

[www.slovenščina.eu](http://www.slovenščina.eu)  
**sporazumevanje**



**Nataša Logar Berginc, Miha Grčar, Marko Brakus,  
Tomaž Erjavec, Špela Arhar Holdt in Simon Krek**  
Korpusi slovenskega jezika Gigafida, KRES, ccGigafida  
in ccKRES: gradnja, vsebina, uporaba

Zbirka *Sporazumevanje*  
Urednik zbirke *Simon Krek*

Recenzentki *Irena Stramljič Breznik, Darja Fišer*  
Prevod povzetka *Iztok Kosem*  
Oblikovna zasnova zbirke *Tomato Košir*  
Prelom *Roman Ražman*  
Naslovnica *Tomato Košir*  
Avtor črkovne vrste »BadNews« *Samo Ačko*

Založili *Trojina, zavod za uporabno slovenistiko*  
in *Fakulteta za družbene vede*  
Za založnika *Iztok Kosem* in *Hermina Krajnc*

Ljubljana, 2012

Naklada: 100 izvodov  
Prva izdaja, prvi natis

Tisk *Collegium Graphicum, d. o. o.*  
Natisnjeno na papirjih *Hello* in *Vega*

© Nataša Logar Berginc, Miha Grčar, Marko Brakus,  
Tomaž Erjavec, Špela Arhar Holdt in Simon Krek

Vse pravice pridržane. Brez pisnega dovoljenja lastnika avtorskih pravic je prepovedano reproduciranje, distribuiranje, dajanje v najem, javna priobčitev, dajanje na voljo javnosti, predelava ali vsaka druga uporaba tega avtorskega dela ali njegovih delov v kakršnemkoli obsegu ali postopku, vključno s fotokopiranjem, tiskanjem ali shranitvijo v elektronski obliki.



CIP – Kataložni zapis o publikaciji  
Narodna in univerzitetna knjižnica, Ljubljana

811.163.6'322

KORPUSI slovenskega jezika Gigafida, KRES, ccGigafida in ccKRES : gradnja, vsebina, uporaba / Nataša Logar Berginc ... [et al.] ; [prevod povzetka Izток Kosem]. - 1. izd., 1. natis. - Ljubljana : Trojina, zavod za uporabno slovenistiko : Fakulteta za družbene vede, 2012. - (Zbirka Sporazumevanje)

ISBN 978-961-92983-6-7 (Trojina)  
ISBN 978-961-235-596-8 (Fakulteta za družbene vede)

1. Logar Berginc, Nataša  
262711040

*Korpusi slovenskega jezika  
Gigafida, KRES, ccGigafida  
in ccKRES: gradnja, vsebina,  
uporaba*

avtorji:

Nataša Logar Berginc

Miha Grčar

Marko Brakus

Tomaž Erjavec

Špela Arhar Holdt

Simon Krek



# Kazalo vsebine

13	<b>1</b>	<b>Zbiranje besedil in vsebina korpusa Gigafida</b>
13	1.1	Uvod
13	1.1.1	Cilj
14	1.1.2	Namen
14	1.2	Merila gradnje
15	1.2.1	Standard za zbiranje gradiva
18	1.2.2	Taksonomija
19	1.2.2.1	Tisk in internet
19	1.2.2.2	Knjižnost, periodičnost in drugo
19	1.2.2.2.1	Leposlovje in stvarna besedila
20	1.2.2.2.2	Časopisi in revije
21	1.3	Zbiranje besedil
21	1.3.1	Podatki za zbiranje
21	1.3.1.1	Nacionalna raziskava branosti
22	1.3.1.2	Izposoja v knjižnicah
23	1.3.1.3	Knjižne nagrade
24	1.3.1.4	Naklada
24	1.3.1.5	Spletne strani: obiskanost, uglednost
24	1.3.1.6	AJPES: izdajatelji knjig; udeleženci knjižnega sejma
25	1.3.1.7	Besedilodajalci in besedila pri FidiPLUS
25	1.3.2	Evidence besedil in besedilodajalcev, stik z besedilodajalci
26	1.3.3	Pogodba z besedilodajalci
27	1.4	Priprava besedil za vključitev
27	1.5	Označitev
29	1.6	Kolofon korpusnih dokumentov: <i>Vrsta besedila in Vir</i>
29	1.6.1	<i>Vrsta besedila in Vir</i> : internet
30	1.6.2	<i>Vir</i> : založba oz. naslov besedila
31	1.6.3	<i>Vir</i> : RTV Slovenija, Državni zbor Republike Slovenije
31	1.7	Vsebina korpusa
31	1.7.1	Taksonomija, čas in besedilodajalci
31	1.7.1.1	Obseg in delež besed po taksonomiji
34	1.7.1.2	Število besed po letih
36	1.7.1.3	Avtorji in založbe
36	1.7.2	Uspešnost zbiranja
37	1.7.2.1	Časopisi in revije
38	1.7.2.2	Leposlovje in stvarna besedila
43	1.8	Zbiranje po Gigafidi
43	1.9	Zaključek
45	<b>2</b>	<b>Spletna besedila korpusa Gigafida</b>
45	2.1	Uvod
46	2.2	Merila izbire in izbrane spletne strani
46	2.2.1	Besedila novinarskih portalov

