



SLOVENSKO-ANGLE[KI KORPUS *ELAN*

Ve-jezi-ni korpusi so pomemben vir podatkov prevodoslovnim raziskavam in tehnolo-gijam strojnega in strojno podprtega prevajanja. ^lanek predstavi dvojezi-ni korpus, ki smo ga na IJS zbrali v okviru projekta EU ELAN. Korpus vsebuje milijon besed, sestavlja pa ga 15 sodobnih terminolo{ko bogatih besedil in njihovih prevodov v slovenskem in an-gle{kem jeziku. Besedila so stav-no poravnana ter ozna-ena v skladu s predpisi TEI (Guidelines for Text Encoding and Interchange). Korpus in vsako besedilo v njem je ozna-eno z glavo TEI, poravnana dvobesedila pa so shranjena podobno kot pomnilniki prevodov. Celoten korpus je dostopen na mre' nem naslovu <http://nl.ijs.si/elan/>, kjer so na voljo tudi dvojezi-ne konkordance korpusa.

Uvod

Pomen ve-jezi-nih korpusov se ka'e skozi veliko {tevilu projektov, ki se jih trudijo zagotoviti: med evropskimi so zglede MLCC (Armstrong et al. 1998) (devet jezikov EU), Crater (McEney et al. 1997) ({pansko-francosko-angle{ko) in ENPC (Johansson et al. 1996) (angle{ko-norve{ko). Za slovenski jezik je bil doslej edini ve-jezi-ni korpus objav-ljen na CD-ROM-u v okviru evropskega zdru'enja TELRI (Erjavec et al. 1998); vzporedni korpus na tej zgo{-enki vsebuje Platonovo Republiko, delo In{tituta za slovenski jezik pri ZRC SAZU, in korpus projekta MULTTEXT-East (Erjavec, Ide 1998). ^eprav CD-ROM TELRI ponuja veliko {tevilu jezikov, je zapisan v skladu z mednarodnimi priporo-ili in bogato ozna-en, pa vsebuje le malo slovenskih besedil (pribl. 500.000 besed iz {tirih be-sedil), ki hkrati niso najbolj primerna za terminolo{ke raziskave, torej za podro-je, na katerem so vzporedni korpusi {e najhitreje uporabni.

S projektom EU ELAN (the European Language Activity Network) se je ponudila prilo'nost, da se izboljša ponudba korpusov za slovenski jezik. V okviru projekta smo zbrali, ozna-ili in uredili korpus, ki vsebuje milijon besed v 32.000 prevodnih enotah, od katerih vsaka vsebuje segment (ve-inoma poved) v slovenskem in angle{kem jeziku. Skladno z nameni projekta je korpus prosto dostopen tako za prepis na ra-unalnik kot za iskanje. ^lanek predstavi korpus IJS-ELAN in ima naslednjo zgradbo: poglavje 2 opi(e na-in zbiranja in obdelave besedil, poglavje 3 na{teje 15 besedil, zbranih v korpusu, po-glavje 4 poda ra-unalni{ki zapis (oznake) korpusa, poglavje 5 mo'nosti dostopa do korpusa in poglavje 6 zaklju-ke.

Zagotovitev besedil in obdelava

Ker smo imeli za izdelavo korpusa IJS-ELAN na voljo samo pol leta, je bil eden od osnovnih ciljev zbrati ~im ve-ji korpus s tem, da zaobidemo najbolj zamudne korake v izdelavi korpusa. Proces izdelave (urejanja), ki je podrobneje opisan v (Erjavec 1999b), je pri korpusu IJS-ELAN zajemal: pridobitev in za{-ito avtorskih pravic do izvirnih besedil in korpusom kot celoto, zagotovitev digitalnih izvirkov besedil, segmentacijo in poravnavo povedi, raz-lenitev na besede (tokenizacija), zapis v standardizirani format; zapis glav kor-pusa in besedil in zagotovitev dostopnosti na svetovnem spletu (WWW). Klju-nega pomena za uspe{nost projekta so bili na-in pridobivanja besedil, izdelave poravnanih dvobesedil ter njihov ra-unalni{ki zapis.

