

Eesthetic: A Paralex Lexicon of Estonian Paradigms

LREC-COLING 2024

Sacha
Beniamine



Mari
Aigro



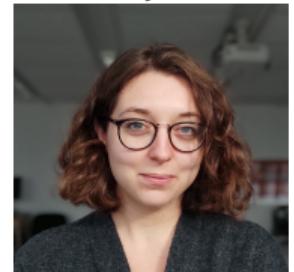
Matthew
Baerman



Jules
Bouton



Maria
Copot



Eesthetic

- a **large lexicon** of Estonian inflection

	lexemes	cells	forms
Nouns	5475	28	193 931
Verbs	5076	51	258 954
total	10551	79	452 885

- **standardized** in the Paralex format (csv)
- with rich **linguistic information**
- for **quantitative & qualitative** morphology

Estonian paradigms

case / number	SG	PL
NOM	sõna	sõnad
PART	sõna	sõnu / sõnasid
GEN	sõna	sõnade
ILL	sõnna / sõnasse	sõnusse / sõnadesse
INESS	sõnas	sõnus / sõnades
ELAT	sõnast	sõnust / sõnadest
ALL	sõnale	sõnule / sõnadele
AD	sõnal	sõnul / sõnadel
ABL	sõnalt	sõnult / sõnadelt
TRANS	sõnaks	sõnuks / sõnadeks
TERM	sõnani	sõnadeni
ESS	sõnana	sõnadena
ABESS	sõnata	sõnadeta
COM	sõnaga	sõnadega

sõna, N., 'word', class 17u

Estonian paradigms

case / number	SG	PL
NOM	sõna	sõnad
PART	sõna	sõnu / sõnasid
GEN	sõna	sõnade
ILL	sõnna / sõnasse	sõnusse / sõnadesse
INESS	sõnas	sõnus / sõnades
ELAT	sõnast	sõnust / sõnadest
ALL	sõnale	sõnule / sõnadele
AD	sõnal	sõnul / sõnadel
ABL	sõnalt	sõnult / sõnadelt
TRANS	sõnaks	sõnuks / sõnadeks
TERM	sõnani	sõnadeni
ESS	sõnana	sõnadena
ABESS	sõnata	sõnadeta
COM	sõnaga	sõnadega

sõna, N., 'word', **class 17u**

case / number	SG	PL
NOM	aas	aasad
PART	aasa	aasu / aasasid
GEN	aasa	aasade
ILL	aasa / aasasse	aasusse / aasadesse
INESS	aasas	aasus / aasades
ELAT	aasast	aasust / aasadest
ALL	aasale	aasule / aasadele
AD	aasal	aasul / aasadel
ABL	aasalt	aasult / aasadelt
TRANS	aasaks	aasuks / aasadeks
TERM	aasani	aasadeni
ESS	aasana	aasadena
ABESS	aasata	aasadeta
COM	aasaga	aasadega

aas, N., 'loop', **class 22u**

Estonian paradigms

Quantity 1

Quantity 2

Quantity 3

case / number	SG	PL
NOM	aas /'a:s/	aasad
PART	aasa /'a:sa/	aasu / aasasid
GEN	aasa /a:sa/	aasade
ILL	aasa /'a:sa/	aasusse / aasadesse
aas, N., 'loop', class 22u		

case / number	SG	PL
NOM	sõna /s'yna/	sõnad
PART	sõna /s'yna/	sõnu / sõnasid
GEN	sõna /s'yna/	sõnade
ILL	sõnnna /s'ynna/	sõnusse / sõnadesse
sõna, N., 'word', class 17u		

Estonian paradigms

case / number	SG	PL
NOM	sõna	sõnad
PART	sõna	sõnu / sõnasid
GEN	sõna	sõnade
ILL	sõnna / sõnasse	sõnusse / sõnadesse
INESS	sõnas	sõnus / sõnades
ELAT	sõnast	sõnust / sõnadest
ALL	sõnale	sõnule / sõnadele
AD	sõnal	sõnul / sõnadel
ABL	sõnalt	sõnult / sõnadelt
TRANS	sõnaks	sõnuks / sõnadeks
TERM	sõnani	sõnadeni
ESS	sõnana	sõnadena
ABESS	sõnata	sõnadeta
COM	sõnaga	sõnadega

sõna, N., 'word', class 17u

Lexicon creation (I)

The screenshot shows a web browser window for the Ekilex dictionary. The URL in the address bar is <https://sonaveeb.ee/dlall/dsall/sõna/1>. The page header includes the Sõnaveeb logo (a blue speech mark icon), the text "Sõnaveeb Eesti Keele Instituut", the word "Ekilex", and a search bar containing "sõna". There are also icons for keyboard, microphone, and search. A language dropdown shows "EST". Below the header, there are links for "See on ka vorm sõnast sõnad sõnama", "Keel Kõik keeled", "Sõnakogud Kõik sõnakogud", and "Ehk mul veab".

