

DĚLAT  
DOBRÝ SOFTWARE  
NÁS BAVÍ

PROFINIT

**NDBI047**

# Aplikace Big Data technologií v Data Science

Petr Paščenko

22. 2. 2019

The background of the slide is composed of numerous overlapping, translucent, light-gray geometric shapes. These shapes, which include rectangles, squares, and irregular polygons, are scattered across the white background, creating a complex, layered, and three-dimensional effect. The shapes vary in size and orientation, some appearing to float above others, which adds depth to the overall design.

Co?

# Cíle předmětu

- › Představit technologie pro zpracování velkých dat
  - Architektonický pohled
    - systém jako celek a jeho součásti
  - Uživatelský pohled
    - jednotlivé technologie a jak je správně používat
  - Programátorský pohled
    - jak vyvíjet aplikace v prostředí velkých dat
- › Data Science aplikace
  - Jaké úlohy můžeme počítat, když máme cluster
- › Ukázat skutečný život na skutečných projektech
- › Těžiště na seminární práci
  - Předmět formou Scientific Lab
    - Jak analyzovat velká data
    - Jak sestavit modely
    - Jak napsat výzkumnou zprávu



The background of the slide is composed of numerous overlapping, translucent, light-gray geometric shapes. These shapes, which include rectangles, parallelograms, and trapezoids, are oriented in various directions, creating a complex, layered, and three-dimensional effect. The shapes are semi-transparent, allowing the ones beneath them to be visible, which adds depth to the overall composition.

Jak?

# Osnova předmětu

- › Úvod, přehled, organizace, motivace
- › Hadoop Cluster
  - Storage
    - HDFS, formáty ukládání a komprese dat, HIVE
  - Map-reduce
    - softwarové paradigma a implementace algoritmů, vztah k sql
  - Apache Spark
    - Distribuované výpočty v RAM
- › Data Science metody a nástroje
  - Python Technology Stack
    - Pandas, NumPy, Scikit-learn, PySpark, matplotlib
  - Analýza velkých dat
    - Analytické notebooky, explorace, vizualizace, statistická analýza, report
- › Samostatná úloha
  - Analýza dat
  - Modelování a predikce
  - Dokumentace formou výzkumné zprávy



# Za co bude známka

- › **Samostatně zpracovaná výzkumná zpráva**
  - Malá bakalářka
  - Trocha teorie, popis experimentů
  - Analýza dat
  - Vyhodnocení experimentů – predikce, výpočetní náročnost
  - Zdrojové kódy experimentů na githubu
  - Popis a vyhodnocení experimentů
- › **Odevzdání na dvě iterace**
  - První iterace
  - Připomínky
  - Druhá iterace
  - Hodnocení
- › **Kritéria hodnocení**

– Věcná správnost	50%
– Rozsah použitých metod	30%
– Forma	20%



The background of the slide is composed of numerous overlapping, translucent, light-gray geometric shapes. These shapes, which include rectangles, trapezoids, and irregular polygons, are scattered across the entire frame, creating a complex, layered, and three-dimensional effect. The shapes vary in size and orientation, with some appearing more prominent than others due to their position and transparency.

Kdo?

# Petr Paščenko

## › Zaměření

- Data Science a Machine Learning
- Big Data Byznys Aplikace

## › Praxe

- V Profinitu 6 let, konzultant, R&D
  - Data Science
    - vývoj vztahových a podobnostních modelů nad velkými daty
  - Grantový výzkum zaměřený na transakční mining
    - Webmining, SNA, text-mining, grafová analytika, zabezpečení online kanálů
  - Big Data
    - aplikace a školení
  - Vývoj Java
- ČEZ, 2 roky
  - Analytik trhů
- Eccam, 2 roky
  - Vývoj c++





# Jan Hučín

- › Zaměření
  - Data Science a Machine Learning
  - Big Data Byznys Aplikace
- › Praxe
  - V Profinitu 4 roky, konzultant
    - Data Science
      - klasifikační a prediktivní modely, metodika
    - Big Data
      - popularizace a školení
  - Scio
    - Psychometrika, datová analytika
  - Institut informatiky
    - Výuka matematiky a statistiky



# Stamenov Sergii

- › Zaměření
  - Data Science, Machine learning, Big Data
  - Vyvoj .NET
- › Praxe
  - V Profinitu 2 roky
  - Data science pro T-Mobile
  - Spark školení



# Dominik Matula

- › Zaměření
  - Data Science, Big Data a Machine Learning
  
- › Praxe
  - Profinit
    - 2 roky, konzultant
    - Data Science, Big Data
    - Podobnostní modelování, klasifikace,
  - Median
    - 2 roky, statistik
    - Analýza mediální konzumace; vývoj a optim. algoritmů pro detekci zásahů billboardy či poslechů reklamních spotů

The background of the slide is a complex, abstract composition of numerous overlapping, translucent geometric shapes. These shapes, which include various polygons and rectangles, are rendered in different shades of light gray and white. They are arranged in a way that creates a sense of depth and movement, as if they are floating or shifting in a three-dimensional space. The overall effect is a modern, minimalist aesthetic that suggests a digital or data-driven environment.

**Big Data?**

# Co jsou Big Data?



› Je to složité.

# Co jsou Big Data?

- › Technologie velkých dat jsou kulminací trendů v několika oblastech
- › Data
  - Od malých dat k velkým
  - od jednoduchých ke složitým
- › Databáze
  - Od souborů přes relační paradigma k distribuovaným clusterům
- › Programování
  - Od procedurálního programování
  - k funkčním frameworkům
- › Data Science
  - Od výběrových statistik
  - k detailní kontextové analýze
- › Postupně si je projdeme





The background of the slide is composed of numerous overlapping, translucent, light-gray geometric shapes. These shapes, which include rectangles, trapezoids, and irregular polygons, are scattered across the white background, creating a complex, layered, and three-dimensional effect. The shapes vary in size and orientation, some appearing to float above others, which adds depth to the overall design.

Data?