Sodelavci projekta so besedila zbirali sami ter jih nato tudi pretvorili iz izvirnega zapisa, jih segmentirali in poravnali. Ker je korpus namenjen ~im bolj svobodnemu raz{ir-



janju, smo izbrali besedila, kjer avtorske pravice niso problematične, bodisi ker besedila sodijo med javne vladne dokumente bodisi med dokumentacijo programja GNU, kjer je nadaljnje razširjanje celo zaželeno, za nekatera besedila pa so bile urejene že prej v okviru drugih projektov. Reprodukcijski besedil v celoti na osnovi korpusa (kar je po navadi glavna ovira nadaljnjemu razširjanju) pa je ostala z našim zapisom korpusa.

Našim pretvorbe in poravnave se je med sodelavci razlikoval; uporabili smo orodja UNIX in pomnilnik prevodov Deja Vu podjetja Atril. V obeh primerih smo kot rezultat dobili poravnana dvobesedila, ki so v precejšnji meri ostala izvornih oznak (npr. HTML, RTF) in zapisana v enostavnem tabelarnem formatu, po ena prevodna enota (tj. poved v izvorniku in njen prevod) v vsaki vrstici. Čeprav prevodne enote v splošnem ustrezajo povedim, so včasih tudi daljše ali pa krajše, saj se zgodi, da eni povedi v izvorniku ustrezata dve ali več v prevodu ali obratno.

Tako dobljena dvobesedila smo nato ostilili s filtri, napisanimi v programskem jeziku Perl, s katerimi smo normalizirali zapis nabora znakov in odstranili ostanke formatiranja. Naslednji korak je bila tokenizacija, tj. identifikacija besed in slovničnih, za kar smo uporabili orodje MULTTEXT 'mtseg' (Di Cristo 1996) s pravili MULTTEXT-East za slovenski in angleški jezik (Dimitrova et al. 1998). Tokenizacija označi tudi besedila, okrajšave itd. Tudi ta korak zaradi nepopolnih orodij in pravil v korpusu prinese napake, ki so spet pretežno popravljene s filtri Perl. Tokenizirana poravnana dvobesedila smo nato pretvorili v standardizirano obliko TEI, kjer se jim je tudi dodal opis v obliki glav TEI. Zadnji korak je bil pakiranje korpusa za prenos preko mreže in ter pretvorba iz standardne oblike v takšno, ki je primerna za WWW in služi kot vir informacij o korpusu, ter v zapis za konkordančni, ki je ravno tako priklopljen na WWW. Tu smo uporabili prosto dostopen program Omnimark Lite, ki zna ravnati z dokumenti, zapisanimi v skladu s standardom SGML.

Zvrstnost korpusa

[tevilno, velikost in vrste besedil v korpusu je kompromis med (pripravljeno) uporabnostjo in enostavnostjo zagotovitve pravic in digitalnega originala. K uporabnosti prispeva predvsem dejstvo, da gre za sodobna besedila (90. leta), ki so terminološko bogata in z zanimivih in dinamičnih področij. Zaradi enostavnosti smo vključili besedila, ki nimajo posebnih omejitev nad nadaljnjim razširjanjem in so nam bila dostopna v digitalni obliki na WWW. S Službo Vlade RS za evropske zadeve, kjer so nam prepustili večjo količino še ne objavljenih besedil, pa smo podpisali posebno pogodbo.

Korpus je sestavljen iz petnajstih enot, ki so vključena integralna dvobesedila, vendar brez pretežno nebesedilnih delov (npr. tabel s slikami). Vsako dvobesedilo ima pripisan identifikacijski niz in skupaj s svojo glavo predstavlja element korpusa. Dvobesedila so razdeljena na tista s slovenskim izvornikom in angleški prevodom, in tista z angleškim izvornikom in prevodom v slovenski jezik. Poleg razlike po smeri prevoda imata ta dva dela tudi precej različno zvrstnost.

