



# ORP & KARP

## **Språkbanken's infrastructures for corpora and lexicons**

**EUROPEAN  
NETWORK OF  
e-LEXICOGRAPHY**

Elena Volodina & Ildikó Pilán  
COST ENeL WG3 meeting  
Vienna, 12 February 2015



# KORP & KARP origin

SweFN++ (2009-2015), a project where one of the objectives is to create, curate, and integrate free Swedish lexical resources --> for researchers and language technology applications

META-NORD (2011-2013), a project aiming at establishing an open linguistic infrastructure in the Baltic and Nordic countries.

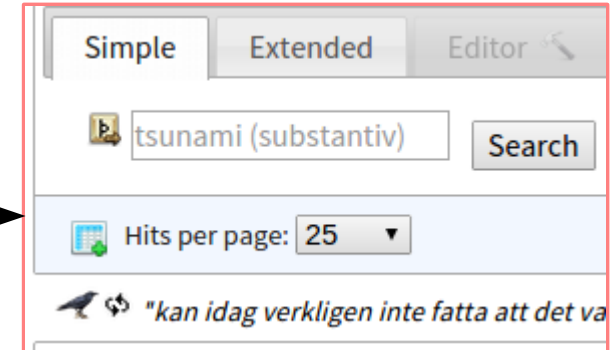
<http://spraakbanken.gu.se/korp>

<http://spraakbanken.gu.se/karp>



# Architecture principles

**Front-end** (user interface): search, present, edit



**Back-end** (web services): select from databases, deliver in machine-readable format

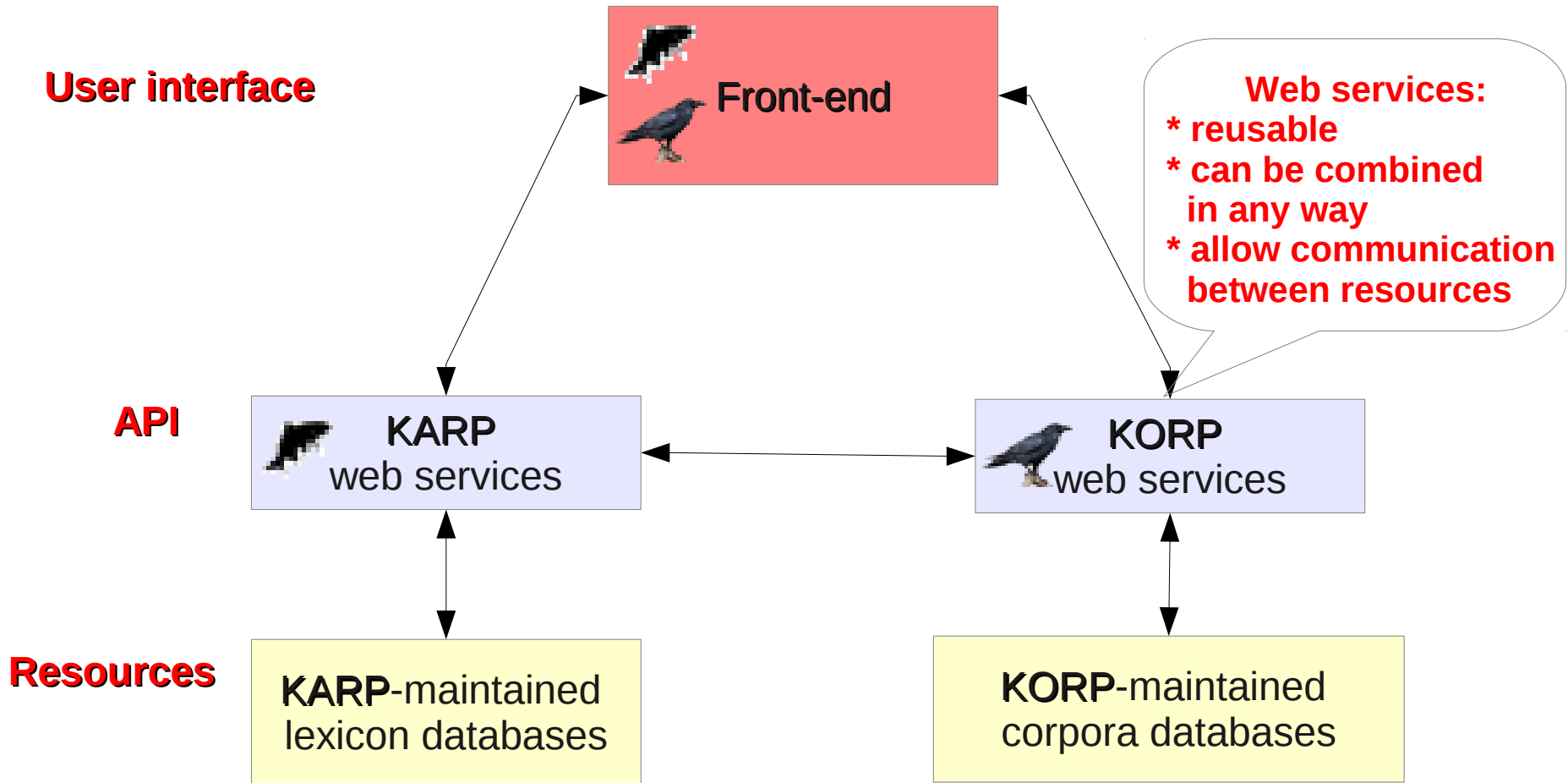
```
"Lemma": {  
  "FormRepresentation": {  
    "feat": [  
      {"att": "lemgram",  
        "val": "huvuddrag..nn.1" },  
      {"att": "partOfSpeech",  
        "val": "nn" },  
      {"att": "paradigm",  
        "val": "nn_6n_bord" }, ]  
    }  
  },  
}
```