# A Very Short History Of (Big) Data

Gil Press pro [forbes](#)



- › **1944:** Fremont Rider, Wesleyan University Librarian
  - „**American university libraries were doubling in size every sixteen years.**“
  - „Yale Library in 2040 will have approximately 200,000,000 volumes, which will occupy over 6,000 miles of shelves“
- › **1961:** Derek Price, Science Since Babylon
  - „**the number of new journals has grown exponentially** rather than linearly, doubling every fifteen years and increasing by a factor of ten during every half-century.“
  - „Each scientific advance generates a new series of advances at a reasonably **constant birth rate**, so that the number of births is strictly proportional to the size of the population of discoveries at any given time. “
- › **1975:** The Ministry of Posts and Telecommunications in Japan
  - „**information supply is increasing much faster than information consumption**“
  - „the demand for information provided by mass media, which are one-way communication, has become stagnant, and the demand for information provided by personal telecommunications media, which are characterized by **two-way communications**, has **drastically increased.**“

# A Very Short History Of Big Data

Gil Press pro [forbes](#)

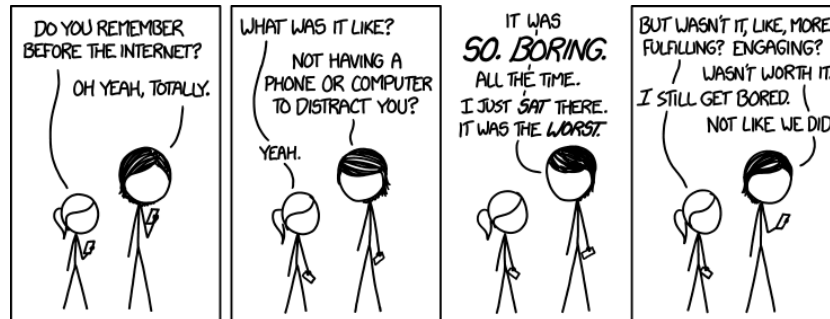


- › **1980:** I.A. Tjomsland, Fourth IEEE Symposium on Mass Storage Systems
  - „Parkinson’s 1st Law paraphrased: **Data expands to fill the space available.**“
  - „The penalties for storing obsolete data are less apparent than are the penalties for discarding potentially useful data.“
- › **1986:** Hal B. Becker, Can users really absorb data at today’s rates?
  - „The recoding density achieved by Gutenberg was approximately 500 symbols per cubic inch – 500 times the **density** of [4,000 B.C. Sumerian] clay tablets. By the year 2000, semiconductor random access memory should be storing  $1.25 \times 10^{11}$  bytes per cubic inch.“ – po pravdě o řád přestřelené v roce 2017
- › **1996:** B.J. Truskowski, The Evolution of Storage Systems
  - **Digital storage becomes more cost-effective for storing data than paper.**
- › **1997** Michael Cox and David Ellsworth
  - „Data sets are generally quite large, taxing the capacities of main memory, local disk, and even remote disk. We call this the problem of **big data**. When data sets do not fit in main memory (*in core*), or when they do not fit even on local disk, the most common solution is to **acquire more resources**. “

# A Very Short History Of Big Data

Gil Press pro [forbes](#)

- › **1997** Michael Lesk, How much information is there in the world?
  - „In only a few years, (a) we will be able [to] save everything—no information will have to be thrown out, and (b) **the typical piece of information will never be looked at by a human being.**“
- › **1998:** K.G. Coffman and Andrew Odlyzko
  - „the growth rate of traffic on the public Internet, while lower than is often cited, is still about 100% per year, much higher than for traffic on other networks.“



## › **2000: Big Data Era**

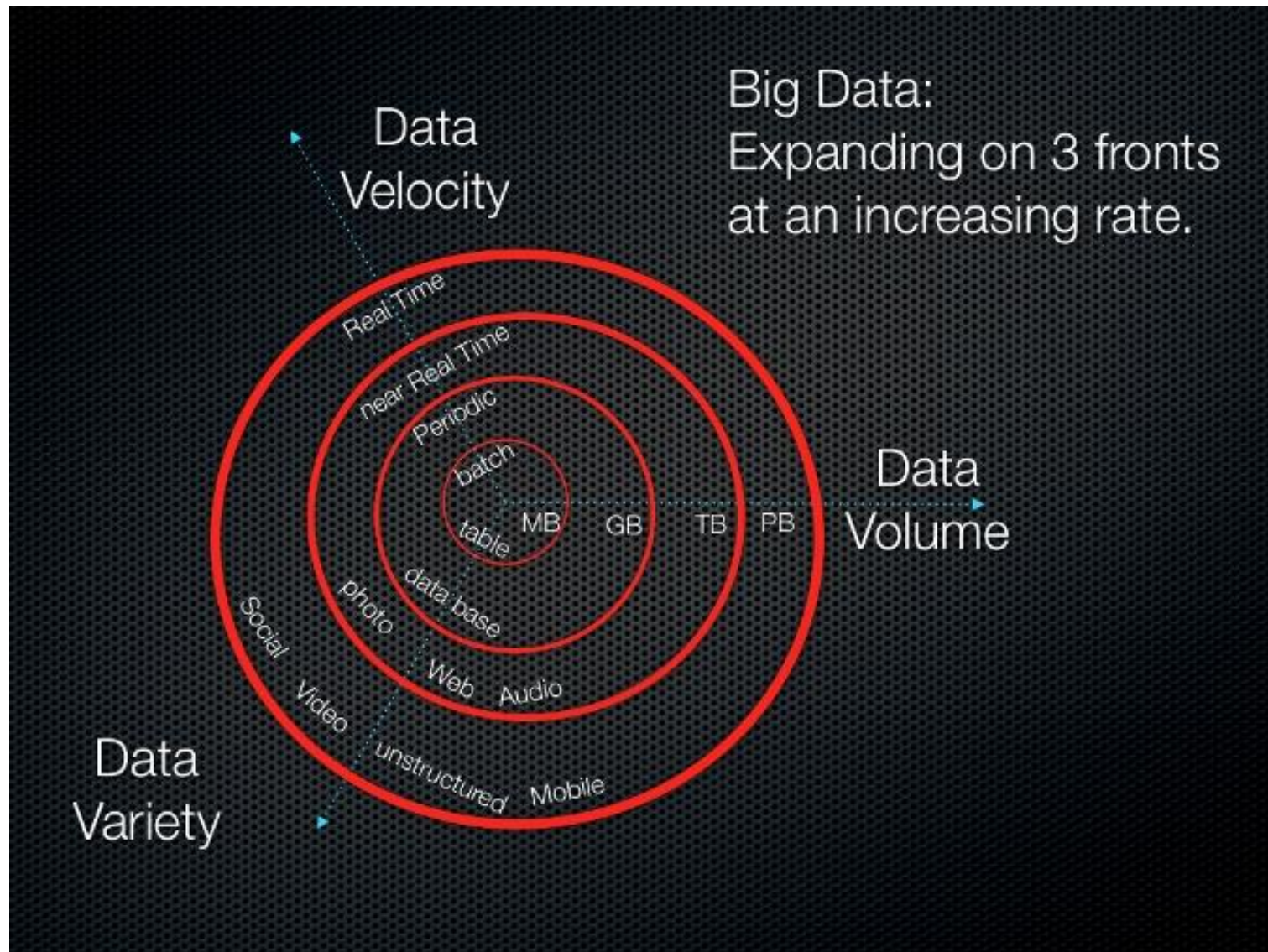
- „the world produced about 1.5 exabytes of unique information, or about **250 megabytes for every man, woman, and child on earth.** It also finds that “a vast amount of unique information is created and stored by individuals” (what it calls the “democratization of data”) and that “not only is digital information production the largest in total, it is also the most rapidly growing.“