Polovica s slovenskim izvornikom je vključena vladnega izvora, sestavlja pa jo enajst enot. Te enote skupaj s svojim identifikatorjem, približno velikostjo v kB in številom besed v tisočih, letnico izida, kratkim naslovom ter zalogom so naslednje:

usta: 364 Kb, 20 kW, 1997

Ustava Republike Slovenije; Ustavno sodišče Republike Slovenije

kuca: 1102 Kb, 69 kW, 1990-95

Govori predsednika RS, M. Kučana; Urad predsednika Republike Slovenije



parl: 325 Kb, 20 kW, 1998

Delovanje Dr'avnega zbora; Dr'avni zbor Republike Slovenije

ecmr: 4056 Kb, 239 kW, 1998/1999

Ekonomsko ogledalo; 13 {tevil; Urad Republike Slovenije za makroekonomske analize in razvoj

ekol: 1222 Kb, 70 kW, 1999

Nacionalni program varstva okolja; Republika Slovenija, Ministrstvo za okolje in prostor, Uprava RS za varstvo narave

spor: 589 Kb, 34 kW, 1996

Evropski sporazum; Slu'ba Vlade RS za evropske zadeve

anx2: 483 Kb, 25 kW, 1996

Evropski sporazum – Priloga II; Slu'ba Vlade RS za evropske zadeve

stra: 1511 Kb, 89 kW, 1997

Strategija Slovenije za vklju-evanje v EU; Slu'ba Vlade RS za evropske zadeve

kmet: 543 Kb, 29 kW

Dr'avni program za prilagajanje zakonodaje – kmetijstvo; Slu'ba Vlade RS za evropske zadeve

ekon: 394 Kb, 23 kW

Dr'avni program za prilagajanje zakonodaje – gospodarstvo; Slu'ba Vlade RS za evropske zadeve

vade: 471 Kb, 24 kW, 1995

Vademecum Lekove doma-e lekarn; Lek d.d.; OTC Division

^eprav del z angle{kim izvornikom vsebuje skoraj polovico besed v korpusu, je sestavljen iz samo {tirih enot, od katerih sta dve knjigi. Vsebuje tudi druga-ne besedilne vrste in podro-ja kot prva polovica: dve enoti se ukvarjata z ra-unalni{tvom, ena pa z vizijo totalitarne dru'be:

vino: 1182 Kb, 69 kW, 1994

EC Council Regulation No 3290/94 – agriculture / Uredba sveta ES {t. 3290/94 – kmetijstvo; Slu'ba Vlade RS za evropske zadeve

ligs: 3044 Kb, 173 kW, 1999

Linux Installation and Getting Started / Namestitev in za-etek dela z Linuxom; Linux Documentation Project (-en: Specialized Systems Consultants; -sl: Linux User Group of Slovenia, LUGOS)

gnpo: 353 Kb, 13 kW, 1999

GNU PO localisation files / GNU PO lokalizacije datoteke; Free Software Foundation, Linux Documentation Project

orwl: 6698 Kb, 195 kW, 1983

G. Orwell: Nineteen Eighty-Four / 1984; projekt MULTEXT-East, slovenski prevod: knji'nica Kondor, Mladinska knjiga (prevajalka: Alenka Puhar).



Kot je razvidno iz seznama, vsebuje korpus raznovrstna, ter večinoma terminološko bogata besedila, primerna za avtomatsko identifikacijo terminov in njihovih prevodov. Za etne raziskave (Vintar, 1999) z uporabo pretežno statističnih metod kažejo, da je korpus močno uporabiti v terminološke namene.

Zapis korpusa

Korpus je zapisan v skladu standardom SGML (Standard Generalized Markup Language, ISO 8879) in uporablja definicijo tipa dokumentov, ki je parametrizacija priporočila TEI (Sperberg-McQueen, Burnard 1994). Čeprav vsebujejo priporočila TEI tudi predlog za zapis vzporednih poravnanih korpusov, se nam ta niso zdela primerna za naš korpus. Namesto tega smo parametrizirali TEI tako, da je zapis bolj podoben tistemu, ki se uporablja pri pomnilnikih prevodov (Erjavec 1999a).

