

Korpusová lingvistika

PLIN059

Mgr. Dana Hlaváčková, Ph.D.

Korpusová lingvistika

- využívá pro studium jazyka velké soubory elektronických textů
 - texty odrážejí a dokládají reálné užívání jazyka
 - korpusy jsou **deskriptivní** (vs. preskriptivní)
 - **korpusové manažery** umožňují data prohlížet a třídit a poskytují statistické údaje
1. podstatná část počítačové lingvistiky – korpusy poskytují **zdroj jazykových dat**
 2. studium jazyka založené na jeho **přirozeném kontextovém užívání**
 3. **metodologický přístup** ke zkoumání jazyka

Jazykový korpus

Rozsáhlý soubor elektronicky uložených jazykových dat, obvykle označovaný, organizovaný se zřetelem k využití pro určitý cíl, vůči němuž je také považován za reprezentativní.

Čermák, F. Jazykový korpus: Prostředek a zdroj poznání. In *Studie z korpusové lingvistiky*. Praha: Karolinum, 2000, s. 15–38.

Přednosti korpusů

- velký **rozsah** s možností dalšího rozšiřování
- jazyková data v **přirozené** kontextové podobě
- převaha **typických** jazykových jevů nad **okrajovými**
- reprezentativní korpus je schopen zachytit **variabilitu** jazyka
- zrychlení a usnadnění lingvistické práce
- **morfologické** a **syntaktické** značkování korpusů zvyšuje jejich informační hodnotu

<S>
Náměstí
republiky
je
přímo
jejich
skanzenem
<g/>

.

</S>
<S>
Průčelí
je
tvořeno
divadlem
Antonína
Balšánka
<g/>

,

vystavěno
bylo

v
letech
1906
až
1909
<g/>

.

</S>

Základní pojmy

- **token, pozice** – řetězec znaků oddělený z obou stran mezerami
- **tokenizace** – proces rozdělení textu na tokeny
- **vertikál** – textový soubor (.vert), ve kterém je text rozdělen na tokeny
- **strukturní atributy** – informace, které se vztahují k textu (hranice dokumentů, hranice vět, ...)
 - **korpusová metadata** (žánrové dělení textů, autor, název díla, ...)
- **poziční atributy** – informace, které se vztahují k jednomu tokenu (word, lemma, tag, ...)
- **korpusový prohlížeč, korpusový manažer** (Bonito, Bonito2, Sketch Engine, KonText)

čisté jako z prádelny . Když se oblékáš , strakatá vrhla na trávník a protahovala si ruce i nohy jako , Regina , Leopard , Jupiter , SmČlandský lev , podařilo odplout . Kapitán Tönnes Speck na lodi Kattan (závod vyhraje . V červnu 1649 opustila loď Kattan (holubího . Malý mužík pozoroval celých pět minut kabelku jako mnou zacházeli , jako kdybych byl něco , co přitáhla mně chovala trochu slušně , byla upírka . " " ale v pádu přešel do salta nazad . Dopadl jako " Černocha . " " Nabral jsem ho hned před píchaly ho do vědomí . Talibe mhouřil oči jak divoká Jindřišky . Za košili . Přesně tím pohybem , jakým Jsi si jistý , že žádnému z mých sousedů neschází vrčení gazíku . Stáli tak , dokud se odněkud neobjevila dva fosforeskující světelné body , jako by se tam skrývala vítr hnul před sebou kupu suchého listí , vyděšená černá okna . Futaki je pořád uvnitř . K řediteli šla , ještě sem ji tu neviděl , co tu taková . „ Micur ! " Na kuchyňském stole seděla černá Šem silnější ! " – hlesklo ií hlavou .

kočka vyskočí na umyvadlo a tře se ti hlavou o bok
kočka . Ležet bez hnutí na lehátku , daleko od svěží , Tygr , Měsíc , Koruna , Klič , Stockholm
Kočka) se o to chtěl pokusit . Se svými dvaadvaceti **Kočka**) Švédsko a zamířila do Severní Ameriky . Loď byla pohybující se klubíčko a pouhá představa jejího obsahu ho hypnotizovala **Kočka** z ulice , a jediná osoba , která se ke **Kočka** ne , " řekla lady Sibyla . " Cože ?
Kočka , vrhl se zpět k užaslému Karotkovi a udeřil ho **Kočka** Barem , na Padesátý čtvrtý a Broadway , a on **Kočka** , ale i tak mu bělostné biče světla pronikaly skrz zdvívá kočata za kůži na hřbetě . " Co si **Kočka** nebo pes ? " vyptával se dál . " Znám , sedla si mezi ně a začala si olizovat tlapky **Kočka** . „ Běž ven a pošli ji do prdele !
Kočka se hbitě protáhla plotem u ředitelova domu . Odstrčil knihu **Kočka** , ještě sem ji tu neviděl , co tu taková **Kočka** k čertu hledá ? ! Zřejmě se něčeho lekla , **Kočka** a z červeného hrnce veselé chlemtala zbytek paprikáše od oběda **Kočka** k ní nízhéhla a třela se ií o nohv

konkordance, konkordanční řádek, konkordanční seznam

KWIC – key word in context (hledaný výraz v korpusu)

Typy korpusů

- druh zachycené komunikace
 - psané (written corpora)
 - mluvené (spoken corpora)
 - multimodální
- časový záběr
 - diachronní
 - synchronní
- účel
 - všeobecné
 - specializované
- způsob vytvoření
 - tradiční
 - webové
- jazyk
 - jednojazyčné
 - paralelní
 - srovnatelné
- možnost rozšíření
 - uzavřené (referenční)
 - otevřené (nereferenční)
- značkování
 - tagging (POS tagging, morfologie)
 - parsing (syntax, treebank)
 - alignment (párování)

Reprezentativnost korpusů

- v závislosti na účelu korpusu (kvantita a kvalita)
- národní korpusy – obraz užívání jazyka
- malý vzorek vzhledem k celku jazyka, nezobrazuje užití jazyka v celé šíři
- snaha zachytit **variabilitu** textů (beletrie, odborná lit., publicistika)

	SYN2000	SYN2005, SYN2010	SYN2015, SYN2020
publicistika	60 %	33 %	33,33 %
odborná lit.	25 %	27 %	33,33 %
beletrie	15 %	40 %	33,33 %

Tvorba korpusů

- korpusy tradiční a webové
- sběr dat
 - poskytovatelé textů
 - webové korpusy – stahování textů (crawler)
- sjednocení formátu a kódování
- odstranění netextového obsahu (boilerplate)
- odstranění duplicitních textů (webové korpusy)
- interní anotace
- tokenizace (vertikál) – lemmatizace – externí anotace (značkování)
- mluvené korpusy – nahrávky, přepis, synchronizace textu se zvukem

Korpusové manažery v ČR

- ÚČNK – ČNK – KonText
 - <http://kontext.korpus.cz>
- FI MU – Sketch Engine
 - <https://www.sketchengine.eu/>
- Český národní korpus
 - <https://www.korpus.cz/>