# A Very Short History Of Big Data – Shrnutí

- › Množství dat exponenciálně roste, stejně jako přenosová kapacita
- › Nabídka dat roste rychleji než poptávka
- › V komunikaci převládá obousměrnost a aktivní role uživatelů
- › Je levnější data skladovat než je třídit a vyhazovat
- › Opačná perspektiva: roste úložná kapacita a data ji jen celou zaplňují – poprvé můžeme uložit (skoro) všechny informace.
- › Datové soubory překračují hranici jednoho zařízení / média
- › Průměrná informace není nikdy přečtena člověkem



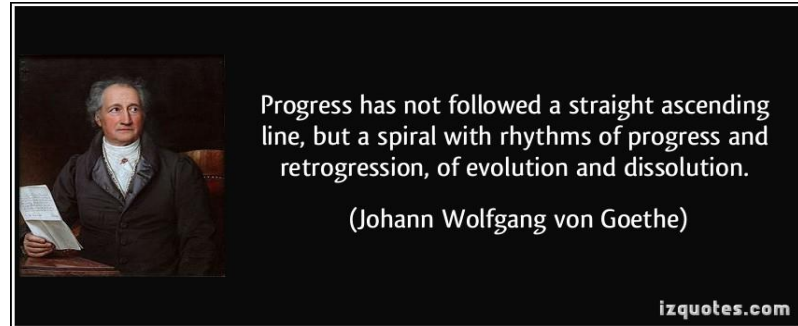
# Big Data Definice – 3V



The background of the slide is a complex, abstract composition of numerous overlapping, translucent geometric shapes. These shapes, which include rectangles, parallelograms, and polygons of various sizes and orientations, are rendered in different shades of light gray. They are layered in a way that creates a sense of depth and movement, with some shapes appearing to float above others. The overall effect is a modern, architectural, and somewhat crystalline aesthetic.

Databáze?

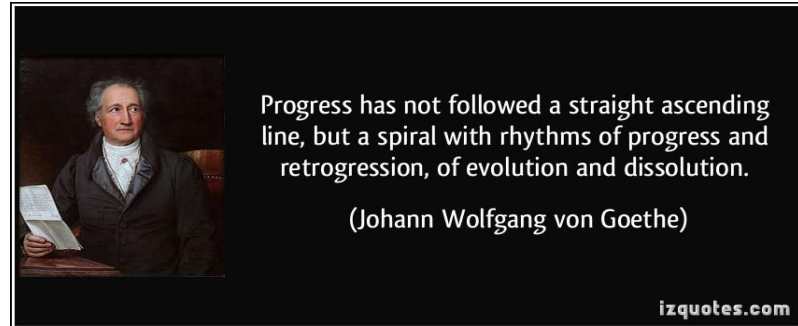
# Prehistorie pojmu Databáze



- › Původně – datová základna, slovník faktů
- › **40. léta – Ad hoc proudy dat**
  - Děrné štítky a pásky
- › **50. léta – File System**
  - Soubory a složky s hierarchickou strukturou a unikátní cestou
  - Nezávislost na použitém médiu (páska, disk, ...)
- › **60. léta – DBMS**
  - Tabulky a Indexování: hashování, B-stromy. Klient-server architektura
- › **70. léta – R-DBMS**
  - Relační paradigma: normální formy, relační algebra, selekce, projekce atd.



# Historie pojmu Databáze



- › **80. léta – SQL**
  - Jednotný dotazovací jazyk, rozvoj proprietálních databází
- › **90. léta – Datové sklady**
  - Datová Integrace, jedna pravda, analytický reporting.
- › **0. léta – NoSQL**
  - Webové programování, rozvolňování normálních forem, grafové databáze, cloud – částečný návrat ke vzdáleným architekturám
- › **10. léta – Big Data**

# RDMS vs Big Data

	Relační DB	Hadoop
Velikost	GB, TB	TB, EB, PB
Přístup	interaktivní i batch	jen batch
Dotazy	SQL a program. nastavby	Map-Reduce a SQL emulace
Úpravy	opakované čtení i zápis	zápis jednou, čtení opakovaně
Struktura	statické databázové schéma	dynamické schéma – volba při načtení
Integrita	ACID	není, durability pomocí redundance
Výkon	limitovaný shora optimalizace na straně DB	lineární škálování zrychlování přidáváním výkonu
Latence	minimální (ms)	značná (desítky sekund i více)
HW	špičkově vyladěné stroje	běžný hardware
Licence	komerční, velmi drahé	open source + podpora
Paralelizace	limitovaná a extrémě drahá cena per jádro	základní kámen architektury

The background of the slide is a complex, abstract composition of numerous overlapping, translucent geometric shapes. These shapes, which include various polygons and rectangular prisms, are rendered in different shades of light gray and white. They are arranged in a way that creates a sense of depth and three-dimensional space, with some shapes appearing to float above others. The overall effect is a modern, architectural, and somewhat crystalline aesthetic.

