



GÖTEBORGS  
UNIVERSITET

Språk  
BANKEN

CLT

# Diabase: Towards a diachronic BLARK in support of historical studies

Lars Borin, Markus Forsberg, Dimitrios Kokkinakis

Språkbanken • Centre for Language Technology  
University of Gothenburg

LREC 2010, Valletta, 19th May, 2010



GÖTEBORGS  
UNIVERSITET

**Språk**  
BANKEN

CLT

# topics

1. text and speech as historical research data, and language technology
2. our work with LT for historical studies: two examples
3. methodological musings on BLARKs and language variation



GÖTEBORGS  
UNIVERSITET

**Språk**  
BANKEN

CLT

# LT in historical studies

In historical studies, text – and speech, i.e., **language** – are central as both primary and secondary research data sources.

In today's world, the normal mode of access to text, speech, images and video is in digital form. Modern material is born digital and older material is being digitized on a vast scale in cultural heritage and digital library projects.

LT can help historians and other researchers make effective use of this flood of language data from all historical periods.



<http://spraakbanken.gu.se/eng/start/>

GÖTEBORGS  
UNIVERSITET

Språk  
BANKEN

CLT

The screenshot shows a Mozilla Firefox browser window displaying the homepage of Språkbanken (the Swedish Language Bank) in English. The browser's address bar shows the URL <http://spraakbanken.gu.se/eng/start>. The page features a navigation menu with links for 'Research', 'Information', 'Staff', 'Publications', and 'Resources'. A search bar is located in the top right corner. The main content area is divided into three columns:

- Språkbanken (the Swedish Language Bank)**: A paragraph describing the bank's history, its creation in the 1960s, and its current role in collecting and processing Swedish text corpora. A [Read more...](#) link is provided.
- Språkbanken tidbits**: A table with columns for 'word', 'description', and 'inflection'. It lists the words 'grönbete', 'bete', and 'gräs', with 'grönbete (nn)' highlighted. Below the table is a 'corpus search' input field.
- Corpora**: A list of corpora including 'PAROLE/SUC', 'Corpus collection', 'ORDAT', and 'Domain specific corpus'.
- Lexical resources**: A list of lexical resources including 'SALDO', 'Söderwall och Schluter', 'Dallins ordbok', 'Svenska ord\_LEXIN', and 'SAOB'.

The footer of the browser window shows the URL <http://www.gu.se/english> and the Zotero logo.



GÖTEBORGS  
UNIVERSITET

Språk  
BANKEN

CLT

# 19th c. fiction in Litteraturbanken

Litteraturbanken.se - Mozilla Firefox

Arkiv Bedigera Visa Historik Bokmärken Verktyg Hjälp

http://litteraturbanken.se/#forfattare/FlygareCarlenE/titlar/PalVarning/sida/3/etext

Mest besökta Welcome to CLT Google Lingvistbloggen Institutionen för ... Language Log The LINGUIST L... OREL - All links by...

Litteraturbanken.se

Litteraturbanken

Hem > författare > FlygareCarlen > titlar > sidan 3 etext

Författare  
Titlar  
Presentationer  
Aktuellt  
Rättigheter  
Om LB  
Kontakt  
Hjälp  
Sök i etext

SAOB :: Dalin

Emilie Flygare-  
Carlén  
Pål Varning

... Bibliografisk info.

« « « » » »

Sida 3

[Start]  
Inledning  
Pål Varning  
Tillkomst och m...  
Forskning

FÖRSTA KAPITLET.

SAKERNAS STÄLLNING FÖRE OCH EFTER PÅLS ANKOMST TILL VERLDEN, SAMT HANS FÖRSTA BESLUT PÅ EGEN HAND.

Seglatsen inomskärs mellan Strömstad och Göteborg utmärker sig genom ett egensinne och en nyckfullhet i lynnet, som mer än en gång under färden pröfvat icke blott passagerarnes, utan äfven de värdige vedskutspatronernas lugna tålmod och bekanta trygghet. Der är ett odrägligt kryssande fram och åter: än motvind och storm, än god vind, men antingen med så hastigt påkommande stiltje, att den goda vinden tjanat till intet, eller ock med så starka orkaner och sjögång, att skutan darrande brakar i sina fogningar och far lik en spån upp och ned på de vreda böljorna. Och alla dessa särskilda fall smyga så tätt och lömskt hvarandra i spåren, att hela resan närmast kan jämföras med ett smuglare-äventyr, der man endast genom tusen beräkningar och svårigheter slingrar sig fram