**Databases:** access to corpora and lexicons





# Secret of KORP-2-KARP bi-directionality





# CORP - the corpus infrastructure of SB

- **8,5 billion words** new and old texts – downloadable, scrambled
- **KWIC** vs **context** display
- Linguistic **annotation**: POS, lemgrams, msd, dependencies ...
- **Statistic** functions: pie charts, word lists, trend diagrams ...
- **Word pictures**
- **Save** or **bookmark** searches
- Different **views**: modern, parallel, Old Swedish ...





190 av 197 korpusar valda — 8,49G av 8,50G token

1850 1900 1950 2000

- Akademiska texter (2)
- August Strindberg (2)
- Finlandssvenska texter (56)
- Skyddade korpusar (7)
- Medicinska texter (12)
- Skönlitteratur (6)
- Sociala medier (60)
- Tidningstexter (37)
- Tidskrifter (1)
- Dramawebben (demo)
- LäSBarT - Lättläst svenska och barnbokstext
- PAROLE
- Psalmboken (1937)
- SNP 78-79 (Riksdagens snabbprotokoll)
- SUC 2.0
- SUC 3.0
- SALT svenska-nederländska
- Europarl svenska
- Svenska partiprogram och valmanifest 1887-2010
- Svensk författningssamling
- Svenskt frasnät (SweFN)
- Svenska Wikipedia (januari 2015)

Enkel Utökad Av

även som  förled  efterled

KWIC: träffar per sida: 25

KWIC Statistik

Antal träffar: 815

Föregående 1 2 3

a fel hur de än gör: de förvänt  
v Iris som en " extremt farlig  
veckan, enligt den tyska nyh  
u tillbaka i äventyrsfilmen nä  
n är oron stor i en del EU-län  
Tv-kändisar, intelle  
ra städer vid Nordsjön och Ö  
ar tvingat fram öppningar av

Markera alla  
Avmarkera

rbild

ontext

GP 2001

I dag väntas en flodvåg nå Jakutsk  
ts oemotståndliga flodvåg ", men om de i stället överlever o  
orsaka livshotande flodvåg och lerskred i bergig terräng.  
vakueras sedan en flodvåg från Wisla närmade sig.  
det handlar om en flodvåg av privata alternativ.  
en jättekondor, en flodvåg och en jaguar.  
are besatta av den flodvåg som åter skulle hemsöka jorden.  
ar slussarna för en flodvåg av arbetslösa invandrare.

