

KORPUS A KORPUSOVÁ LINGVISTIKA

FRANTIŠEK ČERMÁK

KAROLINUM

Korpus a korpusová lingvistika

František Čermák

Recenzenti:

prof. PhDr. Eva Hajičová, DrSc.

PhDr. Jitka Šonková, CSc.

PhDr. Věra Schmiedtová

Vydala Univerzita Karlova

Nakladatelství Karolinum

Redakce Lenka Ščerbaničová

Grafická úprava Jan Šerých

Sazba DTP Nakladatelství Karolinum

Vydání první

© Univerzita Karlova, 2017

© František Čermák, 2017

ISBN 978-80-246-3710-5

ISBN 978-80-246-3776-1 (online : pdf)



Univerzita Karlova
Nakladatelství Karolinum 2017

www.karolinum.cz
ebooks@karolinum.cz

Slovo úvodem

Text knihy zachycuje vznik a rozvoj korpusů a korpusové lingvistiky nutně k určitému datu, vývoj jde však rychle dál. Vyjadřovalo se k němu více lidí, můj hlavní dík však patří jejím recenzentům, kterými byli prof. PhDr. Eva Hajičová, DrSc. (MFF UK), PhDr. Jitka Šonková, CSc. (University of Iowa) a PhDr. Věra Schmiedtová (FF UK).

František Čermák

OBSAH

I. ČÁST: TEORIE ----9

1	ÚVOD: KORPUS A KORPUSOVÁ LINGVISTIKA ---- 10
2	DATA A INFORMACE ---- 12
2.1	Elektronický text ---- 13
2.1.1	Obecná povaha, vlastnosti a zpracování textů ---- 13
2.2	Korpusová data ---- 16
2.2.1	Povaha a zdroje elektronických dat ---- 17
2.2.2	Recepce a produkce textu ---- 19
2.2.3	Reprezentativnost korpusových dat jako zrcadlo reality ---- 21
2.3	Korpusový kontext ---- 25
2.4	Korpus a empirie ---- 27
2.5	Technické a další aspekty přípravy korpusu ---- 28
2.5.1	Příprava dat pro korpus ---- 29
2.5.2	Dotazy a dotazovací jazyk ---- 34
3	KORPUS A KORPUSY ---- 37
3.1	Korpus, jeho povaha a vytváření ---- 37
3.1.1	Korpus a jazyk ---- 37
3.1.2	Tvorba korpusu a jeho možnosti ---- 39
3.1.3	Hledání v korpusu a jeho analýza ---- 41
3.1.3.1	Možnosti hledání a analýzy ---- 41
3.1.3.2	Souhrnný pohled na úzus slova v kontextu: Konkordance ---- 51
3.1.3.3	Kombinace textové i systémové: kolokace a koligace ---- 55
3.1.3.4	Postup od formy k významu a funkci. Typický a zvláštní ---- 62
3.1.4	Hlavní aspekty a charakteristiky korpusu ---- 67
3.2	Český národní korpus (ČNK) ---- 69
3.3	Typy korpusů ---- 74
3.3.1	Synchronní korpus ---- 76
3.3.2	Mluvený korpus ---- 78
3.3.3	Diachronní korpus ---- 80
3.3.4	InterCorp ---- 82
3.3.5	Webové korpusy ---- 83
3.4	Hlavní světové korpusy ---- 84

4	KORPUSOVÁ LINGVISTIKA ---- 90
4.1	Korpusová lingvistika jako disciplína ---- 90
4.2	Korpusová metodologie ---- 92
4.3	Korpusová statistika ---- 100
4.3.1	Základní pojmy a operace ---- 101
4.3.2	Frekvence ---- 102
4.3.3	Zipfův zákon ---- 106
5	KORPUS VE VÝZKUMU A APLIKACI ---- 109
5.1	Výzkum spojený s budováním korpusu ---- 110
5.2	Lingvistický výzkum ---- 111
5.2.1	Výzkum obecně ---- 111
5.2.2	Variabilita jazyka. Mluvená povaha jazyka ---- 112
5.2.3	Lexikon a frazeologie. Kolokace ---- 114
5.2.4	Morfologie a gramatika. pedagogické aplikace ---- 117
5.3	Lingvistický kontrastivní výzkum více jazyků (InterCorp) ---- 118
6	ZÁVĚR: PERSPEKTIVY KORPUSOVÉ LINGVISTIKY ---- 120
II. ČÁST: SONDY, STUDIE, ANALÝZY ---- 123	
A	Korpusy včera, dnes a zítra (2011) ---- 125
B	Some of Current Problems of Corpus and Computational Linguistics or Fifteen Commandments and General Truths (2007) ---- 145
C	Spoken Corpora Design: Their Constitutive Parametres (2009) ---- 156
D	InterCorp: jeho povaha a možnosti (2014) ---- 165
E	Lexical Collocability: The Case of Verbs and Adverbs (2010) ---- 182
F	Abstract Nouns Collocations: Their Nature in a Parallel English-Czech Corpus (2005) ---- 196
G	Statistical Methods for Searching Idioms in Text Corpora (2006) ---- 207
H	Povaha a úzus interjekcí: případ češtiny (2007) ---- 216
I	Konference: Staré známé nebo neznámé slovo? (2001) ---- 225
III. ČÁST: PŘÍLOHY ---- 229	
Příloha A	Jazyková analýza výběrové konkordance (linout se) ---- 231
Příloha B	Frekvence slov ve Frekvenčním slovníku (prvních 100 lexémů) ---- 234
Příloha C	Jeden pohled na českou morfologii (Procentuální rozložení rodu, čísla a pádu u substantiv) ---- 237
Příloha D	Román jako korpus: Zbabělci (Josef Škvorecký) ---- 239
KORPUSOVÁ BIBLIOGRAFIE ---- 247	
GLOSÁŘ ZÁKLADNÍCH POJMŮ A TERMÍNŮ ---- 261	

