

# Aktuálne trendy v **informatike**

## Spracovanie dát (use cases)

Genči

Katedra počítačov a informatiky

2017

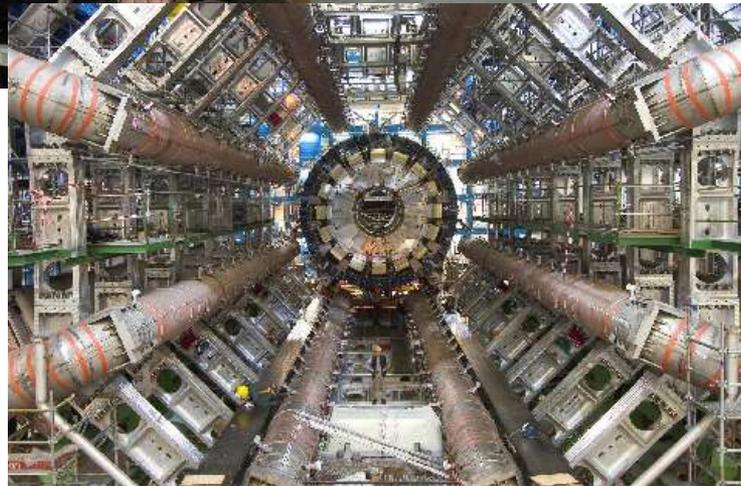
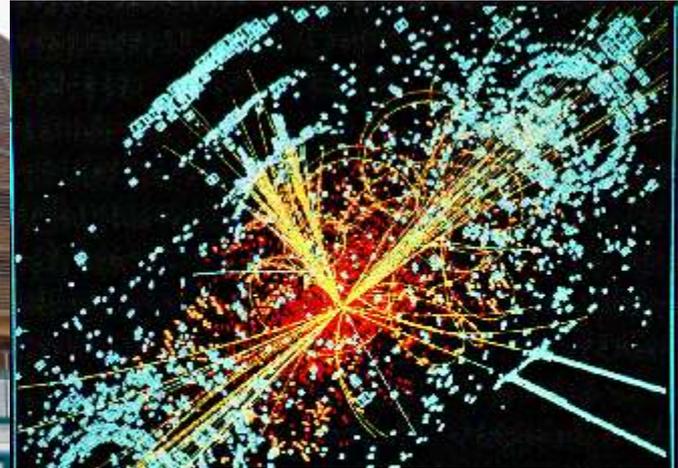
# Galileo Galilei

## Čo to vlastne robil?



# CERN – LHC

Čo to vlastne robia?



# Teda – čo vlastne robí-li/a?

- Naplánovali experiment
- Nazbierali dáta
- Vyhodnotili dáta (premenili ich na informácie)
- Urobili závery (získali znalosti)

# Čo zvyčajne robíme dnes?

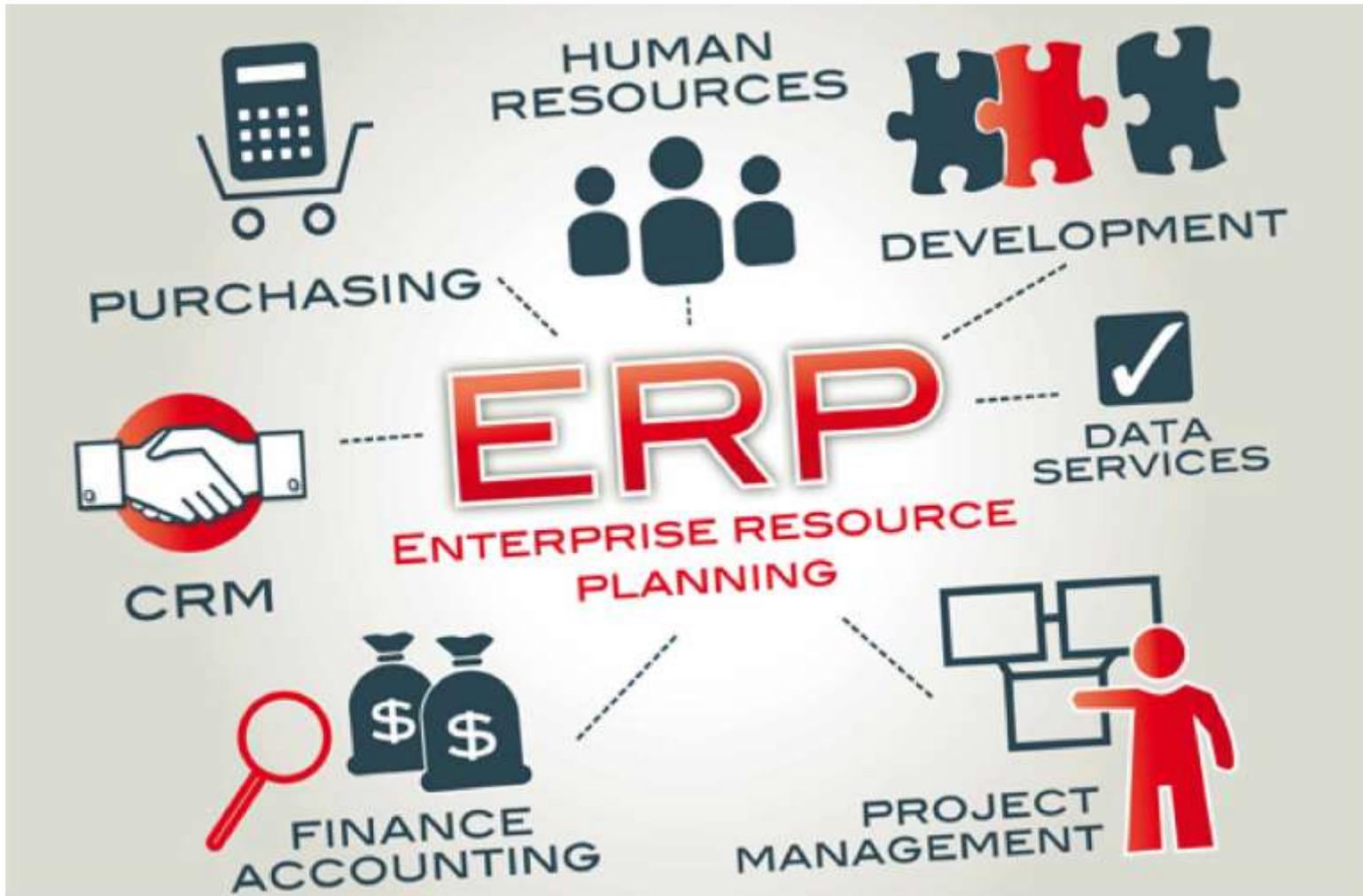
The image shows a screenshot of the website for the Modular Academic Information System (MAIS) at the Technical University of Košice. The website header includes the logo "MAIS MODULÁRNY AKADEMICKÝ INFORMAČNÝ SYSTÉM" and a language selector for "SK". Navigation links for "STÚDIUM" and "KONTAKT" are visible. The main content area features a search bar and a link to "viac informácií" (more information). Three photographs are overlaid on the website: a woman at a computer workstation, a call center operator, and a doctor in a white coat using a computer in a clinical setting. At the bottom of the website screenshot, the contact information "Otázky a pripomienky: [mais@helpdesk.tu](mailto:mais@helpdesk.tu)" is displayed.

# Zvyčajne ...

- Zbierame dáta (v nesmiernom rozsahu)
- Bez nejakého explicitného plánu
- Tieto dáta predstavujú informácie a obsahujú znalosti

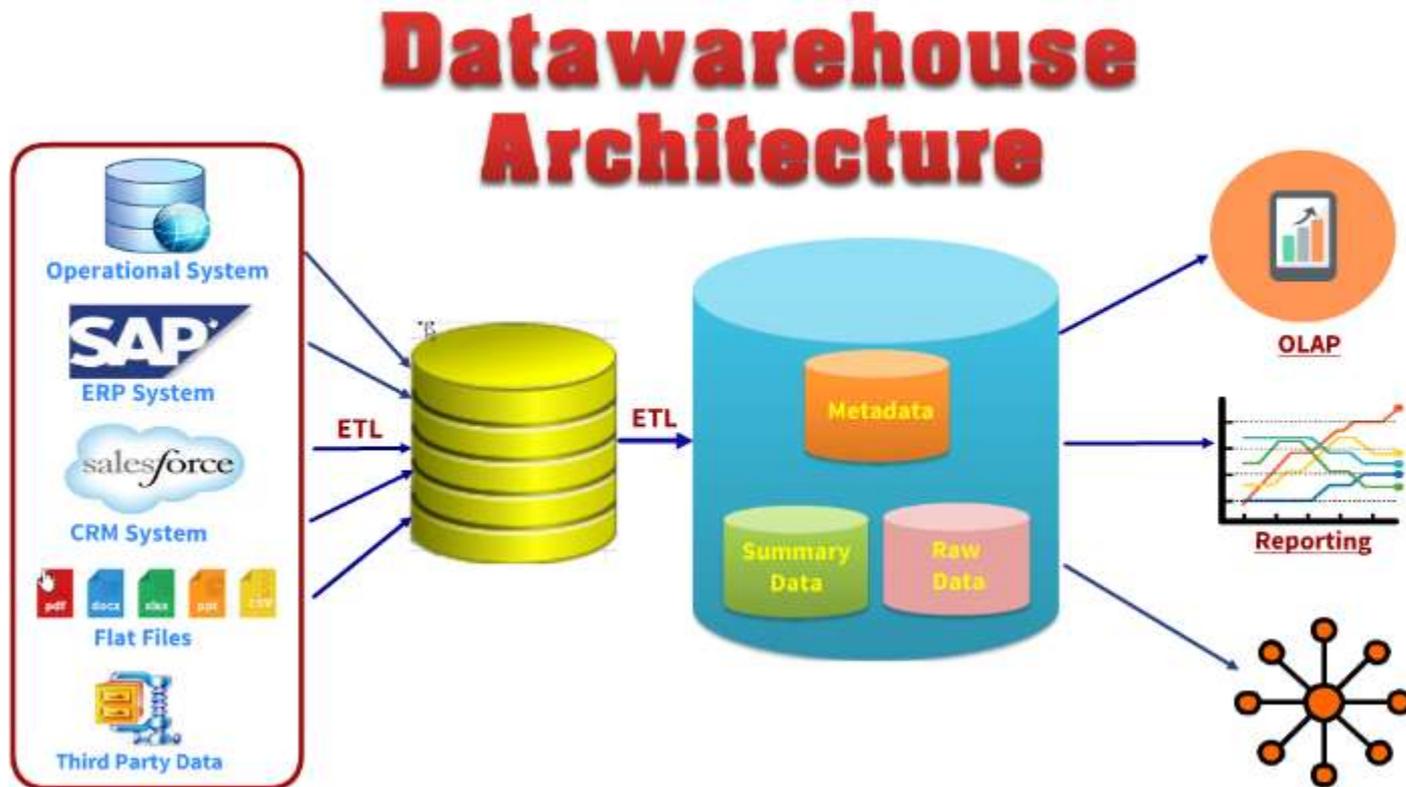
**Ešte „nedávno“**

# ERP (OLTP)





# Data warehouse

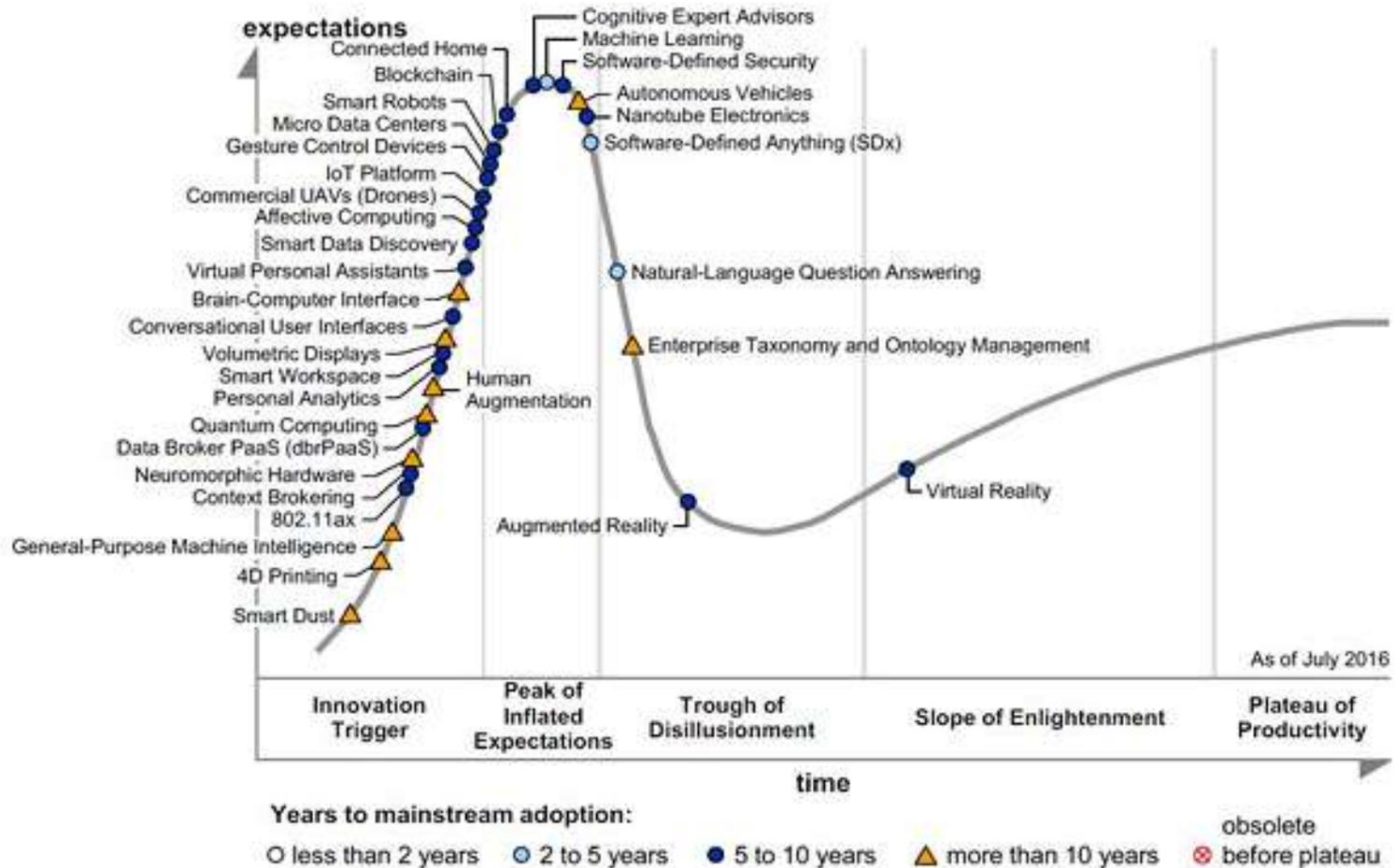




# Business Intelligence

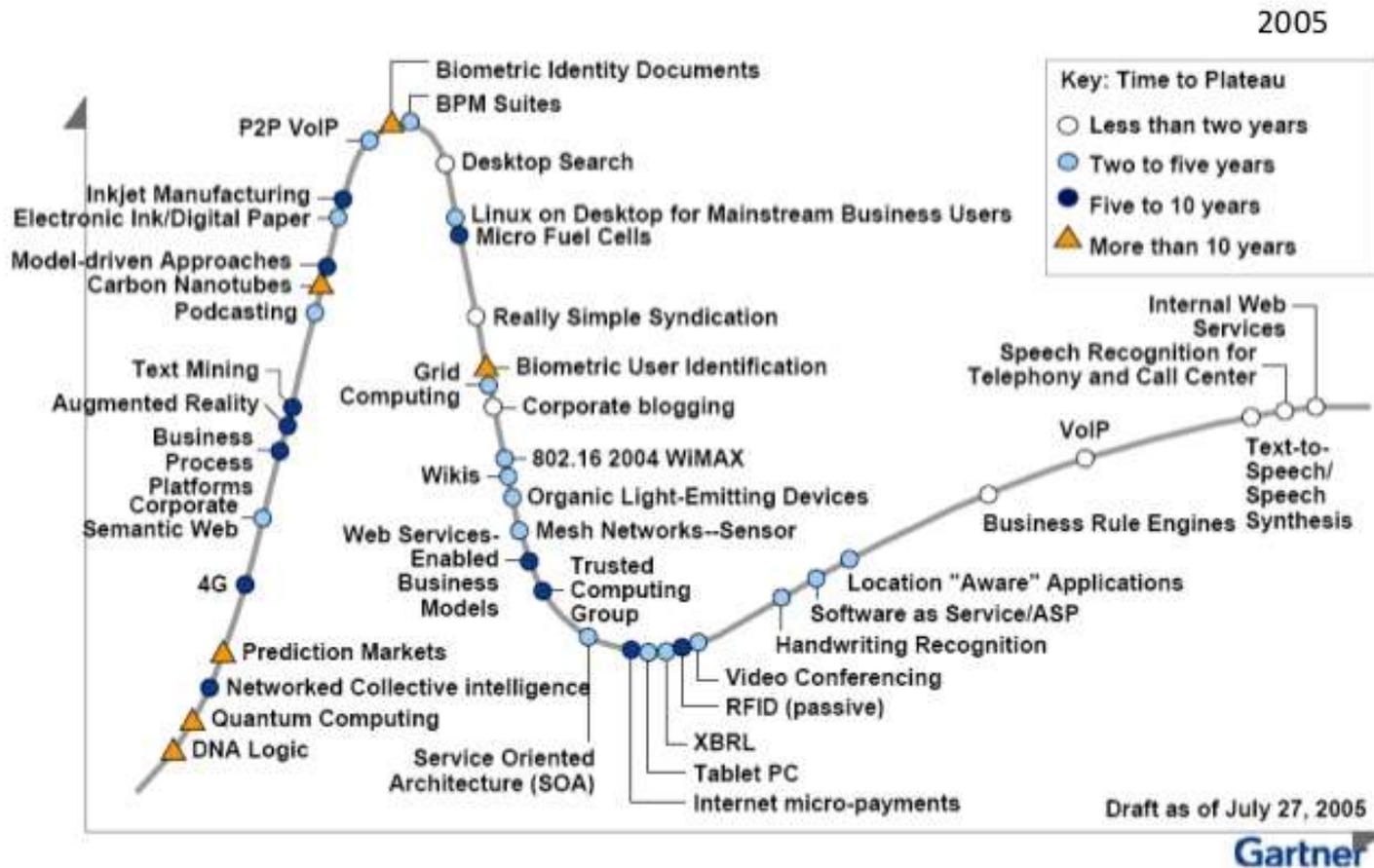


# Gartner Hype Cycle (kde je ERP, DM, BI?)

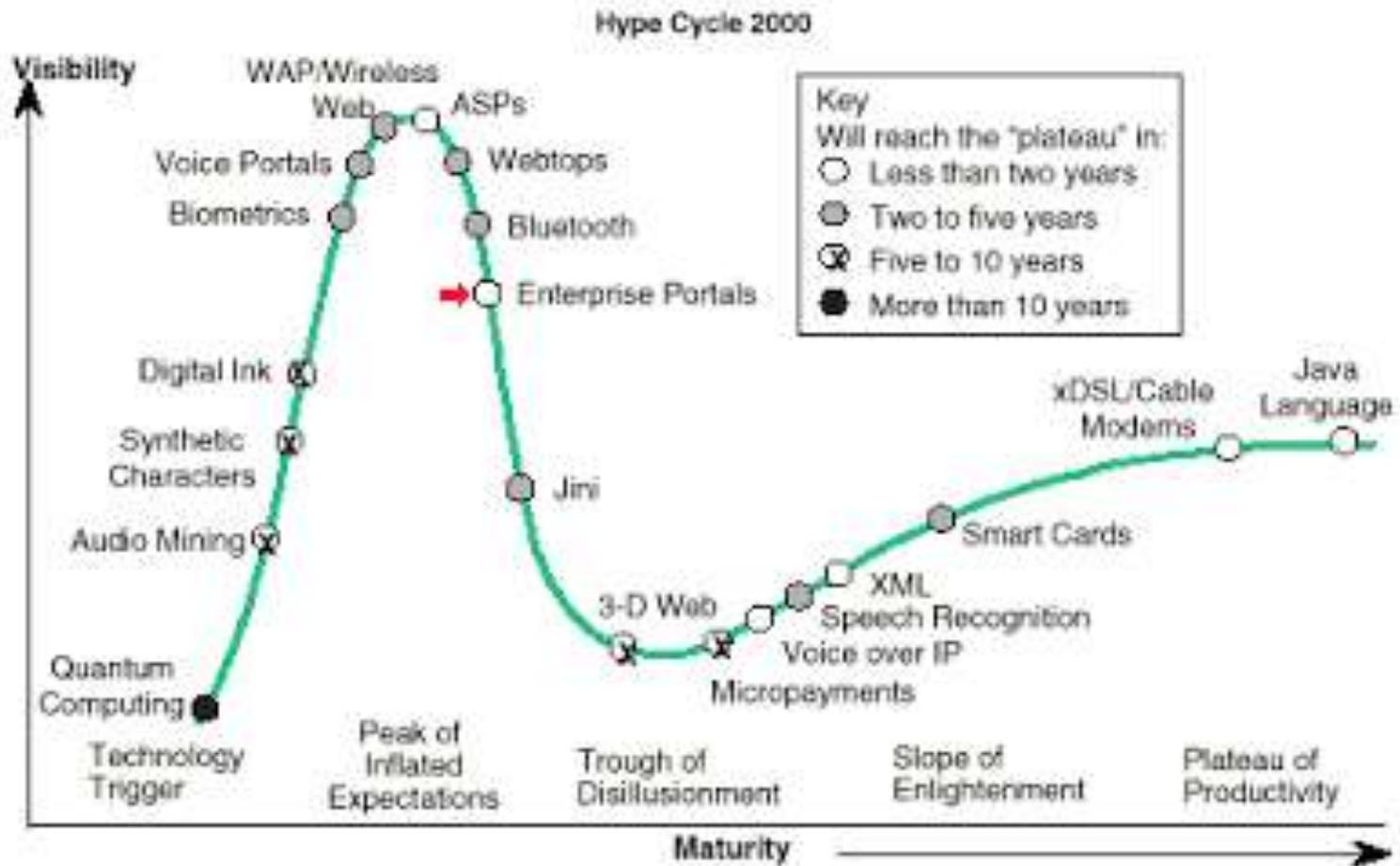


Source: Gartner (July 2016)