GP 2002

v som en fascistisk flodvåg över Marseille.  
ler dukade under i flodvågorna av kulor och granater som svepte  
nad sedan om den flodvåg av prostataremitter som på sena  
a skyddsvallar mot flodvågen som de närmaste dagarna vänta  
a där försäkrar att flodvågen kommer att ha avtagit innan den  
er sedan en enorm flodvåg sköljde med sig bilar, hus och mä  
äntades en formlig flodvåg nå Prag.  
vinkade vad det är för plötslig flodvåg av invandrare Persson är rädd fö



24 of 197 corpora selected — 321.94M of 8.50G tokens

Search history

Simple Extended Advanced Compare 2

ord (noun) Search

also as  initial part  final part and  case-insensitive

Relaterade ord (SWE-FN)  
ord, term, fackterm  
glosa...

KWIC: hits per page: 25 Sort within corpora: not sorted Statistics: compile based on: lemgram  Show word picture

KWIC Statistics Word picture

Results: 88,127

Previous 1 2 3 4 5 6 7 8 9 10 11 .. 3525 3526 Next Show context

ASPAC SVENSKA (does not support extended context)

(Den här gången var hon glad att ingen kunde höra henne, för **ordet** lät inte riktigt rätt.  
a fall på sig och började läsa, men hennes huvud var så fullt av hummerkadriljen att hon knappt visste vad hon sade, och **orden** blev verkligen ganska konstiga:  
Jag har **glömt orden** .  
Jag har hört vartenda **ord** ni sade.  
Låt mig slippa höra det **ordet** mer!  
Dessa **ord** följdes av en lång tystnad, endast avbruten av ett och  
Hon urskilde **orden** : - Var är den andra stegen?



KWIC

Statistics

Word picture

Results: 88,127

Previous

1

2

3

4

5

6

7

8

9

10

11

..

3525

3526

Next

Show context

ASPAC SVENSKA (does not

(Den här gången var hon glad att ingen kunde höra henne, för **ordet** lät int

a fall på sig och började läsa, men hennes huvud var så fullt av hummerkadriljen att hon knappt visste vad hon sade, och **orden** blev v

Jag har glömt **orden**.

Jag har hört vartenda **ord** ni sad

Låt mig slippa höra det **ordet** mer!

Dessa **ord** följde

Hon urskilde **orden** : - Var

Men rösten lät hes och besynnerlig, och **orden** var in

Det sista **ordet** kom r

Hon sade det sista **ordet** ett pa

Hon sade de sista **orden** högt,

" och den fanns absolut inte här förut ", sade Alice ) och vid flaskhalsen satt fastbunden en papperslapp, och på den var **orden** " DRIC

( Alice hade ingen aning om vad latitud och longitud var för något, men hon tyckte det var roligt att säga långa, konstiga **ord** .

Jag begriper inte vad hälften av de där långa konstiga **orden** betyd

Sitt ned båda två, men säg inte ett **ord** förrär

Gripen svarade med nästan exakt samma **ord** som f

stora, tårdränkta ögon men inte sade ett **ord** .

Den börjar med **orden** " När

En del **ord** var no

nnar ckarn tjöt så att Alice knappt kunde höra **orden** :

iktigt, som om han velat ta reda på vilket **ord** som li

" för kålmasken och det blev alldeles fel **ord** .

st fanns en nyfiken liten best på vilken orden " XT

### Corpus

ASPAC svenska

### Text attributes

title: Alice i Underlandet

author: Lewis Carrol

language: swe

description: Översättning av Harry Lundin

### Word attributes

part-of-speech: noun

baseform:

ord

lemgram:

ord (noun)

sense:

ord

final part: [empty]

initial part: [empty]

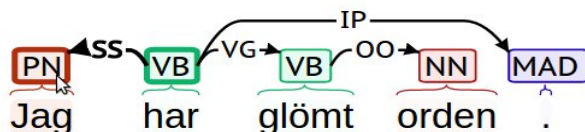
dependency relation: Direct object

msd: NN.NEU.PLU.DEF.NOM ⓘ

Show Dependency Tree

### Dependency Tree

part-of-speech: pronoun









## Word picture for the noun *ord* (word):

Word's grammatical and collocational behavior

Based on MI, statistics from the selected corpora

ord (noun)

Search

also as  initial part  final part and  case-insensitive

Relaterade ord (SWE-FN)  
ord, term, fackterm  
glosa...

