



# Latviešu valodas teksta un runas korpusi vietnē *korpus.lv*

**Baiba Saulīte, Ilze Auziņa**

Praktiskais seminārs  
2022. gada 12. janvāris



## Filtrēt

teksta (16)

morfoloģija (13)

reprezentatīvs (8)

vispārīgs (8)

specializēts (5)

manuāli pārbaudīts (4)

sintakse (3)

runas (3)

apguvēju (3)

kļūdas (2)

semantika (1)

tīmekļa (1)

autora (1)

paralēls (1)

parlamentārs (1)

diahronisks (1)

## Populārākie korpusi

### LVK2018

Līdzsvarotais mūsdienu latviešu valodas tekstu korpus  
2016–2018, 10 milj. vārdlietojumu (12 milj. tekstvienību)

[Vairāk informācijas](#)

### MuLa

Mūsdienu latgaliešu tekstu korpus  
2011–2013, 1 milj. vārdlietojumu (1,3 milj. tekstvienību)

[Vairāk informācijas](#)

### LaVA

Latviešu valodas apguvēju korpus  
2018–2021, 192 000 vārdlietojumi (241 000  
tekstvienību)

[Vairāk informācijas](#)

### LVTB

Latviešu valodas sintaktiski marķētais korpus  
2010–2019, 13 643 teikumi (220 116 tekstvienību) (v2.5)

[Vairāk informācijas](#)

# Korpusa vizītkarte

LaVA

## Latviešu valodas apguvēju korpus

Korpusā iekļauti to Latvijas augstākajās mācību iestādes studējošo ārvalstnieku darbi, kuri latviešu valodu apgūst kā svešvalodu pirmo vai otro semestri. Teksti ir automātiski morfoloģiski marķēti, un tajos ir manuāli marķētas valodas apguvēju kļūdas.

### Publikācija, uz kuru atsaukties:

R. Dargis and I. Auzina and K. Levane-Petrova and I. Kaija

#### Quality Focused Approach to a Learner Corpus Development

*Proceedings of The 12th Language Resources and Evaluation Conference (LREC), 392-396, 2020*

[PDF](#)

teksta (16)

apguvēju (3)

morfoloģija (13)

kļūdas (2)

manuāli pārbaudīts (4)

Korpusa apjoms	192 000 vārdlietojumi (241 000 tekstvienību)
Izstrādes periods	2018–2021
Izstrādātājs	LU MII
Finansējuma avots	LZP Fundamentālo un lietišķo pētījumu programmas projekts Nr. Izp-2018/1-0527
Citas publikācijas	<p>I. Kaija and I. Auziņa <b>Data collection for learner corpus of Latvian: copyright and personal data protection</b> <i>Selected papers from the CLARIN Annual Conference 2019, 41-47, 2020</i> <a href="#">PDF</a> <a href="#">DOI</a></p> <p>I. Auzina and I. Kaija and K. Levane-Petrova <b>Mērķhipotēžu izvirzīšana latviešu valodas apguvēju korpusā</b> <i>Valoda: nozīme un forma, 11, 7-26, 2020</i> <a href="#">PDF</a> <a href="#">DOI</a></p>

# Korpusu iedalījums

**Dati:** teksta vs. runas

**Struktūra:** reprezentatīvs vs. nav reprezentatīvs

**Tips:** vispārīgs vs. specializēts

**Marķējums:** marķēts vs. nemarkēts

## Filtrēt

teksta (16)

morfoloģija (13)

reprezentatīvs (8)

vispārīgs (6)

specializēts (5)

manuāli pārbaudīts (4)

sintakse (3)

runas (3)

apguvēju (3)

kļūdas (2)

semantika (1)

tīmekļa (1)

autora (1)

paralēls (1)

parlamentārs (1)

diahronisks (1)

# Korpusu iedalījums: dati

**Teksta korpuss** – liels un strukturēts mašīnlasāmu tekstu kopums, kas paredzēts rakstītās valodas analīzei.

teksta (16)

LVK2018, LVK2013, LVTB, UDLV, FullStack-LV, LiLa, MuLa, LaVA, VVPP, Pārspriedumi, Saeima, Senie, Rainis, Emuāri, Barometrs, Tīmeklis 2007

**Runas korpuss** – audio ierakstu un to transkripciju kopums, kas paredzēts runātās valodas izpētei.

runas (3)

LVR100, LaRko, LAMBA, Subtitri

# Korpusu iedalījums: struktūra

**Reprezentatīvā** korpusā ir ietverti daudzveidīgi dati noteiktās proporcijās.

Vispārīgā un specializētā korpusa reprezentatīvitate tiek sasniegta un mērīta dažādos veidos.

reprezentatīvs (8)

LVK2018, LVK2013, LVTB, UDLV, FullStack-LV, LiLa, MuLa

LVR100

# Korpusu iedalījums: tips

**Vispārīgais korpuss** nav ierobežots tematiski vai pēc kādas citas pazīmes.

vispārīgs (6)

