

# Získávání a zpracování textových korpusů z internetu

Crawling, deduplikace, odstraňování boilerplate

Vít Suchomel

Centrum zpracování přirozeného jazyka  
Fakulta informatiky, Masarykova univerzita

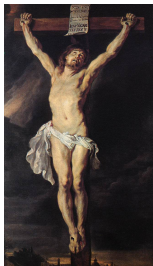
13. března 2017

# Korpus

Co je to korpus?

# Korpus

Co je to korpus?



# Jazykový korpus

- ▶ Rozsáhlý soubor
- ▶ *autentických* textů (psaných nebo mluvených)
- ▶ převedený do *elektronické podoby*,
- ▶ v němž je možné jednoduše vyhledávat *jazykové jevy* (zejména slova a slovní spojení)
- ▶ a zobrazovat je v jejich *přirozeném kontextu*.<sup>1</sup>

---

<sup>1</sup>Zdroj: Český národní korpus, <http://wiki.korpus.cz/doku.php/pojmy:korpus>

# Velikost jazykových korpusů

- ▶ Počet unikátních slov – velikost slovníku z korpusu, *types*
- ▶ Počet výskytů slov – velikost korpusu, *tokens*



# Velikost jazykových korpusů

- ▶ Počet unikátních slov – velikost slovníku z korpusu, *types*
- ▶ Počet výskytů slov – velikost korpusu, *tokens*
- ▶ *Type-Token Ratio*:  $\frac{types}{tokens}$



# Velikost jazykových korpusů

- ▶ Počet unikátních slov – velikost slovníku z korpusu, *types*
- ▶ Počet výskytů slov – velikost korpusu, *tokens*
- ▶ *Type-Token Ratio*:  $\frac{types}{tokens}$
- ▶ Je TTR větší u češtiny nebo angličtiny? Kolikrát?



# Velikost jazykových korpusů

- ▶ Počet unikátních slov – velikost slovníku z korpusu, *types*
- ▶ Počet výskytů slov – velikost korpusu, *tokens*
- ▶ *Type-Token Ratio*:  $\frac{types}{tokens}$
- ▶ Je TTR větší u češtiny nebo angličtiny? Kolikrát?
- ▶ TTR češtiny ku angličtině  $\doteq 1,87$  (český web 2012, anglický web 2013)





# Korpusová lingvistika

- ▶ Věda zkoumající *jazyk* (lingvistika)
- ▶ pomocí *jazykových korpusů*.

Query **korpus** 17,177 > GDEX 17,177 (3.39 per million) 

Page  of 1,718  [Next](#) | [Last](#)

<a href="#">sportovniobchod.cz</a>	Ke <b>korpusu</b> pouze připevníte 4 nohy a stolní fotbal je okamžitě provozu schopný.
<a href="#">kultura21.cz</a>	Ti manuálně zruční mohou sestavit <b>korpus</b> stoličky z dvanácti stejných připravených dílů.
<a href="#">umimeudelatdomov.cz</a>	Také je možná varianta smrkového <b>korpusu</b> s bílými předními dvířky.
<a href="#">ceskahospodynka.cz</a>	Krémem bez ořechů potřeme vrchní <b>korpus</b> a boční stěny.
<a href="#">styleve-kuchyne.cz</a>	<i>&lt;p&gt;</i> Vyšší tloušťka používaných desek <b>Korpusy</b> skříněk jsou vyrobeny z konstrukčních desek o síle 18 mm.
<a href="#">jacques.cz</a>	Všechny dřevěné díly postele jsou z bukového masivu, <b>korpus</b> má tloušťku 30 mm.
<a href="#">ok1blog.blog.cz</a>	Tělo je vyrobeno z hliníkového <b>korpusu</b> , a části mechanismu jsou skutečně kovové. <i>&lt;/p&gt;</i>
<a href="#">apetonline.cz</a>	<b>Korpus</b> dejte znovu minimálně na hodinu chladit. <i>&lt;/p&gt;</i>
<a href="#">korpus.cz</a>	Celková velikost <b>korpusu</b> InterCorp se tak již blíží k 50 milionům slov. <i>&lt;/p&gt;</i>
<a href="#">mix.cz</a>	Tmavý <b>korpus</b> je silný a intenzivní a s jedním soustem roste chuť na další.

## Ukázka dat v korpusu – XML vertikál

```
<dokument zanr="blog"
  navez="Dovolená v Paříži" datum="2011-10-28"
  url="http://karel.bloguje.cz/dovolena-v-parizi">
<odstavec nadpis="1">
<veta>
Po
sedmi
letech
v
kouzelné
Paříži
!
</veta>
</odstavec>
...
</dokument>
```

## Ukázka anotací v korpusu – XML vertikál s anotacemi

```
<dokument zanr="blog" nazev="Dovolená v Paříži">
<veta nadpis="1">
Po          po          k7c6          0  8
sedmi       sedm       k4c6          1  7
letech      léto      k1gNnPc6     2  7
v           v          k7c6          3 10
kouzelné    kouzelný  k2eAgFnSc6d1 4  9
<entita druh="město">
Paříži     Paříž     k1gFnSc6     5  9
</entita>
!          !          kx           6 11
<NP>          7  8
<PP>          8 11
<NP>          9 10
<PP>         10 11
<S>          11  -
</veta>
</dokument>
```

# Použití korpusů

- ▶ Obecně: data ke studiu přirozeného jazyka
- ▶ Lexikografové: slovníky
- ▶ Lingvisté: jazykové analýzy, změny jazyka
- ▶ Sociologové: jak a o čem píšeme, která témata jsou aktuální
- ▶ Marketingoví experti: hodnocení značek a výrobků v textech
- ▶ Statistické nástroje ZPJ: jazykové modely pro značkovače, analyzátory, překladové systémy, prediktivní psaní, . . .

# Použití jazykového korpusu v lexikografii

**korpus** czTenTen [2012] freq = 17,176 (3.38 per million)

<u>a modifier</u>	<u>gen 1</u>	<u>prec verb</u>	<u>post verb</u>
<u>7,932</u> 0.46	<u>1,923</u> 0.11	<u>784</u> 0.05	<u>1,413</u> 0.08
dortový + <u>213</u> 9.31 dortový korpus	skříňka + <u>203</u> 6.44 Korpusy skříněk	upéct <u>41</u> 5.40	prokrojit <u>14</u> 8.20 korpus prokrojíme
piškotový + <u>145</u> 8.90 piškotový korpus	skříň + <u>150</u> 4.71 korpus skříně	potřít <u>19</u> 4.30	rozkrojit <u>32</u> 7.85 korpus rozkrojíme
vychladlý + <u>139</u> 8.56 . Vychladlý korpus	buben <u>34</u> 4.67 korpus bubnu	péct <u>24</u> 3.29 peču korpus	rozříznout <u>46</u> 7.64 korpus rozřízneme
upečený + <u>193</u> 8.46 . Upečený korpus	komoda <u>9</u> 4.67	disponovat <u>16</u> 0.66 disponuje korpusem typu	proříznout <u>9</u> 5.34
závislostní <u>65</u> 7.94 Pražský závislostní korpus	Kristus <u>87</u> 4.63 s korpusem Krista	tvořit <u>48</u> 0.11 tvoří korpus	vyklopit <u>14</u> 5.29 korpus vyklopíme
anotovaný <u>40</u> 7.27 anotovaný korpus	skříňka <u>6</u> 4.48	<u>657</u> 0.04	potřít <u>38</u> 5.27 korpus potřeme
předpečený <u>32</u> 7.01 na předpečený korpus	dort <u>41</u> 4.05 korpus dortu	rozetřít <u>21</u> 6.43 rozetřeme na korpus	pomazat <u>8</u> 5.16
diachronní <u>33</u> 6.96 Diachronní korpus	věvec <u>20</u> 3.91 korpus věnce ,	navržit <u>6</u> 5.88	rozetřít <u>9</u> 5.07
synchronní <u>55</u> 6.90	svítidlo <u>16</u> 3.35	nalít <u>44</u> 4.80	pokapat <u>7</u> 4.90
	trezor <u>8</u> 3.25	natřít <u>25</u> 4.77 natřeme na korpus	postlat <u>14</u> 4.80
	varhany <u>8</u> 3.06		upéct <u>14</u> 3.83
			vychladnout <u>7</u> 3.71

