

K možnostem počítačového zpracování literárního textu

— Petr Pořízka —

Jednu z možností, jak dnes efektivně analyzovat literární text s použitím počítačů, nabízí korpusová lingvistika – obor, jehož hlavním předmětem zájmu jsou tzv. korpusy. Ty lze definovat jako metodologicky jednotné soubory textů (nejčastěji v elektronické verzi), sloužící k analýze prostřednictvím speciálních softwarových programů, tzv. konkordančních programů či korpusových manažerů.¹ Možnosti softwarových analýz jsou přitom určeny či podmíněny nejen počítačovými programy, ale (a to je třeba zdůraznit) i způsobem zpracování samotného textu. V zásadě lze pracovat s textem prostým (s původním

¹ Jako korpusové manažery označujeme softwarové nástroje, jež umožňují komplexní práci s korpusem. Text je prostřednictvím tohoto nástroje zpracován a načten do implicitní databáze programu. Texty lze poté podle různých kritérií (filtrováním) prohledávat a analyzovat. Vyhledávat lze jednotlivé výrazy či slovní spojení, jež jsou zobrazeny i v omezeném kontextu (výsledkem jsou tzv. konkordance), vyhledávat lze také podle provedené anotace (viz pozn. 2) nebo s pomocí tzv. masek – speciálních metaznaků dotazovacího jazyka korpusového programu: blíže viz např. <http://korpus.cz/bonito/regular.php>. Provádět lze i statistické analýzy, vytvářet frekvenční či abecední slovníky ad. Standardně tyto nástroje nabízejí i možnost uložit výsledky analýz do externího samostatného souboru v některém z běžných textových formátů (.txt, .rtf, .doc).

textem, do něž nejsou tvůrcem korpusu vnášeny dodatečné metatextové informace) nebo anotovaným (originální text je obohacen o meta-informace vnětového, ale i vnitrotextového charakteru).² Poměrně efektivně lze pracovat i s neanotovaným korpusem, provedená anotace ale může výrazným způsobem rozšířit možnosti analýz.

Během více než patnáctileté historie korpusové lingvistiky v České republice vznikla řada především jazykových/jazykovědných korpusů.³ Postupně se ale metody a nástroje korpusové lingvistiky začínají aplikovat i na literárněvědné projekty. Nejznámějším a dosud bezesporu nejvýznamnějším literárněvědným korpusovým projektem je *Česká elektronická knihovna* (ČEK), fulltextová databáze české poezie 19. a počátku 20. století.⁴

Právě projekt ČEK nejlépe dokumentuje možnosti, jež počítačové zpracování textů a softwarové nástroje korpusové lingvistiky nabízejí. Uživatel může nejen vyhledávat výrazy prostým vepsáním do dotazového řádku webového rozhraní, ale lze vytvářet strukturované dotazy, tj. vyhledávat jak ve vybraných sbírkách, tak v jejich strukturách: básních, strofách, verších [1].

Aplikace, jež spravuje textovou databázi, provádí i statistické analýzy – nabízí údaje o minimální, průměrné a maximální délce slova, verše či strofy [2], umožňuje generovat frekvenční a abecední seznamy [3] a vyhledávat vazby slov a motivů.

Kromě toho lze volit mezi diplomatickou a ediční verzí textu nebo vpisovat uživatelské poznámky. Každý dokument je doplněn o ediční poznámky a ilustrace, jež jsou součástí knižního vydání. K textům je možno přistupovat skrze webové rozhraní s implementovaným korpusovým manažerem od firmy inSophy [4].⁵

Jedním z projektů, jež by mohly být zajímavé pro literární vědce, je *Intercorp* – databáze paralelních korpusů obsahující literární

2 Anotací, její charakteristikou, možnostmi a typy jsme se blíže zabývali ve studii Pořízka – Schäfer 2010b.

3 V roce 1994 byl založen Ústav Českého národního korpusu (ÚČNK), ale aktivity v oblasti počítačového zpracování jazyka vznikaly již řadu let předtím – mezi iniciátory patřili např. prof. František Čermák, doc. Karel Pala, doc. Vladimír Petkevič ad. Informace o stávajících korpusech ÚČNK lze nalézt na <http://www.korpus.cz/struktura.php>.

