

SPRÅKBANKENTEXT







Swedish morphological profile

Elena Volodina, Språkbanken Text, Sweden (presenter)

Based on work by: Elena Volodina, Yousuf Ali Mohammed, Therese Lindström Tiedemann and annotators

Demo

https://spraakbanken.gu.se/larkalabb/svlp

Login: demo



Swedish L2 Morphological profile

Grammatical profile	Morphological profile		
æ	Word family Morpheme family		
	Grammatical profile		

(will be) freely available at: <u>https://spraakbanken.gu.se/larkalabb/svlp</u>

Morpheme family 🕒



Word formation mechanisms





Word formation	Receptive
Abbreviation	17
Compound	14344
Derivation	13871
Lexicalized form	4367

Word formation	Productive
Abbreviation	4
Compound	1629
Derivation	2448
Lexicalized form	634

Morpheme categories



	Morpheme Category	Receptive
	Binding morpheme	990
	Inflectional ending	2739
	Derivational prefix	1939
Word	Root morpheme	21444

Morpheme Category	Productive
Binding morpheme	138
Inflectional ending	436
Derivational prefix	338
Root morpheme	3108

Statistics: type-token ratio

CEFR level

CEFR level	Receptive (Type)	Productive (Type)	Receptive (Token)	Productive (Token)	Receptive (TTR)	Productive (TTR)
A1	2166	464	25015	4083	0.09	0.11
A2	5134	714	16367	1866	0.31	0.38
B1	10327	1184	27224	2829	0.38	0.42
B2	10916	1310	22569	2748	0.48	0.48











Search morpheme			Search word			ational prefix	Word Cla	
A1	A2	B1	B2	C1	C2		Table	
Add colu	imn for relate	ed roots	Show on	ly first occu	rrence Tab	les - description 🗗	Filters - desci	

Total rows : 428

Morpheme ↓↑	CEFR level ↓↑	Morpheme Category ↓↑	Lemgram ↓↑	Sense ↓ ↑
0	A1	Derivational prefix	obekvämav.1	obekväm1
0	A1	Derivational prefix	obestämdav.1	obestämd1
0	A1	Derivational prefix	ogiftav.1	ogift1
0	A1	Derivational prefix	olikav.1	olik1
0	A1	Derivational prefix	olyckann.1	olycka2
0	A1	Derivational prefix	oregelbundenav.1	oregelbunden1
0	A1	Derivational prefix	orördav.1	orörd1
0	A1	Derivational prefix	ovädernn.1	oväder1

Root morpheme ↓↑	CEFR level ↓↑	Lemgram ↓↑	Sense ↓↑	Word Class ↓↑	Saldo Word Class ↓↑	Word formation ↓↑	Receptive ↓↑ 🔂	Productive 🕼 🗗
jul	A1	julnn.1	jul1	Noun (NN)	Noun (nn)	Root lexeme	0.39 (1)	0.00 (0)
jul	A2	julbordnn.1	julbord1	Noun (NN)	Noun (nn)	Compound	0.15 (1)	0.00 (0)
jul	A2	juldagsmorgonnn.1	juldagsmorgon1	Noun (NN)	Noun (nn)	Compound	0.15 (1)	0.00 (0)
jul	A2	julklappnn.1	julklapp1	Noun (NN)	Noun (nn)	Compound	0.29 (2)	0.00 (0)
jul	A2	jullovnn.1	jullov1	Noun (NN)	Noun (nn)	Compound	0.87 (6)	0.44 (1)
jul	A2	julottann.1	julotta1	Noun (NN)	Noun (nn)	Compound	0.15 (1)	0.00 (0)
jul	A2	jultomtenn.1	jultomte1	Noun (NN)	Noun (nn)	Compound	0.15 (1)	0.00 (0)
jul	B1	julfestnn.1	julfest1	Noun (NN)	Noun (nn)	Compound	0.08 (1)	0.00 (0)
jul	B2	julaftonnn.1	julafton1	Noun (NN)	Noun (nn)	Compound	1.06 (16)	0.00 (0)
jul	B2	julbocknn.1	julbock1	Noun (NN)	Noun (nn)	Compound	0.13 (2)	0.00 (0)
jul	B2	juldagnn.1	juldag1	Noun (NN)	Noun (nn)	Compound	0.26 (4)	0.00 (0)
jul	B2	julfirandenn.1	julfirande1	Noun (NN)	Noun (nn)	Derivation	0.13 (2)	0.00 (0)
jul	B2	julgåvann.1	julgåva1	Noun (NN)	Noun (nn)	Compound	0.07 (1)	0.00 (0)
jul	B2	julgrannn.1	julgran1	Noun (NN)	Noun (nn)	Compound	0.33 (5)	0.00 (0)