Programování?

# Programování – stručná historie aneb stále větší abstrakce



- › Procedurální – Imperativní programování
  - Programátor tvoří program tím, že přímo specifikuje posloupnost příkazů
- › **Nestrukturované paradigma** (cca do 60.let)
  - Asemblerové instrukce a podmíněné skoky – goto era
- › **Strukturované paradigma** (cca 60-80. léta)
  - Ústřední prvek programu jsou funkce sdružené do knihoven
- › **Objektové paradigma** (cca 90. léta)
  - Propojení funkcí a dat do objektů, dědičnost a polymorfismus
- › **Virtualizace** (cca 0. léta)
  - Oddělení programátora od specifického OS a HW
- › Nosná myšlenka
  - Zvyšování úrovně abstrakce (os, kompilátor, linker, vm)
  - Odstraňování complexity ležící dole – knihovny obsaťují nízkoúrovňové úlohy
- › Problém: komplexita neleží jen dole, ale i nahoře.

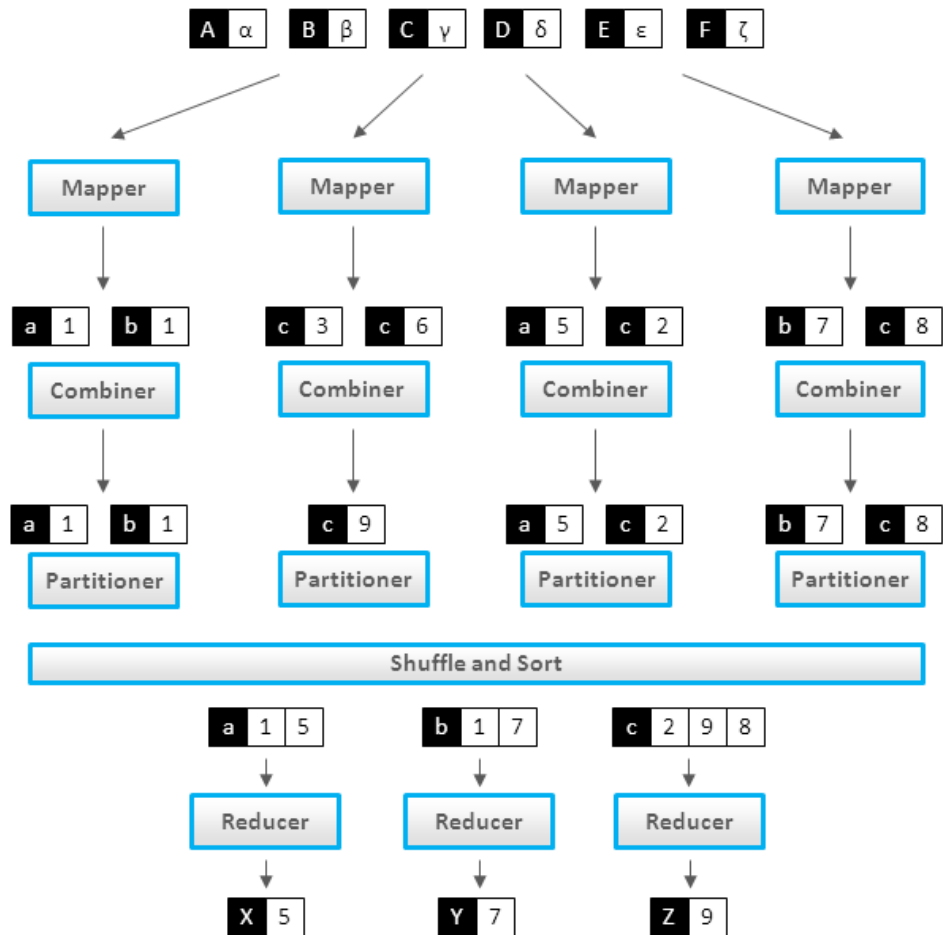
# Paralelní framework



- › Problém paralelního programu – komplexita přichází svrchu
  - Je potřeba formulovat algoritmus rovnou jako paralelní a to je velmi těžké
  - Ještě těžší je program udržovat, debugovat, optimalizovat a ladit
- › **Rozvoj frameworků** (současnost)
  - Framework je programová konstrukce popisující hrubou strukturu algoritmu
    - Například quick-sort s uživatelským komparátorem, struktura rozděl a panuj,...
  - Programátor tvoří program tím, že vkládá vlastní kód na předem připravená místa – obvykle se využívá dědičnosti, šablonování, lambda funkcí atd.
  - Odstínění programátora od komplexity problému, zejména paralelizace
  - Událostmi řízené programování, webové frameworky (apache tomcat), Tensorflow pro neuronové sítě, Map-Reduce paradigma atd.
- › Paralelní knihovny (MPI, BSP, OpenMP ...)
  - Univerzální frameworky mají plnou sílu, ale neřeší komplexitu algoritmu
- › **Map-Reduce** paradigma
  - Limitovaný rozsah funkcionality na zpracování dat
  - výměnou za jednoduchou funkcionální strukturu

# Map-Reduce

- › Schéma algoritmu dáno
- › Doplnujeme implementaci metod
- › Map
  - sestaví páry <key,val>
- › Combine\*
  - lze sloučit páry se stejným klíčem pro omezení síť. přenosu
- › Shuffle and Sort
  - Předun dat mezi uzly
- › Reduce
  - Agregace výsledků
- › Funkce mohou být obecně komplexní
- › Algoritmus = zřetězení Map a Reduce fází

