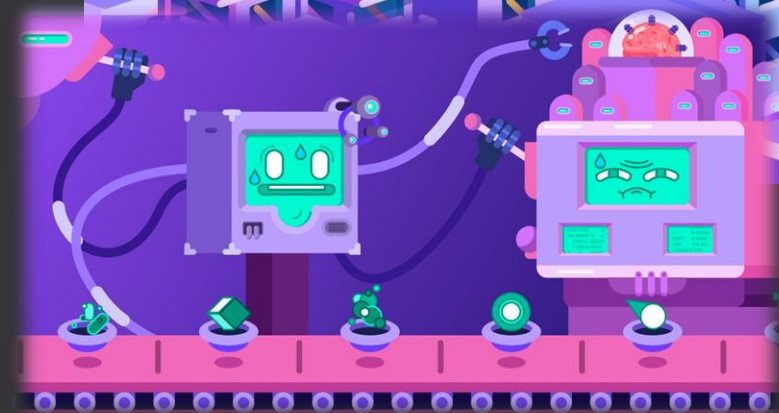


Metody i narzędzia *Big Data*

Modele przetwarzania danych



Co to znaczy, że dane są
duże?



Co to znaczy, że dane są **duże**?

Objętość

nie da się ich załadować na
pojedynczą maszynę i przetwarzać



Duże
dane

Duża objętość

Przykład

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



Co to znaczy, że dane są **duże**?

Objętość

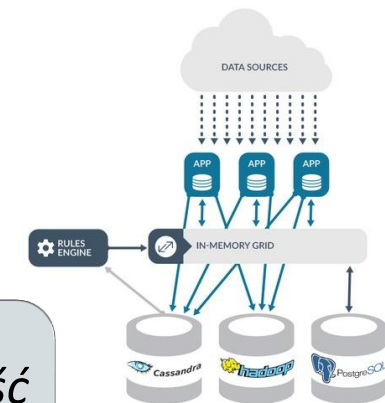
nie da się ich załadować na pojedynczą maszynę i przetwarzać