4 Iniciátory projektu byli Vladimír Macura a Pavel Janoušek; bližší charakteristika a popis projektu viz <http://www.ceska-poezie.cz/cek/>. Nutno ocenit, že stejně jako v případě korpusů ÚČNK je i k *České elektronické knihovně* poskytován bezplatný přístup. Databáze ČEK je o to cennější, že se dostaneme k plně verzi textů, na rozdíl od praxe ÚČNK, kdy je možno zobrazit v korpusech prostřednictvím korpusového manažeru jen hledaný výraz s omezeným kontextem.

5 [Http://www.insophy.cz/](http://www.insophy.cz/).

Strukturované vyhledávání

Vyhledej všechny ,
 které obsahují ,
 ze slov ,
 ,
 ze slov .

uzavřít toto okno po vyhledání [Nápověda](#)

Poznámka: Lze použít zástupné symboly (? = libovolný znak, * = žádný nebo více znaků).

[1] Strukturované vyhledávání České elektronické knihovny

Česká elektronická knihovna
 Pinotextová databáze české poezie 19. a počátku 20. století

Statistika (1 sbírka):
 Březina, Otokar: Tajemné dálky, 1895

[Celkový přehled](#) | [Abecední slovník](#) | [Frekvenční slovník](#)

Celkový přehled ?

	Počet	Počet 1-N písmenných slov:
Počet básní:	30	1 903
Počet strof:	191	2 619
Počet veršů:	883	3 819
Počet slov:	6683	4 1098
Délka strofy (min.):	1	5 1006
Délka strofy (prům.):	4,62	6 881
Délka strofy (max.):	18	7 587
Délka verše (min.):	1	8 386
Délka verše (prům.):	7,57	9 254
Délka verše (max.):	14	10 88
Délka slova (prům.):	4,51	11 26
		12 10
		13 4
		14 1
		15 1



© 2005-2007 Ústav pro českou literaturu AV ČR

[2] Statistika *Tajemných dalek*
 Otokara Březiny

Česká elektronická knihovna
 Pinotextová databáze české poezie 19. a počátku 20. století

Statistika (1 sbírka):
 Březina, Otokar: Tajemné dálky, 1895

[Celkový přehled](#) | [Abecední slovník](#) | [Frekvenční slovník](#)

Frekvenční slovník ?

Slovo	Počet výskytů
a	363
v	345
se	170
jak	114
z	103
na	83
mi	51
do	50
mě	46
jem	45
jež	42
ve	36
nad	33
duše	30
jenž	30
kde	29
kyž	29
duši	27
tvých	26
mých	24
ti	24

[3] Frekvenční slovník téže sbírky

(beletristické) texty.⁶ *Intercorp* vzniká v ÚČNK jako součást *Českého národního korpusu*. V současné době obsahuje texty z 22 jazyků, přičemž čeština má v korpusu pozici tzv. pivotu – česká verze (originál nebo překlad) je vztažena k jedné nebo více verzím cizojazyčným [5].

6 Detailnější informace o projektu *Intercorp*, jeho struktuře, koncepci, užitých aplikacích, tvůrcích apod. lze nalézt na <http://korpus.cz/intercorp-info.php>, příp. <http://www.korpus.cz/intercorp/>. Na tvorbě tohoto paralelního korpusu se velkou měrou podíleli krom spolupracovníků ÚČNK pedagogové a studenti FF UK Praha.

