



GÖTEBORGS  
UNIVERSITET

**Språk**  
BANKEN

CLT

# Towards a knowledge-based culturomics

Lars Borin

Språkbanken • Swedish Language • University of Gothenburg

kickoff meeting • 31st Jan 2013



GÖTEBORGS  
UNIVERSITET

Språk  
BANKEN

CLT

“culturomics”

# Quantitative Analysis of Culture Using Millions of Digitized Books

Jean-Baptiste Michel,<sup>1,2,3,4,5\*†</sup> Yuan Kui Shen,<sup>2,6,7</sup> Aviva Presser Aiden,<sup>2,6,8</sup> Adrian Veres,<sup>2,6,9</sup>  
Matthew K. Gray,<sup>10</sup> The Google Books Team,<sup>10</sup> Joseph P. Pickett,<sup>11</sup> Dale Hoiberg,<sup>12</sup>  
Dan Clancy,<sup>10</sup> Peter Norvig,<sup>10</sup> Jon Orwant,<sup>10</sup> Steven Pinker,<sup>5</sup>  
Martin A. Nowak,<sup>1,13,14</sup> Erez Lieberman Aiden<sup>1,2,6,14,15,16,17\*†</sup>

We constructed a corpus of digitized texts containing about 4% of all books ever printed. Analysis of this corpus enables us to investigate cultural trends quantitatively. We survey the vast terrain of ‘culturomics,’ focusing on linguistic and cultural phenomena that were reflected in the English language between 1800 and 2000. We show how this approach can provide insights about fields as diverse as lexicography, the evolution of grammar, collective memory, the adoption of technology, the pursuit of fame, censorship, and historical epidemiology. Culturomics extends the boundaries of rigorous quantitative inquiry to a wide array of new phenomena spanning the social sciences and the humanities.

*(Science 331, 176 (2011))*



# knowing your data, 1

GÖTEBORGS  
UNIVERSITET

Språk  
BANKEN

CLT



("funk" vs. "sunk" in Google n-grams)



# knowing your data, 2

GÖTEBORGS  
UNIVERSITET

Språk  
BANKEN

CLT

Antal träffar: 781

Föregående 1 2 3 4 5 6 7 8 9 10 11 .. 31 32 Nästa Visa kontext

	DALPILSEN
i Björnsberg kl. 4 em. och i gamla c.is, ionfrgkand dat	<b>8ven</b> Skall jö.
	<b>8ven</b> » l'ct 0llw.
	<b>8ven</b> » l l tolk fal!!!
Sondagen » kollekt:	<b>8ven</b> » ka foli.
	<b>8ven</b> Johan LindtrSm eoh B«nhiid E Isabat Bngtrflm, bida fr. Korania
att je, nle det trvelclrideton s, >rtsarande i lieflerin^ » ser, nen sördelällog	<b>8ven</b> » Ka l ' » dket.
	Kollekt till <b>8ven</b> » ka Bofiafrdsamlingen i Paris.
> ,,,,,	<b>8ven</b> » K l, nr, *IKI> la as <!!.
l>>ulites,	<b>8ven</b> » l<.i, 1') sliir ne l> ^' rnn^li. r l ) m<l>^ » Ol » s ,', m<' l l<' i '.
2 mecl unclerlagcl »	<b>8ven</b> » Ka nrä 8amt upplvsancle text, ul^ilvet al ^?.
SMagn kollkt:	<b>8ven</b> ' k »; kyrkans Diako nittyreleei lånefond.
	<b>8ven</b> Johan Salén från Helsingborg, 64 år, 5 rn, 21 d.
	<b>8ven</b> " Bo: i
	<b>8ven</b> Holger, son till Edvard Be' glund, Fortby.
Apr: S. rjuird Gullan Viktoria, dotter till	<b>8ven</b> Gerhard Soneaton, i kålfä, 8 mån
	<b>8ven</b> Olof, gon till 8ven Olof Ölander i Envikabyn.
	<b>8ven</b> Olof, gon till 8ven Olof Ölander i Envikabyn.
navalet i Östra härad, Blekinge, hvilket i tisdags egde rum, återvaldes hemmansegaren	<b>8ven</b> Aruoelssou.
) u » oeb mörk, oräl när oob	<b>8ven</b> » K;
	<b>8ven</b> kl 7812ll<ier,
Hs	<b>8ven</b> » K » legeringen Kelu lImäktig » s Dtvansrings-^gent.
Skrftddertillskäraien,	<b>8ven</b> Johan Nilsson-Nejbert fr. Vagnmakaren 3, 69 år.

