# Processing of Very Large Text Collections

## Miloš Jakubíček

Natural Language Processing Centre, Faculty of Informatics
Masaryk University, Brno, Czech Republic

jak@fi.muni.cz

Workshop of Natural Language Processing Centre &
Seznam.cz, a.s.
May 16, 2014

# Why to process natural language texts?

- **lots** of information, growing every day (web)
- need for **fast** and continuous knowledge mining
- **no time** for human intervention
- **large** data make statistical processing possible
- **real** data instead of false assumptions

# Text collection = a text corpus

- text collection: usually referred to as **text corpus**
- **humanities** $\rightarrow$ corpus linguistics, language learning
- **computer science** $\rightarrow$ effective design of specialized database management systems
- **applications** $\rightarrow$ usage of *any text* as information source

## goal

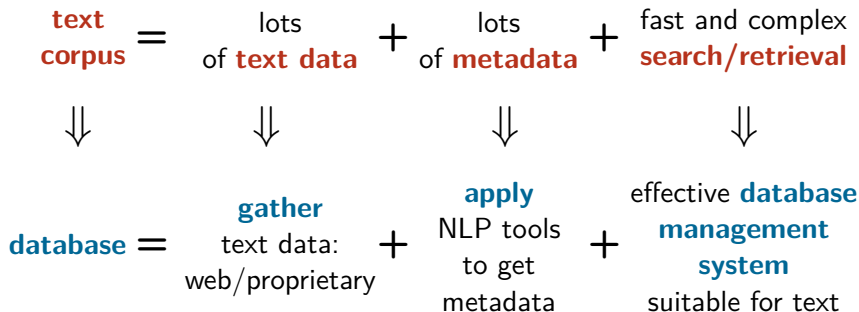| object of | 78390 | 3.0 |
|-----------|-------|-----|
| score | 8390 | 8.8 |
| achieve | | |
| concede | | |
| accomplish | 585 | 7.9 |
| reach | 1924 | 7.57 |
| net | 337 | 7.4 |
| pursue | 648 | 7.35 |
| grab | 406 | 7.33 |
| attain | 400 | 7.32 |
| pull | 504 | 6.69 |

| subject of | 25451 | 2.3 |
|------------|-------|-----|
| score | 903 | 8.8 |
| concede | 204 | 7.47 |
| gape | 105 | 7.3 |
| kick | 100 | 6.5 |
| orientate | 94 | 6.23 |
| rule | 78 | 5.5 |
| come | 175 | 5.21 |
| cap | 65 | 4.32 |
| beat | 20 | 3.69 |

| modifier |
|----------|
| actual |
| argue |
| winning |
| primary |
| secondary |
| strategic |
| common |
| realistic |
| achievable |

...loyer will seek to **achieve** three **goals** once employment...
...m in order to **achieve** the agreed **goal** of" sustainable de...
...ually distribution was to **achieve** its **goals** of peasant mobili...
...ction of how you can **achieve** those **goals** . So it's unlike a la...
... how you, how you could **achieve** the **goal** I mean it just says...
...order to **achieve** stated organisational **goals** . This definition (S...
...s concerned with **achieving** a specific **goal** in a given time us...
...t to be directed towards **achieving** the **goals** of the organisatio...
...ought to be structured to **achieve** their **goals** (Abrahamsson 19...
...maximizing profit and **achieving** other **goals** with which other p...
...derstand' the input. To **achieve** such a **goal** it is necessary to...
...ow you manage to **achieve** the same **goals** ,' said T.E. (he is a...
...ueezed. In a move to **achieve** these **goals** we have merged...
...ot stop until we have **achieved** this **goal** . In our sincere p...

# So what is a corpus?

$$\text{text corpus} = \text{lots of text data} + \text{lots of metadata} + \text{fast and complex search/retrieval}$$

$$\Downarrow \qquad \Downarrow \qquad \Downarrow \qquad \Downarrow$$

$$\text{database} = \text{gather text data: web/proprietary} + \text{apply NLP tools to get metadata} + \text{effective database management system suitable for text}$$

# Corpora

- **text type**
  - *general language* (gather domain independent information: common sense knowledge, global statistics, information defaults)
  - *domain specific* (gather domain specific information: terminology, in-domain knowledge, contrast to common texts)
- **timeline**
  - *synchronic*: one time period / time span ($\rightarrow$ what is up now?)
  - *diachronic*: different time periods / time spans ($\rightarrow$ what are the trends?)
- **language, written/spoken, metadata annotation type, . . .**

So is there any property one should aim at for all corpora?

So is there any property one should aim at for all corpora?

Yes – the size. The bigger, the better.

# Why does size matter so much?



Natural language phaenomena are not distributed evenly.