V našem zapisu uporabljamo generične elemente TEI za zapis glav ter za oznake znotraj segmentov (povedi), medtem ko zapis besedil spremenimo tako, da ta neposredno vsebujejo prevodne enote, tj. so poravnana dvobesedila. Nabor znakov v korpusu je definiran opisno, z mehanizmom entitet SGML. Tako je npr. ~ v korpusu zapisan kot č,] kot ´, & pa kot &, definicija tipa korpusa pa vsebuje nabor in opis uporabljenih entitet. Celoten korpus je sestavljen iz glave korpusa, ki vsebuje informacije o korpusu kot celoti in iz petnajstih elementov korpusa; vsak od teh spet vsebuje glavo in telo (dvobesedilo).

Glava TEI vsebuje podatke o datoteki, o viru ali virih besedila, o zapisu besedila in seznam sprememb. Čeprav je jezik projekta, v okviru katerega je korpus nastajal, angleški, smo poleg angleščine v glavah uporabljali tudi slovenski jezik. Za vtis o tem, kakšne podatke vsebujejo glave v korpusu, podamo nekaj primerov. Prvi je začetek korpusa in njegove glave:

```
<tei Corpus. 2>
<tei header type="corpus" lang="sl" id="ijs-elan" creator="et"
status="update" date.created="1999-04-14" date.updated="1999-06-22">
<filedesc>
<titlestmt>
<title lang="en">The IJS-ELAN Slovene/English Parallel Corpus</title>
<title lang="sl">Slovenskoangleščanski vzporedni korpus IJS-ELAN</title>
```

Naslednji primer iz glave korpusa poda deklaracijo oznak, ki se uporabljajo v korpusu:

```
<tagsdecl>
<tagusage gi="text" occurs="15">Element 'Text'.
Attributes are LANG and ID.</tagusage>
<tagusage gi="body" occurs="15">Element 'Body'.
Content model: TU+</tagusage>
<tagusage gi="tu" occurs="31900">Element 'Translation unit'.
Attributes are LANG and ID.</tagusage>
<tagusage gi="seg" occurs="63800">Element 'Translation segment'.
Attributes are LANG.</tagusage>
<tagusage gi="s" occurs="13386">Element 'Sentence'; only in 'orwl' text.
Attributes are ID (value identical to original MTE bundle.</tagusage>
<tagusage gi="w" occurs="1091745">Element 'Word'.
Attributes are TYPE (IMPLIED/971313, dig/26818, abbr/2662, comp/179)
and (only in 'orwl' text) LEMMA, FUNCTION.</tagusage>
<tagusage gi="c" occurs="167243">Element 'Punctuation'.
Attributes are TYPE (IMPLIED/131019, open/18115, close/18109).</tagusage>
</tagsdecl>
```

Izjava o odgovornosti iz ene od glav besedil:

```
<respstmt>
<name>Jasna Belc, SVEZ</name>
<resp lang="sl">Zagotovitev digitalnega originala</resp>
<resp lang="en">Provision of digital original</resp>
<name>&Scaron;pel a Vintar, FF</name>
<resp lang="sl">Poravnava</resp>
```


- <w type=comp> Ve-besedna leksikalna enota, npr. *medtem ko, vice versa, New York*
- <w type=di g> Beseda, ki vsebuje {tevilke, npr. *1984, 3., IV, 20%, 1993-1996, 25/76, 16MB*
- <w type=abbr> Kratica (beseda, ki se konča s piko), npr. *tar., et al., S.u.S.E., dipl.*
- <w> privzeti tip za "navadne" besede, npr. Slovenije, market, 's, Article, 'ivnorejo, INAVGURACIJSKI, Hurt-Andreatta, Hrup51, E-po{tni, D'you

^eprav je korpus označen po TEI berljiv s poljubnim urejevalnikom besedil, uporabniku ni prav prijazen. Ena od odlik zapisa SGML naj bi bila enostavna pretvorba v format, ki je primeren konkretni aplikaciji. Za bolj{o preglednost smo zato naredili pretvorbo iz zapisa TEI v zapis HTML, tako da so vse glave korpusa ter primeri besedil na voljo preko standardnih mre{nih brkljalnikov.