Search morpho	eme		Search word		Root	eme Category morpheme	Word Class	•
A1	A2	B1	B2	C1	C2		Tables	Graphs
Add colu	mn for relat	ed roots	Show on	ly first occu	rrence Tab	les - description 📴 F	Ilters - description	电

Total rows : 202

Morpheme ↓↑	CEFR level ↓↑	Morpheme Category ↓↑	Lemgram ↓ ↑	Sense ↓↑	Word Class ↓↑
land	A1	Root morpheme	Ångermanlandpm.1	Ångermanland1	Proper noun (PM)
land	A1	Root morpheme	Jämtlandpm.1	Jämtland1	Proper noun (PM)
land	A1	Root morpheme	Svealandpm.1	Svealand1	Proper noun (PM)
land	A1	Root morpheme	Thailandpm.1	Thailand1	Proper noun (PM)
land	A1	Root morpheme	Tysklandpm.1	Tyskland1	Proper noun (PM)
land	A1	Root morpheme	Upplandpm.1	Uppland1	Proper noun (PM)
land	A1	Root morpheme	Värmlandpm.1	Värmland1	Proper noun (PM)
land	A1	Root morpheme	landnn.1	land1	Noun (NN)





Some definitions

Lexical units as the main entry

- word forms → running word forms (walked)
- lemma → base form (walk)
- lemgram → base form + POS (walk, noun)
- flemma → homonymous forms (walk, noun & walk, verb)
- word family → word grouped around the same root (walk, walk-through)
- morpheme family \rightarrow

Presuppose different assumptions on vocabulary acquisition. Some discussion see in

Brown, Dale. 2018. Examining the word family through word lists. *Vocabulary Learning and Instruction* 7, no. 1: 51-65.
Stoeckel, Tim, Tomoko Ishii, and Phil Bennett. 2020. Is the lemma more appropriate than the flemma as a word counting unit? *Applied Linguistics* 41, no. 4 (2020): 601-606.

• Words that share morphological components form so-called **morphological families**, which are sets of words that are *related morphologically* (sharing a morpheme) and *semantically* (sharing a meaning)

Nikolaev et al. (2019) [cognitive studies]

• The **derivation family** is defined as a network of *derivationally related lexemes.*

Körtvélyessy et al. (2020) [typological studies]

• **Derivational families**, i.e. clusters of lemmas in derivational relationships.

Zeller et al. (2013) [computational lexicography]

Suffix -skap (-skap family)