Klar

zotero



GÖTEBORGS  
UNIVERSITET

Språk  
BANKEN

CLT

# NER in Litteraturbanken

Alla författare	Samtliga titlar	Visa 20 träffar per sida
Visar 1021-1040 av 43126 träffa	Ny sökning	efter förekomst
Platsnamn	Samtliga undertyper	

Berger Bendel & Co s. 284	ig och passerades. Bortom <b>Adams street</b> street blef den oerhörda avenyen
Berger Bendel & Co s. 284	street blef den oerhörda <b>avenyen</b> mörkare, och vid van Burens hörn
Berger Bendel & Co s. 284	en loge. - Kommo de till <b>Paris</b> i sommar så skulle de köpa ett par
Berger Bendel & Co s. 284	hon visste en adress vid <b>Boulevard Males-</b> Malesherbes, där de såldes i
Berger Bendel & Co s. 284	elge någonsin har varit i <b>Paris</b> , tillade hon. Van Buren var mörk och
Berger Bendel & Co s. 284	r som - ligt. Från Lister <b>Building</b> hade Helge sett några skepnader
Berger Bendel & Co s. 285	igt och lyfte lätt på sin <b>panama</b> , det är som vore vi i gay Paree -
Berger Bendel & Co s. 286	Reuter, damerna äro från <b>Norge</b> eller, nej, hur var det? Sverige,
Berger Bendel & Co s. 286	eller, nej, hur var det? <b>Sverige</b> , tror jag. Men, herre gud, det är
Berger Bendel & Co s. 291	nda ned till barriären åt <b>Clark street</b> street och bekväma liggstolar, s
Berger Bendel & Co s. 292	rande, svindlande höjder. <b>Michigansjöns</b> ofantliga yta mötte österut
Berger Bendel & Co s. 293	Snedt till vänster syntes <b>Masonic Temples</b> Temples takkrön som ett strål
Berger Bendel & Co s. 293	hvilken var reflexen från <b>Milwaukee</b> . Men själfva sjön tedde sig som
Berger Bendel & Co s. 294	kajmurar, hörde han också <b>Michigans</b> brusande andhämtning och därefter
Berger Bendel & Co s. 296	och se efter - gossarna i <b>England</b> . Men några dagar hinner jag alltid
Berger Bendel & Co s. 296	g öfver Kanalen. Bor du i <b>Paris</b> på Grand Louvre som förut? - Ja, svarad
Berger Bendel & Co s. 296	analen. Bor du i Paris på <b>Grand Louvre</b> Louvre som förut? - Ja, svarade
Berger Bendel & Co s. 298	hvars syster är gift med <b>Indiens</b> vicekonung. Jag måste gå. Han
Berger Bendel & Co s. 298	talade de om den franska <b>rivieran</b> och spelbanken i Monte Carlo. -
Berger Bendel & Co s. 298	rivieran och spelbanken i <b>Monte Carlo</b> Carlo. - Jag skall hvila ut där i



GÖTEBORGS  
UNIVERSITET

Språk  
BANKEN

CLT

# semantic search in 19th c. fiction

## CONPLISIT – the components:

- ▶ SALDO – a modern semantic lexicon with inflectional morphology (~73.000 senses)
- ▶ Dalin – a large 19th century lexicon (~63.000 lemmas)
- ▶ an orthographic mapping database SALDO–Dalin
- ▶ a morphology for regular 19th c. open parts of speech
- ▶ Litteraturbanken (~100 19th c. novels)
- ▶ a research question: can 19th c. fiction throw light on the emergence of consumer society in Sweden?



GÖTEBORGS  
UNIVERSITET

Språk  
BANKEN

CLT

# Dalin–SALDO round trip

vagn..1 - Mozilla Firefox

Arkiv Bedigera Visa Historik Bokmärken Verktyg Hjälp

http://spraakbanken.gu.se/ws/dalin-ws/md1/html/vagn..1

Mest besökta Welcome to CLT ... Google Lingvistbloggen Institutionen for ... Language Log The LINGUIST L... OREL - All links by...