# Gartner Hype Cycle 2005



# Gartner Hype Cycle 2000



# Salary comparision

- <http://www.computerworld.com/article/3169664/it-careers/13-tech-jobs-that-pay-200k-salaries.html>
- <http://www.computerworld.com/salarysurvey/breakdown/2016/joblevel/3>

# Hot Skills

**Top 10 skills**  
respondents **plan**  
**to hire** for in the  
next 12 months:

Source: Computerworld's  
Forecast 2017 survey  
of 196 IT managers,  
directors and executives.

Base: 57 respondents  
who expect to increase  
IT head count in the  
next 12 months.

**Programming/  
application  
development** 35%

**Help desk/  
tech support** 35%

**Security/  
compliance/  
governance** 26%

**Cloud/SaaS** 26%

**Business  
intelligence/  
analytics** 26%

**Web  
development** 26%

**Database  
administration** 25%

**Project  
management** 25%

**Big  
data** 25%

**Mobile  
applications  
and device  
management** 21%

# Harvard Business Review

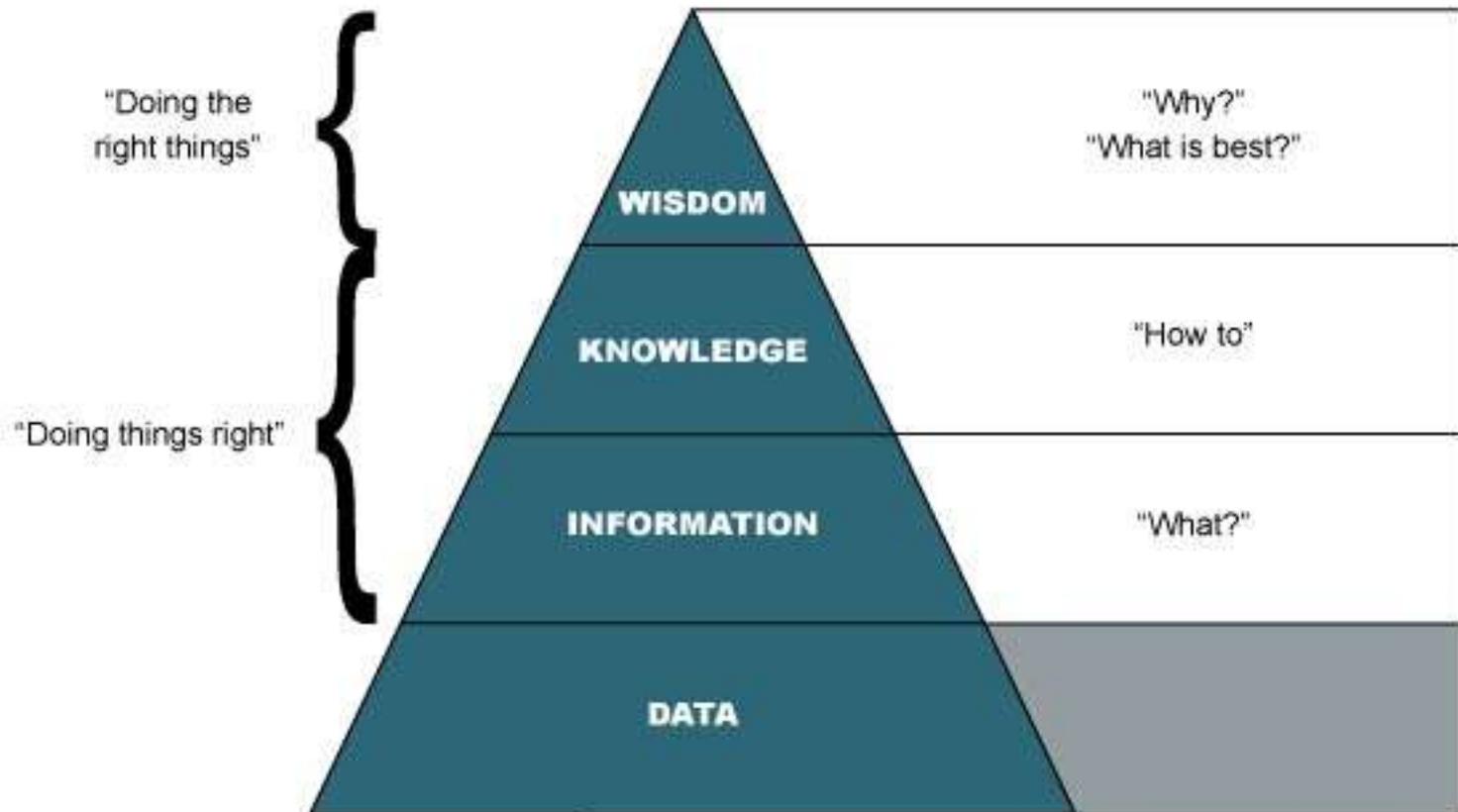
- Sexiest Job of the 21st Century?

The screenshot shows the Harvard Business Review website interface. At the top, there is a navigation bar with the HBR logo, a search bar, and links for 'SUBSCRIBE', a shopping cart, and 'SIGN IN'. Below the navigation bar is a utility bar with icons for 'SUMMARY', 'SAVE', 'SHARE', 'COMMENT', 'TEXT SIZE', 'PRINT', and 'BUY COVER' for \$8.95. The main content area features the article title 'Data Scientist: The Sexiest Job of the 21st Century' in large, bold black text, with the authors 'by Thomas H. Davenport and D.J. Patil' listed below. A small red 'DATA' tag is positioned above the title. To the right of the title, there is a 'WHAT TO READ NEXT' section with three links: 'Big Data: The Management Revolution', '5 Essential Principles for Understanding Analytics', and 'Data Scientists Don't Scale'. Below the article title, the text begins with a large 'W' followed by 'hen Jonathan Goldman arrived for work in June 2006 at LinkedIn, the business networking site, the place still felt like a start-up. The company had just under 8 million accounts, and the number was growing quickly as existing members invited their friends and colleagues to join. But users weren't seeking out connections with the people who were already on the site at the rate executives had expected. Something was apparently missing in the social'. On the right side, there is a 'VIEW MORE FROM THE October 2012 Issue' section featuring a thumbnail of the magazine cover with the headline 'GETTING CONTROL OF BIG DATA'. At the bottom of the page, a red banner contains the text '3/4 FREE ARTICLES LEFT → REGISTER FOR MORE | SUBSCRIBE + SAVE!'.

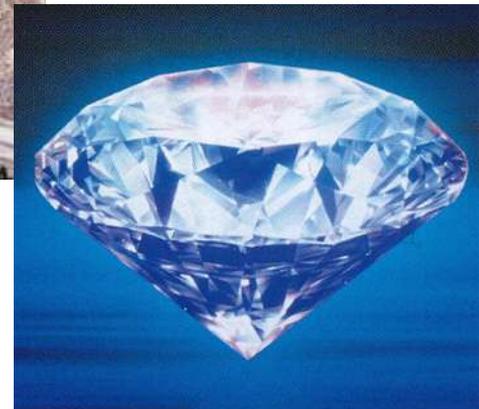
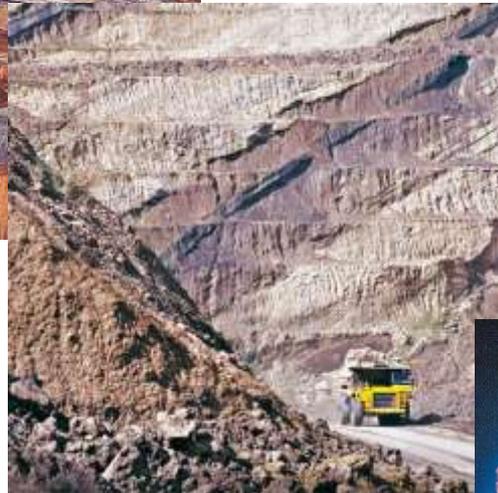
# Wikipedia

- **Data Science** is an interdisciplinary field about processes and systems to extract [knowledge](#) or insights from [data](#) in various forms, either structured or unstructured, which is a continuation of some of the data analysis fields such as [statistics](#), [data mining](#), and [predictive analytics](#), similar to [Knowledge Discovery in Databases](#) (KDD).

# Podstata



# Podstata (analógia)

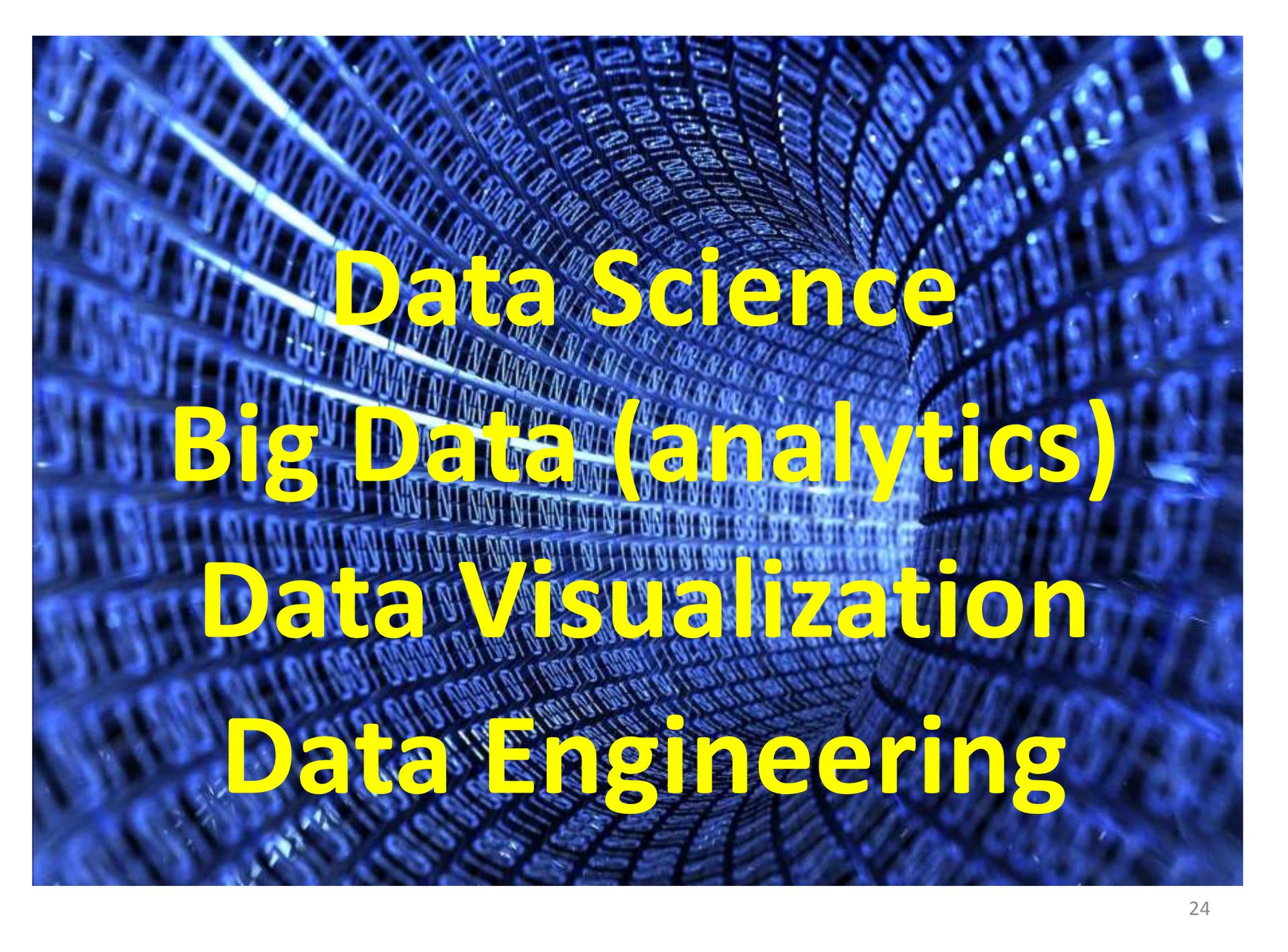




## Audience Data Mining

# Wikipedia (pokr.)

- Data science employs techniques and theories drawn from many fields within the broad areas of [mathematics](#), [statistics](#), [information science](#), and [computer science](#), including [signal processing](#), [probability models](#), [machine learning](#), [statistical learning](#), [data mining](#), [database](#), [data engineering](#), [pattern recognition and learning](#), [visualization](#), [predictive analytics](#), [uncertainty modeling](#), [data warehousing](#), [data compression](#), [computer programming](#), [artificial intelligence](#), and [high performance computing](#).



**Data Science**  
**Big Data (analytics)**  
**Data Visualization**  
**Data Engineering**

# Big Data

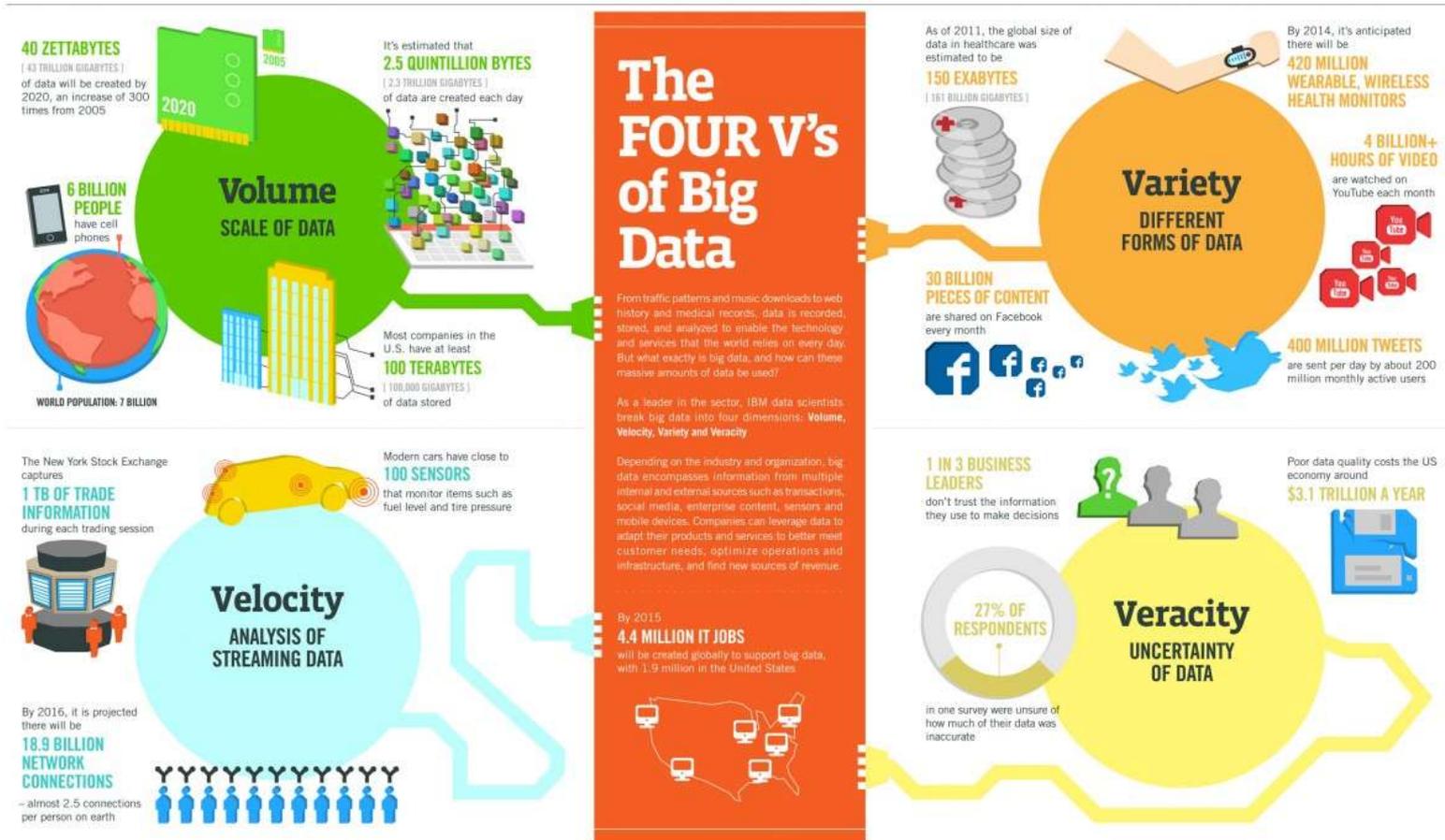
**3V = Volume, Velocity** (stream)

and **Variety** (structured, semi-structured and unstructured)