I. ČÁST: TEORIE

1 ÚVOD: KORPUS A KORPUSOVÁ LINGVISTIKA

Tato kniha je stručným shrnutím a přehledem základních pojmů, zásad a metod **korpusové lingvistiky**, tj. nauky o korpusu a práci s ním, i nového stejnojmenného oboru. Už podle svého jména se tu staví nutně na **korpusu**, tj. velkém elektronickém, systematicky budovaném a organizovaném úhrnu textů, který je *conditio sine qua non* pro jeho disciplínu, tj. korpusovou lingvistiku. Kniha zároveň přináší i malou ilustraci možností výzkumu s modelovými výsledky.

Korpus se objevil v lingvistice, spolu s nástupem počítačů, poměrně nedávno a výzkum jazyka v ní v důsledku existence korpusu zcela převrátil poměry a etablované postupy, především z hlediska rozsahu informací a možností, které data i počítače nabízejí. Dalo by se tu, jakkoliv se to běžně nedělá, mluvit o korpusové revoluci v lingvistice, která jde mj. ruku v ruce s rozvojem počítačů, aspoň na svém počátku.

Jestliže lingvista tradičně donedávna strádal zásadním nedostatkem dat, a tedy, metaforicky řečeno, jistou informační podvýživou, pak se s nástupem korpusů situace radikálně změnila. Tam, kde v celé dlouhé minulosti oboru lingvista pracně, dlouho a nedokonale shromažďoval svá manuální excerpta v podobě lístečků a výpisků, než je mohl začít třídit a dávat do své kartotéky (a později i většího archivu), kterými se obv. pak sekvenčně probíral, tam nastoupil počítač a jeho možnosti. Na místě někdejších lístečků a výpisků dnes korpus poskytuje doslova záplavu informací a dat, ve kterých je sice často obtížné se vyznat přímo a snadno, ve kterých se ale dá už hledat více způsobů, vždy pomocí počítače. To vedlo a vede mimo jiné k nutnosti nalezení (a dalšímu hledání) potřebných technik a metodologie, jak k tomuto množství dat smysluplně přistupovat.

Korpusové poznání přináší v mnoha ohledech překvapivá zjištění. Tato kniha výběrově přibližuje možnosti mnohostranného poznání všeho základního, co se událo za cca dvacet a více let v této nové disciplíně, **korpusové lingvistice**, jejíž

výsledky a aplikace začínají i u nás už otřásat tradiční důvěrou ve spolehlivost stávajících lingvistických prací, starých slovníků, mluvnic a dalších příruček. Mimo jiné se korpusová lingvistika postupně stala už i novým lingvistickým oborem studovaným na univerzitách.

Knihy postupuje od širšího rámce, potřeb moderní lingvistiky, povahy korpusu a jeho možností až k některým základním a konkrétním autorovým výsledkům z jeho studia a výzkumu v daném oboru. Svým systematickým důrazem na vyčerpávající přístup k datům a jejich popisu, který v lingvistice přináší a důsledně uplatňuje poprvé korpusová lingvistika, se tato kniha snaží korpusová data a přístupy k nim aspoň základním způsobem charakterizovat a narýsovat tak zde zároveň i půdorys této rychle se rozvíjející disciplíny. Stejně tak se snaží ukázat na průnik i návaznosti na lingvistické disciplíny tradiční a ovšem přitom i naznačit otevřené otázky a problémy. Jakkoliv se hlavní zkušenost opírá o práci s budováním a analýzou **Českého národního korpusu** (viz blíže na adrese korpus.cz), v pozadí stojí i zkušenost s většinou korpusů větších a řadou menších, především evropských, a kontakty s týmy, které je budují, která se prezentuje výběrově také.

Tato stručná oborová příručka a malý sborník zároveň se pokoušejí komplexně zachytit tedy jak povahu korpusových dat, tak způsoby jejich zpracování, možnosti práce s nimi i všechny hlavní souvislosti a návazné aspekty tak, aby podala stručný komplexní a ucelený, avšak zároveň i základní obraz **oboru** korpusová lingvistika (**I. ČÁST: TEORIE**).