47	2.2.2	Predstavitvene strani podjetij in ustanov
51	2.3	Tehnologije za zajemanje spletnih besedil
52	2.3.1	Zajemanje spletnih vsebin
53	2.3.1.1	Spletni pajki
54	2.3.1.2	Spletno pajkanje v projektu ssj
55	2.3.2	Odstranjevanje spremnih in vnaprej pripravljenih besedil
56	2.3.2.1	Obstoječi pristopi
60	2.3.2.2	Odstranjevanje spremnih in vnaprej pripravljenih besedil v projektu ssj
60	2.3.3	Detekcija jezika
63	2.3.3.1	Detekcija jezika v projektu ssj
64	2.3.4	Detekcija dvojnikov in približnih dvojnikov
65	2.3.4.1	Detekcija dvojnikov in približnih dvojnikov v projektu ssj
66	2.3.5	Nekaj zanimivih statistik
67	2.4	Zaključek
68	<b>3</b>	<b>Zapis korpusa Gigafida</b>
68	3.1	Zapis znakov Unikod
70	3.2	Jezik za označevanje XML
71	3.3	Priporočila za označevanje besedil TEI
72	3.3.1	Kolofon TEI
75	3.4	Besedilne oznake Gigafide
76	3.5	Zaključek
77	<b>4</b>	<b>Gradnja ter vsebina korpusov KRES, ccGigafida in ccKRES</b>
77	4.1	Reprezentativnost, uravnoteženost
79	4.2	KRES
79	4.2.1	Taksonomski deleži
80	4.2.2	Izbira besedil in njihovega obsega
81	4.2.2.1	Tisk
81	4.2.2.1.1	Knjižno
81	4.2.2.1.1.1	Leposlovje
81	4.2.2.1.1.2	Stvarna besedila
81	4.2.2.1.2	Periodično
81	4.2.2.1.2.1	Časopisi
85	4.2.2.1.2.2	Revije
89	4.2.2.1.3	Drugo
89	4.2.2.2	Internet
89	4.2.2.2.1	Novičarski portali
90	4.2.2.2.2	Podjetja in ustanove
90	4.2.3	Končno število besed in število besed po letih
92	4.3	ccGigafida in ccKRES
94	4.4	Postopek vzorčenja
95	4.5	Primerjava pogostosti lem v KRES-u in ccGigafidi
97	4.6	Zaključek
98	<b>5</b>	<b>Konkordančnik ssj z vmesnikom korpusa Gigafida</b>
98	5.1	Od Konkordančnika ASP32 do Konkordančnika ssj
99	5.2	»Splošni« uporabnik besedilnega korpusa

100	5.2.1	Starost in poklic korpusnih uporabnikov
100	5.2.2	Namen uporabe korpusa
101	5.2.3	Pogostost uporabe korpusa
102	5.2.4	Uporaba korpusu sorodnih jezikovnih virov
102	5.2.5	Način seznanitve s korpusom
103	5.3	»Splošna« uporaba besedilnega korpusa
104	5.3.1	Enostavno in razširjeno iskanje
104	5.3.2	Iskanje po kanalih
106	5.3.3	Napredne možnosti izdelave iskalnega pogoja
107	5.3.4	Obdelava konkordančnega niza
108	5.4	Novosti Konkordančnika ssj z vmesnikom Gigafida
109	5.4.1	Začetek dela s korpusom
109	5.4.2	Pomoč pri delu s korpusom
109	5.4.3	Vmesniška navigacija
110	5.4.4	Enostavno iskanje
111	5.4.5	Napredno iskanje
112	5.4.6	Izdelava seznama kolokatorjev
112	5.4.7	Izdelava besednega seznama
113	5.4.7	Podatkovni filtri
114	5.4.9	Tiskanje in izvod podatkov
114	5.5	Prikaz jezikovni podatkov v vmesniku Gigafida
114	5.5.1	Konkordančni niz
115	5.5.2	Seznam kolokatorjev
117	5.5.3	Besedni seznam
118	5.6	Zaključek
119	<b>6</b>	<b>FIDA in FidaPLUS kot predhodnika korpusa Gigafida</b>
119	6.1	Korpus FIDA
119	6.1.1	Zgodovina
121	6.1.2	Sestava
122	6.1.2.1	Taksonomija prenosnik
124	6.1.2.2	Taksonomija zvrst
126	6.1.2.3	Taksonomija lektorirano
127	6.1.2.4	Število besed po letih
128	6.1.3	Besedilodajalci
131	6.1.4	Format in metapodatki
137	6.2	Korpus FidaPLUS
137	6.2.1	Zgodovina
139	6.2.2	Taksonomija in število besed po letih
142	6.2.3	Besedilodajalci
142	6.2.4	Format in metapodatki; konkordančnik
143	6.3	Zaključek
144	<b>7</b>	<b>Povzetek</b>
147	<b>8</b>	<b>Summary</b>
150	<b>9</b>	<b>Literatura</b>
155	<b>10</b>	<b>Priloge</b>

# Kazalo tabel

- 17 Tabela 1.1: Zgradba FidePLUS glede na zvrst
- 18 Tabela 1.2: Taksonomija FidePLUS
- 20 Tabela 1.3: Predvideni delež besed po taksonomiji v Gigafidi
- 22 Tabela 1.4: Cobiss: najbolj izposojane knjige slovenskih avtorjev v letu 2009
- 28 Tabela 1.5: Natančnost statističnega označevalnika Obeliks
- 30 Tabela 1.6: Dvajset založb oz. naslovov besedil, ki so v Gigafido prispevali največ besed
- 32 Tabela 1.7: Število besed po taksonomiji v Gigafidi
- 32 Tabela 1.8: Končni in predvideni delež besed po taksonomiji v Gigafidi
- 34 Tabela 1.9: Delež besed po taksonomiji: primerjava med Gigafido in FidoPLUS
- 35 Tabela 1.10: Število in delež besed po letih v Gigafidi
- 37 Tabela 1.11: Časopisi, ki niso na lestvici NRB 2010, so pa vključeni v Gigafido
- 38 Tabela 1.12: Revije, ki niso na lestvici NRB 2010, so pa vključene v Gigafido
- 39 Tabela 1.13: Najbolj brani avtorji v letu 2009 (po Cobissovem seznamu najbolj izposojanih in največkrat rezerviranih knjig), kateri besedila so vključena v Gigafido
- 40 Tabela 1.14: Najbolj izposojani slovenski avtorji v letu 2009, katerih besedila so vključena v Gigafido
- 41 Tabela 1.15: Stvarna besedila v Gigafidi (naključni izbor)
- 47 Tabela 2.1: Pajkanje: novičarske strani
- 48 Tabela 2.2: Pajkanje: predstavitvene strani podjetij
- 49 Tabela 2.3: Pajkanje: predstavitvene strani ustanov
- 79 Tabela 4.1: Načrtovani delež in število besed po taksonomiji v KRES-u
- 79 Tabela 4.2: Delež besed po besedilnih zvrsteh v nekaterih tujih referenčnih korpusih
- 82 Tabela 4.3: Pridobljeni in nepridobljeni dnevniki, večdnevnik, tedniki ter brezplačniki iz NRB 2010 po branosti
- 84 Tabela 4.4: Časopisi: načrtovano število besed za KRES po branosti
- 84 Tabela 4.5: Število besed najbolj branih prilog v Gigafidi
- 85 Tabela 4.6: Pridobljeni in nepridobljeni tedniki, dvotedniki ter mesečniki iz NRB 2010 po branosti
- 87 Tabela 4.7: Revije: načrtovano število besed za KRES po branosti
- 89 Tabela 4.8: Najpogosteje obiskane novičarske spletne strani: število in delež prikazov za Slovenijo po merjenju Moss (julij 2010)
- 90 Tabela 4.9: Načrtovano število besed z novičarskih portalov za KRES
- 90 Tabela 4.10: Internetna besedila, ki so v Gigafido prišla iz korpusa FIDA
- 91 Tabela 4.11: Število besed po taksonomiji v KRES-u
- 91 Tabela 4.12: Število in delež besed po letih v KRES-u
- 122 Tabela 6.1: Taksonomija prenosnik v korpusu FIDA
- 123 Tabela 6.2: Število in delež dokumentov ter besed po taksonomiji prenosnik v korpusu FIDA