KWIC: hits per page: 25 | sort within corpora: not sorted | Statistics: compile based on: lemgram |  Show word picture

KWIC | Statistics | Word picture

Show word class

Preposition	Pre-Modifier	ord	Post-Modifier	ord	Verb	Verb	ord
1. med	9191	1. fri	2353	1. betyda	358	1. använda	994
2. utan	494	2. ful	314	2. återkomma	57	2. skräda	256
3. bakom	235	3. vacker	376	3. beskriva	58	3. säga	1219
4. mot	484	4. hård	395	4. använda	80	4. sätta	585
5. bortom	50	5. tom	208	5. som bäst beskriva	28	5. höra	456
6. i brist på bättre	23	6. ont	96	6. i gp	62	6. uttala	180
7. mellan	125	7. vänlig	120	7. i ki-duks	20	7. yttra	141
8. ute efter	14	8. engelsk	179	8. av kärlek	30	8. växla	175
9. trots	57	9. fin	229	9. för dag	41	9. byta	273
10. om och om igen	6	10. klok	100	10. på guldväg	20	10. hitta	332
11. me	6	11. latinsk	59	11. i mun	34	11. väga <sup>3</sup>	146
12. af	6	12. fel	111	12. kunna	111	12. väga	149
13. för övrigt	9	13. förflugen	37	13. dyka upp	32	13. få	804
14. till sist	7	14. ond	117			14. lära sig	140
15. till intet förpliktigande	3	15. grekisk	84			15. förstå	150




37 of 38 corpora selected — 354.58M of 356.10M tokens

Search history

Extended

English

26 languages paired with Swedish

baseform is  
sun Aa  
or  
+ 

## Parallel corpora view

Add lang Search

KWIC: hits per page: 10 sort within corpora: not sorted Statistics: compile based on: word

KWIC Statistics

Results: 135

Previous 1 2 3 4 5 6 7 8 9 10 11 .. 13 14 Next Show context

ASPAC SVENSKA-ENGELSKA (does not support extended context)  
er sat still just as she left her, leaning her head on her hand, watching the setting sun, and thinking of little Alice and all her wo  
kvar där hon hade lämnat henne. Hon satt med huvudet lutat mot handen och tittade på solnedgången och tänkte på lilla Alic  
They very soon came upon a Gryphon, lying fast asleep in the sun.  
De kom snart till en grip, som låg och sov i solskenet.  
- timidly at first - to push a charming little spig inoffensively upward toward the sun.

### Corpus

ASPAC svenska-engelska

### Text attributes

title: Alice in Wonderland (1865)  
author: Carroll, Lewis  
language: eng  
description: [empty]

### Word attributes



KWIC: hits per page: 10 ▾ sort within corpora: not sorted ▾ Statistics: compile based on: word ▾

KWIC Statistics

Results: 135

Previous 1 2 3 4 5 6 7 8 9 10 11 .. 13 14

Next

Show context ←

ASPAC SVENSKA-ENGELSKA (d)

ter sat still just as she left her, leaning her head on her hand, watching the setting sun , and

t kvar där hon hade lämnat henne . Hon satt med huvudet lutat mot handen och tittade på

They very soon came upon a Gryphon, lying fast asleep in the sun

De kom snart till en grip , som låg och sov i solskenet .

- timidly at first - to push a charming little sprig inoffensively upward toward the sun .

ett näpet och vänligt skott mot solen .

" I should like to see a sunset... Do me that kindness... Order the sun to s

- Jag skulle så gärna vilja se en solnedgång ... Var snäll och befall solen att g

But if you tame me, it will be as if the sun cam

Om du tämjde mig skulle det lysa upp min tillva

Just so. Everybody knows that when it is noon in the United States the sun is se

Ja , sannerligen - när klockan är tolv i Amerika , går ju solen som alla vet ned i Fra

" And I was born at the same moment as the sun ... "

Och jag föddes just som solen ...