Doprovázejí ji a jen volně na ni navazují reprinty vybraných autorových studií a analýz představující některé výsledky analýzy korpusu (**II. ČÁST: SONDY, STUDIE, ANALÝZY**), které na konci doplňuje ilustrativním způsobem několik konkrétně orientovaných dodatků (**III. ČÁST: PŘÍLOHY**).

Svým pojetím je kniha určena jak studentům a lingvistům či odborníkům z jiných oborů, především humanitních, tak i širší zainteresované veřejnosti. Kniha je vybavena dvojí bibliografií, obecnou, vztahující se k celé problematice disciplíny, a specifickou, vztaženou přímo k tématům v podobě stručných odkazů za jednotlivými kapitolami. Bibliografie tak do značné míry vyvažuje i specifikuje stručnost formulací a popisů zde obsažených a orientuje uživatele na další možnosti studia. Příspěvky v části II mají původní bibliografii vlastní.

ZÁKLADNÍ LITERATURA

Čermák 1995a, 1998, 2001b; McEnery – Wilson 1996; McEnery – Hardie 2012; Wynne 2005.

Speciální a další literaturu viz v Korpusové bibliografii

2 DATA A INFORMACE

Informace se chápe různě a ve více smyslech, intuitivně však především jako nějaký údaj, zjištění o něčem, resp. už i přímo zobecněné poznání či vědění. Informace se odborně někdy vymezuje jako to (tedy zjištění, poznatek), co relevantního lze získat z podkladových dat (údajů) pro analýzu jevu či problému, a to obvykle v nějakém rámci, pro lingvistiku pak z dat textových, a to zvláště v rámci a kontextu sociálním. Ten se v případě jazyka chápe především jako rámec sémiotický a je řízený pragmatickými pravidly.

Současný informační věk a jeho povahu obecně i specificky reflektuje i korpus: obojí, tj. náš věk i korpus sám, nabízí nebyvalou záplavu informace. Její množství, které je zásadně vítané, však i hrozí zahlcením uživatele, a představuje proto také relativně nový problém. Pravda je, že náš informační věk se už dobře neobejde bez počítačů a globálního webu, na tyto nástroje a **hromadné** zdroje informace spoléháme dnes zcela samozřejmě. Této novodobé pravdě ovšem předcházela a předchází pravda starší a stále zcela základní, totiž že **valná většina informace** (významová, formální i pragmatická) **je kódovaná a přenášená (přirozeným) jazykem** a je jen málo oborů, které se při výměně informací bez jazyka obejdou (srov. část matematiky, popř. i malbu apod.); na to se často zapomíná.

Tyto informace, často po určitém zpracování a formalizaci, se obvykle předávají dál, dalším uživatelům, zájemcům. Pak se role toho, kdo informaci objevil, získal (zvl. odborník, lingvista), mění a stává se z něj ten, kdo ji předává, sděluje v nějaké podobě jiným lidem. Tyto informace, tj. poznatky zvláště o aktuální povaze a stavu toho, co nebo kdo nás zajímá, se dnes výrazně sdělují elektronicky, a to často **hromadně**.

Stále se však dosud a hlavně sdělují i po staru také **individuálně**, mj. ústním a soukromým kontaktem mezi lidmi, osobním pozorováním, zkoumáním, poznáním něčeho, popř. i zčásti zobecněnou a zprostředkovanou školní zkušeností, a tedy vlastní deduktivní či induktivní činností. Taková informace, např. článek, kde lingvista zformuluje své poznání a zjištění, či ho odpřednáší studentům, se tak dostává dál. Je však třeba lišit, zda takové poznání autor považuje, po ověření jeho platnosti, tj. zda skutečně odpovídá datům a obvyklým přijímaným normám, popř. i stupni poznání, za zobecnitelné, za obecně přijatelné a považuje ho tedy za poznání a informaci **objektivní**. Proti poznání objektivnímu jde však často i o poznání, zkušenost i informaci **subjektivní**, vázanou na jednotlivce apod., kde míra zobecnitelnosti nebývá jasná.

Oba zdroje poznání a informace z něj odvozené, zdroj individuální i hromadný, jsou ovšem komplementární a jeden nevylučuje druhý. Nicméně individuální a často subjektivní zkušenost a informace nabývaná jedincem obvykle nemívá

možnost **potvrzení** své správnosti ani obecnosti, a tedy zjištěné **platnosti** při opakovaném výskytu, tj. ověřované opakovaně, a použitelné tedy i ostatními. Taková povaha informace se naopak očekává, ba vyžaduje u informací veřejných, obecných, a zvláště vědeckých.

Data, informační jednotky a základní stavební kameny poznání, jsou podle oboru a druhu vnímání a zapojeného smyslu pochopitelně různá, (téměř) všechna jsou však převoditelná do slov, a tedy jazyka, psaného a mluveného, který nám zprostředkuje zrak a sluch. A až tady se zúročuje možnost, zprostředkovaná komputery, dostávat se k zobecnění, popř. zobecnitelné, a tedy i zpravidla **objektivnější a opakovatelné**, resp. opakované informaci, tj. takové, kterou lze nalézt až ve větších informačních jazykových souborech, resp. textech, tj. kterou v nich lze ověřit a kterou individuální a jednorázová zkušenost a poznání z vlastního výzkumu v zásadě nenabízí.

