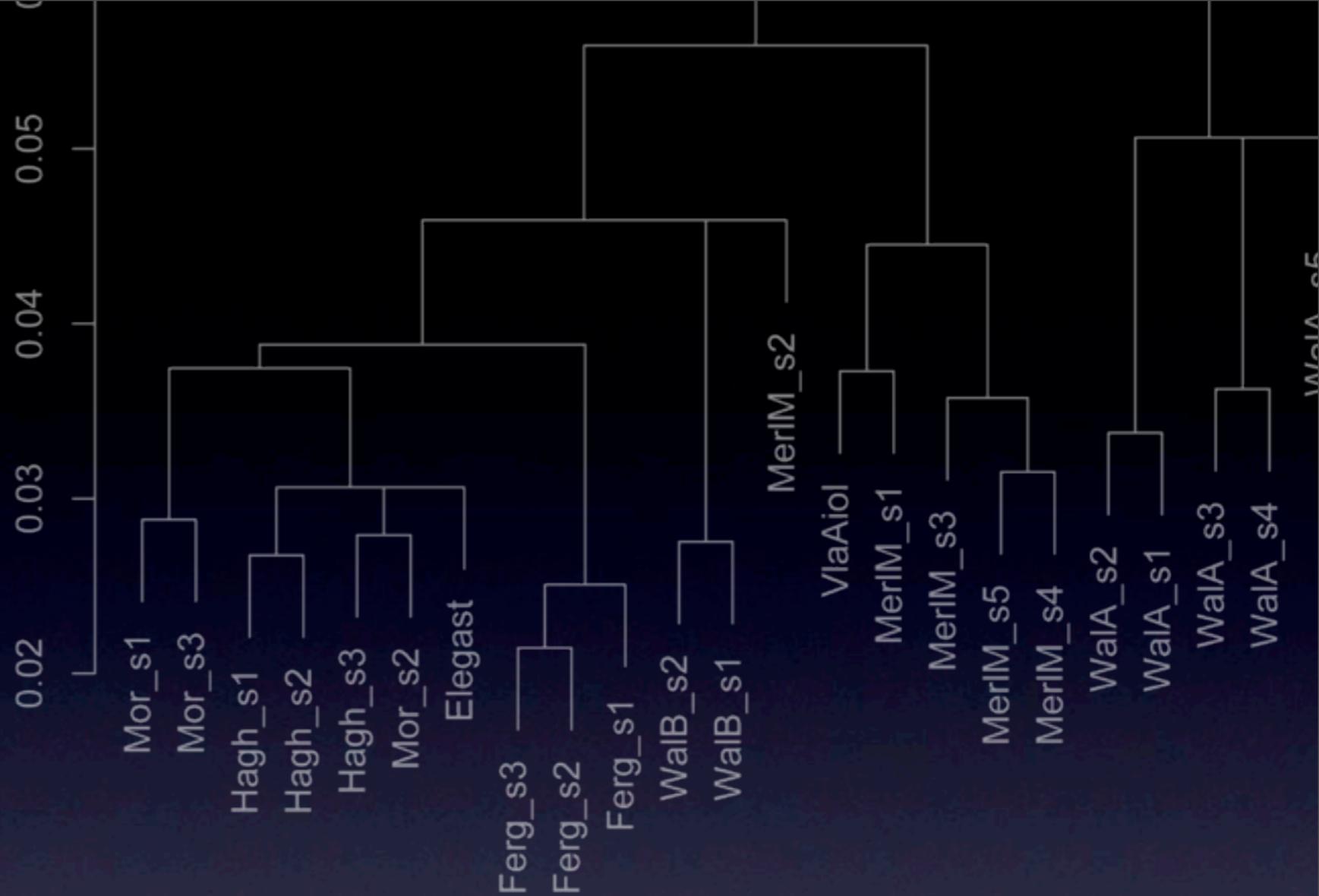


Lezen vanop afstand

Digital Humanities en de computationele analyse van middeleeuwse literatuur

www.mike-kestemont.org | www.fwo.be | www.ua.ac.be



Digital Humanities

- Digitale geesteswetenschappen
- R. **Busa** (1913-2011)
- IBM's *Index Thomisticus*
- Busa award door **ADHO**
- Mediëvist!



Distant reading

- Veelgeciteerde, ‘holle’ term
- **F. Moretti** (2000)
- *Close reading*
- Analyse van grote verzameling teksten “without a single direct textual reading”
- Uitdaging...

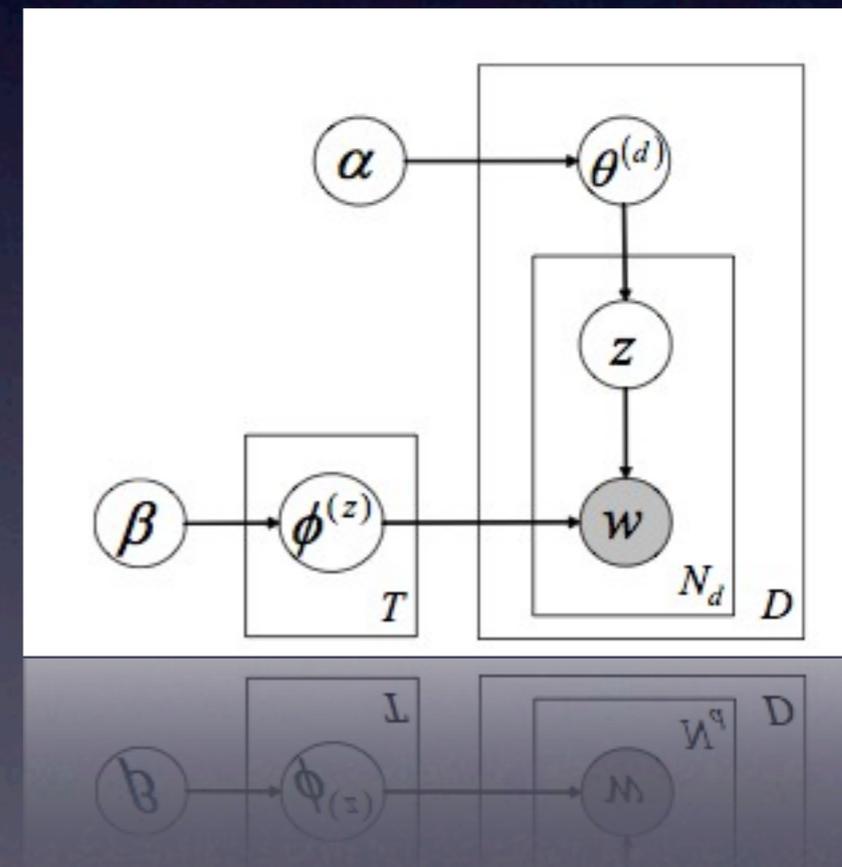


Topic modeling

- Recent veel aandacht in DH
- Automatische detectie van **topics** in groot corpus documenten
- a **synecdoche** of digital humanities. It is distant reading in the most pure sense: focused on corpora and not individual texts, treating the works themselves as unceremonious ‘buckets of words’, and providing seductive but obscure results in the forms of easily interpreted “topics” (Meeks & Weingart 2012)

Latent Dirichlet Allocation

- Dimensie-reductie (multivariate statistiek)
- *Latent Semantic Analysis, Non-negative Matrix Factorization, ...*
- Distributionele semantiek
- Monte Carlo-methode (via Gibbs sampler)
- Niet geometrisch, maar probabilistisch



Topics?

Topic 247

word	prob.
DRUGS	.069
DRUG	.060
MEDICINE	.027
EFFECTS	.026
BODY	.023
MEDICINES	.019
PAIN	.016
PERSON	.016
MARIJUANA	.014
LABEL	.012
ALCOHOL	.012
DANGEROUS	.011
ABUSE	.009
EFFECT	.009
KNOWN	.008
PILLS	.008

Topic 5

word	prob.
RED	.202
BLUE	.099
GREEN	.096
YELLOW	.073
WHITE	.048
COLOR	.048
BRIGHT	.030
COLORS	.029
ORANGE	.027
BROWN	.027
PINK	.017
LOOK	.017
BLACK	.016
PURPLE	.015
CROSS	.011
COLORED	.009

Topic 43

word	prob.
MIND	.081
THOUGHT	.066
REMEMBER	.064
MEMORY	.037
THINKING	.030
PROFESSOR	.028
FELT	.025
REMEMBERED	.022
THOUGHTS	.020
FORGOTTEN	.020
MOMENT	.020
THINK	.019
THING	.016
WONDER	.014
FORGET	.012
RECALL	.012

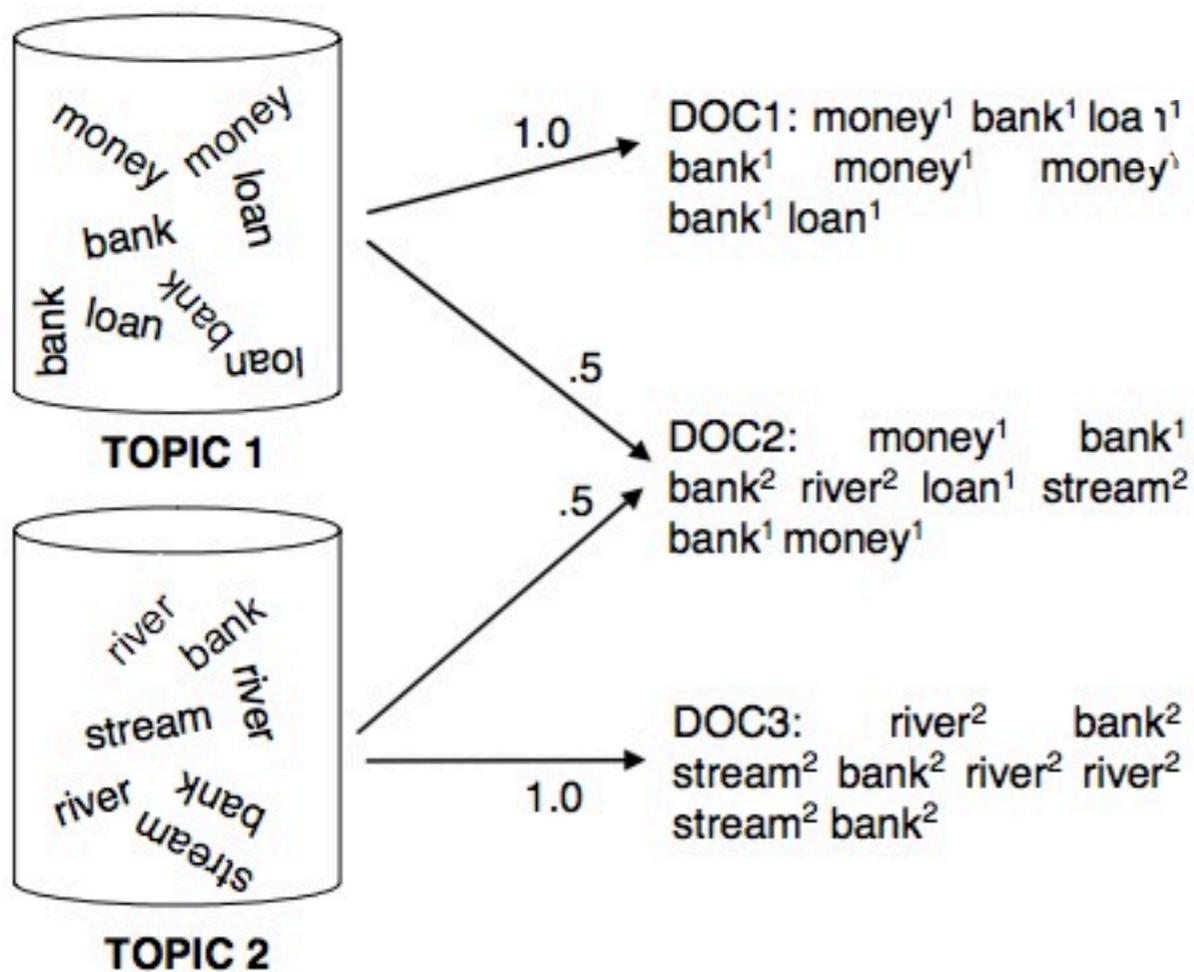
Topic 56

word	prob.
DOCTOR	.074
DR.	.063
PATIENT	.061
HOSPITAL	.049
CARE	.046
MEDICAL	.042
NURSE	.031
PATIENTS	.029
DOCTORS	.028
HEALTH	.025
MEDICINE	.017
NURSING	.017
DENTAL	.015
NURSES	.013
PHYSICIAN	.012
HOSPITALS	.011

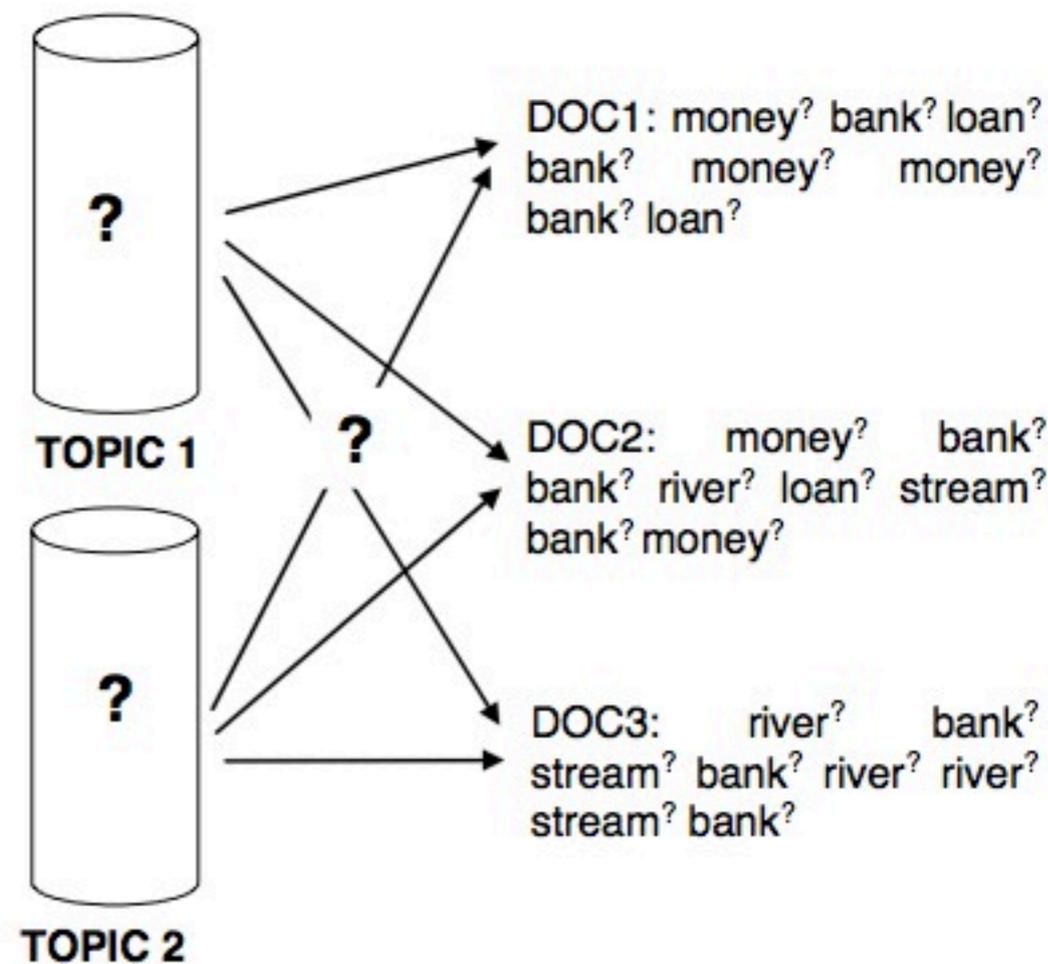
Figure 1. An illustration of four (out of 300) topics extracted from the TASA corpus.