et sõna nimisõna 18.03.2024

EKI ÜHENDSÖNASTIK 2024

1 et iseseisva tähindusega keeleüksus, mida kasutatakse oma mõtete väljendamiseks tekstis ja kõnes (nt maailm, alla kirjutama, Rootsli laud)

Sünonüümid lekseem, sõnake

fr mot

Sõnavormid

	Muuttüüp	17u
sõna	🔊	sõn
sõna	🔊	sõn
sõna	🔊	sõn

[Näita tabelina](#)

Lexicon creation (I)

The screenshot shows a web browser displaying the META-SHARE platform. The URL in the address bar is <https://metashare.ut.ee/repository/browse/estonian-national-corpus-2021-vert/4547c7bea0>. The page title is "Estonian National Corpus 2021 - vert". A blue sidebar on the left contains the text "Estonian national corpus". The main content area includes a "View resource name in all available languages" link, a citation "Cite as: Koppel, Kristina; Kallas, Jelena (2022). Eesti keele ühendkorpus 2021. DOI: 10.15155/3-00-0000-0000-0000-08E60L", and a note about subcorpora. There are also "Read More" and "View resource description in all available languages" links. The bottom navigation bar includes "Browse Resources", "Community", "Statistics", "Help", and "About" buttons.

Estonian National Corpus 2021 - vert

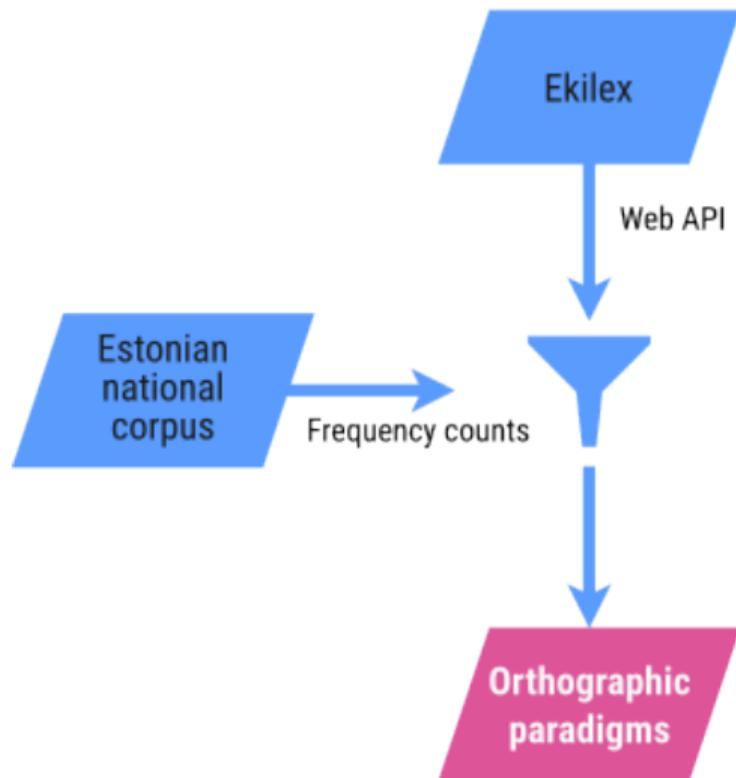
View resource name in all available languages

Cite as: Koppel, Kristina; Kallas, Jelena (2022). Eesti keele ühendkorpus 2021. DOI: 10.15155/3-00-0000-0000-0000-08E60L

Subcorpora are Reference Corpus, incl. Balanced Corpus; Estonian Web 2013, 2017, 2019, 2021; Wikipedia 2021, Wikipedia Talk 20
Estonian Feeds 2014–2021; Literature.
[... Read More](#)

View resource description in all available languages

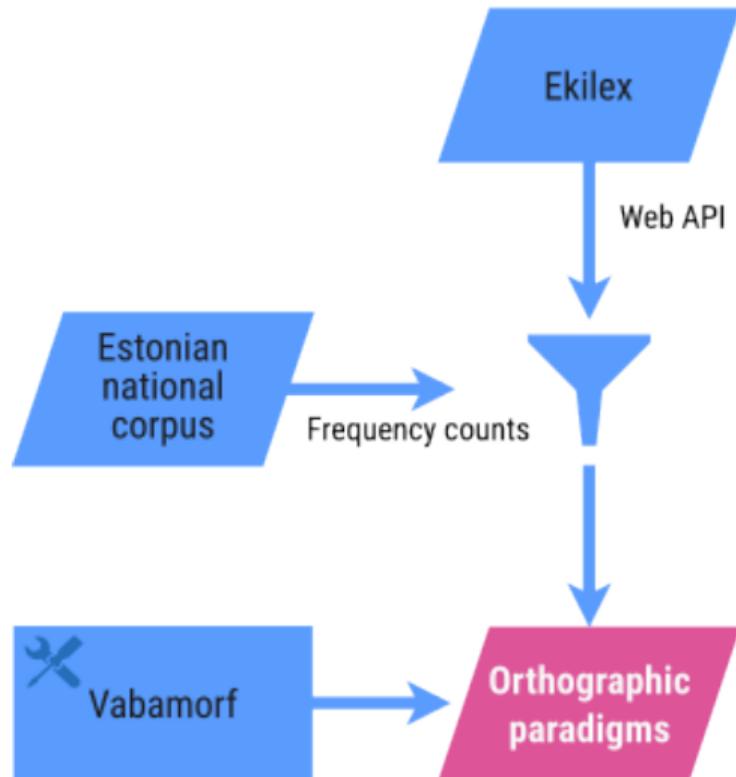
Lexicon creation (I)