## Srovnání korpusů podle původu textů – tradiční korpusy

- ▶ Příklady: British National Corpus, Corpus of Contemporary American English, Google Books, Český národní korpus
- ▶ Vznik: na objednávku, obsahová komise
- ▶ Zdroje: nejčastěji tištěná média, potom skenování knih, přepisy rozhovorů

## Srovnání korpusů podle původu textů – tradiční korpusy

- ▶ Příklady: British National Corpus, Corpus of Contemporary American English, Google Books, Český národní korpus
- ▶ Vznik: na objednávku, obsahová komise
- ▶ Zdroje: nejčastěji tištěná média, potom skenování knih, přepisy rozhovorů
- ▶ Výhody: úplná kontrola nad zdroji (kvalitní a bohaté informace o datech: autor, název, rok vydání, žánr, styl, oblast), známe rozložení typů textů v korpusu, zaručená kvalita textů (možnost opravy chyb)

## Srovnání korpusů podle původu textů – tradiční korpusy

- ▶ Příklady: British National Corpus, Corpus of Contemporary American English, Google Books, Český národní korpus
- ▶ Vznik: na objednávku, obsahová komise
- ▶ Zdroje: nejčastěji tištěná média, potom skenování knih, přepisy rozhovorů
- ▶ Výhody: úplná kontrola nad zdroji (kvalitní a bohaté informace o datech: autor, název, rok vydání, žánr, styl, oblast), známe rozložení typů textů v korpusu, zaručená kvalita textů (možnost opravy chyb)
- ▶ Nevýhody: nákladnost (jednání s vlastníky dat), omezená velikost (nedostatečná pro některé účely)



## Srovnání korpusů podle původu textů – tradiční korpusy

- ▶ Příklady: British National Corpus, Corpus of Contemporary American English, Google Books, Český národní korpus
- ▶ Vznik: na objednávku, obsahová komise
- ▶ Zdroje: nejčastěji tištěná média, potom skenování knih, přepisy rozhovorů
- ▶ Výhody: úplná kontrola nad zdroji (kvalitní a bohaté informace o datech: autor, název, rok vydání, žánr, styl, oblast), známe rozložení typů textů v korpusu, zaručená kvalita textů (možnost opravy chyb)
- ▶ Nevýhody: nákladnost (jednání s vlastníky dat), omezená velikost (nedostatečná pro některé účely)
- ▶ Shrnutí: reprezentativní vyvážený korpus daného jazyka

## Srovnání korpusů podle původu textů – internetové korpusy

- ▶ Příklady: Web as Corpus, ClueWeb, TenTen corpora, Corpora from the Web
- ▶ Vznik: opakované stahování internetu
- ▶ Zdroje: texty na internetu

## Srovnání korpusů podle původu textů – internetové korpusy

- ▶ Příklady: Web as Corpus, ClueWeb, TenTen corpora, Corpora from the Web
- ▶ Vznik: opakované stahování internetu
- ▶ Zdroje: texty na internetu
- ▶ Výhody: aktuální psaná podoba jazyka, velikost – pokrývá lépe/více jazykových jevů (ClueWeb 09: 70 mld. anglických slov), větší rozmanitost textů

# Srovnání korpusů podle původu textů – internetové korpusy

- ▶ Příklady: Web as Corpus, ClueWeb, TenTen corpora, Corpora from the Web
- ▶ Vznik: opakované stahování internetu
- ▶ Zdroje: texty na internetu
- ▶ Výhody: aktuální psaná podoba jazyka, velikost – pokrývá lépe/více jazykových jevů (ClueWeb 09: 70 mld. anglických slov), větší rozmanitost textů
- ▶ Nevýhody: malá kontrola nad zdroji (neuspořádanost, nevíme, co stahujeme), nezaručená kvalita textů (ale množství správných tvarů převáží chyby), nežádoucí obsah, duplicity, spam

## Srovnání korpusů podle původu textů – internetové korpusy

- ▶ Příklady: Web as Corpus, ClueWeb, TenTen corpora, Corpora from the Web
- ▶ Vznik: opakované stahování internetu
- ▶ Zdroje: texty na internetu
- ▶ Výhody: aktuální psaná podoba jazyka, velikost – pokrývá lépe/více jazykových jevů (ClueWeb 09: 70 mld. anglických slov), větší rozmanitost textů
- ▶ Nevýhody: malá kontrola nad zdroji (neuspořádanost, nevíme, co stahujeme), nezaručená kvalita textů (ale množství správných tvarů převáží chyby), nežádoucí obsah, duplicity, spam
- ▶ Shrnutí: velký korpus daného jazyka

## Proč je velikost korpusů důležitá

- ▶ Většina jazykových jevů podléhá Zipfově rozložení

## Proč je velikost korpusů důležitá

- ▶ Většina jazykových jevů podléhá Zipfově rozložení
- ▶ „There is no data like more data“ (Mercer, 1985)
- ▶ „More data usually beats better algorithms“

## Proč je velikost korpusů důležitá

- ▶ Většina jazykových jevů podléhá Zipfově rozložení
- ▶ „There is no data like more data“ (Mercer, 1985)
- ▶ „More data usually beats better algorithms“
- ▶ Větší seznamy slov (více unikátních slov)  
⇒ lepší pokrytí slov jazyka
- ▶ Více vět s výskytem daného slova  
⇒ lepší příklady použití slov v kontextu
- ▶ Lepší pokrytí řídkých jazykových jevů  
⇒ více podkladů pro studium jazyka
- ▶ Více dat pro jazykové modely  
⇒ přesnější (?) jazykové modely s větším pokrytím



## Ukázka: Slova rozvíjející frázi „deliver speech“

- ▶ BNC (96 M words): major (8), keynote (6).
- ▶ ukWaC (1,32 G words): keynote (125), opening (12), budget (8), wedding (7).
- ▶ enTenTen12 (11,2 G words): keynote (813), acceptance (129), major (127), wedding (118), short (101), opening (97), famous (80).
- ▶ enClueWeb09 (70,5 G words): keynote (3802), acceptance (1035), opening (589), famous (555), commencement (356), impassioned (335), inaugural (333).

