Latvian National Corpora Collection – Korpuss.lv

Baiba Saulīte, Roberts Darģis, Normunds Grūzītis, Ilze Auziņa, Kristīne Levāne-Petrova, Lauma Pretkalniņa, Laura Rituma, Pēteris Paikens, Artūrs Znotiņš, Laine Strankale, Kristīne Pokratniece, Ilmārs Poikāns, Guntis Bārzdiņš, Inguna Skadiņa, Anda Baklāne, Valdis Saulespurēns, Jānis Ziediņš

> Institute of Mathematics and Computer Science, University of Latvia National Library of Latvia Culture Information Systems Centre

roberts.dargis@lumii.lv, anda.baklane@lnb.lv, janis.ziedins@kis.gov.lv

Abstract. Latvian National Corpora Collection (LNCC) is a diverse collection of corpora representing both written and spoken language. LNCC covers various use cases and all the important text types and genres. It is a continuous multi-institutional and multi-project effort, supported by the Digital Humanities and Language Technology communities in Latvia. All corpora of LNCC are annotated with a uniform morpho-syntactic annotation scheme which enables federated search and consistent linguistics analysis in all the LNCC corpora and allows to select and mix various corpora for pre-training large Latvian language models like BERT and GPT.

(1) The current compilation of LNCC

Code name	Full name	Size	Туре	Release	
Written langua	ge text corpora				
	Balanced Corpus of				
LVK2018	Modern Latvian	12M tokens	text, general, representative	2016–2018	
	(Levane-Petrova and Dargis, 2018)				
	Latvian UD Treebank,	266k tokens	text general representative		
UDLV-LVTB	subset of LVK2018,	200K tokens	manually annotated	2015–2021	
~	part of UD v2.9 (Zeman et al., 2021)	(TOK Selit.)	manually annotated		
Hugo.lv	Hugo.lv Parallel Corpora	10.5M tokens	text, general, culture	2018	
Tīmeklis2007	Latvian Web Corpus	123 5M tokens	text web	2006_2007	
T IIIICKII32007	(Džeriņš and Džonsons, 2007)	125.5101 tokens		2000-2007	
Tīmeklis2020	Latvian Web Corpus	492.6M tokens	text, web	2020-2022	
Vikipēdija	Latvian Wikipedia	27.7M tokens	text, specialised	2022	
Emuāri	Latvian Blog Corpus	8M tokens	text, specialised	2014–2015	
Barometrs	Corpus of News Portal Comments	447.3M tokens	text, specialised	2011-2021	
	Corpus of Latvian				
Saeima	Parliament Debates	24M tokens	text, specialised:parliamentary	2013–2018	
	(Auziņa et al., 2018)				
Likumi	Corpus of Legal Acts	116 2M tokens	text specialised	2022	
Likuini	of the Republic of Latvia	110.21vi tokelis	text, specialised	2022	
	Lithuanian-Latvian-Lithuanian				
LiLa	Parallel Text Corpus	5.7M tokens	text, parallel, representative	2011–2013	
	(Utka et al., 2013)				
	Corpus of Contemporary		text specialised dialect		
MuLa	Latgalian Texts	1.3M tokens	representative	2011-2013	
	(Sperga et al., 2013)		representative		
	Latvian Language Learner Cornus		text, specialised:learner,		
LaVA	(Auzina et al. 2021)	241k tokens	manually annotated,	2018-2021	
	(Auziția et al., 2021)		error annotation		
Dārspriedumi	Corpus of Students Essays	226k tokons	text specialized	2018 2021	
Faisprieduini	(Levāne-Petrova et al., 2021)	220K IOKEIIS	text, specialised	2010-2021	
Disertācijas	Corpus of Latvian PhD Theses	23.4M tokens	text, specialised	2022	
LatSenRom	Corpus of Latvian Early Novels	3.3M tokens	text, specialised:literary	2019–2021	
Painis	Corpus of Texts Written by Rainis	2 3M tokens	text specialised literary	2018	
Kaillis	(Spektors et al., 2018)		text, specialised.itterary	2018	
	Corpus of Early Written				
Senie	Latvian Texts	2.7M tokens	text, specialised:diachronic	2002-2021	
	(Andronova et al., 2002)				
Spoken languag	ge text corpora				
I DK2012	Latvian Speech Paccamition Comus	975k tokens	snoken general rangeantative	2012	
LKK2015	Latvian Speech Recognition Corpus	(100 hours)	spoken, general, representative	2013	
Subtitri	Latvian Subtitles of	10.8M tokens		2020 2022	
Subuut	Public Broadcasting	(1200 hours)	spoken, specialised	2020-2022	
IVMED	Latvian Medical Speech Comme	157k tokens	spoken specialized	2022	
	Latvian Medical Speech Corpus	(35 hours)	spoken, specialised	2022	
LNCC:	21 corpora	1.3B tokens		Jun 2022	