Grabbelton

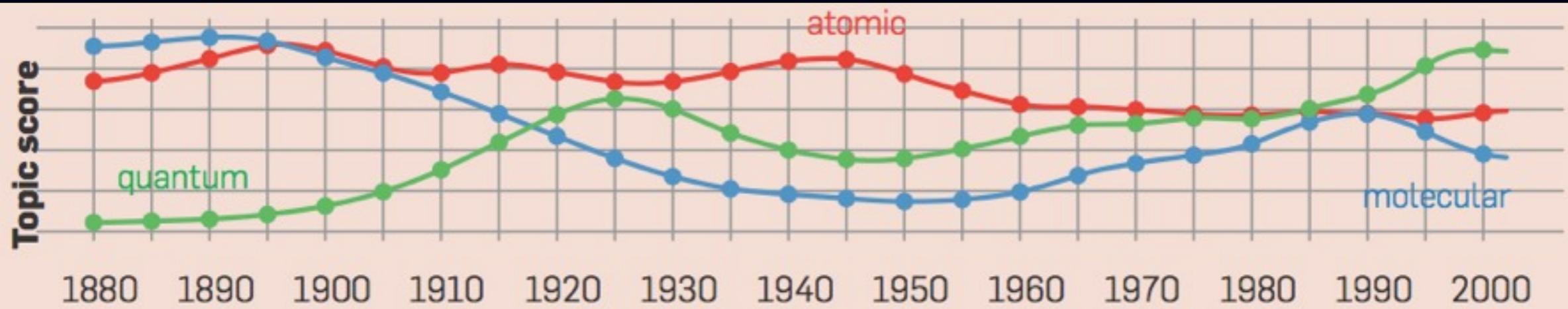
PROBABILISTIC GENERATIVE PROCESS



STATISTICAL INFERENCE

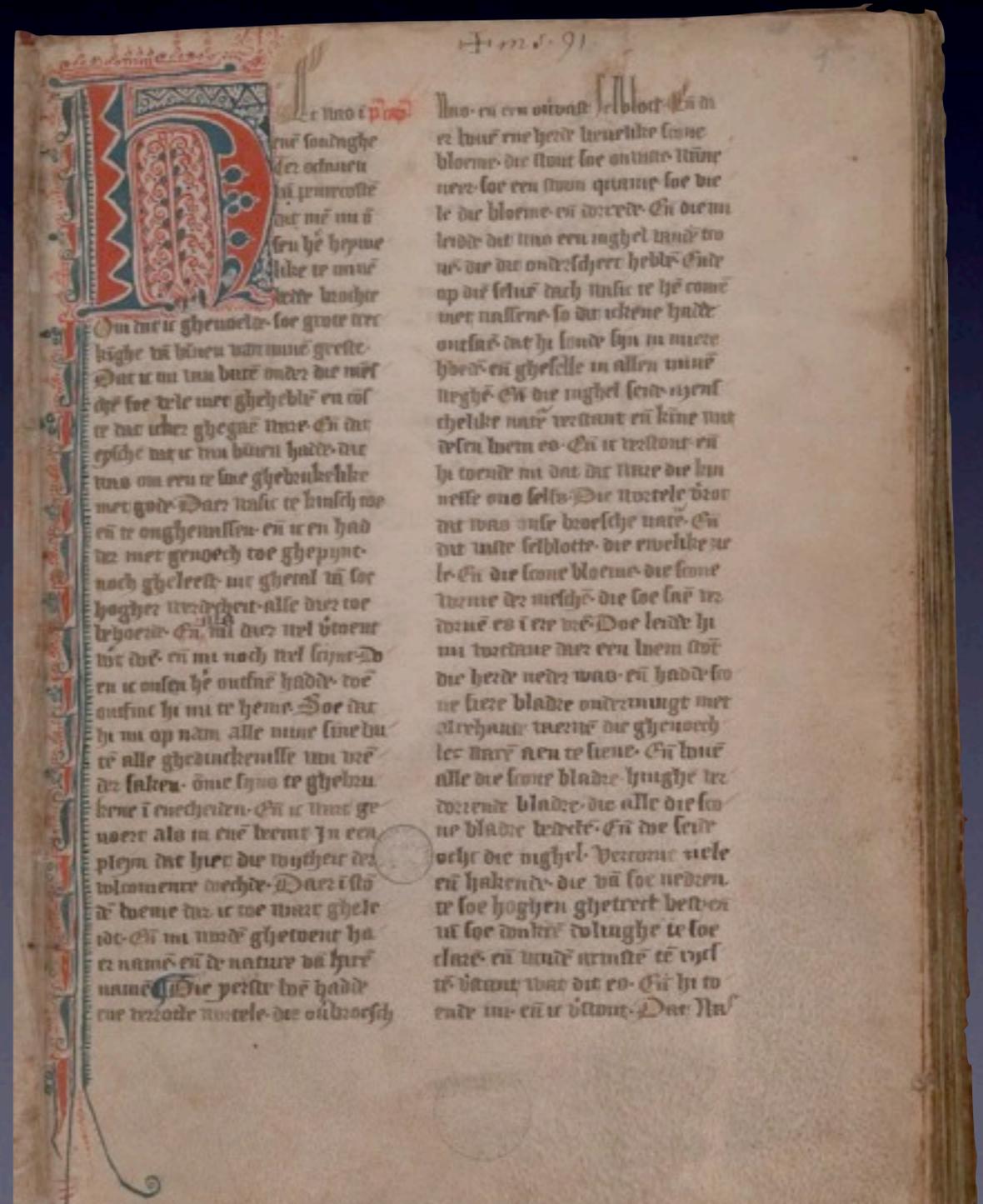


Thematische evoluties

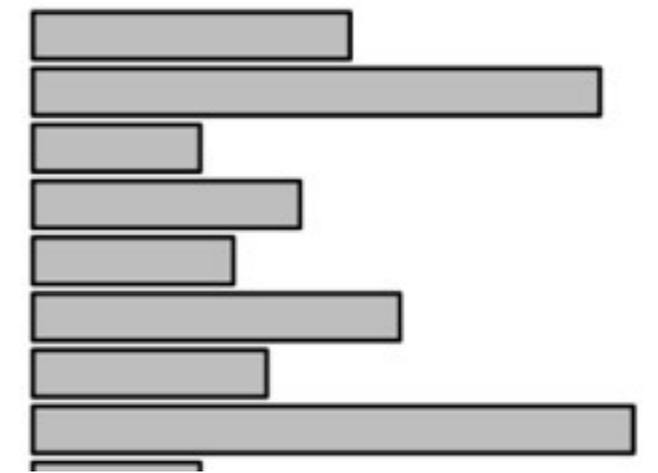


Historische letterkunde?

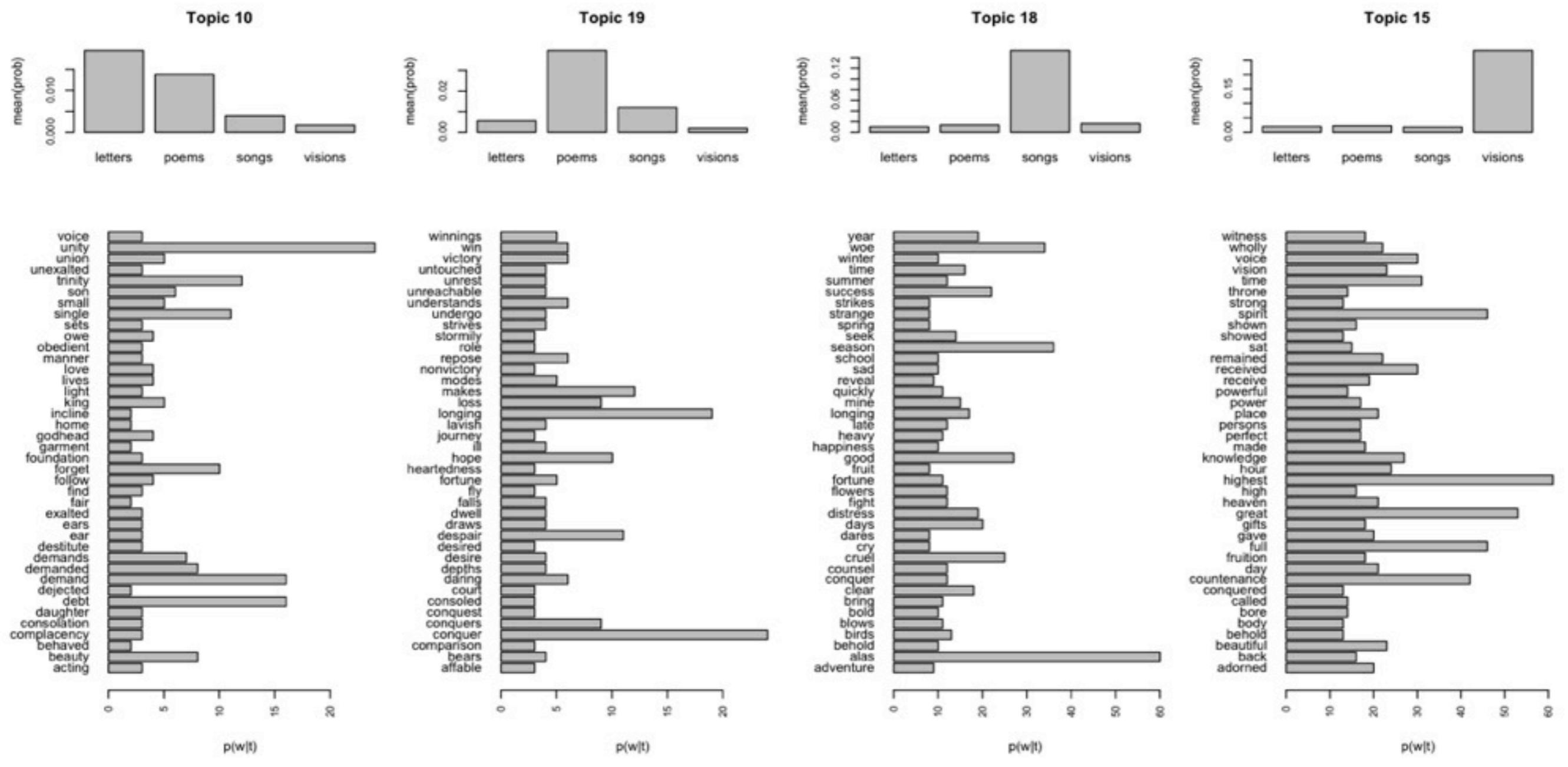
- **Hadewijch** (13e E)
- Brabantse mystica
- Middelnederlands
- 4 genres:
 - visioenen
 - brieven
 - gedichten
 - liederen
- Vert. C. Hart (1980)
- Traditioneel *close reading*

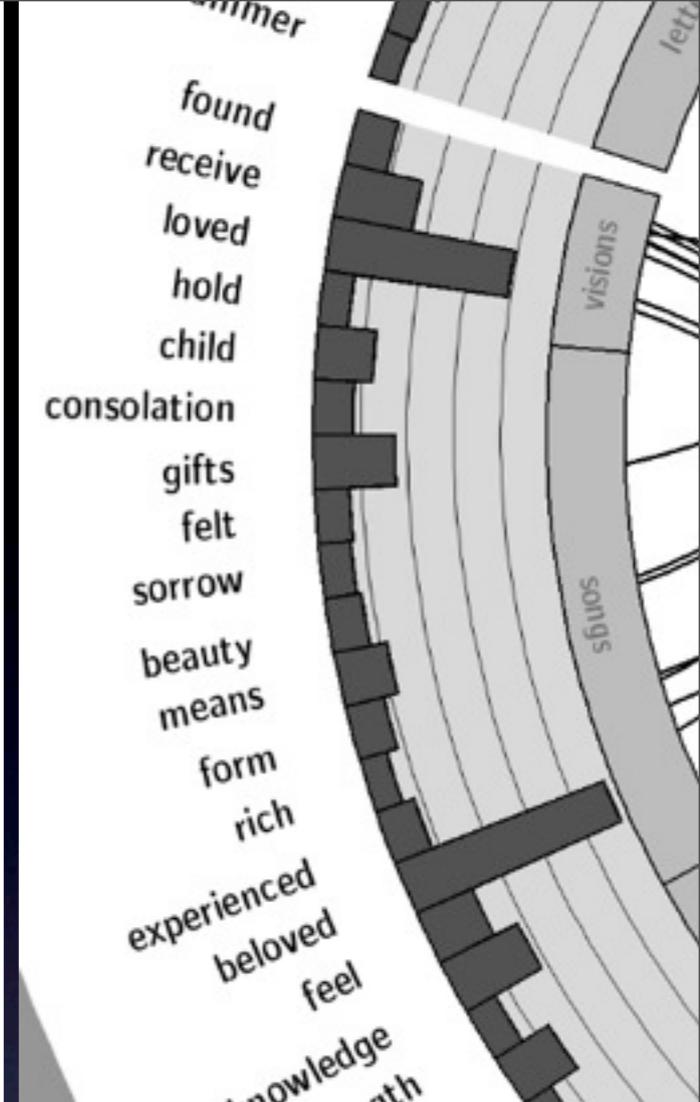
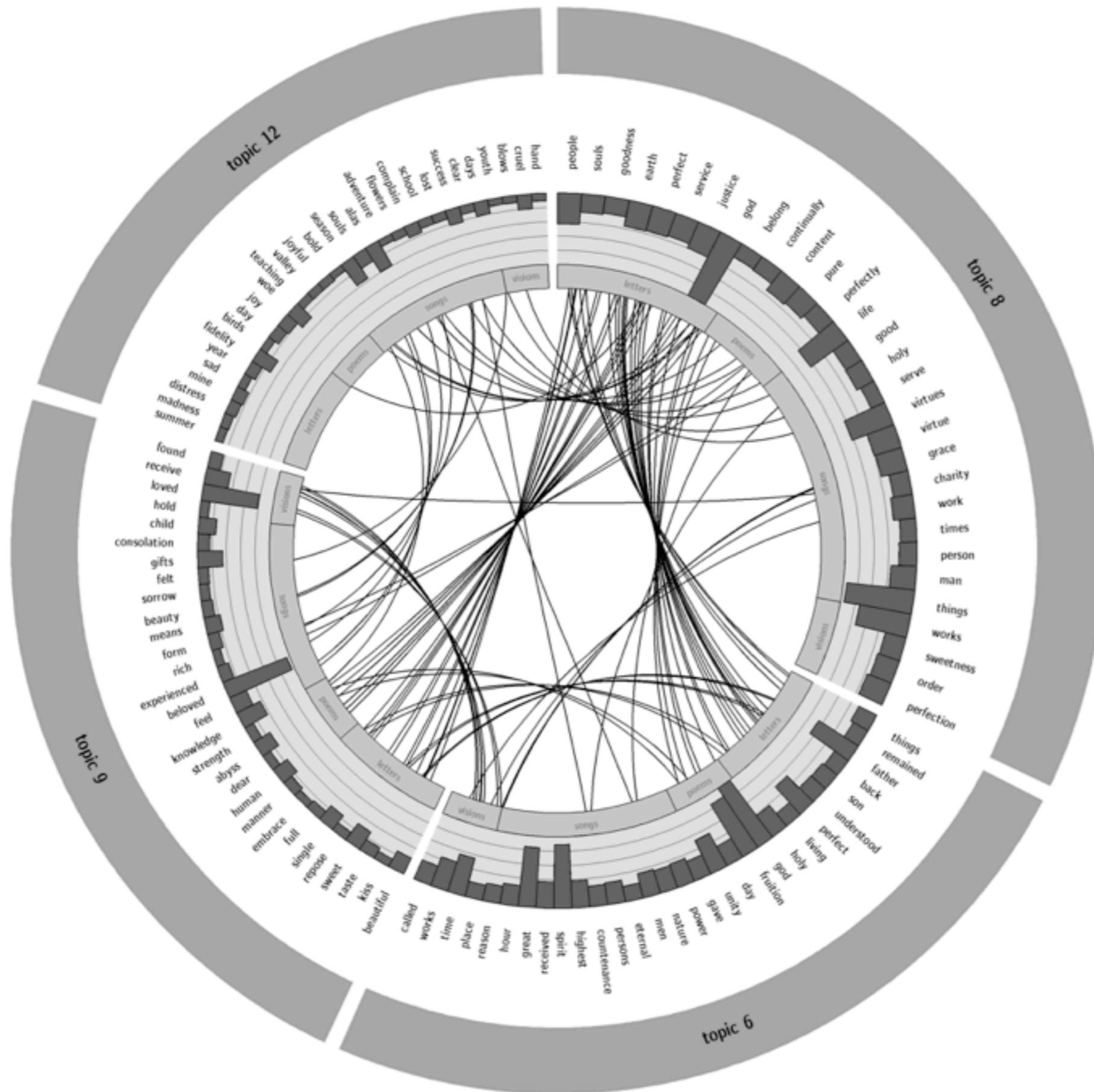


year
 woe
 winter
 time
 summer
 success
 seek
 season



Topics in Hadewijch's Oeuvre
 English translation by C. Hart (1980)





Circos
 Circulaire visualisatie
 van genomische data

“Where’s the beef?”