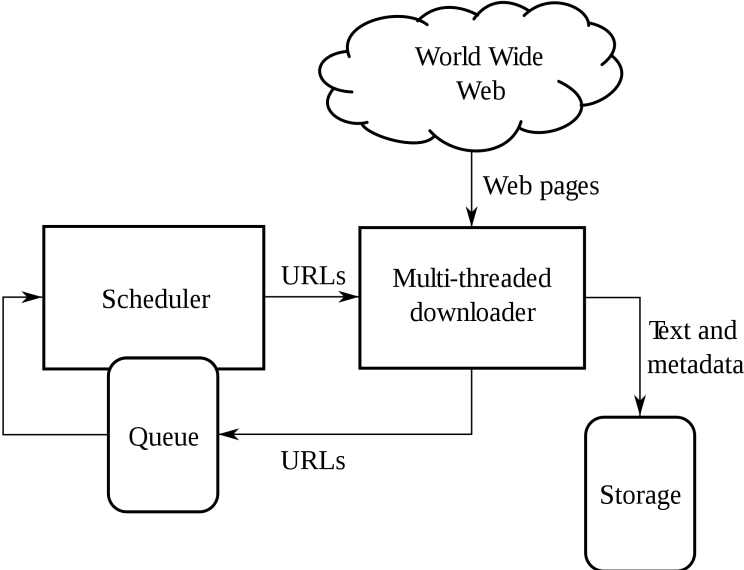
# Web crawler

- ▶ Traverses the internet (graph of pages and links).
- ▶ Downloads documents (content & meta information).
- ▶ Stores documents (or their parts) in various formats for further use.
- ▶ Crawlers for various purposes:
  - ▶ GoogleBot – web indexing,
  - ▶ Linkcrawler – links, broken link checking,
  - ▶ Heritrix – general crawler, (Java, multiple threads),
  - ▶ SpiderLing – text corpora, (Python, multiple sockets).

## Features a crawler should provide

- ▶ **Distributed:** Executable in a distributed fashion across multiple machines.
- ▶ **Scalable:** Scaling up the crawl rate by adding extra machines and bandwidth.
- ▶ **Performance and efficiency:** Efficient use of system resources (processor, memory, storage and network bandwidth).
- ▶ **Quality:** Biased towards fetching “useful” pages first.
- ▶ **Freshness:** Operate in continuous mode: obtain fresh copies of previously fetched pages, i.e. with a frequency that approximates the rate of change of that page. *Search engine crawler → the index contains a fairly current representation of each indexed web page.*
- ▶ **Extensible:** Cope with new data formats, new fetch protocols, various data processing needs. Modular architecture.

# Basic crawler design

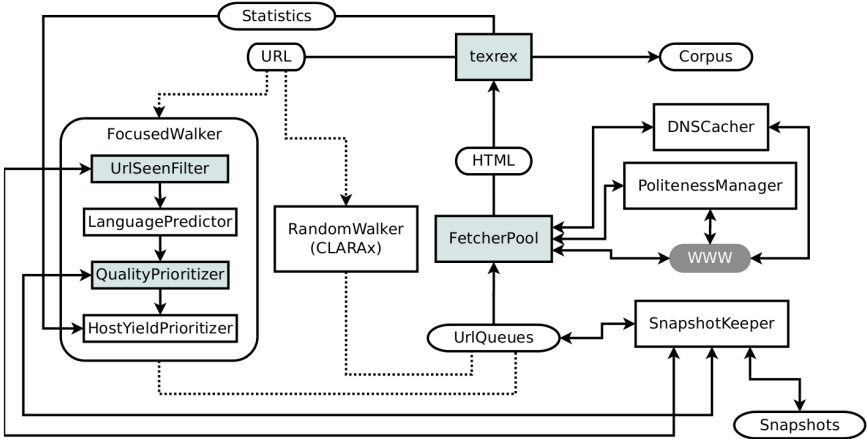


Source: [http://en.wikipedia.org/wiki/Web\\_crawler](http://en.wikipedia.org/wiki/Web_crawler)

# Advanced crawler implementation details

- ▶ Distributed vs. extensible.
- ▶ Multi-threaded synchronous vs. multi-socketed asynchronous.
- ▶ Web traversal policy:
  - ▶ depth vs. breadth,
  - ▶ domain selection,
  - ▶ domain distance,
  - ▶ focused crawling (topic oriented) vs. general crawling,
  - ▶ yield ratio.

# Focused crawler design

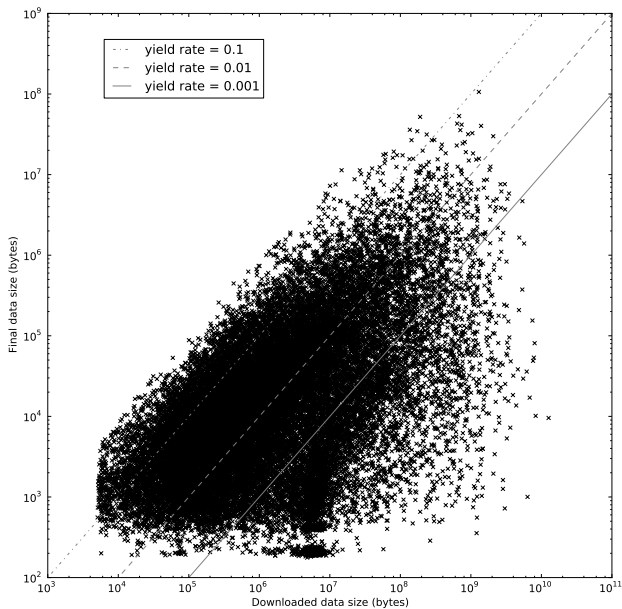


Source: Roland Schafer, Adrien Barbaresi, Felix Bildhauer. Focused Web Corpus Crawling. 9th Web as Corpus Workshop, 2014.

## SpiderLing – crawler pro textové korpusy

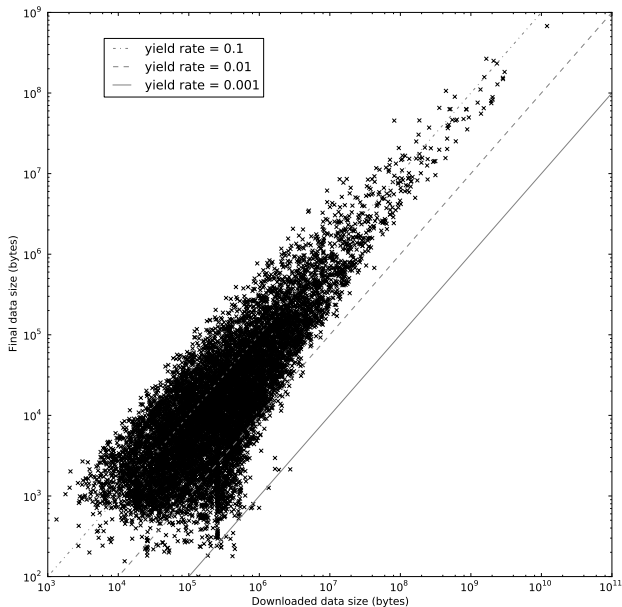
- ▶ důraz kladen na efektivitu stahování
- ▶ *míra výtěžnosti* =  $\frac{\text{velikost výsledných dat}}{\text{velikost stažených dat}}$
- ▶ crawler průběžně vyhodnocuje výtěžnost webových domén, zaměřuje se na „textově bohaté“ a odkládá stahování (nebo vůbec nestahuje) z neperspektivních webů
- ▶ cílem je sestavit korpusy velikosti  $\geq 10^{10}$  slov pro všechny významné jazyky