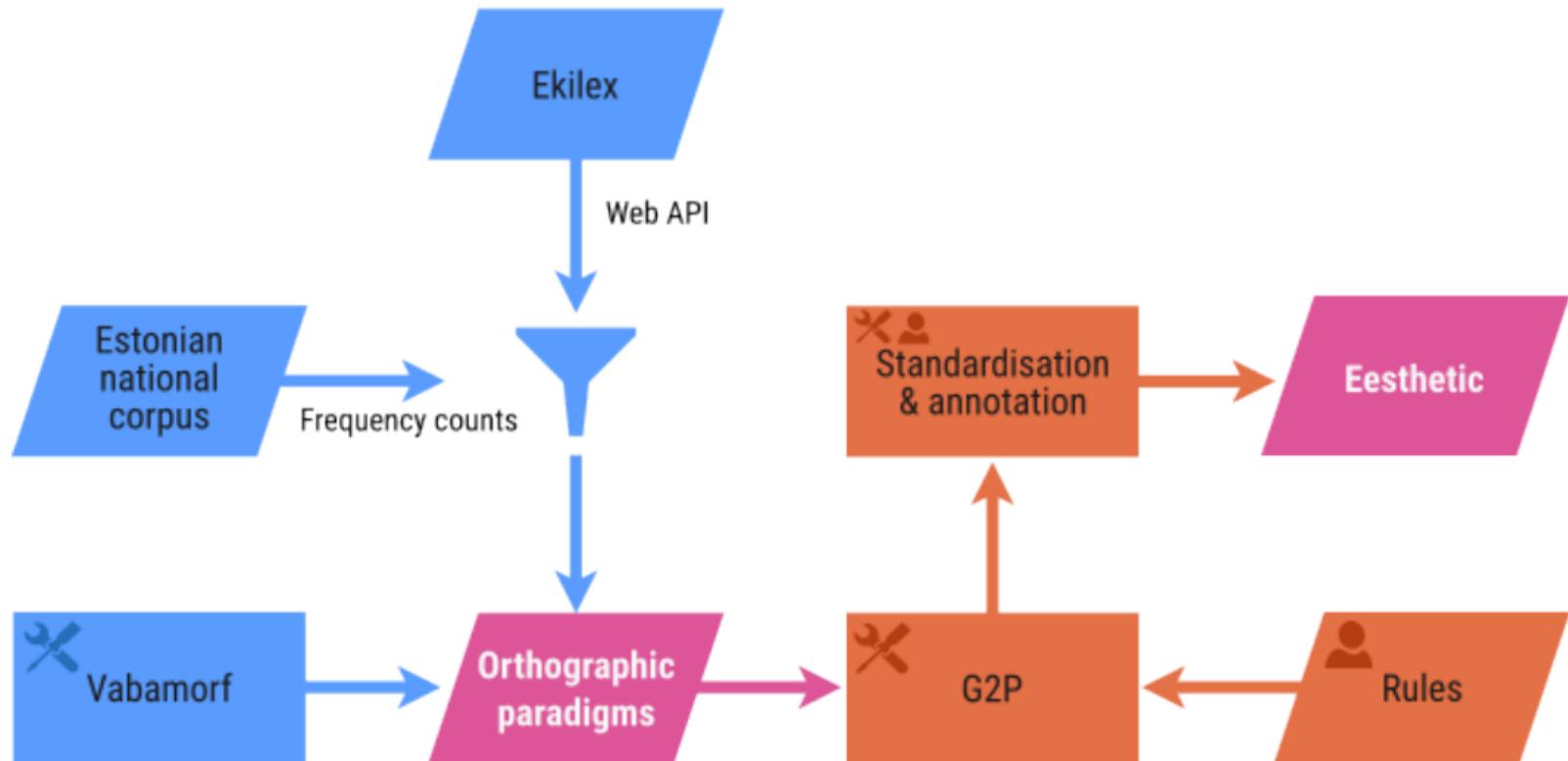
Annotated orthography

- | | | | |
|-----|-------------|-----------------|----------------|
| (1) | <`oman`ik> | ['oman'ik:] | 'owner' NOM.SG |
| (2) | <sõna[sse]> | [sõnas:e] | 'word' ILL.SG |
| (3) | <h`oo+`aeg> | [h'o::'æ:k] | 'season' |
| (4) | <k`ot't> | [k'ot::] | 'bag' |
| (5) | aasa | | 'loop' |
| | a. | <aasa> [a:sa] | GEN.SG |
| | b. | <`aasa> [a::sa] | PART.SG |

Lexicon creation (II)



Lexicon creation (II)