**4V = 3V + Veracity** (IBM)

**5V = 4V + Value** (ORACLE?)

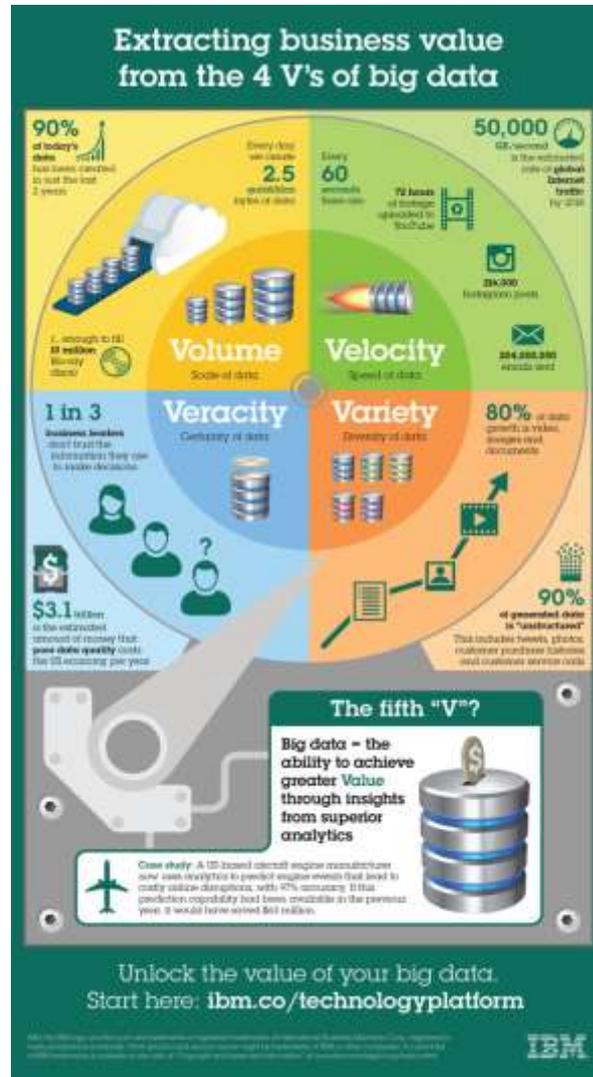
# Big Data (IBM's 4V)



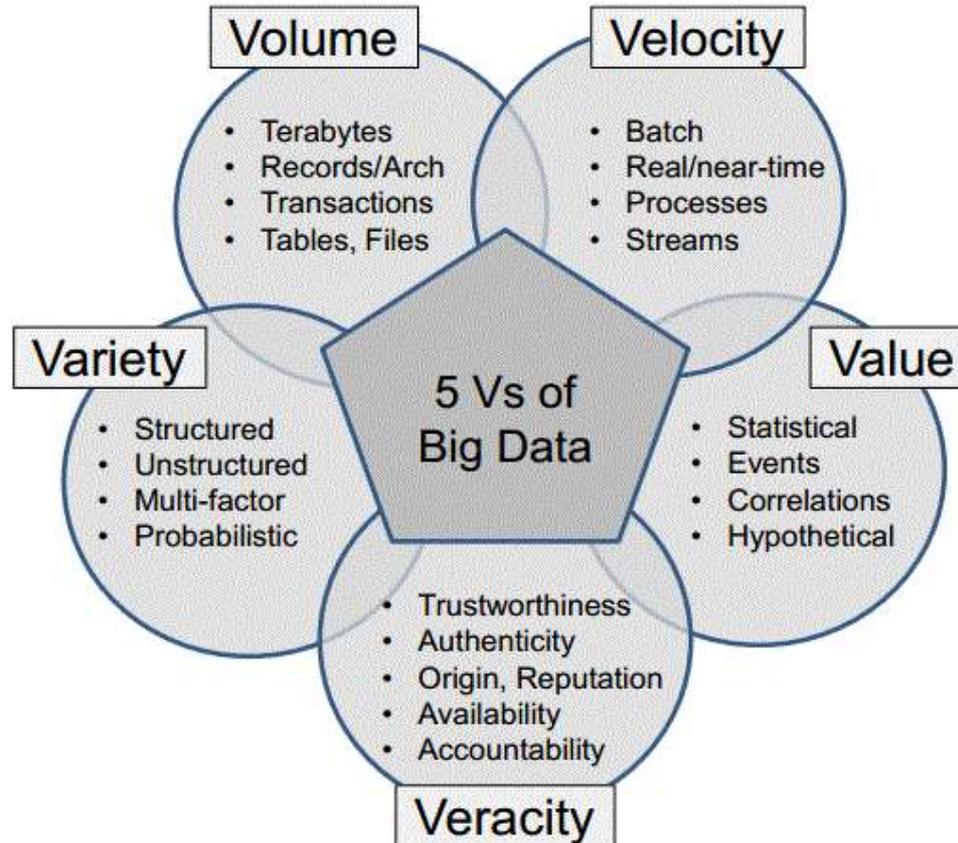
Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTec, GSA



# Big Data (IBM's 4V)



# Big Data 5V



# Big Data 7V

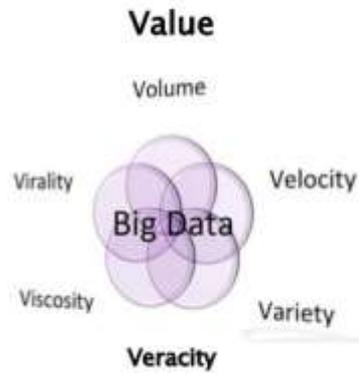
## From the traditional 3-4 V's towards the 5-7 V's

Data Pioneers  
10 april 2013

**Viscosity** - Viscosity measures the resistance to flow in the volume of data. This resistance can come from different data sources, friction from integration flow rates, and processing required to turn the data into insight. Technologies to deal with viscosity include improved streaming, agile integration bus', and complex event processing.

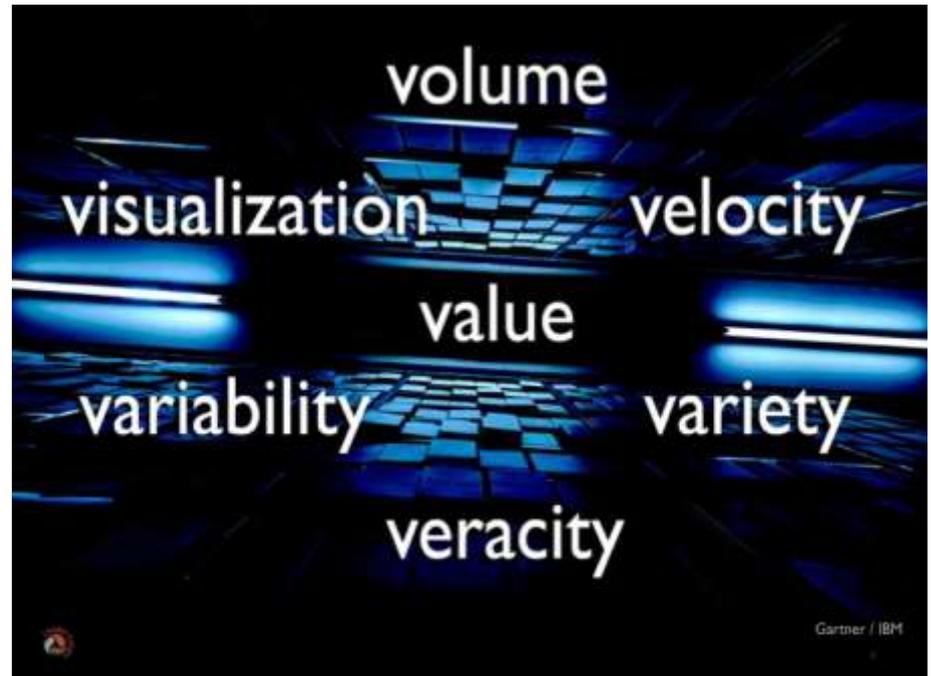
**Virality** - Virality describes how quickly information gets dispersed across people to people (P2P) networks. Virality measures how quickly data is spread and shared to each unique node. Time is a determinant factor along with rate of spread.

**Veracity**: Trust & Quality



Your business technology. Powering progress.

AtoS



Gartner / IBM



# NoSQL a NewSQL

- NoSQL (nie ~~No SQL~~ ale **Not only** SQL !!! )
  - Key value
  - Document Store
  - Column store
  - Graph
- NewSQL
  - class of modern relational database management systems that seek to provide the same scalable performance of NoSQL systems for online transaction processing (OLTP) read-write workloads while still maintaining the ACID guarantees of a traditional database system (<https://en.wikipedia.org/wiki/NewSQL>)

# NoSQL CAP theorem

It is impossible for a [distributed computer system](#) to simultaneously provide more than two out of three of the following guarantees:

- **Consistency** - every read receives the most recent write or an error
- **Availability** - every request receives a (non-error) response – without guarantee that it contains the most recent write
- **Partition tolerance** - the system continues to operate despite an arbitrary number of messages being dropped (or delayed) by the network between nodes

Zdroj: [https://en.wikipedia.org/wiki/CAP\\_theorem](https://en.wikipedia.org/wiki/CAP_theorem)

# Data Lakes

The method of storing [data](#) within a system or repository, in its **natural format**, that facilitates the collocation of data in various schemata and structural forms, usually object blobs or files.

Why natural format?

[https://en.wikipedia.org/wiki/Data\\_lake](https://en.wikipedia.org/wiki/Data_lake)



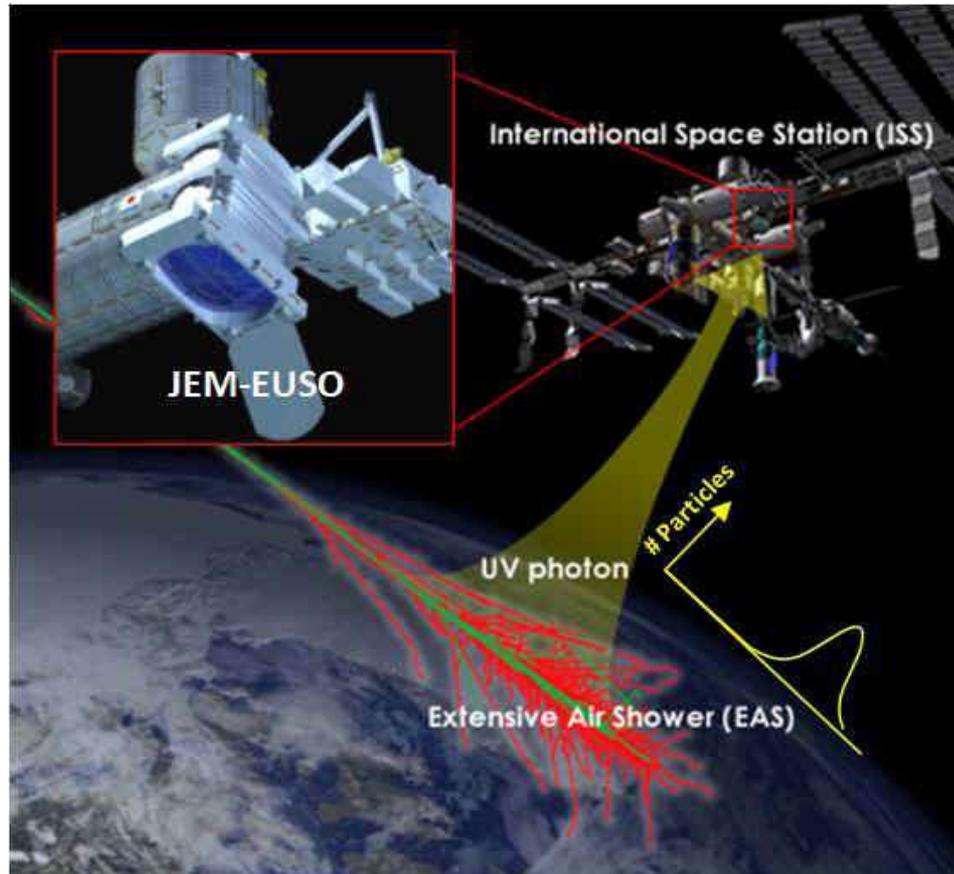
# Príklady našej z praxe

Čo aktuálne robíme

# Kozmická fyzika

Stručně, přídu porozprávět zo SAV

# JEM-EUSO





# Vplyv kozmického počasia na ...

- sentiment obyvateľstva (globálny sentiment);
- výskyt diagnóz (infarkty, presnosť dát);
- výskyt udalostí
- ...

# Intelligentné siete

## Smart Grid (PowerEng+IT)

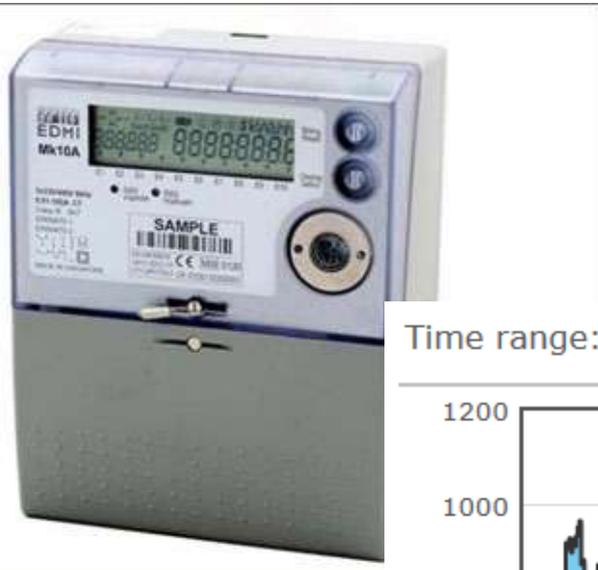
# Čo sme zažili



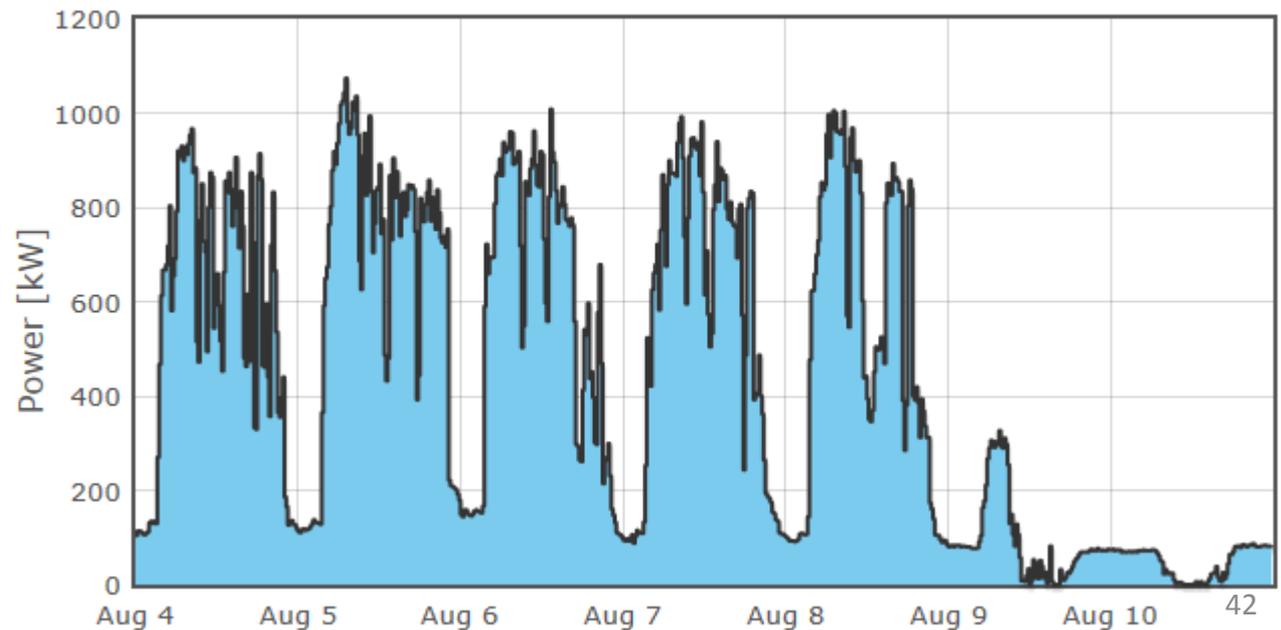
# Čo zažívame



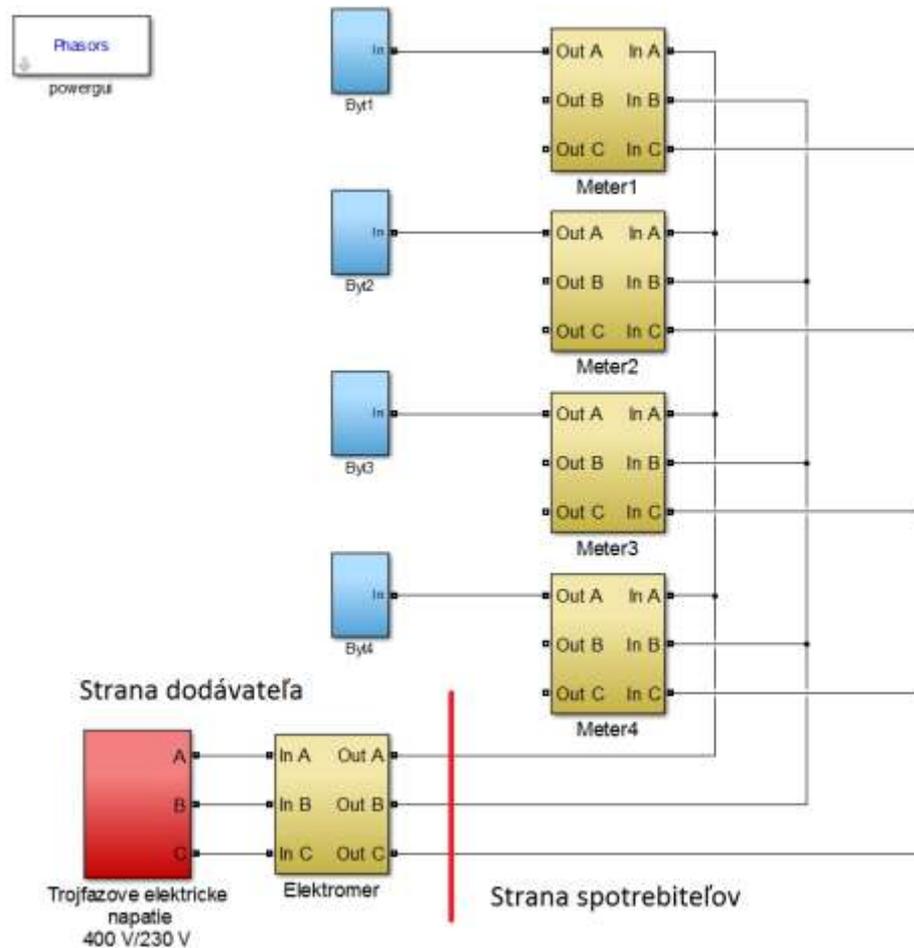
# Simulačné modely pre analýzu inteligentných sietí (Smart Grid)(1)



