

## KORPUS GOS IN NJEGOVA UPORABA V RAZISKOVALNE, DIDAKTIČNE IN LJUBITELJSKE NAMENE

Ana Zwitter Vitez

Trojina, Zavod za uporabno slovenistiko, Ljubljana

UDK 808=163.6:81'42:004.738.52

Korpus GOS je zbirka avtentičnih posnetkov govorjene slovenščine, prosto dostopna na spletnem naslovu [www.korpus-gos.si](http://www.korpus-gos.si). Gre za interdisciplinarni rezultat projekta Sporazumevanje v slovenskem jeziku, ki sloni na analizi diskurza, sociolinguističnih in medijskih raziskavah, programiranju in procesiranju govornih signalov. Namenjen je poučevanju slovenščine, jezikoslovnim in jezikovno-tehnološkim raziskavam, poklicnim govorcem in piscem besedil, pa tudi, nenazadnje, uporabnikom, ki doslej o jeziku sploh niso razmišljali.

govorni korpus, uravnoveženost korpusov, transkribiranje posnetkov, uporaba korpusov

The Spoken Corpus of Slovene GOS is a collection of authentic recordings including speakers from all over Slovenia and available online on the website [www.korpus-gos.net](http://www.korpus-gos.net). This achievement represents one of the results of the project Communication in Slovene and is based on the interdisciplinary research (sociolinguistic, media and discourse analysis, computer processing and sound engineering). The corpus GOS can be used for language teaching, linguistic analysis, language technologies, professional speakers and, nevertheless, for the wider audience.

spoken corpus, balanced corpus, speech transcription, use of language corpora

### 1 Uvod

Od oktobra 2010 je na naslovu [www.korpus-gos.net](http://www.korpus-gos.net) mogoče najti korpus govorjene slovenščine z imenom korpus GOS. Elektronska zbirka transkribiranih posnetkov, zgrajena v okviru projekta Sporazumevanje v slovenskem jeziku, zajema milijon besed oziroma 120 ur posnetkov. Pri tem smo želeli zajeti:

- demografsko reprezentativni vzorec govorcev,
- najpogosteje govorne situacije.

Korpus GOS je plod združenih moči sociolinguističnih in medijskih raziskav, analize diskurza, računalniškega programiranja in procesiranja govornih signalov. Kot edinstven vir avtentičnih primerov govorjenih besedil iz najrazličnejših komunikacijskih situacij je korpus namenjen spoznavanju

realne spontane govorne komunikacije v slovenščini, kar je pomembno za poučevanje slovenščine, poklicne govorce in snovalce besedil, jezikoslovne raziskave na področju sociolinguistike, pragmatike in analize diskurza ter jezikovnotehnoške raziskave za potrebe sinteze in avtomatskega razpoznavanja govora. Poleg omenjenih namenov uporabe ima korpus GOS še en nadvse pomemben cilj: zbuditi zanimanje uporabnika, ki sicer o jeziku sploh ne razmišlja, ter prispevati k njegovemu veselju in ustvarjalnosti pri vsakodnevnom govornem sporazumevanju.

V nadaljevanju prispevka bomo predstavili, kako so različna področja sodelovala pri gradnji korpusa GOS, in razmislili o tem, kako bodo cilji korpusa govorjene slovenščine dejansko uresničeni.

## 2 Zajem gradiva

Osnovni cilj zajema govorjenih besedil sta bili čim boljša uravnoveženost in reprezentativnost končne elektronske zbirke. Zato smo kombinirali demografske in besedilno-vrstne kriterije.

Demografske kriterije predstavljajo dejavniki, ki najverjetneje (vendar še ne dokazano) najbolj vplivajo na razlike v govoru: spol, starost, izobrazba in regionalna pripadnost govorcev (ne pa tudi socialni status, saj je za slovenske razmere težko določljiv).<sup>1</sup> Podatke za določitev demografskih kriterijev smo pridobili pri Statističnem uradu Republike Slovenije.<sup>2</sup>

Pri besedilovrstnih kriterijih smo upoštevali javnost diskurza in prenosnik. Javni diskurz je odprt za širšo javnost in naslavlja veliko skupino ljudi. Nejavni diskurz predstavljajo zasebni pogovori (v okviru družine, prijateljev) in nezasebni pogovori (v uradih, trgovinah in profesionalnem življenju). V javni diskurz smo vključili razvedrilne in informativno-izobraževalne medijske vsebine, pa tudi pedagoški diskurz v različnih izobraževalnih ustanovah.