Toto **zobecnění** jednotlivého *faktu*, tj. poznání, ať je induktivní či deduktivní, a jeho *platnosti* je ovšem možné jen v **kontextu** a na základě či proti pozadí široké znalosti dané vzděláním obecným či oborovým; izolovaný fakt nám naproti tomu neřekne (téměř) nic.

Data a informace v nich však ještě nevedou k tomu, co je zvláště v případě seriózního studia na jeho konci to nejdůležitější, to jest k **věděni**, popř. znalosti, jak informace používat a využívat. To je však už věcí integrace poznatků do platného kontextu a jejich nazírání v něm, ale i studia a učení, opřené o opakování a postaveného kvůli souvislostem do vhodného širšího rámce, zvláště celého jazyka i disciplíny.

2.1 ELEKTRONICKÝ TEXT

2.1.1 OBECNÁ POVAHA, VLASTNOSTI A ZPRACOVÁNÍ TEXTŮ

Základní a nejčastější podoba informace je obsažená v *textu*, dnes většinou už elektronickém, tj. v textu vytvořeném pomocí počítače (a jím i dále zpracovávaném). Každý text je smysluplný řetězec znaků, pro lingvistiku a uživatele jazyka řetězec hlavně slov, který má lineární povahu. Psaný text má ale i povahu lineárně prostorovou, kde platí dimenze *vlevo-vpravo* (podle druhu písma i jiná), nebo lineárně časovou, kde platí v podstatě dimenze *napřed-potom* (u mluveného jazyka se jeho povaha lineárně časová přepisem mění na lineárně prostorovou taky). Paralelní jiná informace (vedle základní psané), zvláště intonace u mluvené komunikace (suprasegmentální rysy), se v písmu nezachycuje, a jen mluvený text je tudíž vlastně vícelineární, skládá se z více souběžných informačních linií (více viz 2.3); mluvený text je však také obecně primární a historicky starší.

Elektronický text je pro potřeby korpusové lingvistiky text, který je počítačově čitelný a dále zpracovatelný. V lingvistice se za text považuje libovolný souvislý text, získaný obvykle jen z autentických zdrojů, v praxi konkrétně (smysluplný) text psaný nebo mluvený. Míra a způsob zpracování a přípravy elektronického textu závisí na cíli, technických podmínkách a potřebách. Zpracování textu psaného a mluveného se liší. Podoby psaného textu se však liší také.

Obecně je při zpracování textů třeba lišit tři základní pojmy, jejichž míra obecnosti se postupně zužuje, a to nepřímou úměrně k nárůstu jejich specifčnosti: typ dat, kód dat a formát (obv. souboru). **Typ dat** je pojem nejširší a patří sem jen okrajově. Tvoří ho proměnlivý *způsob ukládání dat různého druhu*, dat vizuálních, obrázků (např. *jpeg, gif*) či akustických (např. *wav, wma* či kompresní *mp3*). Pro korpusy jsou však zásadní druhé dva pojmy, kód (kódování) dat a jejich formát (obv. formát souboru).

Texty se zachycují, resp. znázorňují v různých **kódech**, které v elektronické oblasti označují *způsoby jejich převodu (a informace v nich) do způsobu jiného*; patří sem např. morseovka (písmena převáděná na tečky a čárky), (dříve) dírky a jejich uspořádání na děrném štítku, a dnes konečně a hlavně čísla (na která se převádějí písmena). V této poslední oblasti (užívající číselného kódu) jsou dnes běžné dva typy kódů a kódování, ASCII a Unicode. ASCII, které prodělalo více proměn a postupné rozšiřování, dnes čítá jen pro češtinu až pět různých verzí, a protože obecně stále není plně použitelný pro všechny jazyky, přechází se proto, resp. přešlo se do univerzálního *Unicode* (viz víc v 2.5.1 a 2.1.1).

Naproti tomu **formát** je *způsob technického záznamu a zpracování informace v určitém kódu technických dat*, který je značně proměnlivý. Formátů může být pro uložení téže informace víc než jeden (vedle netextových formátů, viz výše např. *mpeg*), specificky mj. *doc, T602, ODF pro Linux, html*, či obecný *XML*, i podle druhů operačního systému, pro Windows, Unix či Mac OS). V tomto smyslu tedy formát (resp. způsob kódování znaků, viz dále) je *takové uspořádání dat, které slouží k jejich znázorňování, ukládání a znovuužívání*. Známý je mj. formát *pdf* (portable document format) užívaný pro různé typy textů jako snadno převoditelný, určený zvláště pro tisk (s informací o charakteristikách takového textu, ale i obrázky), který je známý z běžné práce s elektronickými texty.

