

The Balanced Corpus of Modern Latvian (LVK)

Kristīne Levāne-Petrova
kristine.levanepetrova@lumii.lv

The Institute of Mathematics and Computer Science, University of Latvia

Design principles of LVK

LVK has been designed since 2007 based on the on the Latvian Language Corpus Conception (2005).

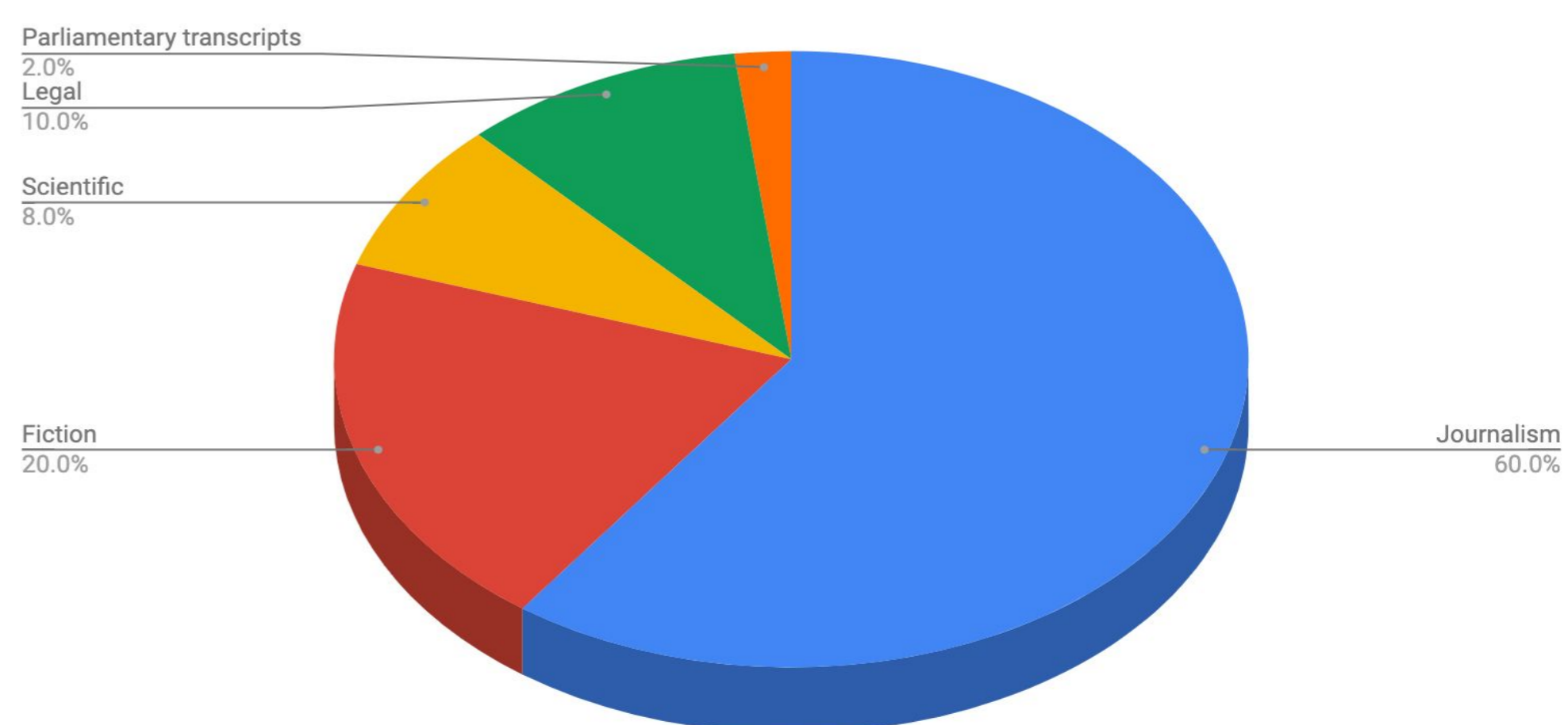
LVK2018 is designed as a general-language, representative and publicly available corpus. It is a monolingual, fully morphologically and partly syntactically and semantically annotated corpus. Presently, it consists of 10 million tokens. Characteristics of LVK2018:

- **General** – the corpus includes sources from different domains, styles, genres, etc.
- **Balanced** – the corpus that aims to cover the variety of existing texts in estimated proportions.
- The corpus represents the **synchronic state** of the language.
- **Originality** – the corpus should only contain texts originally written in Latvian.
- The corpus is **representative**, it contains texts from all language styles, major domains and many subdomains.