# General unfocused crawling efficiency (Heritrix)





# Domain yield ratio optimised efficiency (SpiderLing)



# Konkordance – autentické použití slova v textu

questions our archives. Explore the links we provide in each category. Use a search engine like Google. Questions should be specific . Questions like Hello the best wordpress :p I have a request but it's not about WG => I shook his head, "No." "Jump up and down." "What?" "If I gotten from searching sources to answer my curiosity. Using Google Scholar helps us to understand and trade the market swing. For example, we may eventually lose interest in this type of price action and jump ship in numbers, astrological dates and prayer wheels have all been enlisted in the geographic area covered by each publication. Google News Archive publication. Google News Archive Search -- Google's "News archive Archive Search -- Google's "News archive search provides an easy way to provides an easy way to search and explore historical archives. Users overview of the results by browsing an automatically created timeline. Search Google has developed to put related stories together in the same news related stories together in the same news search 's database of over 1.9 weblogs and get up-to-date information on your monitor breaking news from over 3,000 sources, 24 hours a day. You can also from numerous news sources in your area. NewspaperArchive.com -- Search over 12.3 million newspaper pages. Not a complete resource to morale and possibly paralyze you from taking necessary actions in your job save you a lot of hassle later on and you should do it early in your job on their own in the shortest time possible. Discover more insider job quickest way to locate the information you want is to use the Label and the item of information that you require. You can use the very powerful section that you are currently viewing. Searching You can very easily section that you are currently in. To use, you simply enter

engine like Google. Questions should be specific . Questions like Davichi - love & war mv gif but i don't find =/ I ask you because maybe you, all I find I keep? "You ain't searching me, man." Marco engine, with key words: winter health tips, I compiled and present to the chart for a reversal or breakout pattern that spells opportunity, of a more exciting trading vehicle. The market loses broad sponsorship for that elusive trading edge. Most traders believe Fibonacci fits -- Google's "News archive search provides an easy way to search and explore provides an easy way to search and explore historical archives. Users and explore historical archives. Users can search for events, people, for events, people, ideas and see how they have been described over time results include both content that is accessible to all users and content result. Excellent place to search for news articles by keywords. Technorati --Search Technorati's terms. British Library Newspaper Library -- The only large, integrated for news by zipcode and receive a grouping of news from numerous news over 12.3 million newspaper pages. Not a complete resource to search but a good place to start research. Today's Front Pages -- "Today . If that isn't bad enough, it can also stop you from being hired! Here . In fact, this is something we all should do at least once a year: Find secrets by visiting <http://www.jobchangesecrets.com> Art Canvas Features. The Label pulldown is by far the easiest and quickest, just engine. Remember that each section on shows 'X' number of posts the section that you are currently in. To use, you simply enter item and hit Enter or click the Search Button . Your results are then

# Konkordance – co je špatně?

Homepage | Create Your Own Homepage | Change My Search To: Google Search MSN Search Yahoo! Search Ask this? | Create Your Own Homepage | Change My Search To: Google Search MSN Search Yahoo! Search Ask.com Search Wikipedia Own Homepage | Change My Search To: Google Search MSN Search Yahoo! Search Ask.com Search Wikipedia English | Change My Search To: Google Search MSN Search Yahoo! Search Ask.com Search Wikipedia English Yahoo! Answers | Google Search MSN Search Yahoo! Search Ask.com Search Wikipedia English Yahoo! Answers Answers.com | a z , zaz DIRECT Everything You're Looking For Browse, Search , Find... Get Found (List Yourself) \* Bookmark Your Favorites france24.com | Google News Google News Aggregated headlines and a engine of many of the world's news sources. news.google.com | Gold Suppliers All Products All Suppliers Advanced Search Browse Categories Popular Searches : Urea Fertilizer Dyes Intermediates Infomation Online Inquiry JIANGSU XINYE CHEMICAL CO., LTD. Search products of this supplier Founded in 1997, JIANGSU XINYE CHEMICAL | Gold Suppliers All Products All Suppliers Advanced Search Browse Categories Popular Searches : Urea Fertilizer Dyes Intermediates Infomation Online Inquiry JIANGSU XINYE CHEMICAL CO., LTD. Search products of this supplier O-Fluorobenzoyl Chloride Enquiry | Gold Suppliers All Products All Suppliers Advanced Search Browse Categories Popular Searches : Urea Fertilizer Dyes Intermediates Infomation Online Inquiry JIANGSU XINYE CHEMICAL CO., LTD. Search products of this supplier Flutriafol Enquiry | Gold Suppliers All Products All Suppliers Advanced Search Browse Categories Popular Searches : Urea Fertilizer Dyes Intermediates Infomation Online Inquiry JIANGSU XINYE CHEMICAL CO., LTD. Search products of this supplier O-Fluorobenzoyl Chloride Enquiry | Gold Suppliers All Products All Suppliers Advanced Search Browse Categories Popular Searches : Urea Fertilizer Dyes Intermediates Infomation Online Inquiry JIANGSU XINYE CHEMICAL CO., LTD. Search products of this supplier Flutriafol Enquiry | Gold Suppliers All Products All Suppliers Advanced Search Browse Categories Popular Searches : Urea Fertilizer Dyes Intermediates | Contact Infomation Online Inquiry P-Chlorophenol Search products of this supplier Product Type: Agrochemicals | Gold Suppliers All Products All Suppliers Advanced Search Browse Categories Popular Searches : Urea Fertilizer Dyes Intermediates | Contact Infomation Online Inquiry O-Chlorophenol Search products of this supplier Product Type: Agrochemicals | Gold Suppliers All Products All Suppliers Advanced Search Browse Categories Popular Searches : Urea Fertilizer Dyes Intermediates | Contact Infomation Online Inquiry 2,4-Dichlorophenol Search products of this supplier Product Type: Agrochemicals

# Odstraňování nežádoucího obsahu

## Nežádoucí obsah

- ▶ html značky, styly, poznámky
- ▶ negramatické věty: navigace, reklamy, tabulky, příliš krátké úseky,...

## Používáme nástroj jusText

(<http://nlp.fi.muni.cz/projects/justext>)

- ▶ rozdělení na odstavce
- ▶ slovník častých slov v daném jazyce
- ▶ klasifikace odstavce podle délky, hustoty slov ze slovníku, hustoty odkazů, třídy okolních odstavců
- ▶ demo na stránce nástroje

# Ukázka boilerplate

The image shows a screenshot of a university website with several annotations. A red 'X' labeled 'boilerplate' is placed over the top navigation bar and header image. A yellow checkmark labeled 'heading' is placed over the main title 'Studijní plány - základní informace'. A green checkmark labeled 'content' is placed over the main text area. Another green checkmark labeled 'content' is placed over the left sidebar. A red 'X' labeled 'boilerplate' is placed over the footer area.