Dostopnost

Domača stran IJS-ELAN ima naslov <http://nl.ijs.si/elan/>. Poleg informacij o korpusu (npr. glave TEI v zapisu HTML) je tam dostopen tudi celoten korpus v zapisu TEI. Glave korpusa in posameznih delov določajo, da je korpus prosto dostopen pod pogojem, da se citira njegove vire, dokumentirane v glavi.

Za uporabo besedil označenih po TEI je potrebno vsaj osnovno znanje programiranja in dostop do orodij za izkori{anje korpusa. Da bi olaj{ali uporabo korpusa, omogočamo tudi mre{ne konkordance nad korpusom. Konkordance so eden osnovnih načinov uporabe korpusov, ki omogočajo iskanje in izpis besed in besednih zvez skupaj s sobesedilom. Na IJS imamo licenco za uporabo konkordančnika CQP (Christ 1994), ki je hiter, ima bogat iskalni jezik ter podpira prikaz poravnanih segmentov. Na Univerzi v Gothenburgu so razvili spletni vmesnik za CQP, ki smo ga nato izboljšali in priredili za delo s slovenskim naborem znakov ter za prikaz poravnanih segmentov. Korpus IJS-ELAN smo pretvorili v format, ki ga zahteva CQP in ga vključili v nabor preko mre{e dostopnih korpusov.

S konkordančnikom, ki ga najdemo na <http://nl2.ijs.si/corpus/index-bi.html>, je mogoče iskati po angleškem ali slovenskem delu korpusa, pri čemer lahko dodatno podamo omejitve glede na poravnani segment. Tako lahko npr. iz{emo vse pojavitve besede 'drevo', kjer se v prevedenem segmentu ne pojavi beseda 'tree'. Iskalni izrazi so lahko enostavni (npr. 'corpus' ali 'besed*') ali pa uporabljajo polno (in razmeroma kompleksno) sintakso CQP. Slednja dopu{ča poljubne regularne izraze nad nizi, iskanje po med seboj oddaljenih delih segmentov ter iskanje po tipih besed oz. ločil. Mre{ni vmesnik ponuja tri načine izpisovanja zadetkov: privzeti način je dvojezični, kjer zadetku s sobesedilom sledi tudi poravnani segment v drugem jeziku; enojezični izpis dobimo v načinu KWIC (key-word in context) ali pa kot frekvenčni seznam zadetkov.

Zaključki

V članku smo predstavili slovensko-angleški vzporedni korpus IJS-ELAN. Korpus je uporaben kot vir prevodov terminov, predvsem skozi možnost mre{nih konkordanc, pa tudi kot podatkovna zbirka za raziskave in razvoj jezikovnih tehnologij, predvsem tistih vezanih na terminologijo in prevajanje. Predstavljeni korpus je prosto dostopen v upanju, da dodatno spodbudi razvoj korpusnega jezikoslovja za slovenski jezik. ^eprav je za slovenski jezik 'e na voljo oziroma v delu nekaj korpusov, npr. korpus FIDA (Krek et al.), so le ti enojezični; IJS-ELAN tako predstavlja prvi večji dvojezični korpus za naš jezik.

Nadaljnje delo s korpusom bo usmerjeno predvsem v zagotovitev bogatejšega nabora oznak. Tu je na prvem mestu lematizacija in oblikoslovno označevanje besed v korpusu. V



daljši perspektivi pa bi bilo seveda najbolj koristno povežati količino besedil v korpusu, kar pa bi bilo možno samo v okviru širšega projekta.