LVK2018, LVK2013, LVTB, UDLV, FullStack-LV, LRK100

**Specializētais korpuss** ir ierobežots tematiski, ģeogrāfiski, pēc runātāju vecuma, laikā vai pēc kādas citas pazīmes. Kopā 13 tādu korpusu.

specializēts (5)

LAMBA, LaRko, MuLa, Emuāri, Barometrs

Saeima, LaVa, VVPP, Pārspridumi, Senie, Tīmeklis 2007, LiLa, Rainis

parlamentārs (1)

apguvēju (3)

diahronisks (1)

tīmekļa (1)

paralēls (1)

autora (1)

# Korpusu iedalījums: marķējums

**morfoloģija (13)** – morfoloģiski marķēts

Emuāri, FullStack-LV, LAMBA, LaVA, LVK2013, LVK2018, LVTB, LRK100, VVPP, Pārspriedumi, Rainis, Saeima, Subtitri, Tīmeklis2007, UDLV

**sintakse (3)** – sintaktiski marķēts FullStack-LV, LVTB, UDLV,

**semantika (1)** – semantiski marķēts (*FrameNet* un *PropBank*) FullStack-LV

**kļūdas (2)** – marķētas valodas apguvēju kļūdas LaVA, VVPP

**manuāli pārbaudīts (4)** – manuāli labotas automātiskā marķējuma kļūdas

LVTB, UDLV, FullStack-LV, LaVA



# Piekļuve korpusiem

Konkrētā korpusa mājaslapā (korpuss vai plašāka informācija par to):

Barometrs	<a href="http://barometrs.korpuss.lv/">http://barometrs.korpuss.lv/</a>
LAMBA	<a href="runa.lamba.lv">runa.lamba.lv</a>
LaRko	<a href="larko.ailab.lv">larko.ailab.lv</a>
LaVA	<a href="lava.korpuss.lv/search">lava.korpuss.lv/search</a>
Saeima	<a href="saeima.korpuss.lv">saeima.korpuss.lv</a>
Senie	<a href="senie.korpuss.lv">senie.korpuss.lv</a>
LVTB, LVUD	<a href="http://sintakse.korpuss.lv/">http://sintakse.korpuss.lv/</a>
LRK100	<a href="http://runa.korpuss.lv/">http://runa.korpuss.lv/</a>

Daži korpusi **nav** publiski pieejami (**VVPP**) vai ir pieejami tikai daļēji (**LRK100, Subtitri**)

# Pieklūve korpusiem *NoSketchEngine* platformā

## ***SketchEngine***

- korpusu pārlūks un teksta analīzes rīks
- ar sarežģītiem vaicājumiem ļauj pētniekiem meklēt valodas parādības apjomīgās tekstu kolekcijās
- korpusi vairāk nekā 90 valodās

***NoSketchEngine*** – *SketchEngine* bezmaksas versija

## **AiLab korpusi (14)**

Emuāri, LaVA, LiLa, LVK2013, LRK100, LVK2018, Pārspriedumi, Rainis, Saeima, Subtitri, Tīmeklis2007, UDLV

MuLa, Senie

# Morfoloģisko pazīmju kopa



Latviešu valodas teksta un runas korpusi

Meklēšana

Morfoloģisko pazīmju kopa

Par korpusiem

Morfoloģiskajā marķēšanā izmantotā pazīmju kopa, kurā uzskaitītas vārdšķiras (11) un citas tekstvienību grupas (3), katrai no tām:

- morfoloģisko pazīmju un to apzīmējumu (tagu) kopums (1–13)
- katras pazīmes vērtība
- skaidrojums
- piemēri

Kur nepieciešams, pie vārdšķiras dotas īsas piezīmes, piemēram, par pamatformas izvēli vai par atšķirībām no tradicionālās latviešu valodas gramatikas.

# Meklēšana vairākos korpusos vienlaikus



Latviešu valodas teksta un runas korpusi

Meklēšana

Morfoloģisko pazīmju kopa Par korpusiem

??kūkot|

?

Meklēt

## Meklēšanas instrukcija

Ierakstot pamatformu, tiks atrastas visas vārdformas:

*iet* → *ēju, iet, ejam, gāja* utt.

Ierakstot konkrētu vārdformu, tiks meklēti lietojumi tikai ar šo vārdformu:

*gāja* → *gāja, ziemā* → *ziemā, kūku* → *kūku* u. c.

Ierakstot divu vai vairāku vārdu pamatformas, tiks atrastas konkordances ar šiem vārdiem dažādos locījumos:

*lasīt grāmata* → *lasu grāmatu, jālasa grāmata, jālasa grāmatu, lasīju grāmatu* u. c.