**boilerplate**

**heading**

**content**

**content**

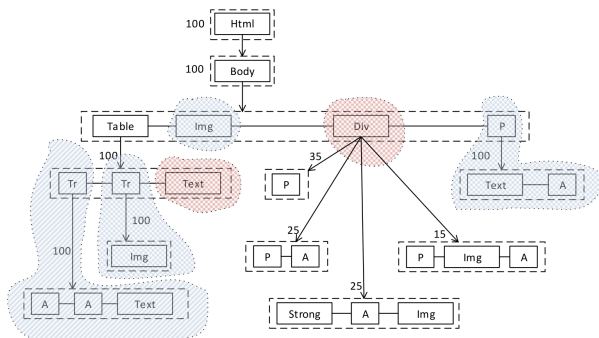
**boilerplate**

# Boilerplate removal approaches

- ▶ Machine learning (SVM, CRF, neural networks, n-gram models):
  - ▶ Annotated web pages required for training.
  - ▶ Victor (CRF),
  - ▶ Ncleaner (n-grams).
- ▶ Heuristics:
  - ▶ Rules for including/excluding sections of text.
  - ▶ BTE (tag density),
  - ▶ Boilerpipe (link/text ratio),
  - ▶ jusText (link/text ratio, frequent words, context sensitive – smoothing).

# Site Style Tree (Yi, Lan, Bing Liu, Xiaoli Li, 2003)

- ▶ Represents both layout and content of a web page.
- ▶ Node importance = node entropy over the whole Site Style Tree.



Context sensitive paragraph classification:



Demo: <http://nlp.fi.muni.cz/projects/justext/>



# Konkordance – co je ještě špatně?

Homepage </p><p> What is this? | Create Your Own Homepage </p><p> Change My Search To: </p><p> Google Search </p><p> MSN Search </p><p> Yahoo! Search </p><p> Ask this? | Create Your Own Homepage </p><p> Change My Search To: </p><p> Google Search </p><p> MSN Search </p><p> Yahoo! Search </p><p> Ask.com Search </p><p> Wikipedia Own Homepage </p><p> Change My Search To: </p><p> Google Search </p><p> MSN Search </p><p> Yahoo! Search </p><p> Ask.com Search </p><p> Wikipedia English </p><p> Change My Search To: </p><p> Google Search </p><p> MSN Search </p><p> Yahoo! Search </p><p> Ask.com Search </p><p> Wikipedia English </p><p> Yahoo! Answers </p><p> Google Search </p><p> MSN Search </p><p> Yahoo! Search </p><p> Ask.com Search </p><p> Wikipedia English </p><p> Yahoo! Answers </p><p> Answers.com </p><p> a z , </p><p> zaz DIRECT Everything You're Looking For </p><p> Browse, Search , Find... </p><p> Get Found (List Yourself) \* </p><p> Bookmark Your Favorites france24.com </p><p> Google News Google News </p><p> Aggregated headlines and a Search engine of many of the world's news sources. </p><p> news.google.com </p><p> Gold Suppliers </p><p> All Products </p><p> All Suppliers </p><p> Advanced Search Browse Categories </p><p> Popular Searches : Urea Fertilizer Dyes Intermediates Infomation </p><p> Online Inquiry </p><p> Jiangsu Xinye Chemical Co., Ltd. </p><p> Search products of this supplier </p><p> Founded in 1997, Jiangsu Xinye Chemical </p><p> Gold Suppliers </p><p> All Products </p><p> All Suppliers </p><p> Advanced Search Browse Categories </p><p> Popular Searches : Urea Fertilizer Dyes Intermediates Infomation </p><p> Online Inquiry </p><p> Jiangsu Xinye Chemical Co., Ltd. </p><p> Search products of this supplier </p><p> O-Fluorobenzoyl Chloride </p><p> Enquiry </p><p> Gold Suppliers </p><p> All Products </p><p> All Suppliers </p><p> Advanced Search Browse Categories </p><p> Popular Searches : Urea Fertilizer Dyes Intermediates Infomation </p><p> Online Inquiry </p><p> Jiangsu Xinye Chemical Co., Ltd. </p><p> Search products of this supplier </p><p> Flutriafol </p><p> Enquiry </p><p> ... </p><p> Gold Suppliers </p><p> All Products </p><p> All Suppliers </p><p> Advanced Search Browse Categories </p><p> Popular Searches : Urea Fertilizer Dyes Intermediates Infomation </p><p> Online Inquiry </p><p> Jiangsu Xinye Chemical Co., Ltd. </p><p> Search products of this supplier </p><p> O-Fluorobenzoyl Chloride </p><p> Enquiry </p><p> Gold Suppliers </p><p> All Products </p><p> All Suppliers </p><p> Advanced Search Browse Categories </p><p> Popular Searches : Urea Fertilizer Dyes Intermediates Infomation </p><p> Online Inquiry </p><p> Jiangsu Xinye Chemical Co., Ltd. </p><p> Search products of this supplier </p><p> Flutriafol </p><p> Enquiry </p><p> ... </p><p> Gold Suppliers </p><p> All Products </p><p> All Suppliers </p><p> Advanced Search Browse Categories </p><p> Popular Searches : Urea Fertilizer Dyes Intermediates </p><p> Contact Infomation </p><p> Online Inquiry </p><p> P-Chlorophenol </p><p> Search products of this supplier </p><p> Product Type: </p><p> Agrochemicals </p><p> Gold Suppliers </p><p> All Products </p><p> All Suppliers </p><p> Advanced Search Browse Categories </p><p> Popular Searches : Urea Fertilizer Dyes Intermediates </p><p> Contact Infomation </p><p> Online Inquiry </p><p> O-Chlorophenol </p><p> Search products of this supplier </p><p> Product Type: </p><p> Agrochemicals </p><p> Gold Suppliers </p><p> All Products </p><p> All Suppliers </p><p> Advanced Search Browse Categories </p><p> Popular Searches : Urea Fertilizer Dyes Intermediates </p><p> Contact Infomation </p><p> Online Inquiry </p><p> 2,4-Dichlorophenol </p><p> Search products of this supplier </p><p> Product Type: </p><p> Agrochemicals </p>

# Duplicitní texty na internetu

- ▶ Jak duplicitní texty vznikají?