Total rows : 53

Morpheme ↓↑	CEFR level ↓↑	Morpheme Category ↓↑	Lemgram ↓↑	Sense ↓↑	Word Class ↓↑	Saldo Word Class ↓↑	Word formation ↓↑	Receptive ↓↑ 🔁	Productive 🕼 🔁
skap	B1	Derivational suffix	äktenskapnn.1	äktenskap1	Noun (NN)	Noun (nn)	Derivation	0.45 (6)	8.33 (19)
skap	B2	Derivational suffix	kunskapnn.1	kunskap1	Noun (NN)	Noun (nn)	Derivation	1.72 (26)	0.90 (4)
skap	B1	Derivational suffix	partnerskapnn.1	partnerskap1	Noun (NN)	Noun (nn)	Derivation	0.15 (2)	1.31 (3)
skap	B1	Derivational suffix	vänskapnn.1	vänskap1	Noun (NN)	Noun (nn)	Derivation	0.60 (8)	0.88 (2)
skap	C1	Derivational suffix	föräldraskapnn.1	föräldraskap1	Noun (NN)	Noun (nn)	Derivation	0.57 (9)	0.19 (1)
skap	B2	Derivational suffix	medborgarskapnn.1	medborgarskap1	Noun (NN)	Noun (nn)	Derivation	0.26 (4)	0.22 (1)
skap	B2	Derivational suffix	kunskapsmässigav.1	kunskapsmässig1	Adjective (JJ)	Adjective (av)	Derivation	0.00 (0)	0.22 (1)
skap	B2	Derivational suffix	språkkunskapnn.1	språkkunskap1	Noun (NN)	Noun (nn)	Compound	0.20 (3)	0.22 (1)

Related resources (a few examples)

- CELEX (Dutch, German, English) lexicons based on L1 corpora
 - <u>https://catalog.ldc.upenn.edu/LDC96L14</u>
 - contain derivational and compositional structure, inflectional paradigms (among others)
 - protected by paywall
- DErivBase (German) resource with morphological families
 - <u>https://www.ims.uni-stuttgart.de/forschung/ressourcen/lexika/derivbase/</u>
 - based on German CELEX
- Derivational networks (40 EU languages)
 - <u>https://www.degruyter.com/document/doi/10.1515/9783110686630/html</u>
 - 3 word classes á 30 core root lexemes across languages
- Word formation networks (Polish, Czech, French, Spanish)
 - Lango et al. (2021); DeriNet (lexical network for Czech)
- Word families (English)
 - Bauer & Nation (1993) <u>https://academic.oup.com/ijl/article-abstract/6/4/253/941347</u>
 - Coxhead (1998) An Academic Word List (570 word families)

• Inclusion of a related form of a word within a **word family** depends on criteria involving *frequency, regularity, productivity and predictability*. These criteria are applied to English *affixes* so that the *inflectional affixes* and the most *useful derivational affixes* are arranged into a **graded** set of seven levels.

Bauer and Nation (1993) [L2 acquisition studies]

word forms		Word families	
2 lemgram (inflectional paradigms)	develop develops developed developing	wood wood's woods wooded	bright brighter brightest
word formation (affixation)	developable undevelopable developer(s) undeveloped	woody woodiest woodier woodiness	brightly brightish brightness
4	development(s) developmental developmentally		
5 Bauer & Nation	developmentwise semideveloped antidevelopment	wooden	brighten
(1993) 6	redevelop predevelopment	anti-wood	

Coxhead (1998)

Academic Word List

The following is a list of the words from the Academic Word List that are highlighted in this dictionary. Words shown in **bold** are one of the 'parent words'.

$abandon^1v$

abandoned *adj* abnormal *adj* **abstract**¹*adj* **abstract**²*n* **abstract**³*v* abstraction *n* academic¹ *adj* academic²*n* **academy** *n* **access**¹*n* adjacent adj adjust v adjustment n administration n administrative adj adult¹n adult² adj advocacy n advocate¹v advocate² n affect v

. 1

analogy n analyse v analysis n analyst n analytical adj analyze v annual¹adj anticipate v anticipation n apparent adj append v

assist¹v assistance n assume v assuming conj assumption n assurance n assurance n assure v attach v attachment n attain v attainment n

Word formation networks: Polish, Czech, Spanish, French



Lango, M., Žabokrtský, Z., & Ševčíková, M. (2021). Semi-automatic construction of word-formation networks. *Language Resources and Evaluation*, *55*(1), 3-32. <u>https://link.springer.com/article/10.1007/s10579-019-09484-2</u>

Our definition

Word family is a group of words sharing

- (1) the same root/base, including potential variations of that (e.g. go, bygone), compounds (e.g. light, light-hearted), derived forms (e.g. teach, teacher) and multi-word expressions (e.g. day off, day by day)
- (2) the same meaning (X sunlight, light-hearted)

How did we do it?

