

Wstęp do Big Data



University
of Bielsko-Biala



iBigWorld:
Innovations for Big Data in a Real World

Prof. dr. Dragan Stojanovic, **UNI**

Prof. dr. Natalija Stojanovic, **UNI**

Disclaimer: Funded by the European Union. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the National Agency (NA). Neither the European Union nor NA can be held responsible for them.



Big Data definitions

- “Big data odnosi się do zbiorów danych, których rozmiar przekracza możliwości typowych narzędzi oprogramowania bazodanowego do gromadzenia, przechowywania, zarządzania i analizowania.” **The McKinsey Global Institute, 2012**
- “Big data to dziedzina, która zajmuje się sposobami analizowania, systematycznego wydobywania informacji lub innego postępowania ze zbiorami danych, które są zbyt duże lub złożone, aby mogły być obsługiwane przez tradycyjne aplikacje do przetwarzania danych.” **Wikipedia, 2020**

Big Data definitions

- “Big data polega głównie na zbieraniu liczb i wykorzystywaniu tych liczb do przewidywania przyszłości. Im większy masz zestaw danych, tym dokładniejsze będą prognozy dotyczące przyszłości.” **Anthony Goldbloom, Kaggle’s founder**
- “Big data to zbiory informacji o dużej objętości, dużej zmienności lub dużej różnorodności, które wymagają nowych form przetwarzania w celu wspomaganie podejmowanie decyzji, odkrywanie nowych zjawisk oraz optymalizacji procesów.” **Gartner**

Big Data przedrostki

- Międzynarodowy Układ Jednostek

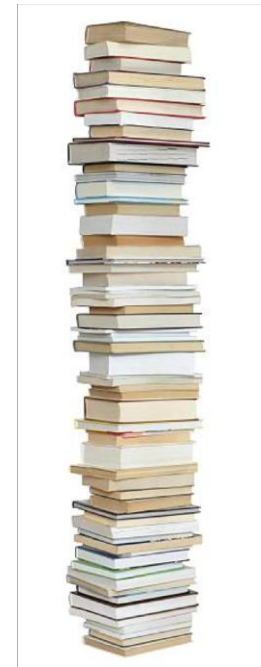
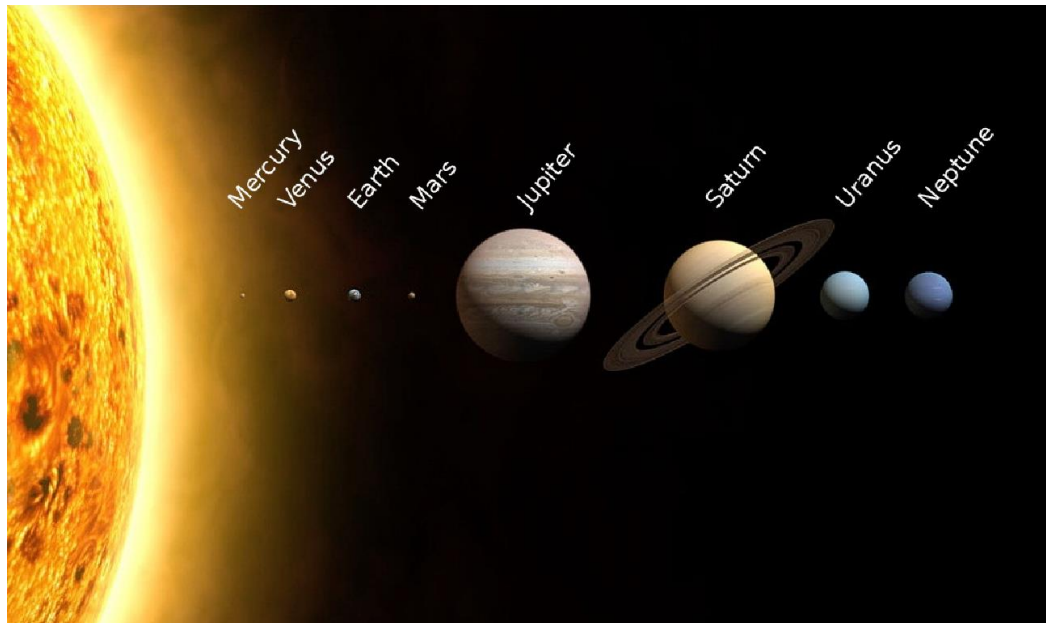
kilo (k)	1,000 (3 zeros)
Mega (M)	1,000,000 (6 zeros)
Giga (G)	1,000,000,000 (9 zeros)
Tera (T)	1,000,000,000,000 (12 zeros)
Peta (P)	1,000,000,000,000,000 (15 zeros)
Exa (E)	1,000,000,000,000,000,000 (18 zeros)
Zetta (Z)	1,000,000,000,000,000,000,000 (21 zeros)
Yotta (Y)	1,000,000,000,000,000,000,000,000 (24 zeros)

Big Data - liczby

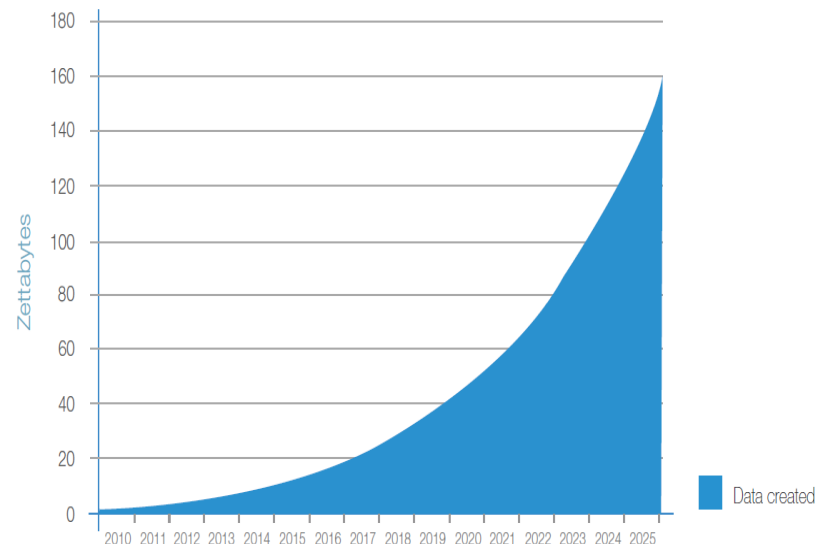
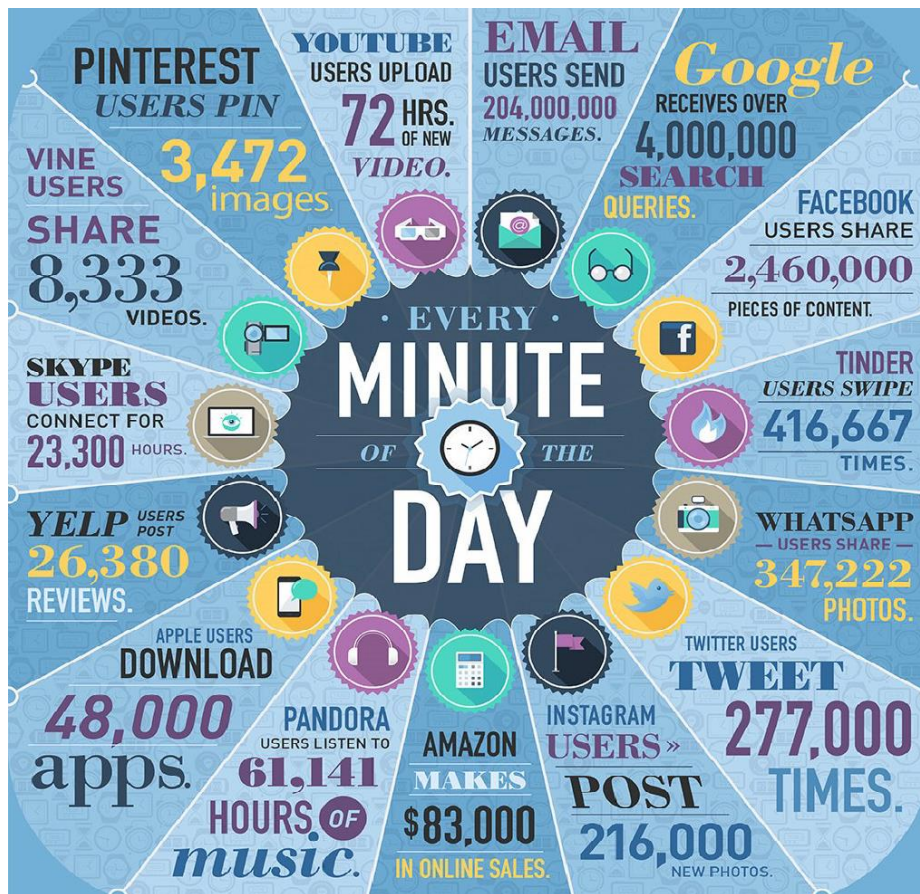
- Ile zgromadzono danych na świecie? Jakie są szacunki?
 - 800 Terabytes, 2000
 - 160 Exabytes, 2006 (1EB = 10^{18} B)
 - 4.5 Zettabytes, 2013 (1ZB = 10^{21} B)
 - 44 Zettabytes by 2020
 - 163 Zettabytes by 2025
- Jak dużo danych generowanych jest w ciągu dnia?
 - 2.5 Exabytes
 - 8 TB, Twitter
 - 50 TB, Facebook
- 90% danych na całym świecie wygenerowano w ciągu ostatnich 2 lat!

