASK**  
Distributed Summit 2021

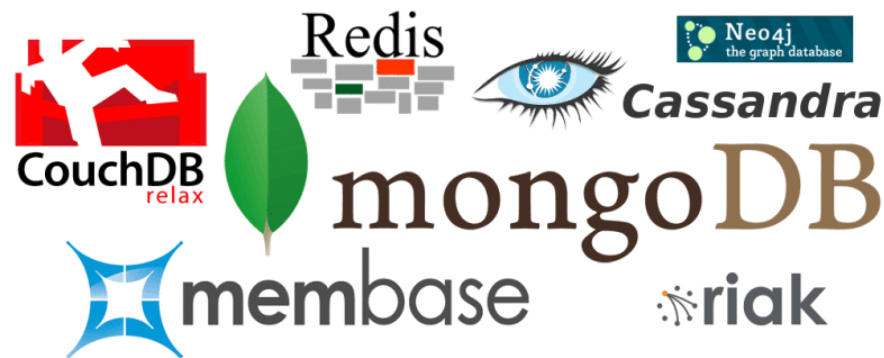
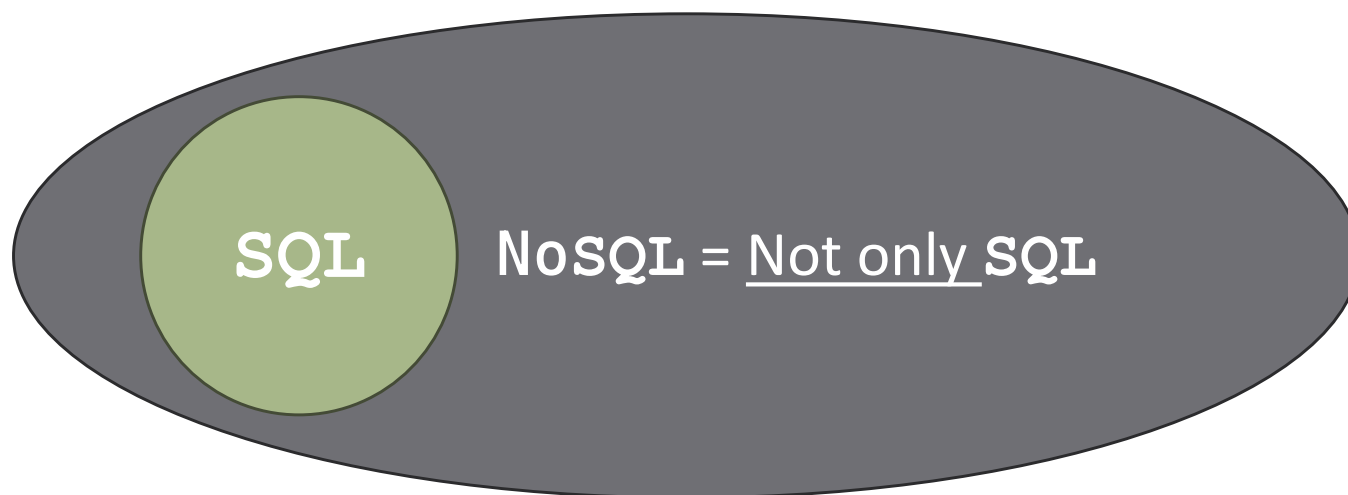


vaex 4.0.0-alpha.13  
documentation



# Duży bałagan (ang. *complexity*)

*Technologie*



# Modele 3V i 5V

Volume



Czy jestem w stanie **zmieścić** dane na dostępnych nośnikach?

Velocity



Jak **szybko** nadchodzą nowe dane?

Variety



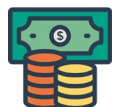
Jaki jest format / struktura / **źródła** danych?

Veracity



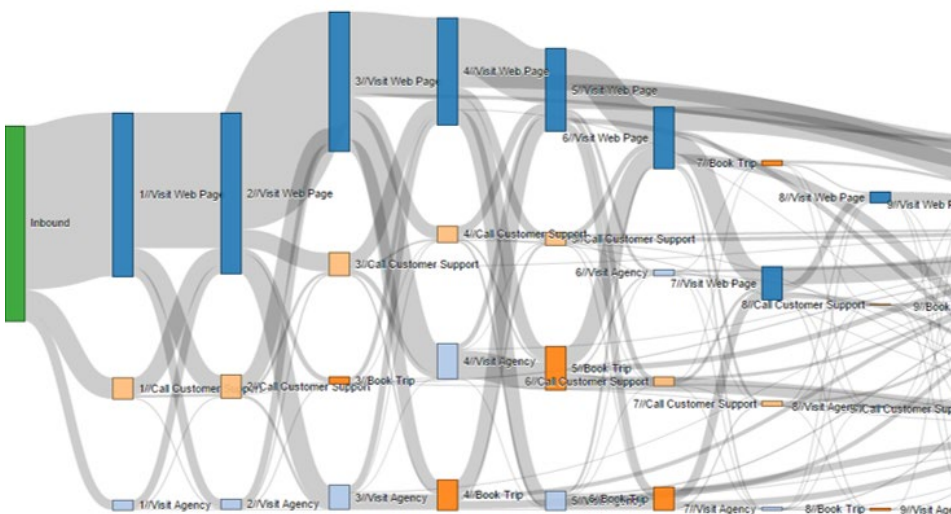
Czy dane są **wiarygodne** / kompletne / spójne?