CoDeRooMor dataset:

Compounding, Derivation, Root, Morphology ... and more



Elena Volodina, Yousuf Ali Mohammed, and Therese Lindström Tiedemann. (2021). CoDeRooMor: A new dataset for non-inflectional morphology studies of Swedish. *NoDaLiDa 2021* : Linköping Electronic University Press, Vol.178. [pdf] [[video-presentation]]

Annotation categories: word formation

Word formation	Definition	Example
Abbreviation	words consisting of the initial components of a word or	AB (aktiebolag) (cf. Eng. 'ltd' = 'lim-
	several words, including chemical abbreviaions and some	ited'), Au (Sw. guld, Eng. 'gold')
	blends	
Compound	words formed by adding together two stems	skol+bok ('school book')
Derivation	words formed by adding a prefix or a suffix to a stem	sorglig ('sad')
Lexicalized form	words that cannot be reduced to baseforms, e.g. MWEs	Aftonbladet (name of a tabloid), järnspik ar (a swearword)
Root lexeme	words consisting of a root only or a root and an inflec-	bok ('book'), adjö ('goodbye'), ande
	tional suffix	('spirit')
Unknown	reserved for difficult or uncertain cases including most	alzheimers (name of a disease), kalen-
	first names	der ('calendar')

Annotation categories: morpheme types

Morph. category	Explanation	Example
р	derivational prefix	för djupa
r	root (orthographic)	kaot isk
rr	real root	kaos (kaotisk)
S	derivational suffix	kaot isk
f	infix*	kedjebrev
i	inflectional suffix	i_höst as
?	unknown	ir on i

Annotation tool: Legato

Lexicographic Annotation Tool (LEGATO)



Quick jump to:

A B C D E F G H I J K L M N O P Q R S T U V W X Y Z Å Ä Ö

			Jump to:	Jump
			1 - 16230	
SALDO lemgram	Saldo POS	Part-of-Speech	Progr CEFR level	ress: 11457/16230
samvetenn.1	noun (nn)	Noun (NN)	B1	
Saldo sense: samvete1 Saldo primary descriptor: känsla. Saldo secondary descriptor: rätt.	.2 2			

Examples:

Jonas : Har du aldrig fått dåligt ** samvete ** ?

Kanske hade han trots allt dåligt ** samvete ** för att han hade lämnat oss och skaffat ny familj

Birgitta erkänner att hon borde ha riktigt dåligt ** samvete ** för det här, men hon menar att det inte drabbar någon enskild.

Previous morphology1:	Previous morphology2:	SO analysis:		SAOL analysis:	
p: sam r: vet s: e	,	sam=vet·et	11.	sam=vete	
Сору		Сору		Сору	

^				I	
սս	rre	nτ	va	iue	

r: sam r: vet s: e

Select word form:

- ROOT LEXEME
- COMPOUND
- C LEXICALIZED FORM
 - UNKNOWN

https://spraakbanken.gu.se/larkalabb/legato

Annotation quality

Annotation type	1-100	-200	-300	-400
Segmentation	0.87	0.87	0.89	0.93
Labeling	0.86	0.86	0.89	0.86
Segmentation+Labeling	0.85	0.85	0.87	0.88
Word formation	0.89	0.89	0.94	0.91

(Krippendorff's alpha for Inter-Annotator Agreement)

Annotators: Beatrice Silén, Stellan Petersson, Maisa Lauriala +Therese Lindström Tiedemann (researcher)

Lemgram view

Lemgram	Sense	POS	Analysis	Segment.	Pattern	RealRoot	WordForm	CEFR
adekvatav.1	adekvat1	JJ	p:ad r:ekv s:at	ad-ekv-at	p:r:s		derivation	C1
adlavb.1	adla1	PC	r:adl s:a	adl-a	r:s	rr:adel	derivation	B2
adelnn.1	adel1	NN	r:adel	adel	r		root_lexeme	B1
adelsmannn.1	adelsman1	NN	r:adel f:s r:man	adel-s-man	r:f:r		compound	B1
adjektivnn.1	adjektiv1	NN	p:ad r:jekt s:iv	ad-jekt-iv	p:r:s		derivation	A2
adjöin.1	adjö1	IN	r:adjö	adjö	r		root_lexeme	A2