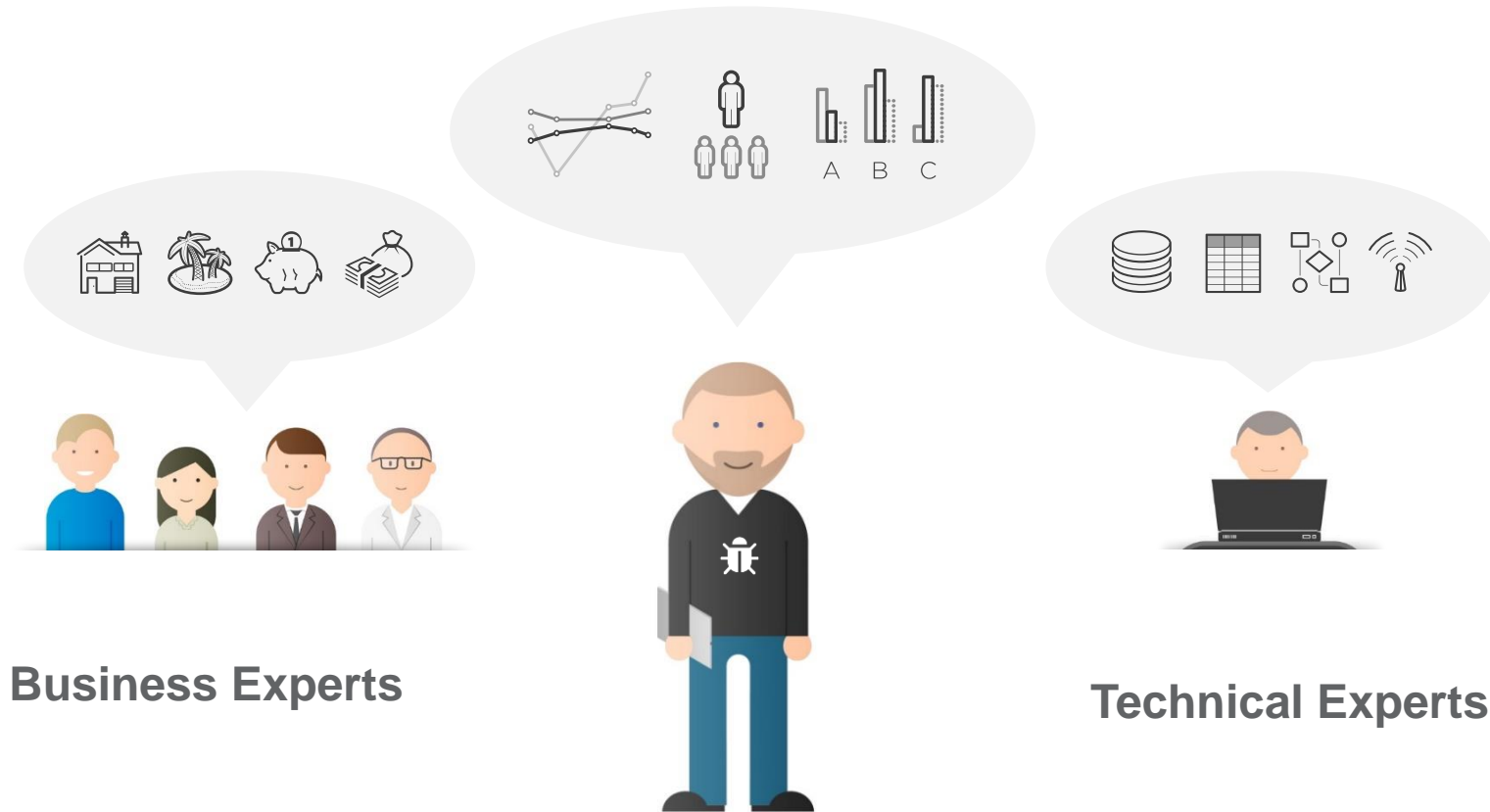


The background of the slide is a complex, abstract composition of numerous overlapping, translucent geometric shapes. These shapes, which include rectangles, squares, and polygons, are rendered in various shades of light gray and white. They are arranged in a way that creates a sense of depth and movement, as if they are floating or falling from the top right towards the bottom left. The overall effect is a modern, digital aesthetic that suggests data and technology.