Value

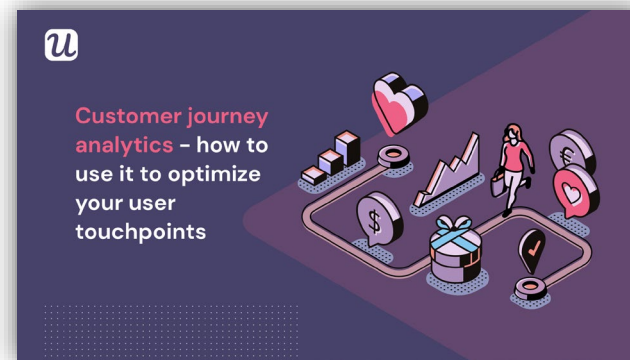


Czy **koszt** przetwarzania danych **się zwróci**?

# Scieżki zakupowe

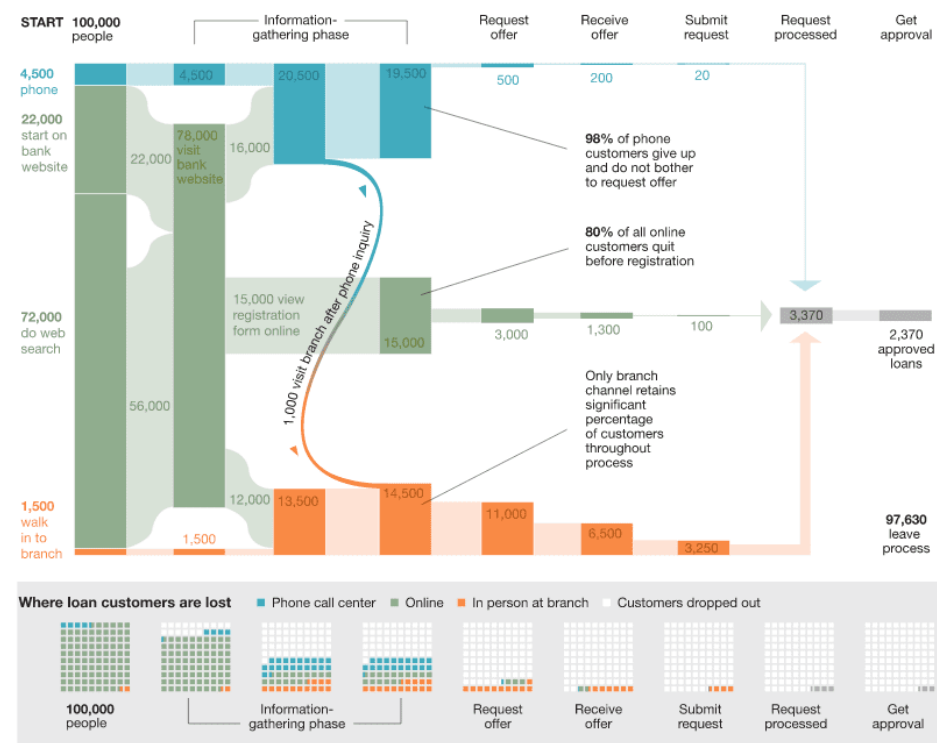


- Porzucone koszyki
- Zakupy krzyżowe, czego nie kupują
- Prognozowanie zamówień – stany magazynowe
- Segmentacja klientów, profilowanie (jakich filtrów używa)
- Odejście klienta – skąd pozyskany, niezbalansowane klasy
- Deduplikacja klientów (różne loginy Neo4j)
- Wielkość sprzedaży towaru – istotność towaru dla całości sprzedaży



Mapping customer flows highlights important pain points.

Average monthly customer flows for loan products by channel,<sup>1</sup> indexed to 100,000



<sup>1</sup>Preapproved loans excluded.

McKinsey&Company | Source: Call-center data; Google Analytics; interviews; McKinsey analysis

# Modele przetwarzania danych w ekosystemie Big Data

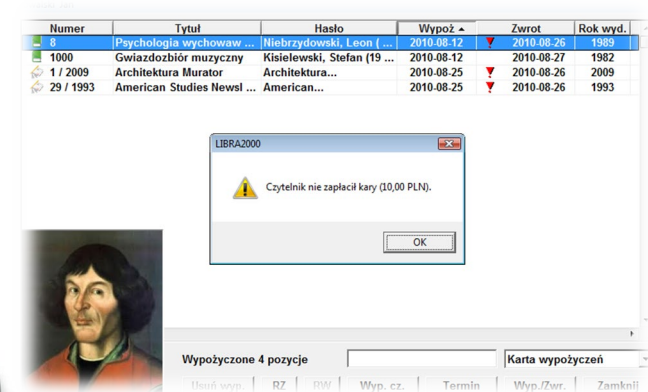
Kiedyś to było...