Zahvale

Pri delu na korpusu, predstavljenem v tem članku, so sodelovali Roman Maurer, Andrej Skubic in Jelena Vintar. Besedila za korpus so prispevali uradi in službe Republike Slovenije, posebej služba Vlade RS za evropske zadeve. Besedila so tudi prispevali Linux Users Group of Slovenia, LUGOS in Lek d.d., OTC Division. Delo na projektu je deloma financirala pogodba z Institut fuer deutsche Sprache v okviru projekta MLIS-ELAN 121 in pogodba MZT L2-0461-0106.

LITERATURA

- Susan ARMSTRONG, Masja KEMPEN, David MCKELVIE, Dominic PETITPIERRE, Reinhardt RAPP, Henry THOMPSON, 1998: Multilingual corpora for cooperation. *Proceedings of the First International Conference on Language Resources and Evaluation*. LREC'98. Granada: ELRA. 579–980.
- Oliver CHRIST, 1994: A Modular and Flexible Architecture for an Integrated Corpus Query System. *Proceedings of COMPLEX '94: 3rd conference on Computational Lexicography and Text Research*. Budimpešta. URL: <http://www.ims.uni-stuttgart.de/projekte/CorpusWorkbench/>
- Philippe di CRISTO, 1996: *Mtseg: The multext multilingual segmenter tools*. MULTEXT Deliverable MSG 1. Version 1.3.1. CNRS. Aix-en-Provence. URL: <http://www.lpl.univ-aix.fr/projects/multext/MtSeg/>
- Ludmila DIMITROVA, Toma' ERJAVEC, Nancy IDE, Heiki-Jan KAALEP, Vladimir PETKEVIČ, Dan TUFIS, 1998: Multext-East: Parallel and Comparable Corpora and Lexicons for Six Central and Eastern European Languages. *COLING-ACL '98*. Montreal, Quebec. 315–319.
- Toma' ERJAVEC, 1999a: A TEI encoding of aligned corpora as translation memories. *Proceedings of the EACL-99 Workshop on Linguistically Interpreted Corpora (LINC-99)*. Bergen: ACL.
- – 1999b: Making the ELAN Slovene/English Corpus. *Proceedings of the Workshop on Language Technologies – Multilingual Aspects*. Ljubljana: Univerza v Ljubljani. 23–30.
- Toma' ERJAVEC, Nancy IDE, 1998: The MULTTEXT-East corpus. *First International Conference on Language Resources and Evaluation, LREC'98*. Granada: ELRA. 971–974.
- Toma' ERJAVEC, Ann LAWSON, Laurent ROMARY (ur.), 1998. *East meets West: A Compendium of Multilingual Resources*. CD-ROM, TELRI Association e.V. URL: <http://www.ids-mannheim.de/telri/cdrom.html>
- Stig JOHANSSON, Jarle EBELING, Knut HOFLAND, 1996. Coding and aligning the English-Norwegian parallel corpus. *Languages in Contrast*. Ur. K. Aijmer, B. Altenberg, M. Johansson. Lund: Lund University Press. 87–112. URL: <http://www.hit.uib.no/enpc/>
- Simon KREK, Marko STABEJ, Vojko GORJANC, Toma' ERJAVEC, Miro ROMIČ, Peter HOLOZAN. *FIDA: korpus slovenskega jezika*. URL: <http://www.fida.net>



- Tony McENERY, Andrew WILSON, Fernando SANCHEZ-LEON, Amalio NIETO-SERRANO, 1997: Multilingual Resources in European Languages: Contributions of the CRATER Project. *Literary and Linguistic Computing* 12/4.
- C. M. SPERBERG-MCQUEEN, Lou BURNARD (ur.), 1994: *Guidelines for Electronic Text Encoding and Interchange*. Chicago, Oxford. URL: <http://www-tei.uic.edu/orgs/tei/>
- [pela VINTAR, 1999: A Lexical Analysis of the ELAN Slovene-English Corpus. *Proceedings of the Workshop on Language Technologies – Multilingual Aspects*. Ljubljana: Univerza v Ljubljani. 63–70.

Tomaž Erjavec
Institut Jožef Stefan, Ljubljana
Odsek za inteligentne sisteme