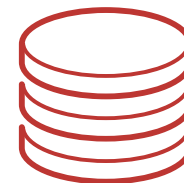
**Data Science?**



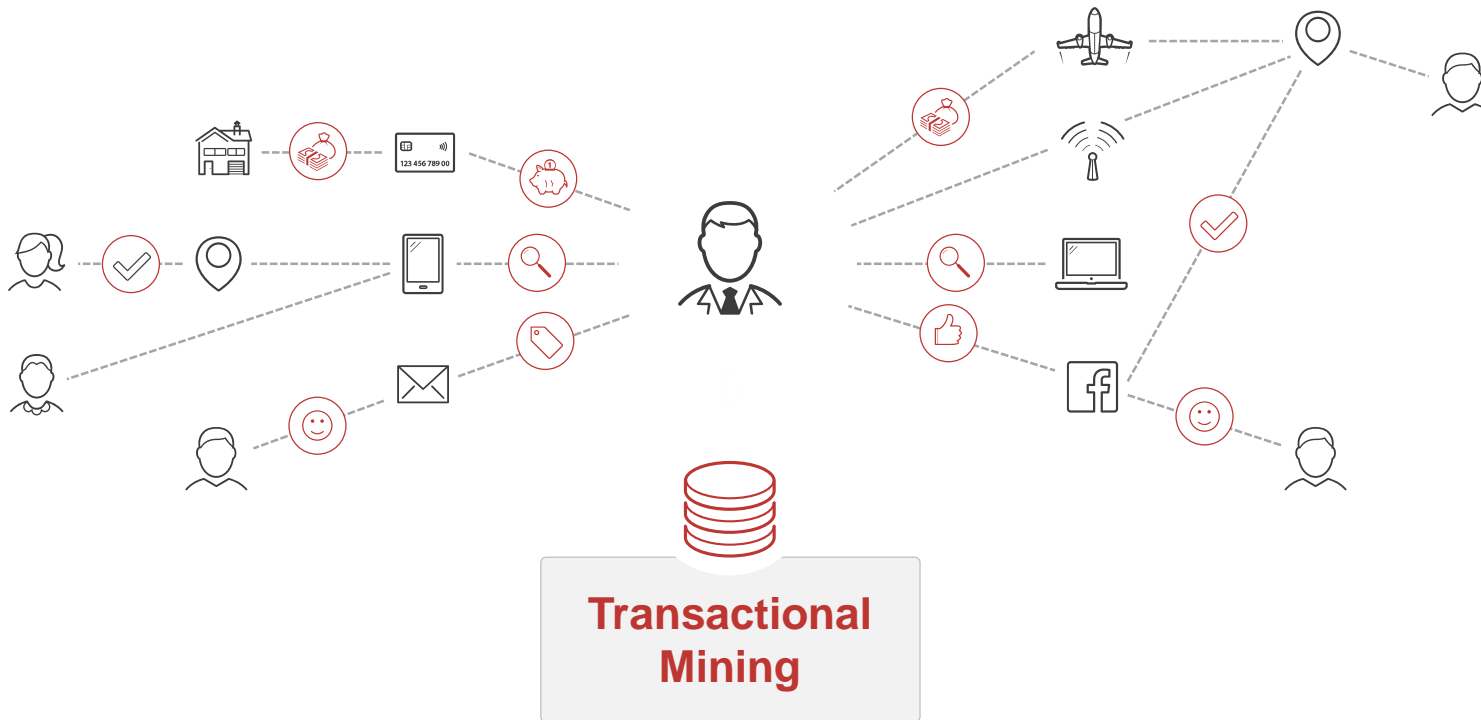
# Data Science in corporations



?



# Where do all the Big Data come from?



## Sources of Big Data

### 1. Technical – generated by machines

- Technical sensors, network traffic, audiovisual streams, astronomical, scientific,...

&

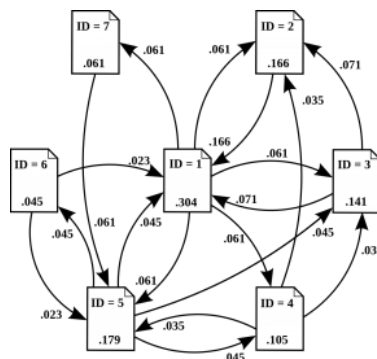
### 2. Transactional – generated by people

- Interactions with other objects (other people or things)
- Digital footprint

# Big Data Technologie

## › Pocházejí ze světa technologických gigantů

- Google
- Facebook
- Netflix
- Amazon



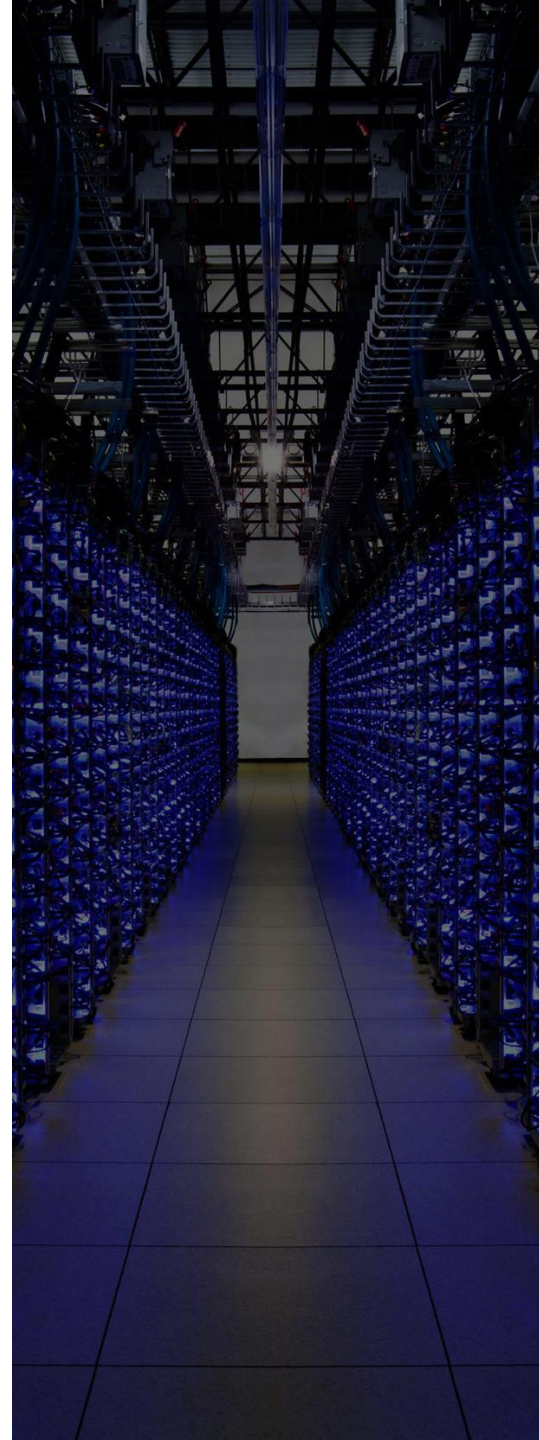
## › V současnosti se šíří do

- Bankovníctví
- Pojišťovnictví
- Telekomunikací
- Průmyslu
- Utilit

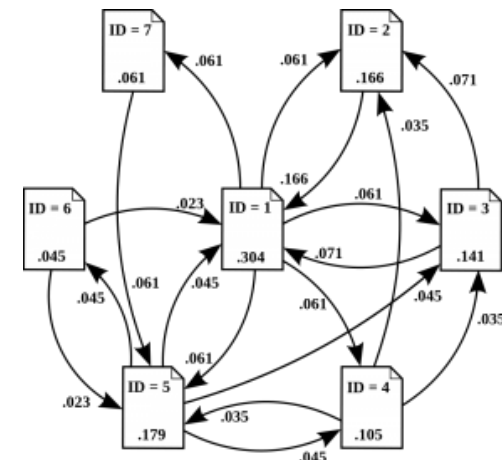
$$\mathbf{A} = \mathbf{U} \mathbf{L} \mathbf{V}^T$$