(3) Federated search result for the query sirds* ('heart*')

sirds*			Search
5230 occurrences (1541 per million)	LatSenRom Corpus of Latvian Early Novels	565 occurrences (426 per million)	MuLa Corpus of Contemporary Latgalian Texts
1734 occurrences (305 per million)	LiLa Lithuanian-Latvian-Lithuanian Parallel Text Corpus	9 occurrences (37 per million)	LaVA Latvian Language Learner Corpus
2128 occurrences (256 per million)	Emuāri Latvian Blog Corpus 2015	157 occurrences (137 per million)	LRK2013 Latvian Speech Recognition Corpus
3321 occurrence (270 per million)	LVK2018 The Balanced Corpus of Modern Latvian	510 occurrences (2257 per million)	Pārspriedumi Corpus of Students' Essays
13 occurrences (5 per million)	Senie Corpus of Early Written Latvian Texts	2752 occurrences (1196 per million)	Rainis Corpus of Texts Written by Rainis
2922 occurrences (121 per million)	Saeima Corpus of the Saeima (the Parliament of Latvia)	1230 occurrences (114 per million)	Subtitri Latvian Subtitles of Public Broadcasting
19 970 occurrences (162 per million)	Tīmeklis2007 Latvian Web Corpus 2007	164 504 occurrences (334 per million)	Tīmeklis2020 CommonCrawl of Latvian 2020

(4.1) Common tagset: morphology

pateicas	Dievam	,	grozās	,	sirdās	par	zeņķiem	,	ו	(as	spļa	t uj
0an/pateikties	ncmsd1/dievs	zc/, v	/myip_330an/grozīties	zc/,	vmyip_330an/sirdīties	sppd/par	ncmpd2/zeņķ	is za	⊳/, ן	pr000nn/kas	vmni	pi130an/spļaut
	S , t zc/ pr000p	kas	gadījās	, 70/	sirdīgi	pārmet	a : Dan/pārmest z	- /o/	_ zd/-	Tu pp20spp/tu	, 70/	Mairita

(2) The Korpuss.lv website

Latvian National Corpora Collection Search Materials - About Korpuss.lv

EN 🤊



\$\vec{2}{2}\$cnotieksirdsdarbībasapstāšanās\$\vec{2}{2}\$cvmnipi130an/notiktncfsg4/sirdsdarbībancfsgr/apstāšanās	Sirds ncfsn6/sirds	darbojas , pateicoties tam vmyip_230an/darboties zc/, vmypu0000000n/pateikties pd3msd
Šamilgibijaapslēpts?ı/patsrp_/ilgivmnis_i30an/būtvmnpdmsnpsnpn/apslēptzs/?	Sirdsbalss	gan sacīja priekšā , bet Dāvis cc/gan vmnist330an/sacīt r0_/priekšā zc/, cc/bet npmsn2/Dāv
ensībās , kuras ir aktiera	sirdslieta	. Patriks piekritis sievas not
'sacensība zc/, pr0fpnn/kura vcnipii30an/būt ncmsg2/aktieris	ncfsn4/sirdslieta	zs/. npmsn1/Patriks vmnpdmsnasnpn/piekrist ncfsg4/sieva ncm
eet-heart (" Ar labu nakti ,	sirdsmīļā	"), – vai šoreiz mājās vajadzēs
weet-heart zb/(zq/" spsa/ar affsanp/labs ncfsa6/nakts zc/,	affsnyp/sirdsmīļš	zq/" zb/) zc/, zd/- cc/vai r0_/šoreiz ncfpl4/māja vmnift330an/\
acis . — Visuspēcīgais Dievs ,	sirdsžēlīgais	Tēvs , Nu es zinu , ka
ncfpa6/acs zs/. zd/- afmsnyp/visuspēcīgs ncmsn1/dievs zc/,	afmsny_/sirdsžēlīgs	ncmsn1/tēvs zc/, q/nu pp10snn/es vmnipt31san/zināt zc/, cs/