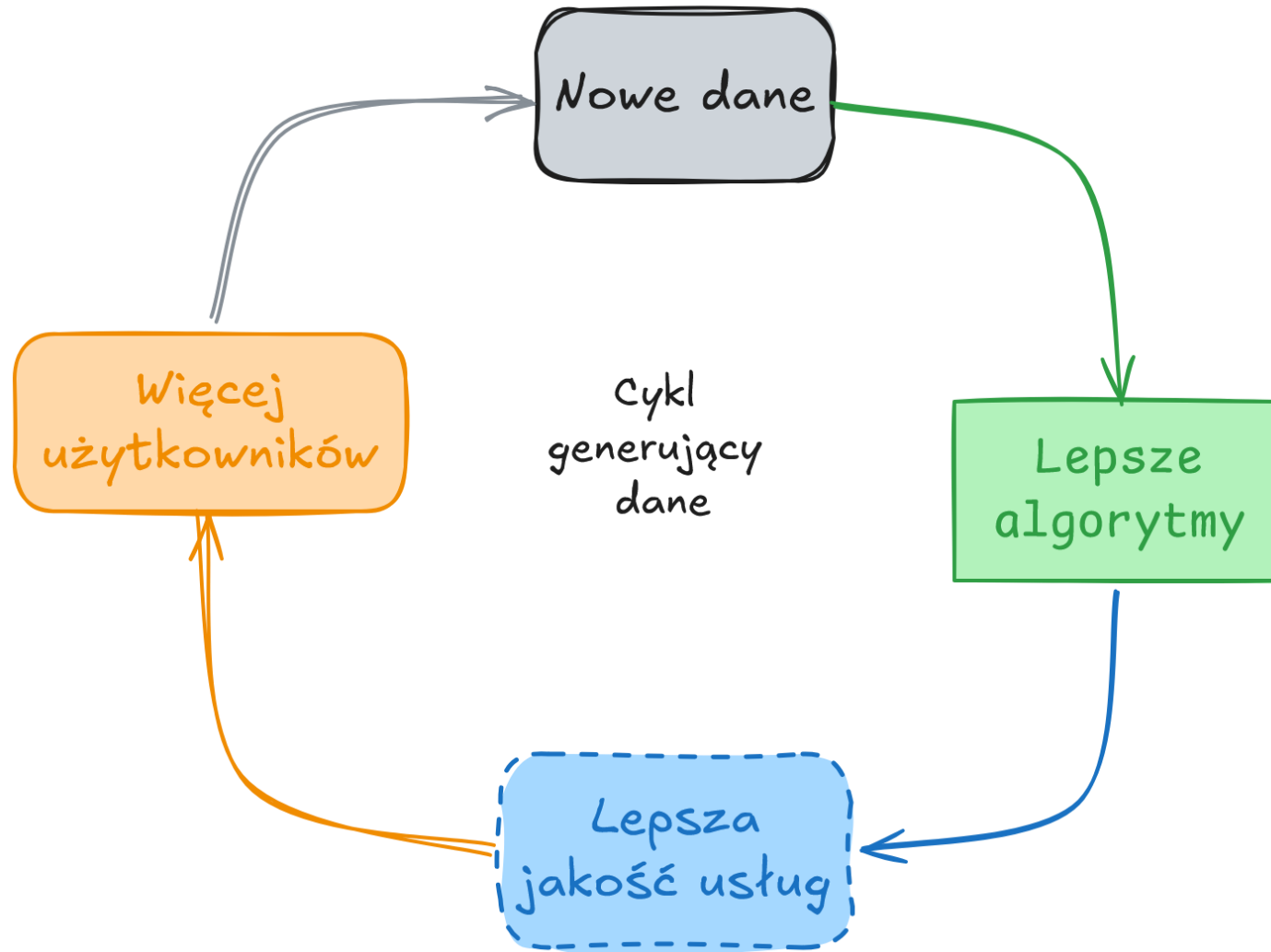
Duże dane



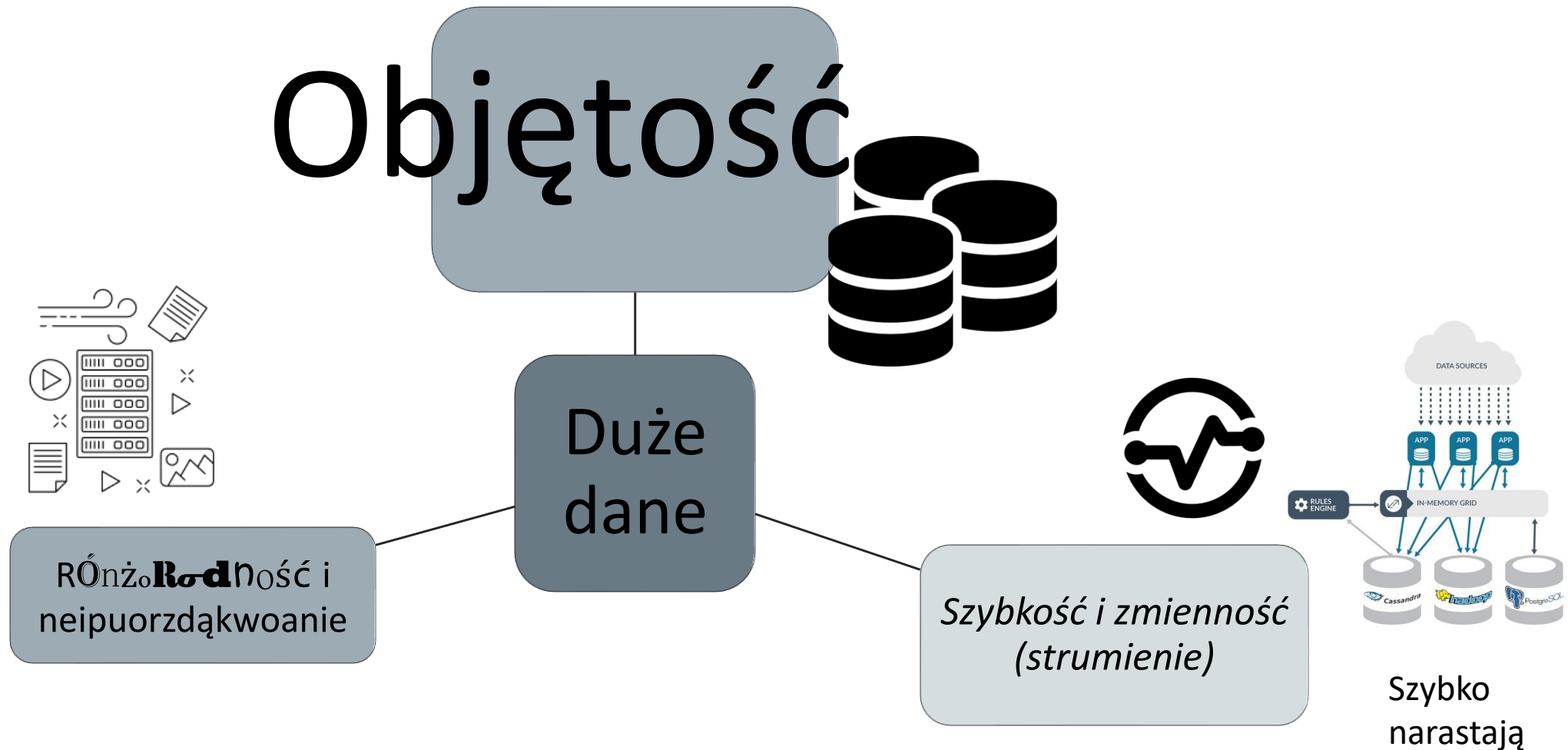
Szybkość i zmienność
(strumienie)