Zajeli smo štiri različne prenosnike govorjenih besedil: radio in televizija (javni diskurz) ter osebni stik in telefon (pedagoški in nejavni diskurz).

## 3 Označevanje in transkribiranje posnetkov

Vsek posnetek je opremljen z naslednjimi podatki:

1. **podatki o govorcih** (spol, starost, izobrazba, regijska pripadnost, prvi jezik),
2. **podatki o posnetku** (vrsta diskurza, prenosnik, regija, leto snemanja),
3. **zapis govora** (transkripcija).

<sup>1</sup> Več o kriterijih zajema besedil je v Zemljarič Miklavčič idr. (2009).

<sup>2</sup> www.stat.si

<sup>3</sup> http://ola.zrc-sazu.si

### 3.1 Zapis govora

Pri pregledu obstoječih praks transkribiranja smo ugotovili, da pri nas ni poenotenega načina zapisovanja, vendar obstajajo tri ključne tendence:

- dialektologija večinoma sledi fonetični transkripciji z diakritičnimi znamenji,<sup>3</sup>
- pragmatične raziskave (Smolej 2006) sicer uporabljajo slovenski črkopis, vendar zapisujejo pojave moderne vokalne redukcije ter druge pogovorne in narečne prvine govorjene slovenščine (npr. *maš čevle*),
- v jezikovnotehnički praksi (Zemljarič Miklavčič 2007) se uporablja poknjiženi zapis govorjenega jezika, iz katerega pogovorne prvine niso več razvidne.

Pri transkribiranju posnetkov korpusa GOS smo želeli najboljše razmerje med hitrostjo postopka zapisovanja, berljivostjo končnih transkripcij in učinkovitostjo iskanja želene besedne oblike. Pri tem smo upoštevali, da je slovenščina močno dialektalno in socialnozvrstno razčlenjen jezik, ter želeli ohraniti čim več različnih oblik določene besede. Zato smo govor zapisali na dva različna načina: pri »pogovornem zapisu« skušamo kar se da zvesto ujeti dejansko akustično podobo izgovorjenih besed; pri »standardiziranem zapisu« dobi več različic neke besedne oblike (npr. *imam, mam, jemam*) krovno standardizirano obliko (npr. *imam*).

S pogovornim zapisom tako omogočimo dober vpogled v besedje in oblike govorjenega jezika, s standardiziranim zapisom omogočimo boljše iskalne možnosti, za povsem natančno zvočno podobo govora pa je vedno na razpolago tudi posnetek izjave.

Primer 1: Besedilo, zapisano s pogovornim (Pog.) in standardiziranim zapisom (Stan.)

Pog. a *lohk mal pošpegamo*

Stan. a *lahko malo pošpegamo*

## 4 Elektronska baza

V zadnji fazi gradnje korpusa GOS je velik del odgovornosti slonel na programiranju. Tekstovni del korpusa je tako zapisan v računalniškem standardu XML in v skladu s priporočili TEI (Text Encoding Initiative). Zahtevnejši uporabniki si ga lahko v tej obliki snamejo s spletnne strani [www.korpus-gos.net](http://www.korpus-gos.net). Opis sheme XML je vključen v paket za prenos korpusa oziroma dostopen na strani [http://nl.ijs.si/ssj/gos/schema/tei\\_gos\\_doc.pdf](http://nl.ijs.si/ssj/gos/schema/tei_gos_doc.pdf).

## 5 Primeri rabe korpusa GOS

Primeri rabe korpusa GOS bodo osredotočeni na tiste situacije, v katerih se govorni korupsi še ne uporabljajo: pouk slovenščine in študij govorjenih žanrov za poklicne govorce. Uporabe korpusa za potrebe jezikoslovnih raziskav in jezikovnih tehnologij ne bomo posebej predstavljalji. Predpostavljamo namreč, da jezikoslovci, programerji in elektroinženirji dobro vedo, kako uporabljati govorne korupse za potrebe sociolinguističnih, pragmatičnih, dialektoloških in

prevodoslovnih analiz ali pri sintezi in avtomatskem razpoznavanju govora.<sup>4</sup>