(4.2) Common tagset: Universal Dependencies

pateicas Dievam , grozās ,	sirdās	par zeņķiem , kas spļauj sau
VERB/root NOUN/iobj PUNCT/punct VERB/root PUNCT/punct	VERB/root	ADP/case NOUN/iobj PUNCT/punct PRON/nsubj VERB/acl NOU
ozdams	sirdīgi	pārmeta: – Tu ,
ERB/acl PUNCT/punct PRON/nsubj VERB/acl PUNCT/punct	ADV/advmod	VERB/root PUNCT/punct PUNCT/punct PRON/vocative PUNCT/
. Kāpēc notiek sirdsdarbības apstāšanās	Sirds	darbojas, pateicoties tam ,
unct ADV/advmod VERB/root NOUN/nmod NOUN/nsubj	NOUN/nsubj	VERB/root PUNCT/punct VERB/advcl PRON/iobj PUNCT/punct
pašam ilgi bija apslēpts ?	Sirdsbalss	gan sacīja priekšā , bet D
RON/iobj ADV/advmod VERB/aux VERB/ccomp PUNCT/punct	NOUN/nsubj	CCONJ/cc VERB/root ADV/advmod PUNCT/punct CCONJ/cc Pl
Icensībās, kurasiraktieraNOUN/oblPUNCT/punctPRON/nsubjAUX/aclNOUN/nmod	sirdslieta NOUN/nsubj	. Patriks piekritis sievas note PUNCT/punct PROPN/nsubj:pass VERB/root NOUN/nmod NC
" Ar labu nakti ,	sirdsmīļā	"), – vai
CT/punct ADP/case ADJ/amod NOUN/discourse PUNCT/punct	ADJ/conj	PUNCT/punct PUNCT/punct PUNCT/punct CCON
. — Visuspēcīgais Dievs , nct PUNCT/root ADJ/amod NOUN/nsubj PUNCT/punct	sirdsžēlīgais ADJ/amod	TēvsNueszinuNOUN/aclPUNCT/punctPART/discoursePRON/nsubjVERB/roc

Benefits:

LVK2018

The Balanced Corpus of Modern Latvian 2016–2018, 10M words (12M tokens) Developers: IMCS UL

More
Q Search

LaVA

Latvian Language Learner Corpus 2018–2021, 192k words (241k tokens) Developers: IMCS UL

More Q Search



Latvian Language Learner Corpus

PDF DOI

The corpus includes more than 1000 texts created by foreign Latvian language learners studying at Latvian higher education institutions for the first or second semester. The morphologically annotated texts have been checked manually; the language learners' errors have been manually annotated.

Publication to be cited:

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

text (17) learner (3) morphology (16) error annotation (2) manually annotated (4)

Corpus size	192k words (241k tokens)
Development period	2018–2021
Developers	Institute of Mathematics and Computer Science UL
Funding	The project "Development of Learner corpus of Latvian: methods, tools and applications" (project No. lzp-2018/1-0527) is financed by Latvian Council of Science.
Homepage	http://lava.korpuss.lv/lv/
CLARIN	http://hdl.handle.net/20.500.12574/49
Other publicatoins	I. Kaija and I. Auzina Data collection for learner corpus of Latvian: copyright and personal data protection Selected papers from the CLARIN Annual Conference 2019, 41-47, 2020

- An open-ended national reference corpus
- Easily accessible otherwise separate and diverse corpora
- Quick quantitative comparison across the whole corpora collection

Future tasks:

- Full-fledged Latvian speech corpora in LNCC
- A new, extended (100M tokens) version of LVK
- Other endpoints and selected corpora
- 10B token corpora collection



http://www.korpuss.lv

This work has received financial support from the State Research Programmes under grant agreements No. VPP-IZM-DH-2020/1-0001 and No. VPP-LETONIKA-2021/1-0006, from the Latvian Language Agency under grant agreement No. 4.6/2019-029, and from ERDF under the grant agreement No. 1.1.1.1/18/A/153.