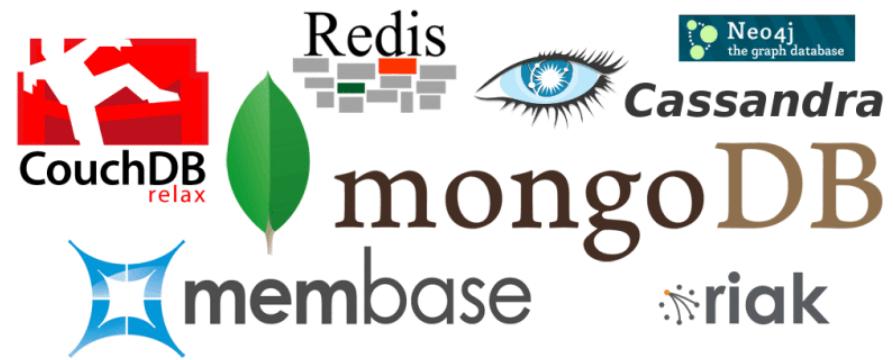
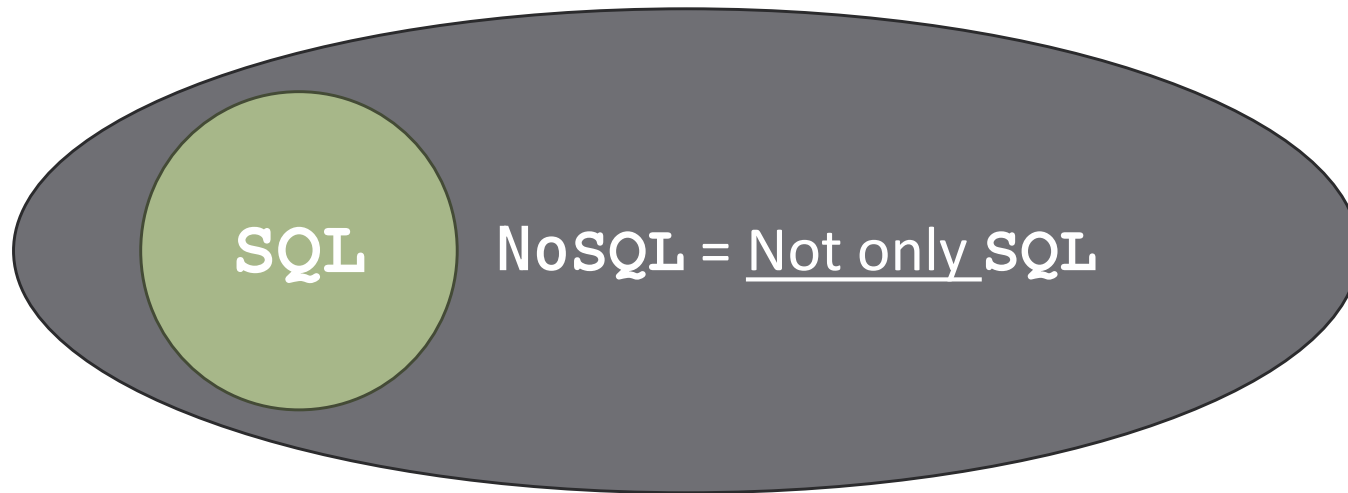
Szybko
narastają



Co to znaczy, że dane są **duże**?



Duży bałagan (ang. *complexity*)



Modele 3V i 5V

Volume



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

Velocity



Jak **szybko** nadchodzą nowe dane?

Variety



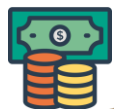
Jaki jest format / struktura / **źródła** danych?

Veracity



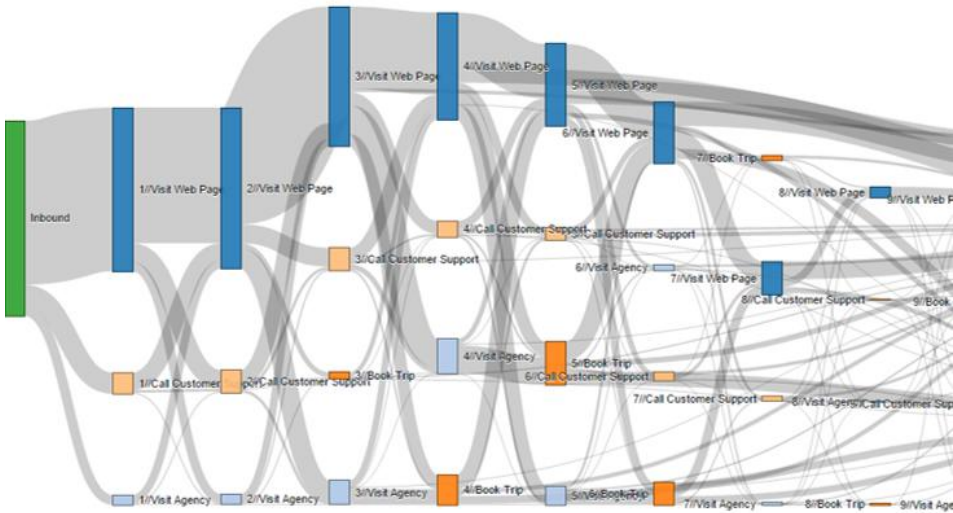
Czy dane są **wiarygodne** / kompletne / spójne?

Value

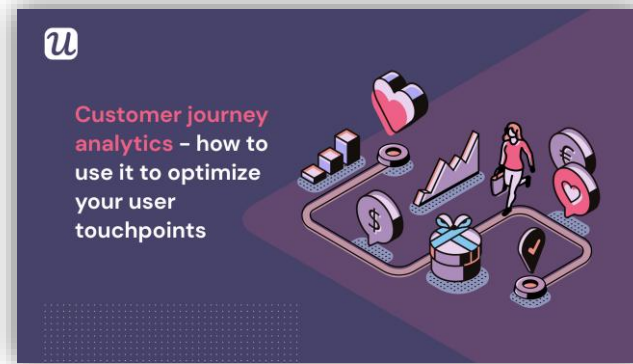


Czy **koszt** przetwarzania danych **się zwróci**?