### 5.1 Učenje jezikov

Eden izmed pomembnejših ciljev je uporaba korpusa GOS pri pouku slovenščine kot neprvega jezika, saj »vnaprej napisani dialogi redko odsevajo nepredvidljivost in dinamiko konverzacije, zato imajo učenci manj možnosti, da bi usvojili ustrezni nabor jezikovnih sredstev, ki jih bodo potrebovali v neznanih situacijah izven učilnice« (Burns 2001: 21). Tudi I. Ferbežar (2003: 35–36) ugotavlja, da je jezikovna raba »odvisna od različnih situacijskih dejavnikov: prostora, časa, vloge (odnosa med udeležencema pogovora, njune starosti, statusa ipd.), teme, sporazumevalnega namena in seveda sobesedila«, vendar »obstoječi učbeniki za slovenščino kot neprvi jezik vse te dejavnike premalo upoštevajo«.

Pri pouku slovenščine za materne govorce je lahko korpus GOS koristen vir avtentičnih pogovorov pri usvajjanju razlik med formalnimi in neformalnimi situacijami. Predstavljajmo si, da petnajst minut po začetku ure v razred stopi učenec, ki svojo zamudo pojasni

The screenshot shows the 'Napredno iskanje' (Advanced search) page of the GOS corpus. It features two main search panels. The top panel is for the word 'moč' (power), with options for 'Oznaka pred' (Marker before) and 'Oznaka za' (Marker after). The bottom panel is for the word 'Glagol' (verb), also with 'Oznaka pred' and 'Oznaka za' options. Both panels have dropdown menus for 'Način iskanja' (Search mode) and 'Besedna vrsta' (Word type). Below the panels are sections for 'V okolici' (In the neighborhood) and 'Dodatek beseda v okolini' (Additional word in the neighborhood), both with similar search parameters. At the bottom right is a 'Najdi' (Find) button.

Slika 1

<sup>4</sup> Sintesa in avtomatsko razpoznavanje govora je do zdaj večinoma slonelo na idealiziranem gradivu dekontekstualiziranih izjav, posnetih v idealnem studijskem okolju in z (domnevno) idealnimi govorci.

## Simpozij OBDOBJA 30

The screenshot shows the GOS search interface. At the top, there are buttons for 'gos' (Search), 'Iskanje' (Search), and 'Seznam' (List). Below these are two radio buttons: 'Iskanje po pogovornem zapisu' (Search by conversational text) and 'Iskanje po standardiziranem zapisu' (Search by standardized text). A large button labeled 'biti' is highlighted, and next to it is a 'Najdi' (Find) button. Below these buttons are links for 'Uporabljš enostavno iskanje' (Use simple search) and 'Napredno iskanje' (Advanced search). The main search results area shows a list of items under 'Tip diskurza' (Type of discourse) and 'Kanal' (Channel). Under 'Tip diskurza', there are categories like 'Nejavni zasebni' (29.126), 'Javni informativno-izobraževalni' (28.756), 'Javni razvedrini' (18.885), 'Nejavni nezasebni' (12.198), and 'Več'. Under 'Kanal', there are categories like 'Osebni stik' (45.323), 'Televizija' (17.813), 'Radio' (17.183), and 'Telefon' (8.646), followed by a 'Več' link. To the right of the search results, there is a list of search results snippets. The first snippet reads: 'nasledil očeta ampak ko se je to zgodilo še vedno je simbolično šel sedet na ta knežji kamen in je sprejel'. The second snippet reads: 'vedno je simbolično šel sedet na ta knežji kamen in je sprejel pač eee voljo ljudstva da jim bo da ga'. The third snippet reads: 'kamen in je sprejel pač eee voljo ljudstva da jim bo da postavlja za kneza in on jim je potem'. The fourth snippet reads: 'jim bo da postavlja za kneza in on jim je potem objubil da dobro vladal zdej jasno klučno temeljno'. The fifth snippet reads: 'postavlja za kneza in on jim je potem objubil da dobro vladal zdej jasno klučno temeljno vprašanje v kakšnem jeziku'. The sixth snippet reads: 'je to potekalo a ne se na žalost ne v zarad'. The seventh snippet reads: 'dobro vladal zdej jasno klučno temeljno vprašanje v kakšnem jeziku'. The eighth snippet reads: 'a ne se na žalost ne v zarad tega ker ni ohranjenih zapisov zej mi z veseljem sklepamo in si želimo'. The ninth snippet reads: 'zapisov zej mi z veseljem sklepamo in si želimo da je bilo to v nekem slovenskem arhaičnem jeziku predvidoma ja ker'. The tenth snippet reads: 'zapisov zej mi z veseljem sklepamo in si želimo da je bilo to v nekem slovenskem arhaičnem jeziku predvidoma ja ker se'. At the bottom of the search results area, there are buttons for 'prejšnja stran' (Previous page), 'naslednja stran' (Next page), and a link to 'Prikazujem 251-300 od 88.965 konkordanc (0.248 sekund)'.