: hurt any more, he still had the delicious taste of cinnamon bun in his mouth, the sun was

längre , han hade fortfarande den härliga kanelbullesmaken kvar i mun , solen sken in gen

### Corpus

ASPAC svenska-engelska

### Text attributes

title: Alice in Wonderland (1865)

author: Carroll, Lewis

language: eng

description: [empty]

### Word attributes

part-of-speech: NN

baseform: sun

Parallel corpora  
search results





# ORP Annotation Laboratory

Language of analysis: English

Load example: [Example](#)

- 1 Lexicography
- 2 Practical lex
- 3 Theoretical l
- 4 syntagmatic a
- 5 developing th
- 6 the needs for
- 7 the data inco
- 8 'metalexicrog
- 9 focuses on th
- 10 provide a des
- 11 dictionary or
- 12 compilation,
- 13 (relatively r
- 14 e.g. legal le
- 15 and following
- 16 dictionaries.

- Bulgarian
- Dutch
- English
- Estonian
- Finnish
- French
- German
- Italian
- Latin
- Polish
- Portuguese
- Russian
- Slovak
- Spanish
- Swedish
- Swedish-dev

Positional attribute g

word pos msd

Show advanced setti

separate but equally important groups:  
 or craft of compiling, writing and editing dictionaries.  
 holarly discipline of analyzing and describing the semantic,  
 onships within the lexicon (vocabulary) of a language,  
 components and structures linking the data in dictionaries,  
 in specific types of situations, and how users may best access  
 nd electronic dictionaries. This is sometimes referred to as  
 ed to lexicography is called a lexicographer. General lexicography  
 , use and evaluation of general dictionaries, i.e. dictionaries that  
 age in general use. Such a dictionary is usually called a general  
 uage for General Purpose). Specialized lexicography focuses on the design,  
 specialized dictionaries, i.e. dictionaries that are devoted to a  
 guistic and factual elements of one or more specialist subject fields,  
 ictionary is usually called a specialized dictionary or LSP dictionary  
 ized dictionaries are either multi-field, single-field or sub-field

 Run!

<sentence id="en847b0-en80d1f"> [Show XML]

word	pos	msd	lemma
Lexicography	noun	NN	lexicography






# ORP Annotation Laboratory

Hide advanced settings

Run!

Corpus name:

Root tag Structural attributes

Tag:

Show Makefile

Show JSON Settings

XML

Show tags

Run!

<sentence id="en847b0-en80d1f"> [Show XML]

word	pos	msd	lemma
Lexicography	noun	NN	lexicography
is	verb	VBZ	be
divided	verb	VBN	divide
into	preposition / subordinating conjunction	IN	into
two	numeral	Z	2
separate	adjective	JJ	separate
but	coordinating conjunction	CC	but
equally	adverb	RB	equally
important	adjective	JJ	important
groups	noun	NNS	group
:	punctuation	Fd	:
Practical	adjective	JJ	practical

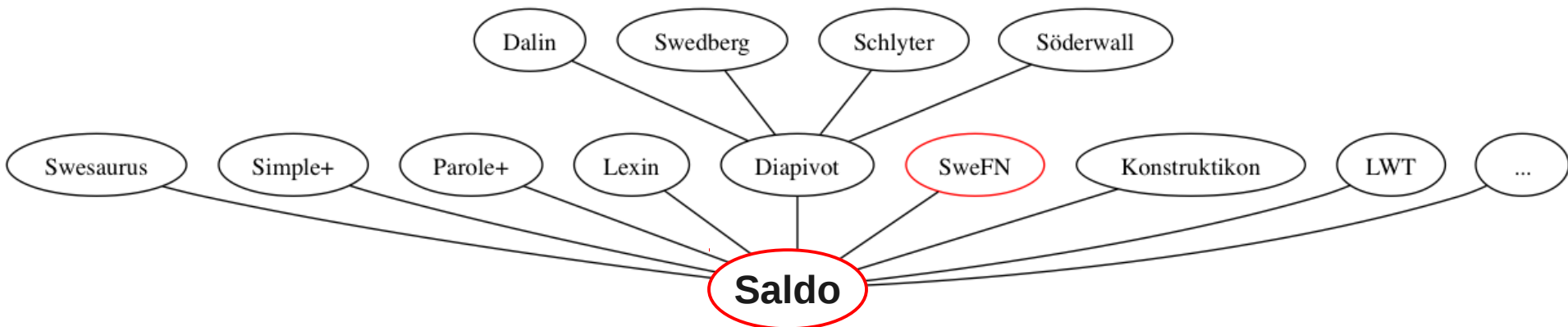


# KARP - Lexical infrastructure of SB

Modern (13)

Historical (9)

Diachronic pivot



All resources (24) are **inter-linked** via Saldo, a Swedish association lexicon, that functions as a pivot for all lexicon resources.

