



LATVIJAS
UNIVERSITĀTE
ANNO 1919



Latviešu valodas
aģentūra

Latviešu valodas apguvēju korpusa izveide

Ilze Auziņa, Roberts Dargis,
Kristīne Levāne-Petrova

Latvijas Universitātes 76. zinātniskās konferences Datorlingvistikas sekcija
2018. gada 1. martā

Mērķis

- Izveidot valsts valodas prasmes pārbaudes darbu datubāzi – t.s. latviešu valodas apguvēju korpusu.
- Noskaidrot latviešu valodas prasmes kvalitāti, analizējot valsts valodas prasmes pārbaudes kārtotāju darbus un to rezultātus.

Korpusā iekļautie teksti

- 900 valsts valodas prasmes pārbaudes darbu rakstītprasmes testi: no katra valodas apguves līmeņa (A1, A2, B1, B2, C1, C2) datubāzē iekļauti 150 darbi.
- Tikai tie rakstītprasmes testa uzdevumi, kuros ir saistīts teksts.
- Apjoms: 1496 teksti, 142 684 teksta vienības.
- Kļūdu skaits tekstos: 22,2 %.
- Datubāzē iekļauti arī valsts valodas prasmes pārbaudes darbu runātprasmes testi (no katra valodas apguves līmeņa 17 audioieraksti).

Korpusa izveides posmi

- Datu digitalizācija
- Tekstu normalizēšana
- Automatizēta morfológiskā anotēšana, tostarp vārdšķiru noteikšana, lemmatizēšana
- Oriģinālā un normalizētā teksta sastatīšana
- Automatizēta kļūdu anotēšana un manuāla pārskatīšana

Datu digitalizācija un tekstu normalizēšana

Text

Labdien mana miļa draudzene Gunita. Es gribu ar tevi satakties. Es gribu tevi aicināt uz cirku. Mes kopā brauksim uz Rīga. Pastaigāsimies pa vecrīgu un parunāt. Biļeti cenu 12 eiro. Es zinu ka manu draudzenei interese cirk. Gaidīšu tavu atbildi
Visu labu.

Correct text

Labdien, mana mīļā draudzene Gunita! Es gribu ar Tevi satikties. Es gribu Tevi aicināt uz cirku. Mēs kopā brauksim uz Rīgu. Pastaigāsimies pa Vecrīgu un parunāsim. Biļetes cena ir 12 eiro. Es zinu, ka manai draudzenei interesē cirks.
Gaidīšu Tavu atbildi.
Visu labu!

Kļūdu tipi (1/3)

Tips	Apakštips	Piemēri
Pareizrakstības kļūdas	Lielie /mazie sākumburti	<i>latvija</i> – <i>Latvija</i> ; <i>Kā tev iet?</i> – <i>Kā Tev iet?</i>
	Diakritiskās zīmes	<i>tēv</i> – <i>tev</i> ; <i>bēt</i> – <i>bet</i> ; <i>bilete</i> – <i>biļete</i>
	Kopā / šķirti rakstāmi vārdi	<i>Uzredzēšanos!</i> – <i>Uz redzēšanos!</i> <i>vis lielākais</i> – <i>vislielākais</i>
	Izlaisti burti	<i>30. maja</i> – <i>30. maijā</i>
	Lieki burti	<i>uotraja augustā</i> – <i>otrajā augustā</i>
	Citas pareizrakstības kļūdas	<i>lajme</i> – <i>laime</i> <i>Es gaigu atbildi!</i> – <i>Es gaidu atbildi!</i>
Formveidošanas kļūdas	Nepareiza vārda forma (piem., lietv. – locījums, skaitlis; īpaš. v. – dzimte, noteiktā /nenoteiktā galotne; verbiem – izteiksme, persona)	<i>Es dzīvoju Latvija</i> – <i>Es dzīvoju Latvijā</i> <i>Es gribas apmeklēt šo pasākumu</i> – <i>Man gribas apmeklēt šo pasākumu</i>
	Skaņu mijas esamība/ neesamība	<i>vīreti</i> – <i>vīrieši</i> ; <i>angla valoda</i> – <i>angļu valoda</i>

Kļūdu tipi (2/3)

Tips	Apakštips	Piemēri
Leksikas kļūdas	Nozīmes ziņā neatbilstoša vārda lietojums	<p><i>Izeja</i> ir bez maksas. – <i>leeja</i> ir bez maksas. Es tev <i>pazvanīšu</i>. – Es Tev <i>piezvanīšu</i>. Aptauja piedalījās Latvijās <i>dzivotāji</i>. – Aptaujā piedalījās Latvijas <i>iedzīvotāji</i>.</p> <hr/> <p><i>Ka</i> klajas? – <i>Kā</i> klājas? Nopirkt lūdzu mani 2 biļetus <i>no</i> cirku. – Nopirksi, lūdzu, man 2 biļetes <i>uz</i> cirku. Biļešu maksa 1,5 eiro, <i>kāpēc</i> es esmu pensioniere. – Biļešu maksa ir 1,5 eiro, <i>tāpēc ka</i> es esmu pensionāre.</p> <hr/>
	Nesaprotams vārds / latv. val. neeksistējošs vārds	<p>Tas ir smuk <i>pūtik</i> un dažados krāsa. – Tās ir smukas <i>puķes</i> un dažādās krāsās.</p> <hr/>
	Kalkēts vārds	<p><i>gimnasti</i> – <i>vingrotāji</i>; <i>nedārgs</i> – <i>nav dārgs</i>; <i>surpriziņš</i> – <i>pārsteigums</i>; <i>papraši</i> – <i>lūgt</i>; <i>importanta</i> – <i>svarīga</i></p>