# Duplicitní texty na internetu

- ▶ Jak duplicitní texty vznikají?
  - ▶ Hlavičky, patičky, navigace, copyright

# Duplicitní texty na internetu

- ▶ Jak duplicitní texty vznikají?
  - ▶ Hlavičky, patičky, navigace, copyright
  - ▶ Kopírování – informace na více místech

# Duplicitní texty na internetu

- ▶ Jak duplicitní texty vznikají?
  - ▶ Hlavičky, patičky, navigace, copyright
  - ▶ Kopírování – informace na více místech
  - ▶ Levné získání obsahu vykradením textu z původních zdrojů

# Duplicitní texty na internetu

- ▶ Jak duplicitní texty vznikají?
  - ▶ Hlavičky, patičky, navigace, copyright
  - ▶ Kopírování – informace na více místech
  - ▶ Levné získání obsahu vykradením textu z původních zdrojů
  - ▶ Média přejímají agenturní zprávy

# Duplicitní texty na internetu

- ▶ Jak duplicitní texty vznikají?
  - ▶ Hlavičky, patičky, navigace, copyright
  - ▶ Kopírování – informace na více místech
  - ▶ Levné získání obsahu vykradením textu z původních zdrojů
  - ▶ Média přejímají agenturní zprávy
  - ▶ Minimální změny: sportovní výsledky, předpověď počasí

# Duplicitní texty na internetu

- ▶ Jak duplicitní texty vznikají?
  - ▶ Hlavičky, patičky, navigace, copyright
  - ▶ Kopírování – informace na více místech
  - ▶ Levné získání obsahu vykradením textu z původních zdrojů
  - ▶ Média přejímají agenturní zprávy
  - ▶ Minimální změny: sportovní výsledky, předpověď počasí
  - ▶ Licenční ujednání



# Duplicitní texty na internetu

- ▶ Jak duplicitní texty vznikají?
  - ▶ Hlavičky, patičky, navigace, copyright
  - ▶ Kopírování – informace na více místech
  - ▶ Levné získání obsahu vykradením textu z původních zdrojů
  - ▶ Média přejímají agenturní zprávy
  - ▶ Minimální změny: sportovní výsledky, předpověď počasí
  - ▶ Licenční ujednání
- ▶ Jak podobné texty jsou příliš podobné?

# Duplicitní texty na internetu

- ▶ Jak duplicitní texty vznikají?
  - ▶ Hlavičky, patičky, navigace, copyright
  - ▶ Kopírování – informace na více místech
  - ▶ Levné získání obsahu vykradením textu z původních zdrojů
  - ▶ Média přejímají agenturní zprávy
  - ▶ Minimální změny: sportovní výsledky, předpověď počasí
  - ▶ Licenční ujednání
- ▶ Jak podobné texty jsou příliš podobné?
  - ▶ Členění textu: dokumenty, odstavce, věty, slova

# Duplicitní texty na internetu

- ▶ Jak duplicitní texty vznikají?
  - ▶ Hlavičky, patičky, navigace, copyright
  - ▶ Kopírování – informace na více místech
  - ▶ Levné získání obsahu vykradením textu z původních zdrojů
  - ▶ Média přejímají agenturní zprávy
  - ▶ Minimální změny: sportovní výsledky, předpověď počasí
  - ▶ Licenční ujednání
- ▶ Jak podobné texty jsou příliš podobné?
  - ▶ Členění textu: dokumenty, odstavce, věty, slova
  - ▶ Jak můžeme měřit podobnost textů?

# Duplicitní texty na internetu

- ▶ Jak duplicitní texty vznikají?
  - ▶ Hlavičky, patičky, navigace, copyright
  - ▶ Kopírování – informace na více místech
  - ▶ Levné získání obsahu vykradením textu z původních zdrojů
  - ▶ Média přejímají agenturní zprávy
  - ▶ Minimální změny: sportovní výsledky, předpověď počasí
  - ▶ Licenční ujednání
- ▶ Jak podobné texty jsou příliš podobné?
  - ▶ Členění textu: dokumenty, odstavce, věty, slova
  - ▶ Jak můžeme měřit podobnost textů?
  - ▶ Jak nastavit práh podobnosti?

# Ukázka agenturní zprávy – Bahrain News Agency



Fri 10 Feb 2017

[Home](#)

[About BNA](#)

[Photo Archive](#)

[TV](#)

[Radio](#)

[Contact Us](#)

عربي

Search



0

Main

Home

Latest News

Local News

Arab News

World News

Business

Sports

INA Reports

INA video

## Explosion Strikes German Convoy in Afghanistan

12 : 47 PM - 19/06/2011

**Kabul, June. 19 (BNA) --** Afghan officials say an explosion has struck a German military convoy on a main road in northern Afghanistan, killing two Afghan civilians who were nearby.

A reporter at the scene saw at least one overturned armored vehicle and what appeared to be a civilian car that was mangled in the blast Sunday.

Provincial spokesman Muhbobullah Sayedi says the two civilians who were killed were pedestrians who were caught up in the blast. Sayedi said it appeared that bomb came from a suicide attacker.

### *Explosion strikes German convoy in Afghanistan*

AP — PUBLISHED Jun 19, 2011 08:35am

**KUNDUZ, Afghanistan: A suicide car bomber struck a German military convoy in northern Afghanistan on Sunday, detonating explosives that killed three Afghan civilians and overturned at least one armored vehicle, according to officials and witnesses.**

Taliban spokesman Zabiullah Mujahid claimed responsibility for the attack.

The bomber blew his vehicle up shortly before 10 a.m. local time on a busy road on the edge of Kunduz city, near the airport.

## Explosion strikes German convoy in Afghanistan

By **KATHY GANNON**, Associated Press

JUNE 18, 2011, 11:37 PM | KUNDUZ, AFGHANISTAN

**A** suicide attacker blew up his explosives-laden car next to a German military convoy in northern Afghanistan on Sunday, killing three Afghan civilians, officials and witnesses said.

[Taliban](#) spokesman Zabiullah Mujahid claimed responsibility for the attack.

The bomber detonated his vehicle shortly before 10 a.m. local time on a busy road on the edge of Kunduz city, near the airport, the Afghan Interior Ministry said in a statement. An Associated Press reporter at the scene saw at least one overturned armored vehicle and what appeared to be a civilian car that was mangled in the blast.

# Taliban Targets German Troops

June 19th, 2011 at 12:18 pm [FRUMFORUM NEWS](#) | [No Comments](#) | [Share](#) | [Print](#)



The Associated Press **reports:**

*A suicide attacker blew up his explosives-laden car next to a German military convoy in northern Afghanistan on Sunday, killing three Afghan civilians, officials and witnesses said.*

*Taliban spokesman Zabiullah Mujahid claimed responsibility for the attack.*

*The bomber detonated his vehicle shortly before 10 a.m. local time on a busy road on the edge of Kunduz city, near the airport, the Afghan Interior Ministry said in a statement. An Associated Press reporter at the scene saw at least one overturned armored vehicle and what appeared to be a civilian car that was mangled in the blast.*

*Three civilians were killed and 11 were wounded in the attack, the ministry said.*

*Germany's military said two German soldiers were lightly wounded and treated at a nearby base. Two vehicles were damaged, according to a German military spokesman, who declined to be named in line with department policy. The spokesman said the military could not immediately confirm whether the attack was a suicide assault or a roadside bomb.*



# Deduplikace

Odstranění opakujících se úseků textu

- ▶ Zcela identické texty – snadné?