Jak duże są zbiory Big Data?

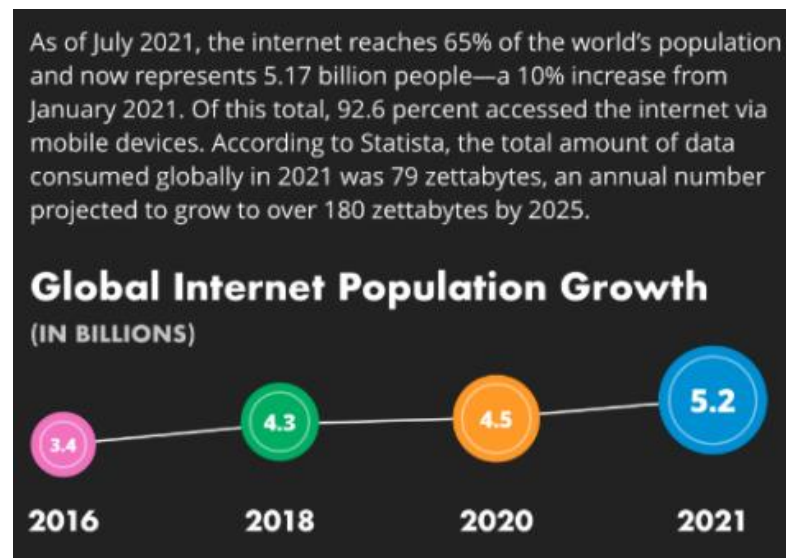
- Spróbujmy zrobić stos książek zawierający Zettabajt danych
- 1 Zettabyte = stos książek, których wysokość jest PIĘCIOKROTNĄ odległością między Ziemią a Plutonem



Więcej, coraz więcej danych...



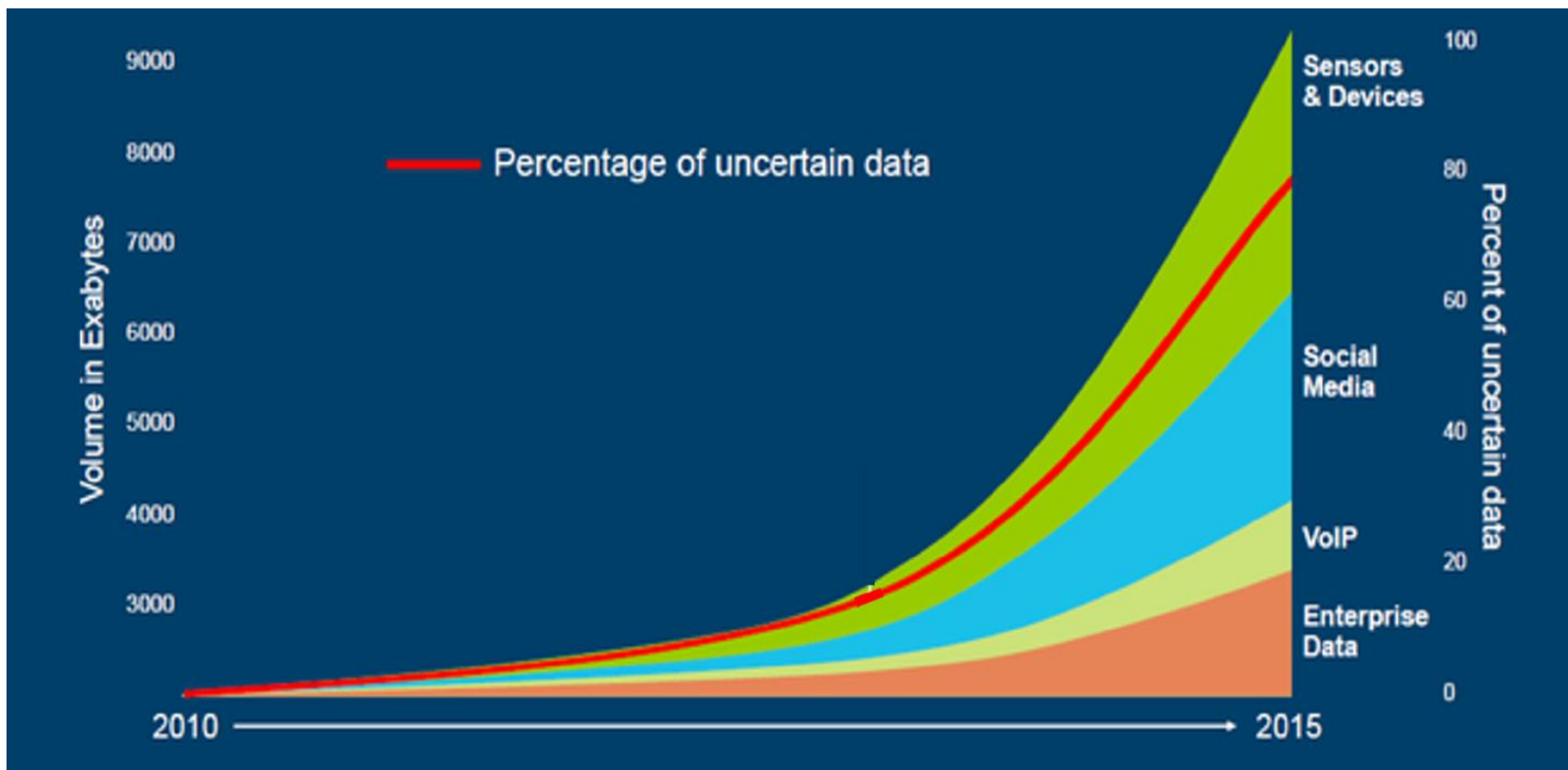
Dane nigdy nie śpią (*domo.com*)