Grapheme to phoneme

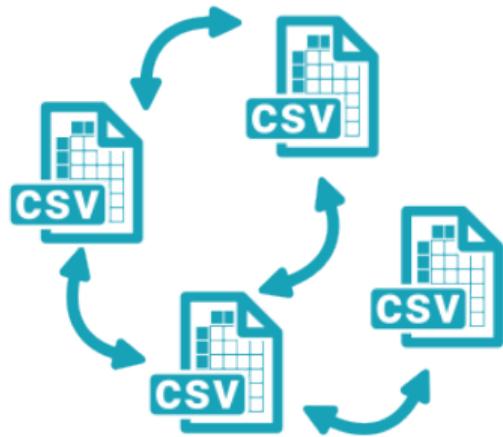
Custom rules for Epitran (Mortensen et al., 2018)

n -> η / _ [kg]

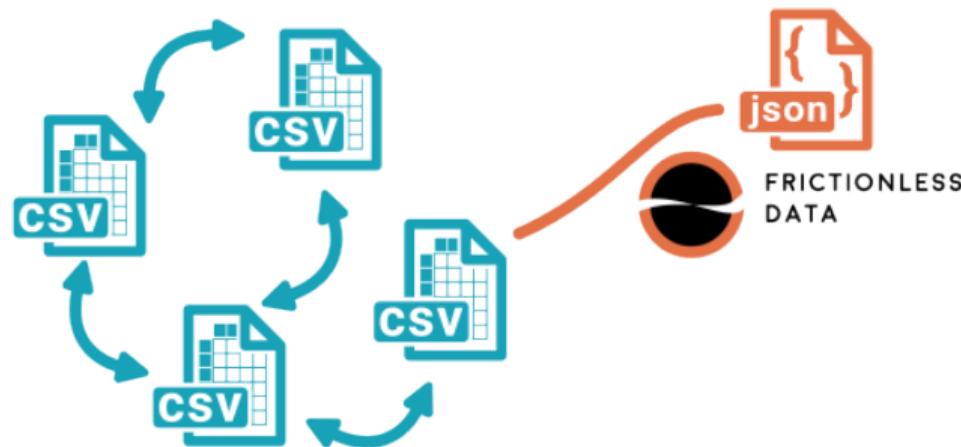
- Pre-processing rules:
 - Transcribe orthographic conventions (quantity, length, palatalization)
 - Apply regular phonological rules
- Non-contextual mapping
- Post-processing rules:
 - Add ligatures to diphthongs
 - Identify which sounds to lengthen under Q3
 - Cleanup typography
- Adjustments