- 124 Tabela 6.3: Taksonomija zvrst v korpusu FIDA
- 125 Tabela 6.4: Število in delež dokumentov ter besed po taksonomiji zvrst v korpusu FIDA
- 126 Tabela 6.5: Taksonomija lektorirano v korpusu FIDA
- 126 Tabela 6.6: Število in delež dokumentov ter besed po taksonomiji lektorirano v korpusu FIDA
- 127 Tabela 6.7: Število in delež besed po letih v korpusu FIDA
- 128 Tabela 6.8: Pokritost tematskih sklopov v korpusu FIDA
- 129 Tabela 6.9: Besedilodajalci (institucije) korpusa FIDA in število besed, ki so jih prispevali v korpus
- 137 Tabela 6.10: Atribut @msds pridevnika *mladinski* v korpusu FIDA
- 139 Tabela 6.11: Število in delež besed po letih v FidiPLUS
- 140 Tabela 6.12: Taksonomija prenosnik: število besed v FidiPLUS ter razmerja med deleži besed v FidiPLUS in korpusu FIDA
- 141 Tabela 6.13: Taksonomija zvrst: število besed v FidiPLUS ter razmerja med deleži besed v FidiPLUS in korpusu FIDA
- 142 Tabela 6.14: Največji besedilodajalci FidePLUS ter število in delež besed, ki so jih prispevali v korpus

# Kazalo slik

- 13 Slika 1.1: Povezanost ciljev projekta ssj  
19 Slika 1.2: Taksonomija Gigafide  
22 Slika 1.3: Cobiss: najbolj izposojane knjige v letu 2009  
23 Slika 1.4: Cobiss: slovenski avtorji najbolj izposojanih knjig v letu 2009  
29 Slika 1.5: Del konkordančnih vrstic besede *sodelovati* v Gigafidi s filtroma *Vrsta besedila* in *Vir*  
32 Slika 1.6: Število besed po taksonomiji v Gigafidi  
34 Slika 1.7: Število besed iz besedil, izdanih do leta 2005 in pridobljenih pri novem zbiranju za Gigafido  
36 Slika 1.8: Število besed po letih v Gigafidi  
51 Slika 2.1: Cevovod za zajem besedil v projektu ssj  
53 Slika 2.2: Visokonivojska arhitektura tipičnega spletnega pajka  
56 Slika 2.3: Tipična novičarska spletna stran  
58 Slika 2.4: Prvih nekaj nivojev odločitvenega drevesa za odstranjevanje spremnih in vnaprej pripravljenih besedil  
62 Slika 2.5: Referenčna jezikovna profila za slovenski (levo) in angleški jezik (desno)  
62 Slika 2.6: Korelacija med slovenskim besedilom in slovenskim jezikovnim profilom (levo) ter korelacija med slovenskim besedilom in angleškim jezikovnim profilom (desno)  
66 Slika 2.7: Pajkanje: nekaj zanimivih statistik  
70 Slika 3.1: Primer dokumenta XML  
72 Slika 3.2: Struktura dokumenta TEI v Gigafidi  
73 Slika 3.3: Primer kolofona TEI v Gigafidi (1. del)  
74 Slika 3.4: Primer kolofona TEI v Gigafidi (2. del)  
75 Slika 3.5: Oznake besedila v Gigafidi  
91 Slika 4.1: Število besed po taksonomiji v KRES-u  
92 Slika 4.2: Število besed po letih v KRES-u  
97 Slika 4.3: Frekvenčni profil lem KRES-a in ccGigafide  
115 Slika 5.1: Del konkordančnega niza za iskalni pogoj *medvedki*  
116 Slika 5.2: Del seznama kolokatorjev za iskalni pogoj *medved*  
117 Slika 5.3: Del besednega seznama za iskalni pogoj *medved\**  
121 Slika 6.1: Vstopna stran spletnega konkordančnika korpusa FIDA  
124 Slika 6.2: Delež besed po taksonomiji prenosnik v korpusu FIDA  
125 Slika 6.3: Delež besed po taksonomiji vrst v korpusu FIDA  
128 Slika 6.4: Število besed po letih v korpusu FIDA  
131 Slika 6.5: Začetek dokumenta v formatu SGML v korpusu FIDA  
132 Slika 6.6: Primer kolofona TEI v korpusu FIDA  
140 Slika 6.7: Število besed po letih v FidiPLUS

# Spremna beseda

*»Pripravi naj se par strani dolgo besedilo o namenu korpusa in razlogih zbiranja materialov – namenjeno ljudem in institucijam, od katerih bomo skušali dobiti material.«*

**T**ako se je začelo. Bil je 24. januar 1997, skupaj so sedeli Tomaž Erjavec, Vojko Gorjanc, Simon Krek in Marko Stabej ter navedeno sprejeli kot sklep. V naslednjih mesecih je nastal še SGML/TEI »muštr za DZS proto korpus«, seznam meril za uravnoteževanje, glava korpusnih dokumentov, poskusno »tagiranje«, zametek konkordančnika in še marsikaj, pa seveda tudi – in to natanko na današnji dan pred 15 leti – pogodbeni zaveza k izvedbi projekta ter uradno poimenovanje cilja: Korpus slovenskega jezika FIDA.