Elektronické texty, které dnes vstupují jako jeho stavební kameny do korpusu, se tedy vyskytují v různých podobách, *formátech*, a pro potřeby počítače se musejí konvertovat do formátu jednotného, tj. sjednotit. Vedle formátů spojených s operačním systémem (*Windows, Unix/Linux*, dříve *DOS*) se texty vyskytují ve formátu, který jim dávají různé textové editory a procesory (jako *Word, Writer* v *LibreOffice* či *TextEdit* v *Mac OS X* aj.), existují ve formátech vznikajících i ve webovém prostředí (zvl. *html, Hypertext Markup Language*) a ty je třeba podrobovat *konverzi* (viz 2.5.1).

V praxi se po konverzi a převodu ze specifického formátu výsledný formát označuje obvykle jako *txt* (prostý text). Zachycuje ho rozšířený kód ASCII (pův. *American Standard Code for Information Interchange*), obvykle už v univerzálním kódování *Unicode* (viz 2.5.1), specificky v UTF-8. Tak se tvoří základ pro další korpusové zpracování a budování vlastního korpusu (srov. dál ještě následnou úpravu textu jazykem XML, zahrnujícím i další zpracování, 2.5.1).

Vybraný korpusový text, získaný různými konverzemi textů z textových procesorů užívajících různých formátů, je ve výsledku zásadně text *autentický*, a má tu podobu, kterou mu dal autor, resp. vydavatel. Je tedy nepředstavitelné, že korpusový text, který by měl ideálně zachovávat záměr autora, by se měl opravovat, dál editovat, pospisovňovat či jinak pravopisně, nebo dokonce cenzorně, a tedy krajně sporně „vylepšovat“ (což je věcí konkrétní a vždy problematické jazykové politiky a inženýrství). Ochranu autora a autentičnost jeho textu je třeba dodržet za každou cenu (jakkoliv i autorská formální úprava textu může být nevyrovnaná a naznačovat problémy tvůrce textu s jazykem). Předpokládá se tedy, že korpusový text zachovává i různé individuální překlepy a omyly autora, což pro masové korpusové hledání nevádí. Takové, v tradičním pohledu „spisovníků“ a tradicionalistů, nenáležitosti bývají totiž řídké a neovlivňují zásadně celkový výsledek hledání v korpusu (viz 2.5.1). Korpusy mají tedy mj. i „konzervační“, archivní roli tím, že uchovávají (relativně) trvale texty v neměnné podobě (jakkoliv ale přitom z hlediska typografického už není jejich původní podoba zpětně rekonstruovatelná). Při přípravě textů pro korpus se však vyskytují i (často hromadné) chyby textu vzniklé technickým zpracováním či sazbou textu (a vzniklé tedy bez, anebo proti přání autora), a ty je naopak třeba zjistit a řadou procedur opravit, jakkoliv rozlišení obou typů nebývá snadné.

Protože zásadní podobou korpusových dat je stále text psaný, spadají sem nutně i *elektronické přepisy* mluvených textů (viz 3.3.2), které mají také textovou podobu. Pro zachování jednotnosti s texty psanými mají být s nimi ideálně srovnatelné, avšak hledání jejich optimální formy (v zásadě formy transkripce) při zachování jejich věrné původní podoby je zatím otevřená otázka (zvl. s ohledem např. na nespisovné či individuální podoby slov). Takové texty se můžou dodatečně ještě vybavovat i fonetickým či prozodickým přepisem (zachycujícím jejich suprasegmentální rysy).

Proti dřívější koncepci budovat korpus z homogenních (náhodných) **vzorků** textů (ty vyhovují zvl. statistickým potřebám někdy lépe) se dnes dává, mj. i vzhledem k dostupnosti elektronických textů, přednost tomu, do korpusu začleňovat **texty celé**, které umožňují optimální zkoumání kontextu, stejně tak ale i povahu odlišných částí téhož textu (jako je specifický začátek apod., při textové analýze). Problém, kdy se může přístup založený na vzorcích textů jevit výhodnější, se např. projevuje tehdy, když jde o malou textovou kategorii či (pod)obor, jehož naplnění

více různými malými vzorky je lepší a vhodnější než jediným a úplným textem větším. Ani velké množství náhodných vzorků však nemusí zachytit některé zvláštnosti textu.

2.2 KORPUSOVÁ DATA

Protože většina korpusů je veřejných a zájemcům nekomerčně přístupných, bývá získávání elektronického textu (*akvizice*) od poskytovatelů spojeno s nemalými problémy autorských práv, **copyrightu**, který je podmíněn souhlasem různých agentur (zastupujících autory, tj. vlastníky textů) udělujících souhlas k využití textů na základě zvláštních smluv. Někdy je dokonce jejich využívání korpusu i zpoplatněno (stranou zůstávají soukromé či podnikové korpusy, veřejnosti nedostupné, vyvíjené zvl. v zahraničí například některými nakladatelstvími, kde může být autorský souhlas delegován na tato nakladatelství apod.). Ať už jde přímo o nakladatelství, vydavatelství či texty zprostředkující organizace a agentury, všechny jsou vždy odpovědné autorům za to, že se dodržuje autorský zákon v případě textů, se kterými pracují, i za jejich ochranu, tj. zajišťuje dodržování licenčních podmínek, za nichž byly poskytnuty; neautorizované šíření textů je právně postižitelné. Korpusová centra přebírají tyto texty sice z různých zdrojů, ale vždy za podmínky respektování domluvených podmínek, především ve smyslu ochrany před dalším šířením, a to na základě právně závazných smluv.