DATA



DATA METADATA



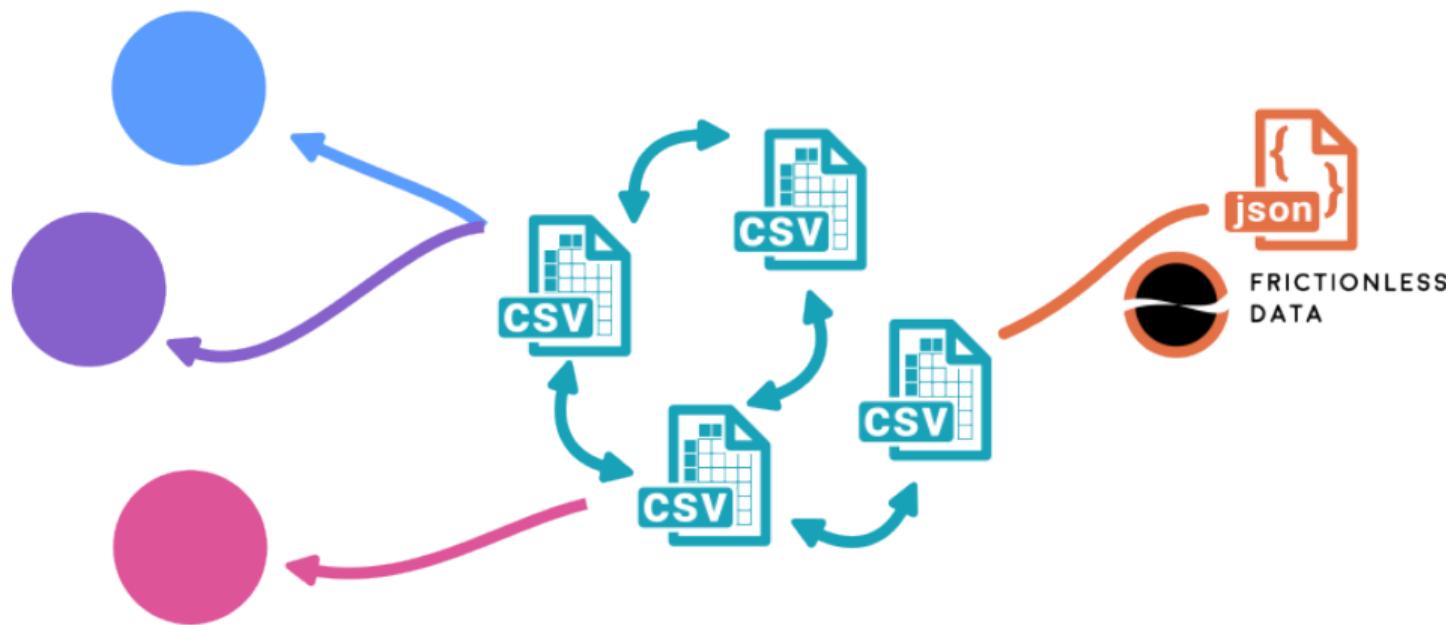


Paralex standard

LINKED

DATA

METADATA



Paralex standard

The screenshot shows the header of the Paralex standard website. It features a dark blue header bar with the Paralex logo (a speech bubble icon) and the text "Paralex". To the right is a search bar with a magnifying glass icon and the word "Search". Further right is a user profile section showing "sbeniamine/paralex", a star icon, "0", and "1".

Paralex

[Home](#)

Standard

Specs

Ontology

Background

Principles

What are metadata ?

Why the long form ?

How to

Link sources

Structure paradigms

Use comments and tags

Express sounds

Tutorial

FAQ

Datasets

Paralex: lexicons of morphological paradigms

Table of contents

1 Contributors

2 Version

Paralex is a standard for morphological lexicons which document inflectional paradigms.

It strives to provide data which is FAIR, so it can be used automatically, CARE, so it respects and empowers language communities, and DeAR (our own set of principles), so we can create a virtuous data ecosystem. It was inspired by the [Cross-Linguistic Data Formats \(CLDF\)](#) standard, and adheres to a similar philosophy and the same design principles. We aim to keep the two standards compatible in order to facilitate inter-operability.

A **paralex** lexicon is a set of tables written as [comma separated value \(csv\)](#) files. It follows a relational model, tables are written in [long form](#), [metadata](#) is written using the [frictionless standard](#), and the tables respect [pre-defined conventions](#). An [ontology](#) is also provided to allow converting paralex lexicons into RDF [lemon](#)/[ontolex](#) lexicons.

The standard is meant for sharing and interfacing, but not necessarily for data input. The expectation is for data creators to first input data through any convenient means, then convert the result into the standardized structure for publishing and sharing.

<http://www.paralex-standard.org>

DeAR Principles

To publish high quality, easily citable, scientifically impactful data, useful for the long term:

- **Decentralized**

- Centralized data is neither robust nor long-lasting
- International collaboration must be incentivised

- **Automatically validated**

- Manual curation of large datasets is necessary but error-prone
- Automatic quality control of structure and content is crucial

- **Revisable**

- Data evolves, its presentation must be updateable
- Seamless updates of showcases can be generated from a single data source