## › Naše mise

- Adaptace přístupů big data science ve financích



- › **Co je to Google?**
- › 1G stránek (webů) s řádově 10x více odkazy
- › Fulltextové vyhledávání
  - Roboti procházejí web a indexují obsah – inženýrská úloha
  - Uživatel zadá dotaz (slovo, frázi) a výsledkem je seznam stránek
- › V jakém pořadí zobrazit výsledky vyhledávání
  - Většinou je hledaný výraz hned ten první nebo na první stránce
  - Jak je to možné?
- › Google Pagerank
  - Iterativní metoda pro ohodnocení stránek
  - Stránky jsou uzly, odkazy jsou hrany
  - Markovovské řetězce s okrajovými podmínkami
  - Násobení řádkových matic o hraně 1G
  - Implementace v map-reduce



# Amazon, Netflix, YouTube

- › Doporučování obsahu
  - Viděl jste těchto deset filmů, podívejte se na jedenáctý
  - Kdo si koupil Babičku, ten si koupí Broučky

## › Dva přístupy

- Podobnost zboží
- Podobnost zákazníků

## › Kolaborativní filtrování

- Singular Value Decomposition
- Alternating Least Squares (Spark)

## › A – matice klient x produkt

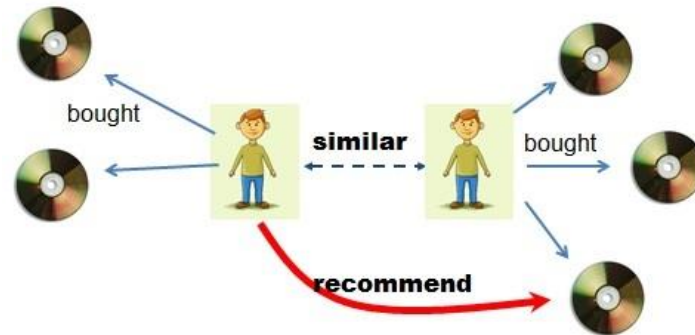
## › U – matice klient x faktor

## › L – matice faktorů

## › V – matice faktor x produkt

## › Opětovné vynásobení

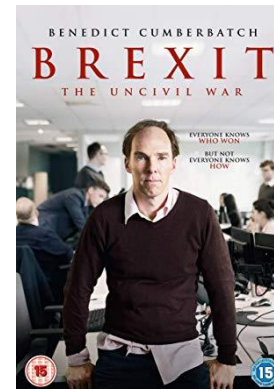
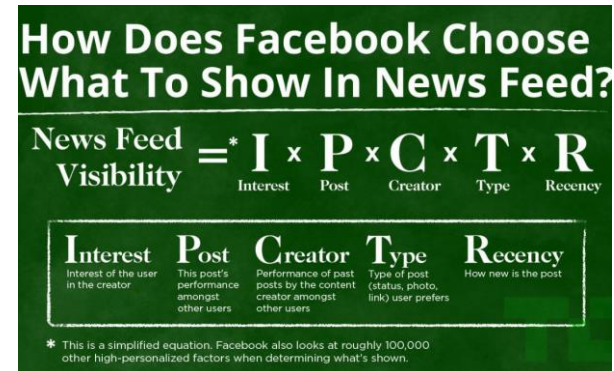
- Doporučení chybějících



$$\begin{array}{c} \mathbf{A} \end{array} = \begin{array}{c} \mathbf{U} \\ \begin{array}{|c|} \hline \text{gray bar} \\ \hline \end{array} \end{array} \begin{array}{c} \mathbf{L} \\ \begin{array}{|c|} \hline \text{gray bar} \\ \hline \end{array} \end{array} \begin{array}{c} \mathbf{V}^T \\ \begin{array}{|c|} \hline \text{gray bar} \\ \hline \end{array} \end{array}$$

# Facebook

- › Doporučování přátel
  - SNA, detekce vztahů a komunit
- › Doporučování obsahu
  - Personalizovaný výběr zpráv ve feedu
  - Zobrazování zpráv
    - které se líbí
    - které se líbí vašemu okolí
    - od zajímavých lidí
    - ...
  - Efektivní cílená reklama
    - Čím více víme, tím přesnější jsme
  - Real-time testování hypotéz
  - Vytváření bublin (echo chambers)
    - Efektivně funguje samoorganizace
    - Konfirmační bias
    - Mediální zkreslení
    - Cambridge Analytica (Brexit, Trump)
    - ...