<https://www.visualcapitalist.com/from-amazon-to-zoom-what-happens-in-an-internet-minute-in-2021/>

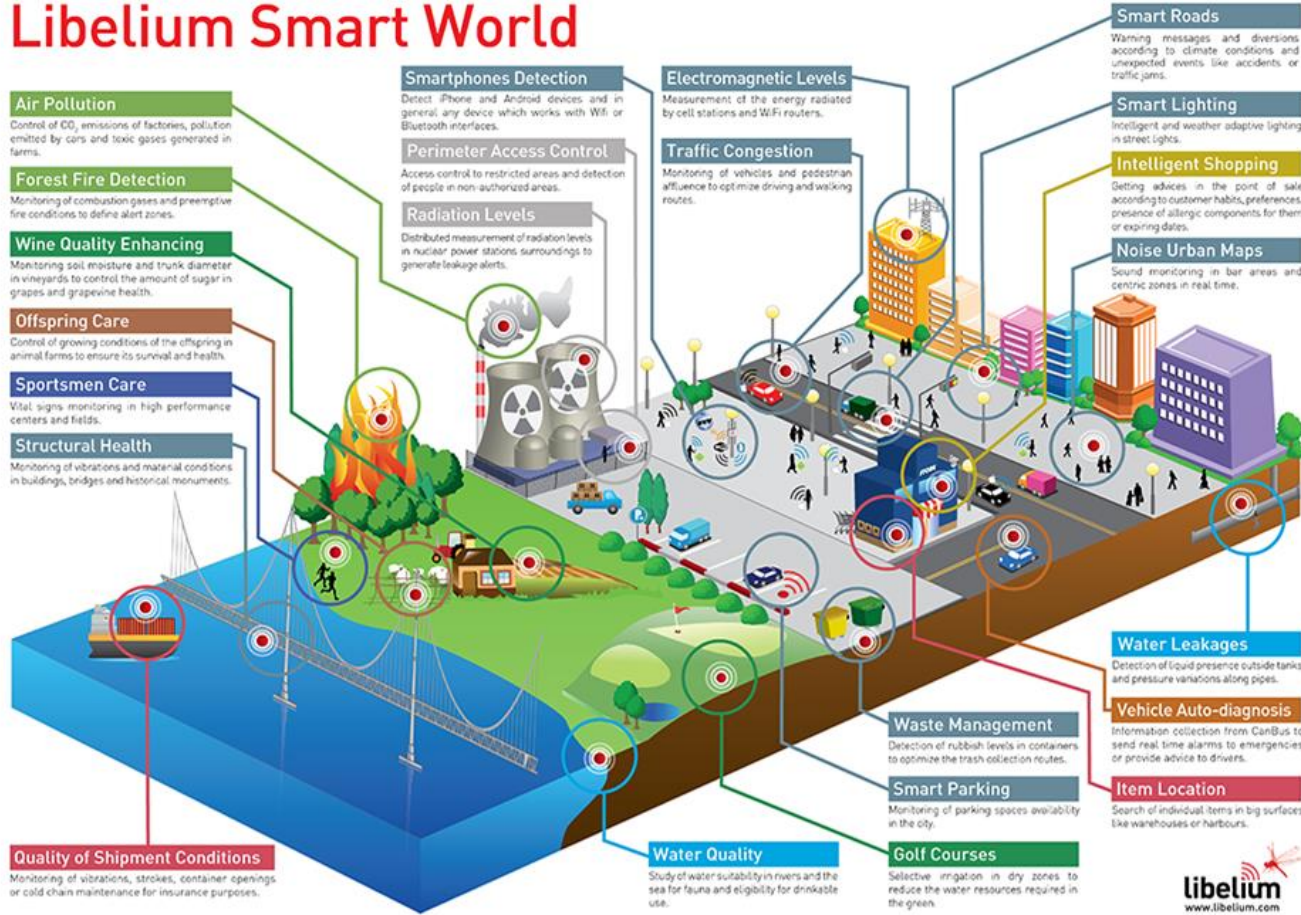


Źródła danych



Źródła danych - Inteligentny świat

Libelium Smart World



Internet rzeczy (IoT)

'There will be as many as **40 TO 80 BILLION** connected objects by 2020.

There will be **10** connected objects for every man, woman, and child on the **PLANET.**