- + Visuele DH (“mooie plaatjes”): bevrijdend?
- + Innovatie als drijfveer: vruchtbare confrontatie met exacte wetenschappen?
- - Wat levert het op? Nieuwigheidswaarde?
- - Verwijt: Oppervlakkigheid als Achilleshiel

Stylometrie

- Kwantitatieve studie van schrijfstijl
- Stijl \Leftrightarrow meta-data
 - Auteurschap (*authorship attribution*)
 - Datering (*stylochronometry*)
 - Tekstsoort (*genre studies*)
 - ...

The Evolution of Stylometry in Humanities Scholarship

DAVID I. HOLMES
The College of New Jersey, USA

Abstract

This paper traces the historical development of the use of statistical methods in the analysis of literary style. Commencing with stylometry's early origins, the paper looks at both successful and unsuccessful applications, and at the internal struggles as statisticians search for a proven methodology. The growing power of the computer and the ready availability of machine-readable texts are transforming modern stylometry, which has now attracted the attention of the media. Stylometry's interaction with more traditional literary scholarship is also discussed.

1. Introduction

The successful uncovering of Joe Klein as author of *Primary Colors* by Don Foster, the coverage by CBS TV of the UCLA conference in 1996 discussing the attribution of *A Funeral Elegy*, Robert Matthews' thoughtful articles in *Newsweek* and *The Sunday Telegraph* (see, for example, Matthews, 1994), and the request by the BBC for the Bristol Stylometry Research Unit to investigate the authorship of the 'Cassandra' letters in *Tribune* followed by footage of Richard Foesyth and myself on their *Newsnight* programme are examples of how stylometry—the statistical analysis of literary style—has begun to intrigue the media. Indeed, the age of 'pop' stylometry could now be with us.

Not all such publicity has been necessarily good, however. Channel 4's *Street-legal* programme unmasked flaws in the reliability of Andrew Morton's so-called 'cusum' technique, which had been used in UK courts in a number of high-profile criminal cases, and the Shakespeare scholar Stanley Wells recently was moved to write in the *Times Literary Supplement*:

'Most people who make a profession of the study of literature do so because they have an artistic rather than a scientific bent. They are, to put it simply, better at English than at Maths. But the investigation of authorship of disputed and apocryphal works relies increasingly on statistical tests that take their exponents into the realms of higher mathematics. Their works are packed with tables, charts and statistical analyses which many of those interested in the works investigated have neither the inclination nor even the intellectual ability to understand, let alone to assess.' (Wells, 1996).

Correspondence: David I. Holmes, Department of Mathematics and Statistics, The College of New Jersey, PO Box 7718, Ewing, NJ 08628-0718, USA.

Literary and Linguistic Computing, Vol. 13, No. 3, 1998

What, then, does stylometry seek to do? It certainly does not seek to overturn traditional scholarship by literary experts and historians, rather it seeks to complement their work by providing an alternative means of investigating works of doubtful provenance. At its heart lies an assumption that authors have an unconscious aspect to their style, an aspect which cannot consciously be manipulated but which possesses features which are quantifiable and which may be distinctive.

Bailey (1979) lists the general properties which quantifiable features of a text should possess: 'They should be salient, structural, frequent and easily quantifiable, and relatively immune from conscious control.' By measuring and counting these features, stylometrists hope to uncover the 'characteristics' of an author.

The two primary applications of stylometry are attributional studies and chronological problems, and Laan (1995) has provided a sound exposition of the rationale behind such applications. Laan points out that these applications seem to be based on seemingly contradictory premises, with attributional studies claiming that the unconscious aspect of an author's style remains fixed whilst chronological studies claim that stylistic features develop rectilinearly during the course of an author's life. These claims may not be incompatible; the choice of features is the overriding concern. With attributional problems, it is certainly wise to work within the same genre when faced with a list of candidate authors for a disputed work and to work within as close a time period as possible when selecting appropriate 'control' authors. Genre effects generally will supersede authorial features in the discrimination process.

The major problem inhibiting stylometry's acceptance within humanities scholarship is that, as yet, there is no consensus as to correct methodology or technique. Rudman (1997) provides a perceptive analysis of this problem and illustrates well how, for every method that 'works', there soon appear counter-arguments pointing out crucial flaws. A methodology successful for one attributional problem does not necessarily 'work' for another. Practitioners search for the 'holy grail' of stylometry, a technique beyond reproach which may be applied successfully to all genres, languages, and eras.

The historical development of stylometry is reflected in the choice of quantifiable features used as authorial discriminators. Lexical features have predominated, yet this decade has seen the application of syntactic and semantic features to attributional problems, allied to the enormous growth in computing power and the

© Oxford University Press 1998

Auteursherkenning

- Populairste toepassing
- Stylome Hypothesis
 - Unieke vingerafdruk
 - Kwantitatief meten

Journal of Quantitative Linguistics
2005, Vol. 12, No. 1, pp. 65–77
DOI: 10.1080/09296170500055350



New Machine Learning Methods Demonstrate the Existence of a Human Stylome*

Hans van Halteren¹, R. Harald Baayen², Fiona Tweedie³,
Marco Haverkort⁴ & Anneke Neijt⁵

¹Department of Language and Speech, Radboud University Nijmegen ²Max-Planck-Institut für Psycholinguistik, Nijmegen ³School of Mathematics, University of Edinburgh
⁴Department of Linguistics, Radboud University Nijmegen (also Department of Linguistics, Boston University) ⁵Department of Dutch, Radboud University Nijmegen

ABSTRACT

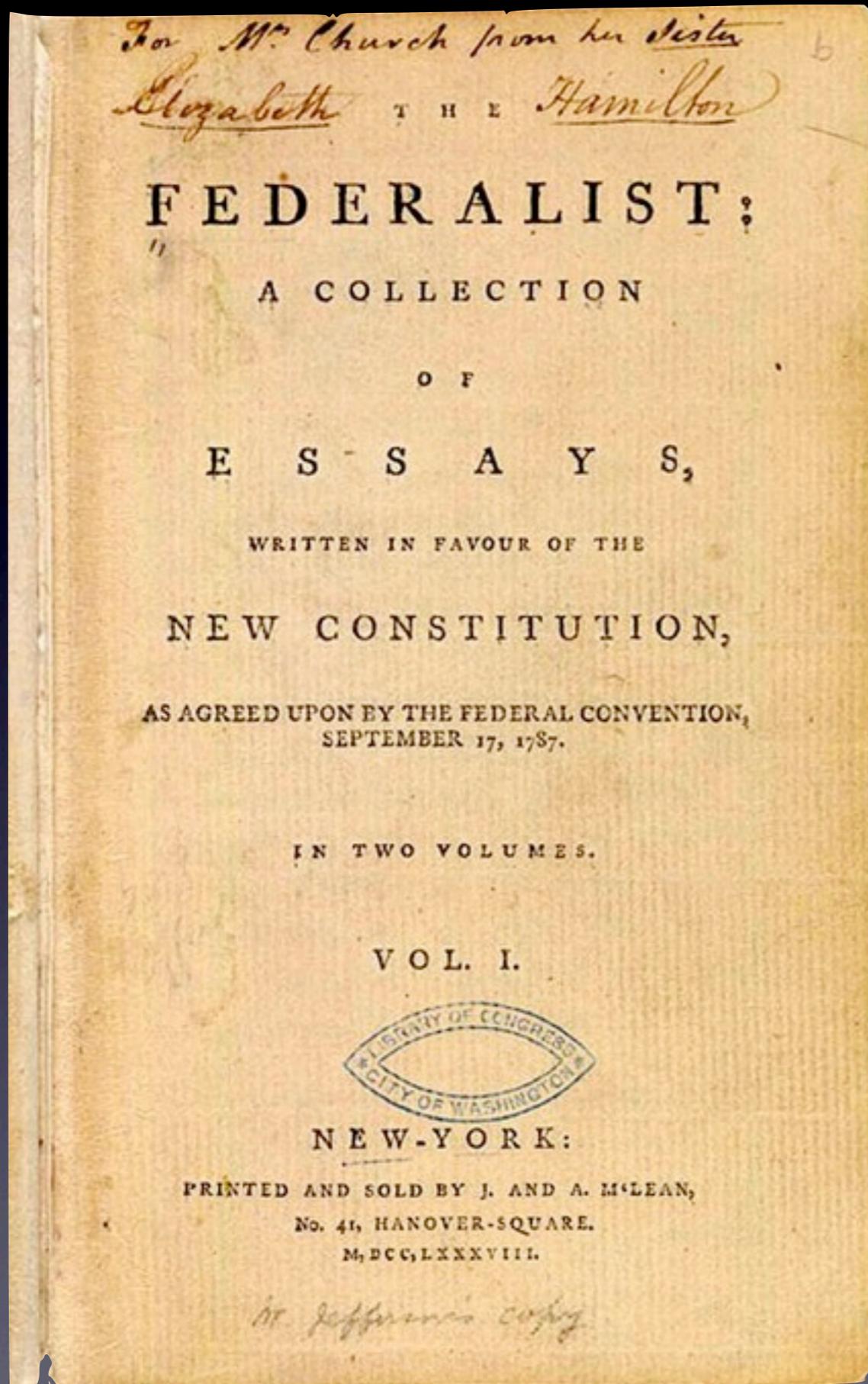
Earlier research has shown that established authors can be distinguished by measuring specific properties of their writings, their stylome as it were. Here, we examine writings of less experienced authors. We succeed in distinguishing between these authors with a very high probability, which implies that a stylome exists even in the general population. However, the number of traits needed for so successful a distinction is an order of magnitude larger than assumed so far. Furthermore, traits referring to syntactic patterns prove less distinctive than traits referring to vocabulary, but much more distinctive than expected on the basis of current generativist theories of language learning.

INTRODUCTION

We all make extensive use of a natural language to communicate with others. The ease with which we do this might give the impression that all the users of a particular natural language are using the exact same language. Now, specific languages do not come completely hard-wired in the brain, although the brain is thought to contain a set of hard-wired expectations about the structure of natural language, the so-called Universal Grammar. Whether or not there is indeed some support from structures in the brain, we always have to learn a language on the basis of

*Address correspondence to: Hans van Halteren, Department of Language and Speech, Radboud University Nijmegen, P.O. Box 9103, NL-6500 HD Nijmegen. Tel. + 31 24361 2836. E-mail: hvh@let.kun.nl.

0929-6174/05/12010065\$16.00 © Taylor & Francis Group Ltd.



- Jong paradigma (1960s)
- **Mosteller & Wallace** (US)
- *Federalist papers* (1780s)
- Twee innovaties:
 - Kwantitatieve aanpak
 - “Functoren”

Traditioneel

- **Natte vinger...**
- Opvallende kenmerken
 - bv. zeldzaam werkwoord
- “Checklist”
- Maar:
 - scholen, ateliers, ...
 - tradities
 - vervalsing, imitaties ...
 - ...

Mosteller & Wallace

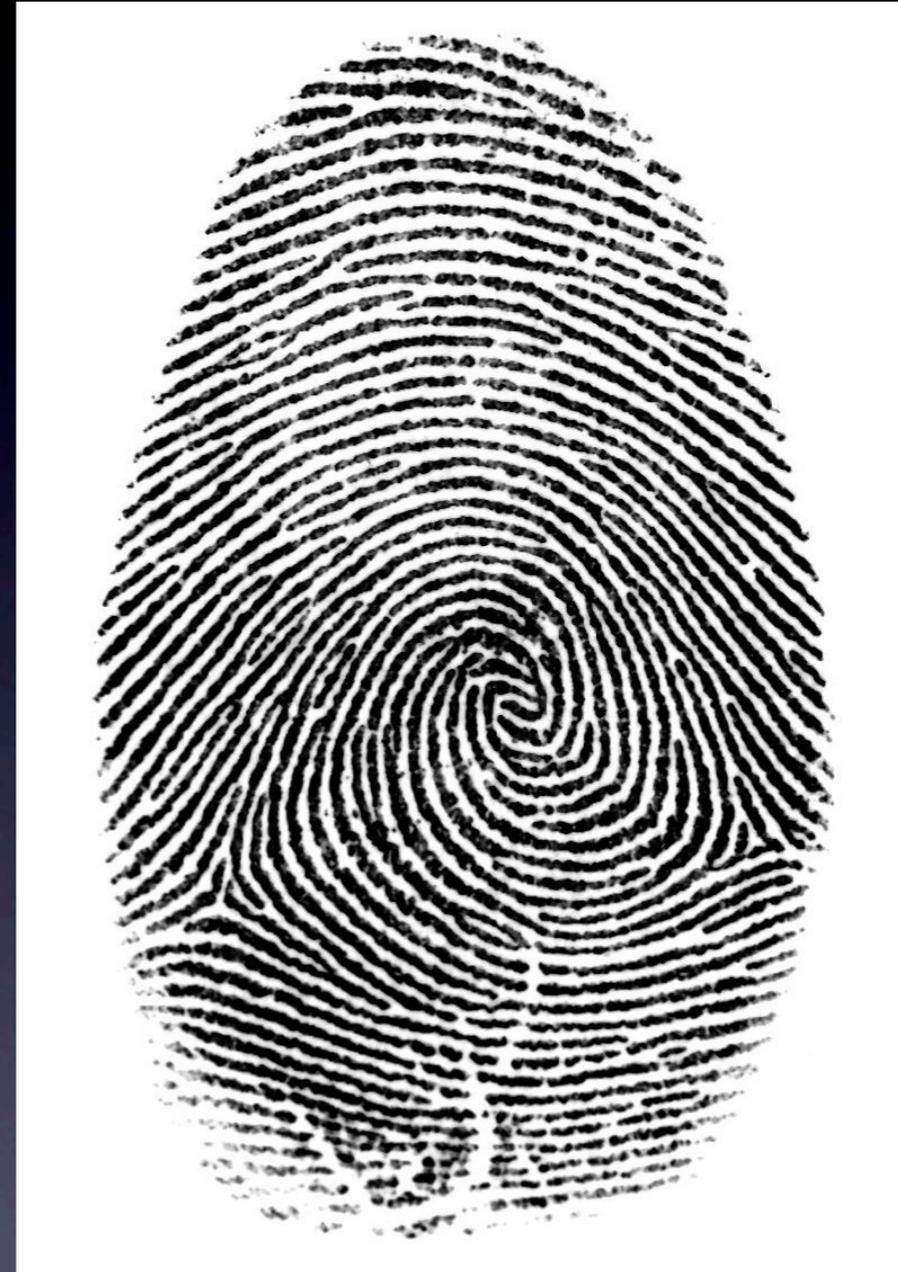
- Onopvallend kenmerken
- **Functiewoorden of functoren**
 - lidwoorden
 - voorzetsels
 - voornaamwoorden
 - [naamvallen?]

Voordelen?

Veel observaties

Alle auteurs, zelfde set

Relatief inhoudsonafhankelijk



Aantal letters *f* op volgende
slide?

Finished files are the result
of years of scientific study
combined with the experience
of many years.

Hoeveel?

Verwerken wij functoren 'onbewust'?

Finished files are the result
of years of scientific study
combined with the experience
of many years.

Welke tekst staat op de
volgende slide?



**PARIS
IN THE
THE SPRING**

En?

Moeilijk fouten detecteren...