Enak sklep je bil sprejet še dvakrat in bil v svojem »bomo skušali dobiti material« ter vsem, kar še sodi zraven, uresničen najprej kot FidaPLUS, pred kratkim pa še kot Gigafida. Leta 2008 smo se lotili štetja objav izsledkov raziskav, ki so vključevale vsaj vpogled v korpusa FIDA in FidaPLUS. Prišli smo do številke 60, nato pa ugotovili, da je tega še veliko več ... in – zavedajoč se, da količina ni vse, a smo lahko z njo zadovoljni – nehali šteti. Nedvomno lahko zapišemo, da sta oba predhodnika Gigafide v raziskovanje sodobne slovenščine ter razvoj jezikovnih tehnologij za slovenščino prinesla veliko sprememb: raziskovalno ugodje ob merljivosti podatkov, nove poglede na stara jezikovna pričanja, presenetljive menjave v jeziku tipičnega in posebnega, še večjo previdnost pri posploševanju ter interpretaciji vsega slovenističnega in porast statističnih metod, povezanih z jezikom, uresničljivost možnosti strojnega prevajanja slovenskih besedil, samodejnega učenja sistemov na podlagi učnih zbirk ter še in še.

Gradnja korpusa je seveda veliko več kot zbiranje, pretvarjanje in označevanje besedil, je tudi nenehno ponovno premišljanje o stvareh, ki so bile že premišljene, a jim je jeziko(slo)vni in informacijskotehnološki razvoj pokazal nove poti naprej, zato smo se odločili v knjigi predstaviti interno razumevanje stvari, utemeljiti odločitve, ki so bile kdaj tudi subjektivne, predstaviti dejavnike, ki so na naše delo vplivali od zunaj, ter pokazati na spoznanja domačih in tujih strokovnjakov, po katerih smo se zgledovali.

Zahvaljujemo se vsem besedilodajalcem, ki so nam brezplačno odstopili svoja besedila, marsikdaj pa zraven pripisali še prijazen »Z veseljem bomo sodelovali pri vašem delu, ker menimo, da bo koristilo razvoju slovenskega jezika«.

Gradnja vseh štirih korpusov skupaj z vmesniškim delom je bila kompleksna ter ni niti približno rezultat dela samo piscev te knjige – pri gradnji, vsebini in uporabnosti Gigafide, KRES-a, ccGigafide ter cc-KRES-a so sodelovali tudi Miro Romih, Peter Holozan, Rok Rejc, Simon Rigač, Iztok Kosem in Simon Šuster. Hvala vsem, z veseljem še kdaj!

Ljubljana, Kopenhagen, Škofja Loka, 1. julij 2012

Avtorji

# 1 Zbiranje besedil in vsebina korpusa Gigafida

1 Operacija je delno financirala Evropska unija iz Evropskega socialnega sklada ter Ministrstvo za izobraževanje, znanost, kulturo in šport Republike Slovenije. Operacija se je izvajala v okviru Operativnega programa razvoja človeških virov za obdobje 2007–2013, razvojne prioritete: razvoj človeških virov in vseživljenjskega učenja; prednostne usmeritve: izboljšanje kakovosti in učinkovitosti sistemov izobraževanja in usposabljanja 2007–2013.

## 1.1 Uvod

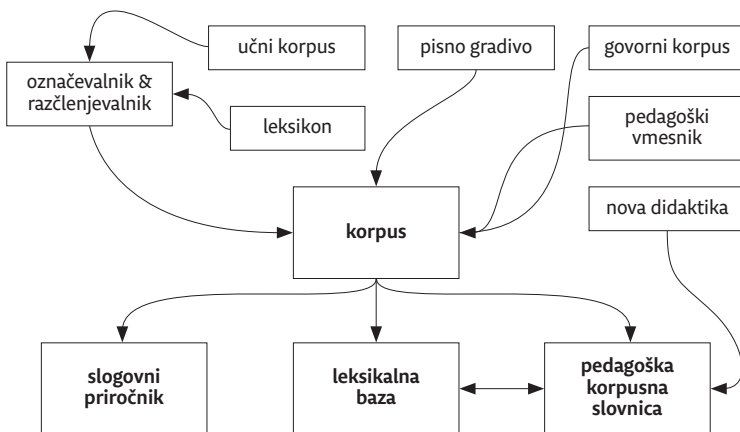
### 1.1.1 Cilj

Izgradnja referenčnega, enojezičnega in pisnega korpus slovensčine Gigafida pomeni uresničitev enega od ciljev projekta *Sporazumevanje v slovenskem jeziku* (<http://www.slovenscina.eu>, <http://www.projekt.slovenscina.eu>, dalje ssj).<sup>1</sup> Cilji projekta so bili sicer trije:

1. referenčni korpus in leksikalna baza slovenskega jezika s slovničnim analizatorjem,
2. jezikovne tehnologije kot del didaktičnih pristopov v vzgojno-izobraževalnih procesih,
3. pedagoška korpusna slovnica in slogovni priročnik.

Na kakšen način je Gigafida kot korpus vpeta v vse tri cilje, prikazuje naslednja slika:

Slika 1.1: Povezanost ciljev projekta SSJ.



Končni cilj gradnje referenčnega korpusa je bil nov javno in prosto dostopen pisni korpus v obsegu do ene milijarde pojavníc oz. besed, ki bo izdelan po zgledu korpusov FIDA in FidaPLUS ter zapisan v formatu XML TEI P5; določeno je bilo, da bo korpus lematiziran, v celoti oblikoskladenjsko označen, v določenem delu skladenjsko razčlenjen

in bo imel orodje za avtomatsko prepoznavo lastnih imen. Kot bo natančneje razvidno v nadaljevanju, je bil cilj v celoti dosežen.

(O drugih ciljih projekta gl. <http://www.slovenscina.eu> ter med drugimi: Gantar 2008; Arhar Holdt 2009; Gantar 2009; Gantar, Krek 2009; Gantar 2010; Krek, Arhar Holdt 2010; Verdonik in dr. 2010; Rozman, Krapš Vodopivec 2010; Gantar, Krek 2011; Kosem, Može 2011; Zwitter Vitez, Krapš Vodopivec 2011; Zwitter Vitez 2011; Arhar Holdt 2011; Verdonik, Zwitter Vitez 2011; Kosem in dr. 2011; Gantar 2011.)