Internet of things

Vehicle, asset, person & pet monitoring & controlling

Agriculture automation

Energy consumption

Security & surveillance

Building management

Embedded Mobile

Everyday things get connected for smarter tomorrow

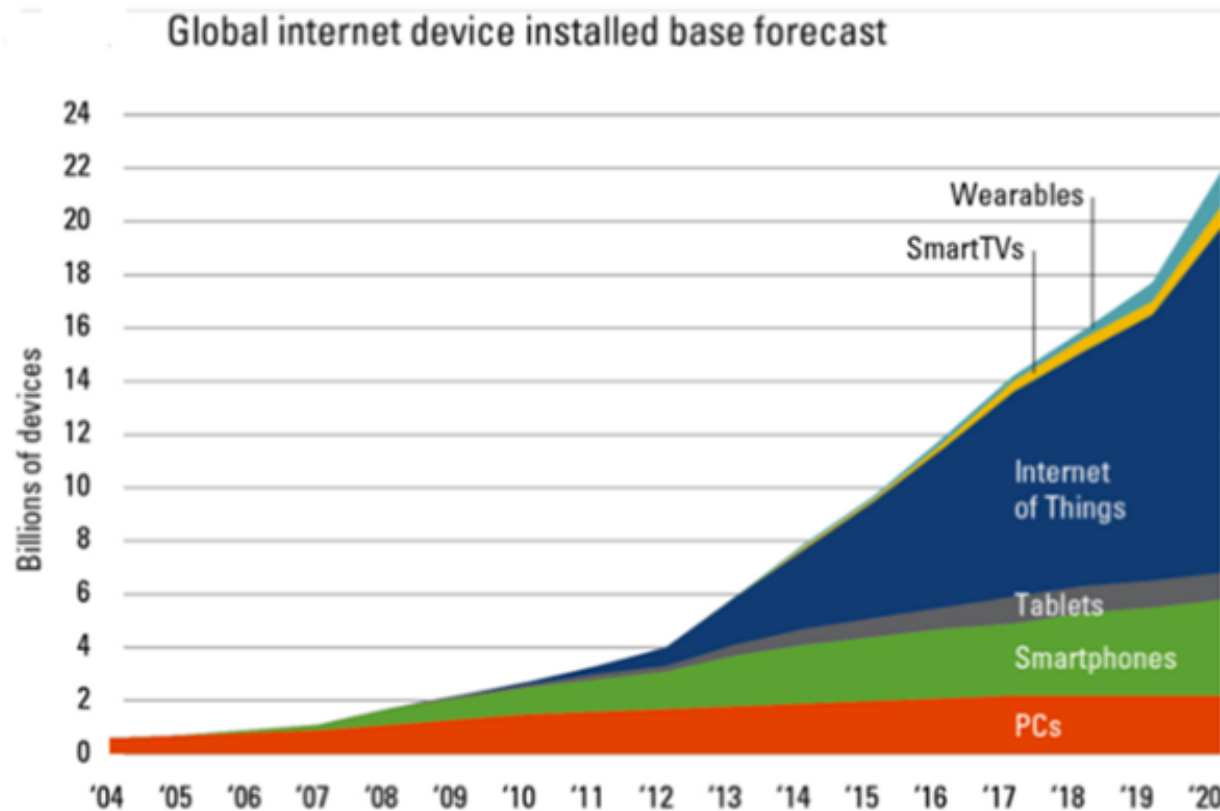
M2M & wireless sensor network

Everyday things

Smart homes & cities

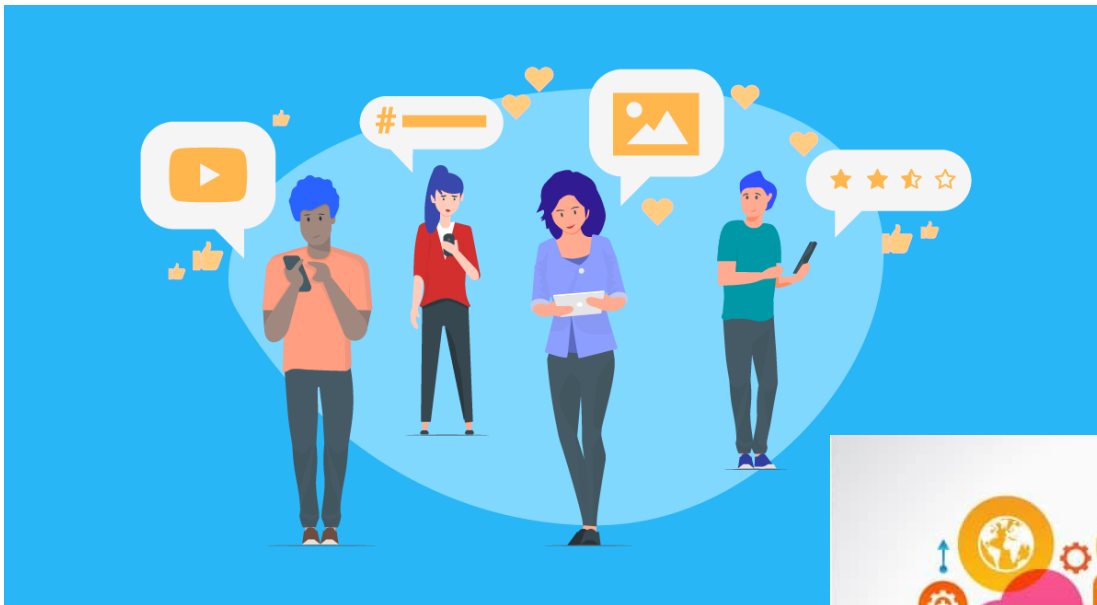
Telemedicine & healthcare

Światowy sieć - rozwój urządzeń



Global Internet Device Forecast (Source: Market Realist, Gartner, IDC, Strategy Analytics.)

Dane generowane przez użytkowników

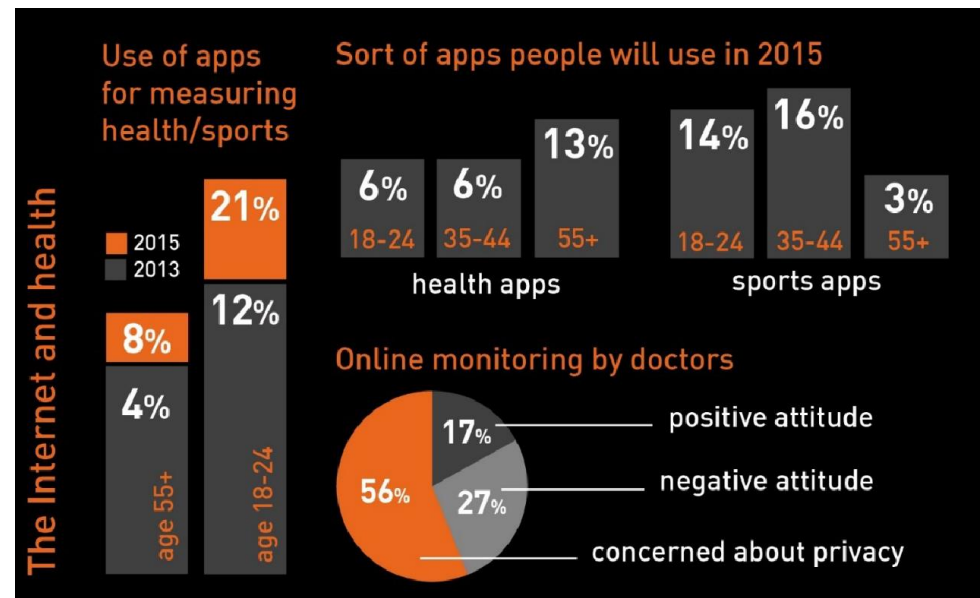
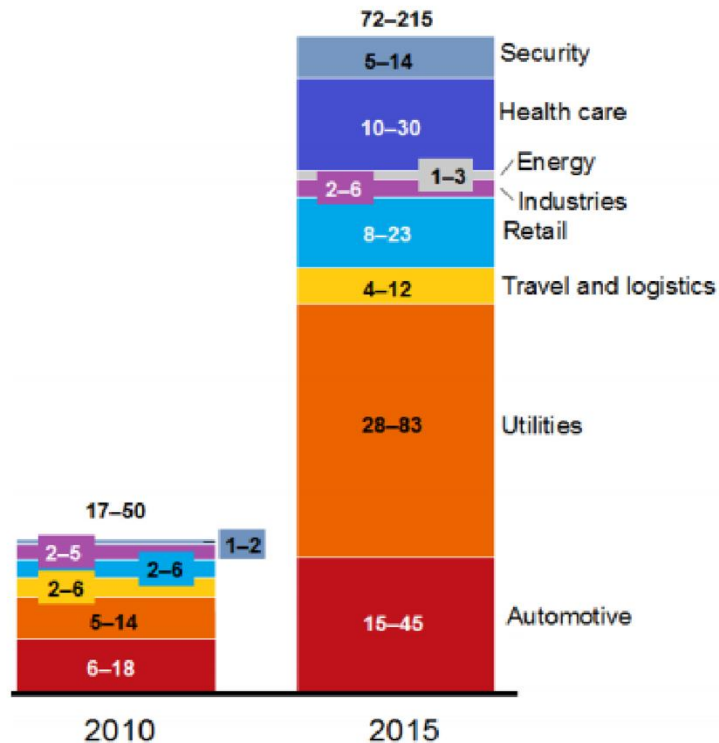


2006



Big Data w różnych dziedzinach życia

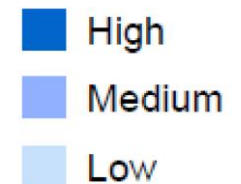
- Dane mają kluczowe znaczenie, niezależnie od tego, czy jesteśmy gigantycznym przedsiębiorstwem, czy indywidualną osobą.