Případně širší, tj. zvláště technické a tedy nelingvistické využití autorských textů, které tvoří korpus, je věcí dohody takového zájemce a poskytovatele textů (zastupujícího autora textů). Avšak i standardní **využití** korpusových textů, které jsou pro korpusové užívání provozovateli korpusu, tedy pro dané využití, jen svěřené (v českém případě je to zvl. Ústav Českého národního korpusu Filozofické fakulty Univerzity Karlovy), je omezené na dohodnutou limitovanou délku citace (zvl. v podobě rozsahu, pozic či vět ap.). Takové užívání je v zásadě otevřené a u konkrétního lingvisty apod. vázané na respektování podmínek užívání. Pro odborný technický výzkum textů se po dohodě může zájemcům nabídnout i *podoba textů s namátkově přeházenými větami* či jinými úseky (*reshuffled texts*), kde už nejde přímo o autorský text v původní podobě, a práva autorů se tak v zásadě nenarušují.

Ekvivalentem, obdobou copyrightu je u mluvených korpusů *dohoda* s nahrávanými osobami o podmínkách zveřejnění nahrávek, kde právně namísto autorského zákona začíná platit občanský zákoník. Nemalý problém, narážející na nestejnou legislativu autorských práv v různých zemích, se musí řešit při výstavbě paralelních (vícejazyčných) korpusů. Více o vlastnictví textů viz též v 2.5.

2.2.1 POVAHA A ZDROJE ELEKTRONICKÝCH DAT

Elektronická data se zpracovávají z více druhů zdrojů z dodaných, do počítače vnesených textových dat (viz 2.5.1). Dnes je u psaných textů *synchronních* nejčastější, nejsnadnější a nejužitečnější ta podoba elektronická, která se získává, často přes nějakého zprostředkovatele (nakladatelství, agenturu ap.), a je už produkovaná přímo v elektronické podobě (většina autorů dnes sama píše na počítači v nějakém textovém editoru).

Naproti tomu se *starší* texty musejí zpravidla ručně a pracně skenovat (programy OCR, optical character recognition), zatímco *mluvené* texty se nejčastěji do počítače napřed a neméně pracně přepisují z akustického záznamu, resp. nahrávky a až pak dál zpracovávají jako jiné elektronické texty.

Korpusová data, prostá či (většinou) různě anotacemi obohacená a technicky zpracovaná, jsou, co do zdroje informace, zhruba řečeno moderním elektronickým ekvivalentem starých a rozsahem omezených manuálních excerpt z různých (většinou psaných) zdrojů (v podobě kartotéčních lístků). Z nich se získávala a někdy ještě i získává prakticky veškerá informace o jazyce, tj. generalizací z více podobných či stejných výskytů (lecko si dokonce i dnes dělá malé kartotéky dosud sám). Korpusová data jsou, jinými slovy a zjednodušeně řečeno, rozsáhlé soubory elektronických textů, které jsou následně uspořádány a propojeny do podoby korpusu. Tato data jsou na rozdíl od soukromých osobních, omezených, a tedy nutně subjektivních excerpt jednotlivce zásadně *objektivní*, protože ve velkém a zvláště reprezentativním korpusu odrážejí v širokém průřezu nejružnější a typologizované druhy textů, a tedy jak většinový a *typický úzus* jazykových forem, tak ale i úzus menšinový (daný nižší frekvencí jejich výskytu), který zrcadlí jejich různé a různé časté zastoupení v korpusu. **Reprezentativní podoba** dat v korpusu se obvykle považuje v nespécifickém a globálním výzkumu jazyka usilujícím o celkový a vyvážený obraz úzu za optimální. Zprůměrování získaných výsledků naznačuje pak i typický úzus.

V kontrastu k tomuto většinovému přístupu korpusové lingvistiky usilujícímu o objektivnost informací a podkladových dat stojí dosud některé přístupy např. N. Chomského a jeho stoupenců. Ten objektivitou informace, kterou korpusy nabízejí, v zásadě pohrdá a upřednostňuje své vlastní izolované a často i spekulativní či konstruované příklady úzu, na kterých zakládá své argumenty a proměnlivé teorie. Vyvozovat zobecnění z jednotlivých a izolovaných příkladů, které takový přístup nabízí, je však krajně problematické a stojí v protikladu k cílům korpusové lingvistiky analyzovat text důsledně celý a systematickým způsobem, bez vynechávek.