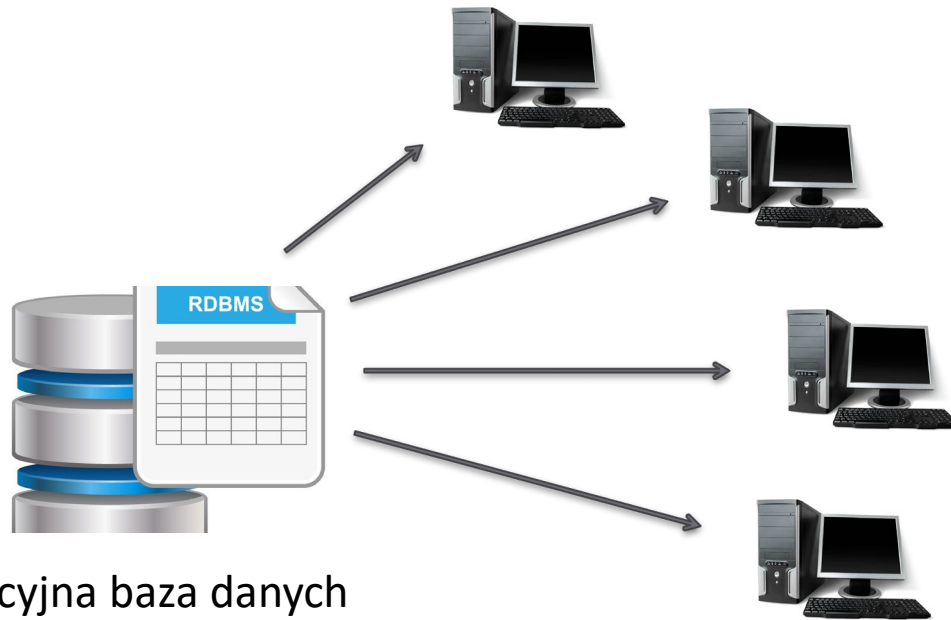
Relacyjna baza danych



Przetwarzanie informacji

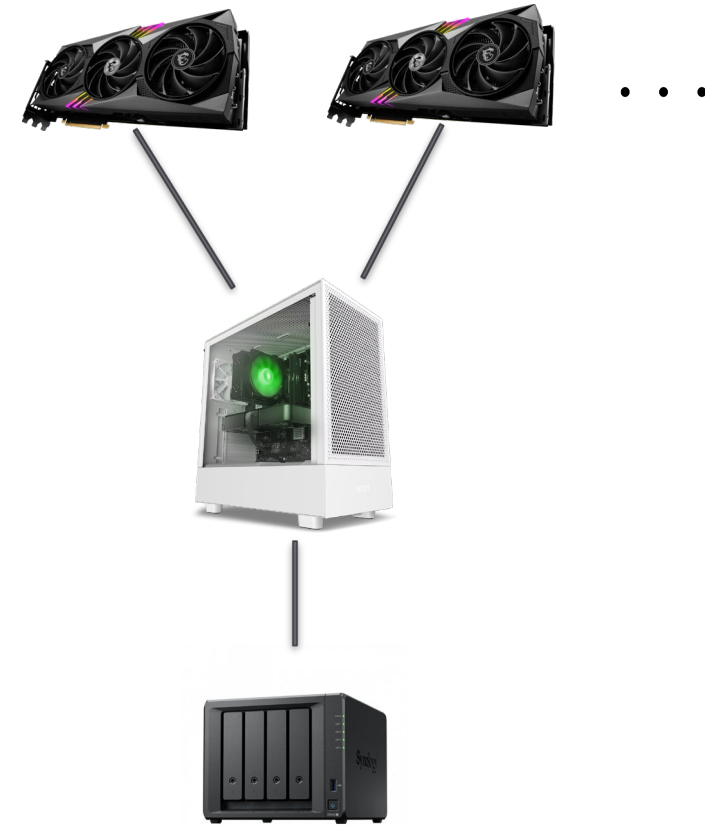


*Kiedyś to było...*



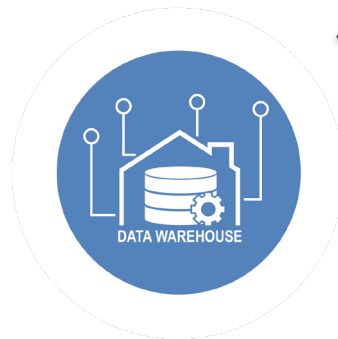
Relacyjna baza danych

...



NAS  
*Network Attached Storage*

*Kiedyś to było...*



Hurtownia danych



...

# Rozproszony system obliczeniowy

- Wiele komputerów widzianych przez użytkownika jako jedna całość

Skalowanie

p  
i  
o  
n  
o  
w  
e



Skalowanie p o z i o m e



W dłuższej perspektywie  
trudno wdrożyć nowe technologie

- wysokie koszty na wejściu w system zarządzania
- łatwo uzyskać niezawodność
- trudno spełnić warunki ACID dla transakcji  
*Atomicity, Consistency, Isolation, Durability*
- warunki BASE  
*Basically Available, Soft state, Eventual consistency*
- język dostępu do danych komplikuje się

# Rozproszony system obliczeniowy

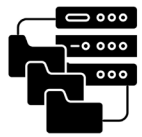


Rozproszony system  
plików

GFS



# Rozproszony system obliczeniowy



Rozproszony system plików

GFS



System nierelacyjnych baz danych



Cloud  
Bigtable



Amazon DynamoDB



mongo DB

# Rozproszony system obliczeniowy



Rozproszony system plików

GFS



System nierelacyjnych baz danych



Cloud  
Bigtable



Amazon DynamoDB



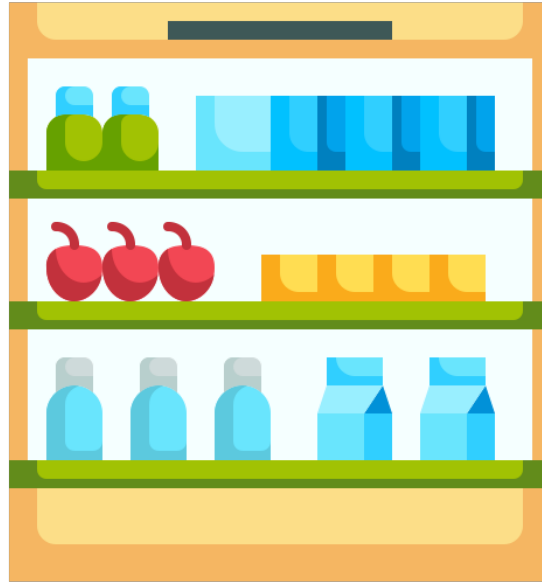
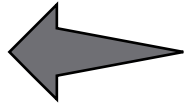
mongo DB