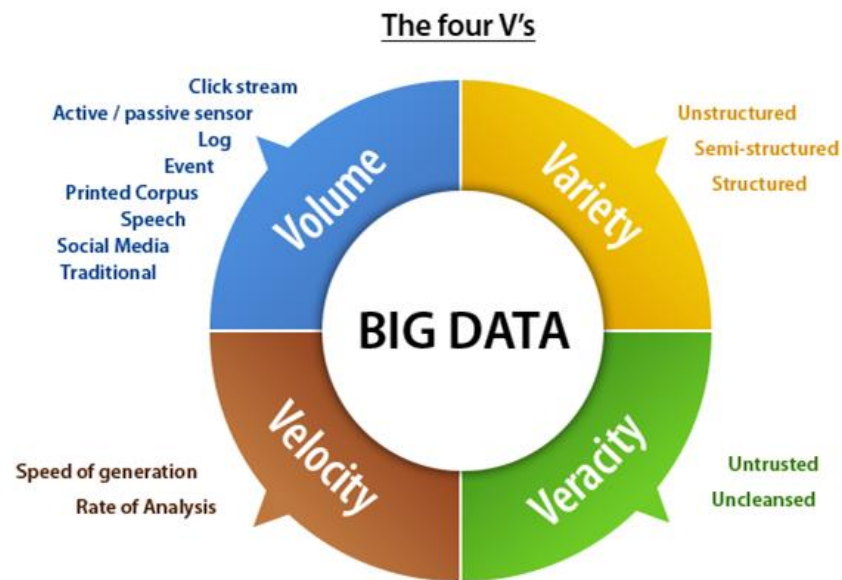
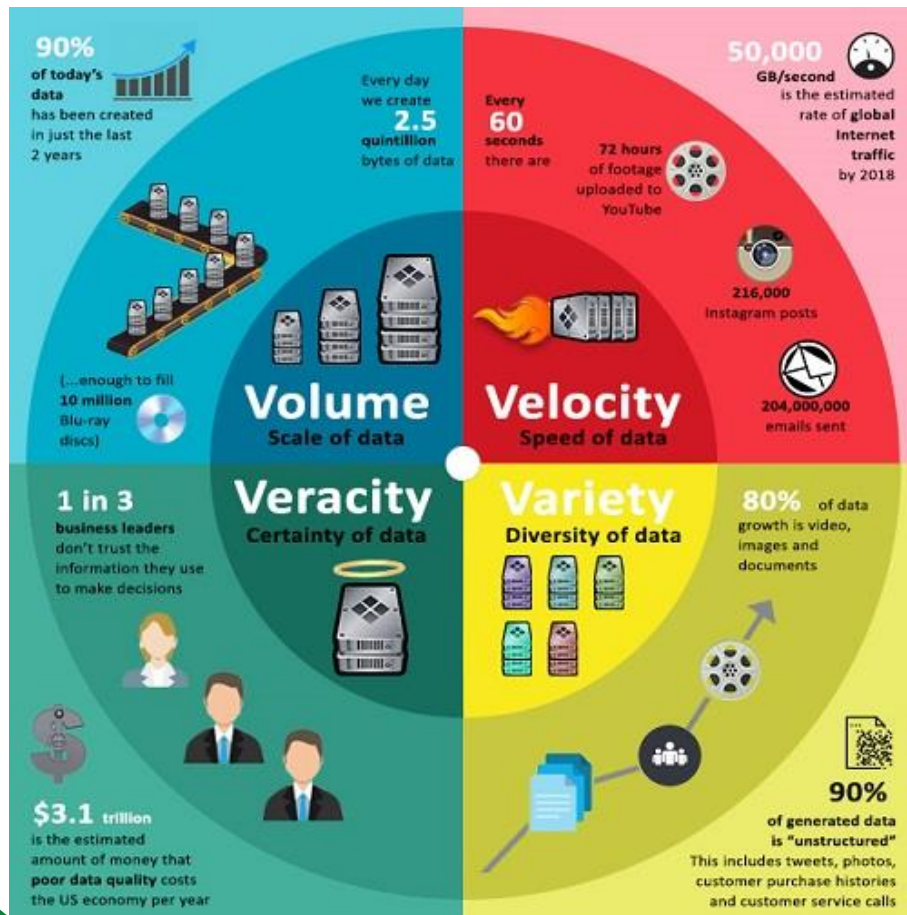
Big Data - typy danych

	Video	Image	Audio	Text/ numbers
Banking	Medium	Medium	Medium	High
Insurance	Low	Low	Low	High
Securities and investment services	Low	Low	Low	High
Discrete manufacturing	Medium	Medium	Low	High
Process manufacturing	Medium	Medium	Low	High
Retail	Medium	Low	Low	High
Wholesale	Low	Low	Low	High
Professional services	Medium	Medium	Medium	High
Consumer and recreational services	Medium	Low	Medium	Medium
Health care	Low	High	Low	High
Transportation	Low	Medium	Low	High
Communications and media ²	High	Medium	High	High
Utilities	Low	Medium	Low	High
Construction	Low	High	Low	Medium
Resource industries	Medium	Medium	Low	High
Government	High	Medium	High	High
Education	High	Medium	High	Medium

Penetration



4V Big Data + Wartość



Big Data: V4+Wartość (Value)

- Rozmiar (Volume)
 - Terabyte (10^{12}), Petabyte (10^{15}), Exabyte (10^{18}), Zettabyte (10^{21})
- Różnorodność (Variety)
 - Structured, semi-structured, unstructured
 - Tekst, obrazy, audio, wideo, Text, image, audio, video, record
- Zmienność (Velocity)
 - Okresowe, bliskie czasu rzeczywistego, czasu rzeczywistego
- Wiarygodność (Veracity)
 - Jakość danych może się znacznie różnić
- **Wartość (Value)**
 - Big Data może generować ogromne przewagi konkurencyjne

Dlaczego?

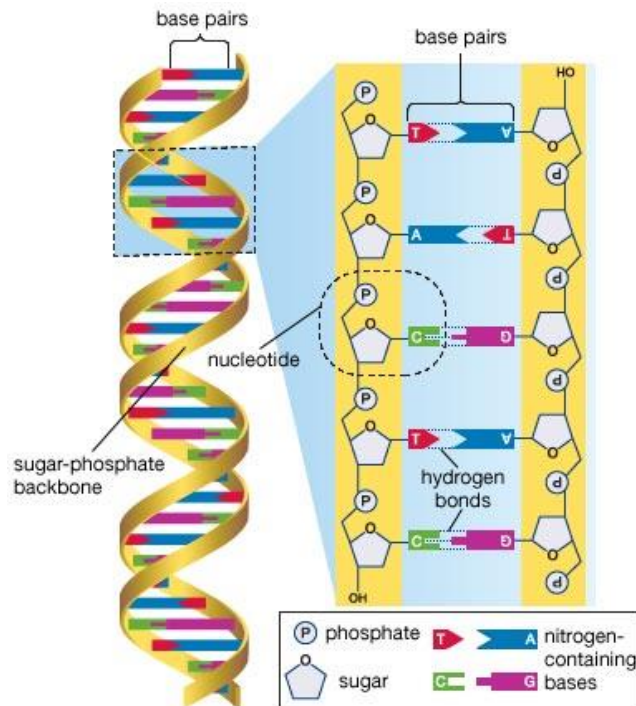
- Ponieważ mamy dane
 - Dane utworzone są już w formie cyfrowej
 - 40% przyrostu danych rocznie
 - 7 lat temu: Internet ma oczy i uszy
- Ponieważ możemy
 - 300 \$ za dysk, na którym można przechowywać całą muzykę świata
 - >40 lat prawa Moore'a → duże zasoby obliczeniowe
 - 40 years of databases → pojemność pamięci masowych
 - 68% firm zainwestowało w Big Data w 2018 roku
 - 57 miliardów \$ zainwestowanych w Big Data w 2018 roku
- “[Software is eating the world](#)”, Marc Andreessen, The Wall Street Journal on August 20, 2011.

Wielki zderzacz hadronów (LHC) CERN

- Zebrał 200 Petabytes danych from 2012, a 1 Petabyte danych jest przetwarzany każdego dnia

Ludzki genom

- DNA pojedynczego osobnika zawiera około 3,2 miliarda par nukleotydów DNA



© 2007 Encyclopædia Britannica, Inc.

