

Korp and Karp – a bestiary of language resources: the research infrastructure of Språkbanken

Malin Ahlberg, Lars Borin, Markus Forsberg, Martin Hammarstedt, Leif-Jöran Olsson, Olof Olsson, Johan Roxendal, Jonatan Uppström

Språkbanken, Dept. of Swedish, University of Gothenburg, Sweden

{malin.ahlberg, lars.borin, markus.forsberg, martin.hammarstedt,
leif-joran.olsson, olof.olsson.2, johan.roxendal,
jonatan.uppstrom}@svenska.gu.se

Abstract

A central activity in Språkbanken, an R&D unit at the University of Gothenburg, is the systematic construction of a research infrastructure based on interoperability and widely accepted standards for metadata and data. The two main components of this infrastructure deal with text corpora and with lexical resources. For modularity and flexibility, both components have a backend, or server-side part, accessed through an API made up of a set of well-defined web services. This means that there can be any number of different user interfaces to these components, corresponding, e.g., to different research needs. Here, we will demonstrate the standard corpus and lexicon search interfaces, designed primarily for linguistic searches: *Korp* and *Karp*.

Keywords: Swedish, corpora, lexical resources, research infrastructure.

1 The research infrastructure of Språkbanken

Språkbanken <<http://spraakbanken.gu.se/eng/start>> is a research and development unit at the University of Gothenburg, and a node in a cross-disciplinary research collaboration at the University of Gothenburg and Chalmers University of Technology formalized under the name of Centre for Language Technology <<http://www.clt.gu.se>>. Språkbanken was established with government funding already in 1975 as a national centre. The main focus of Språkbanken's present-day activities is the development and refinement of language resources and language technology (LT) tools, and their application to research in language technology, in linguistics, and in several other disciplines, notably text-based research in the humanities, social sciences, medicine and health sciences.

The larger context of these activities is the systematic construction of a research infrastructure based on interoperability and widely accepted standards for metadata and data. The two main components of this infrastructure deal with text corpora and with lexical resources. For modularity and flexibility, both components have a backend, or server-side part, accessed through an API made up of set of well-defined web services. This means that there can be any number of different user interfaces to these components, corresponding, e.g., to different research needs. Here, we will focus on the standard corpus and lexicon search interfaces *Korp* and *Karp* – designed primarily for linguistic searches – but the same backend web services are also used, e.g., in a corpus-driven grammar and vocabulary exercise generator (Volodina et al., 2012).

2 The search interface of Korp

The search interface of Korp (Borin et al., 2012b; <<http://http://spraakbanken.gu.se/korp>>) has been inspired by corpus search interfaces such as SketchEngine (Kilgarriff et al., 2008), Glossa (Nygaard et al., 2008), and DeepDict (Bick, 2009).

At first glance, the search interface of Korp is a concordance tool that displays search results in the standard KWIC (*keywords in context*) layout (figure 1), here with the example word *svininfluensa* (*noun*) 'swine flu'. This basic functionality is extended by various visualisations of statistical data, such as the basic table (figure 2) and an interactive trend diagram (figure 4) plotting relative frequency over time. Furthermore, the interface features so called word pictures that provides an overview of a selected set of syntactic relations for a word (figure 3). The purpose is to quickly gain an understanding of the contexts in which a word most commonly appears.

We are also working on increasing the diachronic coverage of the corpora, by including Swedish texts from the 19th century back to the 13th century. Ultimately, our goal is to develop tools for all types of text, at various levels of annotation, such as part-of-speech, morphosyntactic information, and dependency parses (Borin et al., 2010; Borin and Forsberg, 2011; Adesam et al., 2012). Our primary source material for Old Swedish (ca 1225–1526) comes from *Fornsvenska textbanken*, <<http://project2.sol.lu.se/fornsvenska>> a 3 MW collection of around 160 digitized texts, mainly from the 13th to the 16th century. Further, a 1 MW corpus of medieval letters from the Swedish National Archives <<http://riksarkivet.se>> is available. Work in progress concerns newspaper texts (17th–19th century) and a collection of law texts (13th century – present).

A number of issues are problematic for annotation of historical texts. For example, sentence splitting cannot be handled with standard tools, as sentence boundaries are often not marked by punctuation or uppercase letters. Compared to modern Swedish texts, the Old Swedish texts have a different vocabulary and richer morphology, show a more free word order, and Latin and German influences. Finally, the lack of a standardized orthography results in a wide variety of spellings for the same word.