**Saldo** contains around 130 000 entries, is organized by sense and each entry has a persistent identifier (PID).

**Diapivot** is a mapping between Saldo entries and entries (+ inflectional variants) in historical lexicons [SKOS principle; 20 names for *hens*; or “*bakvade*”-example ]

<http://spraakbanken.gu.se/karp>



# KARP

24 lexical resources selected (all) — 720,246 lexical entries

[Search history](#) ▾

Simple Extended Editor

tsunami (substantiv)

Search

through Diapivot

Hits per page: 25 ▾

"kan idag verkligen inte fatta att det var så stort att få sitt första barn att man missar hela **tsunami** katastrofen"

Full Hits (3)

Statistics

Saldo

Saldos morfologi

Lexin

## Saldo (1)

Sense	Lemgram	POS	Primary	Secondary	Children (primary)	Children (secondary)
<i>tsunami</i>	<i>tsunami (noun)</i>	noun	<i>flodvåg</i>	<i>förödande</i>	<i>tsunamikatastrof</i>	<i>tsunamivarning varning</i>
	(21399)					

## Saldo morphology (1)

Lemgram	POS	Inherent	Paradigm	Morphology
<i>tsunami (noun)</i>	noun	u	nn_vu_safari	<i>tsunami (nn_vu_safari) ...</i>
(21399)				


## Lexin (1)



# KARP

[Log in](#) [Download resources](#) [Cite Karp](#) [Svenska](#) | [English](#) [About Karp](#) 

24 lexical resources selected (all) — 720,246 lexical entries

Search history 

Simple

Extended

Editor 


Search for entries that match

or  equals



Sort by

Search

 Hits p

Phonetic form

Wordform

Baseform

Msd

Phonetic form

Saldo Sense

Lemgram

Part of speech

Primary descriptor

Secondary descriptor

Grammar note

Text in definitions

Text in examples

Paradigm

Cefr

Form information

Raw frequency

Wpm

Source

Rank



Simple

Extended

Editor

Search for entries that match



Phonetic form

equals

çö:k



or



Sort by Wordform

ascending

Search



Hits per page: 25



Full Hits (3)



Statistics

## Lexin (3)

Baseform	Phonetic form	POS	Rank	Lexin-ID	Morphology	Gram
kök	çö:k	noun	1637	1078455	kök ...	

[▶ translations](#)

rum där man lagar mat [Lexintema 10 11 ]



en lägenhet på två rum och kök  
det kom härliga dofter från köket  
köks|bord






# KARP – editing



Password protected



Goal: to edit all resources


  **Name** dubbelimperativ

  **Illustration** Fortsätt ställ frågor i frågestunden HÄR!

 Berkeley ID




**Type** ▾   konstruktion

  direktiv huvudsats



  **Category** S

 FrameNet

  **Definition** Konstruktionen omfattar [verb]v som inte kräver infinitivmärke i infinitivfras, i samordning. Konstruktionen har vardaglig prägel. [Det inledande verbet]Modifierarverb har modifierande funktion. 

  **Structure** [V<sub>1,imp</sub> V<sub>2,imp</sub>]

**Inheritance** ▾   pseudosamordning







**Keywords** ▾ 

**Common words** ▾   fortsätta<sup>1</sup>

  sluta<sup>1</sup>





**Construction elements, internal** ▾   name=Modifierarverb, cat=V, msd=imp, other=Modifierarverb

  name=V, cat=V, msd=imp




**Construction elements, external** ▾ 

  **Examples** ▾   [[Fortsätt]Modifierarverb [sälj]v]dubbelimperativ biljetter SJ – och sluta släng av barnen !

  [[Sluta]Modifierarverb [skriv]v]dubbelimperativ så flummigt och redogör hur det verkligen ligger till , tack .