Dalin | Språkbanken vagn..1

*SweFN++*

korgvagn, kolskrinda, räddningsbåt, haktåg, tröskvagn, triumftåg, hisståg, kulvagn, skärbåt, stupkärra, slåde, aflöpningsståg, gigg, täckvagn, själtåg, sågvagn, brödvagn, trossvagn, likvagn, paketbåt, källslåde, kärr, fröskida, båt, brunsvagn, värjskida, lappsråde, rullvagn, stegvagn, sorgetåg, stegkärra, makrillbåt, barnvagn, kasttåg, bondskrinna, färjbåt, postvagn, snällvagn, långbåt, landå, kappsråde, sjötåg, dragkärra, trilla<sup>2</sup>, blockvagn, falltåg, char, bagarkärra, vippkärra, sytåg, bondkärra, basttåg, lastbåt, tältvagn, kalesch, bombvagn, skjutsvagn, törntåg, forvagn, sorgvagn, kanonvagn, jutvagn, familjevagn, bröllopsvagn, långkärra, iltåg, åkarkärra, släptåg, hofbuss, bergningsbåt, brobåt, verkvagn, sillbåt, krutvagn, fodervagn, kupévagn, ångvagn, strömbåt, källvagn, bakvagn, vagnskorg, bjellersråde, enbetsvagn, surrtåg, triumfvagn, kurvagn, skida, landtåg, buss, hyrvagn, buss<sup>2</sup>, sandkärra, fyrbåt, kursråde, fordon, kolryss, grundtåg, diligens, vurstvagn, låringsbåt, statsvagn, segervagn, portschäs, brudvagn, kaross, sådesvagn, krigsvagn, köksvagn, resvagn, gigtåg, rackarkärra, aflöpningsråde, skid, presskärra, åkning, vipptåg, jagtvagn, redskapsvagn, arbetsvagn, vefbåt, stenvagn, roddarbåt, bredschäs, bondvagn, tågsråde, passagerare, kärra, vagn, borttåg, skrinna, ställningståg, genomtåg, durktåg, borrhvagn, packvagn, schäs, bogtåg, rustvagn

klar zotero





GÖTEBORGS  
UNIVERSITET

**Språk**  
BANKEN

CLT

## some issues encountered

- ▶ the 1906 spelling reform – in reality lasting some decades around 1900
- ▶ large synchronic orthographic variation before the 19th century – and again today!
- ▶ slightly different inflectional morphologies 19th–20th century – gradual abandonment of verb-subject agreement during the first half of the 20th century until it officially went out of use around 1950
- ▶ very different inflectional system (and syntax) in the Old Swedish period
- ▶ changes in vocabulary
- ▶ changes in word meanings



GÖTEBORGS  
UNIVERSITET

**Språk**  
BANKEN

CLT

# methodological musings: what's in a BLARK?

- ▶ linguistically annotated text corpora
- ▶ speech databases
- ▶ tools for basic text and speech processing
- ▶ basic lexical resources
- ▶ tools for linguistic annotation of text (POS taggers, chunkers, parsers)
- ▶ text-to-speech and speech-to-text systems

... in interoperable, standardized formats



GÖTEBORGS  
UNIVERSITET

**Språk**  
BANKEN

CLT

## BLARKs are static, . . .

As it is normally conceived of and presented, the BLARK assumes a modern standard language variety as the object of description, at least as far as the written language part of the BLARK is concerned, which is the part that we are competent to make judgements about. Part of the reason for this is certainly historical: The BLARK has been – and continues to be – informed more than anything else by language technology work on modern stable written standard languages.



GÖTEBORGS  
UNIVERSITET

**Språk**  
BANKEN

CLT

## ... language is dynamic

Modern linguistics increasingly recognizes **variation** as a fundamental and essential characteristic of human language. In this regard, the study of history through textual primary sources makes up an interesting and challenging testbed, where the robustness and the generality of existing language technology are subjected to the acid test of messy and multilingual reality, more so than in many other application areas, since we have to deal with, *inter alia*, historical, non-standardized language varieties in addition to a number of modern standard languages.



GÖTEBORGS  
UNIVERSITET

**Språk**  
BANKEN

CLT

# linguistic variation

Language varies (at least),

- ▶ by community (languages, dialects, sociolects)
- ▶ by subject, purpose or medium (topics, genres)
- ▶ by time (historical language stages)

The BLARK attempts to abstract away from all three; it can be thought of as reflecting a modern standard language, which is topic- and genre-neutral.



GÖTEBORGS  
UNIVERSITET

Språk  
BANKEN

CLT

## describing variation in a BLARK?

Our work described above can be seen as the first steps towards the development of *Diabase*, a Swedish BLARK extended along the diachronic – or time – axis. This work also raises some interesting methodological questions about the description of variation in the linguistic resources making up the BLARK:

- ▶ fundamental/norm(al) ~ deviation (e.g., ‘correct’ vs. ‘incorrect’ verb agreement in a novel published in 1935; correct vs. incorrect spelling in an internet page in 2010)
- ▶ one system ~ a mix of systems (e.g., pre- vs. post-1906 spelling)



thank you for listening!

GÖTEBORGS  
UNIVERSITET

**Språk**  
BANKEN

CLT