Ścieżki zakupowe

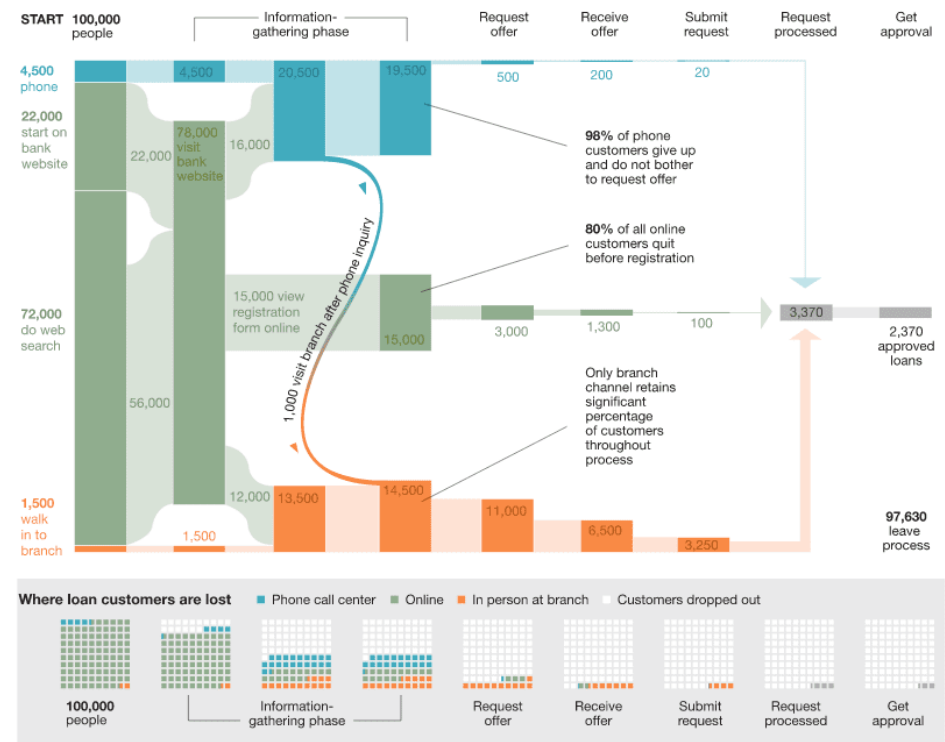


- Porzucone koszyki
- Zakupy krzyżowe, czego nie kupują
- Prognozowanie zamówień – stany magazynowe
- Segmentacja klientów, profilowanie (jakich filtrów używa)
- Odejście klienta – skąd pozyskany, niezbalansowane klasy
- Deduplikacja klientów (różne loginy Neo4j)
- Wielkość sprzedaży towaru – istotność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



Mapping customer flows highlights important pain points.

Average monthly customer flows for loan products by channel,¹ indexed to 100,000



¹Preapproved loans excluded.

Kiedyś to było...




Numer	Tytuł	Hasło	Wypoż.	Zwrot	Rok wyd.
8	Psychologia wychowaw...	Wiebzydowski, Leon (...)	2010.08.12	2010.08.26	1987
1000	Gwiazdzbior muzyczny	Kisielewski, Stefan (19 ...)	2010.08.12	2010.08.27	1982
1 / 2009	Architektura Murator	Architektura...	2010.08.25	2010.08.26	2009
29 / 1993	American Studies Newsl ...	American...	2010.08.25	2010.08.26	1993

LIBRA2000

⚠ Czytelnik nie zapłacił kary (10,00 PLN).