LEMGRAM

fisk (noun)  
 384421

PART OF SPEECH

noun

INHERENT PARADIGM

u nn\_2u\_sten

BÖJNING

fisk ...

sg indef nom	<b>fisk: 215</b>
sg indef gen	<b>fisks: 2</b>
sg def nom	<b>fisken: 84</b>
sg def gen	<b>fiskens: 3</b>
pl indef nom	<b>fiskar: 79</b>
pl indef gen	<b>fiskars: 0</b>
pl def nom	<b>fiskarna: 45</b>
pl def gen	<b>fiskarnas: 2</b>
ci	<b>fisk-: 1 fisk: 215</b>
cm	<b>fisk-: 1 fisk: 215</b>
sms	<b>fisk-: 1</b>

Statistics over form usage in Korp

Number of lemgram occurrences in Korp

fisk<sup>2</sup> (noun)  
 294213

noun

u

nn\_0u\_månsing


fisk ...

som fisk i vattnet (multiword adverb)  
 8

multiword adverb

abm\_i\_till\_exempel

som fisk i vattnet (oböjl.)

som en fisk (multiword adverb)  
 3348

multiword adverb

abm\_i\_till\_exempel

som en fisk (oböjl.)

som en fisk på torra land (multiword adverb)  
 135

multiword adverb

abm\_i\_till\_exempel

som en fisk på torra land (oböjl.)



1 lexicons chosen

Simple Search

Extended Search

Haus

Search

Hits per page: 25

Statistics...



Overview

Hits (1)

Language Comparis

Page: 1 / 1

LWT 1 HIT

SALDO-SENSE LWT-ID UTTRYCK DEF

hus S07.120 Haus deu

Sök på alla språk

svenska

engelska

tibetanska

danska

tyska

franska

finska

grekiska

hindi

italienska

kotgarhi

marathi

nepali

nederländska

portugisiska

spanska

ryska

telugu

jiddish

Search for items in  
other languages





24 lexicons chosen

Simple Search

Extended Search

hus

Search

Hits per page: 25

Statistics...



Overview

Hits (725)

Language Comparison

svenska engelska tibetanska danska tyska franska finska grekiska hindi italienska kotgarhi

marathi nepali nederländska portugisiska spanska ryska telugu jiddish



SWE	ENG	DEU	FRA	FIN	HIN	ITA	NLD	POR	SPA	RUS
hus	the house, house	Haus	maison	talo	घर, g <sup>h</sup> ar	casa	huis, huis	casa	casa	дом



24 lexicons chosen ▾

Simple Search

Extended Search

Find entries where    or  + *or*

Add term

Search

Reset

Hits per page:

Statistics...