Česká elektronická knihovna
 Pínotová databáze české poezie 19. a počátku 20. století

uživatel: porizka | jcd@lilid.cz

Uživatelská nastavení | Správa filtrů | Správa uživatelských pozámek | Historie

Seznam sbírek -7-
 Zobrazeno 112 celkem 1700 sbílek

Filtrování: **aktívní - Března**
 Nový | Vybrat | Bez filtrů

Almanachy a rok NČM, 1999
 ▶ Almanachy secese, 1999
 ▶ BŘEZNA, Otakar: Básnické spisy, 1949
 ▶ BŘEZNA, Otakar: Básnické spisy 2, 1, 19;
 ▶ BŘEZNA, Otakar: Básnické spisy 2, 2, 19;
 ▶ BŘEZNA, Otakar: Prvoloty, 1933
 ▶ BŘEZNA, Otakar: Ruce, 1901
 ▶ BŘEZNA, Otakar: Statistické chrám, 1999
 ▶ BŘEZNA, Otakar: Světlní na západě, 1899
 ▶ BŘEZNA, Otakar: Tajemné dáky, 1895
 ▶ BŘEZNA, Otakar: Věry od půli, 1897

Vybrané sbírky -7-
 Března, Otakar: Ruce, 1901
 Března, Otakar: Statistické chrám, 1999
 Března, Otakar: Světlní na západě, 1899
 Března, Otakar: Tajemné dáky, 1895
 Března, Otakar: Věry od půli, 1897

Zobrazení
 Zobrazit označené

Statistika
 Zobrazit statistiku
 Abecední slovník
 Frekvencní slovník

Strukturované hledání
 duše
 Roziřené vřehledání
 Vyhledat

Fulltextové hledání
 Hledaný text
 Vřv těle Vřv hlasůvkách
 Vyhledat

Kontexty
 slovo pro kontexty
 Vyhledat kontexty

Uložení výběru
 Uložit jako výchozí
 Obnovit výchozí

Práce s výběrem
 Smal všechny
 Smal označené

Třídění podle:
 autora názvu roku
 vzestupně sestupně
 Třídít | Nastavit jako výchozí

Vybrat všechny

[4] Webové rozhraní (pracovní prostředí) České elektronické knihovny

Jazyky | Dokumenty

intercorp_cs Vůbec_vše / Odeber_vše
 intercorp_bg adami-stopanus_přivodc
 intercorp_de angelova-dvoj_svit
 intercorp_en brown-chut_lasky
 intercorp_es brown-stravin_temnoto
 intercorp_fr cermak-zaklady_metod
 intercorp_it dark-epekani_ramou
 intercorp_ja day-cirkus_v_time
 intercorp_lv feldingova-panerla
 intercorp_nl franzen-roztreseni
 intercorp_pl graham-partner
 intercorp_pt halley-konecna_dag
 intercorp_ru irving-rcik_vidovou
 intercorp_sv jrotka-saturan
 intercorp_sy silham-mez_medvedy
 krentz-zajato_snu
 kundera-Nesmetnost
 kundera-zert
 Lindeyova-Zamlouvaný
 Ondastay-Anglicky_Pao
 otomasek-romeo_julie
 Palatrusk-cakrnost
 pavic-chazarsky_slov
 roblant-mikeno_benat
 rowlingova-top_kamen
 searle-mysl_mozek_veda
 Steel-Druha_sance
 Steel-Strazy_andel
 SYNDICATE
 Topol-Kodakatajolechj
 Tuku-tbetskaj_metody
 Velegh-VychovaDveleCR
 woolfova-dalcovayova
 woolfova-mez_alky
 Woolfova-strasdelj

[5] Intercorp: V levém sloupci je možno filtrem zvolit subkorpusy (zde český a anglický). V pravém bloku jsou texty dostupné v obou jazycích, v nichž je možno paralelně vyhledávat.