Big Data – Potencjalne możliwości



US health care

- \$300 billion value per year
- ~0.7 percent annual productivity growth



Europe public sector administration

- €250 billion value per year
- ~0.5 percent annual productivity growth



Global personal location data

- \$100 billion+ revenue for service providers
- Up to \$700 billion value to end users



US retail

- 60+% increase in net margin possible
- 0.5–1.0 percent annual productivity growth

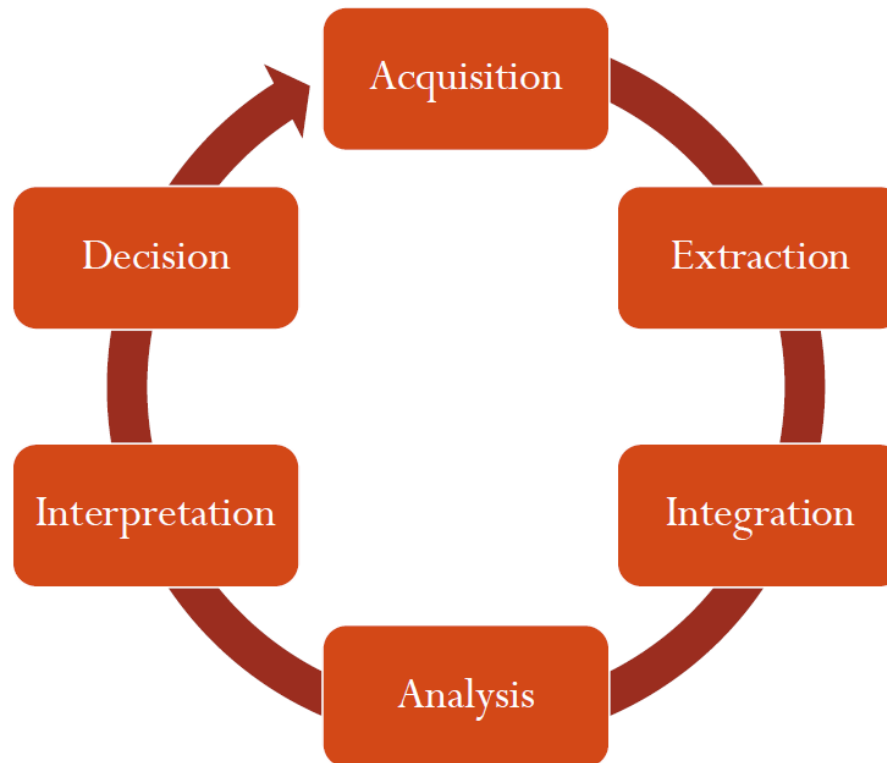


Manufacturing

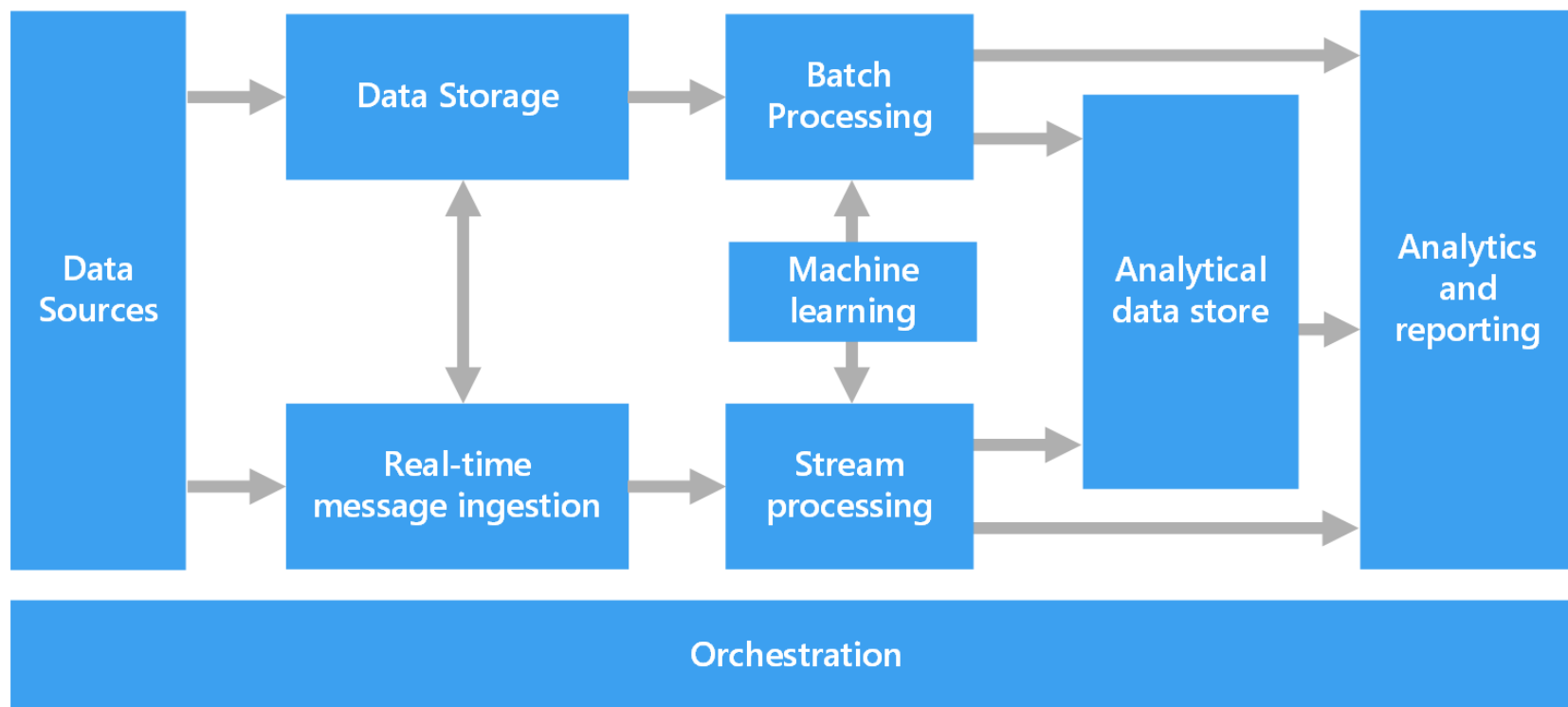
- Up to 50 percent decrease in product development, assembly costs
- Up to 7 percent reduction in working capital

Proces Big Data

- Podejmuj skuteczne decyzje strategiczne, wykorzystując dostępność przetwarzania, analizy i wizualizacji Big Data