OK



Wypożyczone 4 pozycje

Karta wypożyczeń

Obudź wyjm. RZ RW Wyp. cz. Termin Wyp./Zwr. Zamknij

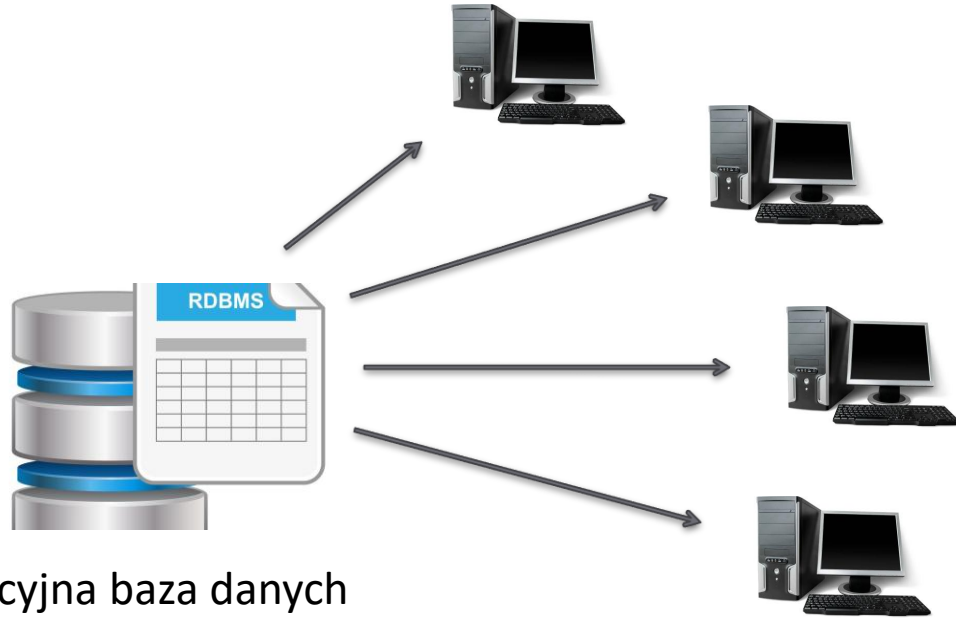


Relacyjna baza danych



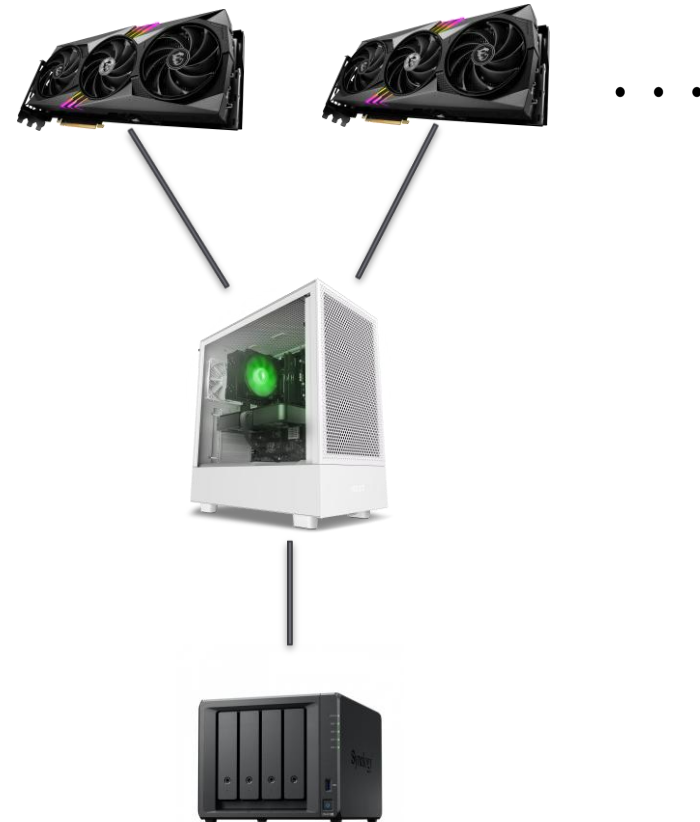
Przetwarzanie informacji

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

- Większa złożoność zarządzania
- Łatwiej uzyskać niezawodność przez replikację
- Trudno zapewnić wszystkie warunki ACID dla transakcji rozproszonych
Atomicity, Consistency, Isolation, Durability
- Warunki BASE
Basically Available, Soft state, Eventually consistent
- Dostęp do danych wymaga myślenia o partycjonowaniu, replikacji i opóźnieniach

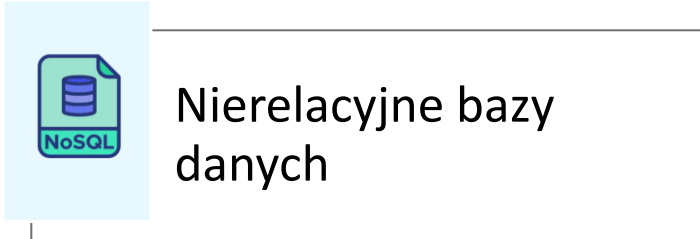
Rozproszony system obliczeniowy



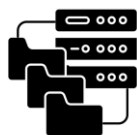
GFS



Rozproszony system obliczeniowy



Rozproszony system obliczeniowy



Rozproszony system plików



Nierelacyjne bazy danych



Platforma przetwarzania rozproszonego

GFS

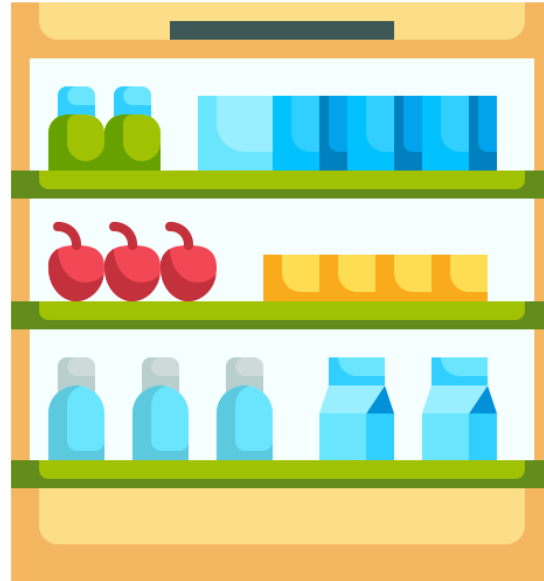
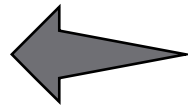