Korpus: intercorp_cs | **Korpus: intercorp_en**

Lemma: nesmetnost | Lemma:

Slovní spojení: | Slovní spojení:

Word Form: match case: | Word Form: match case:

CQL: | CQL:

Default attribute: word | Default attribute: word

Rádi-0 na stránce 10 |

[6] Pracovní prostředí Intercorp – dotazovací okno webového rozhraní Park

intercorp_cs (404051 nalezen)	intercorp_en (470459 nalezen)
Účel: nezobrazit - kóp	intercorp_en
Věra v nesmetnost vnesla relativitu do rozdílu mezi životem a smrtí.	Belief in immortality made the opposition between life and death relativ... intercorp_en
Vědomím vyjádřením tohoto snu je obsažená realita na jistou biologick... smrti - přehledně, že dostatečně velká věra v nádej vřdy ponaci samotnou smrt, neboch vřdy přinese jistou formu nesmetnosti .	The conscious expression of this dream is an obsessive response to the certainty of biological death - the belief that a big enough vein in the game of space will bear death itself by conferring a form of immortality on the winner intercorp_en
Přetváru vědecké nesmetnosti slizce přeměnilo objevu neoddělitel... Hábilj, staršho a vřce neveděckého snu o unikatu vlastní nevymřutelné smrti se nelí nřdění hájka popření.	Nothing but the thinnest membrane of denial separates the notion of scientific immortality through priority of discovery from the deeper, older, and wholly non- scientific dream of escaping one's own mortal death intercorp_en
Třebaže dnešní vědecká torž, že hledá způsob nápravy neoddělitel... přijdu, až přijde mnohá zřho, jako j křdělčuj, se chová tak, jako by jehož skutečným cílem bylo pouze získat nevřdučnou přeměnu, ač už to znamená nebo odřadu stop křděl.	Though today's biomedical science claims to search for repairs of nature's 'defects', too many of its practitioners behave as if their real purpose were only to gain the mythical immortality of precendence, at whatever cost to themselves or others intercorp_en
Sadit se získal výjimečné místo v dějinách a dosáhl nesmetnosti v okamžiku kdy užil z tohoto pohleduho slizce nesmetnosti , jakož i z partnerské sodržady a vypracované vědecké soudržnosti arabůžých summá.	Sadlet gained a privileged place in history and achieved immortality the moment he fled from the comfortable prison of reality - and from the partners' solidarity and hollow metaphorical cohesion of Arab summits intercorp_en
"V dětství, " napsal sám Greene (v Ministerstvu strachu), " žijeme v jasu nesmetnosti - nebo je tak blžže a skutečné jako mrtvost žráh.	In childhood, " Greene himself had written (in the Ministry of Fear), " we live in jasu of immortality - heaven is as near and actual as the seaside intercorp_en
Žle oděle stranou republikovanoj studij, o nř Martin Belvedereho doufá, že mu zápas nesmetnosti , až vřde v nřkterém uznaměnaném vědeckém	He put aside the unpublished paper that Martin Belvedere had no doubt hoped to see immortalized in one of the respectable journals of sleep and dream research and sank intercorp_en

[7] Výsledek vyhledávání v Intercorp: paralelně zaruované konkordance českého a anglického subkorpusu.

Korpus je přístupný přes webové rozhraní *Park* (autor Michal Štourač) [6], jež je nadstavbou nejužívanějšího českého korpusového manažeru *Manatee*.⁷ *Manatee* je komplexním korpusovým nástrojem a vždy záleží na tvůrcích daného korpusu, které možnosti programu využijí. Projekt *Intercorp* je zpracován způsobem, který umožňuje využít všechny základní funkce manažeru: je možné specifikovat prohledávané části korpusu – jazyky i konkrétní texty (použitím filtrů), vyhledávat podle řady kritérií – v jednom či více jazycích současně, podle slovního tvaru, frází či posloupností tvarů, podle dotazovacího jazyka programu *Manatee*,⁸ podle lemmatu a morfosyntaktické značky (tagu).⁹ Výsledky vyhledávání (konkordance) jsou zobrazeny jako paralelně zarovnané úseky textu ve zvolených jazycích. [7]

Statistické analýzy, mezi něž patří například výpočet absolutní a relativní frekvence výrazů, rozložení hledaného výrazu v korpusu apod., lze v programu *Manatee* aplikovat na jakýkoli text – anotovaný či prostý (neanotovaný), neboť jde o implicitní funkce, jejichž využití není přímo podmíněno mírou dodatečného zpracování textu.¹⁰

