

# Lexical Computing

Vojtěch Kovář

Lexical  Computing

Brighton, UK & Brno, CZ

Setkání SPP FI MU  
listopad 2017

# Lexical Computing

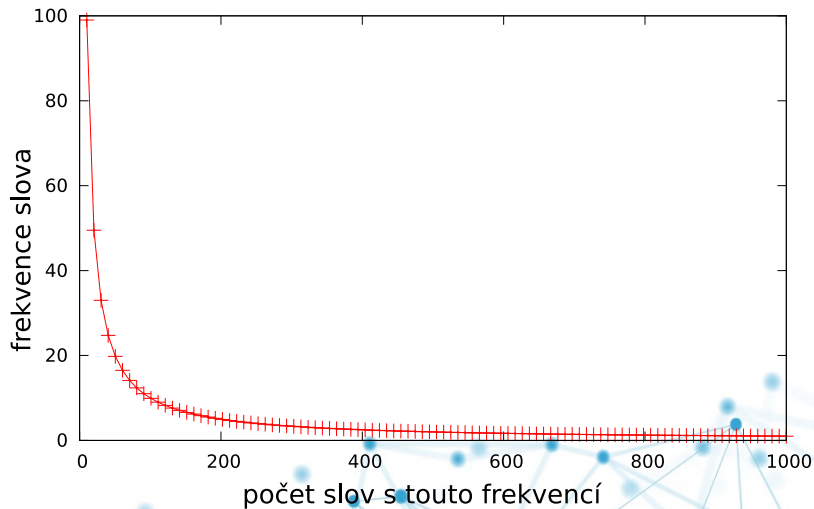
- výzkumná firma založená v Británii v roce 2003
- jazykové technologie, zejména korpusová lingvistika a počítačová lexikografie
- významný podíl činnosti tvoří vlastní výzkum
- partner na FI: Centrum zpracování přirozeného jazyka

# Co je to korpus

## Textový korpus

- rozsáhlý soubor textů
- často jazykově označovaný  
základní tvary, slovní druhy, morfologické značky
- všeobecné i vysoce specializované (oborově, jazykově)
- obsahuje aktuální užití jazyka
- pro reprezentativnost musí být opravdu velký ( $10^9+$  slov)

# Zipfovo rozložení



# Využití korpusů

- komerční:
  - počítačová lexikografie – tvorba slovníků
  - informační systémy, analytické nástroje
  - inteligentní (mobilní) aplikace
- akademické:
  - univerzity pro výzkum a výuku jazyků
  - studenti jazyků

# Uživatelé



# Uživatelé

- lexikografové
- vědci
- učitelé
- studenti
- překladatelé
- terminologové
- textaři, kreativci

# Sketch Engine

- hlavní produkt LCL
- webová aplikace
- nástroj na zpracování a zkoumání textových korpusů
- jazykově nezávislý, velké množství funkcí pro různá použití
- stovky korpusů, desítky tisíc uživatelů
- pro kohokoli z MU volný přístup na <http://ske.fi.muni.cz>



# Sketch Engine

více než **600 korpusů** pro **89 jazyků**:

- 100+ korpusů větších než 100M pozic
- 30+ korpusů větších než 1G pozic
  - od roku 2010 korpusy TenTen ( $10^{10}$ )
- 60+ jazyků s anotací slovních tvarů
- 50+ jazyků s anotací slovních profilů (word sketches)
- vícejazyčné korpusy