Time range:   Mon, 4 August 2014 - Sun, 10 August 2014

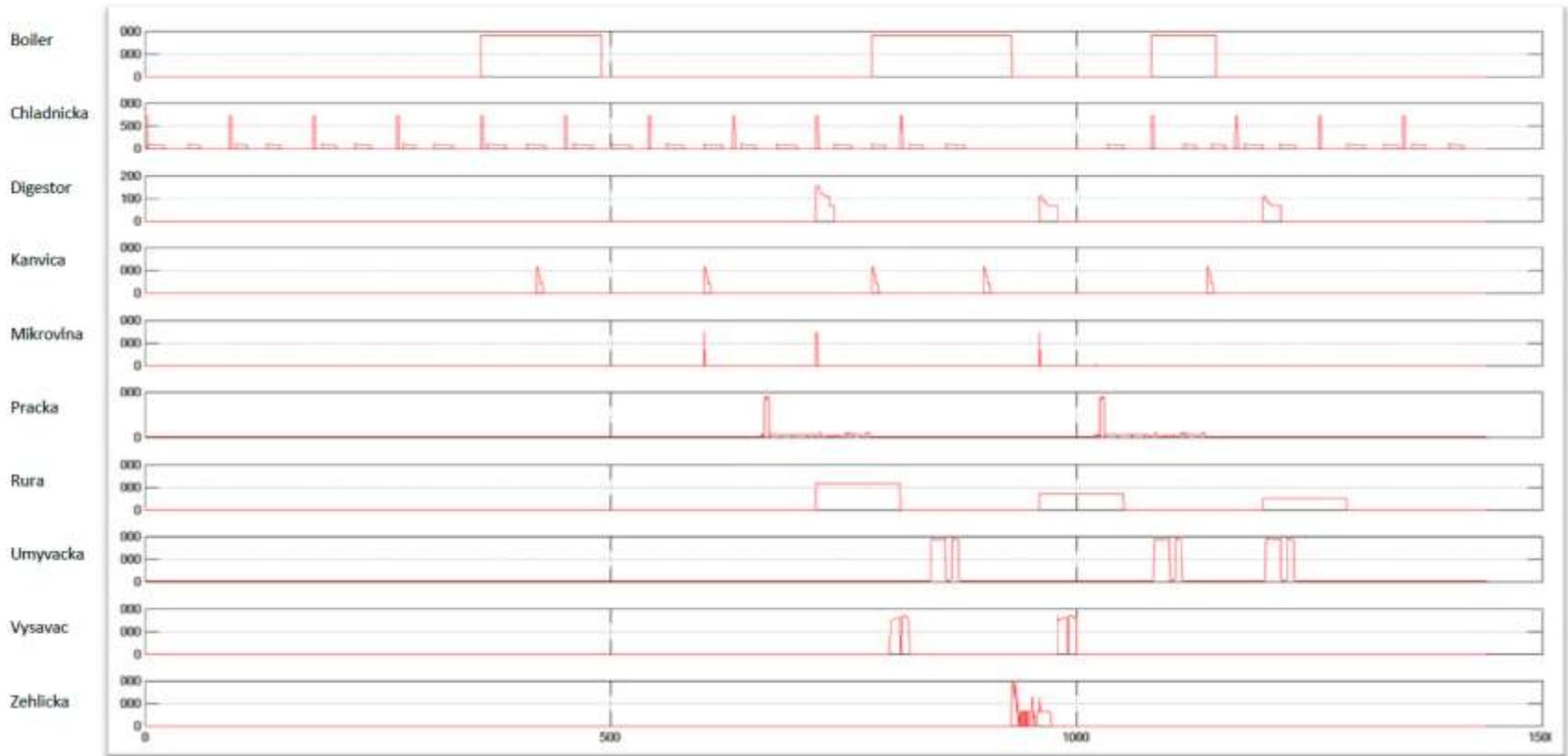


# Simulačné modely pre analýzu inteligentných sietí (Smart Grid)(2)

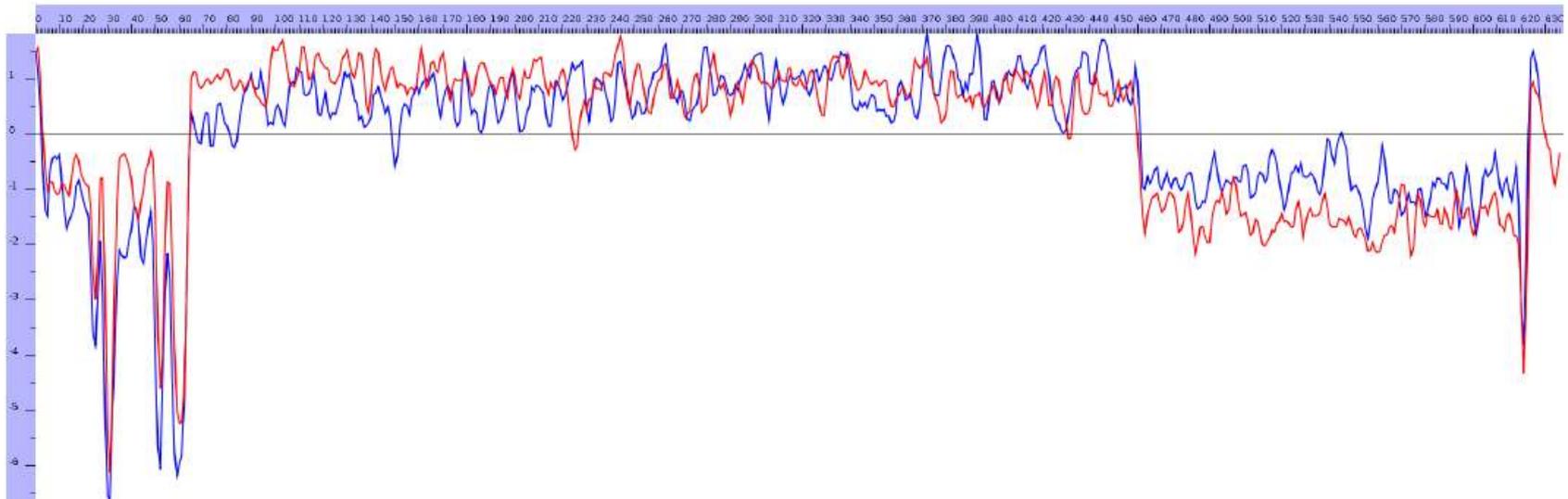


# Simulačné modely pre analýzu inteligentných sietí (Smart Grid)(3)

Jednotlivé spotrebiče (ukážka domácnosti)



# Určenie filmu, podobnosti programov (1)



Obrázok 8 – 6 Dva profily spotreby pre rovnaké video s množstvom tmavých scén

# Určenie filmu, podobnosti programov (2)

		Philips 32PHH4100/88					
		Nahrávky	A	B	C	D	E
Tesla S 7090 TSP2	A	<u>0,48</u>	0,23	-0,01	-0,25	0,23	
	B	0,12	<u>0,53</u>	0,00	-0,07	0,12	
	C	-0,19	-0,09	<u>0,37</u>	0,29	-0,16	
	D	-0,06	-0,07	0,11	<u>0,44</u>	0,02	
	E	0,10	0,12	0,09	-0,17	<u>0,55</u>	

		Samsung UE40D5003					
		Nahrávky	A	B	C	D	E
Tesla S 7090 TSP2	A	<u>0,12</u>	-0,20	0,01	-0,05	<u>0,04</u>	
	B	0,00	0,15	<u>0,06</u>	-0,04	-0,07	
	C	-0,07	<u>0,26</u>	0,00	0,05	-0,09	
	D	-0,08	-0,03	0,00	<u>0,07</u>	0,02	
	E	-0,04	-0,15	-0,12	-0,12	-0,02	

# Ďalšie úlohy

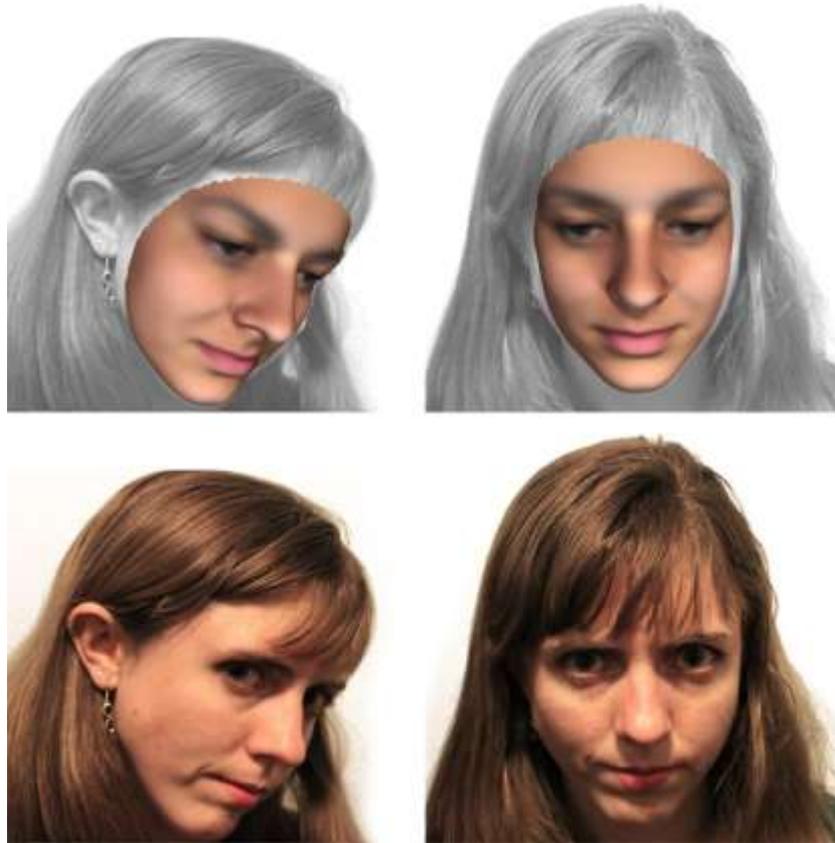
Security

# Bioinformatika

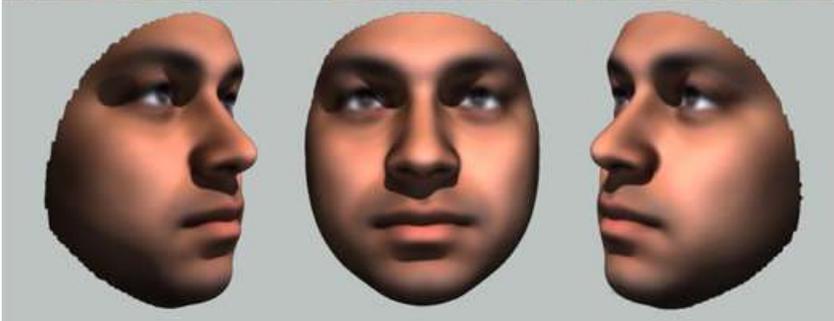
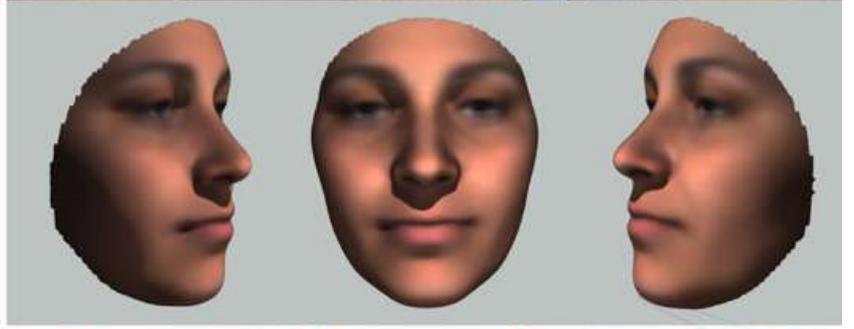
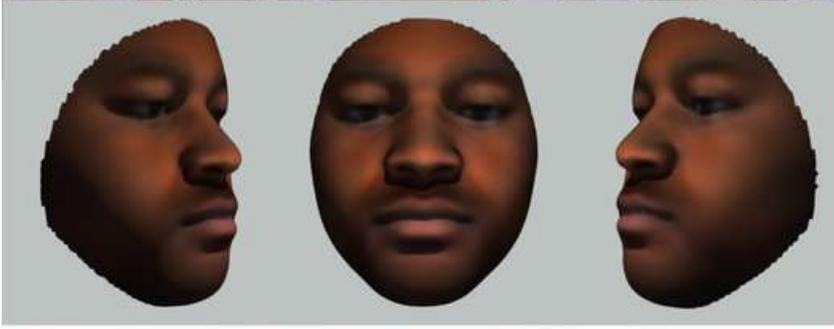
# Genomika



# Čo napr. (už?) vieme získať z genómu?

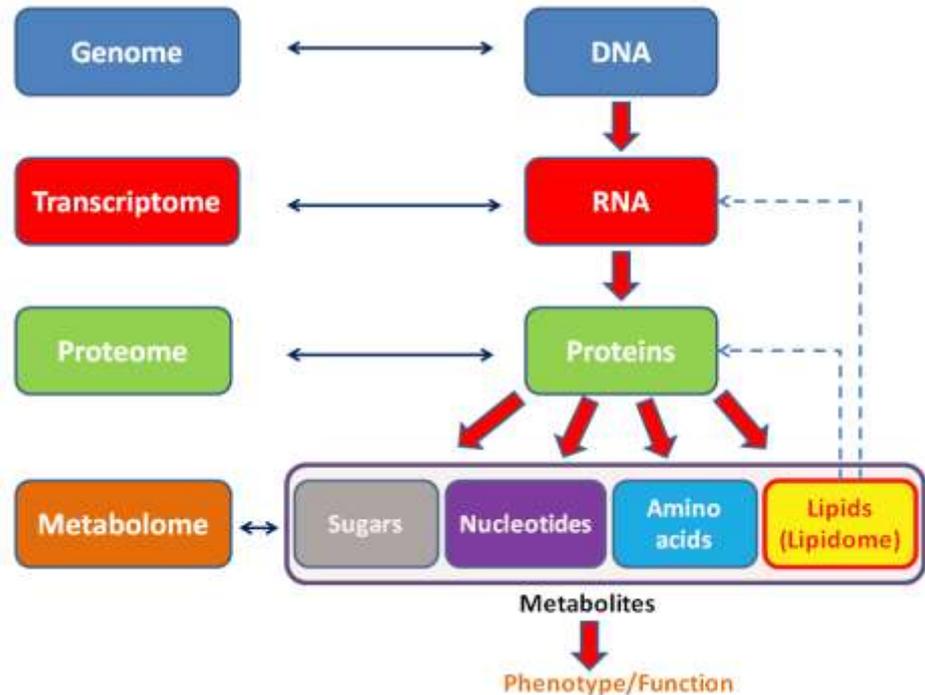


- **Modeling 3D Facial Shape from DNA**
- <https://www.newscientist.com/article/mg22129613-600-genetic-mugshot-recreates-faces-from-nothing-but-dna/>

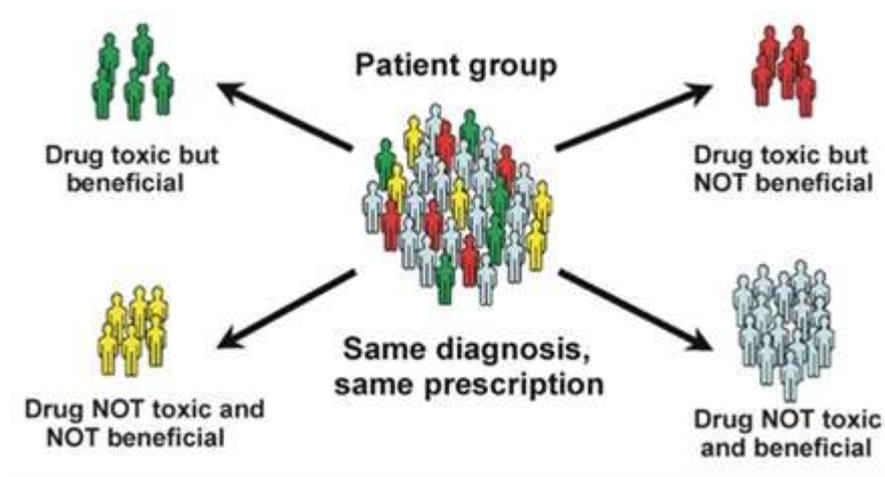
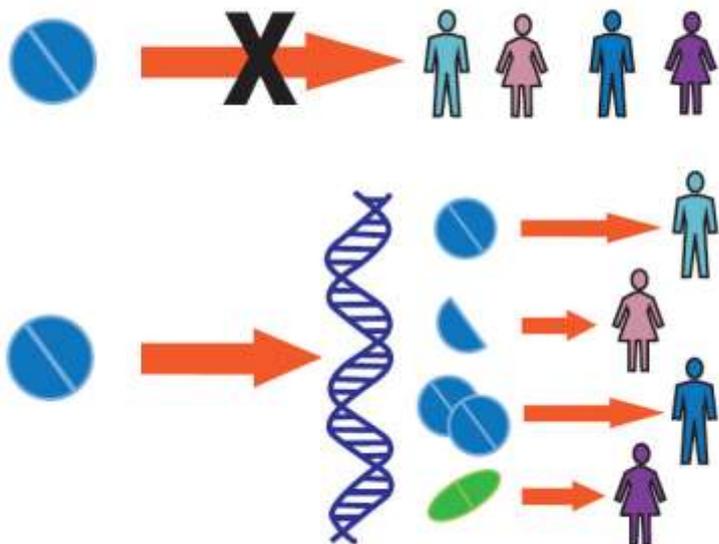


# -omics

- <https://en.wikipedia.org/wiki/Omics>
  - Genomics
  - Transcriptomics
  - Proteomics
  - Lipidomics
  - Metalobomics
  - ...



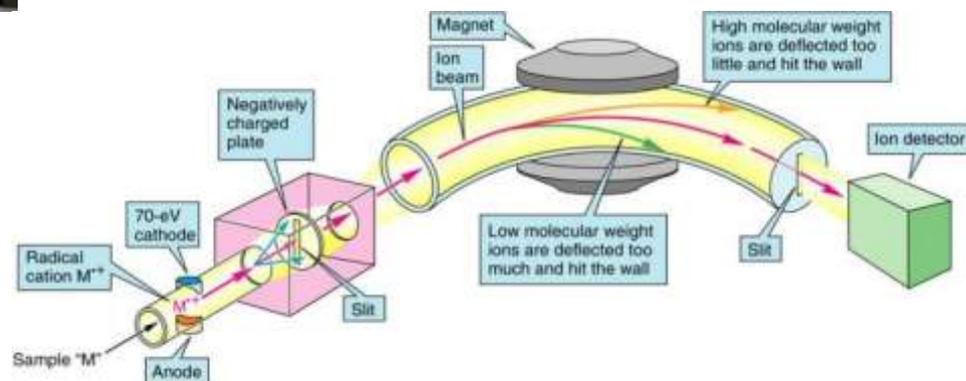
# Načo je to dobré?