Slika 2

z izjavo *Oprostite, sem mogu k zobarju*. Učencu (in razredu) želimo pojasniti nepri-mernost uporabljenih struktur. Uporabimo funkcijo napredno iskanje, ki omogoča iska-nje po posameznih oblikah sestavljenih struk-ture. V naprednjem iskanju pod prvo iskano besedo vpišemo *moči* in pri oznaki besedna vrsta označimo, da gre za glagol. V podrob-nostih označimo, naj bo glagol nezanikan (ker je problematična samo struktura v trdilni obliki). Nato kliknemo na možnost beseda v okolini in označimo, naj prvi besedi sledi glagol v nedoločniku (Slika 1).

Tako iskanje bo vrnilo rezultate tipa *kolkol draže bi mogle biti, združenje bi moglo reč mi ne bomo nič delali, ga je po mojem moglo ful ful bolet*. Hiter pregled situacij, v katerih se ta struktura pojavlja najpogosteje, nam pove, da gre za zasebni diskurz. Glede na dejstvo, da je šola ustanova formalnega diskurza, je lahko to za učenca dober dokaz, da bi bila v njegovem primeru ustreznejša raba strukture *morati + nedoločnik*.

## 5.2 Poklicni govorci in snovalci besedil

Predstavljammo si situacijo, v kateri želi igralec ali pisec literarnega dela usvojiti dis-kurz določenega profila govorca za potrebe svojega lika. V tem primeru bo zagotovo ko-ristno filtriranje besedil glede na tip situacije in profil udeležencev. Smiselno je izbrati zelo pogosto besedo, ki nastopa v vsaki govorni situaciji, na primer vse oblike glagola *biti*. Izberemo enostavno iskanje po standar-diziranem zapisu, vpišemo *biti* in vidimo, da je rezultatov res veliko (okrog 90.000; Slika 2).

Levo od rezultatov so filtri diskurzov in govorcev. Pri tipu diskurza izberemo zasebni diskurz, pri regiji snemanja Dolenjska, pri podatkih o govorcu pa izberemo moške med 35. in 59. letom s končano srednješolsko izobrazbo. Dobimo približno 65 izjav izbranega profila govorca v izbrani govorni situaciji, s katerimi si lahko pomagajo poklicni govorci in snovalci besedil (igralci, scenaristi, novi-narji, prevajalci).

## 6 Prvi odzivi na korpus GOS

Prve odzive uporabnikov med učitelji slovenščine,<sup>5</sup> med raziskovalci s področja jezikovnih tehnologij<sup>6</sup> ter med bodočimi

<sup>5</sup> Predavanje na Filozofski fakulteti Univerze v Mariboru, predavanje na Filozofski Fakulteti Univerze v Ljubljani, konferenca Sirikt 2010.

poklicnimi govorci in snovalci besedil<sup>7</sup> lahko strnemo v naslednje komentarje in vprašanja:

- komaj čakam, da korpus GOS uporabim pri svojem delu,
- kdaj bodo na voljo obsežnejše delavnice za delo s korpusom GOS,
- kdaj bo zbirka nadgrajena z dodatnim naborom besedil (da bo po njej na primer relevantno tudi iskanje kolokacij).

## 7 Zaključek

Korpus govorjene slovenščine kot edinstvena elektronska zbirka avtentičnih govorjenih besedil v najrazličnejših govornih situacijah predstavlja jezikovni vir, ki uporabnika postavi v dejavno vlogo raziskovalca izgovorjave, besedišča in struktur govorjene slovenščine v domačem in zasebnem okolju, pa tudi govorjene slovenščine, ki jo zahtevajo bolj formalni govorni položaji.

