

Type-based Search of Idiomatic Expression

Jan Bušta

December 7, 2013

7th Workshop on Recent Advances in Slavonic Natural Language Processing

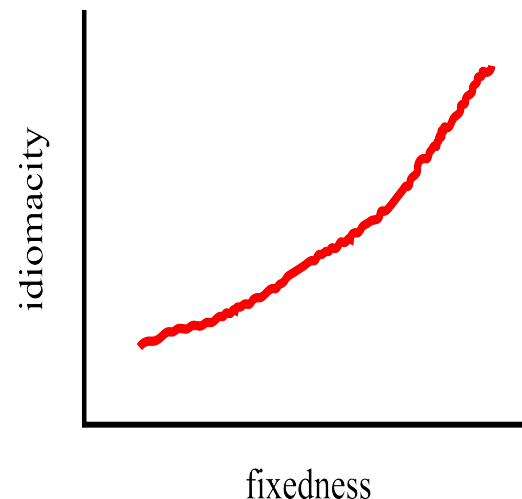
Karlova Studánka

Idiom

- Frege – compositionality
- Čermák – abnormality
- Fixedness – computer processing
 - type-based vs. token-based
- Token based: verb-noun
 - transitive verb (Verbalex: verbs allowing pasivation) + it's direct object (accusative case)

Fixedness

- How much we can change the MWE? Will the change preserve the idiomaticity?
- More fixed more idiomatic
 - Lexical fixedness
 - Syntactic fixedness



Methods

- PMI

$$PMI(v, n_j) = \log \frac{|\mathcal{V} \times \mathcal{N}| f(v, n_j)}{f(v, *) f(*, n_j)}$$

- PMI-based lexical fixedness

$$Fix_{lex}(v, n) = \frac{PMI(v, n) - \overline{PMI}}{s}$$

- Intersegment-based syntactic fixedness

$$Fix_{syn}(v, n) = \frac{\max(f(v, n), f(v, d, n), f(v, d, d, n))}{\sum f(v, n), f(v, d, n), f(v, d, d, n)}$$

Results

Candidate phrases better than idioms.
(SYN 2000)

Lexikální fixnost – extrahované fráze

znát slabina
 nanést krém
 * najít hrob
 vystřelit světlice
 být mašina
 mít svědek
 mít jiskra
 být výstup
 mít opletačka
 vmíchat kůra
 zlomit odpor
 dávat inzerát
 mít čelo
 * mít kůže
 získat klíč
 vypálit rána
 * najít jmenovatel

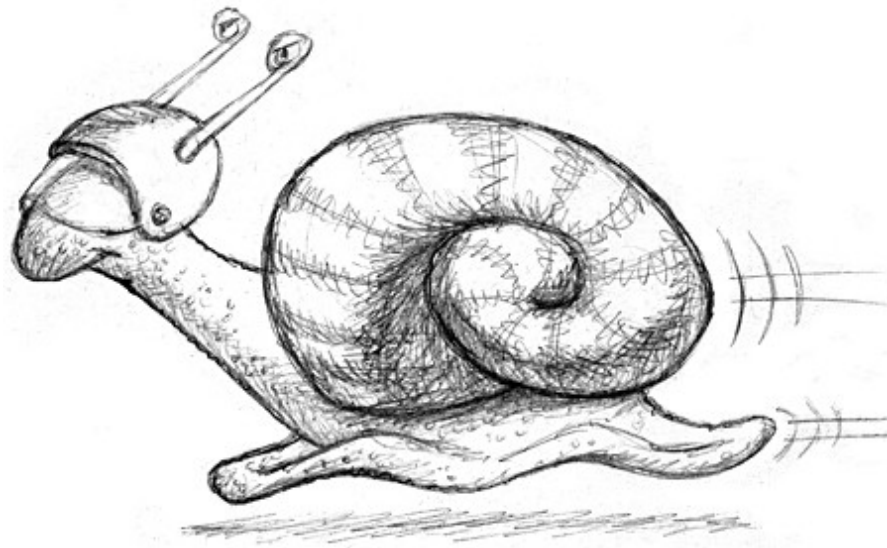
získat kondice
 nabírat vzduch
 * zachránit kůže
 zničit úroda
 zrušit pořadník
 nastartovat kariéra
 zatáhnout břicho
 mít ambice
 vyhrát poločas
 získat výpomoc
 nalévat voda
 vytvářet mapa
 vytvářet semeno
 dělat milost
 zastavit voda
 natočit šnyt
 vmíchat kečup

mít koncese
 mít rám
 * zamést stopa
 vlastnit karta
 mít návratnost
 vybrat kombinace
 mít paluba
 * vidět smrt
 mít důstojnost
 změnit chod
 mít rozchod
 mít království
 mít opice
 zachycovat proměna
 být mládež
 mít bratranec
 mít obojek

Conclusion

- Corpus-based fast generating **idiom candidates**.
- Annotators decide (yes/no decision).
- Faster building the lexicon of idioms.
- Up-to-date language processing (no useless idioms anymore).

Questions



Thank you for your attention.