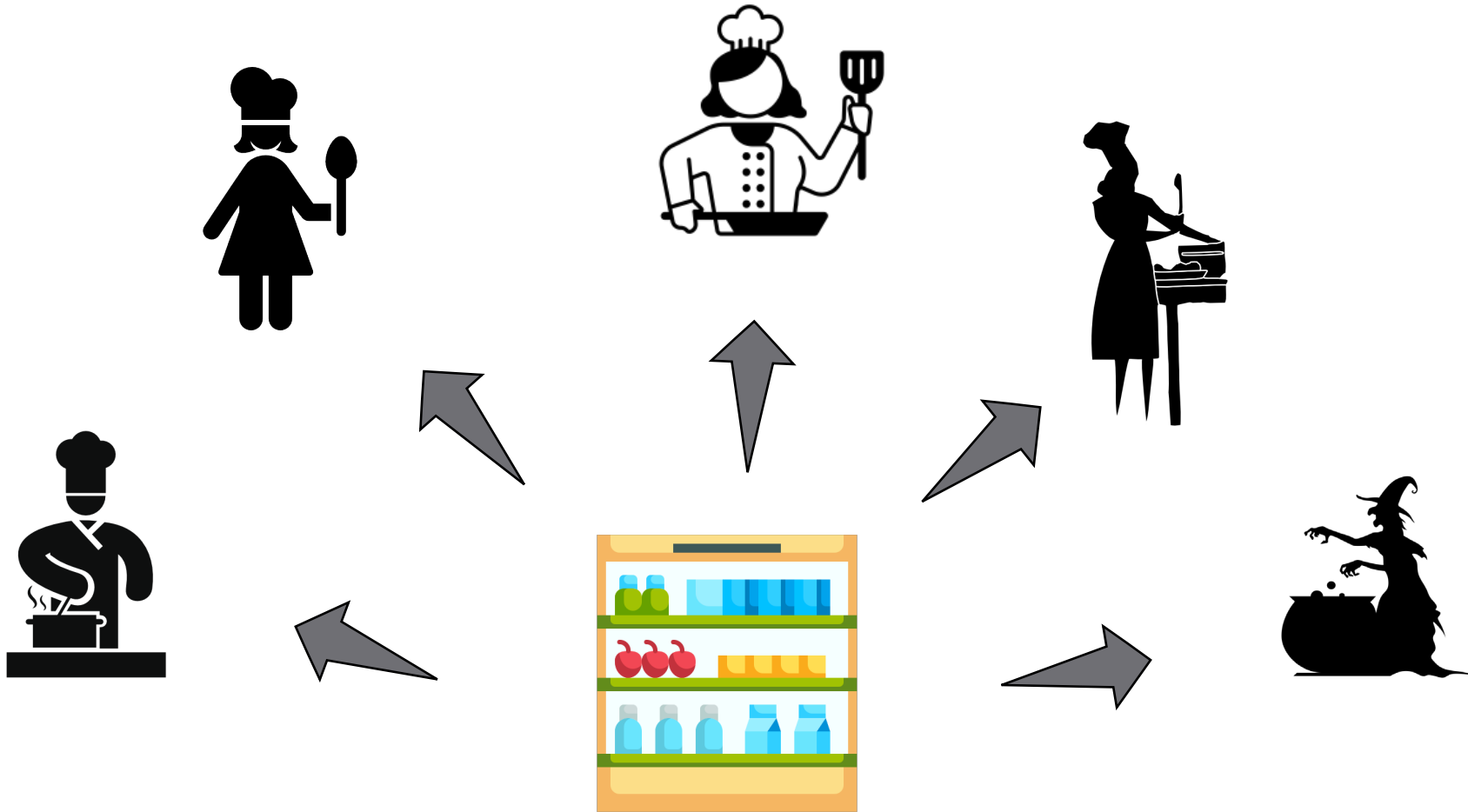


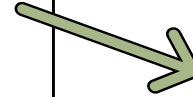
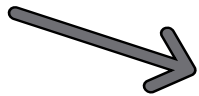
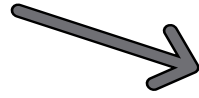
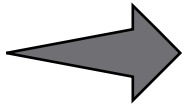
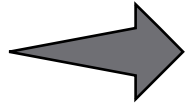
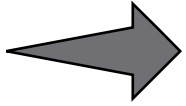
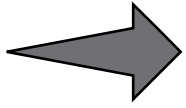
Środowisko przetwarzania rozproszonego



modele programowania



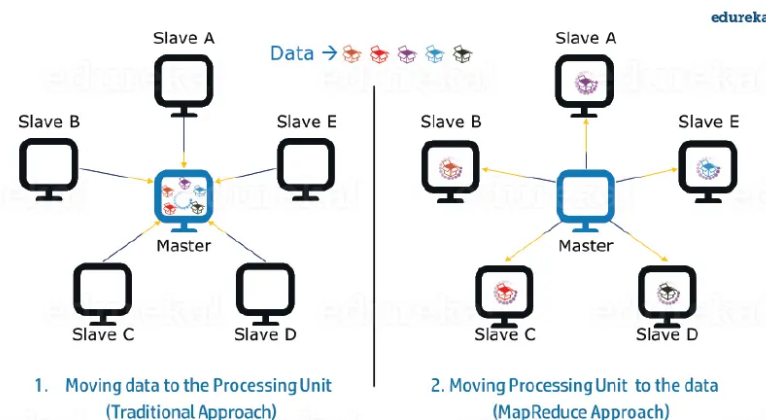
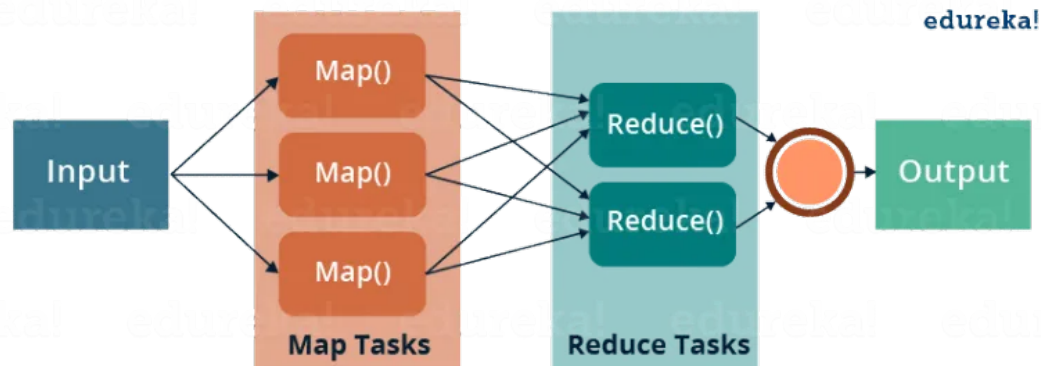
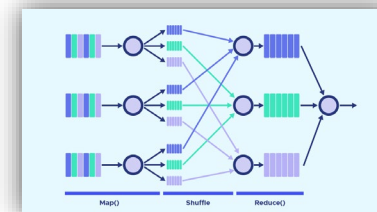