Relational Structure

form_id	lexeme	cell	phon_form	analysed_phon_form	orth_form	analysed_orth_form	overabundance_tag	defectiveness_tag	epistemic_tag
14618895	sōna_238850_1349429	nom.sg	s y n a	s y n a	sōna	sōna			
14618896	sōna_238850_1349429	gen.sg	s y n a	s y n a	sōna	sōna			
14618897	sōna_238850_1349429	part.sg	s y n a	s y n a	sōna	sōna			
14618898	sōna_238850_1349429	ill.sg	s 'y n:: a	s 'y n:: a	sōnna	s'ōnna		aditive	
14618899	sōna_238850_1349429	ill.sg	s y n a s: e	s y n a + s: e	sōnasse	sōn[sse		sse_illative	
14618900	sōna_238850_1349429	iness.sg	s y n a s	s y n a + s	sōnas	sōna[s			
14618901	sōna_238850_1349429	elat.sg	s y n a s t	s y n a + st	sōnast	sōna[st			
14618902	sōna_238850_1349429	all.sg	s y n a l e	s y n a + le	sōnale	sōna[le			
14618903	sōna_238850_1349429	ad.sg	s y n a l	s y n a + l	sōnal	sōna[l			
14618904	sōna_238850_1349429	abl.sg	s y n a l t:	s y n a + lt:	sōnalt	sōna[lt			
14618905	sōna_238850_1349429	trans.sg	s y n a k s	s y n a + ks	sōnaks	sōna[ks			
14618906	sōna_238850_1349429	term.sg	s y n a n'i	s y n a + n'i	sōnani	sōna[ni			
14618907	sōna_238850_1349429	ess.sg	s y n a n a	s y n a + na	sōnana	sōna[na			
14618908	sōna_238850_1349429	abess.sg	s y n a t a	s y n a + t:a	sōnata	sōna[t:a			
14618909	sōna_238850_1349429	com.sg	s y n a k a	s y n a + ka	sōnaga	sōna[ga			
14618910	sōna_238850_1349429	nom.pl	s y n a t	s y n a + t	sōnad	sōna[d			
14618911	sōna_238850_1349429	gen.pl	s y n a t e	s y n a + te	sōnade	sōna[de	de_te_genitive		
14618912	sōna_238850_1349429	part.pl	s y n u	s y n u	sōnu	sōnu	voc_partitive		
14618913	sōna_238850_1349429	part.pl	s y n a s i t	s y n a + s i t	sōnasid	sōna[sid	sid_partitive		
14618914	sōna_238850_1349429	ill.pl	s y n a t e s: e	s y n a + tes: e	sōnadesse	sōna[desse	de_te_plural		
14618915	sōna_238850_1349429	ill.pl	s y n u s: e	s y n u + s: e	sōnusse	sōnu[ss:e	voc_rad_plural		
14618916	sōna_238850_1349429	iness.pl	s y n a t e s	s y n a + tes	sōnades	sōna[des	de_te_plural		
14618917	sōna_238850_1349429	iness.pl	s y n u s	s y n u + s	sōnus	sōnu[s	voc_rad_plural		
14618918	sōna_238850_1349429	elat.pl	s y n a t e s t	s y n a + test	sōnatest	sōna[dest	de_te_plural		
14618919	sōna_238850_1349429	elat.pl	s y n u s t	s y n u + st	sōnust	sōnu[st	voc_rad_plural		
14618920	sōna_238850_1349429	all.pl	s y n a t e l e	s y n a + tele	sōnadele	sōna[dele	de_te_plural		
14618921	sōna_238850_1349429	all.pl	s y n u l e	s y n u + le	sōnule	sōnu[le	voc_rad_plural		
14618922	sōna_238850_1349429	ad.pl	s y n a t e l	s y n a + tel	sōnadel	sōna[del	de_te_plural		
14618923	sōna_238850_1349429	ad.pl	s y n u l	s y n u + l	sōnul	sōnu[l	voc_rad_plural		
14618924	sōna_238850_1349429	abl.pl	s y n a t e l t:	s y n a + tel:t	sōnadelet	sōna[delt	de_te_plural		
14618925	sōna_238850_1349429	abl.pl	s y n u l t:	s y n u + lt:	sōnult	sōnu[lt	voc_rad_plural		
14618926	sōna_238850_1349429	trans.pl	s y n a t e k s	s y n a + teks	sōnadeks	sōna[deks	de_te_plural		
14618927	sōna_238850_1349429	trans.pl	s y n u k s	s y n u + ks	sōnuks	sōnu[ks	voc_rad_plural		
14618928	sōna_238850_1349429	term.pl	s y n a t e n'i	s y n a + ten'i	sōnadeni	sōna[deni	de_te_plural		
14618929	sōna_238850_1349429	ess.pl	s y n a t e n a	s y n a + ten a	sōnadena	sōna[dena	de_te_plural		
14618930	sōna_238850_1349429	abess.pl	s y n a t e t a	s y n a + tet:a	sōnadeta	sōna[deta	de_te_plural		
14618931	sōna_238850_1349429	com.pl	s y n a t e k a	s y n a + teka	sōnadega	sōna[dega	de_te_plural		

Relational Structure

forms	
form_id	string
lexeme	string
cell	string
phon_form	string
analysed_phon_form	string
orth_form	string
analysed_orth_form	string
overabundance_tag	string
defectiveness_tag	string
epistemic_tag	string

Relational Structure

lexeme_id	lemma	wordId	paradigmId	POS	inflection_class	homonymNr
sõna_238850_1349429	sõna	238850	1349429	noun	17u	1
sahkerdama_230449_1387659	sahkerdama	230449	1387659	verb	27	1
gastroleerima_166864_1388631	gastroleerima	166864	1388631	verb	28	1
vastukaja_254709_1479772	vastukaja	254709	1479772	noun	17	1
sinetama_233880_1386716	sinetama	233880	1386716	verb	27	1
servima_232818_1387793	servima	232818	1387793	verb	28	1
pärija_223664_1348672	pärija	223664	1348672	noun	1	1
tipphetk_244967_1477638	tipphetk	244967	1477638	noun	22i	1
ammendama_157058_1385304	ammendama	157058	1385304	verb	27	1
töö_250035_1373819	töö	250035	1373819	noun	26i	1
transleerima_246364_1389340	transleerima	246364	1389340	verb	28	1
ägisema_261520_1386191	ägisema	261520	1386191	verb	27	1
kristalliseerima_186562_1412429	kristalliseerima	186562	1412429	verb	28	1
vaade_252049_1371583	vaade	252049	1371583	noun	6	1
kõpsutama_189238_1384885	kõpsutama	189238	1384885	verb	27	1
jäätis_175576_1288507	jäätis	175576	1288507	noun	11	1
rutiin_1111368_1444852	rutiin	1111368	1444852	noun	22e	2
pesu_215900_1352908	pesu	215900	1352908	noun	17	2
põlistuma_222793_1506098	põlistuma	222793	1506098	verb	27	1
veerg_255264_1297645	veerg	255264	1297645	noun	22e	1
sonkima_235110_1387555	sonkima	235110	1387555	verb	28	1