Memory & Cognition
1977, Vol. 5 (6), 636-647

Detection errors on *the* and *and*: Evidence for reading units larger than the word

ADAM DREWNOWSKI

Rockefeller University, New York, New York 10021

and

ALICE F. HEALY

Yale University, New Haven, Connecticut 06520
and Haskins Laboratories, New Haven, Connecticut 06511

In five experiments, subjects read 100-word passages and circled instances of a given target letter, letter group, or word. In each case subjects made a disproportionate number of detection errors on the common function words *the* and *and*. The predominance of errors on these two words was reduced for passages in which the words were placed in an inappropriate syntactic context and for passages in which word-group identification was disturbed by the use of mixed typecases or a list, rather than a paragraph, format. These effects for the word *and* were not found for the control word *ant*. These results were taken as evidence that familiarity

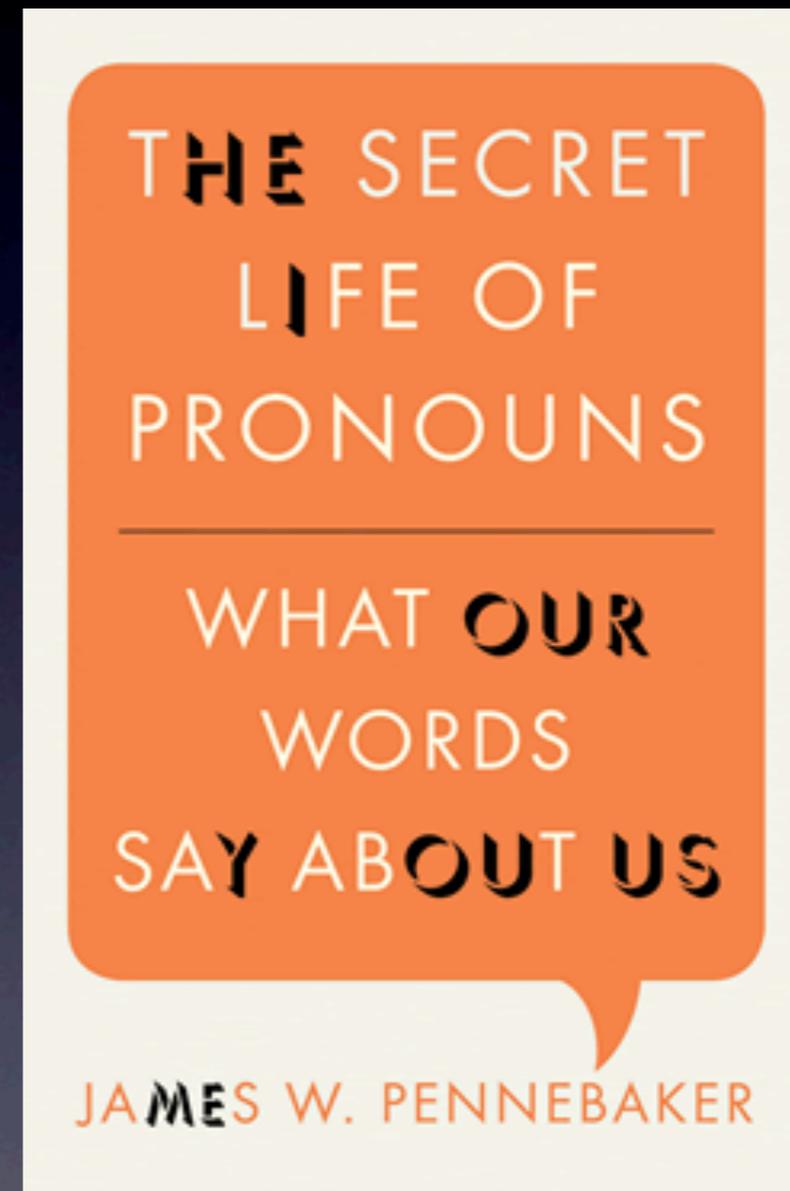
Onbelangrijk?

According to a research at Cambridge University, it doesn't matter in what order the letters in a word are, the only important thing is that the first and last letter be at the right place. The rest can be a total mess and you can still read it without problem. This is because the human mind does not read every letter by itself, but the word as a whole.

According to a research at Cambridge University, it doesn't matter in what order the letters in a word are, the only important thing is that the first and last letter be at the right place. The rest can be a total mess and you can still read it without problem. This is because the human mind does not read every letter by itself, but the word as a whole.

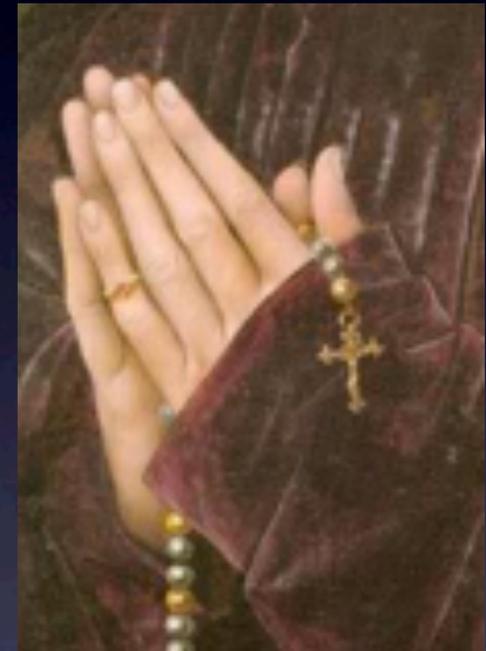
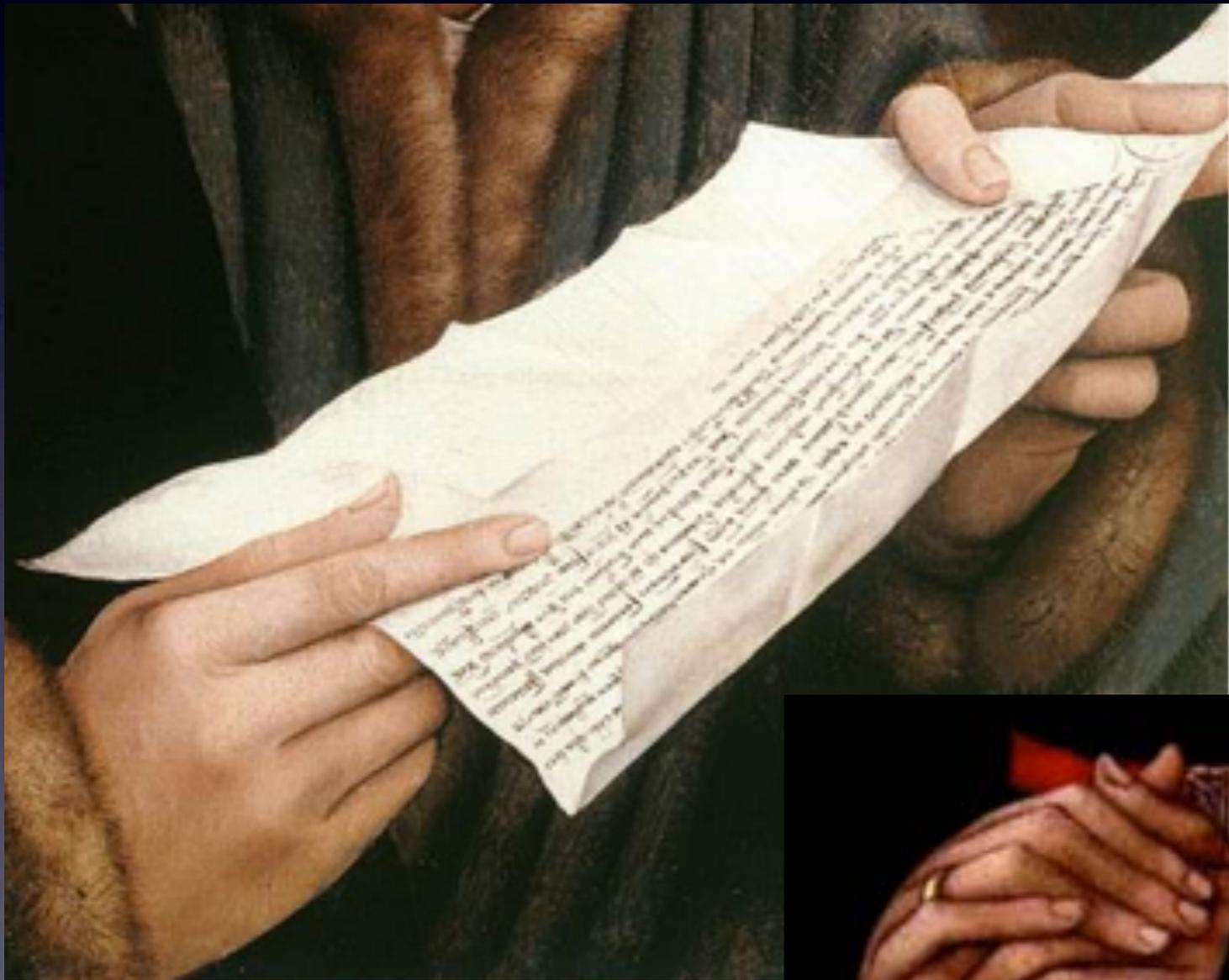
Functoren

- Populair in stylometrie
- Frequenties als input
- Pennebaker (2011)



Parallel in kunstgeschichte

Morelli (1816–1891)



Middelnederlands

- Germaanse volkstaal
- Lage landen
- **Middeleeuwen**
- Ca. 1200-1500
- Literatuur



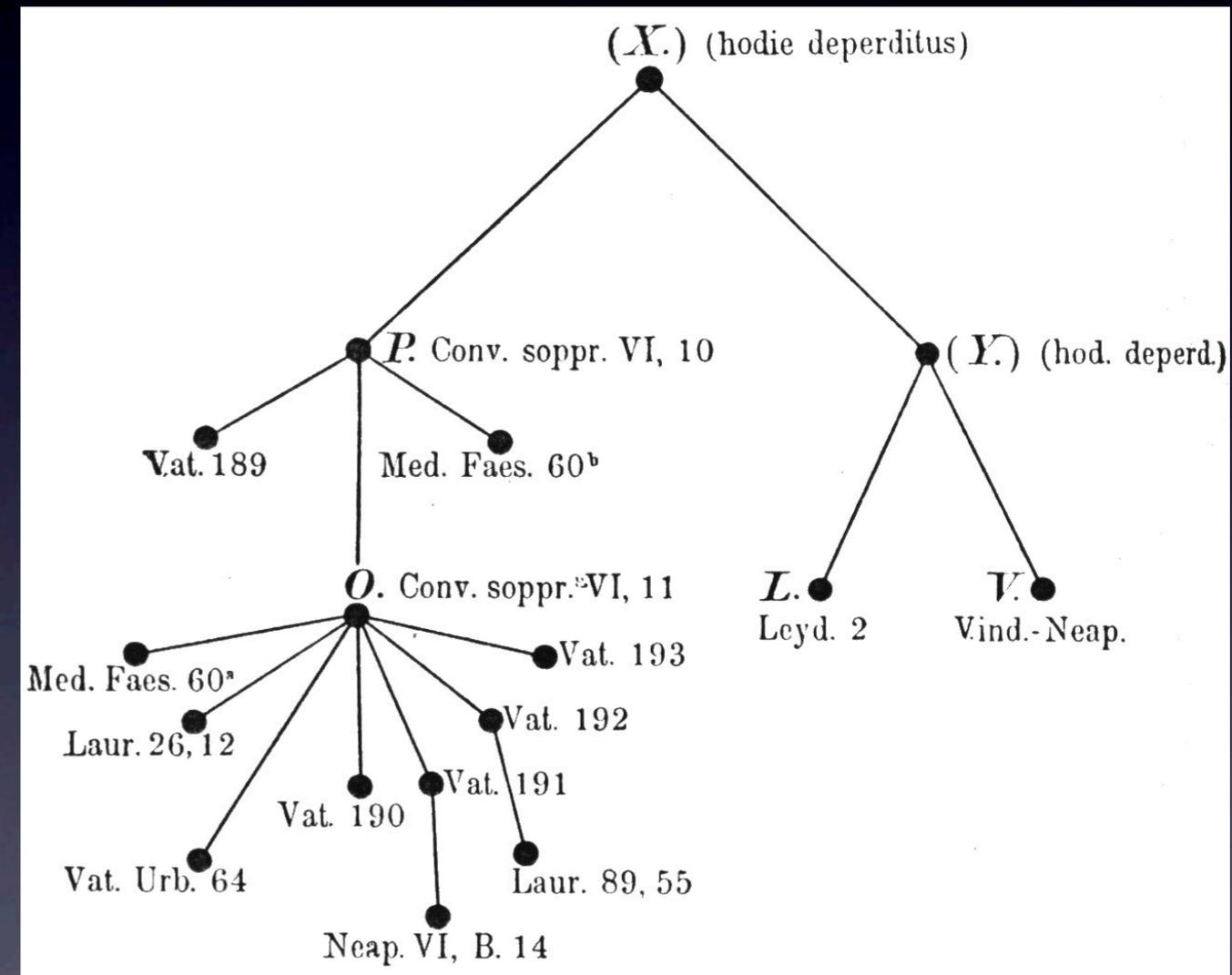
Handgeschreven wereld

- Geen drukpers
- **Manu-scripten**
- Kopiisten scribenten
- Iedere kopie uniek



Variatie als regel

- Geen standaardtaal, -spelling
- **Kopieën wijken af van legger:**
 - Spelling
 - Locale dialecten
 - Stilistische voorkeur?
 - ...
- Tekst oorspronkelijke auteur?



Middelnederlands

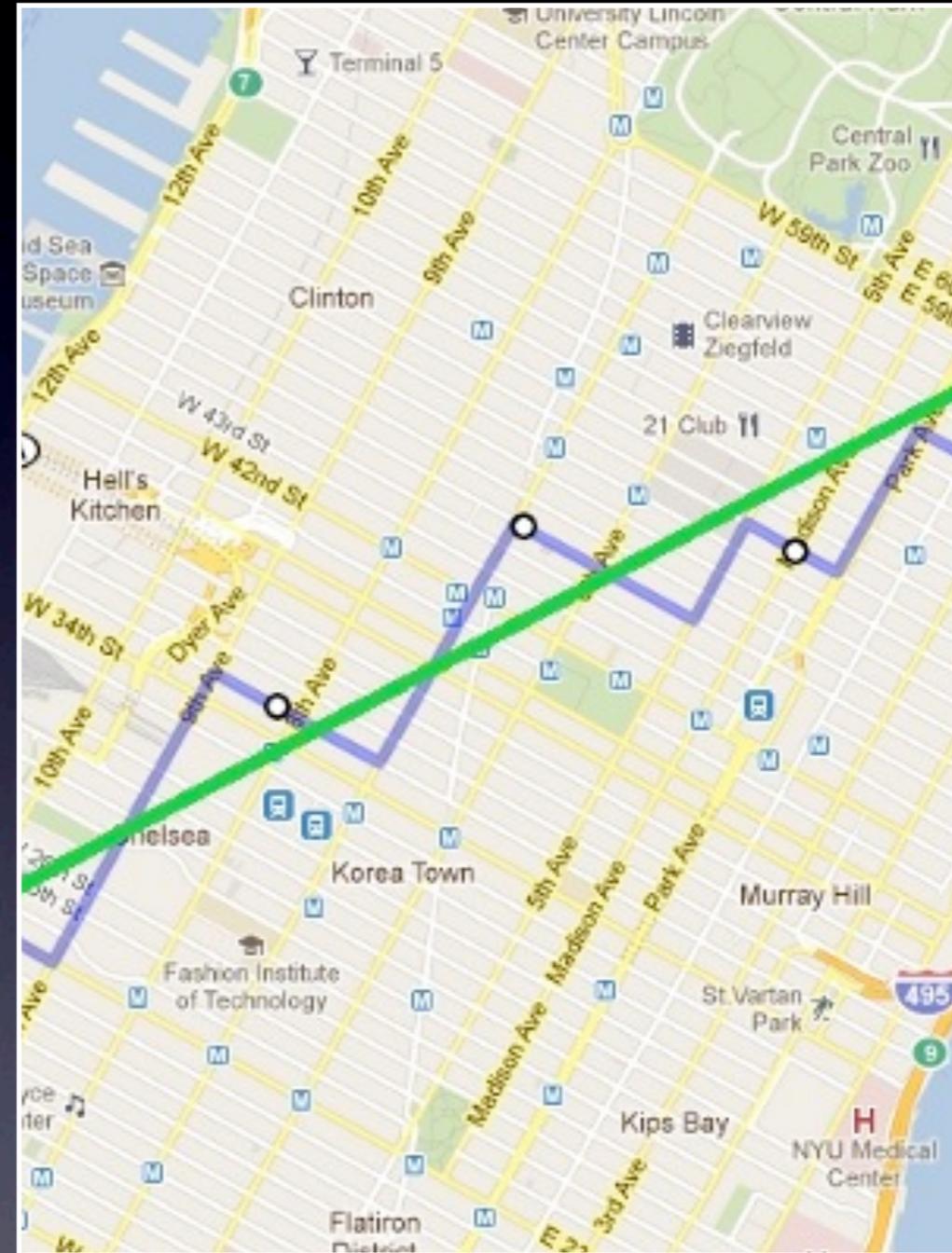
- Van Dalen-Oskam & Van Zundert 2007
- *Literary and Linguistic Computing*
- *Roman van Walewein*
- Auteursovergang
- Leiden, Ltk. 195
- Pionierswerk