Overview

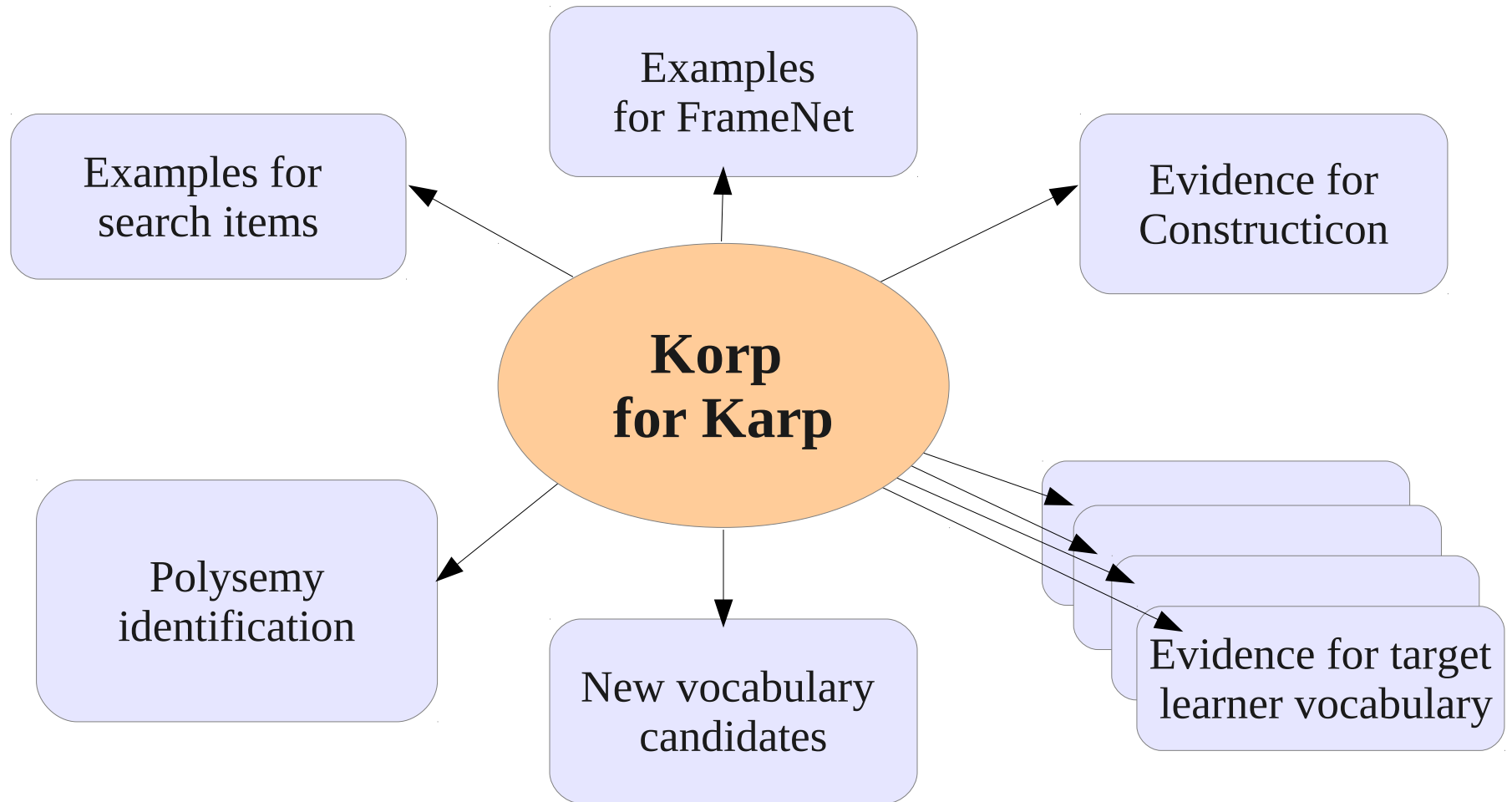
Hits (65)

Language Comparison



SWE	fyra	fem
ENG	four	five
BOD	bźi	lña
DEU	vier	fünf
FIN	neljä	viisi
NEP	cāra	pāñca
RUS	четыре	пять







# Korp & Karp in



- Corpus examples
- Swe FrameNet
- Saldo morphology
- Lexin
- Kelly
- Wikipedia
- Wiktionary

## research - Advanced English

### research

► n

1. systematic investigation to establish facts

**Hypernyms:** [investigation](#)

**Hyponyms:** [operations research](#), [field work](#), [marketing research](#), [microscopy](#), [probe](#), [scientific research](#)

2. **inquiry, enquiry** - a search for knowledge

*"their pottery deserves more research than it has received"*



# SaldoM exercise: training inflection patterns



Inflection paradigm  
(SaldoM as source)

Reference materials

Language learners Train inflections, multiple-choice General purpose vocabulary 1 of 6 learner levels selected  
e.g. Saldo morphology

## Korp sentences

Fully automatic

self-study mode test mode timed test **Generate**

**Result Tracker**

Exercise name	Correct/Total
Learners/inflections, självstudier	1/2

**Train inflections, multiple-choice**

Choose an appropriate word form for the gap

Nr	Sentence
2	Ännu en skatteförvaltning har slarvat bort _____ .
1	Husfasaderna pryds av enorma <b>målningar</b> med tecknade figurer som motiv .

Select a word

- deklaration
- deklarationen
- deklarationens
- deklarationer
- deklarationerna
- deklarationers

Correct answer: **målningar** ✓

**MÅLNING**  
 lemgram: målning..nn.1;  
 part-of-speech: nn;  
 saldo-paradigm: nn\_2u\_mening;  
 inherent: u;

sg indef målning  
 nom

sg indef målnings  
 gen

sg def målningen  
 nom

sg def målnings  
 gen

pl indef målningar  
 nom



# ORP & KARP export

We are exporting our tools, if anyone is interested in using them for their corpora or lexica. They are free of charge, and we help with installation.



Thank you!