## 1.1.2 Namen

Namen Gigafide v okviru projekta ssj je povezan s prikazom realne podobe slovenskega jezika v pedagoški korpusni slovnici, slogovnem priručniku in leksikalni bazi slovenskega jezika, in sicer tako v smislu iz korpusa pridobljenih podatkov ter njihovih interpretacij kot konkretnih zgledov. Sicer pa je Gigafida namenjena raziskovanju jezika na več ravneh. Ob odgovorih na posamezne sprotne poizvedbe je še pomembneje, da daje podatke o celotni podobi jezika, tako da je danes skoraj edini razmeroma zanesljiv vir za izdelavo sodobnih slovníc, slovarjev in različnih jezikovnih priručnikov za slovenščino, uporablja pa se tudi v jezikovnih tehnologijah. Z Gigafido želimo seznaniti ne le znanstvenike in raziskovalce v jezikoslovlju, temveč tudi učitelje slovenščine v osnovnih in srednjih šolah, tiste, ki se slovenščine učijo kot drugega ali tujega jezika, pa tudi vse, ki grede namesto na knjižno polico odgovor na svojo jezikovno zadrego raje iskat na svetovni splet.

## 1.2 Merila gradnje

Gigafida je nadgradnja referenčnega korpusa slovenskega jezika FidaPLUS (<http://www.fidaplus.net>), ki je v obsegu več kot 621 milijonov besed na spletu prosto dostopen od leta 2006 in že vključuje (oz. nadgrajuje) prvi tak korpus slovenščine, tj. v letih 1997–2000 nastali korpus FIDA (več o teh dveh korpusih gl. v 6. pogl. in tam navedeni literaturi).

Da bi dosegli cilj milijardnega korpusa, smo morali FidoPLUS dopolniti s približno 380 milijoni besed, nabor besedil, s katerim smo želeli to doseči, pa so vodila merila gradnje, ki smo jih določili v *Standardu za redno zbiranje pisnega gradiva za referenčni korpus* (december 2008; Kazalnik 1 na projektni spletni strani; dalje Standard za zbiranje gradiva).

## 1.2.1 Standard za zbiranje gradiva

Standard za zbiranje gradiva vključuje izhodiščni premislek lastnosti, ki jih lahko pripišemo besedilom oz. jih prepoznamo v besedilih in na podlagi katerih usmerjamo zbiranje gradiva ter uravnotežujemo korpus. Merila zbiranja smo določili na podlagi v domači in tuji literaturi popisanih spoznanj (npr. Atkins, Clear, Ostler 1992; Gorjanc 2002: 32–33; Arhar Holdt 2004; McEnery, Xiao, Tono 2006), na podlagi izkušenj, pridobljenih pri gradnji korpusov FIDA in FidaPLUS, ter na podlagi pogovorov med člani ožje projektne skupine, povezane s pisnim korpusom (po abecednem vrstnem redu: Š. Arhar Holdt, P. Gantar, V. Gorjanc, P. Kocjančič, S. Krek, N. Logar Berginc, M. Stabej, M. Šorli in S. Šuster). Merila zbiranja besedil so bila naslednja: besedilna zvrst/ vrsta, področje/tema, dolžina besedil, ustroj dokumenta, avtorstvo, ciljna publika, branost, prenosnik, objavljenost/internost/zasebnost, čas izdaje/nastanka, prevedenost in lektoriranost. V nadaljevanju smernice za zbiranje gradiva na kratko povzemamo.

**a) Besedilna zvrst/vrsta:** Za novi korpus smo zbirali javno objavljena pisna besedila raznoterih zvrsti/vrst in žanrov. Glede na to, da je bilo v FidiPLUS umetnostnih besedil 3,48 %, vseh drugih (brez upoštevanja neopredeljenih besedil) pa 96,41 % (Tabela 1.1), smo si prizadevali dobiti čim več umetnostnih besedil. Ta so bila v FidiPLUS razdeljena na prozna, dramska in pesniška, v novem korpusu pa leposlovje ni podrobneje členjeno, saj je bilo utemeljeno pričakovati, da bo podobno kot v FidiPLUS delež dramskih in pesniških besedil izredno majhen (več o tem gl. v točki 1.2.2).

**b) Področje/tema:** Pri zbiranju smo si prizadevali dobiti gradivo z različnih področij in različnih tem:

- aktualni dogodki
- gospodarstvo, politika
- vzgoja in izobraževanje
- narava, dom, hišni ljubljenci
- ljudje, družina, moški, ženske, otroci, mladina
- zdravje, hrana
- posel, finance
- prosti čas, glasba, film, razvedrilo, moda
- šport, turizem
- kultura, umetnost
- religija, duhovnost
- računalništvo, avtomobilizem itd.

**c) Dolžina besedil:** Pri dolžini besedil za vključitev v novi korpus ni bilo omejitev, smo pa za dela, ki bi izstopala po svojem velikem obsegu

ali bi zanje tako željo izrazil avtor oz. založba, predvideli možnost individualne odločitve o skrajšanju ali vključitvi le določenih delov.

**č) Ustroj dokumenta:** Korpusni dokument (z eno besedilno glavo oz. kolofonom) je lahko sestavljen iz enega besedila (npr. cel roman) ali več besedil (časopis, revija, zbirka pesmi ipd.). Naknadna členitev večbesedilnih dokumentov na enobesedilne se pri gradnji korpusa ni izvajala. Veljalo je tudi obratno: besedil, pridobljenih v več dokumentih, nismo združevali v en dokument.

**d) Avtorstvo:** Pri gradivu, pri katerem je avtorstvo razvidno oz. merljivo, smo bili pozorni na to, da posamezni avtorji ne bi bili prekomerno zastopani, kjer števila avtorjev in obsega njihovih besedil ni bilo mogoče na preprost način nadzorovati (časopisi, revije, internet, drugo), pa smo to merilo zanemarili. Pri tem lastnosti, kot so spol, starost, tip (posameznik, ustanova), regijska pripadnost, nacionalnost in prvi jezik avtorja ter število (eden, več) avtorjev, naknadno nismo ugotavljali ter na zbiranje niso vplivale. Le pri časopisih in revijah smo bili pozorni na regijsko razpršenost (lokalno, izseljensko, zamejsko). Ime in priimek avtorja sta del bibliografskih podatkov v kolofonu korpusnih dokumentov pri tistih enobesedilnih dokumentih, ki imajo podatek na voljo brez iskanja.

**e) Ciljna publika:** Spol, starost, regijska pripadnost in raven izobrazbe ciljne publike niso vplivali na zbiranje. Zahtevnih specializiranih besedil za korpus nismo pridobivali.

**f) Branost:** Branost je najpomembnejši kazalnik recepcije pisnih besedil (več gl. nadaljevanju tega poglavja in v 4. pogl.).