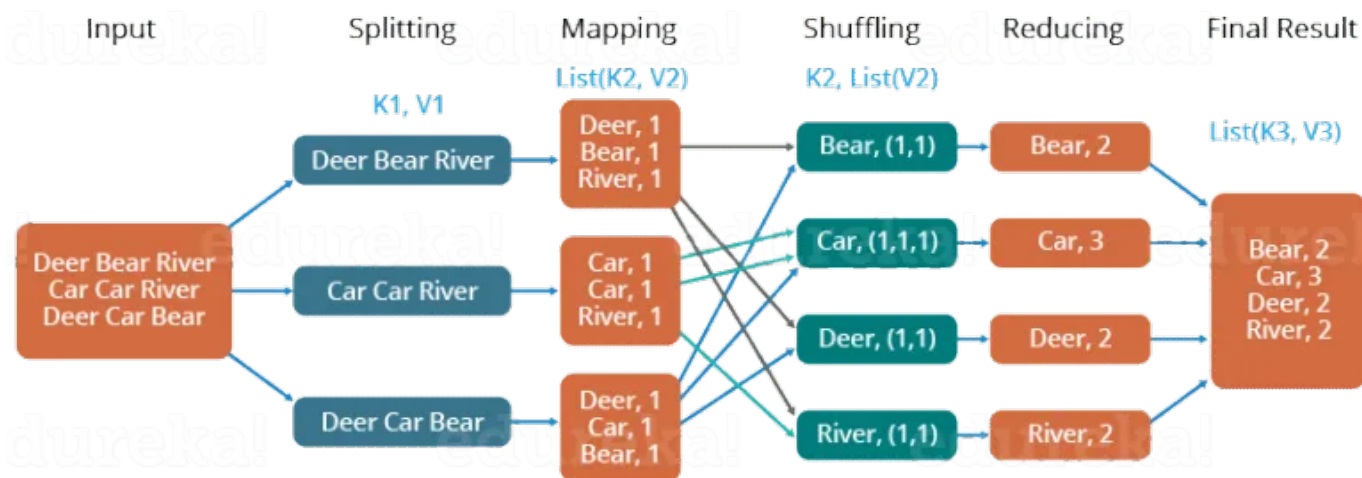
MAP

REDUCE

# Model przetwarzania MapReduce

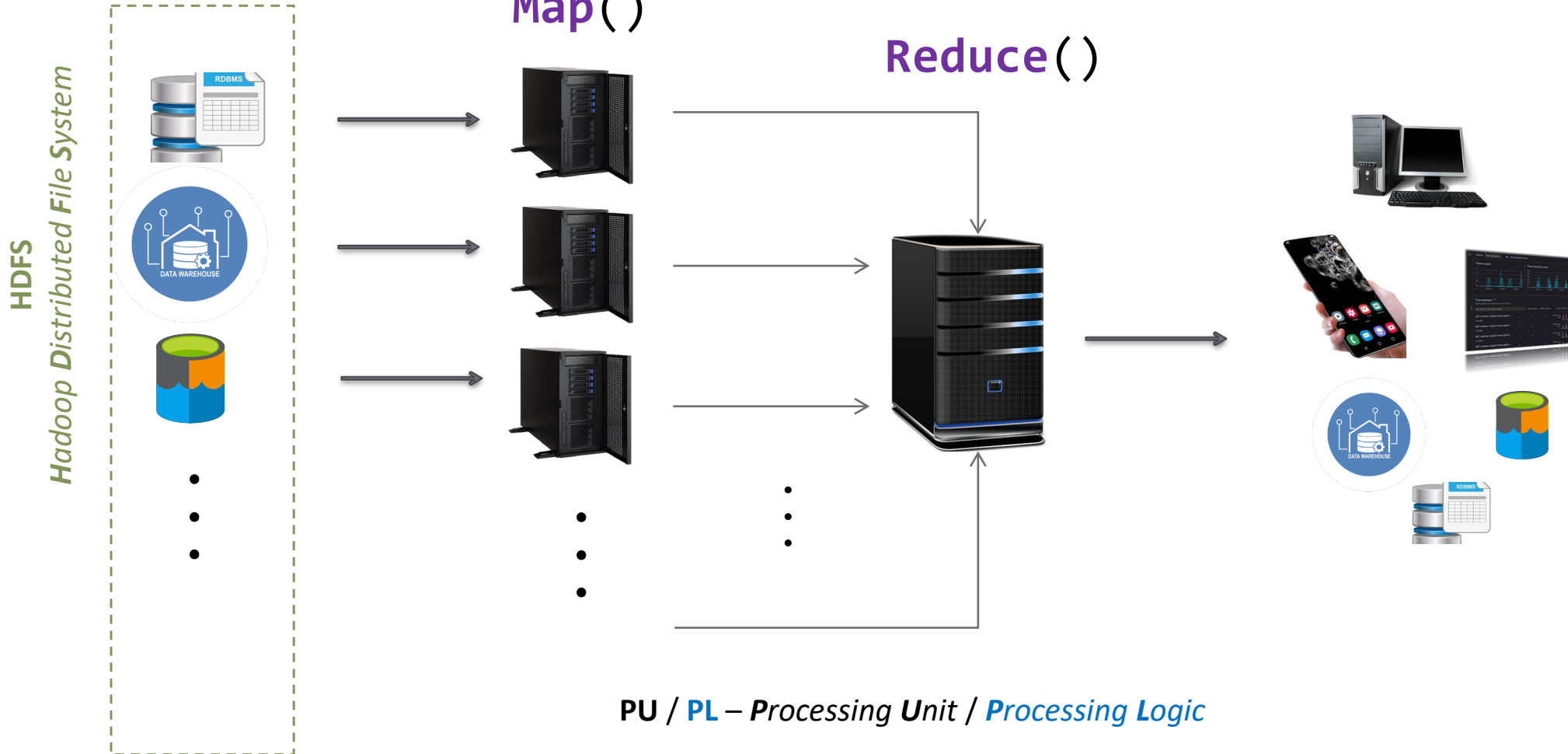


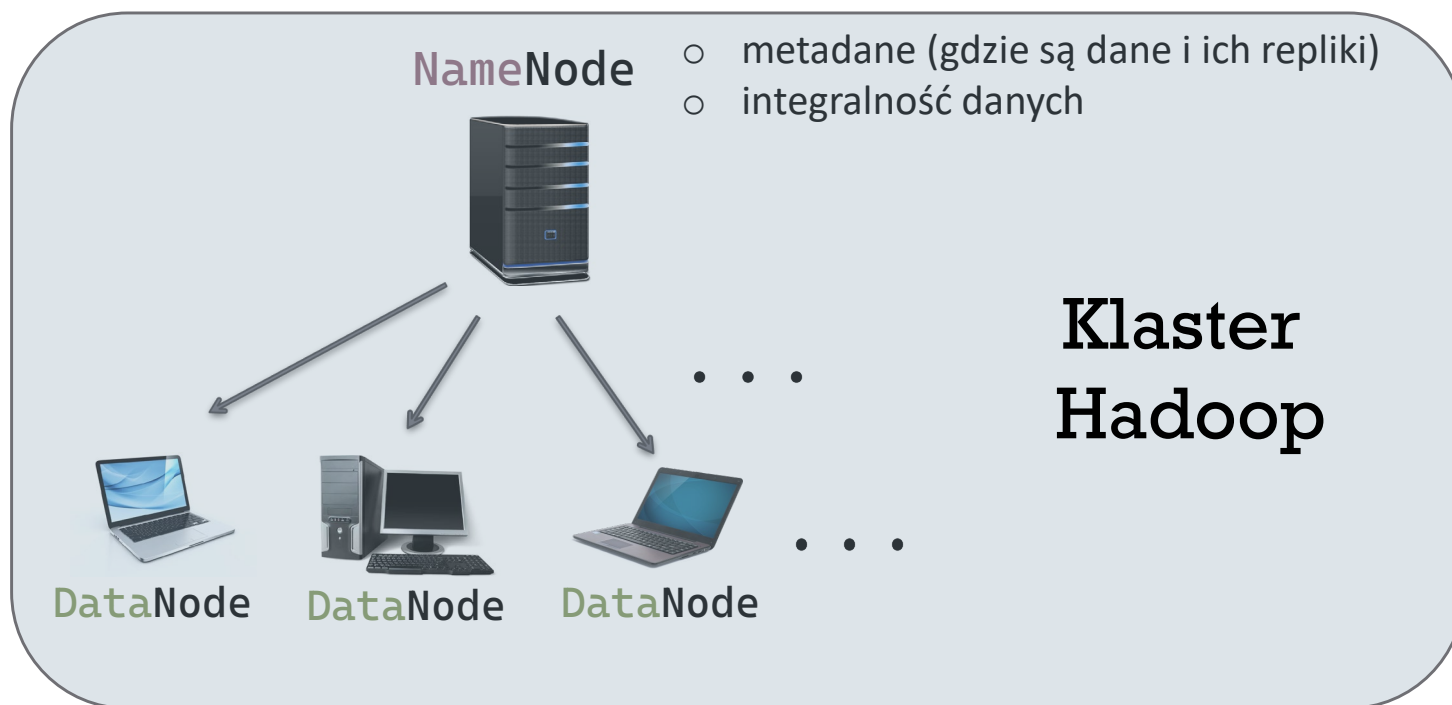
## The Overall MapReduce Word Count Process



Dean, Jeffrey, and Sanjay Ghemawat. "MapReduce: Simplified data processing on large clusters." (2004).