# Word sketch

**water** (*noun*) British National Corpus freq = 34246 (305.3 per million)

| <u>modifier</u> | <u>9591</u> | 1.1   | <u>object_of</u> | <u>5126</u> | 1.6  | <u>subject_of</u> | <u>2835</u> | 1.7  |
|-----------------|-------------|-------|------------------|-------------|------|-------------------|-------------|------|
| hot             | <u>665</u>  | 10.17 | pump             | <u>92</u>   | 8.82 | flow              | <u>113</u>  | 9.29 |
| drinking        | <u>352</u>  | 9.97  | pour             | <u>139</u>  | 8.74 | drip              | <u>36</u>   | 8.33 |
| cold            | <u>459</u>  | 9.63  | drink            | <u>133</u>  | 8.55 | seep              | <u>30</u>   | 8.2  |
| boiling         | <u>237</u>  | 9.58  | heat             | <u>72</u>   | 8.43 | gush              | <u>23</u>   | 7.94 |
| fresh           | <u>231</u>  | 8.81  | boil             | <u>55</u>   | 8.12 | boil              | <u>25</u>   | 7.62 |
| mineral         | <u>173</u>  | 8.76  | tread            | <u>43</u>   | 7.88 | cascade           | <u>17</u>   | 7.48 |
| running         | <u>145</u>  | 8.74  | fetch            | <u>41</u>   | 7.51 | swirl             | <u>18</u>   | 7.45 |
| Thames          | <u>140</u>  | 8.54  | splash           | <u>29</u>   | 7.26 | lap               | <u>17</u>   | 7.43 |

# Formy spolupráce LCL a CZPJ

- bakalářské, diplomové, disertační práce pod vedením/na náměty LCL
- více než 40 obhájených BP a DP, 2 sponzorované PhD
- práce na společných grantech (EU projekt PRESEMT)
- vytváření, úpravy a anotace korpusů
- vývoj algoritmů nad korpusy

# Řešené problémy z pohledu informatika

- jak rychle zpracovat desítky TB dat
- jak taková data efektivně indexovat
- jak v indexech o desítkách GB hledat ve zlomcích vteřiny
- jak spočítat podobnost matic  $1M \times 1M$
- ...

# Vytvořené webové korpusy

enClueWeb09:

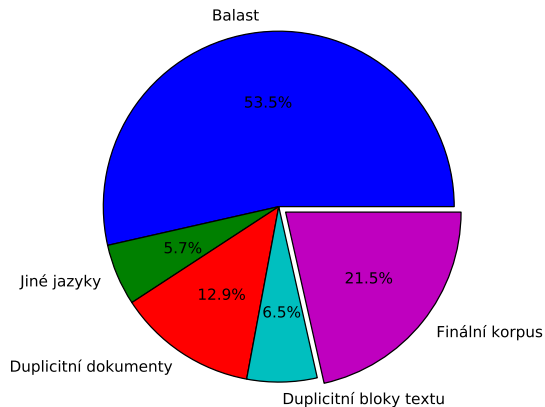
|              | Size [mil] |
|--------------|------------|
| total tokens | 81 990     |
| alphanumeric | 70 330     |
| numbers      | 1 485      |
| punctuation  | 9 849      |
| documents    | 138        |

korpusy řady TenTen:

|             |      |     |      |     |      |     |
|-------------|------|-----|------|-----|------|-----|
| jazyk       | cz   | de  | en   | es  | fr   | it  |
| mld. tokenů | 5,4  | 2,8 | 33,1 | 9,8 | 12,4 | 3,1 |
| jazyk       | jp   | pt  | ru   | sk  | zh   |     |
| mld. tokenů | 11,1 | 0,9 | 20,2 | 0,9 | 2,1  |     |

# Zpracování EnTenTen

původní velikost 14,6 mld slov, výsledek 3,2 mld slov



# SpiderLing

- webový crawler pro vytváření velkých korpusů
- zaměřuje se na webové domény bohaté na text
- paralelní zpracování získaných dat
- napsán v Pythonu

# Závěr

## Lexical Computing

- korpusová lingvistika a jazykové technologie

## Kontakt

- [milos.jakubicek@sketchengine.co.uk](mailto:milos.jakubicek@sketchengine.co.uk)
- [vojtech.kovar@sketchengine.co.uk](mailto:vojtech.kovar@sketchengine.co.uk)
- místnost S214