118 corpora selected — 1,253,702,515 tokens

Search history

Simple Extended Advanced

Search for Search also as initial part final part and case-insensitive

Related words

fågelinfluensa flunsa asiat såsångsinfluensa influensaliknande influensa

KWIC: Statistics:

KWIC Statistics Word picture

Results: 7,427

ÅBO UNDERÅRTELESER 2012	
ÅRSTÄMMINGEN 2012	
Det var strax före jul den hösten då	svinfluensan grass
Efter de stora nordiska vaccinationerna mot	svinfluensa , med
En 22-årig svensk man är svårt sjuk i	svinfluensa .
Tio miljoner människor från EUs länder har vaccinerat sig mot	svinfluensa .
En 22-årig svensk man är svårt sjuk i	svinfluensa komr
örsta gången som någon smittad i Sverige har blivit allvarligt sjuk i	svinfluensa .
En 22-årig svensk man är svårt sjuk i	svinfluensa .
Ännu en svensk har dött i sjukdomen	svinfluensa .
Mannen hade inga sjukdomar när han fick	svinfluensa .
ollsturneringen Gothia Cup i Göteborg har också blivit smittade av	svinfluensa .
Förra året dog 30 svenskar sedan de blivit sjuka i	svinfluensa .
Flera miljoner svenskar har vaccinerats mot	svinfluensan .
Den här veckan skulle Sverige ha fått mer vaccin mot	svinfluensa .
Ireas Heddini tror att omkring femtio svenskar kan komma att dö i	svinfluensa .

Corpus

Åbo Underrättelser 2012

text attributes

date: 2012-12-11

word attributes

part-of-speech: noun

baseform:

svinfluensa

lemgram:

svinfluensa (noun)

sense:

svinfluensa

initial part:

svin (noun)

final part:

influensa (noun)

dependency relation: Other subject

msd: NN.UTRSIN.DEF.NOM

Show Dependency Tree

Figure 1: The Korp KWIC view of *svinfluensa (noun)* ‘swine flu’

Hit	Total	Åbo Und...	Åbo Und...	Astra No...	Ålandsti...	8 SIDOR	Bloggmix...	Bloggmi
<input checked="" type="checkbox"/> Σ	5.9 (7,427)	0.8 (1)			1.0 (1)	100.2 (68)		19.4 (8)
<input type="checkbox"/> svinfluensan	3.0 (3,818)	0.8 (1)				47.1 (32)		9.7 (4)
<input type="checkbox"/> svinfluensa	1.4 (1,785)				1.0 (1)	45.7 (31)		4.8 (2)
<input type="checkbox"/> Svininfluensan	0.7 (878)					7.4 (5)		2.4 (1)
<input type="checkbox"/> Svininfluensa	0.3 (416)							
<input type="checkbox"/> svinfluensavaccinet	0.1 (95)							
<input type="checkbox"/> svinfluensavaccin	0.0 (51)							
<input type="checkbox"/> svinfluensans	0.0 (49)							2.4 (1)
<input type="checkbox"/> svinfluensaviruset	0.0 (26)							
<input type="checkbox"/> SVININFLUENSAN	0.0 (22)							
<input type="checkbox"/> Svininfluensans	0.0 (17)							
<input type="checkbox"/> svinfluensavirus	0.0 (16)							
<input type="checkbox"/> svinfluensasprutan	0.0 (15)							

Figure 2: The Korp statistics of *svinfluensa (noun)* ‘swine flu’ (including compounds)

3 The search interface of Karp

The interface of Karp (Borin et al., 2012a; <<http://http://spraakbanken.gu.se/karp>>) supports searching and editing lexical resources. It currently hosts 21 lexical resources, some of which have been created from scratch using existing free resources, both external and in-house. The resources have been converted to the Lexical Markup Framework format (ISO, 2008) to ensure uniformity and interchangeability. The infrastructure has one primary lexical resource, SALDO (Borin and Forsberg,

Preposition	Pre-Modifier	svininfluensa	Post-Modifier	Svininfluensa	verb	Verb	svininfluensa
1. mot	818	1. mexikansk 9	1. i Sverige 22	1. sprida 59		1. få 157	
2. om	381	2. aktuell 8	2. med argument 5	2. sprida sig 25		2. ha 158	
3. av	593	3. jävlig 8	3. i Norge 7	3. vara 228		3. vaccinera 9	
4. kring	54	4. dödlig 4	4. förpanikartad 3	4. härja 13		4. skita 9	
5. för	244	5. ny 18	5. påvis 5	5. drabba 14		5. kalla 14	
6. på grund av	19	6. eventuell 4	6. i län 5	6. smitta 9		6. stoppa 11	
7. upp mot	8	7. farlig 4	7. i land 7	7. nå 17		7. dra på sig 4	
8. pga	8	8. bli kvitt 2	8. som vara 15	8. slå 20		8. dra på 4	
9. utav	6	9. resistent 2	9. som smitta 3	9. mutera 7		9. slippa 6	
10. inför	14	10. kvitt 2	10. i höst 4	10. klassa 8		10. undvika 5	
11. för hand	3	11. oförarglig 2	11. under pandemi 2	11. komma 44		11. ta 17	
12. angående	5	12. sibirisk 2	12. till influensa 2	12. skörda 7		12. smitta 3	
13. för sig	7	13. livsfarlig 2	13. som sprida 3	13. fortsätta 13		13. diskutera 6	
14. emot	6	14. inrikes 2	14. på gris 2	14. influensa 4		14. sprida 5	
15. p.g.a.	1	15. jäkla 3	15. i Halland 3	15. bryta ut 6		15. hantera 4	