Každý korpus je dnes založený na souborech elektronických **dat psaných**, připravených pro tisk nebo internetovou komunikaci různého druhu (jiná psaná data

je třeba skenovat, viz už výše) či **dat mluvených**. Psané texty se dnes získávají většinou (i když ne vždy) od vydavatelů tištěných textů na základě dohody s nimi a kvůli jejich velikému rozsahu jsou také automaticky zpracovávány, i když mívají různé formáty; v řadě případů je třeba při získávání dat pro korpus autorská práva tvůrců či majitelů textů chránit (copyright, viz i výše).

Náročnější je proces získávání **mluvených dat** pro mluvené korpusy (viz 3.3.2); představuje ho především nahrávání zvláště spontánních promluv. Takové nahrávky se následně přepisují do grafické podoby tak, aby mohly případně existovat paralelně s akustickou podobou. Jiné, **vícemodální korpusy** (zahrnující i obrazy, digitální, někdy i filmové, ve snaze zachytit dynamiku zachycované situace promluvy) pro náročnost jejich získání i zpracování dosud ve významnějším rozsahu neexistují, a jsou proto stále víceméně zbožným přáním. Problém je nejen v metodologii záznamu (propojení obrazu, zvuku a písma), způsobu hledání v nich, ale i takových věcech, jako je způsob anotace gest, posunků, výrazů tváře, typů situace, nálady či postojů mluvčích ap., tedy všech **mimotextových faktorů**, které se spolupodílejí na významu a funkci; žádná shoda napříč jazyky a týmy (pokud se věci vůbec zabývají) není. O takovém korpusu, srovnatelném svou mírou informativnosti a objektivnosti s korpusem psaným, nemůže být proto zatím ani řeč, není tu dostatečné množství slov ani kontextů. Přestože existuje pár pokusných vícemodálních korpusů, které implicitně ovšem v praxi zatím vlastně znamenají *pouze dva zaznamenané mody* podle lidských smyslů uplatňujících se při vnímání nejčastěji, tj. zraku a sluchu (např. film představuje jejich spojení), lze se pro budoucnost ptát na možnost zapojení i smyslů dalších, zvl. čichu a hmatu. Nejde tedy tolik o mody (a název multimodální není vhodný), jako o druhy percepce (a mělo by se event. mluvit o korpusech multiperceptuálních). Nikdo si však otázku, kolik módů komunikace zachycovat a jakým způsobem, vlastně dosud neklade, zřejmě proto, že neví, jak by na ni odpověděl.

Od korpusových dat, korpusů se liší **archivy**, původně manuální, dnes však už většinou také elektronické, které bývají veřejné (jako *Oxford Text Archive*), nebo vázané na lexikografická, popř. jiná centra (jako *Archiv ÚJČ ČSAV*), která svá data shromažďovala už dříve, a to ruční excerpce (na kartotéční lístky); jejich význam už dnešní korpusy v podstatě odsunuly do pozadí.

Archivy, konzervující na rozdíl od korpusů jakákoliv cenná elektronická data, ať už diachronní nebo mluvená, jsou však spíše výjimkou. V druhém případě, tj. u mluvených dat, jde však mnohem víc o výjimku, protože automatický spolehlivý převod akustických dat (z existujících archivů) do psané podoby dosud k dispozici není. V prvním případě, i když je jejich povaha také manuální, archivy se stále udržují a někdy i rostou dál (jako *archiv Oxfordského slovníku* sestavovaný přes 150 let na základě jeho *Reading programme* s pomocí původně dobrovolníků výběrově excerpujících dodané knihy). Starší data se však i zde musejí převádět

do elektronické podoby (viz 3.3.3) většinou náročným ručním skenováním a procházejí pak dalšími procesy.

Významným typem aktivity, obecně sdíleným v mnoha zemích, je ve snaze o zachování národního kulturního dědictví úsilí o naskenování národní literatury minulé i nové v knihovnách, popř. i dalších textů, které jsou (jednoduchým způsobem) pak na webu veřejně a obecně zpřístupňované. Tak postupuje pro českou oblast např. Národní knihovna. Často se však takový archiv skládá, pro snadnější způsob jeho vytvoření (protože kvalita automatické digitalizace starších textů není uspokojivá), jen z „obrázků“, tj. stránek textu naskenovaných jako celek (a tedy obrázek), ve kterých slova obvyklým způsobem hledat nejde, a pro korpusové využití taková data nepostačují.

K neznámějším obecným a velkým mezinárodním zdrojům textů patří výsledky činnosti několika organizací, především *Project Gutenberg* (Open Language Archives Community), *Elsevier.org* či *OTA* (Oxford Text Archive) či u nás *Národní knihovna* (srov. <http://www.nkp.cz/digitalni-knihovna>, zvl. skrze aplikaci *Kramerius*, <http://kramerius.nkp.cz/kramerius/>). Některé další jsou komerční jako *ELRA*, *ELDA* či americká *LDC* (adresy viz v bibliografii).