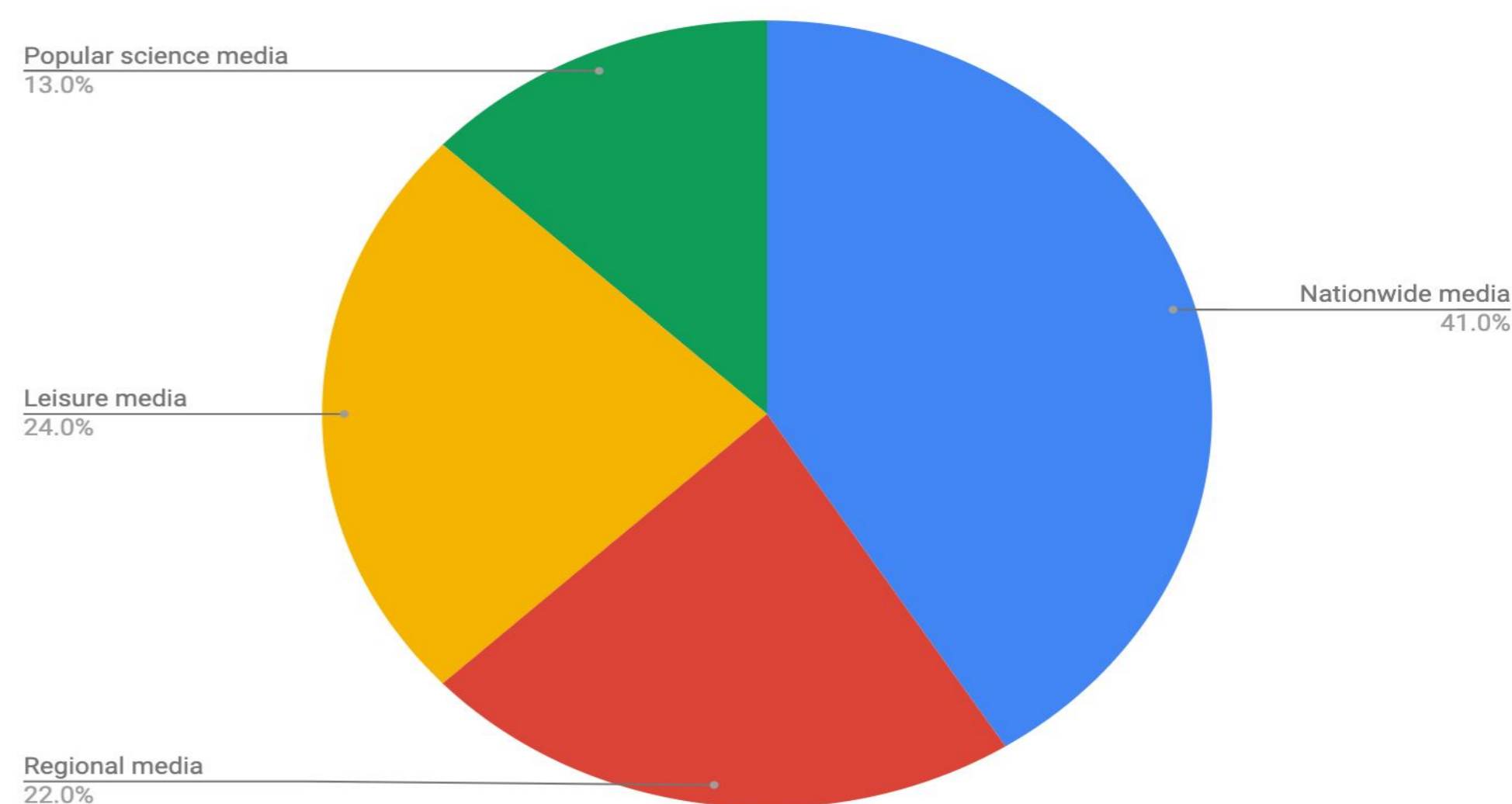
Text selection criteria

- **Time** – the corpus should contain texts created and published after 1991.
- The corpus should contain **full-text**.
- **Diversity** – texts should cover as wide range of topics as possible. The sample cannot exceed more than 5 % and 50,000 words of the particular section of the corpus to not dominate one author or domain in the corpus. If the text is longer than the limit, it should be cut from the end of the sample only.
- **Uniqueness** – the corpus sample should be represented in corpus just once.
- **Quality** – samples should only contain clean text written in literary language with appropriate usage of diacritics and punctuation in Latvian. Tables and other nontext parts should be removed.