(instances of "8ven" i the oldest Digidaily material)



GÖTEBORGS  
UNIVERSITET

Språk  
BANKEN

CLT

## the hitch (Mark Liberman, LL 17 Dec, 2010)

There are also a large number of cases where you'd like to group word-strings into categories: dates, organizations, minerals, place names, novelists, etc., and then treat these categories (rather than words or word strings) as units of analysis. Again, there are well-known techniques for inducing such categories in text collections — but to use these techniques, you need to be able to have the text collection in hand so as to be able to run your algorithms over it.

Many — maybe most — questions about historical texts are like these last few examples: relatively easy to answer if you have a corpus in hand, and not addressed very well (if at all) by a collection of "culturomic trajectories", defined as the year-by-year time-functions of common word sequences. In particular, nearly all questions about the history of the English language fall beyond the grasp of time-functions of n-gram frequencies. This is not to deny the interest and value of such time functions. It's just that they're not nearly enough.



GÖTEBORGS  
UNIVERSITET

Språk  
BANKEN

CLT

## project consortium • strengths

- ▶ Chalmers • machine learning, data mining
- ▶ Lund • semantic parsing and role labeling
- ▶ Gothenburg • language resources, knowledge-based LT



GÖTEBORGS  
UNIVERSITET

**Språk**  
BANKEN

CLT

## our proposal

- ▶ **general:**

- ▶ work on Swedish (and smaller datasets)
- ▶ apply deep linguistic processing to texts
- ▶ combine knowledge-rich and statistical approaches

- ▶ **specifics:**

- ▶ a Swedish Watson
- ▶ discover and track semantic change over time



# why we got funded

GÖTEBORGS  
UNIVERSITET

Språk  
BANKEN

CLT

## Overall grade for the application

6

(7=Outstanding, 6=Excellent, 5=Very good to excellent, 4=Very good, 3=Good, 2=Weak, 1=Poor)

## Motivation for the overall assessment of the application

The panel concluded that this was a highly technical project. The reputations and experience of the team seem matched to the challenge although even this could be a little difficult to judge without more detail as to what will actually be achieved – what success for the challenge is defined to be.

The project therefore would gain by specifying very clear goals or outcomes against which its worth can be considered and its achievements judged. The scientific reviewers agreed that this project falls short on the detail with which it specifies these aspirations. However, the panel were inclined to be supportive to the project overall.

It was noted that in the recent past, Culturomics research was able to retroactively predict the 2011 Arab Spring and successfully estimate the final location of Osama Bin Laden to within 124 miles. This is very impressive, since social science is not usually so precise. It was agreed that finding ourselves in a situation of text-information-overload, more and better applications of semantic processing are in great need. Focusing on the rich treasure of Swedish texts seems adequate and highly relevant for this grant. The researchers also draw from long-standing concepts of linguistics and language technology, which gives it a solid disciplinary grounding. If the researchers achieve to create concrete and user-friendly tool that can be applied by others, the outcomes of this project could have a very high and real impact.





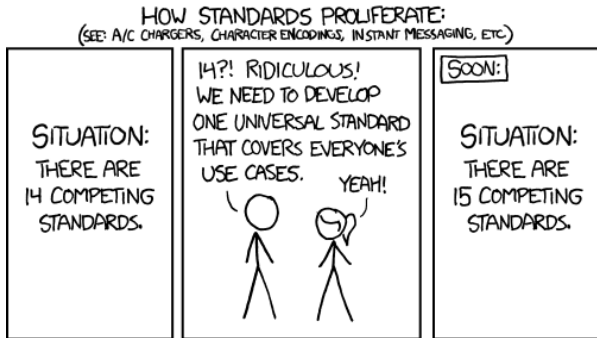
GÖTEBORGS  
UNIVERSITET

Språk  
BANKEN

CLT

# challenges

- ▶ operationalizing and concretizing our goals
- ▶ coordinating our work:
  - ▶ in time
  - ▶ in objectives
  - ▶ in interoperability:



(<<http://http://xkcd.com/927/>>)