Vaicājumā var izmantot zvaigznīti (\*) nenoteikta rakstzīmju skaita norādīšanai, jautājuma zīmi (?) tieši vienas rakstzīmes apzīmēšanai, vertikālu līniju (|) vairāku vārdu vai frāžu meklēšanai.

## Vaicājumu piemēri

Vaicājums	Atrastie lietojumi	Komentāri
inic*	<i>iniciators, iniciatīva, iniciēt, iniciālis</i> utt.	Vārdi, kas sākas ar <i>inic-</i> .
*šana	<i>dibināšana, piekrišana, ražošanas, pagatavošanai, zināšanu</i> utt.	Vārdi, kas beidzas ar <i>-šana</i> , un to formas.
??kūkot	<i>izkūkot, pakūkot, nokūkot, iekūkot</i> utt.	Verbs <i>kūkot</i> ar divburtu priedēkli dažādos locījumos.
f finansu   finanšu	<i>finansu, finanšu</i>	Atrod abas vārdformas.

# Meklēšana vairākos korpusos: piemērs



Latviešu valodas teksta un runas korpusi

Meklēšana

Morfoloģisko pazīmju kopa

Par korpusiem

??kūkot

?

Meklēt

0 lietojumi

(0 uz miljonu)

**Saeima**

LR 5.–12. Saeimas sēžu stenogrammu korpuss



7 lietojumi

(0,84 uz miljonu)

**Emuāri**

Latviešu valodas emuāru korpuss



8 lietojumi

(1,43 uz miljonu)

**LVK2013**

Līdzsvarotais mūsdienu latviešu valodas tekstu korpuss



0 lietojumi

(0 uz miljonu)

**LaVA**

Latviešu valodas apguvēju korpuss



18 lietojumi

(1,46 uz miljonu)

**LVK2018**

Līdzsvarotais mūsdienu latviešu valodas tekstu korpuss



2 lietojumi

(0,87 uz miljonu)

**Rainis**

Raiņa darbu korpuss



0 lietojumi

(0 uz miljonu)

**Pārspridumi**

Skolēnu pārspridumu korpuss



45 lietojumi

(0,36 uz miljonu)

**Tīmeklis 2007**

Latviešu valodas tīmekļa korpuss



# Meklēšana vairākos korpusos: piemērs

## CONCORDANCE



simple ??kūkot 7 (0.84 per million)



KWIC



Details

Left context

KWIC

Right context

1	<input type="checkbox"/>		doc#1486	jad . No otras puses uzrodas kādi pavisam prātu	<b>izkūkojuši</b>	ekstrēmisti , un cenšas šo pasivitāti atsvērt ar mo	
2	<input type="checkbox"/>		doc#7235	ī kompānijā : ) Bija arī dzeguze . Tukšas kabatas	<b>iekūkoja</b>	. Toties zaļumiem pilnu vēderu : ) Nemaz nezinu ,	
3	<input type="checkbox"/>		doc#7237	askatījāties un klusējot vienojāties , ka Vilnis ir	<b>izkūkojis</b>	prātu . Publiski atzīstos , ka bija pagājušas tikai 2	
4	<input type="checkbox"/>		doc#9072	, pielaikot Martai kleitu ( dejām ) , tad nesteidzīgi	<b>pakūkojām</b>	kafejnīcā , tālāk mūsu ceļš veda uz rokdarbu veik	
5	<input type="checkbox"/>		doc#14262	Slaista skolotāja nāves ir vainojams kāds prātiņu	<b>izkūkojis</b>	burvis vārdā Ištvars , kas nolēmis nogalināt visus	
6	<input type="checkbox"/>		doc#18767	un pietiek . Lidostā atgriezāties laicīgi , un jauki	<b>nokūkojām</b>	atlikušas stundas gaidot savu reisu . Arī uz Atēnē	
7	<input type="checkbox"/>		doc#19053	us lokus ap dzelzsceļa staciju , un pēdējo stundu	<b>nokūkoju</b>	stacijā , gaidot svinīgo brīdi iekāpšanai vilcienā .	

# korpuss.lv nākotne

Paplašināt meklēšanu visos korpusos – pievienot meklēšanu arī citu iestāžu korpusos.

Pilnveidot esošos korpusus un pievienot jaunus.

Regulāri ievietot jaunākās versijas.

Pēc iespējas pievienot meklēšanas instrukcijas.



Ieskaties: <http://ailab.lv>



Seko mums: @AiLab.lv

Jautājumi, ieteikumi: [korpuss@ailab.lv](mailto:korpuss@ailab.lv)

# Projekti



Nr. VPP-IZM-DH-2020/1-0001  
Nr. VPP-Letonika-2021/1-0006



Nr. lzp-2018/1-0527 (2018-2021)  
Nr. lzp-2019/1-0464 (2020-2022)



Nr. VPP-IZM-2018/2-0002 (2018-2021)



IEGULDĪJUMS TAVĀ NĀKOTNĒ