ÚČNK vydal i dvě lexikograficky zaměřené monotematické monografie – slovníky Karla Čapka a Bohumila Hrabala (Čermák 2007, Čermák – Cvrček 2009). Oba vyšly knižně s příloženým CD, jež je de facto elektronickou verzí tištěné knihy, CD-ROM dokonce ve srovnání s knihou obsahuje korpus doplněný o lemmatizaci a morfologickou anotaci.¹¹ V obou případech se jedná v podstatě o abecedně uspořádaný frekvenční slovník. Charakteristiku obou slovníků (metodologicky se od sebe neliší) podává Štíchova recenze slovníku Karla Čapka

-
- 7 Autorem programu *Manatee* je Pavel Rychlý. Až na výjimky používají korpusový manažer *Manatee* všechny české korpusy, jde tedy o jakýsi český softwarový standard. Ke korpusům uloženým pod systémem *Manatee* lze přistupovat dvěma způsoby: přes grafické uživatelské rozhraní *Bonito* (<http://www.textforge.cz/products>), nebo skrze novější verzi, webové rozhraní *Word Sketch Engine* (<http://www.sketchengine.co.uk/>, autoři Pavel Rychlý, Adam Kilgarriff a Jan Pomikálek). Zatímco systém *Manatee/Bonito* je poskytován zdarma, užití aplikace *Word Sketch Engine* je zpoplatněno. To je důvod, proč se i nadále ve většině korpusových projektů používá vývojově starší verze *Manatee/Bonito*.
 - 8 Jde o tzv. *Corpus Query Language* (CQL), dotazovací jazyk vyvinutý na univerzitě ve Stuttgartu při práci na konkordančním programu *Xkwoic*. Tento nástroj se stal základem systému *Manatee*. Vývoj projektu *Xkwoic* v současné době pokračuje pod názvem *The IMC Open Corpus Workbench* – viz <http://cwb.sourceforge.net/>.
 - 9 Otázkám lemmatizace a morfologického značkování jsme se věnovali ve studii Pořízka – Schäfer 2010a.
 - 10 Zevrubnější popis základních statistických funkcí *Manatee/Bonita* viz <http://korpus.cz/bonito/stat.php>.
 - 11 Nabízí se tedy i otázka ekonomického charakteru, neboť by slovníky mohly být vydány pouze na CD-ROMu, což by nepochybně přineslo velkou úsporu nákladů.

uveřejněná v *Naší řeči*: „Hlavní částí *Slovníku Karla Čapka* je abecedně uspořádaný slovník všech slov, která Čapek užil ve svém publikovaném díle literárním (próza, drama, poezie), odborném, v rozsáhlé publicistice i ve vydané soukromé korespondenci. U každého z těchto slov, řazených pod sebou ve dvou sloupcích na stránce, je v sedmi sloupcích vedle sebe uvedeno, kolikrát Čapek daného slova užil, a to nejdřív celkem a pak v šesti hlavních žánrech (próza, drama, publicistika, poezie, odborná literatura a korespondence)“ (Štícha 2009: 38–39).

K oběma slovníkům dodejme, že lze poměrně jednoduchým a nepracným způsobem vytvořit stejný abecední či frekvenční slovník libovolného autora s použitím volně dostupných konkordančních nástrojů (viz níže program *AntConc*).

Od roku 2008 se systematicky věnujeme možnostem tvorby malých, specializovaných korpusů, výsledkem těchto aktivit jsou i dílčí korpusové projekty, jež vznikají ve spolupráci se studenty na Katedře bohemistiky Filozofické fakulty Univerzity Palackého v Olomouci.¹² Jako první byl v roce 2008 sestaven korpus esejů Otokara Březiny,¹³ v roce 2008–2009 vznikl korpus esejů Ladislava Klímy¹⁴ a v roce 2010 jsme začali pracovat na výukovém korpusu Karla Čapka.¹⁵ Všechny korpusy jsou stále ve fázi budování a každý sloužil či slouží jinému účelu.¹⁶ Cílem, který spojuje všechny tři korpusy, bylo vypracovat efektivní postup tvorby malých autorských korpusů pro jazykovědné a literárněvědné účely. Budování korpusů ovšem zahrnuje několik etap a vyžaduje jak znalosti filologické, tak technické (kódování znaků, formát dat apod.). Jednotlivé korpusy tedy vznikaly jako pilotní projekty, jež měly prověřit možnosti zpracování textů korpusovými nástroji pro každou z klíčových oblastí tvorby textových korpusů.