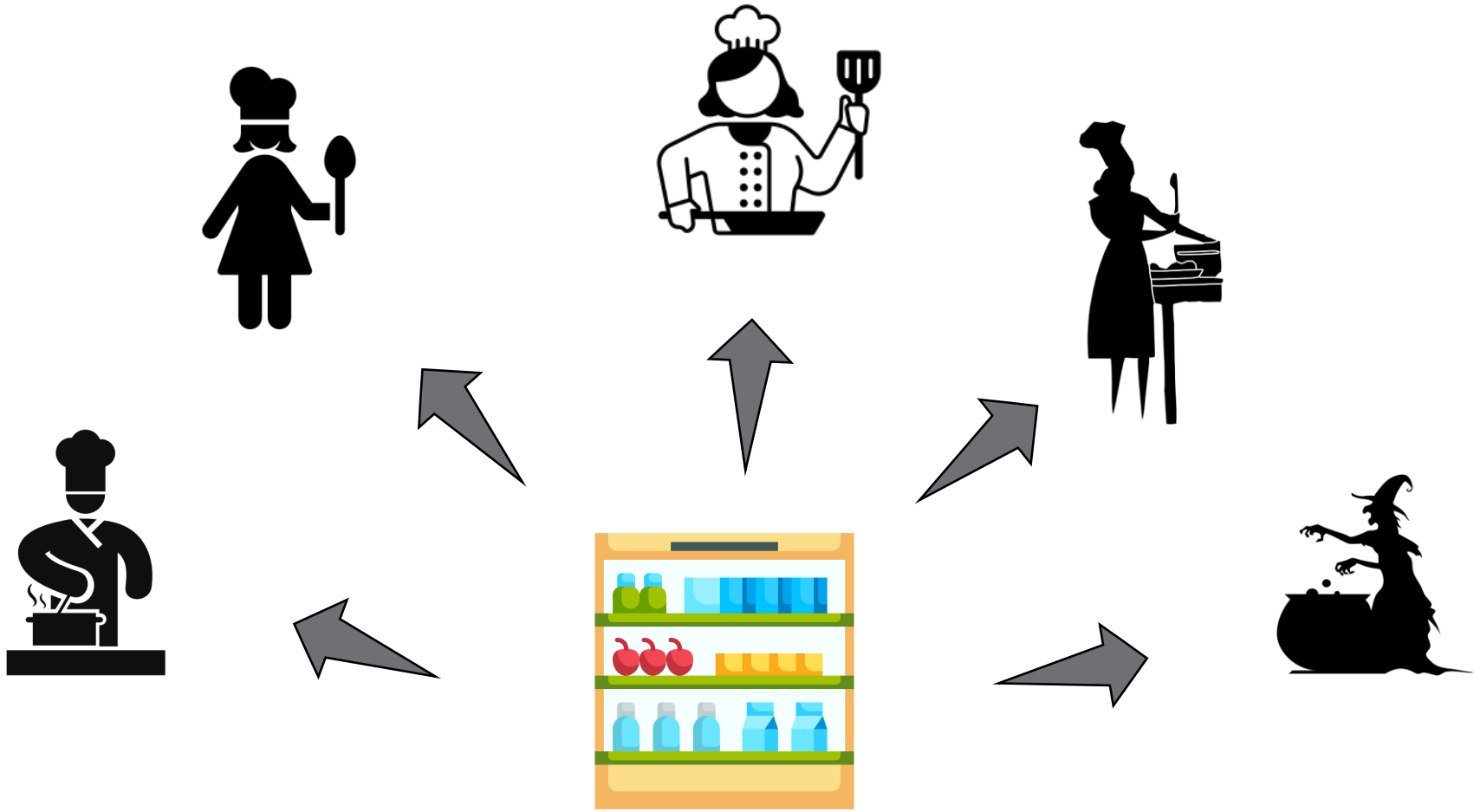
mongo DB

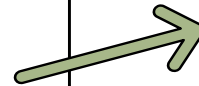
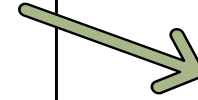
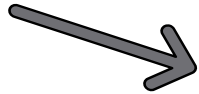
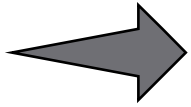
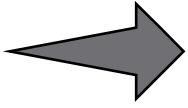
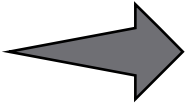
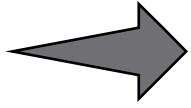
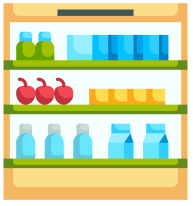