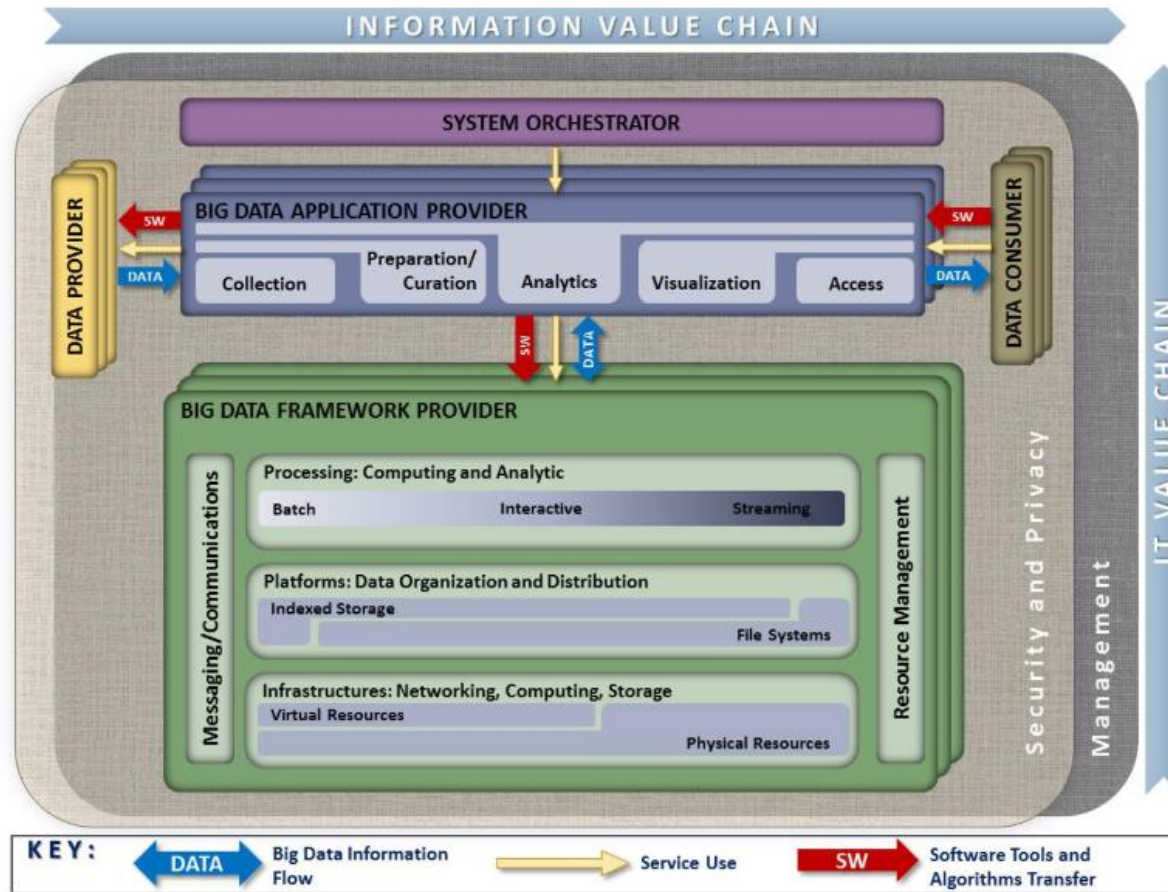


Architektura Big Data



<https://docs.microsoft.com/en-us/azure/architecture/data-guide/big-data/>

Architektura referencyjna NIST Big Data



<https://doi.org/10.6028/NIST.SP.1500-6r2>



Dystrybucja i replikacja

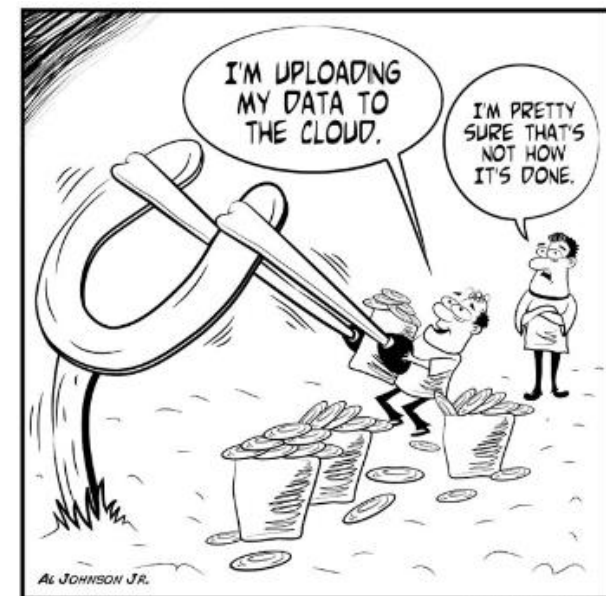
- Architektura rozproszona
 - Wspólne rozwiązanie architektoniczne do przetwarzania Big Data: klaster podstawowych zasobów sprzętowych, również w chmurze
 - **Skaluj w poziomie, a nie w górę (w pionie)!**
 - Wyzwania: elastyczność i przetwarzanie danych na krawędziach sieci
- Przetwarzanie rozproszone
 - Model nie współdzielony
 - Nowe paradygmaty programistyczne, np. programowanie funkcjonalne
- Replikacja zasobów
 - Dobrze znane rozwiązanie pozwalające na osiągnięcie odporności na uszkodzenia
 - Spójność ostateczna (twierdzenie CAP)

Jak implementować Big Data

- Tradycyjny sposób: korzystanie z klastra serwerów w siedzibie
 - Węzły obliczeniowe są przechowywane w szafach rackowych
 - 8-64 węzłów obliczeniowych w szafie
- Może być wiele szaf węzłów obliczeniowych
 - Węzły w szafie połączone są w sieć, zwykle gigabitowym Ethernetem
 - Szafy są połączone na innym poziomie sieci lub przełącznikiem
 - Przepustowość komunikacji wewnątrz szaf jest zwykle znacznie większa niż w przypadku komunikacji między szafami
- Rozważyć:
 - Konieczność zarządzania infrastrukturą sprzętową i
 - Platformy przetwarzania (pozyskiwanie, instalacja, konfiguracja, ...)

Gdzie przetwarzać Big Data

- Ale centra danych w chmurze znajdują się w rdzeniu sieci
- Główne wyzwania:
 - Przenieść dane do chmury
 - Opóźnienie nie jest zerowe (z powodu prędkości światła)!
 - Mniejszy problem: przepustowość sieci
 - Bezpieczeństwo i prywatność danych



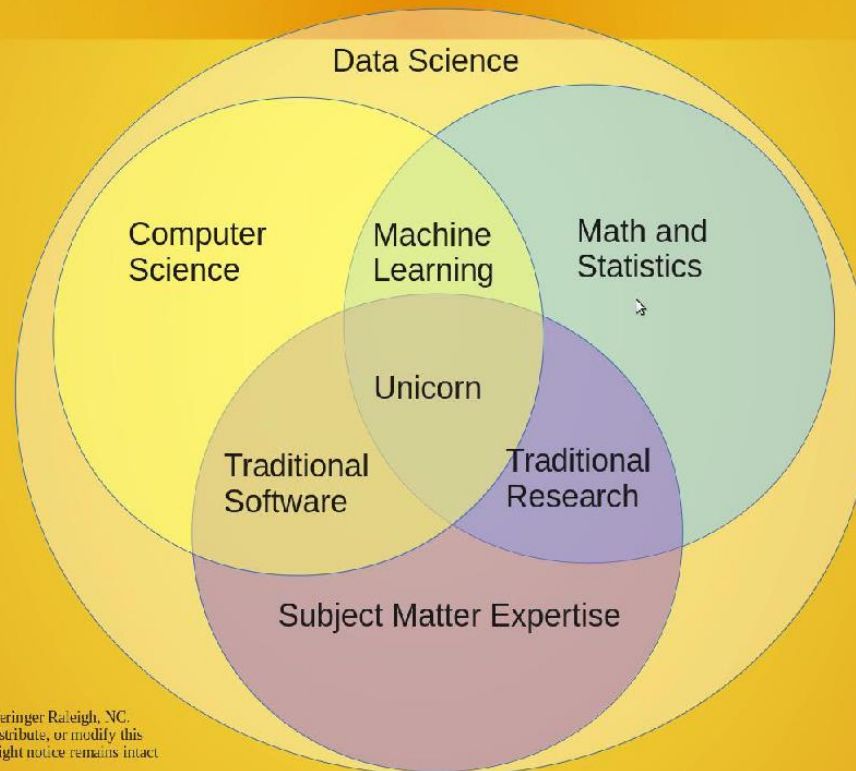
©CloudTweaks.com

Nauka o danych

- Interdyscyplinarna dziedzina, która próbuje wyodrębnić wzorce, spostrzeżenia i wiedzę z ustrukturyzowanych, częściowo ustrukturyzowanych i nieustrukturyzowanych danych
- Eksploracja danych
- Odkrywanie wiedzy w bazach danych
- Nauczanie maszynowe
- Rozpoznawanie wzorców
- Analiza predykcyjna
- Analiza Big Data

Diagram analizy danych

Data Science Venn Diagram v2.0



Copyright © 2014 by Steven Geringer Raleigh, NC.
Permission is granted to use, distribute, or modify this
image, provided that this copyright notice remains intact.