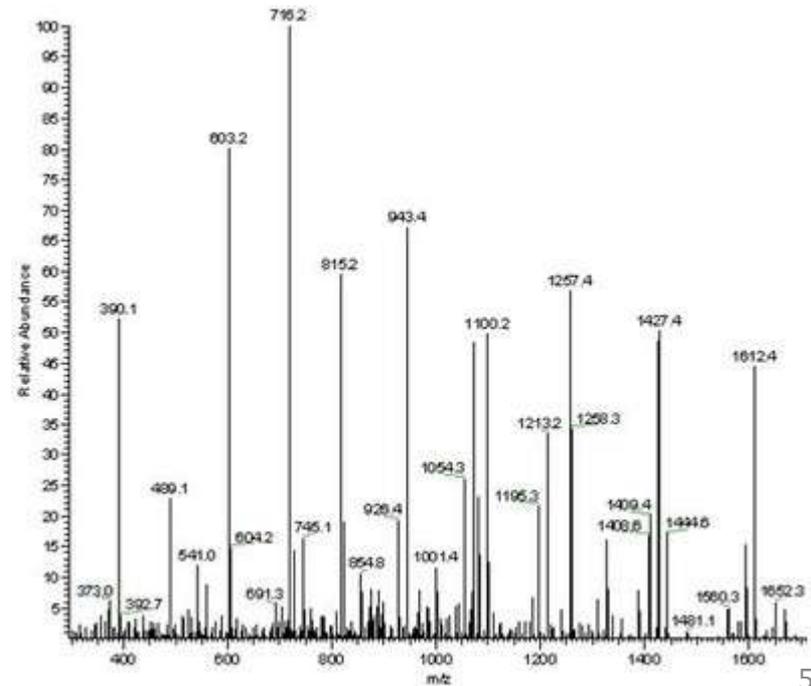
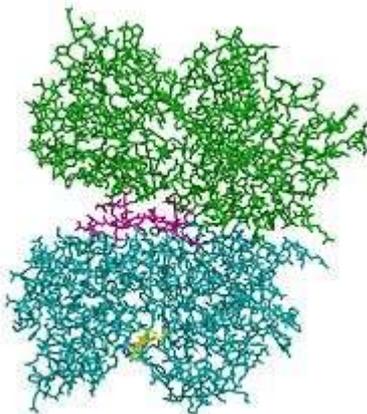
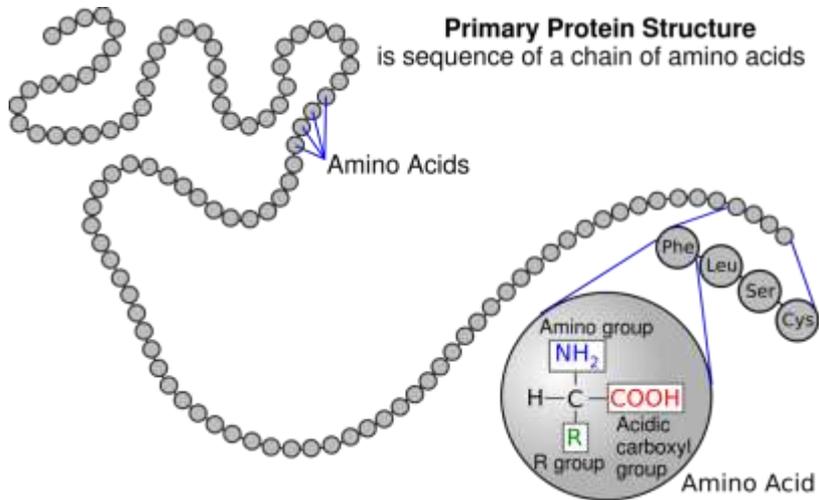
# Personalizovaná medicína

## Data mining v proteomike

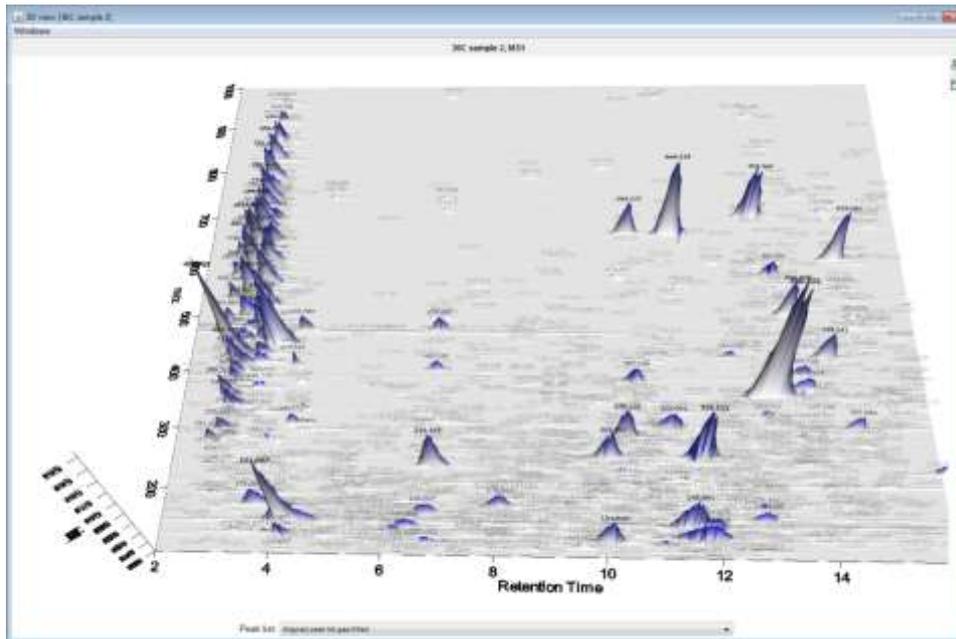
# Hmotnostná spektrometria



# Proteíny a MS



# „Biomarkery“

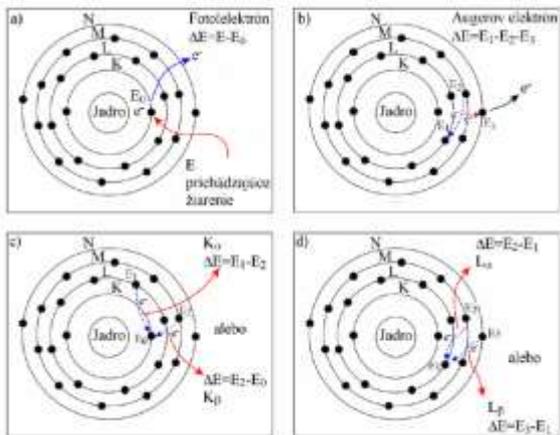


# Reálny život

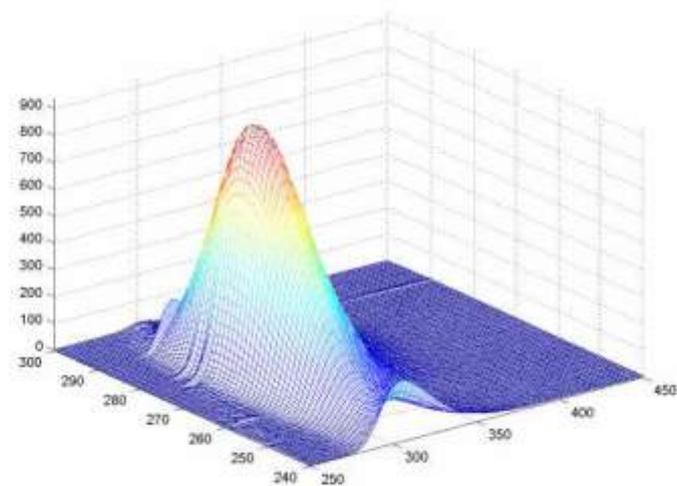
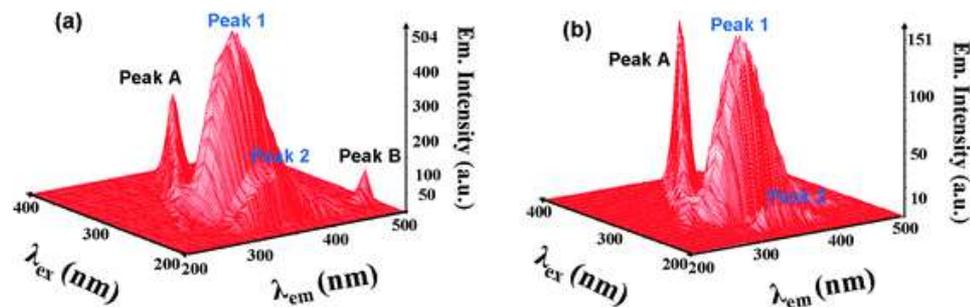
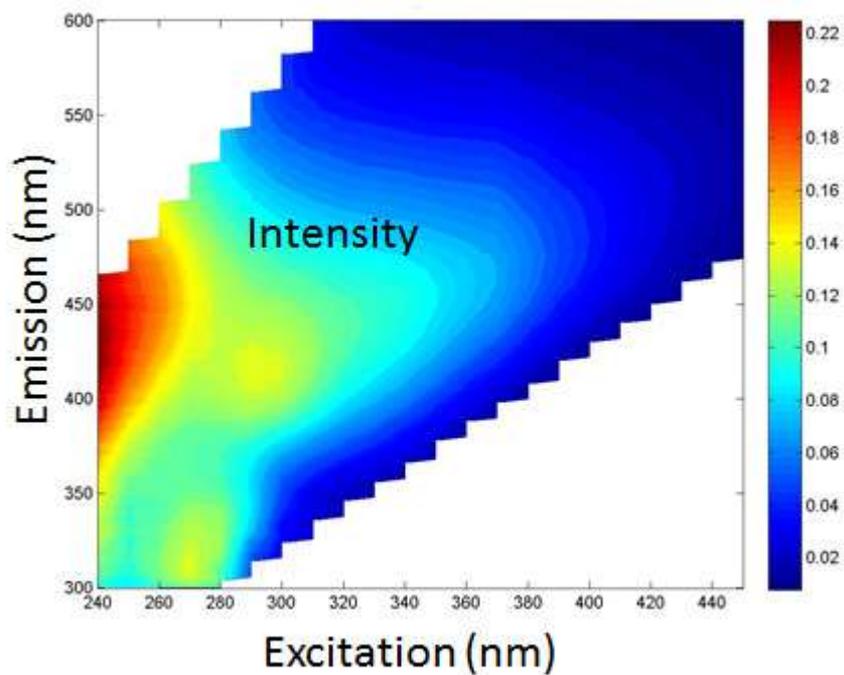
- Podpora spracovania dát – data engineering

# Spracovanie fluorescenčných spektrálnych matric biologického materiálu

# Fluorescencia



# Spektrá - biomarkery



# Počítačová lingvistika

## Computational linguistics

# Vymedzenie pôsobnosti

Aplikácia informačných technológií na  
spracovanie prirodzeného jazyka



natural language  
processing

# Úlohy

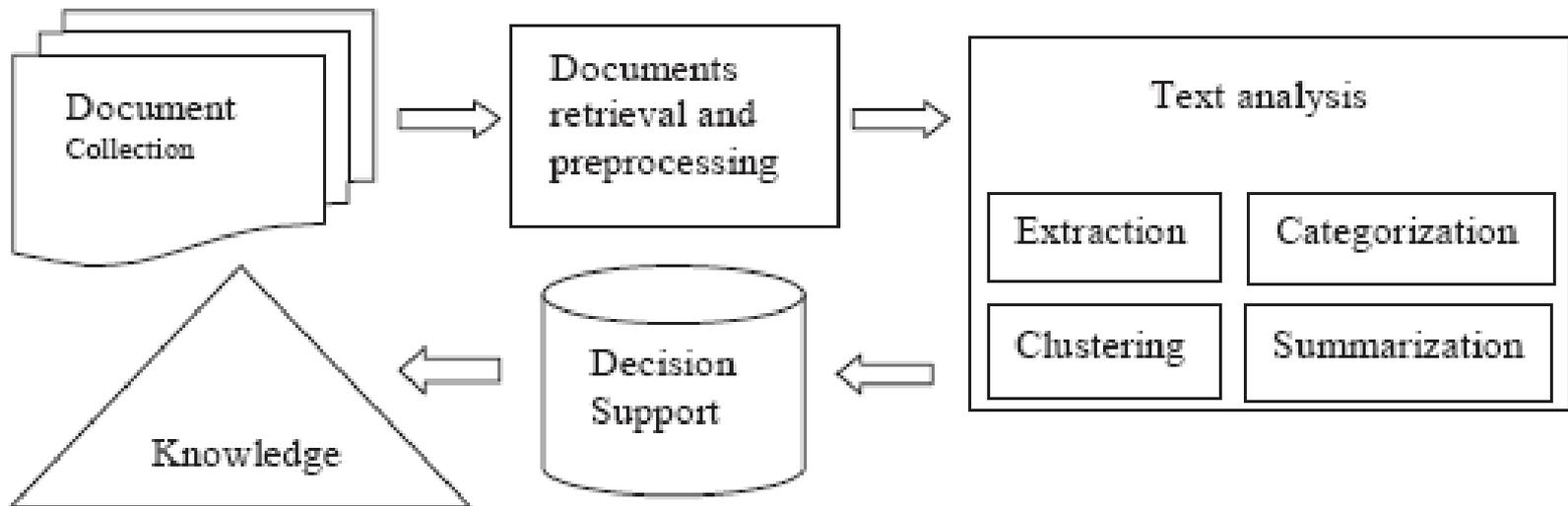
- strojový preklad (asi top doména)
- dolovanie informácií (text mining)
- korektory pravopisu, gramatiky (len prípady)
- budovanie lingvistických zdrojov pre SK (morfologická DB, retrográdny, MWE, WordNetSK)
- Určenie autorstva textu
- ...



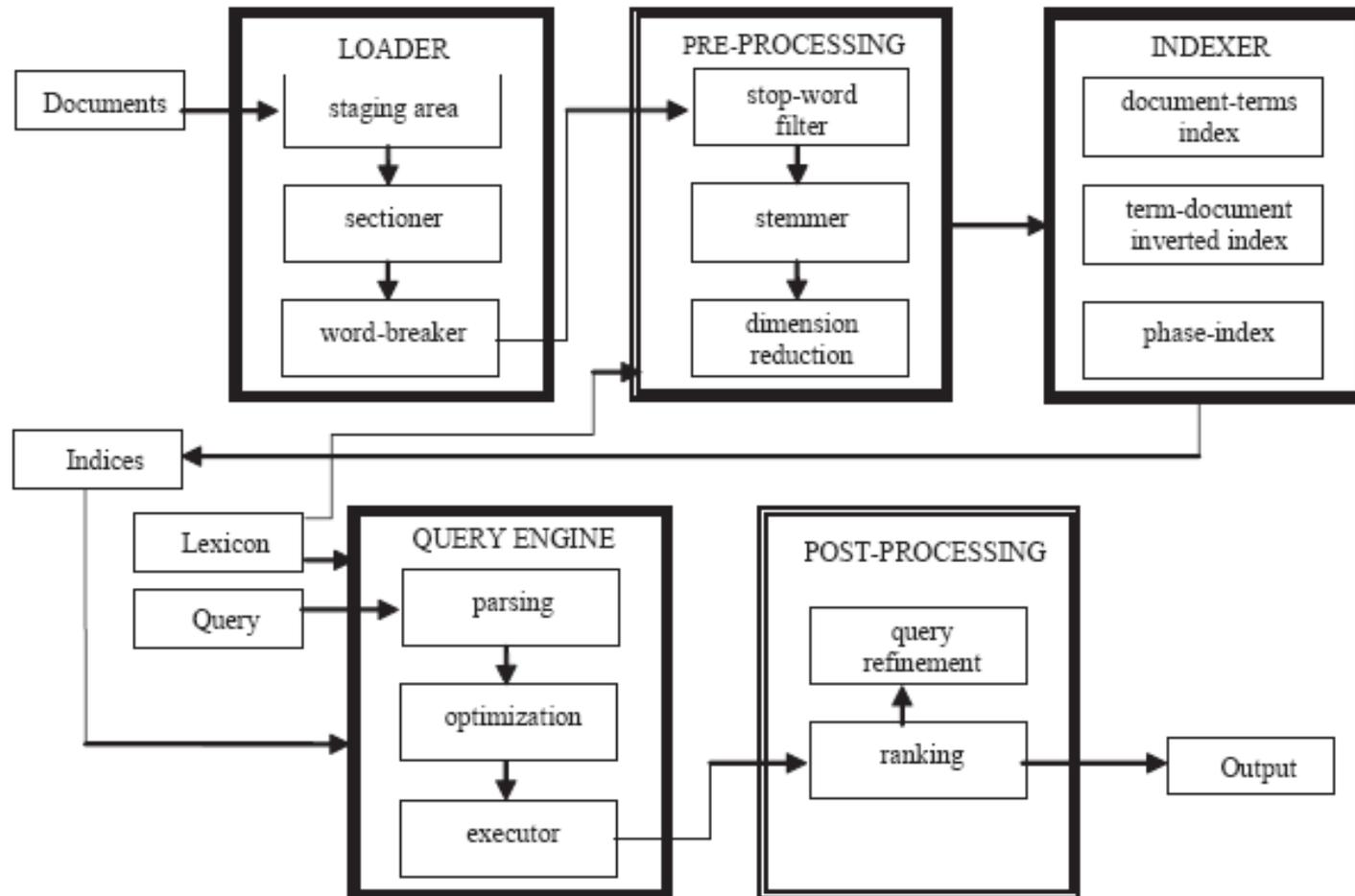
# Témy

- Budovanie morfolologickej databázy:
  - harvesting dát
  - určenie morfologických atribútov
- Slovníky (koreňových morfém, retrográdny)
- Kontrola pravopisu a gramatiky
- Viacslovné pomenovania (MWE)
- WordNetSK
- Prepis numerálov
- Hodnotenie krátkych odpovedí
- Určenie autorstva textu
- Adaptácia nástrojov pre slovenčinu
- Analýza sentimentu („nálady“, postoja)
- ...

# Schéma aplikácie text mining-u



# FULL-TEXT SEARCH (FTS) ENGINES



# Morfologická databáza

otec Nájdí Stratégia: Presne v:  KSSJ4  PSP  SSSJ-AG  SSSJ-HL  SCS  SSS  SSJ  MA  HSSJ-V  Bernolák  Obce  Príezviská  UN  locutio  sk-cs  sk-en  substantiva

Výsledky hľadania pre: otec

## Paradigmy podstatných mien. [Viac informácií](#)

### Otec

mužský rod, životné, jednotné číslo, substantívna paradigma

1	(jeden)	Otec	ruku a obližla si prsty Soľ! <b>Otec</b> jej dokonca priniesol i more.	m++++
2	(bez)	Otca	etnické čistky. A tak potomci <b>Otca</b> vlasti čistia svoju krajinu aj	o++++
3	(k)	Otcu	že láskavo všetko zmieruješ, i <b>Otcu</b> , Duchu Svätému, s ktorými večne	o++
3	(k)	Otcovi	Amerike, kam sa vybral za prácou. <b>Otcovi</b> však jeho matka, moja stará mať	o++++
4	(vidím)	Otca	prekrásnymi svetlými vlasmi. <b>Otca</b> si pamätal neveľmi jasne ako	m++++
5	(nej)	Otče!	dokážeme odniekať bez chýb. „ <b>Otče</b> , nauč ma, čo vieš. Chcem	m++++
5	(nej)	Otecl	nejakom afekte, teda napr. aj <b>Otec</b> ! Výborné! Pohni sa!... Vráťme	o
6	(o)	Otcovi	Červeňák). Lenže Gorkij knihu o " <b>Otcovi</b> " nenapísal, vyhováral sa na	o+++
7	(s)	Otcom	kompozíciu pre syntetizátor Moog <b>Otcom</b> skutočne použiteľných	o++++

### otec

mužský rod, životné, jednotné číslo, substantívna paradigma

1	(jeden)	otec	ovsom okolo hrdla, „vyhral tvoj <b>otec</b> trojitú stávkú?“ Synovec prikývol.	o
2	(bez)	otca	ktorá prýštila zo smútiaceho <b>otca</b> , sa k nemu mlčky pripojili.	o
3	(k)	otcu	Gerlachovského štítu. No pánu <b>otcu</b> prisahám, neraz sa vám tam budú	o
3	(k)	otcovi	je koniec? Marek to vedel, ale <b>otcovi</b> nedokázal povedať pravdu. Báľ	o
4	(vidím)	otca	ho, že ju vidí pri posteli, aj <b>otca</b> , a zdalo sa mu to čudné, veď	o
5	(nej)	otče!	to trest Boží!“ zaúpel. „Začo, <b>otče</b> ?“ „Za to, že sme klamali, pani	o
5	(nej)	otec!	doteraz rád. Prečo odchádzaš, <b>otec</b> môj? Prečo ma práve ty nemáš rád?	o
6	(o)	otcovi	do zelených vln a rozmýšľala o <b>otcovi</b> . A vtom som si uvedomila, že	o
7	(s)	otcom	ružomerský advokát. So starým <b>otcom</b> sa zblížili znova po mojom narodení	o

mužský rod, životné, množné číslo, substantívna paradigma

1	(dvaja)	otcovia	táto kultúra robí. Tak ako ich <b>otcovia</b> , aj oni sa stali obľúbeným terčom	o
2	(bez)	otcov	bábáčka a zlomených, plačúcich <b>otcov</b> nad hrobmi, svet naplnený bolesťou	o
3	(k)	otcom	ale hlavne čas: čo sa nepodarilo <b>otcom</b> , podarilo sa synom a vnukom,	o
4	(vidím)	otcov	pripomínal raných kresťanských <b>otcov</b> , ktorí verili, že pripútanosť	o
6	(o)	otcoch	riadny človek. Podobne sa o svojich <b>otcoch</b> vyjadrovali aj iní. Taký bol vtedy	o
7	(s)	otcami	(1668 – 1744), ale jeho pravými <b>otcami</b> sa stali Johann Gottfried von	o

# Morfologická databáza (2)

abakus abakus SSis1  
abakus abakusu SSis2  
abakus abakusu SSis3  
abakus abakus SSis4  
abakus abakus SSis5  
abakus abakuse SSis6  
abakus abakusom SSis7  
abakus abakusy SSip1  
abakus abakusov SSip2  
abakus abakusom SSip3  
abakus abakusy SSip4  
abakus abakusy SSip5  
abakus abakusoch SSip6  
abakus abakusmi SSip7

zvlúdniet->VKdpa+  
zvlúdniet->VKdpb+  
zvlúdniet->VKdpc+  
zvlúdniet->VKdsa+  
zvlúdniet->VKdsb+  
zvlúdniet->VKdsc+  
zvlúdniet->VLdpah+  
zvlúdniet->VLdpbh+  
zvlúdniet->VLdpcf+  
zvlúdniet->VLdpci+  
zvlúdniet->VLdpcm+  
zvlúdniet->VLdpcn+  
zvlúdniet->VLdsaf+  
zvlúdniet->VLdsai+  
zvlúdniet->VLdsam+  
zvlúdniet->VLdsan+  
zvlúdniet->VLdsbf+  
zvlúdniet->VLdsbi+  
zvlúdniet->VLdsbm+  
zvlúdniet->VLdsbn+  
zvlúdniet->VLdscf+  
zvlúdniet->VLdsci+  
zvlúdniet->VLdscm+  
zvlúdniet->VLdscn+  
zvlúdniet->VMdpa+  
zvlúdniet->VMdpb+  
zvlúdniet->VMdsa+

zvlúdnieme  
zvlúdniete  
zvlúdnejú  
zvlúdniem  
zvlúdnieš  
zvlúdnie  
zvlúdneli  
zvlúdneli  
zvlúdneli  
zvlúdneli  
zvlúdneli  
zvlúdneli  
zvlúdnela  
zvlúdnel  
zvlúdnel  
zvlúdnelo  
zvlúdnela  
zvlúdnel  
zvlúdnel  
zvlúdnelo  
zvlúdnela  
zvlúdnel  
zvlúdnelo  
zvlúdnela  
zvlúdnel  
zvlúdnelo  
zvlúdnejme  
zvlúdnejte  
zvlúdnej

# Transkripcia numerálov

## Vstupný text:

Márii Kvopkovej olympijský limit ušiel tesne, lebo skončila na 40. mieste (+25,20) a zo 78 pretekárov, ktoré sa postavili na štart, bolo potrebné zdolať polovicu.