**g) Prenosnik:** Za novi korpus smo zbirali pisna besedila, ki so (a) tiskana (in sicer periodična, če je zanje značilna rednost ali pogostnost izhajanja, ter knjižna) in (b) internetna. Pri zadnjih smo se omejili na strani z informativnimi vsebinami, in sicer z dveh vsebinskih vidikov – zajeli smo: besedila novičarskih portalov ter predstavljene strani podjetij in državnih, pedagoških, raziskovalnih, kulturnih ipd. ustanov (več o tem gl. v 2. pogl.).

**h) Objavljenost/internost/zasebnost:** Novi korpus vsebuje objavljena besedila, ki jih razumemo kot javno dostopna besedila. Zasebnih in internih besedil, kot so npr. okrožnice v podjetju, zapiski ter vabila, ki so namenjeni ožji, znani skupini ljudi, na novo nismo zbirali.



**i) Čas nastanka/izdaje:** Čas nastanka/izdaje je relevanten z dveh vidikov:

– Vidik produkcije: Besedilodajalce, ki so v FidoPLUS že prispevali besedila, smo prosili za dela, ki so jih izdali po letu 2005; pri novih besedilodajalcih smo skušali dobiti besedila, ki so jih izdali po letu 1995.

– Vidik recepcije: branost in izposoja tiskanih del ter obiskanost spletnih strani niso nujno povezane z novejšim datumom nastanka del. Če npr. podatki o izposoji v knjižnicah kažejo visoko branost starejših del (zlasti t. i. klasikov), smo si ta besedila prizadevali dobiti.

**j) Prevedenost/izvirnost:** V novi korpus so vključena tudi prevedena dela, njihov delež vnaprej ni bil določen. Zaželeno je bilo, da so izvirniki v različnih jezikih.

**k) Lektoriranost:** V FidiPLUS je bila oznaka »nelektorirano« pripisana le zelo majhnemu delu korpusa (0,6 %), oznaka »lektorirano« pa večinoma avtomatsko vsemu periodičnemu in knjižnemu gradivu. V novem korpusu lektoriranosti nismo beležili in je tudi nismo razumeli kot uravnoteževalne lastnosti, saj je njeno naknadno ugotavljanje časovno potratno, pridobitev znatnega deleža nelektoriranih besedil pa prav tako zelo zamudna ter primernejša za specializirani korpus (prim. tudi točko 6.1.2.3 v 6. pogl.).

Premisleku o različnih lastnostih besedil je sledila ocena okvirnih deležev besed, ki jih bodo v novi korpus prinesle posamezne besedilne zvrsti/vrste (prim. Tabela 1.1, v kateri so podatki o zgradbi FidePLUS glede na zvrst). To pomeni, da je bila delu zgoraj predstavljenih lastnosti že v času priprav na zbiranje pripisana količina, ki smo jo želeli vključiti v korpus. Hkrati so nekatere od teh lastnosti postale kategorije korpusove taksonomije.

**Tabela 1.1: Zgradba FidePLUS glede na zvrst.\***

<b>Taksonomija zvrst</b>	<b>FidaPLUS: število besed</b>	<b>FidaPLUS: delež v %</b>	<b>SKUPAJ v %</b>
NI PODATKA	709.344	0,11	0,11
Ft.Z.N (neumetnostna)	368.208	0,06	
Ft.Z.N.N (nestrokovna)	536.314.007	86,34	
Ft.Z.N.P (pravna)	124.817	0,02	
Ft.Z.N.S (strokovna)	4.530.801	0,73	
Ft.Z.N.S.H (humanistična in družboslovna)	19.331.249	3,11	96,41
Ft.Z.N.S.N (naravoslovna in tehnična)	38.202.106	6,15	

Ft.Z.U (umetnostna)	543.750	0,09	
Ft.Z.U.D (dramska)	480.957	0,08	
Ft.Z.U.P (pesniška)	366.215	0,06	3,48
Ft.Z.U.R (prozna)	20.178.021	3,25	
<b>SKUPAJ</b>	<b>621.149.475</b>	<b>100,00</b>	<b>100,00</b>

\* Podatki so iz Erjavec 2008.

## 1.2.2 Taksonomija

Medtem ko je bila taksonomija FidePLUS tridelna (prenosnik, zvrst, lektoriranost; Tabela 1.2) in tudi dalje notranje dokaj podrobno členjena (prim. npr. periodično, ki je imelo podkategoriji časopisno ter revijalno, znotraj druge pa nato še tedensko, štirinajstdnevno, mesečno, redkeje kot na mesec in občasno), smo taksonomijo Gigafide poenostavili v enodelno in členjeno do tretje podravnine (Slika 1.2).

**Tabela 1.2: Taksonomija FidePLUS.**

Ft.P – prenosnik	Ft.P.P.N.J – javno
Ft.P.G – govorni	Ft.P.P.N.I – interno
Ft.P.E – elektronski	Ft.P.P.N.Z – zasebno
Ft.P.P – pisni	
Ft.P.P.O – objavljeno	Ft.Z – zvrst
Ft.P.P.O.K – knjižno	Ft.Z.U – umetnostna
Ft.P.P.O.P – periodično	Ft.Z.U.P – pesniška
Ft.P.P.O.P.C – časopisno	Ft.Z.U.R – prozna
Ft.P.P.O.P.C.D – dnevno	Ft.Z.U.D – dramska
Ft.P.P.O.P.C.V – večkrat tedensko	Ft.Z.N – neumetnostna
Ft.P.P.O.P.C.T – tedensko	Ft.Z.N.S – strokovna
Ft.P.P.O.P.R – revijalno	Ft.Z.N.S.H – humanistična in družboslovna
Ft.P.P.O.P.R.T – tedensko	Ft.Z.N.S.N – naravoslovna in tehnična
Ft.P.P.O.P.R.S – štirinajstdnevno	Ft.Z.N.N – nestrokovna
Ft.P.P.O.P.R.M – mesečno	
Ft.P.P.O.P.R.D – redkeje kot na mesec	Ft.L – lektorirano
Ft.P.P.O.P.R.O – občasno	Ft.L.D – da
Ft.P.P.N – neobjavljeno	Ft.L.N – ne

**Slika 1.2: Taksonomija Gigafide.**

---

tisk
knjižno
leposlovje
stvarna besedila
periodično
časopisi
revije
drugo
internet

---

2 Posredno o večji vplivnosti govorijo podatki raziskave Slovenija in internet 2005–2008 (Raba interneta v Sloveniji): delež gospodinjstev, ki uporabljajo internet, se je s 43 % v letu 2004 povzpел na 58 % v letu 2008, prav tako se je povečal delež dnevnih uporabnikov interneta z 28 % v letu 2005 na 42 % v letu 2008.