# Model MapReduce





Rozwiązuje problemy:

- Volume
- Veracity + Variety

# *Yet Another Resource Negotiator*



**Organizuje** przetworzenie ogromnej ilości danych w sposób efektywny





Otwarta platforma do **przechowywania i przetwarzania** dużych wolumenów danych w **rozproszonych klastrach** bazujących na **ogólnodostępnym sprzęcie** komputerowym

Podział plików

+

Podział zadań

+

Niezawodność

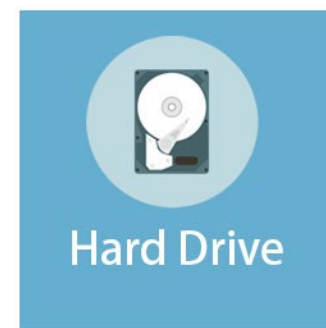
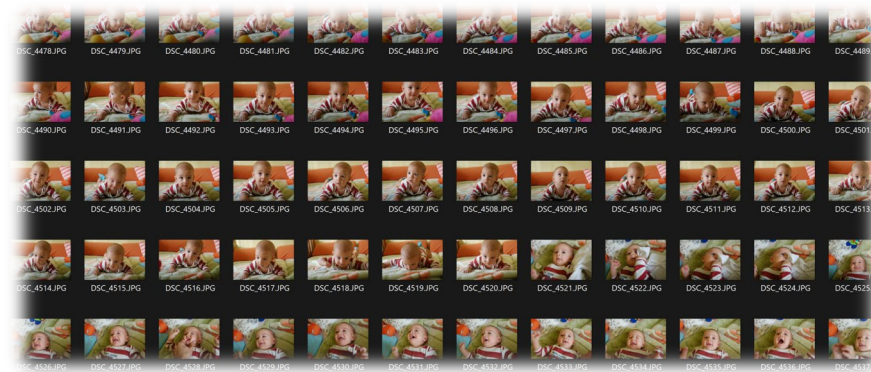
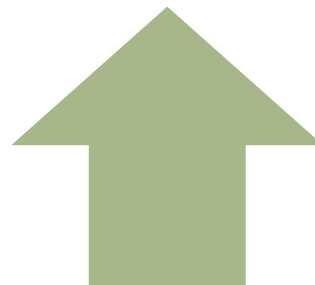
obsługa awarii  
sprzętu i danych



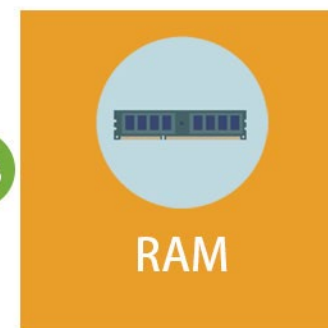
Koszt  
filtrowania i  
usuwania



Koszt  
przechowywania  
bezużytecznych  
danych



vs





Model przetwarzania  
oparty na **danych**  
przechowywanych  
na dyskach



Nieefektywny w zadaniach  
iteracyjnych i interakcyjnych



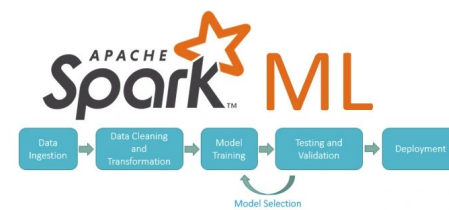
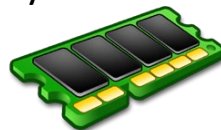
Uczenie  
maszynowe



Przetwarzanie  
strumieni



Model przetwarzania  
oparty na **danych**  
**pośrednich**  
przechowywanych  
w pamięci



DataFrame		
	Website	Visited
0	datagy.io	2023-01-23
1	google.com	2023-01-24
2	bing.com	2023-01-04

Otwarta platforma do przetwarzania w czasie rzeczywistym  
danych w rozproszonym systemie obliczeniowym

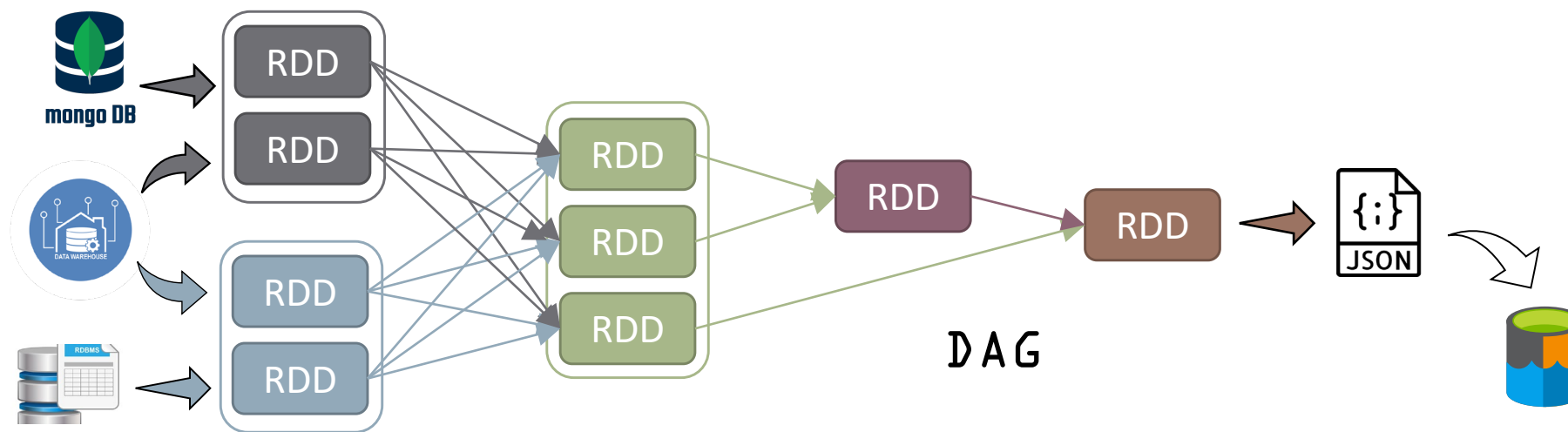


# Model przetwarzania w *Spark*

- Dane lub instrukcje ich wytworzenia są reprezentowane przez **RDD**

*Resilient Distributed Dataset*

- RDD może być większy od dostępnej pamięci



- *DataFrame*
- *Dataset*

# Podsumowanie Hadoop vs Spark



Key features	Apache Spark	Hadoop MapReduce
Speed	About 100 times faster than Hadoop	Faster than non-distributed systems
Data model	Use in-memory model to transfer data between RDD	Use HDFS to read, process and store large amounts of information (always use persistent storage)
Created in	Scala	Java
Style of processing	Batch, real-time, iterative, interactive, graph	Batch
Caching	Store data in memory	Caching data is not supported