## Medzivýsledok:

Márii Kvopkovej olympijský limit ušiel tesne, lebo skončila na <cislovka typ=radove tvar=samostatne rod=s cislo=s pad=5>40.</cislovka> mieste (<cislovka typ=desatinne tvar=samostatne rod= pad=>+25,20</cislovka>) a zo <cislovka typ=cele tvar=samostatne rod=ž pad=2>78</cislovka> pretekárov, ktoré sa postavili na štart, bolo potrebné zdolať polovicu.

## Výstupný text:

Márii Kvopkovej olympijský limit ušiel tesne, lebo skončila na štyridsiatom mieste (plus dvadsaťpäť celých dvadsať stotín ) a zo sedemdesiatich ôsmich pretekárov, ktoré sa postavili na štart, bolo potrebné zdolať polovicu.

# Transkripcia numerálov (2)

## Vstupný text:

Diskusia: 8 - posl. 27.4.2009 08:16 od: Martin

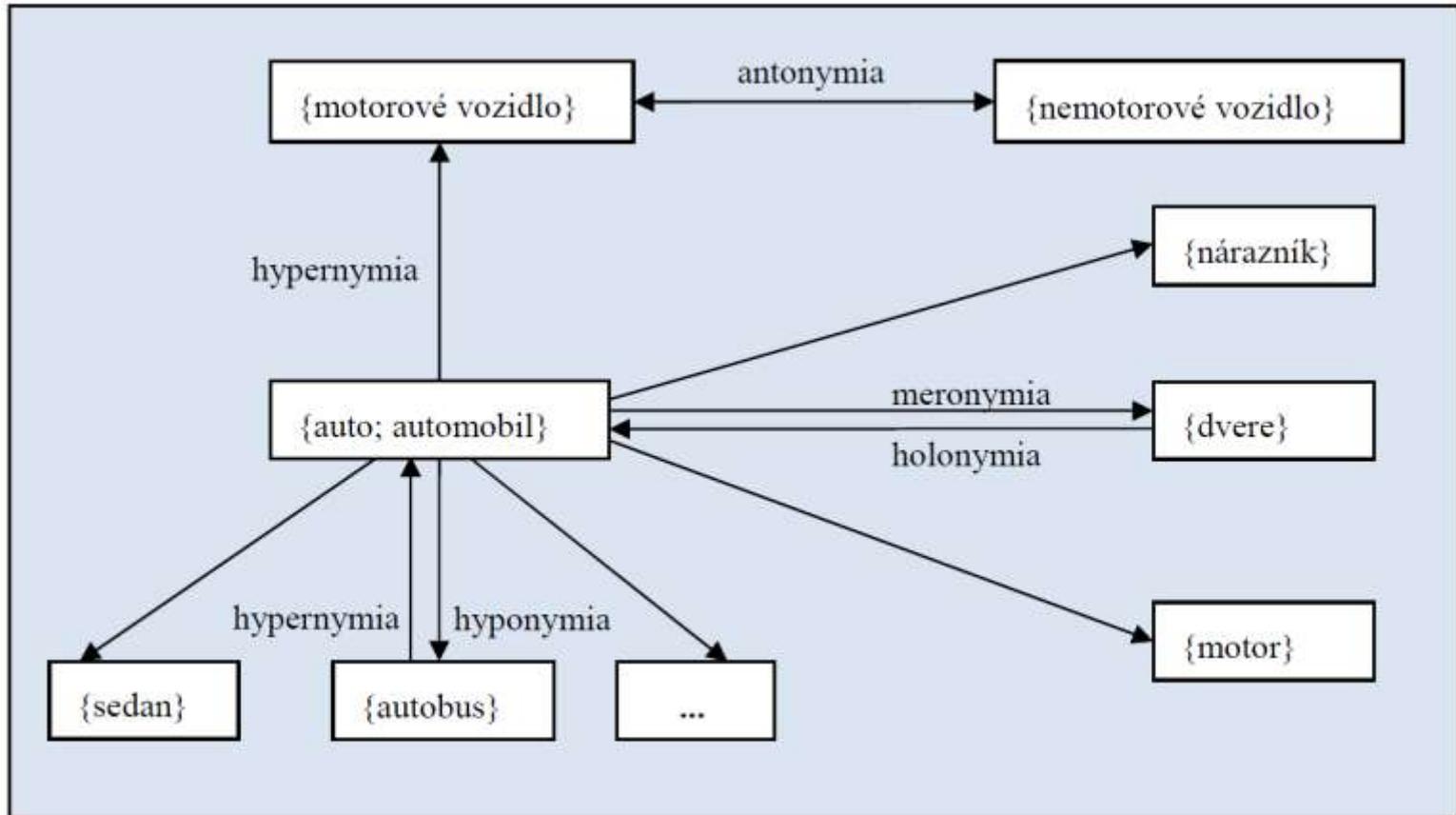
## Medzivýsledok:

Diskusia: <cislovka typ=cele tvar=samostatne rod= pad=>8</cislovka> - posl. <cislovka typ=datum tvar=samostatne rod=m pad=>27.4.2009</cislovka> <cislovka typ=cas tvar=samostatne format=HHMM rod=ž pad=>08:16</cislovka> od: Martin

## Výstupný text:

Diskusia: osem - posl. dvadsiaty siedmy apríl dvetisícdeväť osem hod. šestnásť min. od: Martin

# WordNet SK

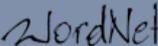


Obr. 1: Základné vzťahy nachádzajúce sa vo WordNet-e

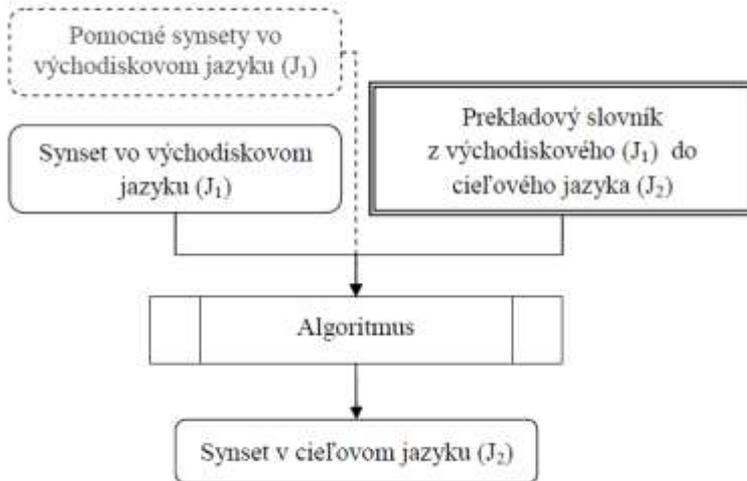
# SynsetBuilder

**Synsets Builder 1.0**

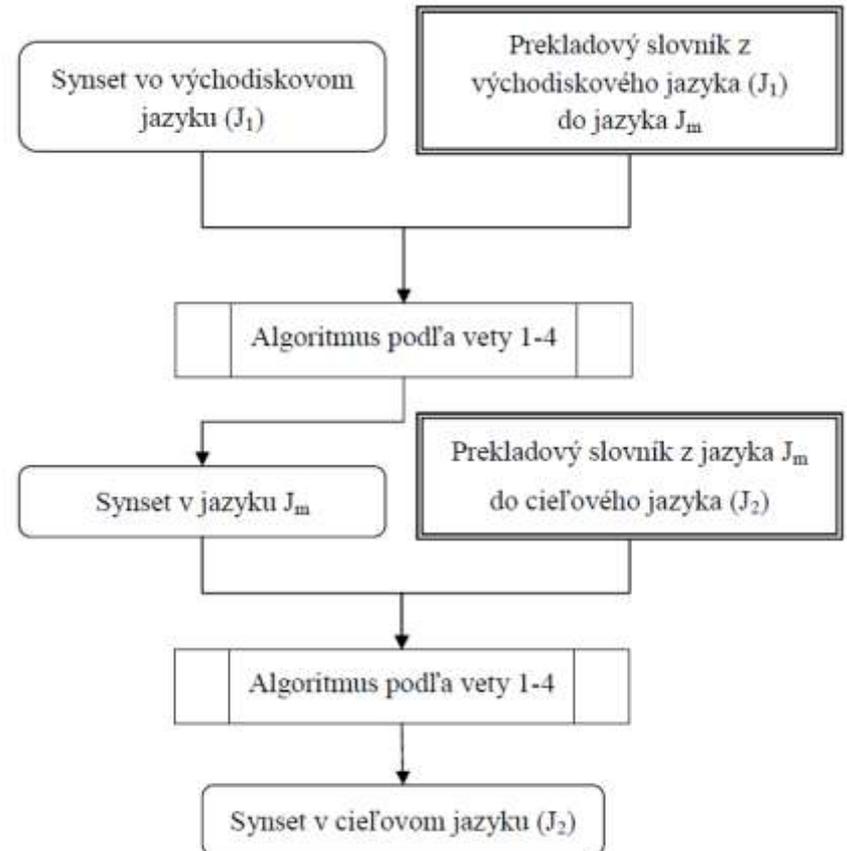
Hľadanie	Štatistika slovníka	Prihlásenie																								
<p><b>Zadanie pojmu</b> </p> <p>Anglický pojem <input type="text"/></p> <p>Slovník: <input type="text" value="slovník.azet.sk   en-&gt;sk"/> </p> <p><input type="button" value="Search"/></p>	<p style="text-align: center;"><b>Počet synsetov</b></p> <table border="1"><thead><tr><th>Slovný druh / jazyk</th><th>anglické</th><th>slovenské</th><th>české</th></tr></thead><tbody><tr><td>Podstatné mená:</td><td>1422</td><td>90</td><td>63</td></tr><tr><td>Slovesá:</td><td>819</td><td>95</td><td>16</td></tr><tr><td>Prídavné mená:</td><td>32</td><td>14</td><td>8</td></tr><tr><td>Príslovky:</td><td>2</td><td>1</td><td>2</td></tr><tr><td>Spolu:</td><td>2275</td><td>200</td><td>89</td></tr></tbody></table> <p style="text-align: right;">Spracované za 0 sek.</p>	Slovný druh / jazyk	anglické	slovenské	české	Podstatné mená:	1422	90	63	Slovesá:	819	95	16	Prídavné mená:	32	14	8	Príslovky:	2	1	2	Spolu:	2275	200	89	<p><b>Prihlásenie</b> </p> <p>Login <input type="text"/></p> <p>Heslo <input type="password"/></p> <p><input type="button" value="Login"/></p> <p><b>Info</b> </p> <p> Dokumentácia  Pomoc</p> <p> Odkazy a kontakt</p>
Slovný druh / jazyk	anglické	slovenské	české																							
Podstatné mená:	1422	90	63																							
Slovesá:	819	95	16																							
Prídavné mená:	32	14	8																							
Príslovky:	2	1	2																							
Spolu:	2275	200	89																							

# SynsetBuilder - princíp



Obr. 5: Všeobecný postup pri priamom budovaní synsetov v cieľovom jazyku



Obr. 7: Postup pri budovaní synsetov v cieľovom jazyku cez medzijazyk

# SynsetBuilder - výsledky

Slovo car – vyplnené 2 synsety z 5:

<p><b>02853224: car; auto; automobile; machine; motorcar</b>  <i>4-wheeled motor vehicle; usually propelled by an internal combustion engine; "he needs a car to get to work"</i></p>	<p>vozidlo; automobil; motorové vozidlo; auto; automobilový; osobný automobil</p>
<p><b>02854760: car; railcar; railway car; railroad car</b>  <i>a wheeled vehicle adapted to the rails of railroad; "three cars had jumped the rails"</i></p>	<p>železničný vozeň</p>

Slovo computer – vyplnené 2 synsety z 2:

<p><b>02971359: computer, computing machine, computing device, data processor, electronic computer, information processing system</b>  <i>a machine for performing calculations automatically</i></p>	<p>počítač</p>
<p><b>09257296: calculator, reckoner, figurer, estimator, computer</b>  <i>an expert at calculation (or at operating calculating machines)</i></p>	<p>kalkulant; počítač; kalkulátor</p>

Slovo information – vyplnené 3 synsety z 5:

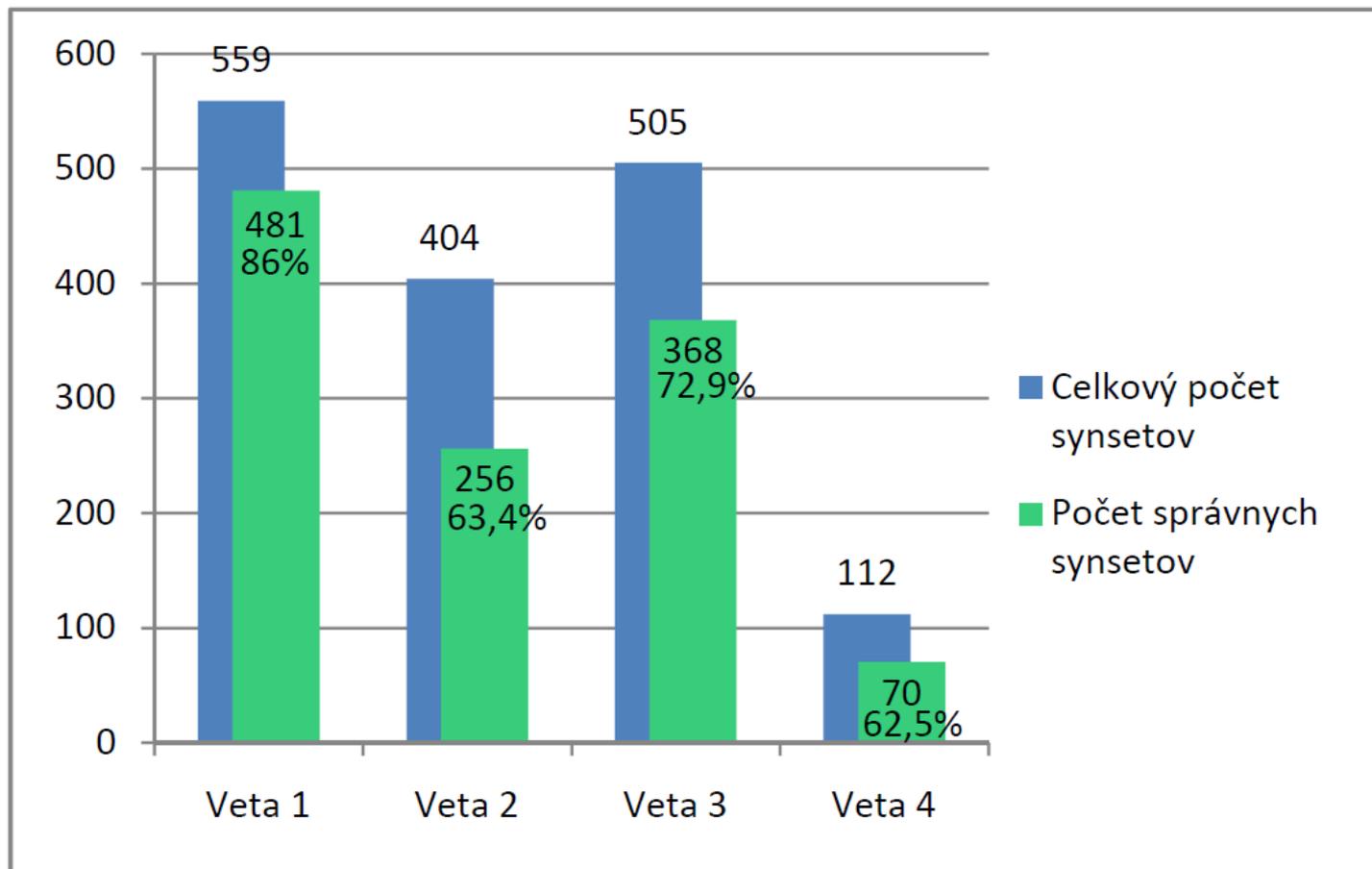
<p><b>06225142: information, info</b>  <i>a message received and understood</i></p>	<p>správa; informácia</p>
<p><b>07949563: data, information</b>  <i>a collection of facts from which conclusions may be drawn; "statistical data"</i></p>	<p>informačný</p>
<p><b>05479334: information</b>  <i>knowledge acquired through study or experience or instruction</i></p>	<p>informácia, vedomosť, znalosť</p>