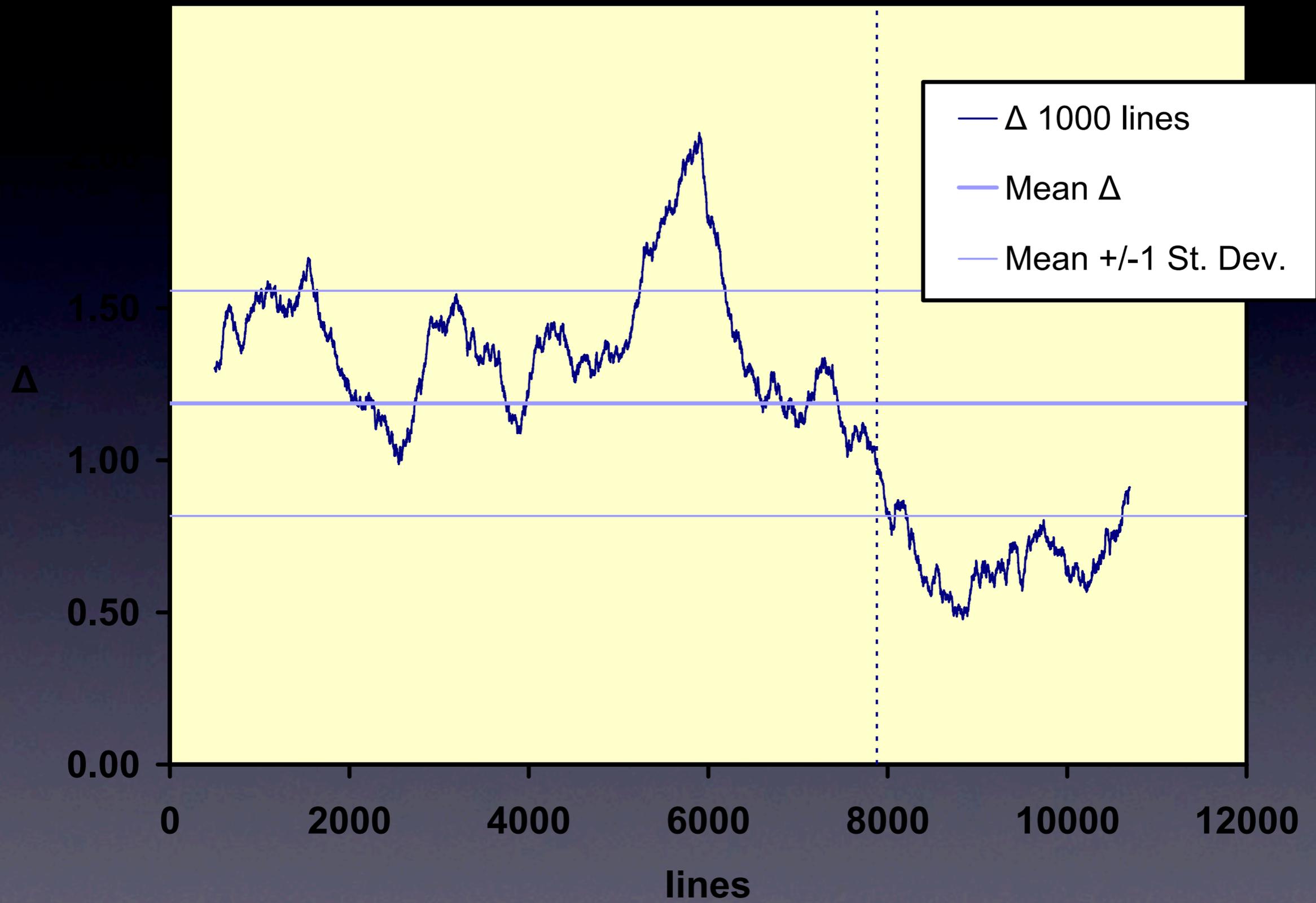


Burrows's Delta

- Stilstatische afstandsrekening
- *Manhattan distance* op MWF
- *Nearest neighbor learning*
- Zowel classificatie als afstand

$$\Delta_B^{(n)}(D, D') = \sum_{i=1}^n \frac{1}{\sigma_i} |f_i(D) - f_i(D')|$$





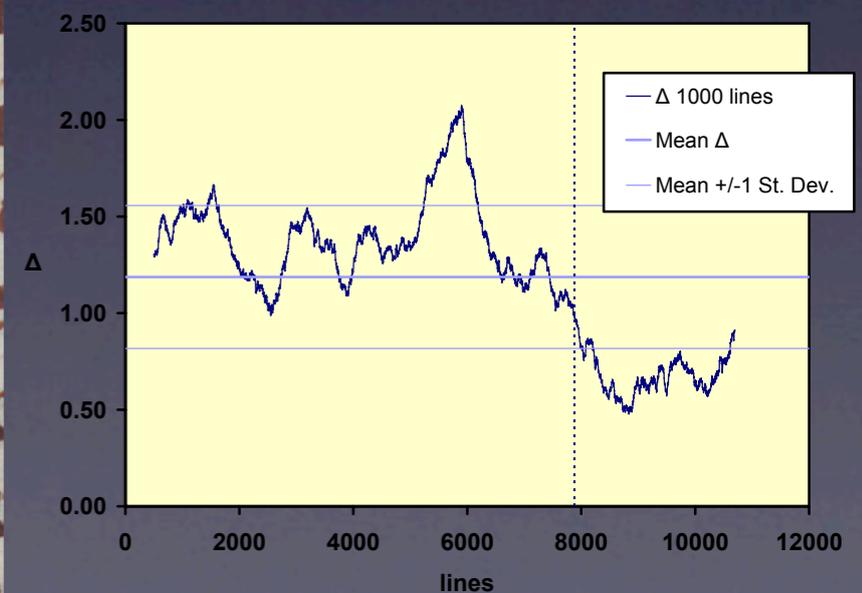
Inden casteel dien ghi ziet staen
Ouer dese riuere in zijn behout
En ware niere also stoue
Dese halen wyde dat wende wel
En ware walewen die vaders suet
Hets drauentuen vader
Hco ware myn gelue al gader
H ordeyt van hare yec gewagen
H me liets hene meo vlagert
Scherlike hinc soude comen
Thunt hem te scaden ok te vromen
En proeuen oec syn gheual
Rv hebbick v berecht al
Serecht in so doet wel
Van wak den vader suet
Dan salic v berechten vore
Ic hebbe dicke wile ghehort
Van wak sueten meo de drom



E
R
h
tu
D
H
S
E
E
G
D
S
D
B
h
D

Leiden

UB, Ltk. 195



Serendipiteit?

Zochten auteurs en
vonden kopiisten...

Delta for Middle Dutch—Author and Copyist Distinction in *Walewein*

Karina van Dalen-Oskam and Joris van Zundert
Huygens Instituut, The Hague, The Netherlands

Abstract

The Middle Dutch Arthurian romance *Roman van Walewein* ('Romance of Gawain') is attributed in the text itself to two authors, Penninc and Vostaert. Very little quantitative research into this dual authorship has been done. This article describes our progress in applying different non-traditional authorship attribution methods to the text of *Walewein*. After providing an introduction to the romance and an overview of earlier research, we evaluate previous statements on authorship and stylistics by applying both Yule's measure of lexical richness and Burrows's Delta. To find out whether these new methods would confirm or even enhance our present knowledge about the differences between the two authors, we applied an adapted version of John Burrows's Delta procedure. The adapted version seems to be able to distinguish the double authorship of the romance. It also helps us to confirm some and to reject other earlier statements about the position in the text where the second author started his work.

Correspondence:

Karina van Dalen-Oskam,
Huygens Instituut, Postbus
90754, NL-2509 LT Den
Haag, The Netherlands.
E-mail:
karina.van.dalen@
huygensinstituut.knaw.nl

1 Introduction¹

Until now, non-traditional authorship attribution techniques have hardly been applied to medieval texts. The fact that we do not and probably never will know the author of many texts, may play a role in this. There are cases, however, in which the researcher does not necessarily want to identify a specific author but wishes to distinguish one (perhaps anonymous) author from another. In these cases, it could indeed be very interesting to turn to modern authorship attribution techniques and find out of how much help they actually can be.

In our research on medieval Dutch literature, we have long been intrigued by a Middle Dutch Arthurian romance that seems to be the perfect case for an experiment like this. In this article, we first present the case and the questions we started out with (Section 2). We then summarize our first results and present the new questions to which these lead (Section 3). In Section 4, we investigate

two specific questions. We then focus on the possibilities of Burrows's Delta for distinguishing the authors as well as the scribes of *Walewein* (Section 5). After presenting our results, we sum up our conclusions and describe how we plan to test our results and to refine the method we are developing.

2 The Case: *Roman van Walewein*

One of the few Middle Dutch romances that were not translated from Old French is the Arthurian romance known as *Walewein* or *Roman van Walewein* ('Romance of Walewein'). The main character of this verse text in rhyming couplets is Arthur's nephew Walewein—the 'Gawain' of English texts and the 'Gauvain' of the French. *Walewein* was probably written around the year 1260, but the only surviving complete manuscript dates from almost a century later. The manuscript was written by two scribes. The first wrote lines

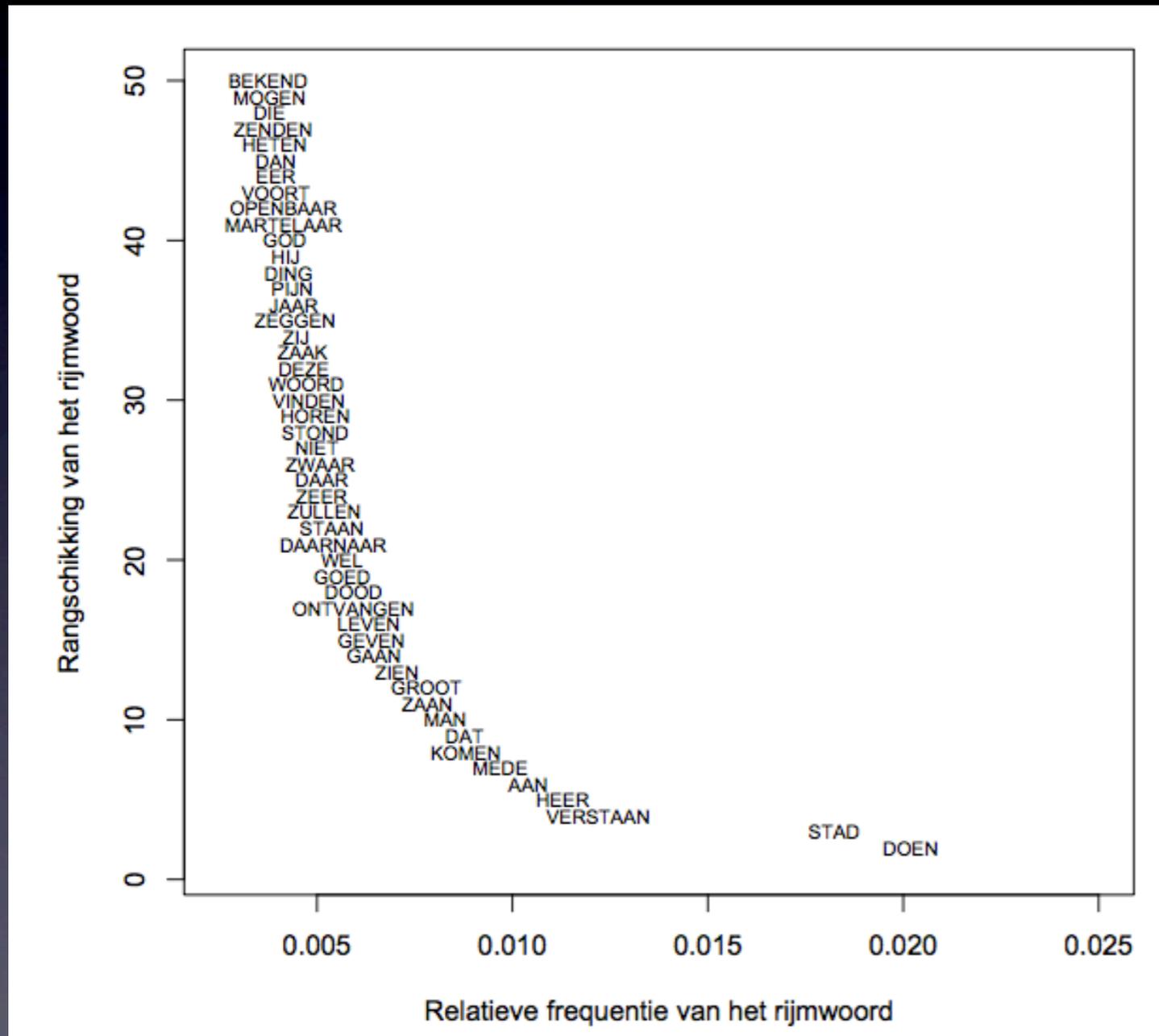
D	Ter stont ende ter seluer vren
E	Tier stont ende ter seluer vren
F	Tlere stont enter seluer vren
G	Tottien stonden en ter uren
H	TEn stonden ende ter seluer vren
I	Tjerst stont ende tier veren
J	Tyer stont ende tier seluer vren
N	TJer stont tier seluer vre

Auteursherkenning vs. **Scribent**herkenning...

Rijmwoord

- Meeste Middelnederlandse literatuur **berijmd**
- Bijbel, encyclopedie, liefdesgedichten, ...
 - Willam die madock **maecte** / Dair hi dicke om **waecte**
- **Skelet** van tekst
- Heel moeilijk aan te passen
- “Oorspronkelijke dichter in het aangezicht staren”?

Rijmwoorden als functoren?



- Eindig # combinaties
- Formules, “stoplap”
- Popliedjes: “De lucht is blauw : Ik hou van ...”
- Zipfianse verdeling
- Surrogaat functoren?

Lemmatiseren

- Spellingvariatie
- Lemmatiser
- Lemma-frequenties

in der ierster partien inden	EIND
mogedi van claudiusse vinden	VINDEN
hoe hi was keyser ende here groot	GROOT
dese heeft gedaen so vor sine doot	DOOD
dat nero keyser na hem blive	BLIJVEN
bi den rade van sinen wive	WIJF
ende der heren daer hi in desen	DEZE
te seer af scheen bedwongen wesen	WEZEN
octavien siere dochter man	MAN
was nero ende also dan	DAN

Spiegel historiael *Speculum historiale* (13e E)





Jacob van Maerlant | Filip Utenbroeke

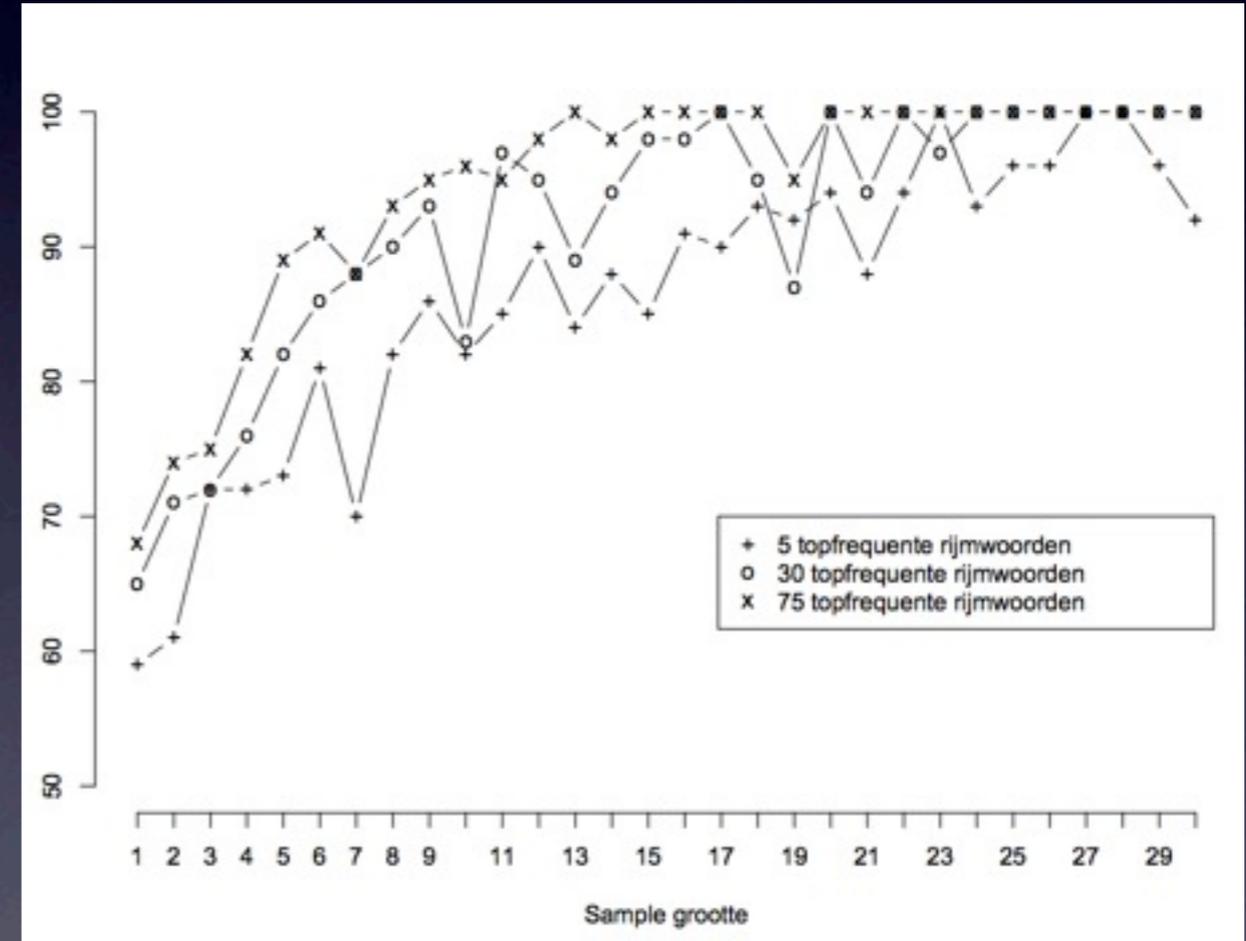
Tekstclassificatie

Spam filtering

 Mail thinks this message is Junk Mail.