redis





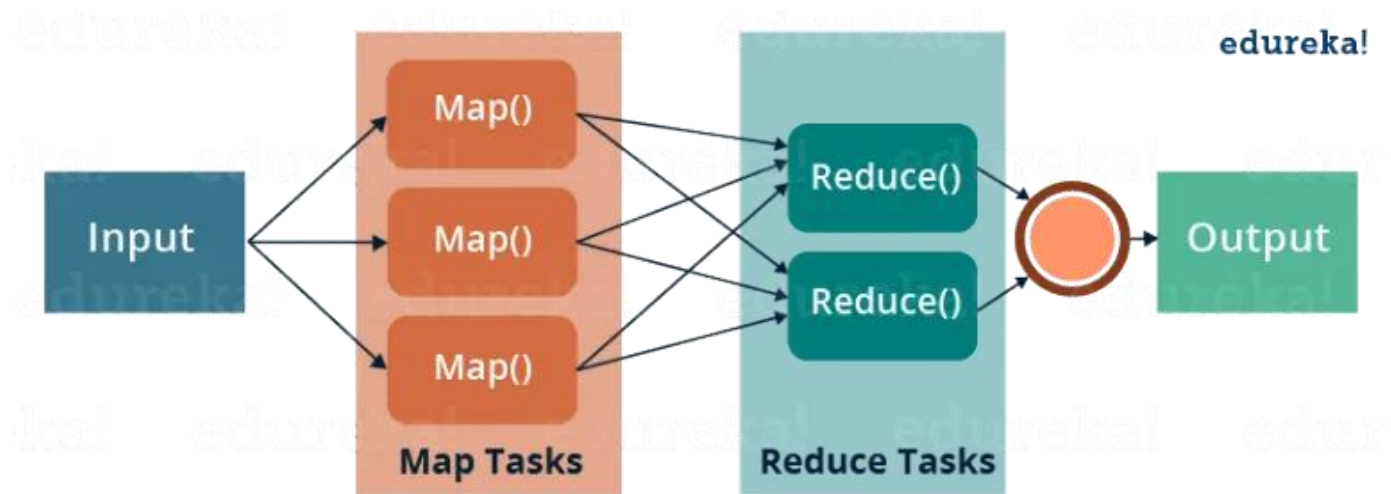




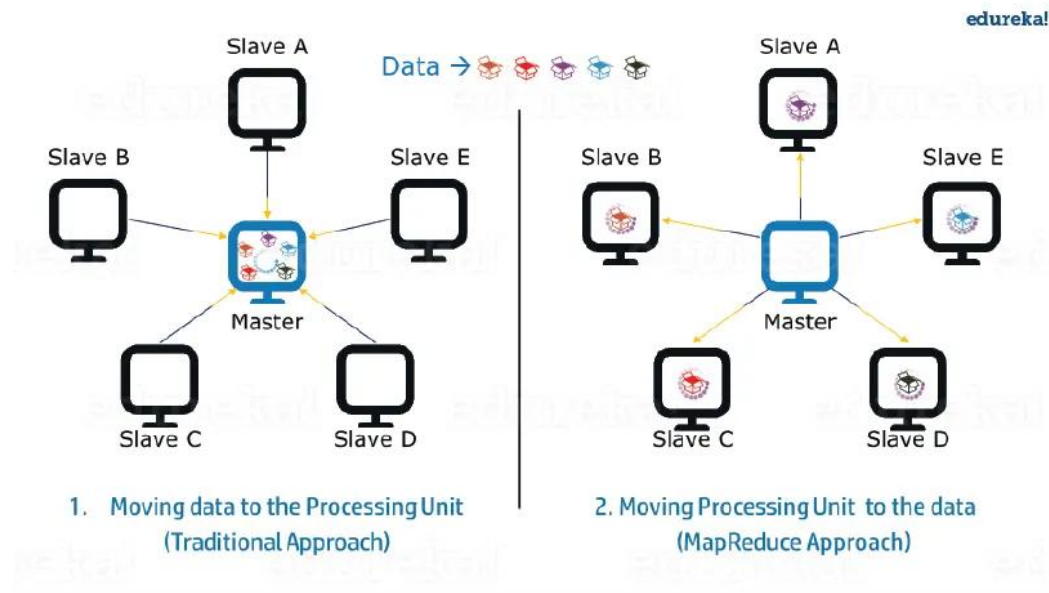
MAP

REDUCE

Model przetwarzania MapReduce



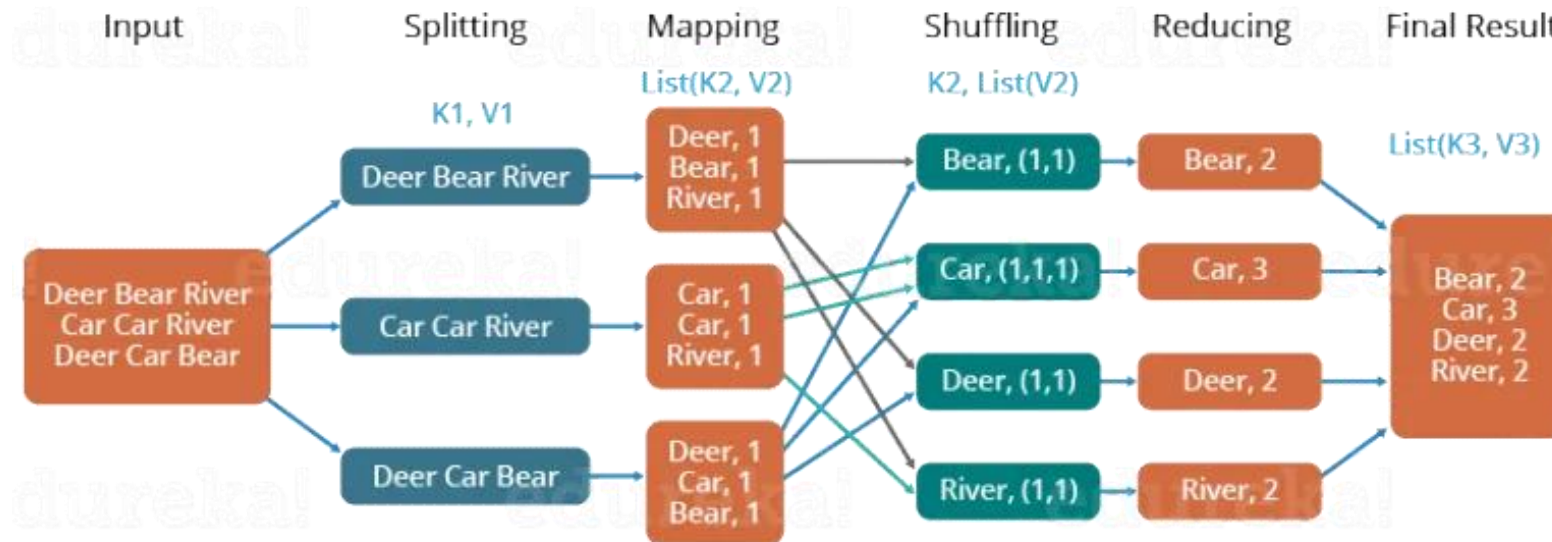
Model przetwarzania MapReduce



Model przetwarzania MapReduce

The Overall MapReduce Word Count Process

edureka!





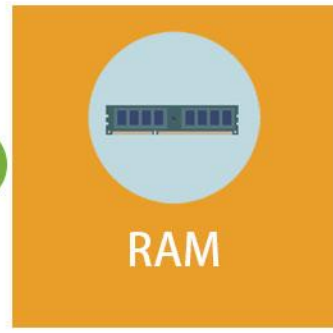
Koszt
filtrowania i
usuwania



Koszt
przechowywania
bezużytecznych
danych



vs





- Dane pośrednie przechowywane **na dyskach**
- Przetwarzanie wsadowe



Nieefektywny w zadaniach iteracyjnych i interakcyjnych



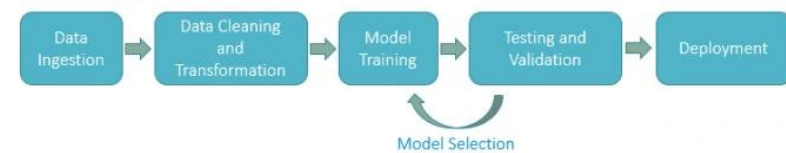
Uczenie maszynowe



Przetwarzanie strumieni



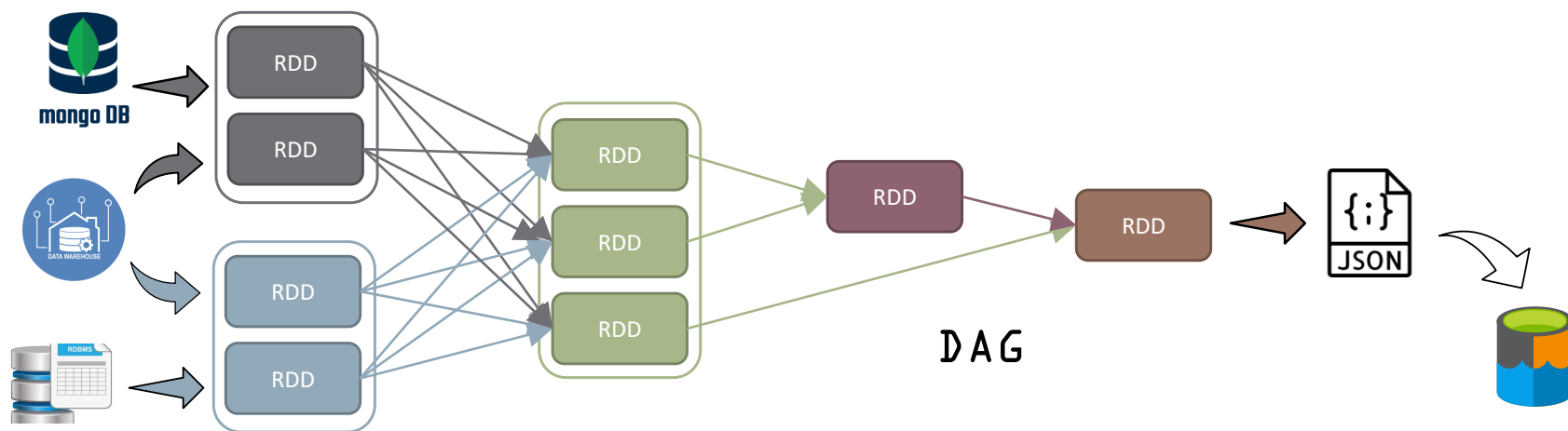
- Dane pośrednie mogą być przechowywane **w pamięci**
- Szybsze iteracje



Otwarta platforma do szybkiego, rozproszonego przetwarzania danych, także **strumieniowych**



Model przetwarzania w *Spark*

- **RDD** (*Resilient Distributed Dataset*) reprezentuje rozproszoną kolekcję danych oraz przepis na ich odtworzenie (ang. *lineage*)
- RDD może być większy od dostępnej pamięci
- Leniwa ewaluacja – DAG jest tworzony przy użyciu *transformacji*, obliczenia wykonuje się *akcjami*



Wybrane kierunki rozwoju



- Łańcuch bloków (ang. *blockchain*) – model rozproszonego rejestru danych ze współużytkowaniem, gdzie użytkownik może mieć większą kontrolę nad własnymi danymi
- *Interplanetary File System* () – rozproszony system przechowywania i **adresowania treści**: sieć peer-to-peer z systemem nazw IPNS oraz wersjonowaniem plików
-  **Filecoin** – wynagradzanie za przechowywanie i udostępnianie przestrzeni dyskowej



THE END