# Deduplikace

Odstranění opakujících se úseků textu

- ▶ Zcela identické texty – snadné?
  - ▶ Text převedeme na číslo pomocí *hašovací funkce*

# Deduplikace

Odstranění opakujících se úseků textu

- ▶ Zcela identické texty – snadné?
  - ▶ Text převedeme na číslo pomocí *hašovací funkce*
  - ▶ Srovnáváme haš nového textu se všemi předchozími

# Deduplikace

Odstranění opakujících se úseků textu

- ▶ Zcela identické texty – snadné?
  - ▶ Text převedeme na číslo pomocí *hašovací funkce*
  - ▶ Srovnáváme haš nového textu se všemi předchozími
- ▶ A co podobné texty – obtížné?

# Deduplikace

Odstranění opakujících se úseků textu

- ▶ Zcela identické texty – snadné?
  - ▶ Text převedeme na číslo pomocí *hašovací funkce*
  - ▶ Srovnáváme haš nového textu se všemi předchozími
- ▶ A co podobné texty – obtížné?
  - ▶ Lidé často zkopírují jen části textu

# Deduplikace

Odstranění opakujících se úseků textu

- ▶ Zcela identické texty – snadné?
  - ▶ Text převedeme na číslo pomocí *hašovací funkce*
  - ▶ Srovnáváme haš nového textu se všemi předchozími
- ▶ A co podobné texty – obtížné?
  - ▶ Lidé často zkopírují jen části textu
  - ▶ Nebo dokument zkopírují a změní

# Deduplikace

Odstranění opakujících se úseků textu

- ▶ Zcela identické texty – snadné?
  - ▶ Text převedeme na číslo pomocí *hašovací funkce*
  - ▶ Srovnáváme haš nového textu se všemi předchozími
- ▶ A co podobné texty – obtížné?
  - ▶ Lidé často zkopírují jen části textu
  - ▶ Nebo dokument zkopírují a změní
  - ▶ Nebo zkopírují a rozšíří vlastním textem

# Deduplikace

## Odstranění opakujících se úseků textu

- ▶ Zcela identické texty – snadné?
  - ▶ Text převedeme na číslo pomocí *hašovací funkce*
  - ▶ Srovnáváme haš nového textu se všemi předchozími
- ▶ A co podobné texty – obtížné?
  - ▶ Lidé často zkopírují jen části textu
  - ▶ Nebo dokument zkopírují a změní
  - ▶ Nebo zkopírují a rozšíří vlastním textem
  - ▶ Problém je třeba řešit na nižší úrovni, než celý dokument



# Deduplikace

## Odstranění opakujících se úseků textu

- ▶ Zcela identické texty – snadné?
  - ▶ Text převedeme na číslo pomocí *hašovací funkce*
  - ▶ Srovnáváme haš nového textu se všemi předchozími
- ▶ A co podobné texty – obtížné?
  - ▶ Lidé často zkopírují jen části textu
  - ▶ Nebo dokument zkopírují a změní
  - ▶ Nebo zkopírují a rozšíří vlastním textem
  - ▶ Problém je třeba řešit na nižší úrovni, než celý dokument
  - ▶ Nápady?

## Dokument jako vektor četností slov

- ▶ Hlavní myšlenka: počet výskytů slova je důležitý

## Dokument jako vektor četností slov

- ▶ Hlavní myšlenka: počet výskytů slova je důležitý
- ▶ Spočítáme *relativní četnost* slov

## Dokument jako vektor četností slov

- ▶ Hlavní myšlenka: počet výskytů slova je důležitý
- ▶ Spočítáme *relativní četnost* slov
- ▶ Dokument reprezentujeme *vektorem četností slov*

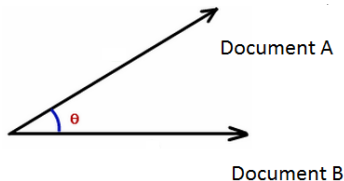
## Dokument jako vektor četností slov

- ▶ Hlavní myšlenka: počet výskytů slova je důležitý
- ▶ Spočítáme *relativní četnost* slov
- ▶ Dokument reprezentujeme *vektorem četností slov*
- ▶ Jak zjistíme podobnost vektorů v n-rozměrném prostoru?

# Dokument jako vektor četností slov

- ▶ Hlavní myšlenka: počet výskytů slova je důležitý
- ▶ Spočítáme *relativní četnost* slov
- ▶ Dokument reprezentujeme *vektorem četností slov*
- ▶ Jak zjistíme podobnost vektorů v n-rozměrném prostoru?
- ▶ Mírou podobnosti dvou vektorů může být *kosinová podobnost*

$$\cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$



# Dokument jako překrývající se n-tice slov

- ▶ Hlavní myšlenka: okolí slova je důležité

## Dokument jako překrývající se $n$ -tice slov

- ▶ Hlavní myšlenka: okolí slova je důležité
- ▶ Rozdělíme dokument na překrývající se  *$n$ -tice slov*  
Hospodin je můj, je můj pastýř, můj pastýř nebudu, pastýř  
nebudu mít, nebudu mít nedostatek



## Dokument jako překrývající se $n$ -tice slov

- ▶ Hlavní myšlenka: okolí slova je důležité
- ▶ Rozdělíme dokument na překrývající se  *$n$ -tice slov*  
Hospodin je můj, je můj pastýř, můj pastýř nebudu, pastýř nebudu mít, nebudu mít nedostatek
- ▶ Dokument reprezentujeme množinou hašů těchto  $n$ -tic

## Dokument jako překrývající se $n$ -tice slov

- ▶ Hlavní myšlenka: okolí slova je důležité
- ▶ Rozdělíme dokument na překrývající se  *$n$ -tice slov*  
Hospodin je můj, je můj pastýř, můj pastýř nebudu, pastýř nebudu mít, nebudu mít nedostatek
- ▶ Dokument reprezentujeme množinou hašů těchto  $n$ -tic
- ▶ Jak zjistíme podobnost množin?

## Dokument jako překrývající se n-tice slov

- ▶ Hlavní myšlenka: okolí slova je důležité
- ▶ Rozdělíme dokument na překrývající se *n-tice slov*  
Hospodin je můj, je můj pastýř, můj pastýř nebudu, pastýř nebudu mít, nebudu mít nedostatek
- ▶ Dokument reprezentujeme množinou hašů těchto n-tic
- ▶ Jak zjistíme podobnost množin?
- ▶ Mírou podobnosti dvou množin může být *Jaccardova podobnost*

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}.$$

		A	
		0	1
B	0	$M_{00}$	$M_{10}$
	1	$M_{01}$	$M_{11}$

## Dokument jako překrývající se n-tice slov

V CZPJ používáme následující parametry deduplikace pomocí překrývajících se n-tic slov:

- ▶ Úroveň jemnosti: odstavce,
- ▶ velikost úseků: sedmice slov
- ▶ hranice úseků: začátek a konec věty,
- ▶ práh podobnosti sedmic: 50 %
- ▶ vyhlazování.

## Dokument jako překrývající se n-tice slov

V CZPJ používáme následující parametry deduplikace pomocí překrývajících se n-tic slov:

- ▶ Úroveň jemnosti: odstavce,
- ▶ velikost úseků: sedmice slov
- ▶ hranice úseků: začátek a konec věty,
- ▶ práh podobnosti sedmic: 50 %
- ▶ vyhlazování.