Kļūdu tipi (3/3)

Tips	Apakštips	Piemēri
Sintakses kļūdas	Vārdu secība	<i>cena ekskursija – ekskursijas cena svetki sporta – sporta svētki</i>
	Lieks vārds	<i>Apsveicu tevi ar svētku dienu. – Apsveicu Tevi svētku dienā!</i>
	Izlaists vārds	<i>Gaidīšu tevi informāciju – Gaidīšu no Tevis informāciju</i>
	Nepilna teikuma struktūra	<i>Rozes sievietēm. Es gribu pasākum. – Es gribu apmeklēt pasākumu.</i>
Interpunkcijas kļūdas	Pieturzīmes trūkums	<i>Sveiki Anna! – Sveiki, Anna!</i>
	Lieka pieturzīme	<i>Labdien, dārgais, draugs – Labdien, dārgais draugs!</i>
	Neatbilstoša pieturzīme	<i>Kā tev iet. – Kā tev iet?</i>
Neskaidrs teksts	Jālabo viss izteikums. Turklāt ir iespējami vairāki varianti.	<i>Olga tu gribet pasakumu man teātra. Es ceru apmeklēt akrobāti un dresēti dzīvnieki.</i>

Meklēšana korpusā (1/2)

Korpusā ir iespējams meklēt, izmantojot dažādas pazīmes:

- vārds / vārdforma
- pamatforma
- morfoloģiskās pazīmes (*Tag*)
- valodas prasmes līmenis
- dzimums
- pilsonība

The screenshot shows a search interface with the following fields and options:

- Token:** An empty text input field.
- Lemma:** An empty text input field.
- Tag:** A dropdown menu with 'n...6' selected.
- Error type:** A dropdown menu with 'Spelling Error' selected.
- Test level:** A dropdown menu with options 'A', 'B', and 'C'.
- Sublevel:** A dropdown menu with options '1' and '2'.
- Gender:** A dropdown menu with options 'Male' and 'Female'.
- Nationality:** A dropdown menu with options 'Amerika', 'Armēnija', 'ASV', 'Austrālija', and 'Apvienotā Karaliste'.

At the bottom of the interface are two buttons: 'Filter' (orange) and 'Reset' (grey).

Meklēšana korpusā (2/2)

Vaicājums: 6. deklinācijas lietvārdi, pareizrakstības kļūdas

Candidate 2132 (MALE, 42, Krievija) , 3. uzdevums, A līmenis

ar jums braukt ekskursijā . Ekskursija būs 14. maijā uz Vinspili vai Liepāju . Cena ekskursijas būs 60 eiro vai 75 eiro
ar Jums braukt ekskursijā . Ekskursija būs 14. maijā uz Ventspili vai Liepāju . Cena ekskursijas būs 60 eiro vai 75 eiro

Candidate 3112 (FEMALE, 41, Krievija) , 3. uzdevums, B līmenis

lesim kopā uz džeza mūzikas festivālu ? festivāls būs Cesu koncertzāli no 17.10 . līdz 19.11.2015 . Biļetes maksā 55 eiro
lesim kopā uz džeza mūzikas festivālu ? Festivāls būs Cēsu koncertzāli no 17.10 . līdz 19.11.2015 . Biļetes maksā 55 eiro

Candidate 9212 (MALE, 30, Latvija) , 2. uzdevums, C līmenis

Latvija ir tik bagata ar dzintaru , ka ne vienā Valsta . Turistiem ļoti interesanti atvest to , ko nav nevienam
Latvija ir tik bagāta ar dzintaru , kā ne neviena valsts . Tūristiem ļoti interesanti un atvest to , kā nav nevienam

Candidate 9261 (MALE, 25, Latvija) , 2. uzdevums, C līmenis

labs pamatēdiens . Labāk pamatēdienā dārzeņus aizvietot ar gaļu vai zivīm .
labs pamatēdiens . Labāk pamatēdienā dārzeņus aizvietot ar gaļu vai zivīm .

Candidate 9454 (MALE, 25, Krievija) , 3. uzdevums, C līmenis

es nezināju , ko tieši par mani domā šie ļaudi , ja viņi visi smaida . Tādējādi , teiciens par smaidu
es nezināju , ko tieši domā par mani domā šie ļaudis , ja viņi visi smaida . Tādējādi , teiciens par smaidu

Kļūdu analīze

Kļūdu tips	Skaitis	Kļūdainās teksta vienības (%)	Kļūdu izplatība korpusā (%)
Pareizrakstības kļūdas	14956	47,13%	10,48%
Formveidošanas kļūdas	8075	25,45%	5,66%
Interpunkcijas kļūdas	5857	18,46%	4,10%
Leksikas kļūdas	1756	5,53%	1,23%
Sintakses kļūdas (vārdu secība)	1703	5,37%	1,19%
Nesaprotams teksts	1546	4,87%	1,08%
Sintakses kļūdas	1321	4,16%	0,93%