# Synset Builder - štatistiky

**Tab. 10: Úspešnosť celkového generovania slovenských synsetov**

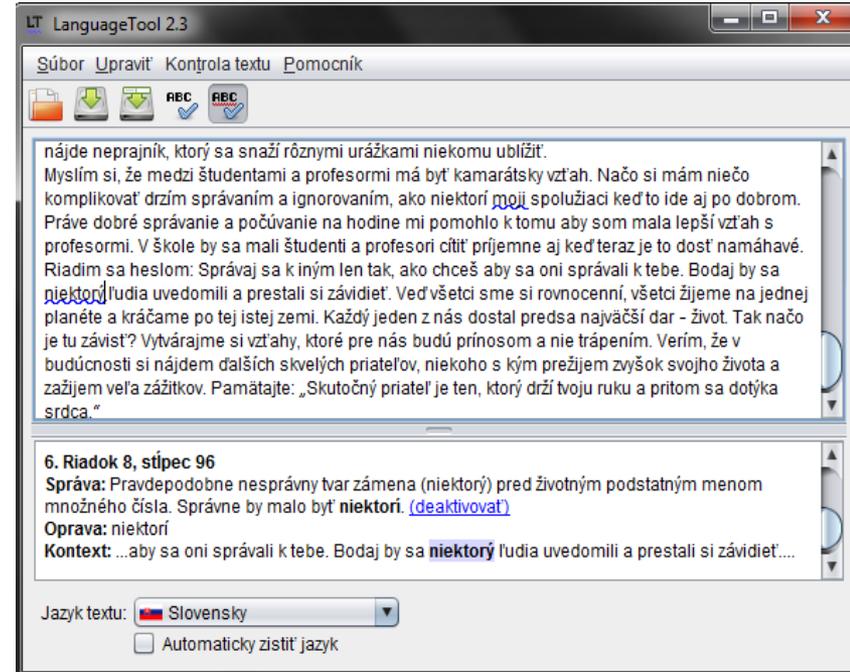
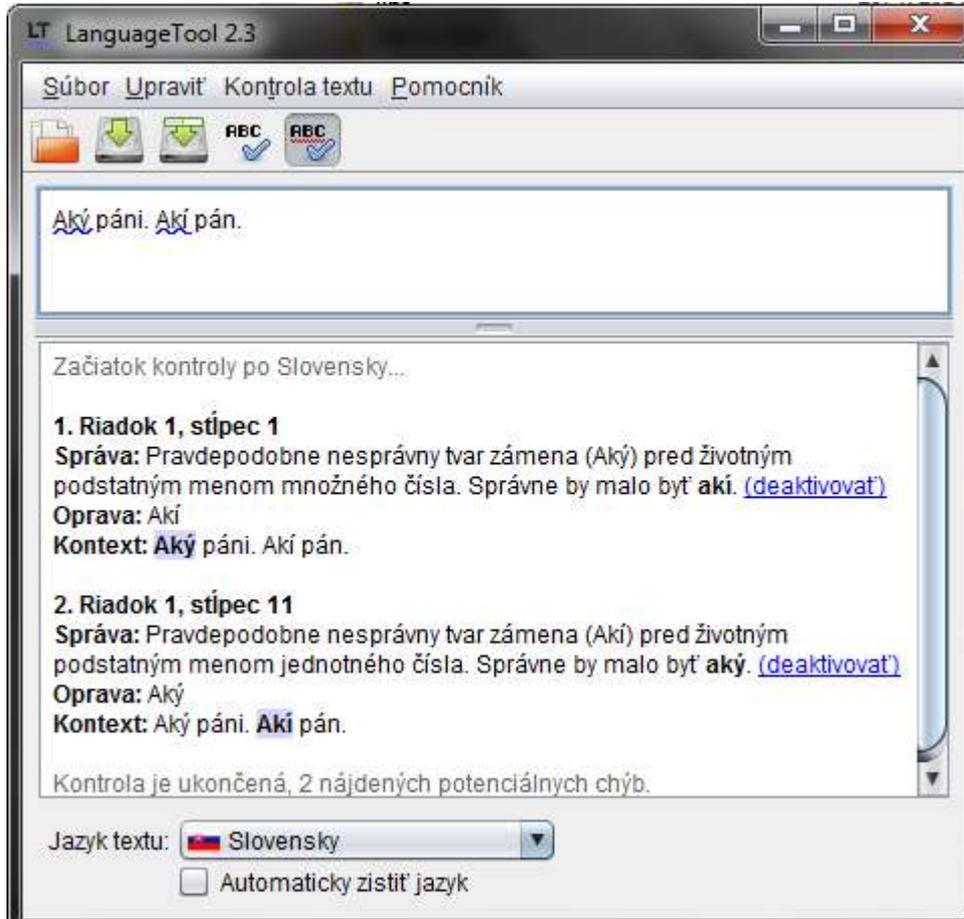
	Všetky	Podstatné mená	Pridavné mená	Slovesá	Príslovky
Počet EN synsetov vo WordNet-e	117 659	82 115	18 156	13 767	3 621
Počet EN synsetov, s nagerovaným slovenským synsetom	40 521 (34,4%)	26 787 (32,6%)	6 859 (37,8%)	5 839 (42,4%)	1 036 (28,6%)
Počet SK synsetov vytvorených vetou 1	10 267 (8,7%)	5 705	2 175	2 109	278
Počet SK synsetov vytvorených vetou 2	30 243 (25,7%)	20 510	6 059	2 715	959
Počet SK synsetov vytvorených vetou 3	11 533 (12%) <sup>3</sup>	8 192	-	3 341	-
Počet SK synsetov vytvorených vetou 4	1 917 (1,4%) <sup>3</sup>	1 348	-	569	-

# Synste Builder - štatistiky



**Obr. 11: Pomer správnych slovenských synsetov k ich celkovému počtu vo vzorke**

# Kontrola gramatiky



# Určenie autorstva textu

	Blaha	Chren	Hudacky	Kanik	Matovic	Micovsky	Mikus	Pollack	Tomanova	Vasecka
Blaha	9.74	14.03	10.50	21.45	19.15	12.58	18.25	10.91	11.76	
Chren	10.60	11.52	7.32	13.18	13.42	9.48	10.09	7.64	10.92	
Hudacky	10.32	13.78	7.80	6.88	12.97	10.03	11.02	11.43	7.32	
Kanik	6.52	11.79	7.33	5.28	12.88	10.36	10.52	10.87	6.39	
Matovic	18.22	18.27	15.75	14.29	9.82	11.30	9.57	20.29	14.35	
Micovsky	10.03	12.21	9.26	7.43	12.05	7.45	8.97	15.10	8.38	
Mikus	11.28	13.53	9.02	8.43	11.03	7.20	9.08	12.98	7.91	
Pollack	11.60	10.01	8.43	9.31	9.15	8.01	6.87	10.24	11.31	
Tomanova	6.95	12.68	7.33	7.09	16.97	14.19	7.43	13.83	10.80	
Vasecka	9.74	12.63	8.36	19.48	17.93	12.41	15.78	9.55	12.27	

Obrázok 2 Frekvencia čiarok, politické rozpravy

	Blaha	Chren	Hudacky	Kanik	Matovic	Micovsky	Mikus	Pollack	Tomanova	Vasecka
Blaha	15.00	15.83	14.42	11.34	12.70	11.39	12.01	11.69	9.91	10.00
Chren	12.19	7.87	14.56	10.33	10.01	12.58	13.37	12.54	11.38	10.00
Hudacky	13.81	13.35	12.79	8.90	12.84	9.92	13.58	12.22	11.03	10.00
Kanik	9.58	12.84	12.88	7.94	8.78	10.93	9.60	10.98	10.57	9.42
Matovic	12.13	13.49	12.92	10.64	9.31	10.24	10.44	9.47	12.65	10.69
Micovsky	12.30	12.95	8.61	9.25	8.36	10.61	8.91	8.20	9.22	10.00
Mikus	10.48	13.35	8.61	10.48	9.90	7.68	11.05	11.35	9.43	10.00
Pollack	10.37	10.36	13.38	11.00	9.00	10.88	10.46	9.93	10.83	10.00
Tomanova	10.10	12.39	10.83	8.97	9.57	11.33	8.93	10.83	9.30	10.00
Vasecka	9.48	11.88	13.38	9.99	11.28	11.53	11.83	10.45	9.63	10.00

Obrázok 4 Dĺžka viet, politické rozpravy

	Blaha	Chren	Hudacky	Kanik	Matovic	Micovsky	Mikus	Pollack	Tomanova	Vasecka
Blaha	3.75	4.71	5.20	4.36	4.65	4.97	4.16	4.65	5.15	3.80
Chren	4.18	3.85	4.58	6.15	4.75	5.11	4.23	4.42	4.49	3.88
Hudacky	6.32	5.22	3.15	7.31	5.35	5.10	5.05	4.90	6.43	5.48
Kanik	4.51	4.30	6.17	5.11	5.24	4.63	5.57	4.98	6.46	3.25
Matovic	5.02	4.34	5.31	6.47	3.74	3.73	4.76	4.61	6.32	4.45
Micovsky	4.75	4.21	5.28	5.25	3.19	3.46	4.95	4.05	5.98	3.93
Mikus	4.84	4.02	4.51	6.41	4.51	4.83	3.88	4.19	4.77	4.01
Pollack	4.49	2.92	3.67	5.32	4.41	4.31	3.32	3.88	4.41	3.19
Tomanova	4.33	4.27	5.52	7.19	6.48	7.21	4.28	5.23	3.88	4.79
Vasecka	3.90	3.67	5.71	5.35	3.64	3.38	4.84	4.41	4.66	3.58

Obrázok 6 Stop slová, politické rozpravy

n =	Blaha	Chren	Hudacky	Kanik	Matovic	Micovsky	Mikus	Pollack	Tomanova	Vasecka
Blaha	1.54	4.18	1.94	2.95	5.31	4.24	1.62	4.50	1.33	2.13
Chren	1.54	1.78	1.06	0.88	2.74	2.04	1.27	1.85	0.84	0.99
Hudacky	1.52	2.16	1.06	1.10	3.13	2.25	0.83	2.33	0.70	0.72
Kanik	1.82	1.40	1.13	1.32	2.49	1.59	1.36	1.68	1.09	1.01
Matovic	3.18	3.54	2.48	1.50	1.00	0.75	2.78		2.51	2.48
Micovsky	2.49	0.79	1.07	1.10	1.90	1.04	2.24	1.25	2.14	2.27
Mikus	0.91	2.22		1.11	3.37	2.20	0.42	2.56	0.61	1.04
Pollack	2.76		2.05	1.09	1.48	0.82	2.34	0.61	2.07	2.06
Tomanova	0.90	3.06	1.01	1.90	4.20	3.23		3.38	0.79	1.04
Vasecka	1.36	2.98	1.08	1.80	3.98	3.15	0.92	-3.20	0.70	0.88

Obrázok 8 Výskyt najfrekventovanejších slov, politické rozpravy

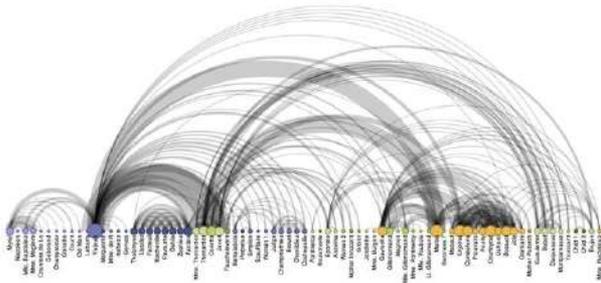
# Triedenie

- Podľa kódu
- Podľa normy
- Retrográdne

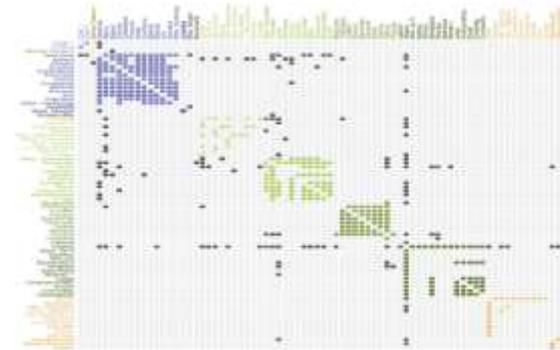
Vizualizácia



# Príklady vizualizácie (2)

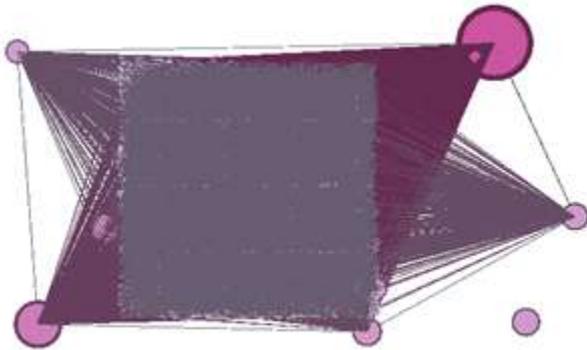
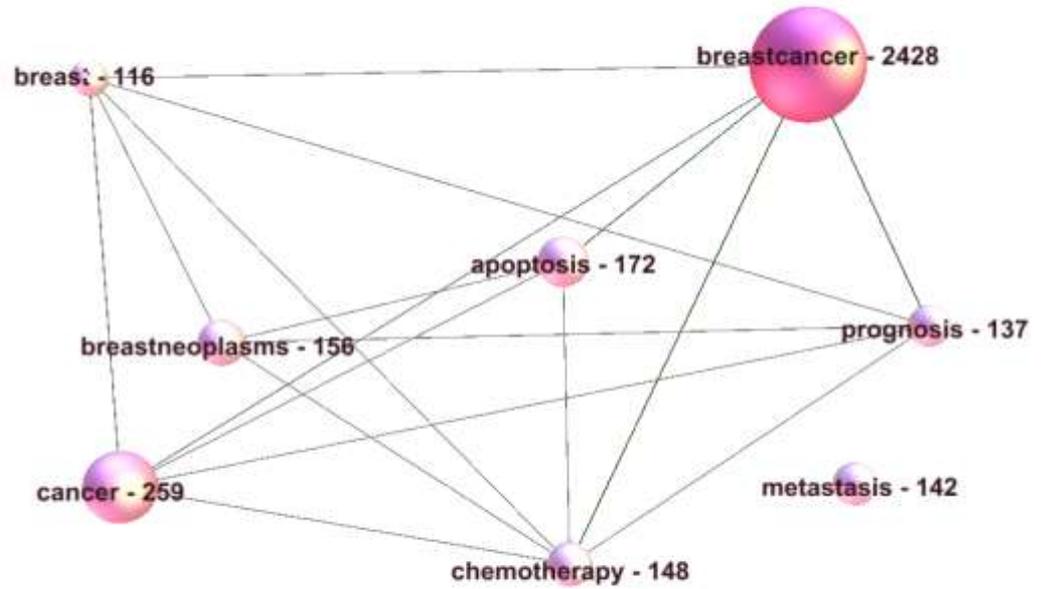
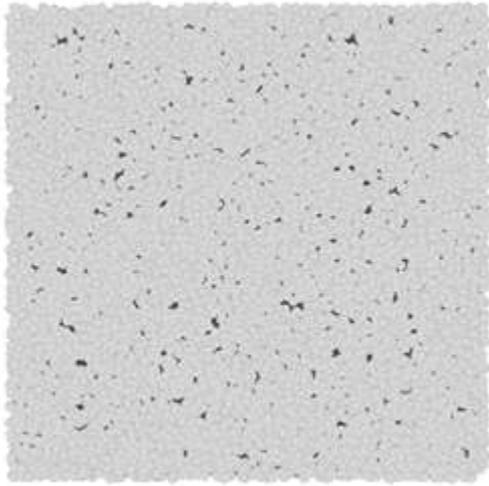


Obr. 5 Príklad vizualizácie slov v kapitole dokumentu pomocou oblúčkového diagramu [8].

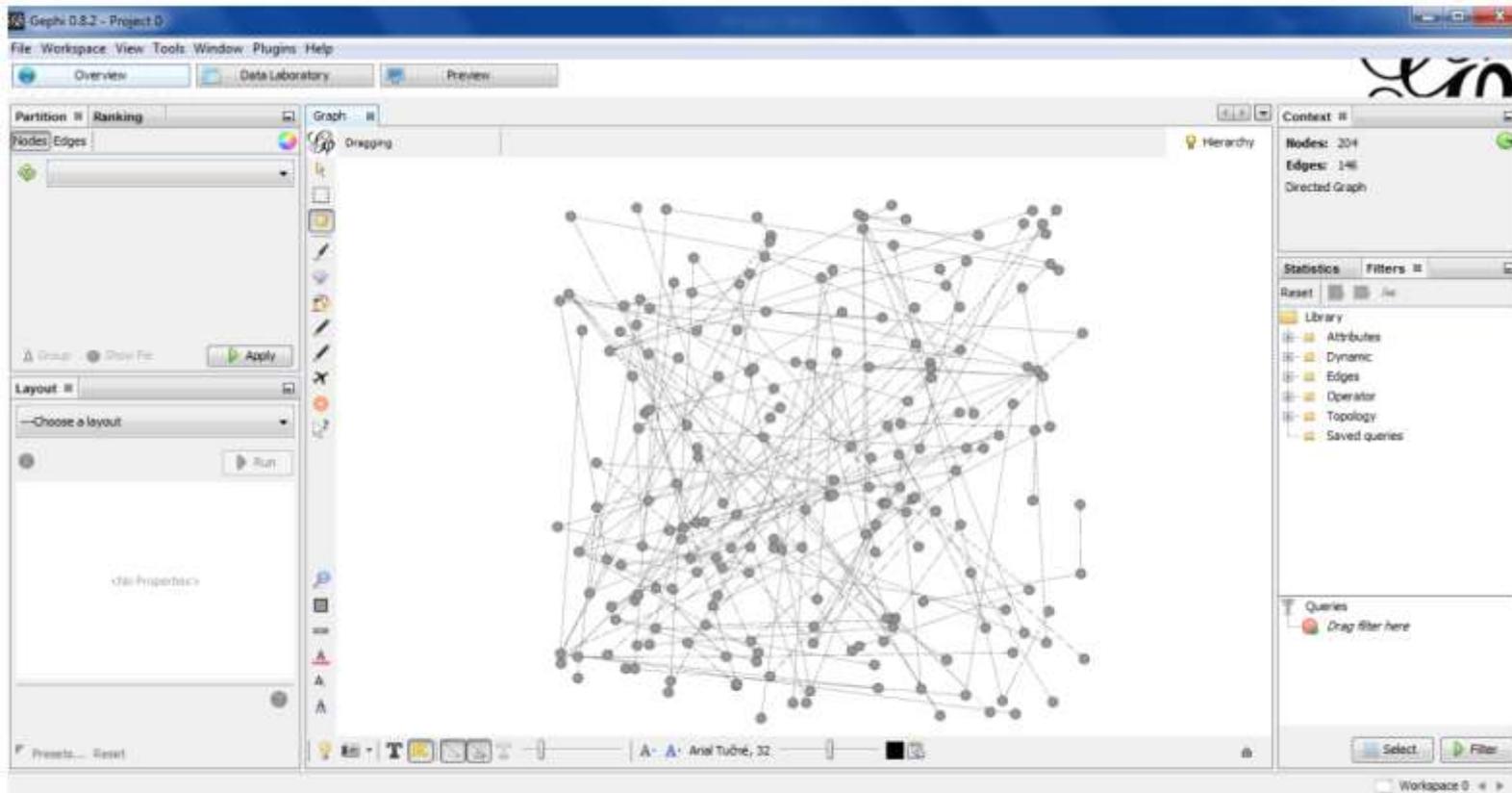


Obr. 6 Príklad maticovej vizualizácie slov v dokumentoch [9].