Odstavec je ponechán,

- ▶ pokud alespoň 50 % sedmic slov v daném odstavci nebylo zaznamenáno dříve,
- ▶ nebo leží bezprostředně mezi dvěma ponechanými odstavci.

## Srovnání a slabá místa metod deduplikace

Zkonstruujte text, kterým uvedené metody obelstíte

- ▶ Metoda založená na relativní četnosti slov

# Srovnání a slabá místa metod deduplikace

Zkonstruujte text, kterým uvedené metody obelstíte

- ▶ Metoda založená na relativní četnosti slov
  - ▶ Stačí zkopírovat dva různé dokumenty za sebe

# Srovnání a slabá místa metod deduplikace

Zkonstruujte text, kterým uvedené metody obelstíte

- ▶ Metoda založená na relativní četnosti slov
  - ▶ Stačí zkopírovat dva různé dokumenty za sebe
- ▶ Metoda založená na překrývajících se n-ticích slov



# Srovnání a slabá místa metod deduplikace

Zkonstruujte text, kterým uvedené metody obelstíte

- ▶ Metoda založená na relativní četnosti slov
  - ▶ Stačí zkopírovat dva různé dokumenty za sebe
- ▶ Metoda založená na překrývajících se n-ticích slov
  - ▶ Stačí změnit pořadí slov ve větách

# Srovnání a slabá místa metod deduplikace

Zkonstruujte text, kterým uvedené metody obelstíte

- ▶ Metoda založená na relativní četnosti slov
  - ▶ Stačí zkopírovat dva různé dokumenty za sebe
- ▶ Metoda založená na překrývajících se n-ticích slov
  - ▶ Stačí změnit pořadí slov ve větách
- ▶ Obelstěte obě metody deduplikace zároveň

# Srovnání a slabá místa metod deduplikace

Zkonstruujte text, kterým uvedené metody obelstíte

- ▶ Metoda založená na relativní četnosti slov
  - ▶ Stačí zkopírovat dva různé dokumenty za sebe
- ▶ Metoda založená na překrývajících se n-ticích slov
  - ▶ Stačí změnit pořadí slov ve větách
- ▶ Obelstěte obě metody deduplikace zároveň
  - ▶ Stačí zaměnit některá slova za synonyma,
  - ▶ případně celý text převyprávět

# Srovnání a slabá místa metod deduplikace

Zkonstruujte text, kterým uvedené metody obelstíte

- ▶ Metoda založená na relativní četnosti slov
  - ▶ Stačí zkopírovat dva různé dokumenty za sebe
- ▶ Metoda založená na překrývajících se  $n$ -ticích slov
  - ▶ Stačí změnit pořadí slov ve větách
- ▶ Obelstěte obě metody deduplikace zároveň
  - ▶ Stačí zaměnit některá slova za synonyma,
  - ▶ případně celý text převyprávět
- ▶ Ukažte, že existuje text s unikátními větami, který s naším nastavením metody překrývajících se  $n$ -tic slov zahodíme

# Srovnání a slabá místa metod deduplikace

Zkonstruujte text, kterým uvedené metody obelstíte

- ▶ Metoda založená na relativní četnosti slov
  - ▶ Stačí zkopírovat dva různé dokumenty za sebe
- ▶ Metoda založená na překrývajících se  $n$ -ticích slov
  - ▶ Stačí změnit pořadí slov ve větách
- ▶ Obelstěte obě metody deduplikace zároveň
  - ▶ Stačí zaměnit některá slova za synonyma,
  - ▶ případně celý text převyprávět
- ▶ Ukažte, že existuje text s unikátními větami, který s naším nastavením metody překrývajících se  $n$ -tic slov zahodíme
  - ▶ Tvoří-li unikátní věty méně než 50 % sedmic slov, celý odstavec zahodíme

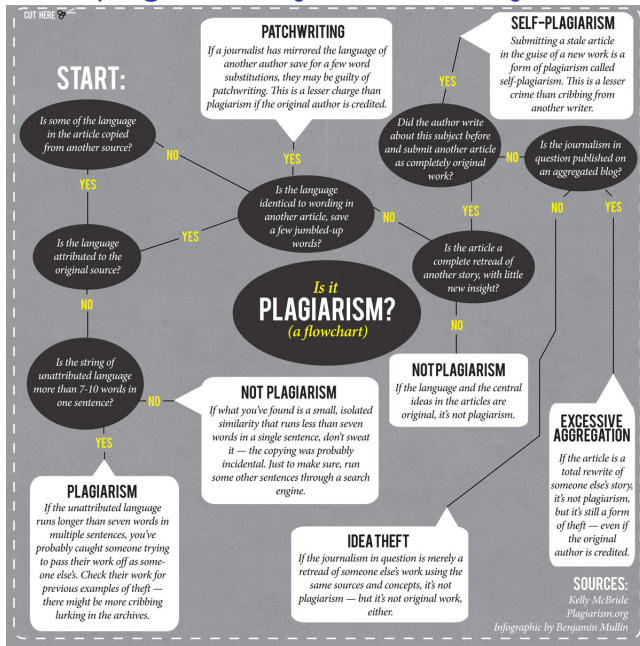
4. ŘÍJNA 2016 8:52 | [LIDOVKY.CZ](#) > [ZPRÁVY](#) > [DOMOV](#)

## Piráti: Chovanec opsal úvod bakalářské práce. Snaha o předvolební skandálek, brání se ministr



PRAHA Ministr vnitra Milan Chovanec (ČSSD) okopíroval značnou část své bakalářské práce, tvrdí Piráti a dokazují to pasážemi bez řádného citování. Práce pochází z roku 2009, kdy Chovanec studoval na Právnické fakultě Západočeské univerzity v Plzni. Ministr se brání, že se Piráti snaží jen o „předvolební skandálek“ a za svým studiem si stojí.

# Problém plagiátorství je však složitější



# Issues of Building Language Resources from the Web

Particular tasks:

- ▶ Language identification,
- ▶ Character encoding detection,
- ▶ Efficient web crawling,
- ▶ Boilerplate removal,
- ▶ De-duplication (removal of identical or nearly identical texts),
- ▶ Fighting web spam,
- ▶ Authorship recognition & plagiarism detection,
- ▶ Storing & indexing of large text collections.

NLPC & Lexical Computing corpus tools:

<http://corpus.tools/>



# Postup získávání webových korpusů v CZPJ

- ▶ příprava jazykově závislých modelů používaných v dalších krocích — učení na dokumentech z Wikipedie
- ▶ spuštění crawleru (SpiderLing)
- ▶ zpracování a vyhodnocování během běhu crawleru
  - ▶ detekce znakové sady dokumentu (Chared)
  - ▶ filtrování jazyka (vektor trigramů znaků)
  - ▶ odstraňování nežádoucího obsahu (Justext)
  - ▶ kontrola duplicitních dokumentů
  - ▶ vyhodnocování průběžné výtěžnosti webových domén
- ▶ zpracování získaných dat
  - ▶ odstranění podobných odstavců (Onion)
  - ▶ tokenizace (Unitok nebo jiný nástroj)
  - ▶ značkování morfologické a syntaktické — externími nástroji, jsou-li dostupné
  - ▶ zakódování a nahrání do korpusového manažeru (Manatee/Sketch Engine)