12 Průběžné výsledky těchto aktivit jsme publikovali v následujících dvou studiích: Pořízka – Schäfer 2010a a Pořízka – Schäfer 2010b.

13 Zpracována byla prozatím první kniha esejů *Hudba pramenů* (1903). Korpus byl prezentován na mezinárodním sympoziu Otokar Březina 2008 v Jaroměřicích nad Rokytou (Pořízka – Schäfer 2010a).

14 Tento korpus obsahuje text knihy *Svět jako vědomí a nic* (1904), prezentován byl na vědecké konferenci Ladislav Klíma konané na FF UPOL v Olomouci (Pořízka – Schäfer 2010b).

15 Vznik čapkovského korpusu je motivován čistě didaktickými účely, zatímco předchozí byly zaměřeny metodologicky. Vzniká v experimentálním semináři, v němž si studenti osvojují potřebné know-how, jak korpusy vytvořit. Obsahuje jen studentské seminární práce – dílčí texty Karla Čapka s různou formou zpracování. Materiálovým zdrojem jsou digitalizované texty e-knihovny Městské knihovny v Praze; srov. <http://www.mlp.cz/karelcapck/> a rovněž http://www.mlp.cz/knihovna_on-line.htm.

16 Tyto dílčí korpusové projekty budou po dokončení zveřejněny na korpusovém portálu <http://corpus.upol.cz>.

Březinovský korpus byl zaměřen na lingvistickou anotaci: byla provedena lematizace (přiřazení reprezentativní formy slovních tvarů) a jednoduchá morfologická anotace přiřazením slovnědruhové interpretace. Použita přitom byla celá řada nástrojů urychlujících a automatizujících práci (počítačové skripty *fsm Tokenize*, *Annot1.pl*, *Annot1_to_Annot2.pl*) (Pořízka – Schäfer 2010a).

Klímovský korpus prověřoval možnosti typograficko-ediční a strukturně-obsahové anotace s pomocí značkovacího (meta)jazyka *XML (eXtensible Markup Language)*, který uživateli umožňuje definovat si vlastní sadu značek se speciálním významem. V *XML* tak byla anotována struktura textu a jeho hierarchizace (kniha a její části: kapitoly, oddíly, odstavce, věty, slova...) a editorské/typografické jevy: u esejů Ladislava Klímy například a) grecismy, latinismy; b) korektury: překlepy, chybějící uvozovky; c) řezy písma: kurziva, bold; d) uvozovky, interpunkce apod. (Pořízka – Schäfer 2010b).

Zároveň byly během sestavování obou korpusů vypracovány a algoritimizovány jednotlivé kroky pro přípravu a importování souboru textů do aplikace *Manatee/Bonito*, včetně přesných formátů zápisu pro použití počítačových skriptů. Tento software totiž pracuje s tzv. binárními texty, jež je třeba s pomocí speciálních nástrojů (skriptů) zkonvertovat z prostého textu (ve formátu *.txt*) do požadovaného formátu – především nástrojem *encodevert.exe* (součást aplikace *Manatee*). Až po této konverzi a řadě dalších úkonů (zápisy do registrů, deklarace kódování znaků, korpusu a jeho atributů apod.) je možno s korpusem pracovat – přitom je nutné, aby byl text připraven v tzv. vertikále (jedno slovo na jeden řádek) a ve formátu *.txt*.