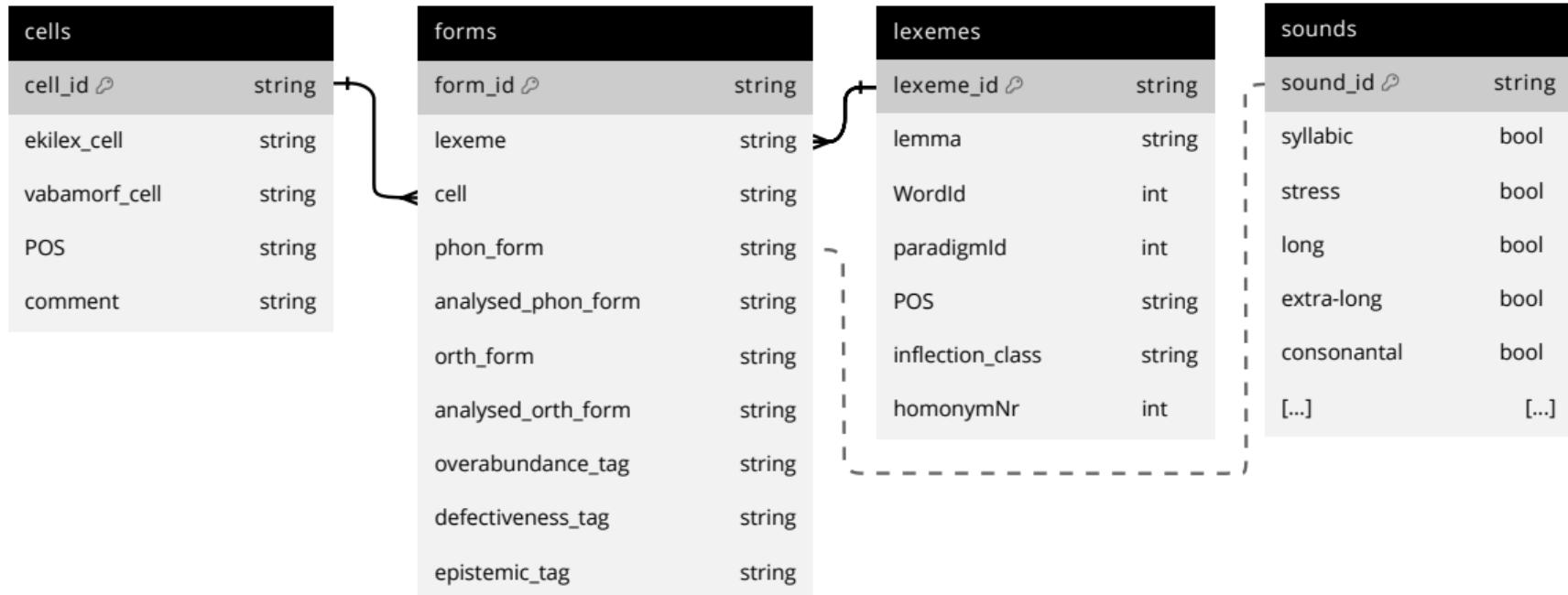
Relational Structure

forms		lexemes	
form_id	string	lexeme_id	string
lexeme	string	lemma	string
cell	string	WordId	int
phon_form	string	paradigmId	int
analysed_phon_form	string	POS	string
orth_form	string	inflection_class	string
analysed_orth_form	string	homonymNr	int
overabundance_tag	string		
defectiveness_tag	string		
epistemic_tag	string		