3 Korpus govornje slovenščine GOS (<http://www.korpus-gos.net/>), ki je prav tako nastal v projektu SSJ, namreč vključuje le spontani govor (Zemljarič Miklavčič in dr. 2009; Verdonik, Zwitter Vitez 2011).

V nadaljevanju bomo predstavili razloge, ki so nas vodili k oblikovanju take taksonomije – v skladu z dejstvom, da gre za nadgradnjo že obstoječega korpusa, so ti razlogi podani primerjalno s FidoPLUS oz. temeljijo na povratnih informacijah v zvezi z njo.

### 1.2.2.1 Tisk in internet

Tradicionalnemu pisnemu prenosniku – tisku – se je v javnih sporočajnskih položajih vsaj v zadnjem desetletju kot vsakodnevni način prenosa sporočil pridružil še elektronski. V FidoPLUS je internetnega gradiva 1,24 %. V nastajajočem korpusu smo se zaradi večje vplivnosti<sup>2</sup> odločili ta delež povečati, ker pa je šlo tudi v tehničnem in metodološkem smislu za prvi večji poskus pridobivanja besedil s svetovnega spleta za referenčni korpus pri nas, smo se – kot smo že zapisali – omejili na strani z informativnimi vsebinami (več gl. v 2. pogl.).

### 1.2.2.2 Knjižnost, periodičnost in drugo

V obliki knjige izdana besedila so v FidoPLUS prinesla slabih 9 % besed, skoraj vse drugo izhaja iz publicistične periodike. Načinu izhajanja – enkrat (z možnostjo ponatisa) : večkrat – smo poskusno pridružili še deloma odprto skupino »drugo«. Zanj smo se odločili zbirati podnapise tujih filmov, nadaljevank in dokumentarnih oddaj (vključno s podnapisi za slušno prizadete) ter besedila, ki so v različnih oddajah brana – t. i. scenarije in postproduksijska besedila.<sup>3</sup>

#### 1.2.2.2.1 LEPOSLOVJE IN STVARNA BESEDILA

Kot je razvidno v Tabelah 1.1 in 1.2, je bila v FidoPLUS uporabljena delitev na umetnostna in neumetnostna besedila. Določitev, ali gre za umetnostna besedila ali ne, je samodejno mogoča le pri knjižnem gradivu (pri časopisju, ki tudi lahko vsebuje besedila umetnostne zvrsti, zaradi večbesedilnosti dokumentov to skoraj ni mogoče (vsekakor pa ni časovno gospodarno)), zato ti dve skupini v novi enodelni

4 Korošec (1976: 106) je znotraj publicistike izrecno ločil le na vsakodnevno izhajanje vezano poročevalstvo – kajti vsakodnevno pisanje o podobnih ali ponavljajočih se situacijah je najpomembnejši objektivni stilotvorni dejavnik časopisnega poročevalstva, ki je od jezika zahteval prilagoditev novi vlogi in s tem nastanek novega, tj. poročevalnega stila (prim. tudi Kalin Golob 2003).

taksonomiji umeščamo kot podravnini v kategorijo knjižno. Namesto sicer na tradiciji slovenske zvrstnosti temelječega poimenovanja ne-umetnostni, ki izraža pravzaprav to, česa v tej skupini *ni* (z izločitvijo publicistike pa postane hkrati tudi preširoko), smo se knjižna besedila z nefikcijsko vsebino odločili poimenovati stvarna besedila (tudi oznaka strokovna besedila je namreč zavajajoča), njej nasprotno skupino pa leposlovje.

### 1.2.2.2.2 ČASOPISI IN REVIJE

Delež časopisne in revijalne periodike je v korpusu FidaPLUS daleč največji – več kot 85-odstotni. Tudi na podlagi odzivov stalnih uporabnikov tega korpusa (sicer zaznanih povsem nesistematično) v smislu, da je – čeprav najvplivnejši – novinarski jezik v korpusu količinsko preveč izpostavljen, smo se odločili, da bomo v uravnoteženem 100-milijonskem delu Gigafide, tj. v korpusu KRES (o njem gl. 4. pogl.), delež publicistike zmanjšali, opustili pa smo tudi delitev na tedensko, štirinajstnevno ipd., ker je raziskave slovenskega poročevalstva kot stilotvorno ali jezikovnorazlikovalno relevantne (še) niso potrdile,<sup>4</sup> za referenčni korpus pa je preveč podrobna.

Tabela 1.3: Predvideni delež besed po taksonomiji v Gigafidi.

Taksonomija	Oznaka	Delež besed v %
tisk	T	50–90
knjižno	T.K	15–35
leposlovje	T.K.L	20–50
stvarna besedila	T.K.S	30–60
periodično	T.P	20–40
časopisi	T.P.C	30–70
revije	T.P.R	30–70
drugo	T.D	5–10
internet	I	10–50

Pri oblikovanju taksonomije z deleži nas je vodilo tudi pravilo, ki smo ga posredno že nakazali: vključili smo le kategorije, za katere je bilo pričakovati, da bomo zanje lahko dobili toliko besedil, da bo obstoj kategorije upravičen, tj. da bo dosegel vsaj 5 % v 100-milijonskem KRES-u. Opustili smo kategorije, ki zahtevajo več notranjega uravnoteževanja in več časa pri zbiranju, saj je zanje bolj smiselna gradnja specializiranih korpusov (npr. korpus zasebnih besedil ali korpus nelektoriranih besedil). Za opustitev nekaterih podravnin taksonomije smo se odločili tudi na podlagi podatkov o načinih iskanja po FidiPLUS. Analiza iskanj, opravljena v novembru 2008, je pokazala, da je bilo kar 93 % izdelav konkordanc v FidiPLUS izvedeno pri osnovnem iskanju, le 7 % zahtev po pridobitvi konkordančnih nizov

pa je potekalo v razširjenem iskanju z izbiro taksonomskih kategorij, časa nastanka dela ali izpisa Cobiss. V teh primerih so bila nekatera iskanja izredno redka, tako so bile npr. podkategorije pri revijalnih in časopisnih besedilih glede na pogostost izhajanja izbrane v manj kot enem odstotku razširjenih iskanj. Sicer pa je bil v okviru razširjenega iskanja prenosnik izbran v 15 %, čas nastanka dela v 35 %, zvrst v 17 %, lektoriranost v 18 % in izpis Cobiss v 4 % (prim. tudi podatke v 5. pogl. ter v Arhar Holdt 2009b in 2010). Kljub na videz manjši izbirnosti vnaprej pripravljenih možnosti razširjenega iskanja zaradi enodelne in poenostavljene taksonomije je uporabnikom novega korpusa še vedno omogočena izdelava podkorpusov na podlagi podatkov v kolofonu korpusnih dokumentov oz. korpusovih filtrov (več o tem v točki 5.4.8 v 5. pogl.).