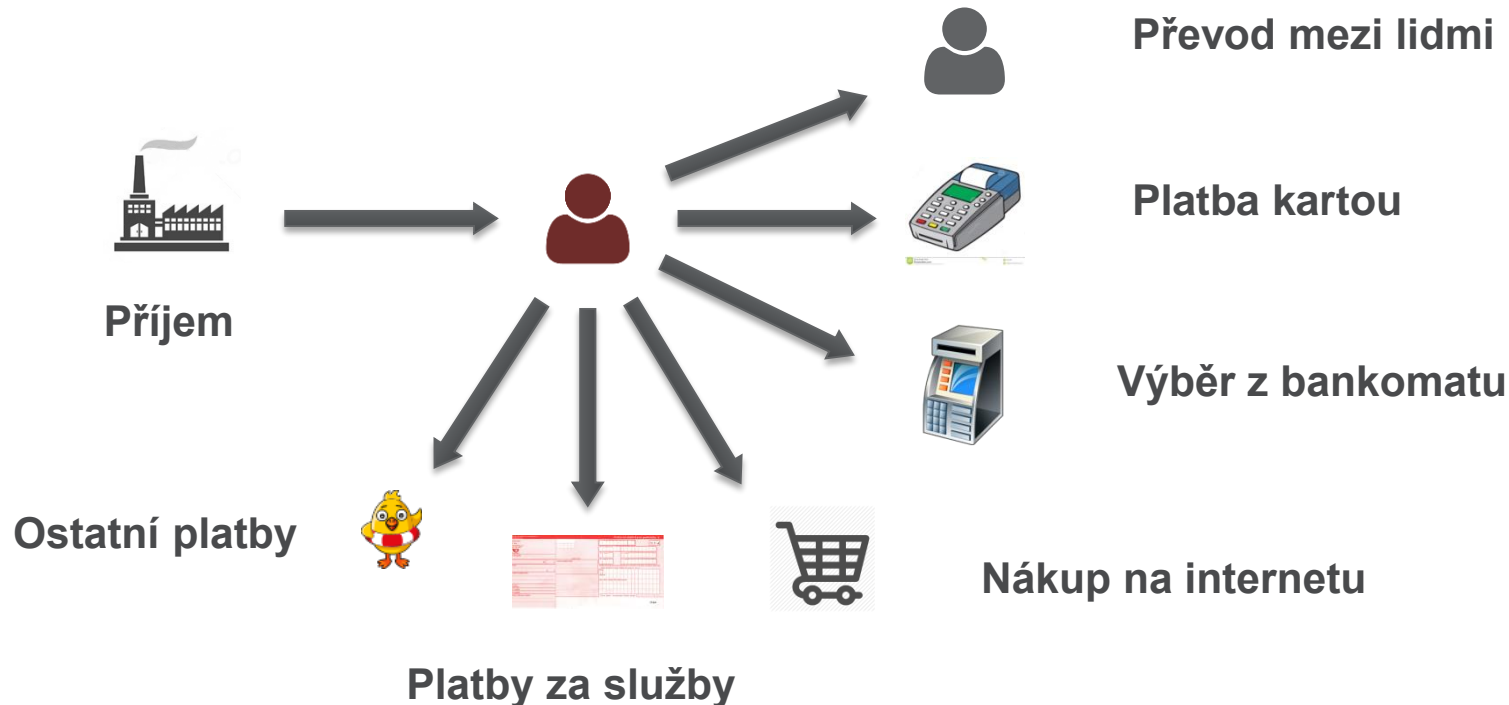
The background of the slide is a complex, abstract composition of numerous overlapping, translucent geometric shapes. These shapes, which include various polygons and rectangular prisms, are rendered in different shades of light gray and white. They are arranged in a way that creates a sense of depth and movement, as if they are floating or falling from the top right towards the bottom left. The overall effect is a modern, minimalist aesthetic.

Co děláme my?

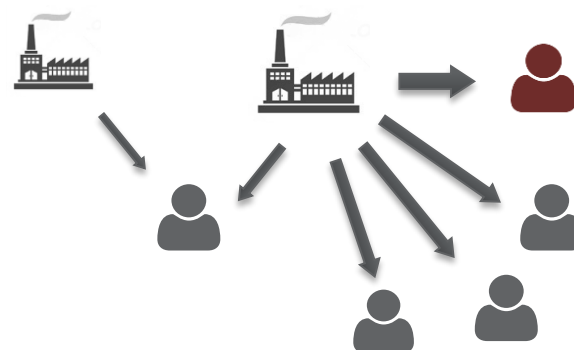


# Analýza Finančních Transakcí pomocí BD

- › Vytváříme vyladěné modely pro retailové banky
- › Vstup – finanční transakce
- › Výstup – využitelné informace o klientovi, příznaky, události,
- › Cílem je obohatit stávající obchodní proces o novou znalost



# Salary detector



- › Vstup
  - Finanční transakce typu firma - klient
- › Výstup: Identifikované vztahy zaměstnavatel – zaměstnanec
- › Business case
  - Rizikové skóre, detekce událostí, podobnosti (c2c/b2b),...
- › Principy
  - Detekce transakčních vzorců, text mining, pokročilá statistika
- › Vysoká přesnost i pro
  - Krátké úvazky – délka nepřesahující 3 měsíce
  - Nestandardní úvazky (částečné úvazky, práce na živnost, atd.)
  - Firmy s malým počtem zaměstnanců

# Detekce domácnosti – Banka/Telco

## › Vstup

- Klientské transakce – banka (c2c, karetní operace,...)
- Informace ze sítě – telco (cdr, lokace, billing)
- Základní demografie (věk, pohlaví, adresa, příjmení,...)

## › Výstup

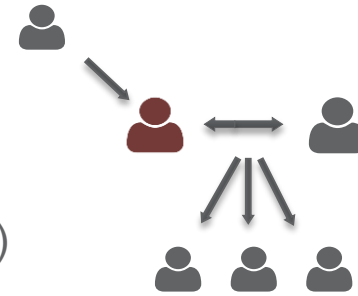
- Identifikace členů domácnosti a rodinných vztahů

## › Obchodní využití

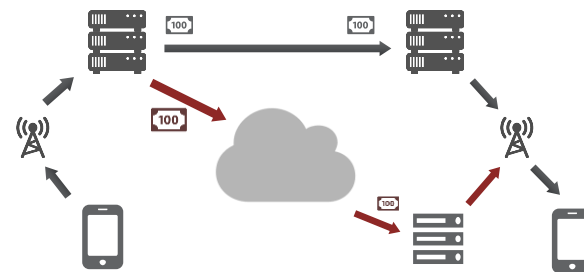
- Rodinný marketing, robustní rizikové skóre,...

## › Principy

- Detekce transakčních vzorců, analýza interakcí, text mining



# Telco Big Data SimBox Fraud



## › Scénář

- Zahraniční operátor/subjekt obchází standardní mezistátní hovor přes internet s cílem ušetřit na mezinárodním propojovacím poplatku

## › Vstup

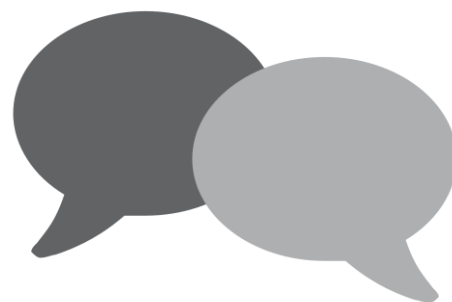
- Telco, síťová data (cdr, location, billing)

## › Výstup

- Identifikované podezřelé sim karty

## › Principy

- Detekce specifických typů neobvyklého chování
- Rozpoznání skupin s podobným chováním
- Automatická detekce pomocí roamingových dat



**Dotazy**