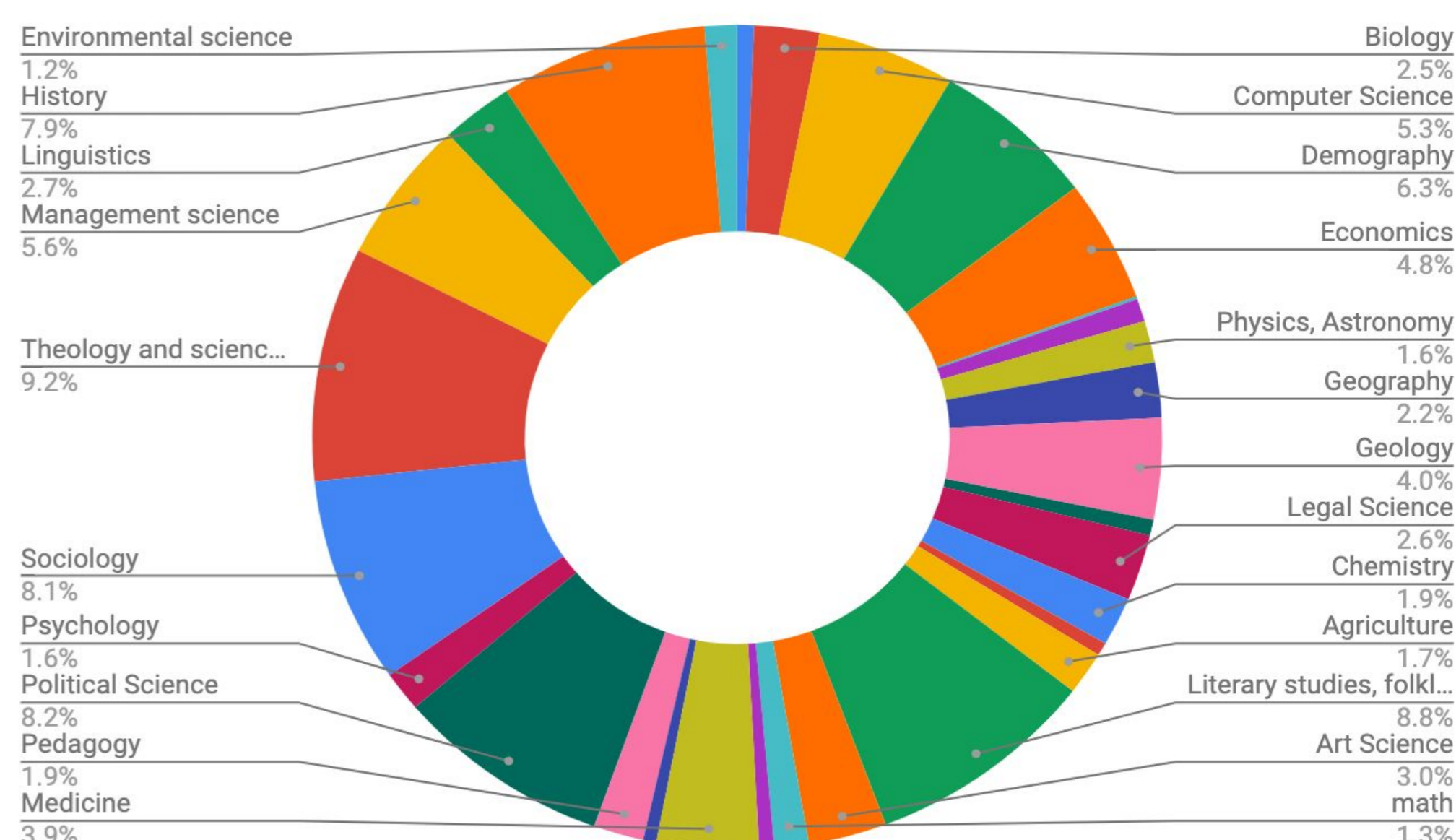
Composition of LVK2018



Composition of Journalism



Composition of Scientific texts



Latvian National Corpora Collection Search Materials About Korpuss.lv EN

general (6) web (2) learner (3) literary (2) parallel (1) parliamentary (1) diachronic (1)
specialised (6) representative (7) text (18) speech (4) morphology (17) syntax (3) semantics (1)
error annotation (2) manually annotated (4)

LVK2018
The Balanced Corpus of Modern Latvian
2016–2018, 10M words (12M tokens)
Developers: IMCS UL

[More](#)

MuLa
Corpus of Contemporary Latgalian Texts
2011–2013, 1M words (1.3M tokens)
Developers: IMCS UL, RAT

[More](#)

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

[More](#)

LVTB
Latvian Treebank
2010–2021, 15984 sentences (265 722 tokens)
Developers: IMCS UL

[More](#)

All LVK corpora and aubcorpora have been released in the framework of Latvian National Corpus and registered at the CLARIN-LV repository. LVK2018 and other corpora are freely available via the corpus query interface NoSketch Engine.



LVK2022

The Latvian National Corpora Collection is available:

<http://www.korpuss.lv>

This work has received support from the "Strengthening of the capacity of doctoral studies at the University of Latvia within the framework of the new doctoral model", identification No. 8.2.2.0/20/I/006