Tyto postupy a procesy při přípravě a zpracování textů mohou být pro začínajícího či méně zkušeného uživatele velmi komplikované a i sebemenší chyba během technického zpracování znamená nefunkčnost celého korpusu. Naštěstí pro běžného (i začínajícího) uživatele není nutno znát zevrubně všechny technické aspekty počítačového zpracování korpusů – pro základní práci s textovými daty postačí prostý (původní, tj. neanotovaný) text a jednoduchý konkordanční software.

Jedním z takových programů je *AntConc* (© Laurence Anthony),¹⁷ aplikace umožňující pracovat s textem podobně jako *Manatee/Bonito* (v mnohém snese srovnání), ale tvorba korpusu je značně jednodušší – v podstatě triviální, neboť stačí text prostě jen importovat do aplikace.

17 Jedná se o výborný freewareový program, který je volně ke stažení ve verzi *AntConc3.2.1w* na webové adrese <http://www.antlab.sci.waseda.ac.jp/software.html>.

AntConc přitom pracuje s celou řadou formátů: *s.txt*, *.html*, *.xml*, importovat lze dokonce i wordovský *.doc* nebo *.rtf*, u nich ale nastanou s největší pravděpodobností problémy s kódováním (tj. se zobrazováním) znaků. *AntConc* umožňuje uživateli velmi snadno vyhledávat v textu jak slovní výrazy v kontextu (konkordance), tak s pomocí masek (speciálních metaznaků), provádět statistické výpočty, generovat abecední a frekvenční slovníky a další funkce.

Pro prezentaci základních možností tohoto konkordančního programu jsme zvolili Čapkovy *Povídky z jedné kapsy*.¹⁸ Text lze načíst do programu standardní cestou: dokument lze otevřít v menu *File* → *Open File(s)* pro jednotlivé texty, nebo lze zvolit cestu *File* → *Open Directory* a importovat do aplikace ze zvolené složky počítače všechny dokumenty najednou. Načtené soubory se objeví v levém sloupci (*Corpus Files*) programu *AntConc*. Aby se texty zobrazovaly korektně, je nutno zkontrolovat či nastavit tzv. kódování znaků: v menu *Global Settings* → *Language Encodings* → *Edit* lze zvolit jedno ze tří kódování, která se dnes pro české texty užívají.¹⁹ Poté již lze pracovat s textem a využívat všechny funkce, jež *AntConc* nabízí. Jak již bylo zmíněno, tou základní je možnost vyhledávat jednotlivé výrazy či slovní spojení (případně je kombinovat se speciálními zástupnými metaznakly). [8]

Obr. [8] zobrazuje kromě grafického prostředí aplikace jednotlivé konkordance (záložka *Concordance*) – výsledek hledaného výrazu „metod.“ (tečka představuje zástupný metasymbol s významem „jakýkoli znak“). Dotazový řádek se nachází v dolním panelu (*Search Term*), kde se také zobrazuje počet nalezených výskytů (*Concordance Hits*). Klíčové slovo je od kontextu barevně odlišeno.

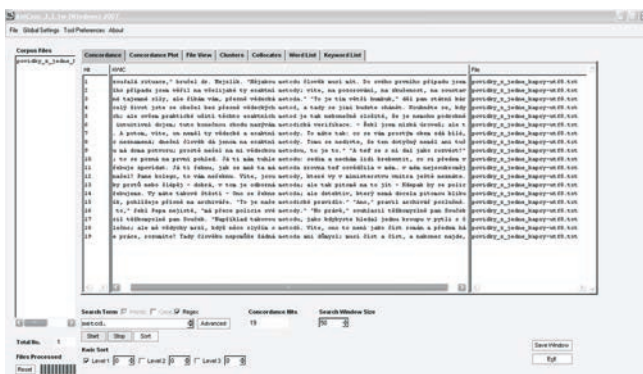
Velkou výhodou oproti systému *Manatee/Bonito* je možnost zobrazit celý text v úplnosti (záložka *File View*), nejen klíčové slovo s omezeným kontextem. [9]