Glede na prve odzive strokovne javnosti lahko zaključimo, da se je korpus GOS v slovenskem raziskovalno-strokovnem prostoru že prijel. Zdi pa se smiselno, da ostanemo previdni pri pričakovanjih glede moči demokratizacije tega novega dosežka jezikovnih tehnologij. V okviru programa mWomen<sup>8</sup> so med muslimanske ženske razdelili na tisoče mobilnih telefonov, da bi jim zagotovili bolj demokratičen status (več stikov z vrstnicami ter boljše možnosti za zaposlitev in izobraževanje), vendar so mobilni telefoni končali v rokah njihovih soprogov, ki jim jih posodijo le, ko morajo na pot, da jih lahko bolje nadzorujejo.

V primeru korpusa GOS seveda zlorabe ne morejo biti tako radikalne, vendar bi bilo žalostno, če bi tak jezikovni vir ostal zgolj v rabi znanstvenoraziskovalne elite. Nam bo s korpusom GOS uspelo zbuditi zanimanje širšega kroga uporabnikov? Bo korpusu GOS uspelo zmanjšati frustracije maternih govorcev slovenščine glede njihovih jezikov-

nih zmožnosti in zbuditi veselje pri njihovem vsakodnevnom govornem sporazumevanju? Če bomo nekega dne dosegli ta cilj – morda tudi s pomočjo težko priborjenega prostega dostopa na spletu in uporabniku prijaznega vmesnika – bo to velika reč.

## Literatura

- BURNS, Anne, 2001: *I see what you mean: Using Spoken Discourse in the Classroom*. Sydney: National Centre for English Language Teaching and Research.
- FERBEŽAR, Ina, 2003: Dialog med ciljem in metodo. *Jezik in slovstvo* 48/1. 31–44
- KRNEL, Dušan, 2008: Uporaba informacijsko-komunikacijske tehnologije (IKT) pri pouku v nižjih razredih osnovne šole. *Naravoslovna solnica* 13/1. 6–9.
- SMOLEJ, Mojca, 2006. *Vpliv besedilne vrste na uresničitev skladenjskih struktur: Primer narativnih besedil v vsakdanjem spontanem govoru. Doktorska disertacija*. Ljubljana: Filozofska fakulteta.
- STABEJ, Marko, VITEZ, Primož, 2000: KGB (korpus govorjenih besedil) v slovenščini. *Informacijska družba IS'2000: Jezikovne tehnologije*. Ljubljana: Institut Jožef Stefan.
- VERDONIK, Darinka, 2007: *Jezikovni elementi spontanosti v govoru*. Maribor: Slavistično društvo Maribor.
- ZEMLJARIČ MIKLAVČIČ, Jana, 2009: *Govorni korpsi*. Ljubljana: Znanstvena založba Filozofske fakultete.
- ZEMLJARIČ MIKLAVČIČ, Jana idr., 2009: Kaj in zakaj v referenčni govorni korpus slovenščine. Marko Stabej (ur.): *Infrastruktura slovenščine in slovenistike. Obdobja* 28. Ljubljana: Znanstvena založba Filozofske fakultete.
- ZWITTER VITEZ, Ana 2011: À la recherche de la parole perdue: l'exemple du corpus national du slovène parlé. Noelle Serpollet (ur.): *Colloque Transcrire, Ecrire, Formaliser II*. Orléans: Presses universitaires d'Orléans.
- ZWITTER VITEZ, Ana idr., 2009: Načela transkribiranja in označevanja posnetkov v refe-

<sup>6</sup> Konferanca Jezikovne tehnologije na Institutu Jožef Stefan 2010, konferanca Transcrire, Ecrire, Formaliser v Orleansu 2011.

<sup>8</sup> www.mwomen.org

ZWITTER VITEZ, Ana idr., 2009: Načela transkribiranja in označevanja posnetkov v referenčnem govornem korpusu slovenščine. Marko Stabej (ur.): *Infrastruktura slovenščine in slovenistike. Obdobja 28.* Ljubljana: Znanstvena založba Filozofske fakultete.

[www.korpus-gos.net](http://www.korpus-gos.net)  
[www.stat.si](http://www.stat.si)  
<http://ola.zrc-sazu.si>