Analiza danych

- Liczba zastosowań dla analizy danych wzrosła o 650% od 2012 roku
- Według amerykańskiego Biura Statystyki Pracy do 2026 r. zostanie utworzonych około 11,5 miliona miejsc pracy
- Ponadto stanowisko Data Scientist należy do najlepszych wschodzących miejsc pracy w LinkedIn
- Wszystkie statystyki wskazują na rosnące zapotrzebowanie na naukowców zajmujących się danymi

Big Data i AI środowisko 2020

DATA & AI LANDSCAPE 2020

INFRASTRUCTURE STORAGE HADOOP CLOUDERA DATA LAKES DATA WAREHOUSES STREAMING / IN-MEMORY				ANALYTICS & MACHINE INTELLIGENCE BI PLATFORMS VISUALIZATION DATA ANALYST PLATFORMS DATA SCIENCE NOTEBOOKS DATA SCIENCE PLATFORMS MACHINE LEARNING				APPLICATIONS - ENTERPRISE SALES MARKETING - B2B MARKETING - B2C CUSTOMER EXPERIENCE / SERVICE HUMAN CAPITAL			
NOSQL DATABASES NEWSQL DATABASES GRAPH DBS MPP DBS SERVER-LESS CLUSTER DBS				COMPUTER VISION HORIZONTAL AI SPEECH & NLP				LEGAL RETECH & COMPLIANCE FINANCE FINANCE - B2B FINANCE - B2C FINANCE - INVESTING INSURANCE			
ETL / DATA TRANSFORMATION DATA INTEGRATION DATA GOVERNANCE DATA QUALITY				SEARCH LOG ANALYTICS SOCIAL ANALYTICS WEB / MOBILE / COMMERCE ANALYTICS				ADVERTISING EDUCATION REAL ESTATE GOVT & INTELLIGENCE COMMERCE FINANCE - LENDING			
MGMT / MONITORING DATA GENERATION & LABELLING AI OPS GPU DBS & CLOUD AI HARDWARE				OPEN SOURCE FRAMEWORKS QUERY / DATA FLOW DATA ACCESS & DATABASES ORCHESTRATION & PIPELINES STREAMING & MESSAGING STAT TOOLS & LANGUAGES AI OPS & INFRA AI / MACHINE LEARNING / DEEP LEARNING				HEALTHCARE LIFE SCIENCES TRANSPORTATION AGRICULTURE INDUSTRIAL OTHER			
DATA MARKETPLACES & DISCOVERY FINANCIAL & ECONOMIC DATA AIR / SPACE / SEA PEOPLE / ENTITIES LOCATION INTELLIGENCE OTHER				DATA SOURCES & APIs DATA SERVICES DATA RESOURCES INCUBATORS & SCHOOLS RESEARCH							

Version 1.0 - September 2020

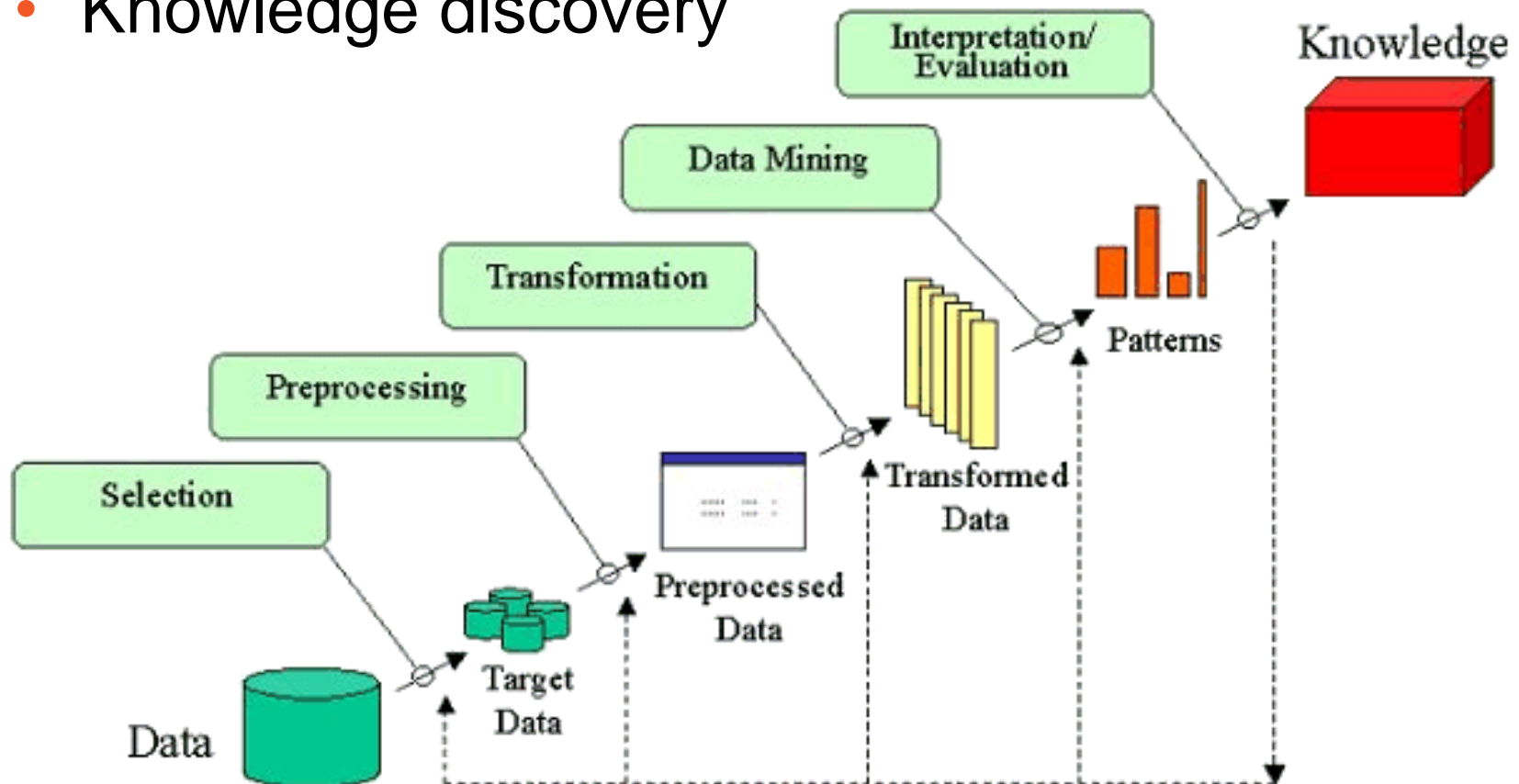
© Matt Turck (@matturck) & FirstMark (@firstmarkcap) matturck.com/data2020

FIRSTMARK
EARLY STAGE VENTURE CAPITAL



Proces przetwarzania danych

- Knowledge discovery



Prywatność i bezpieczeństwo danych

- Prywatność danych
- Bezpieczeństwo danych
- Integralność danych
- Analiza pod kątem bezpieczeństwa

Bibliografia

- Nathan Marz, James Warren, Big Data Principles and best practices of scalable real-time data systems. Manning Publications Co., 2015.
- Martin Kleppmann, Designing Data-Intensive Applications: The Big Ideas Behind Reliable, Scalable, and Maintainable Systems, O'Reilly Media, 2017.
- John D. Kelleher, Brendan Tierney, Data Science, The MIT Press, 2018
- Joel Grus, Data Science from Scratch: First Principles with Python, 2nd Edition, O'Reilly Media, 2019