2.2.2 RECEPCE A PRODUKCE TEXTU

V autorství psaných textů se projevuje jasná disproporce: je mnohem víc pasivních uživatelů (čtenářů, konzumentů) než těch aktivních (autorů, tvůrců), tj. těch, co texty vytvářejí, píšou. Ideální stará představa bývala, že korpus by měl být ve vyvážené proporcii složený z textů odpovídajících (širokému publiku a tedy) recepci i produkci, tj. mělo by být zohledněno i množství autorů textů; tak se měl specificky odrážet i jejich soukromý jazyk (ne nutně však veřejný). V korpusové lingvistice (i jinde) nejde tedy pouze o uživatele v první řadě, nýbrž i o autory, tvůrce (korpusových) textů. Tento rozdíl se nazývá **recepce textu** a stojí proti **produkcí textu**; je to specifický problém především textů psaných. U mluveného jazyka a mluvených textů tento rozdíl do té míry nevyvstává, všichni mluví i poslouchají zároveň, jakkoliv ani zde ne ve zcela stejné míře (ale to už je věc situace a psychologicky i povahy mluvčích; je tu i nejednoznačný vliv rozhlasu). Původní představa, že by u psaného jazyka korpus měl nějak odrážet proporce produkce a recepce, byla zkusmo uplatněná v zásadě jen u dánského korpusu, kvůli chybějícímu výzkumu a dostupnosti reprezentativního výběru textů podle produkce se však dnes (zatím) neuplatňuje.

Pokud chtějí mít reprezentativní podobu, budují se korpusy v praxi tedy především podle kritéria **recepce**, tj. rozsahu a složení čtenářstva (psaných) textů a povahy jimi čtených textů. To je různě velké, největší u novin, nejmenších zřejmě u odborných tisků či tisků administrativních a podle obsahových hledisek je i proměnlivé. Tyto

texty se dají relativně snadno získat, jejich proporce promítající se do reprezentativního složení korpusu jsou však stále věcí diskuze. Proti textům z hlediska recepce stojí texty získávané podle **produkce**, tj. texty ideálně pocházející od (co nejvíce) různých autorů a posuzované i podle množství těchto autorů; tento ohled tvůrci korpusů při výběru textů sledují jen málo a takové texty bývají v menšině a je obtížné je získávat ve vyvážených proporcích.

Obecně jde však také a z jistého hlediska především o *formování, ovlivňování čtenáře* a jeho jazyka skrze texty. Při bližším pohledu je zřejmé, že na jedné straně je relativně velmi málo lidí (autorů, spisovatelů, novinářů), kteří svou tvorbou zásadně a *aktivně* ovlivňují svou produkcí (psaný) jazyk a jeho konzumaci čtenáři (jakkoliv významným a blíže nezkoumaným způsobem). Proti tomu však stojí většinový úzus platný pro obrovskou většinu čtenářů na straně druhé, kteří jsou v roli pasivních čtenářů textů, a tedy i uživatelů korpusu. Platí to v podstatě i obecněji, srov. např. čtenáře bulvárních novin (s nejvyššími náklady) vytvářených malou skupinou novinářů pro obrovskou masu čtenářů (přitom erudice a jazyková kompetence těchto novinářů nemusejí být nejlepší).

Určité omezené možnosti vyvážení této jinak přirozené disproporce (zvl. v dialogu) představují *speciální korpusy*, jako je např. korpus korespondence (viz 3.3), zaměřené jen na jeden speciální druh textu, jehož uživatelé jsou odpovídajícím způsobem specializovaní také; zpravidla to jsou lingvisti a literáti.

Moderní korpusy, zvláště ve své publicistické složce, však nemusejí být zcela jednojazyčné, i když jednojazyčnost je předpoklad samozřejmý. V důsledku reklamy (někdy i rozsáhlejší), popř. **cizojazyčných** citátů (kombinace slov i celé věty) se v novinách apod. objevují i menší texty cizojazyčné, zvl. z jazyků velkých a známých, tj. angličtiny, němčiny, latiny apod. (někdy celé odstavce). Tyto nečeské texty do textu českého bývají ovšem vloženy (zvl. z různých důvodů nakladatelem či vydavatelem) a text jejich dodatečnou přítomností pak ztrácí svou autentičnost. Je pochopitelné, že pro národní jednojazyčný korpus takovéto texty nejsou žádoucí a musejí se odstraňovat. Rysem této cizojazyčné složky je malá či žádná opakovanost a (i pro český korpus) jistá bezkontextovost (resp. slabá, popř. odlišná závislost cizích slov na kontextu jiného jazyka, např. němčiny na českém kontextu). To se ovšem netýká *citátových cizojazyčných slov a výrazů*, ať jde o módní anglicismy či klasické latinské či řecké výrazy, srov. *à propos, acquis communautaire, last but not least* aj., které se do českého úzu svou vysokou frekvencí už zařazují a užívají se i v širším českém jazykovém kontextu. Okrajovým případem ovšem je český text, který autor cizojazyčným textem vybavil sám; ty je pak třeba, v zájmu zachování autentičnosti, v korpusu ponechat.