Čeprav smo pregledali stanje v tujih korpusih (ki pa je zelo različno, prim. Tabela 4.2 v 4. pogl.), so bili deleži v taksonomiji Gigafide v končni fazi subjektivna odločitev sestavljalcev korpusa, zavedali pa smo se, da bo uporabnikom korpusa treba dati možnost prepoznanja teh subjektivnih odločitev v smislu, da je korpus sicer zaznamovan s teoretičnimi prepričanji in odločitvami svojih snovalcev, vendar mora biti uporabnikom omogočeno, da to zaznamovanost razberejo in presežejo (Stabej 1998: 98).

## 1.3 Zbiranje besedil

### 1.3.1 Podatki za zbiranje

V slovenskem prostoru je mogoče podatke, iz katerih lahko okvirno sklepamo o recepciji besedil, dobiti iz več virov.

#### 1.3.1.1 Nacionalna raziskava branosti

V okviru *Nacionalne raziskave branosti* (dalje NRB) se zbirajo podatki o bralnih navadah bralcev časopisov in revij. Raziskavo izvaja družba Valicon, d. o. o., njen naročnik pa je Svet pristopnikov k NRB, ki deluje v okviru Slovenske oglaševalske zbornice. Splošni podatki iz raziskave so objavljeni dvakrat na leto na spletni strani <http://www.nrb.info/>. Pri zbiranju besedil za Gigafido, še bolj pa pri uravnoteževanju korpusa KRES (4. pogl.), smo izhajali iz podatkov NRB za leta 2006, 2007, 2008, 2009 in 2010 (podatki za leto 2010 – valutno obdobje: 2. polovica leta 2009 in 1. polovica leta 2010 – so v Prilogi 1). Tako je npr. iz NRB 2010 razvidno, da je bil najbolj bran časopis (tj. besedilo, ki je v korpusu označeno kot časopis, T.P.C) brezplačnik *Žurnal*, sledil mu je *Nedeljski dnevnik*, nato zopet brezplačnik *Dobro jutro*, na četrtem mestu so bile *Slovenske novice* itd. Med revijami (T.P.R) je bila najbolj brana *Lady*, sledili so ji *Ognjišče*, *Motorevija*, *Zdravje* itd.

### 1.3.1.2 Izposoja v knjižnicah

Drugi vir podatkov o branosti je knjižnična izposoja, ki pove, katere knjige so bile v knjižnicah, ki so vključene v sistem Cobiss ter imajo avtomatsko izposajo, najbolj izposojane in največkrat rezervirane ter kateri slovenski avtorji in njihova dela so bili najbolj izposojani (gre za avtorje, ki so upravičeni do knjižničnega nadomestila). Podatki so na voljo na spletni strani <http://www.cobiss.si/>. Prvih deset najbolj izposojanih knjig v letu 2009 prikazuje Slika 1.3. Na tem seznamu je izmed stotih del 17 del slovenskih avtorjev (Tabela 1.4), vsa druga dela so prevodi. Iz Tabele 1.4 je razvidno, da so najbolj izposojane knjige domačih avtorjev bodisi otroška ali mladinska literatura bodisi obvezno šolsko branje. Prvih 10 slovenskih avtorjev, katerih dela so bila v letu 2009 najbolj izposojana, prikazuje Slika 1.4.

Slika 1.3: Cobiss: najbolj izposojane knjige v letu 2009.

leto:  mesec:  gradivo:

Število knjižnic, ki ustrezajo iskalnemu pogoju: 261

	Naslov	Avtor	Izposoj.	Rezerv.	COBISS/OPAC
1.	Pepel v vetru	Woodiwiss, Kathleen E.	11815	2931	<input type="button" value="COBISS.SI-ID"/>
2.	Vreden ljubezni	Woodiwiss, Kathleen E.	11638	2835	<input type="button" value="COBISS.SI-ID"/>
3.	Antigona	Sophocles	10283	194	<input type="button" value="COBISS.SI-ID"/>
4.	Pride ženska k zdravniku --	Kluun	7551	3641	<input type="button" value="COBISS.SI-ID"/>
5.	Zločin in kazen	Dostoevskij, Fedor Mihajlovič	7347	214	<input type="button" value="COBISS.SI-ID"/>
6.	Varna vožnja : priročnik za voznike		7186	351	<input type="button" value="COBISS.SI-ID"/>
7.	Somrak	Meyer, Stephenie	7121	5539	<input type="button" value="COBISS.SI-ID"/>
8.	Mlada luna	Meyer, Stephenie	6956	4214	<input type="button" value="COBISS.SI-ID"/>
9.	Matilda	Dahl, Roald	6640	285	<input type="button" value="COBISS.SI-ID"/>
10.	Zimska vrtnica	Woodiwiss, Kathleen E.	6626	684	<input type="button" value="COBISS.SI-ID"/>

Tabela 1.4: Cobiss: najbolj izposojane knjige slovenskih avtorjev v letu 2009.

Mesto na lestvici	Avtor	Delo
11.	Goran Vojnovič	Čefurji raus!
19.		Od Ivana Preglja do Cirila Kosmača: izbor novel
26.	Ela Peroci	Muca Copatarica
27.	Prežihov Voranc	Solzice
29.	Desa Muck	Anica in Grozovitež
40.	Ivan Cankar	Na klancu
48.	Ivan Tavčar	Visoška kronika
55.	Tone Seliškar	Bratovščina Sinjega galeba
66.	Ivan Cankar	Hlapci
70.	Frane Milčinski	Zvezdica Zaspanka
71.	Svetlana Makarovič	Sapramiška
77.	Desa Muck	Anica in počitnice
78.	Vid Pečjak	Drežček in trije marsovčki