"Morpheme family" view

Morpheme	Identifier	Category	Frequency	Examples
a	S	suffix	1 605	leverera, lugna_sig, meritera, narkotika, pumpa, rasa
er	s	suffix	577	abdikera, intrigera, politiker, kritiker, motivera, tekniker
tid	r	root	128	arbetstid, nutid, skoltid, livstid, dåtid, deltid
ny	r	root	46	nyinköpt, nykokt, nykomling, nyligen, nymodighet, Nynäshamn
0	р	prefix	240	olaglig, olämplig, olik, olika, olikhet, oljud
re	p	prefix	105	reaktionstid, rebell, rebellisk, recensent, recensera, recension
S	f	infix*	803	fredstid, landsfader, riksbank, tvångsgift, landsdel, riksdag
0	f	infix*	76	vilopaus, sagobok, sannolik, sociolog, vilorum, typografi

freely available from https://spraakbanken.gu.se/en/resources/coderoomor

A few facts

- 4 429 word families
- smallest: 1 member (e.g. adrenalin)
- biggest: 313 members (e.g. utbildning)
- most numerous: function words consisting of root lexelmes (e.g. preposition ut)

Total rows : 22553 Unique roots/level: A1: 849, A2: 1527, B1: 2168, B2: 2355, C1: 1870,

Why did we do it?

Complexity as a proxy of learner language development



Our focus



Complexity measures (infl. morphology)

- Syntactic variety (Forster & Skenann, 1996; Ellis & Yan, 2005)
- Inflectional diversity index (Malvern et al., 2004)
- Normalized mean size of paradigm (Xanthos & Gillis, 2010)
- MCI Morphological Complexity Index (Brezina & Pallotti, 2016):
 - <u>http://corpora.lancs.ac.uk/vocab/process_text_morph.php</u>

Complexity measures (deriv. morphology)

• ...yet to be developed

What can we do with it?

Statistic analysis of L2 Swedish

Morpheme	Unique	Sen*Lex	COCTAILL	SweLL-pilot	Examples
category	count				
root	4429	23 987	471 056	142 381	matbord, kärleksaffär, sagolik
suffix	259	10 062	91 646	28 638	markn ad , kost sam , milit är
prefix	155	2 183	19 828	5 489	konsonant, nyrenoverad
infix	12	1 089	3 441	1 641	kännedom, kvinnorörelse
inflection	32	3 067	88 641	28 810	saker_och_ting, Medelhavet, läsa

Relate statistics and levels of language development



Derivational complexity

Item	A1	A2	B1	B2	C1	Total
word	1369	2689	4518	4440	3211	16230
morpheme	2457	5727	11318	11732	9292	40534
morpheme/word	1.79	2.13	2.51	2.64	2.89	2.50

Within L2 context

- Relation of morphological awareness and language proficiency
- Awareness of derivational morphemes
- Priming effect of roots/affixes on learning of new vocabulary (cf word families)
- Vocabulary testing

. . .

• Pedagogical applications / ICALL (e.g. exercises)

Outside L2 context

- Typological studies of languages
- Computational complexity of languages
- Psycholinguistic studies
- Cognitive studies
- Machine translation
- Cognate identification

Automatic tools for

- ...detection of morpheme boundaries
- ...labeling morphemes for their types
- ...labeling words for their word formation type



• ...with the ultimate aim/hope to add this type of analysis to automatic pipelines, primarily to Sparv



Our major principles

- research oriented, in addition to teacher- & learner focus
- descriptive, transparent, empirically based, provides access to all corpus hits
- splits receptive and productive knowledge
- statistic analysis and graphical representation





SPRÅKBANKENTEXT







Thank you! Comments? Questions?

(or mailto: elena dot volodina at svenska dot gu dot se)