From: first choice <info@yahoo.com>
Date: 21 Apr 2010 19:25:23 GMT+02:00
Subject:
To: undisclosed-recipients::
Reply-To: firstchoice_loanenquiry01@w.cn

Reliable loan offer!!!
Here is an opportunity for those in financial problem and financial up lift in their life, we give out loan at a very rate 3.05%, we give out all kind of loan to help the national stress. Many are suffering and needs help to improve are jobless and need financial help to start a business help to clear their bills and debt. Here is a wise decision



Dit is sonne sprake... met coning... 1. Boec... 1. 2. cap.

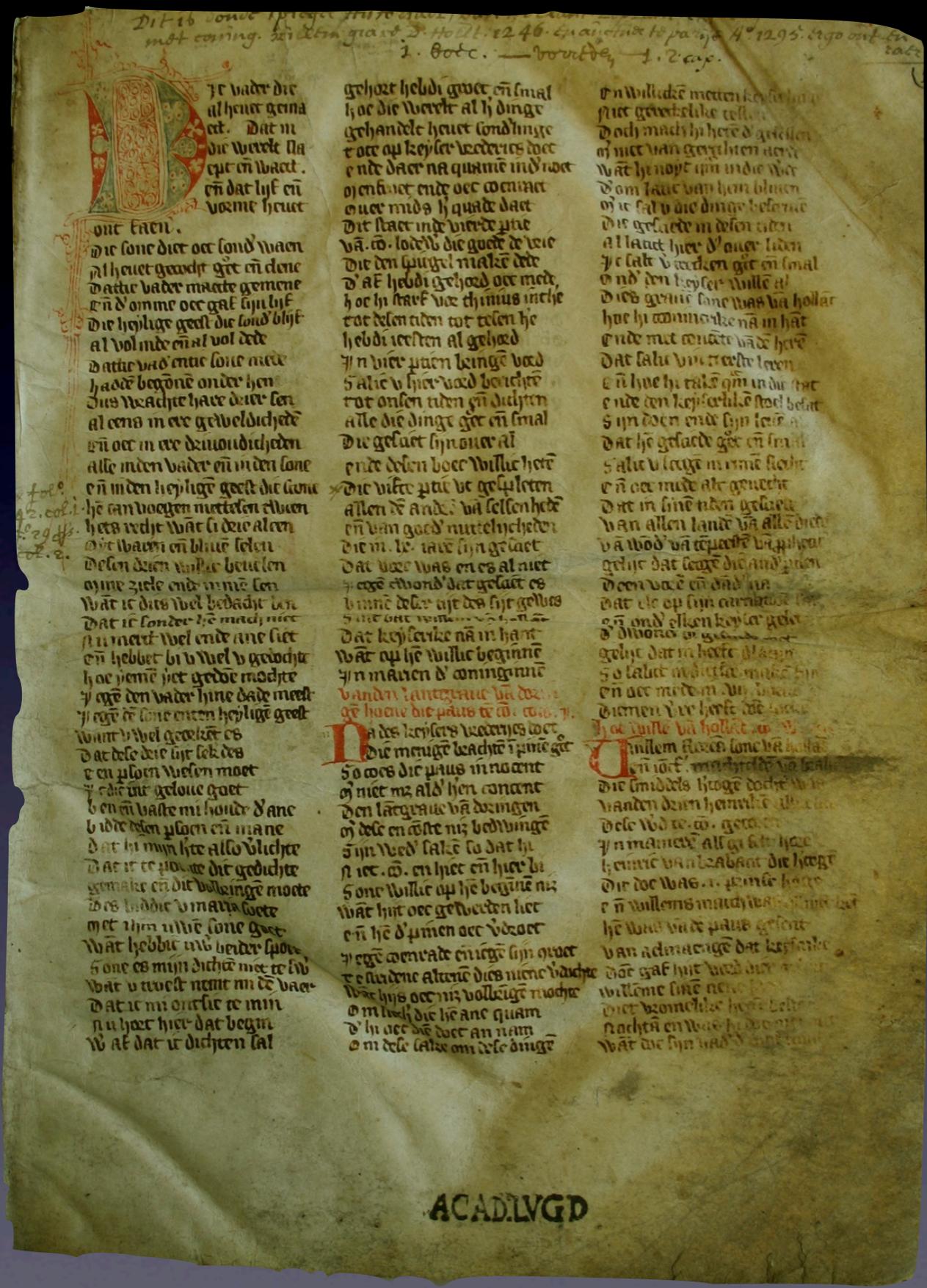
Die vader die al heuet genaet. Dat in die werelt naer en waert. en dat luf en vorne heuet ont faen. Die lone die oec sond waen al heuet geacht get en dene dattie vader maecte gemene en d'omme oec gaf luf luf die heilige geest die lone blif al vol mede en al vol dede dattie vad' enne lone mede hadde begone onder hen dus waachte hare d'werlen al eens in ere gebeldichede en oec in ere d'ruoudicheden alle inden vader en in den lone en in den heilige geest die lone he can voegen metelen elken hees vedt wat si die alen die waen en blau' selen desen dien wil' heu len ayne ziele end' in me' len wat ic dies wel bedacht ben dat ic sonder he mach niet si in mart' wel ende ane liet en hebbet bi v'wel v'gedochte hoe yem' het gedde mochte y'egge den vader h'ne dade meel y'egge de lone enen heilige geest want v'wel gevellet es dat dese die luf des e en y'loen welen moet y'ed' die geloue goet ben en v'aste mi houde d'anc v'ide den y'loen en mane dat hi myn h'ce also v'lichte dat ic te y'owe die gedichte gemake en die volleinge moete des l'iddie v'maria' loete met ihm' u'lie lone goet wat hebbet u'v' heider spoe' s'one es myn dichte met te l'v' wat v'weest neent in de vaer dat ic mi ontfic te min nu haer hier dat begun w'at dat ic dichten sal

gehoez hebdi gweet en smal hoe die werelt al h'dinge gehandele heuet sondlinge t'ore op keijser wederics doet ende daer na quam' in d'woet o'enh'iet ende oec coen'iet o'uec mids h'quade daet dit staet in d'v'erde y'ue' va' co. loede die goede de veie die den l'ugel maect dede d'ae hebdi gehoez oec mede hoe hi staef v'ce th'imus inche t'oe delen aden tot telen he hebdi i'celten al gehoez y'n v'ier y'uen leinge v'oad s'alie v' h'ier v'oad bechere tot onsen tiden en d'ich'ten alle die dingge get en smal die gelact l'yn oner al p'nde delen boec w'illie h'ere die v'ite y'ue' de g'el'eten allen de a'nde' v'alellen h'ede en van goed' nuttelicheden die m. le. i'ace l'yn ge'act dat v'oe' was en es al niet y'egge d'wond' dat gelact es v'inne d'ele' tijt des l'ij' g'el'us s'it o'at v'one' y'n' n' n' dat keijserike na' in h'ere wat op h'e w'illie be'g'inne y'n' m'auen d' coning'one v'anden l'ant'grau' v'ad' d'oe' ge' h'one die p'aus te co. co. 7. Die menige' be'acht' i'p'nie get so coes die p'aus in no'ant of niet m'z ald' h'en con'ant den l'at'grau' v'ad' d'oringen of dese en'aste m'z bed'vinge s'ijn v'eed' sake so dat hi n'iet. co. en h'et en h'ere bi s'one w'illie op h'e be'g'ine n'z wat h'it oec g'el'oe'den h'et en h'e d'p'men oec v'oe'oe' y'egge coen'rade en'ige l'yn g'roet t'el'andene al'ene dies m'ene v'ad'che wat h'ys oec m'z volleige' mochte om' l'och' die he'anc quam d' h' oec die doet an nam om' dese sake om' dese dingge

en w'illie' meeten keijser' l'one niet ge'el'el'ike tellen doch mach hi h'ere d' g'el'eten of niet van g'el'eten ac'oe' wat hi noyt in' in die v'oe' som' laet van hen bl'imen of ic sal v' die dingge be'lo'ne die g'el'acte in desen aden al laet h'ere d'ou'oe' l'iden y'c late v'oe'den get en smal on' den keijser' w'illie al d'ies g'raue lone was v'ad' h'ollat hoe hi co'nn'ike na' in h'at ende met co'acte v'ade h'ere dat salu' v'm' t'ee'de l'oen en h'oe hi r'als q'm in die t'at ende den keij'or'like stoel be'at s'yn' eden ende l'yn' l'oe' n' dat h'e g'el'acte get en l'ra' s'alie v' l'ee'ge m' r'ime' l'oe'ht en oec mede al' g'ene'che dat in l'ine' n'den g'el'act' van allen lande' v'ad' alle' d'ie' v'ad' w'od' v'ad' t'ee'p'che v'ad' p'p'he'at g'el'ic dat l'ee'ge die and' y'uen d'oen v'ad' en d'ad' h'at dat ic op l'yn' car'na'ge' l'oe' d' d'wond' v'ad' g'el'act' g'el'ic dat in h'ere d' l'ee'ge s'o l'at'et in d'ic' l'ee' m' l'yn' en oec mede m' d' m' l'ee' d' m'nen v'ad' h'ere d' l'ee' h'oe' v'ille v'ad' h'ollat. co. 7. d' m' l'ee' m' d' h'ere v'ad' l'ee' die l'middels h'oe'ge do'che' v'ad' v'anden d'een' h'and'ic' al' l'ee' d'ele' w'od' te. co. g'oe' d' y'n' m'ame' d' all' g'at h'ere h'ere v'ad' l'ee' d' h'ere die doe' was. n' p'one' h'ere en w'illems mach' v'ad' l'ee' h'e w'ad' v'ad' p'aus' g'el'oe' van ad'm'ac'inge' dat keij'or'ike d'oe' gaf h'it v'ad' d'ic' v'ad' w'illeme l'ine' n'ee' d'ic' v'rom'el'ic' h'ere l'ee' h'ere' n'acht' en v'ad' h'ere d'ic' v'ad' wat die l'yn' v'ad' d' h'ere

ACAD.LVGD

Laatste deel, ca. 1316:
Lodewijk van Velthem
Vierde en Vijfde Partie



Top 25			
	Utenbroeke	Maerlant	Velthem
Utenbroeke	15	0	0
Maerlant	0	15	0
Velthem	1	0	14
Top 50			
	Utenbroeke	Maerlant	Velthem
Utenbroeke	15	0	0
Maerlant	0	15	0
Velthem	0	1	14
Top 75			
	Utenbroeke	Maerlant	Velthem
Utenbroeke	15	0	0
Maerlant	0	15	0
Velthem	0	1	14

4de boek steeds "fout" toegeschreven...

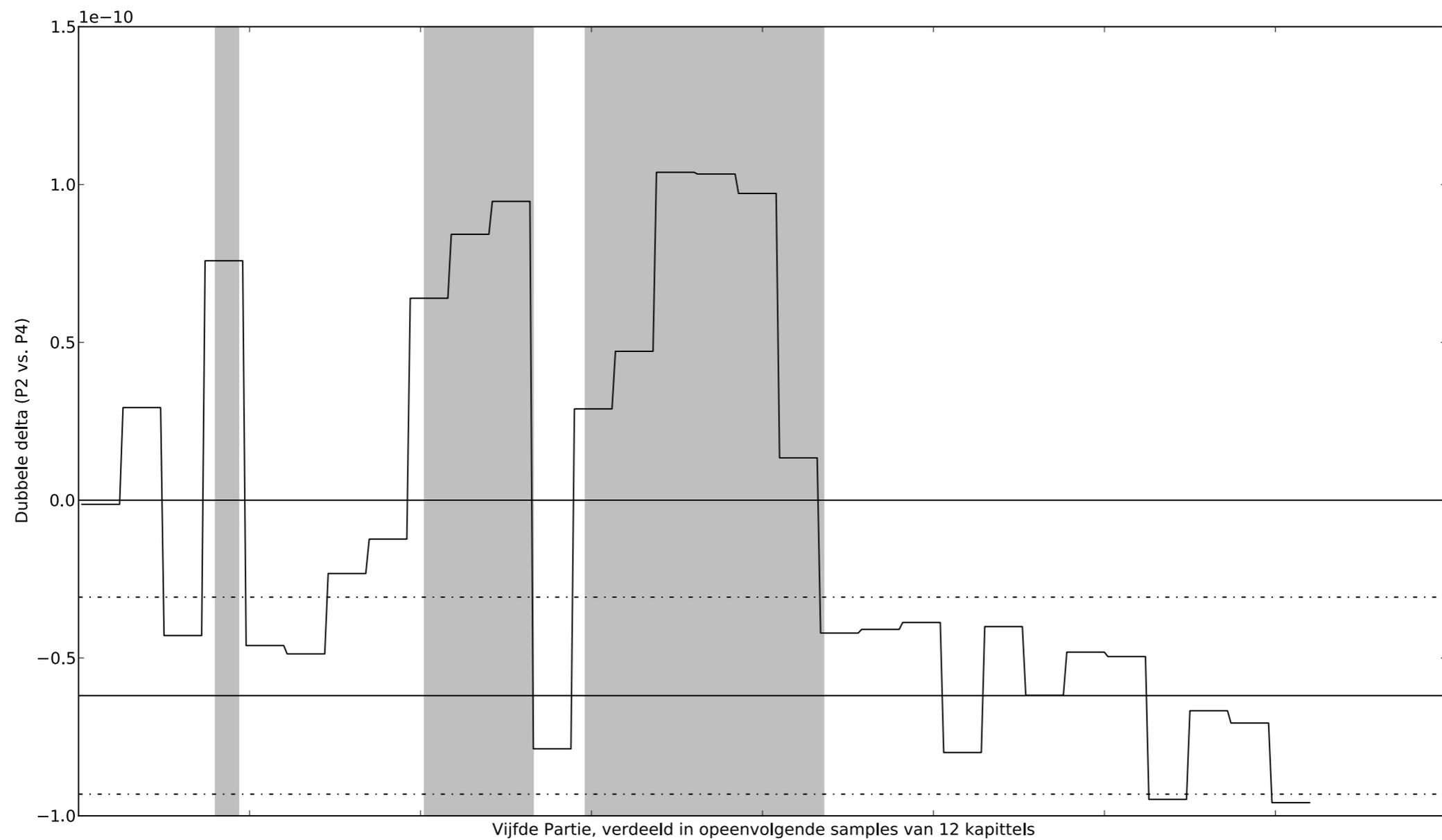
Guldensporenslag, 1302



Jan van Heelu Slag bij Woeringen 1288



Stijlcurve Velthems Vijfde Partie...