# Vizualizácia kľúčových slov (1)

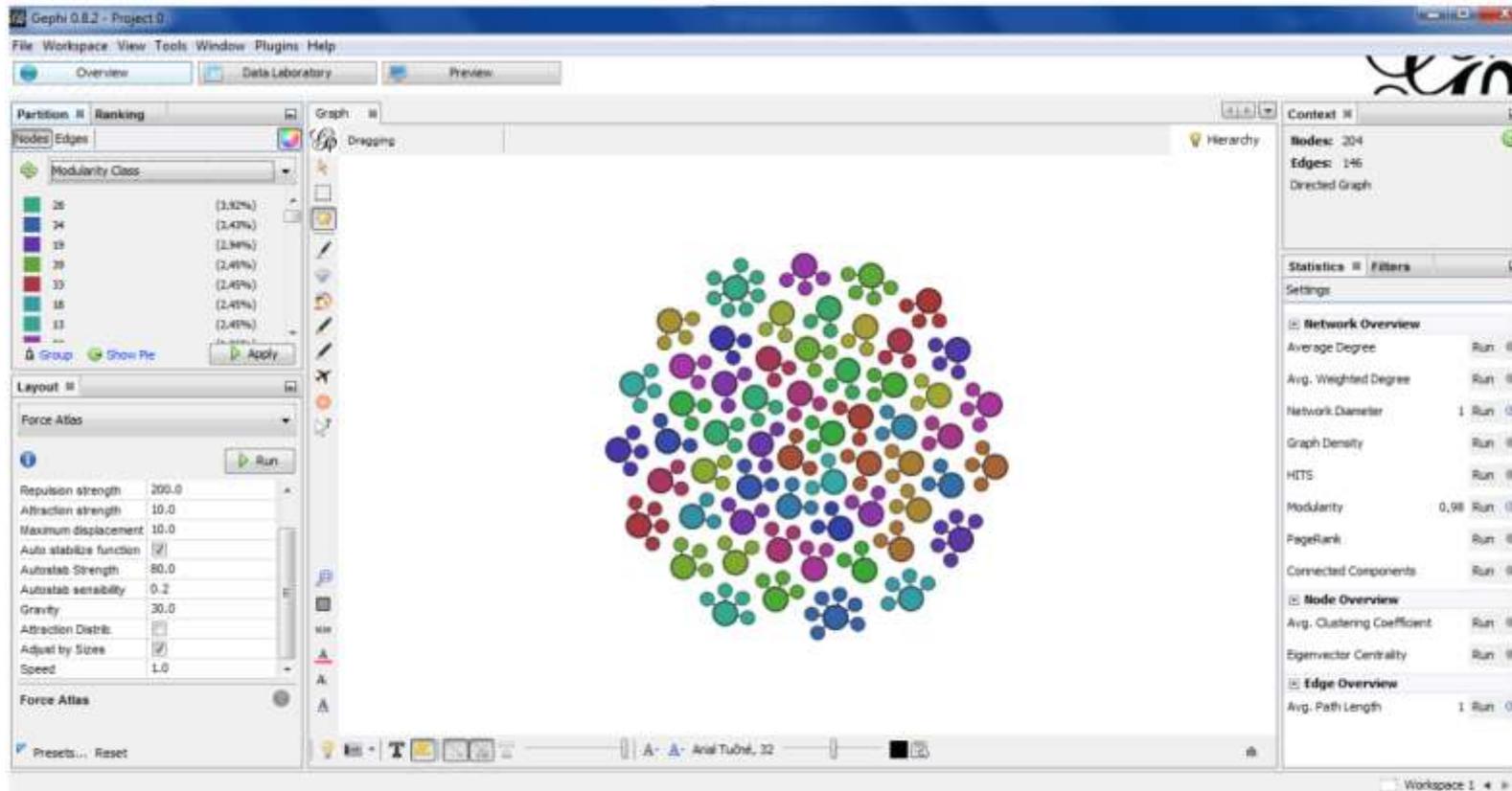


# Vizualizácia kľúčových slov (2)



Obr. 15 Vzhľad grafu po importovaní dát bez akýchkoľvek ďalších úprav.

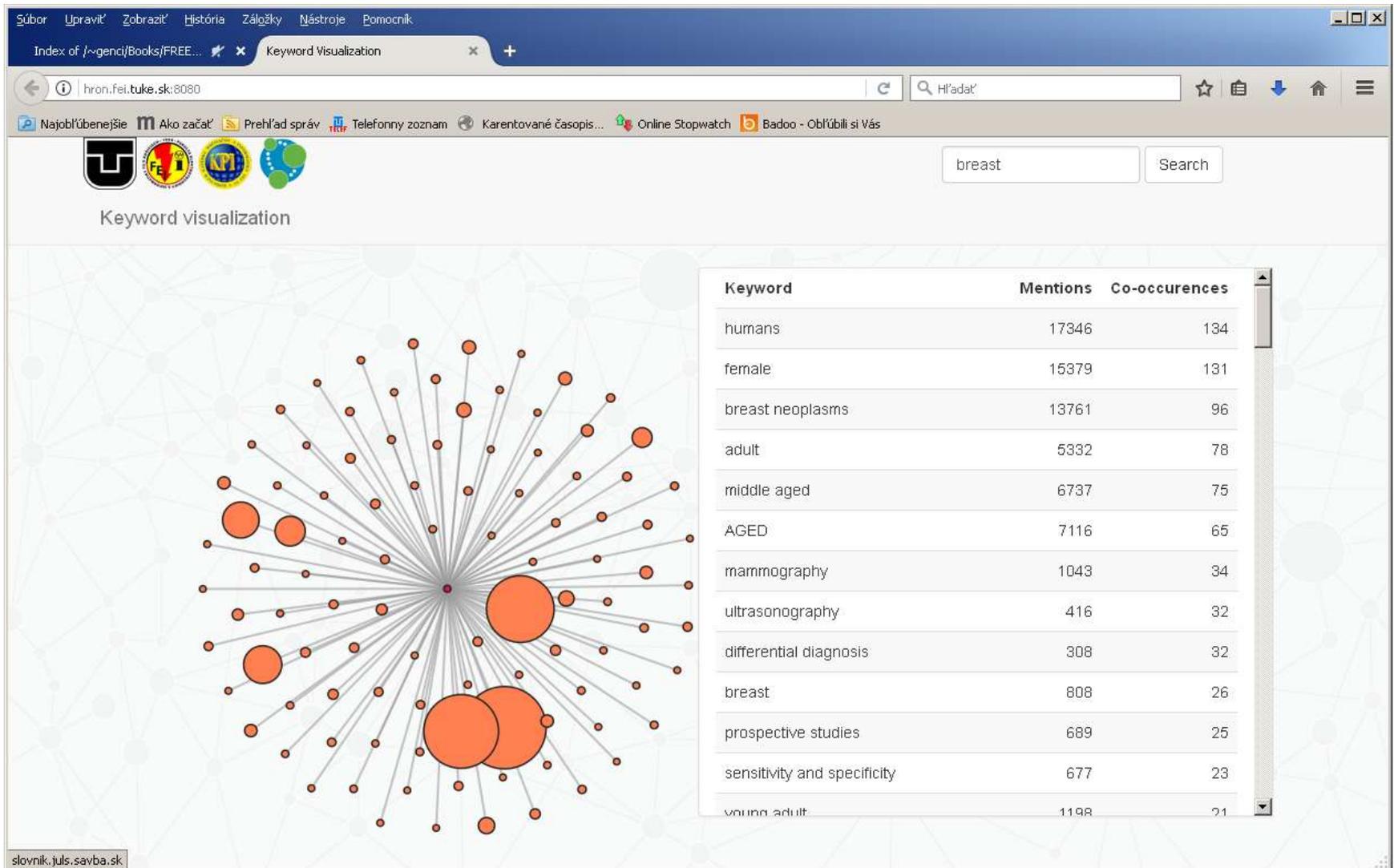
# Vizualizácia kľúčových slov (3)



Obr. 18 Farebne odlišené uzly grafu



# Vizualizácia kľúčových slov (5)



# Viacslovné pomenovania, MWE

- Slovné konštrukcie, ktoré majú iný význam ako slová z ktorých sa skladajú. Problém (niekedy?) pri prekladoch:
  - starý otec -> old father (grandfather)
  - vysoká škola -> high school (university)
  - dať pokoj -> give rest (give peace)
- Problém identifikácie:
  - **daj** mi už dnes, prosím ťa, konečne **pokoj**

# Viacslovné pomenovania (korpusové štatistiky)

Tab. 4 Manuálne spracované údaje pre spojenia obsahujúce slovo otec

<i>bigram</i>	<i>f(xy)</i>	<i>f(x)</i>	<i>f(y)</i>	<i>PMI</i>	<i>AJ-slovník</i>
nebeský otec	1047	3651	93192	2,741057718	N
svätý otec	10885	44369	93192	2,730414679	A
starostlivý otec	226	1271	93192	2,488156216	N
starý otec	11162	77050	93192	2,420370161	A
milujúci otec	215	1897	93192	2,280584256	N
vzorný otec	96	1191	93192	2,191590743	N
nebohý otec	206	3241	93192	2,084992497	N
duchovný otec	706	14789	93192	1,936980839	A
otec biskup	1111	93192	31464	1,833809557	N
hrdý otec	208	10757	93192	1,70948726	N
otec arcibiskup	310	93192	16478	1,578521493	N
nešťastný otec	92	8674	93192	1,316993534	N
vlastný otec	404	58045	93192	1,305078563	N
šťastný otec	138	33018	93192	1,056189777	N
dobry otec	318	144632	93192	0,873237917	N

# Viacslovné pomenovania

(dáta v Webu)

Bigram	Frekvencia	VBI	T-skóre	Dice
stainless steel	383129	5.84485	608.205	0.940286
obchodné podmienky	446433	5.2125	650.136	0.809389
e mail	402408	4.28945	601.916	0.494587
made of	331733	4.5548	551.457	0.513679
práva vyhradené	153738	6.94351	388.909	0.865647
obchodné podmienky	446433	5.2125	650.136	0.809389
stainless steel	383129	5.84485	608.205	0.940286
e mail	402408	4.28945	601.916	0.494587
made of	331733	4.5548	551.457	0.513679
nákupný košík	179520	5.84878	416.346	0.598913

**Tabuľka:** Bigramy - prých 5 má najvyššiu hodnotu súčinu hodnôt, posledných 5 hodnotu T-skóre

# Viacslovné pomenovania

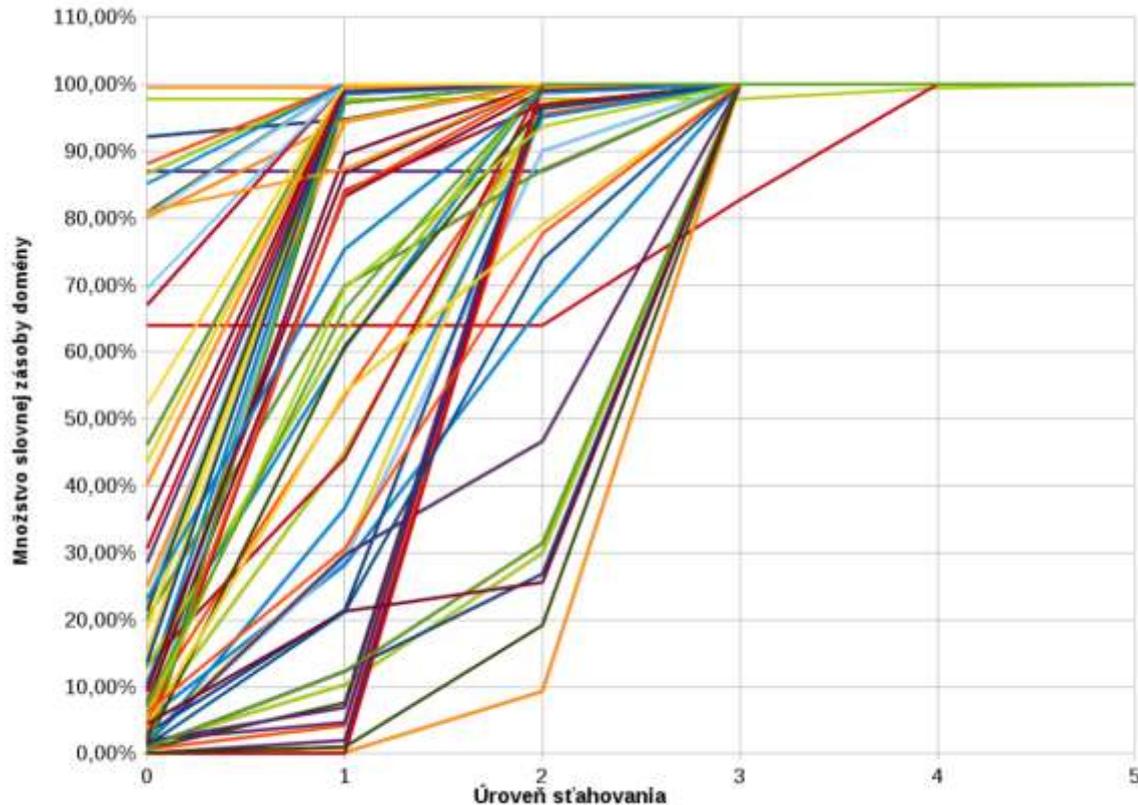
(dáta v Webu)

Trigram	Frekvencia	VBI	T-skóre	Dice
najkrajšej slovenke rozhodovali	3613	1.94527	60.1042	0.95355
petru kvitovú povedie	3237	2.19628	56.8911	0.937355
abdél hakim awyan	2365	3.28054	48.6292	0.979161
receiving personalized recommendations	3178	2.10449	56.3701	0.907395
pomlérun pomlérun pomlérun	2160	3.45609	46.4739	0.96
všetky práva vyhradené	151825	-10.7569	386.805	0.539609
ochrana osobných údajov	86795	-9.84355	293.167	0.491108
vytvorené systémom www	79429	-10.6806	279.196	0.271958
farebný uv gél	72089	-8.67889	267.49	0.596628
aisi stainless steel	72130	-11.8715	262.016	0.238011

**Tabuľka:** Trigramy - prých 5 má najvyššiu hodnotu súčinu hodnôt, posledných 5 hodnotu T-skóre

# Viacslovné pomenovania

(dáta v Webu)



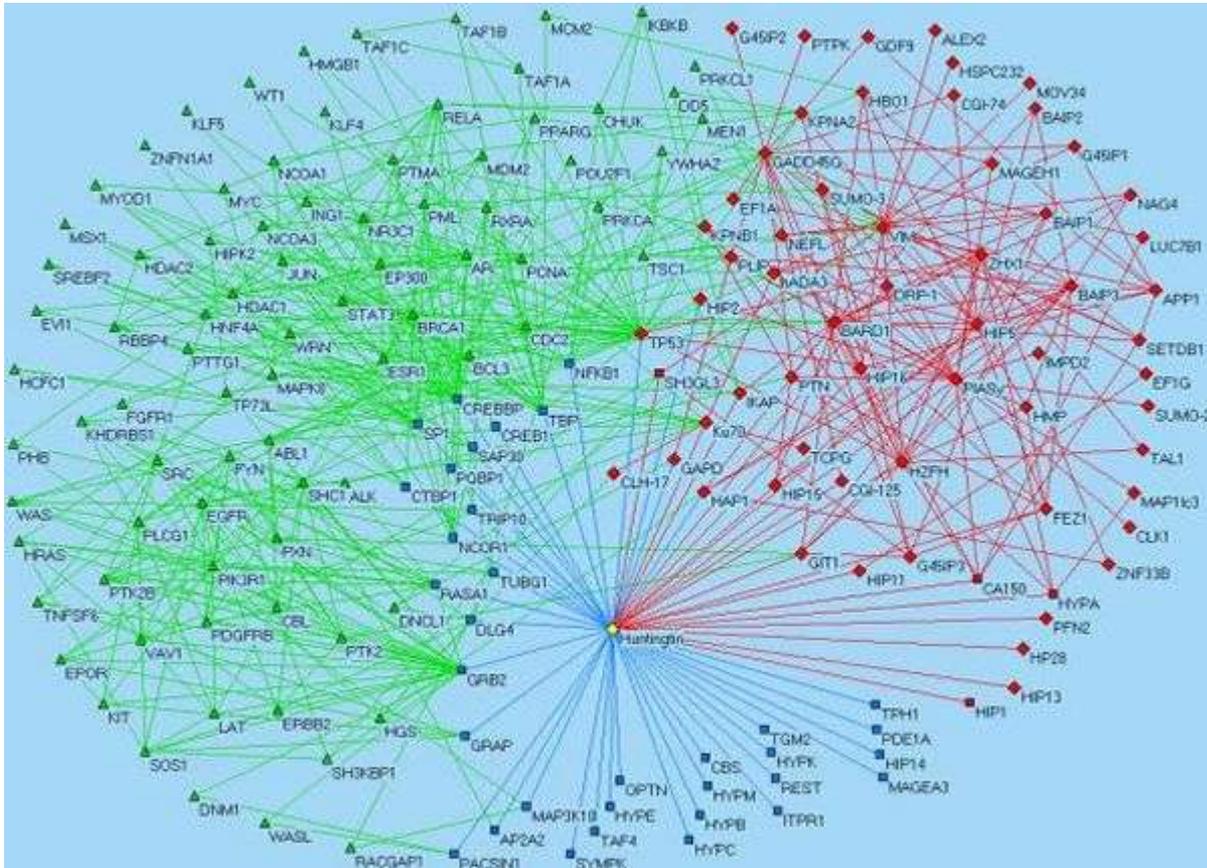
Obr.: Množstvo slovnej zásoby k úrovni sťahovania

# NLP a Bioinformatika

# Budovanie „informačných“ modelov

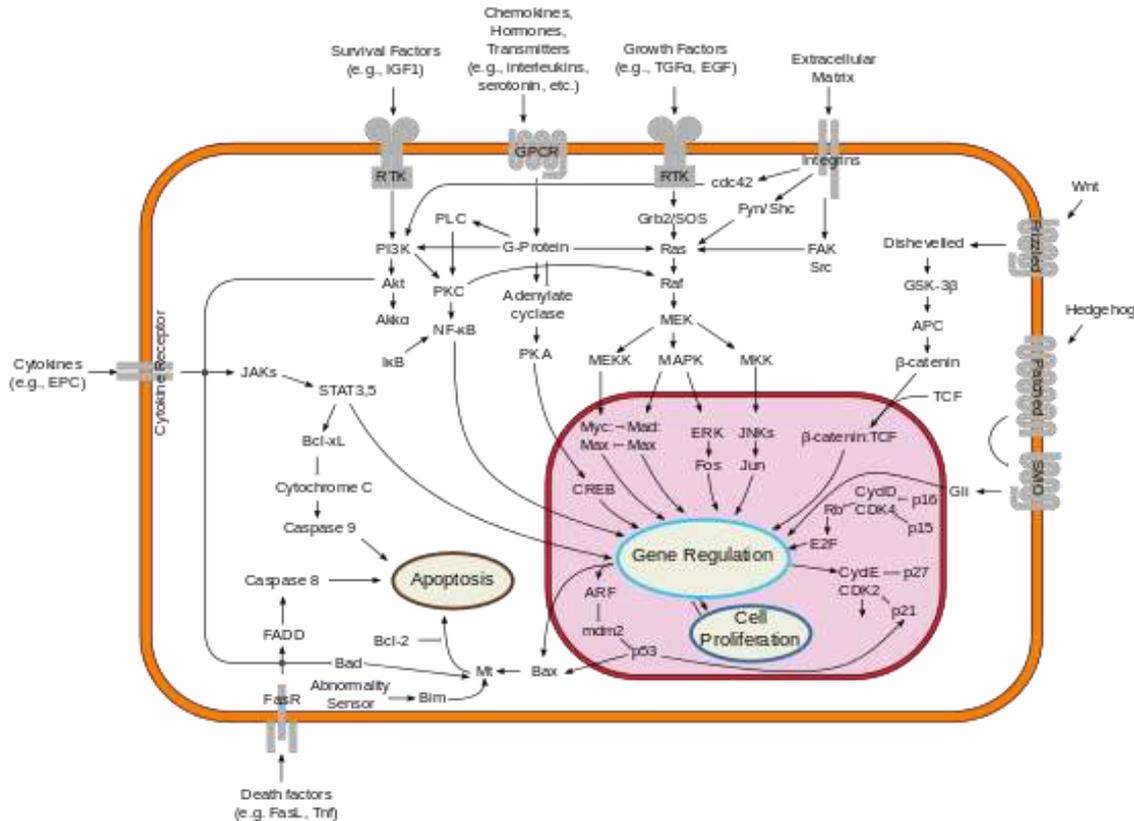
- Spracovanie publikácií
  - breast cancer od 1.1.2014
    - v databáze Pubmed – 336 508 článkov
    - Sciencedirect.com – 72 753 článkov
    - PlosOne – 64 035 článkov
- Modely interakcií proteínov
- Signálne dráhy

# PPI



# Protein signal path

chemical or physical signal is transmitted through a cell as a series of molecular events



Zdroj: [https://en.wikipedia.org/wiki/Signal\\_transduction](https://en.wikipedia.org/wiki/Signal_transduction)

Ďakujem za pozornosť!

Otázky?