Zajímavou funkci nabízí záložka *Concordance Plot*, jež zobrazuje rozložení hledaného výrazu, tj. jeho pozici v celém textu či korpusu – uživatel tak vidí, ve které části se daný výraz vyskytuje nejčastěji. [10]

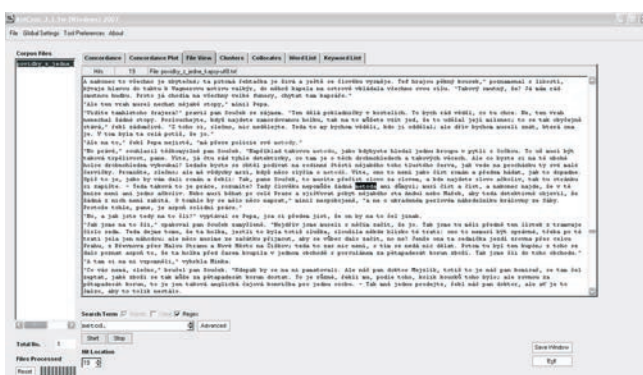
V souvislosti se slovníky Karla Čapka a Bohumila Hrabala jsme

18 Zdroj: http://www.mlp.cz/koweb/00/03/34/75/70/povidky_z_jedne_kapsy.txt.

19 Pro češtinu lze zvolit 1. univerzální kódování *Unicode (utf8)*, 2. *ISO Central Europe (iso-8859-2)* nebo 3. *windowsovské WinLatin2 (cp-1250)*. Při špatném kódování se znaky buď nezobrazují vůbec, nebo jsou zobrazeny nesprávně. Pokud uživatel neví, v jakém kódování je text zpracován (lze to jednoduše zjistit v některém z textových editorů), je třeba vyzkoušet postupně všechna tři kódování, až se všechny znaky zobrazí korektně.



[8] Grafické uživatelské rozhraní programu *AntConc* s konkordancemi



[9] Funkce *File View*, jež zobrazuje barevně odlišená klíčová slova v celém dokumentu, nikoli jako konkordance



[10] Rozložení výrazu „metod.“ v korpusu s pomocí funkce *Concordance Plot*

uvedli, že existuje jednoduchý způsob, jak z libovolného textu vytvořit abecední či frekvenční slovník. Tímto způsobem je tak možno získat slovník kteréhokoli literárního autora, máme-li k dispozici elektronickou verzi příslušných (literárních) textů. Oba typy slovníků, ale i slovník retrogradní lze v programu *AntConc* vygenerovat v záložce *Word List*, přičemž jednotlivé typy výstupů se nastavují v dolní části panelu funkcí *Sort by Freq* (frekvenční slovník), *Sort by Word* (abecední slovník) a *Sort by Word End* (retrogradní slovník). [11]

2010b „Svět jako vědomí a nic Ladislava Klímy v olomouckém korpusu české esejistiky přelomu 19. a 20. století“, *Aluze* (v tisku)

ŠTÍCHA, František

2009 „Nad slovníkem Karla Čapka“, *Naše řeč* 92, č. 1, s. 38–39

Computer processing of literary texts: the opportunities involved

This study adopts an interdisciplinary approach towards text and deals with the technical options involved in text processing, which enables us by means of software tools to provide data retrieval, and to perform statistical analysis and other processes in accordance with preselected criteria and on the basis of an annotation text. The first part is devoted to the most important corpus projects focusing on literary texts: *Czech Electronic Library*, the *Intercorp* corpus and lexicographical dictionaries of Karel Čapek and Bohumil Hrabal. The second part presents the basic possibilities for the creation of small corpora, demonstrated on corpora of Otokar Březina and Ladislav Klíma, and primarily the usage of a corpus concordancer called *AntConc* during analysis of (literary) texts: data retrieval, creating alphabetical and frequency dictionaries, etc.

Keywords

corpus linguistics, text processing, data retrieval, corpus concordances, creation of corpora, Czech corpora of literary texts