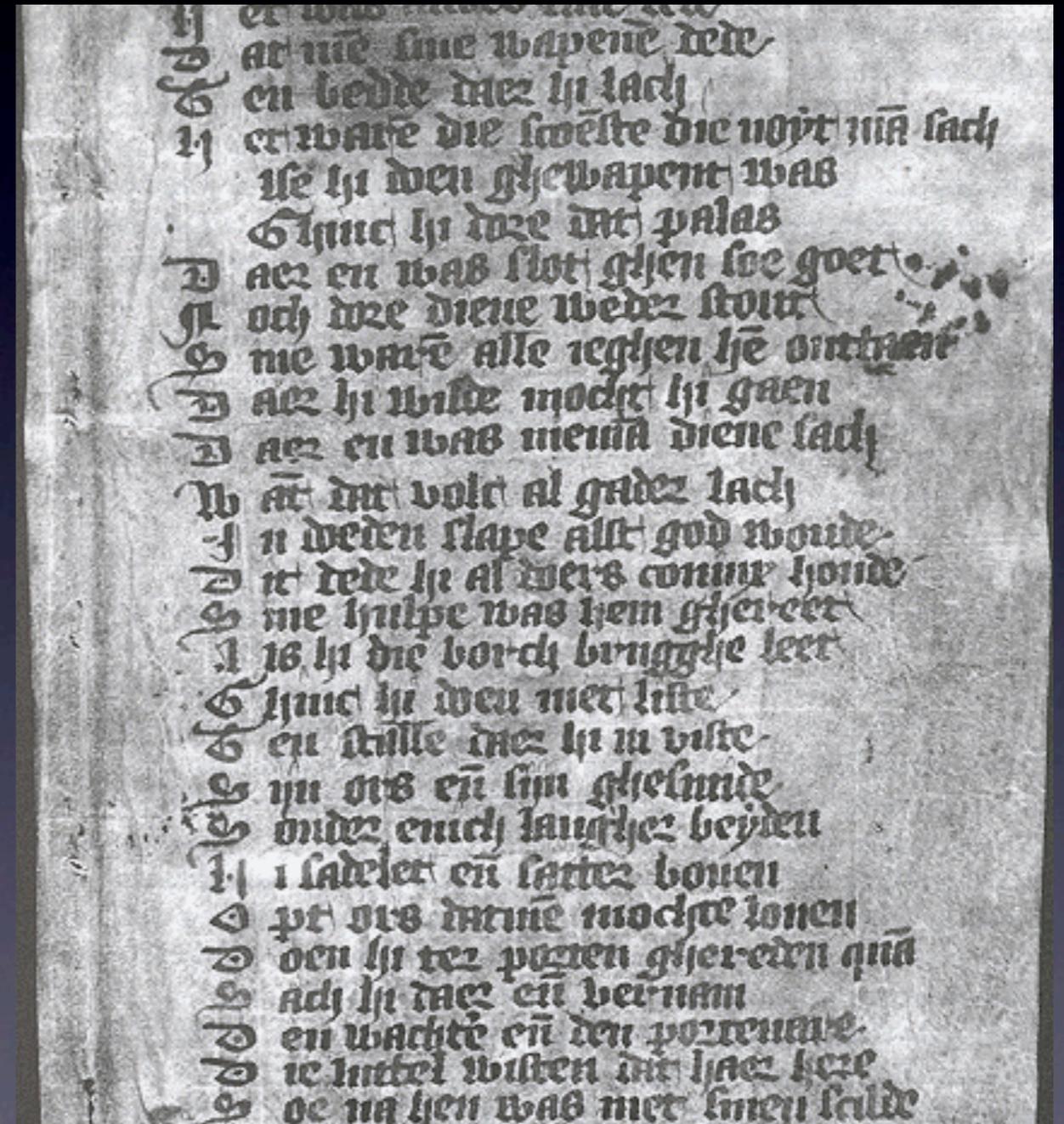


“Ontleende” Velthem een
reeds bestaand Vlaams
ooggetuigeverslag?
(serendipiteit...)

Wij willen Gedenck ons op goud
of dele blaminge laet ons laus
proclant vrey to dat Gyt harte
le die he me behoorden
Hertike volcom

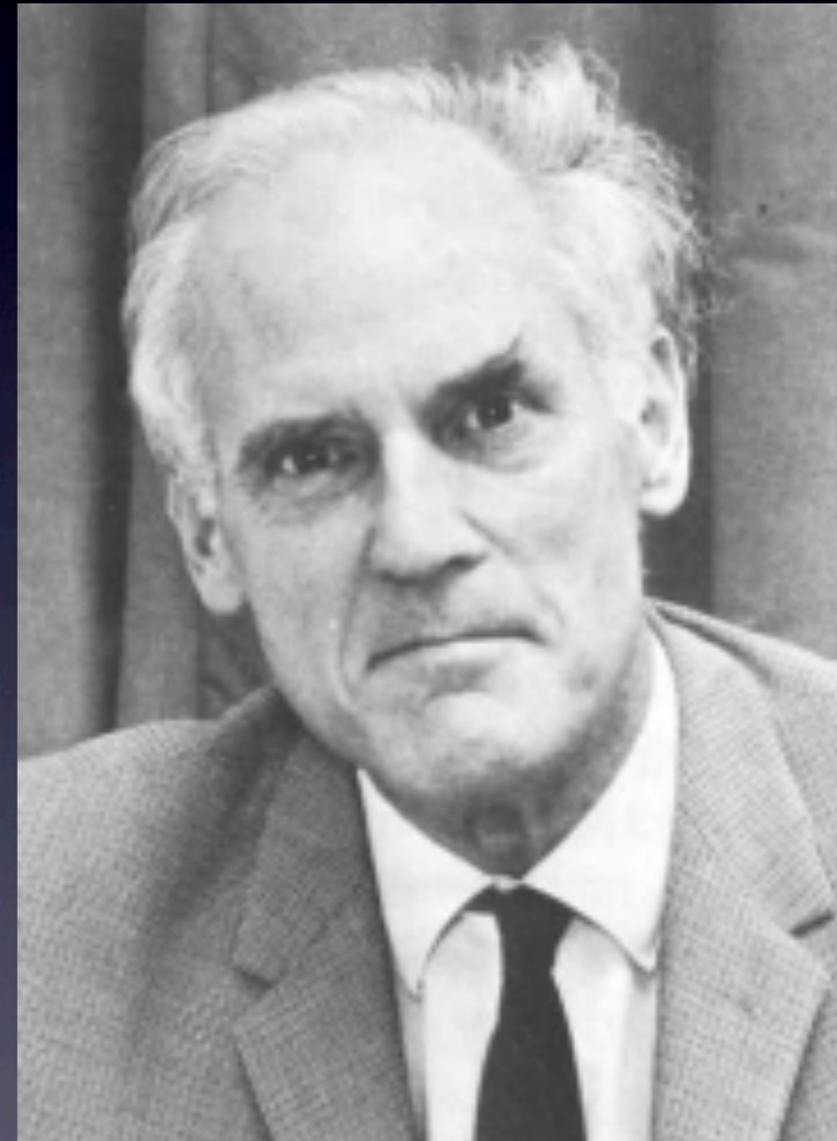
Karel ende Elegast

- Icoon Mnl'se literatuur
- 13e eeuw (?), Vlaanderen
- Anoniem, berijmd verhaal
- *Fraeye historie ende al waer*
- Nachtelijke rooftocht van Karel de Grote en Elegast



K.H. Heeroma (1909-1972)

- Nederlands filoloog, zelf dichter
- “Subjectieve” stilistische analyse
- Controversieel
- “Hoorde” zelfde stem in *Elegast* en *Moriaen*
- Zelfde auteur?



Moriaen

- Koning Artur en Rondetafelridders
- Moriaen, zwarte ridder uit Afrika
- Op zoek naar biologische vader
- Eerste Nederlandse roman met zwart hoofdpersonage
- Ongebruikelijk onderwerp

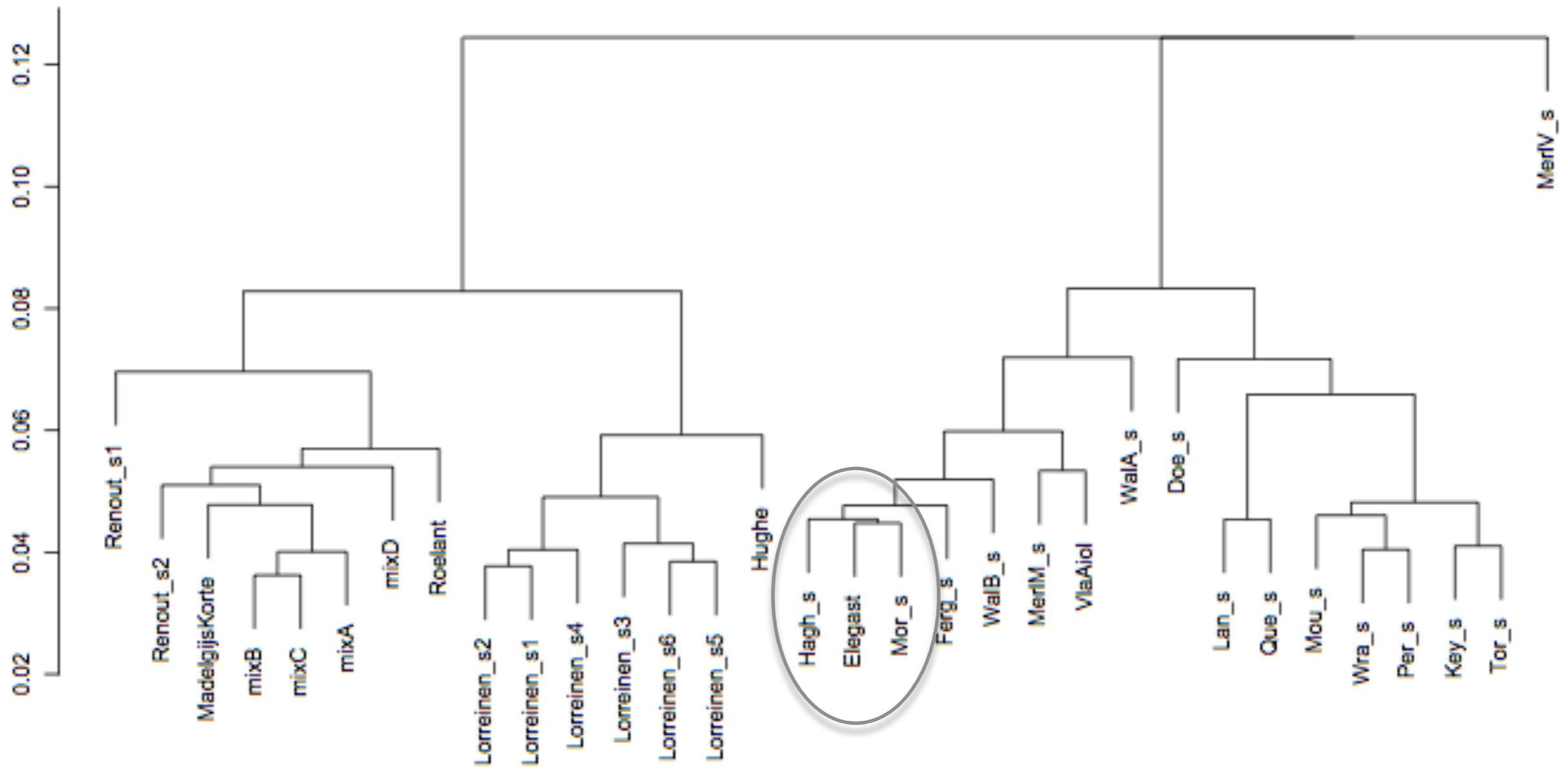


Icarus...

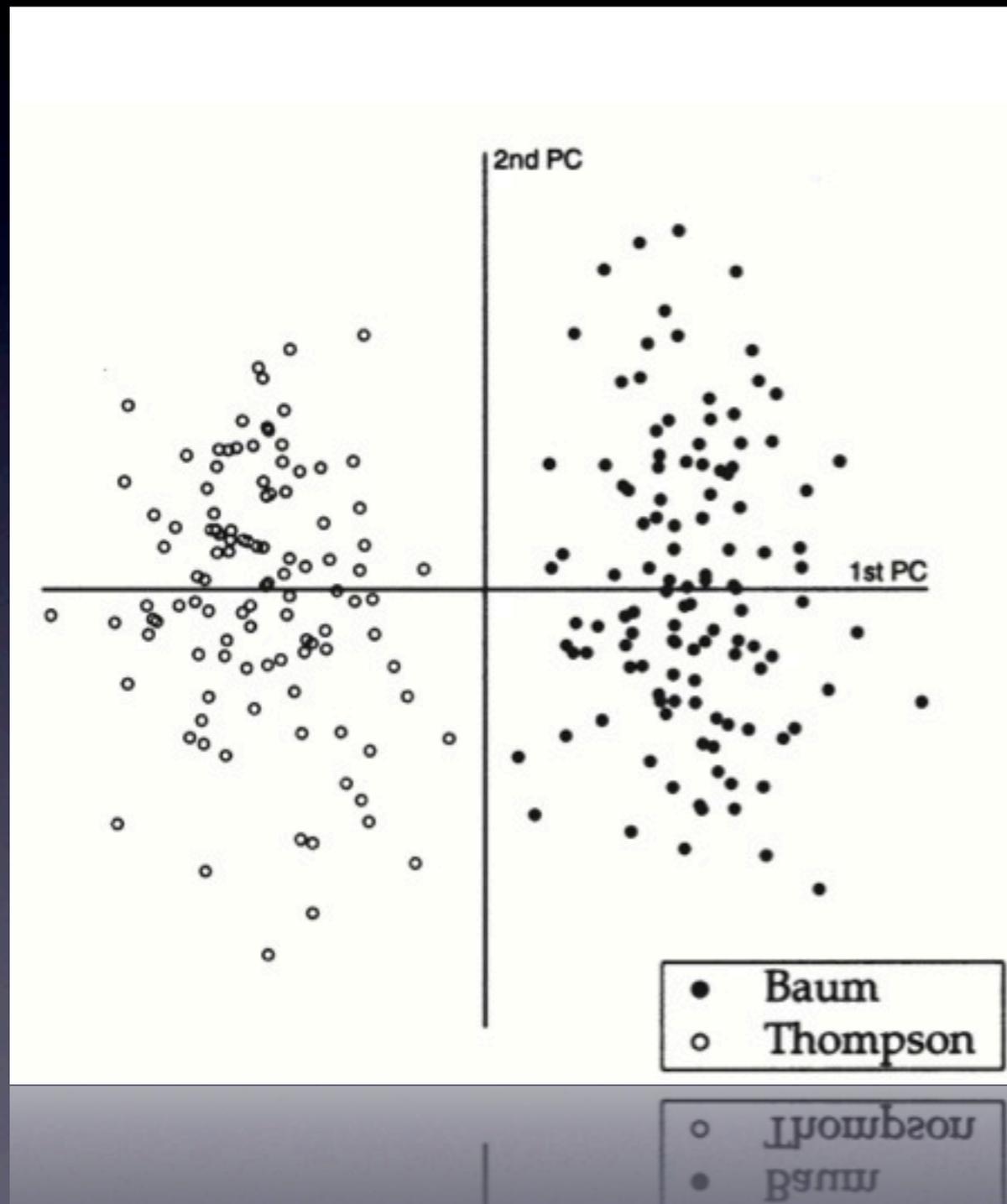


- Toeschrijving afgewezen...
- Té subjectief
- *Ohrenphilologie* ≠ wetenschap
- Bijnaam “Icarus”
- Stylometrie?!

Genres: Serendipiteit?



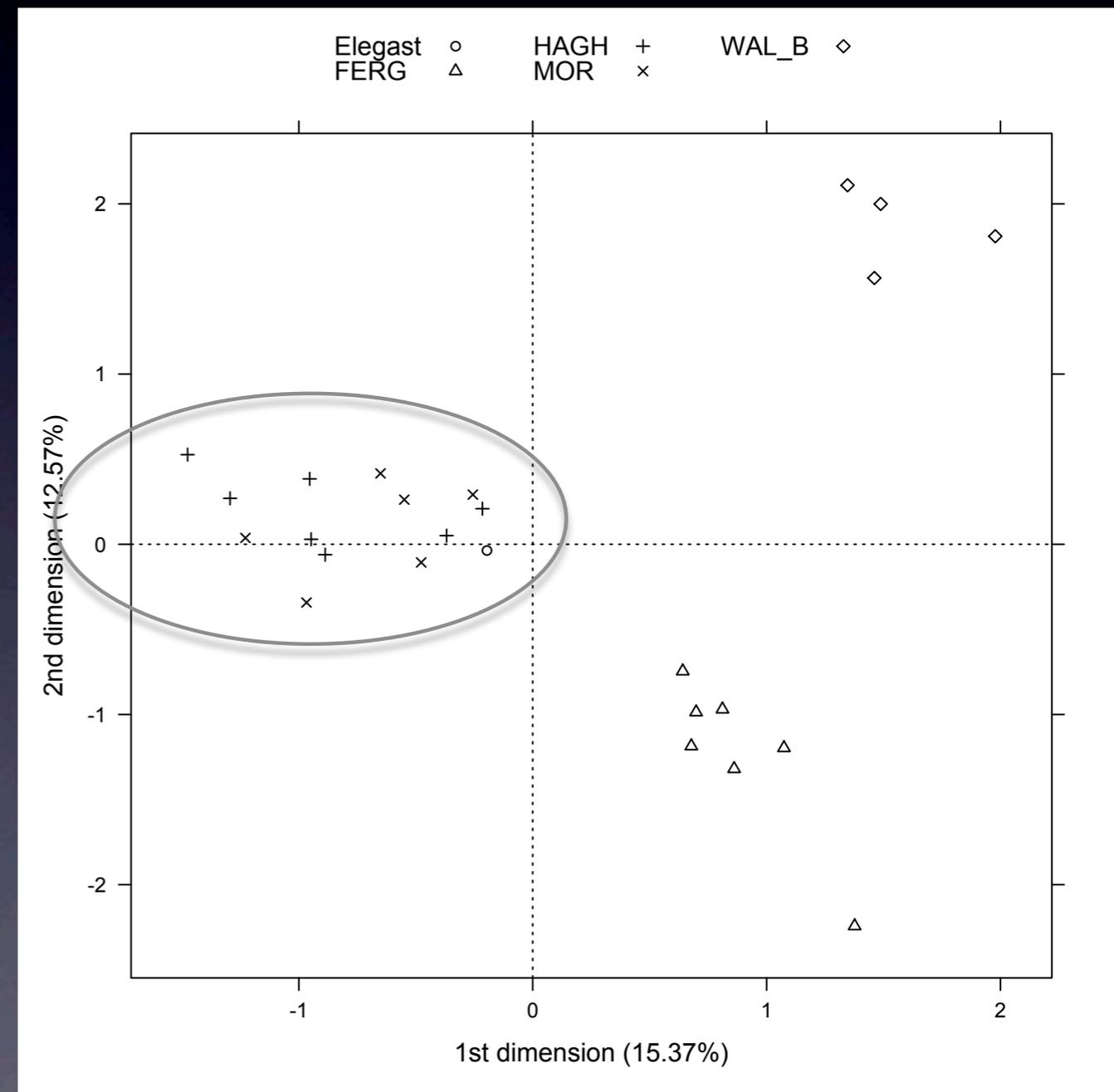
Dimensie-reductie



- Terugdringen # variabelen
- **Abstracte, latente dimensies**
- Heel populair in stylometrie
- Attributie via clusters
- Bv. PCA, CA, MDS, ...