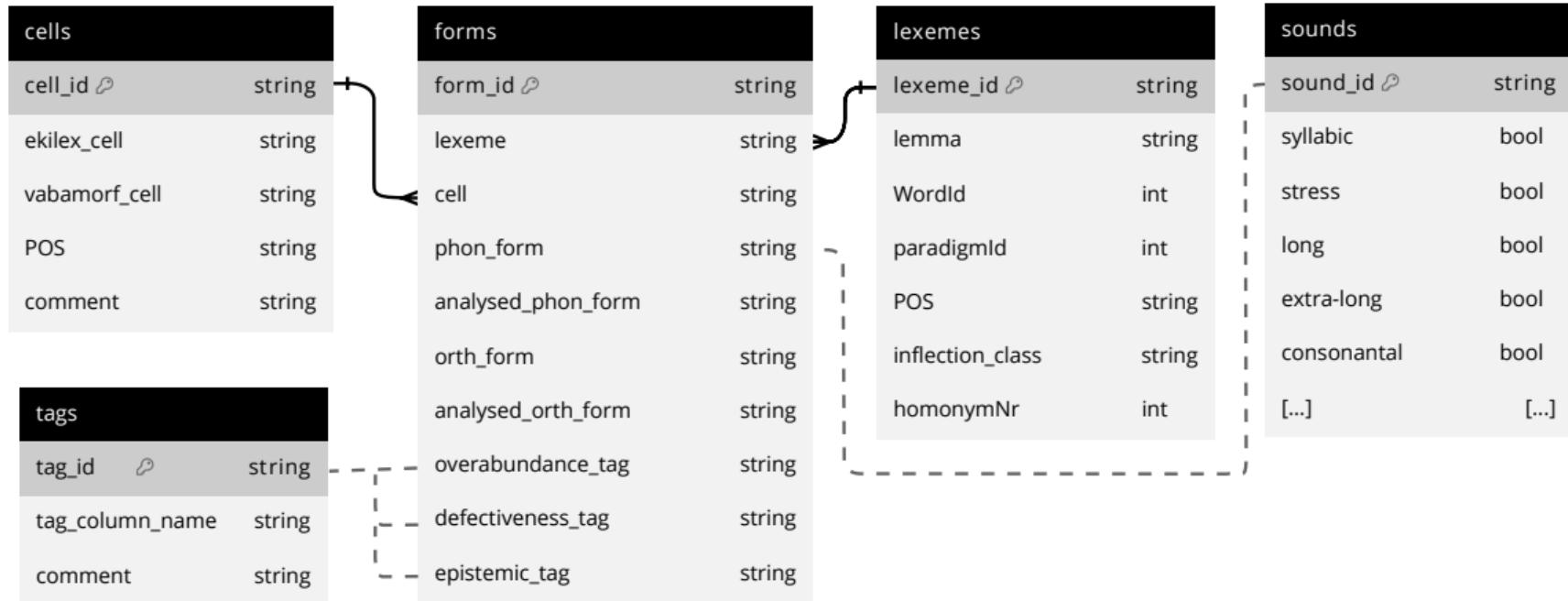
Relational Structure

cells		forms		lexemes	
cell_id	string	form_id	string	lexeme_id	string
ekilex_cell	string	lexeme	string	lemma	string
vabamorf_cell	string	cell	string	WordId	int
POS	string	phon_form	string	paradigmId	int
comment	string	analysed_phon_form	string	POS	string
		orth_form	string	inflection_class	string
		analysed_orth_form	string	homonymNr	int
		overabundance_tag	string		
		defectiveness_tag	string		
		epistemic_tag	string		

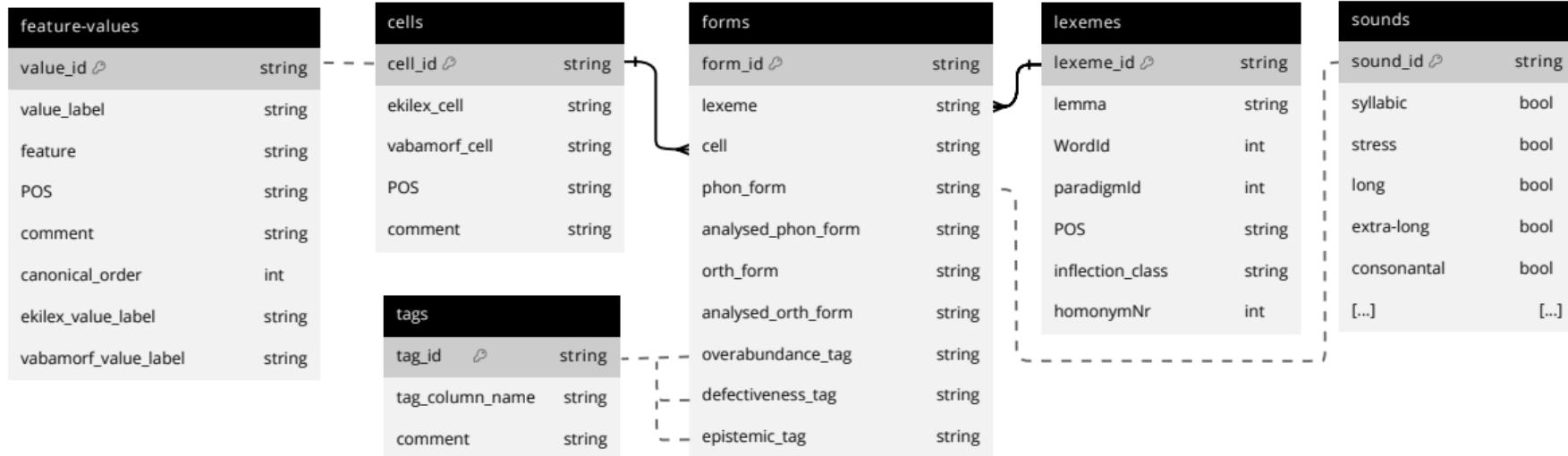
Relational Structure



Relational Structure



Relational Structure



Evaluation

- **Validation:** frictionless & paralex
- **G2P:** development set; targeted manual verifications
- **Paradigms:** Comparison to Unimorph

Conclusion: Eesthetic lexicon

DOI [10.5281/zenodo.8383522](https://doi.org/10.5281/zenodo.8383522)

- DeAR Paralex lexicon for Estonian N & V
- Large (5k lexemes per POS)
- Rich linguistic annotation
- Generated phonological forms
- Linguistically motivated G2P rules

Conclusion: Eesthetic lexicon

DOI [10.5281/zenodo.8383522](https://doi.org/10.5281/zenodo.8383522)

- DeAR Paralex lexicon for Estonian N & V
- Large (5k lexemes per POS)
- Rich linguistic annotation
- Generated phonological forms
- Linguistically motivated G2P rules

Thank you !