AUDYCJA ZAWIERA LOKOWANIE PRODUKTU

# Inne rozwiązania



Photon

The next generation engine for the Lakehouse

Photon is the next generation engine on the Databricks Lakehouse Platform that provides extremely fast query performance at low cost – from data ingestion, ETL, streaming, data science and interactive queries – directly on your data lake. Photon is compatible with Apache Spark™ APIs, so getting started is as easy as turning it on – no code changes and no lock-in.



**Massively scalable storage for demanding applications**

Red Hat® Ceph® Storage, an open, massively scalable, simplified storage solution for modern data pipelines. Engineered for data analytics, AI/ML, and emerging workloads, Red Hat Ceph Storage delivers software-defined storage on your choice of industry-standard hardware.

# Przyszłość

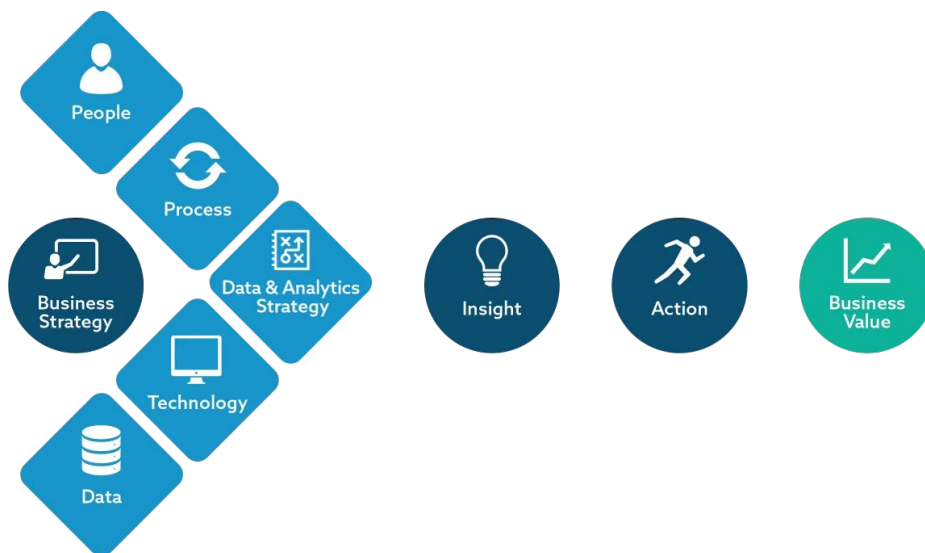


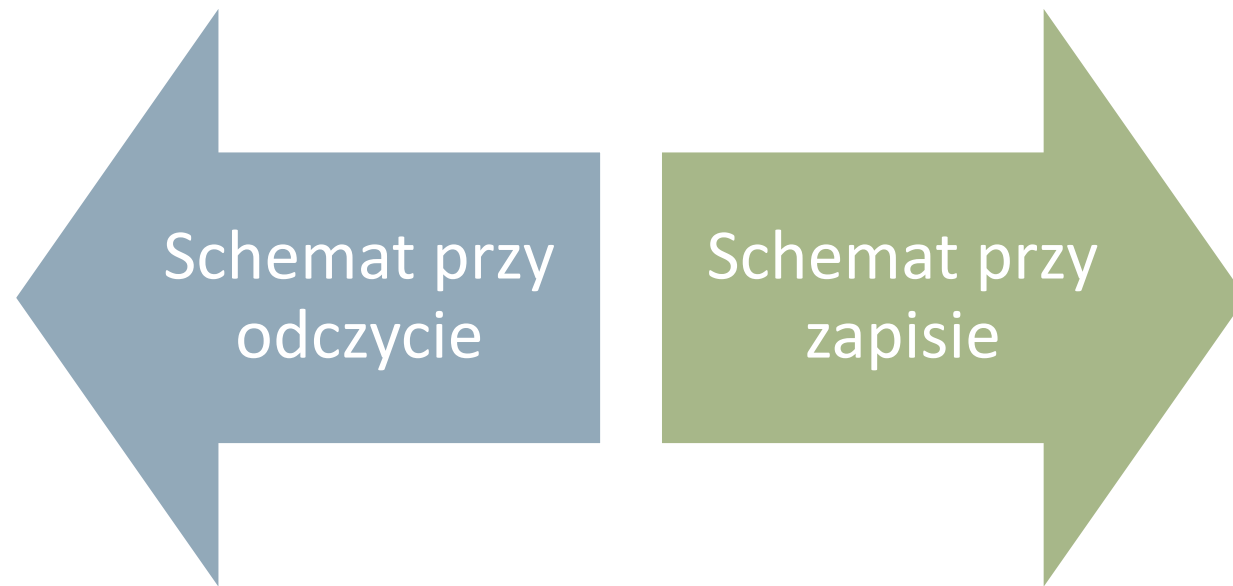
- Nowy model rozproszonej bazy danych (patrz *blockchain*)  
prawa własności ma twórca danych
- *Interplanetary File System* (IPFS) – sieć peer-to-peer z własnym DNS (IPNS) oraz wersjonowaniem plików



Uwagi końcowe

# Przetwarzanie danych w ekosystemie Big Data





*Baza pozbawiona schematów* = baza z elastyczniejszym schematem

# THE NEVER ENDING STORY