Figure 3: The word picture of *svininfluensa* (noun) ‘swine flu’

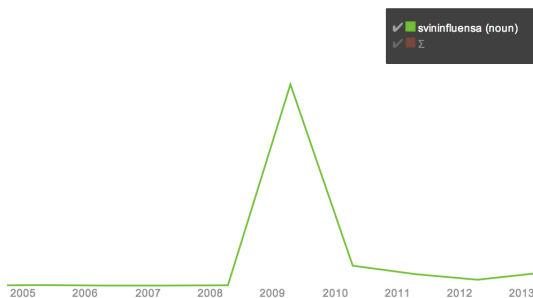


Figure 4: The trend diagram of *svininfluensa* (noun) ‘swine flu’

The screenshot shows the SALDO search results for 'influensa'. The search bar contains 'influensa (substantiv) > \'. The results section shows 'Fullständiga träffar' and 'Träfflistning'. The 'Saldo (1)' section lists related terms: 'asiaten', 'influensaliknande', 'fluensa', 'influensafall', 'influensamedicin', and 'influensadödsfall'. The 'Saldos morfologi (1)' section shows the morphological breakdown of 'influensa (substantiv)'. A tooltip indicates that there are 8733 hits for 'Korp' and provides a link to click for a corpus search.

Figure 5: The search result of *influensa* ‘flu’

2009), which acts as a pivot to which all other modern resources are linked. SALDO is a large freely available morphological and lexical-semantic lexicon for modern Swedish. Moreover, there is a

diachronic pivot resource with links between the modern and the historical morphologies.

In the simple search a user can input either a word form, a lemgram (a form unit), or a sense unit, and the interface will render all information associated to all sense units related to the input. E.g., figure 5 displays the search result of the lemgram *influenta* (*noun*) ‘flu’. In the extended search the user can combine available filters from drop down boxes which are translated to SRU/CQL expressions. E.g. a word form as regular expression and a certain part of speech type.

In addition, the interface supports full text search in the textual parts of the lexical resources, such as examples and definitions. The full text search, beyond extending the search capabilities, also makes the lexical information not linked to SALDO discoverable.

References

- Adesam, Y., Ahlberg, M., and Bouma, G. (2012). *bokstaffua, bokstaffwa, bokstafwa, bokstaua, bokstawa*. . . Towards lexical link-up for a corpus of Old Swedish. In *Proceedings of LTHist 2012*.
- Bick, E. (2009). A graphical corpus-based dictionary of word relations. In *Proceedings of NODALIDA 2009. NEALT Proceedings Series Vol. 4*, Odense. NEALT.
- Borin, L. and Forsberg, M. (2009). All in the family: A comparison of SALDO and WordNet. In *Proceedings of the Nodalida 2009 Workshop on WordNets and other Lexical Semantic Resources – between Lexical Semantics, Lexicography, Terminology and Formal Ontologies*, Odense. NEALT.
- Borin, L. and Forsberg, M. (2011). A diachronic computational lexical resource for 800 years of Swedish. In *Language technology for cultural heritage*, pages 41–61. Springer, Berlin.
- Borin, L., Forsberg, M., and Kokkinakis, D. (2010). Diabase: Towards a diachronic blark in support of historical studies. In *Proceedings of LREC 2010*.
- Borin, L., Forsberg, M., Olsson, L.-J., and Uppström, J. (2012a). The open lexical infrastructure of Språkbanken. In *Proceedings of LREC 2012*, pages 3598–3602, Istanbul. ELRA.
- Borin, L., Forsberg, M., and Roxendal, J. (2012b). Korp – the corpus infrastructure of Språkbanken. In *Proceedings of LREC 2012*, pages 474–478, Istanbul. ELRA.
- ISO (2008). Language resource management – lexical markup framework (lmf). International Standard ISO 24613:2008.
- Kilgarriff, A., Rychlý, P., Smrž, P., and Tugwell, D. (2008). The Sketch Engine. In Fontenelle, T., editor, *Practical Lexicography: A Reader*, pages 297–306. Oxford University Press, Oxford.
- Nygaard, L., Priestley, J., Nøklestad, A., and Johannessen, J. B. (2008). Glossa: a multilingual, multimodal, configurable user interface. In *Proceedings of the Sixth International Language Resources and Evaluation Conference (LREC’08)*, Marrakech. ELRA.
- Volodina, E., Borin, L., Loftsson, H., Arnbjörnsdóttir, B., and Leifsson, G. Ö. (2012). Waste not, want not: Towards a system architecture for icall based on nlp component re-use. In *Proceedings of the SLTC 2012 workshop on NLP for CALL, Lund, 25th October, 2012*, pages 47–58.