# Corpora now

Corpora at NLP Centre:

- **LARGE**: billions ($\sim 10^{10}$) of words
- **COMPLEX**: muti-level multi-value annotation, wide range of languages

Corpora at NLP Centre:

- **LARGE**: billions ($\sim 10^{10}$) of words

Corpora at NLP Centre:

- **COMPLEX**: muti-level multi-value annotation, wide range of languages

A big need for search/retrieval that is:

- **INTELLIGENT**: complex searching involving large amounts of metadata
- **VERY FAST**: parallel and distributed processing
- **ACCESSIBLE**: interfaces for automatic processing via third-party tools

# Applications

- **information systems** (going beyond fulltext search)
- **information analytics** (opinion mining, marketing assessment)
- **intelligent text processing** (predictive and adaptive writing, correction tools, effective writing in mobile devices)
- **computer lexicography** (better dictionaries, larger dictionaries)
- **machine translation** (parallel corpora)
- **statistics** for enhancing NLP tools

# What can we offer?

Ready-made tools for corpus building, management and effective search:

- **Building:** from own data/from the web, crawling, cleaning, deduplication
- **Management:** effective indexing in special DBMS
- **Search:** very fast evaluation of complex queries, keywords extraction, extraction of semantically related words, word sketches

Most of the tools are part of Sketch Engine, a product developed in collaboration with Lexical Computing Ltd.

# Demo: Sketch Engine
compare and contrast words visually



| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| informative | 0 | 4 | – | 6.6 | artifically | 30 | 12 | 5.1 | 3.7 |
| perceptive | 0 | 3 | – | 6.6 | extremely | 14 | 5 | 6.1 | 4.6 |
| cultured | 0 | 3 | – | 6.6 | very | 272 | 74 | 6.8 | 4.9 |
| knowledgeable | 0 | 3 | – | 6.4 | emotionally | 76 | 20 | 6.2 | 4.3 |
| humorous | 0 | 3 | –– | 6.3 | fiercely | 72 | 16 | 5.4 | 3.3 |
| dedicated | 0 | 4 | –– | 6.3 | particularly | 11 | 0 | 5.0 | 3.0 |
| charming | 3 | 10 | 5.8 | 7.2 | rather | 13 | 0 | 5.2 | 2.8 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| being | 70 | 10 | 5.9 | 3.1 |
| robot | 62 | 7 | 6.0 | 2.9 |
| agent | 3 | 0 | 5.7 | 2.8 |
| guess | 5 | 0 | 5.8 | 2.6 |
| conversation | 3 | 0 | 5.8 | 2.5 |
| human | 4 | 0 | 6.3 | 2.2 |
| creature | 3 | 0 | 6.4 | 2.0 |

witty **clever** 6.0 4.0 2.0 0 -2.0 -4.0 -6.0 **intelligent**

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| sexy | 3 | 3 | 6.6 | 6.2 | terribly | 3 | 0 | 5.8 | 2.4 |
| ambitious | 5 | 4 | 6.4 | 5.8 | pretty | 8 | 0 | 5.9 | 2.0 |
| amusing | 5 | 4 | 7.2 | 6.5 | jolly | 3 | 0 | 6.7 | 1.8 |
| clever | 10 | 4 | 7.2 | 5.7 | that | 11 | 0 | 6.8 | 1.1 |
| subtle | 6 | 0 | 6.4 | 5.4 | damn | 5 | 0 | 6.8 | 1.0 |
| brave | 6 | 0 | 6.6 | 5.1 | awfully | 4 | 0 | 7.0 | 0.0 |
| devious | 3 | 0 | 7.1 | 0.0 | extraordinarily | 6 | 0 | 7.4 | 0.0 |
| cunning | 5 | 0 | 7.7 | 0.0 | fiendishly | 5 | 0 | 7.9 | 0.0 |

| | | | | |
|---|---|---|---|---|
| fellow | 11 | 0 | 6. | 1.7 |
| pass | 4 | 0 | 6.6 | 1.5 |
| wordplay | 4 | 0 | 6.7 | 1.2 |
| chap | 10 | 0 | 7.0 | 1.1 |
| snap | 7 | 0 | 7.1 | 0.8 |
| twist | 18 | 0 | 7.2 | 0.6 |
| kick | 8 | 0 | 7.7 | 0.1 |
| trick | 12 | 0 | 8.2 | 0.0 |

# Demo: Sketch Engine
build specialised corpora instantly from the Web

# test *(verb)*   **enClueWeb (full) freq = 6180301** (74.8 per million)

| Lemma | Score | Freq |
|---|---|---|
| evaluate | 0.532 | 4453262 |
| analyze | 0.475 | 3595762 |
| monitor | 0.467 | 4047771 |
| examine | 0.455 | 5078101 |
| investigate | 0.453 | 3907848 |
| utilize | 0.439 | 4047715 |
| maintain | 0.438 | 10975886 |
| introduce | 0.435 | 8263900 |
| assess | 0.43 | 3196297 |
| demonstrate | 0.426 | 5668643 |
| identify | 0.423 | 10722177 |

# Conclusions

- Text corpora represent a **valuable information source** useful for many practical applications

- Corpora as text databases require **special solutions** that are fast and powerful

- There are number of **tools developed in the NLP Centre** for corpus building, management and efficient search