Eerherstel?

- Correspondentie-analyse
- Sterkst verwant in ridderepiek
- **Meting bevestigt buikgevoel**
- Eerherstel Heeroma?
- (Maar Gruuthuse...)

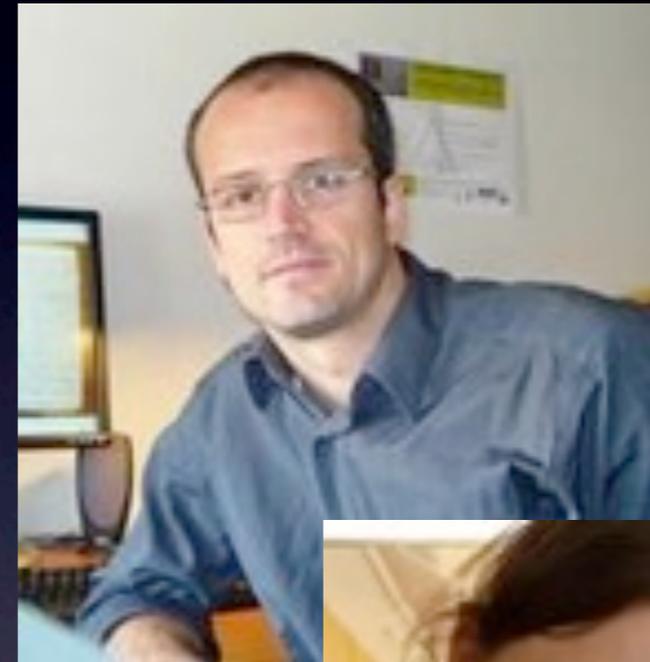


Momenteel...

- FWO-postdoc (2012-2015)
- Uitbreiding stylometrie:
 - proza
 - Latijn
- Veel **groter toepassingsgebied**

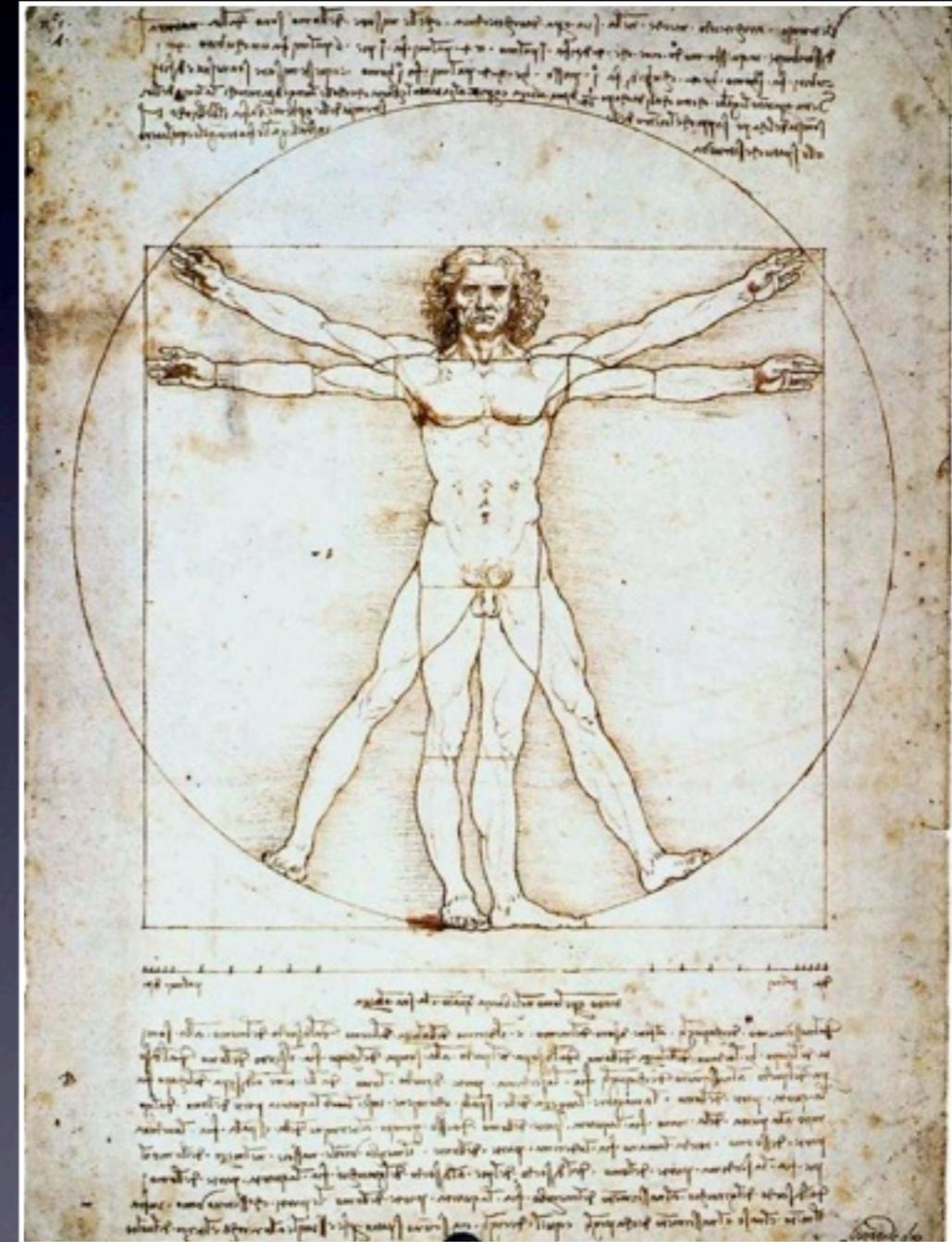
Latijnse mediëvistiek?

- J. Deploige & S. Moens
- Editie 2 korte tekstjes:
 - *Visio de sancto Martino*
 - *Visio ad Guibertum missa*
- Toegeschreven aan Hildegard van Bingen
- Twijfels...



Vitruvische man (LDO)

Renaissance 12e eeuw



Hildegard von Bingen

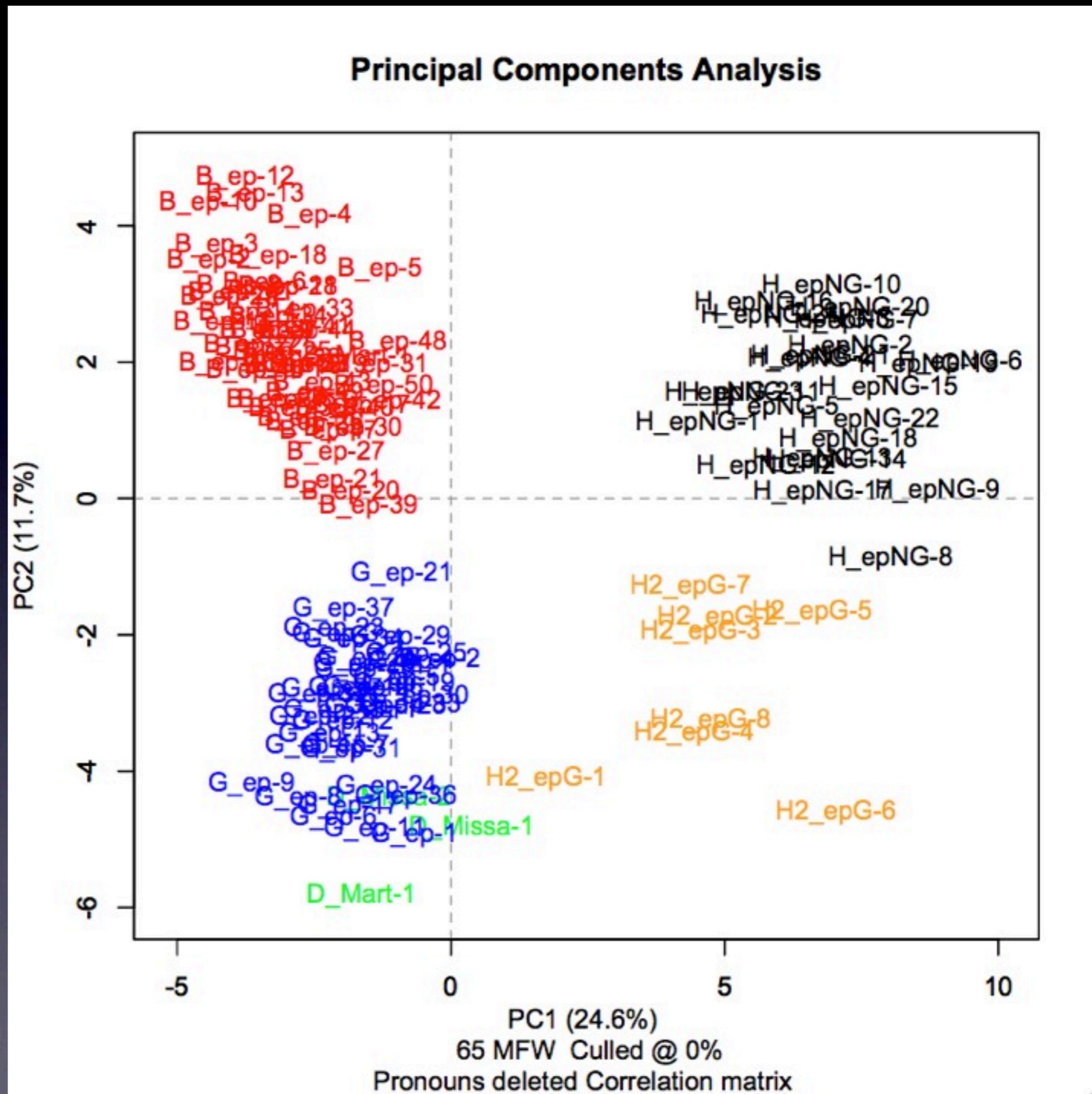
- Rijnlandse mystica (1098-1179)
- Hoog aanzien (Bernardus, paus, ...)
- Mystieke visioenen:
 - Dicteerde aan secretarissen
 - Latijn niet volledig machtig (*indocta*)
- Complex “auteurschap”



Laatste secretaris: Guibert van Gembloux

When you correct [the *Visio de sancto Martino*] you should keep to this rule: that adding, subtracting, and **changing nothing**, you apply your skill only to make corrections where the order or the rules of **correct Latin** are violated. Or if you prefer – and this is something I have conceded in this letter **beyond my normal practice** – you need not hesitate to clothe the whole sequence of the vision in **a more becoming garment of speech**.

- Stylometrische vingerafdruk
- Toeschrijving onder druk...
- Synergy hypothesis



Meerwaarde stylometrie?

- Erg gerichte applicatie in DH
- Meerwaarde:
 - Falsifiëren oude inzichten
 - Genereren nieuwe inzichten
- Serendipiteit als nevenproduct van *distant reading*



Toekomst?

- Breder dan stylometrie en mediëvistiek
- *TIME Magazine* archief (i.s.m. Folgert Karsdorp; 1920s-2000s)
- **Big Data & Humanities**: onderschat probleem
- Recente reeks papers:
 - toptijdschriften
 - “geesteswetenschappelijk”
 - retoriek van de Big Data
 - methodologie unaniem verworpen door vakgenoten

Google Books paper

- Science paper
- Google books corpus
- Michel et al.
- “Culturomics”
- Woordfrequenties
- Diachroon

RESEARCH ARTICLE

Quantitative Analysis of Culture Using Millions of Digitized Books

Jean-Baptiste Michel,^{1,2,3,4,5,6,7} Yuan Kui Shen,^{2,4,7} Ailva Presser Aiden,^{2,4,8} Adrian Vares,^{2,4,9} Matthew K. Gray,¹⁰ The Google Books Team,¹¹ Joseph P. Pickett,¹² Dale Holberg,¹³ Dan Clancy,¹⁴ Peter Norvig,¹⁵ Jon Orwant,¹⁶ Steven Pinker,⁵ Martin A. Nowak,^{1,7,17,18} Erez Lieberman Aiden^{1,2,4,7,11,12,17,19}

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.

Reading small collections of carefully chosen works enables scholars to make powerful inferences about trends in human thought. However, this approach rarely enables precise measurement of the underlying phenomena. Attempts to introduce quantitative methods into the study of culture (1–6) have been hampered by the lack of suitable data.

We report the creation of a corpus of 3,195,769 digitized books containing ~4% of all books ever published. Computational analysis of this corpus enables us to observe cultural trends and subject them to quantitative investigation. ‘Culturomics’ extends the boundaries of scientific inquiry to a wide array of new phenomena.

The corpus has emerged from Google’s effort to digitize books. Most books were drawn from over 40 university libraries around the world. Each page was scanned with custom equipment (7), and the text was digitized by means of optical character recognition (OCR). Additional volumes, both physical and digital, were contributed

by publishers. Metadata describing the date and place of publication were provided by the libraries and publishers and supplemented with bibliographic databases. Over 15 million books have been digitized [~12% of all books ever published (7)]. We selected a subset of over 5 million books for analysis on the basis of the quality of their OCR and metadata (Fig. 1A and fig. S1) (7). Periodicals were excluded.

The resulting corpus contains over 500 billion words, in English (161 billion), French (45 billion), Spanish (45 billion), German (37 billion), Chinese (13 billion), Russian (13 billion), and Hebrew (2 billion). The oldest works were published in the 1500s. The early decades are represented by only a few books per year, comprising several hundred thousand words. By 1800, the corpus grows to 98 million words per year; by 1900, 1.8 billion, and by 2000, 11 billion (fig. S2).

The corpus cannot be read by a human. If you tried to read only English-language entries from the year 2000 alone, at the reasonable pace of 200 words/min, without interruptions for food or sleep, it would take 80 years. The sequence of letters is 1000 times longer than the human genome: if you wrote it out in a straight line, it would reach to the Moon and back 90 times over (8).

To make release of the data possible in light of copyright constraints, we restricted this initial study to the question of how often a given 1-gram or *n*-gram was used over time. A 1-gram is a string of characters uninterrupted by a space; this includes words (“banana”, “SCLUBA”) but also numbers (“3.14159”) and typos (“essess”). An *n*-gram is a sequence of 1-grams, such as the phrases “stock market” (a 2-gram) and “the United States of America” (a 5-gram). We restricted *n* to 5 and limited our study to *n*-grams occurring at least 40 times in the corpus.

Usage frequency is computed by dividing the number of instances of the *n*-gram in a given year by the total number of words in the corpus in that year. For instance, in 1861, the 1-gram “slavery” appeared in the corpus 21,660 times, on 11,687

pages of 1298 books. The corpus contains 386,434,758 words from 1861; thus, the frequency is 5.5×10^{-7} . The use of “slavery” peaked during the Civil War (early 1860s) and then again during the civil rights movement (1955–1968) (Fig. 1B).

In contrast, we compare the frequency of “the Great War” to the frequencies of “World War I” and “World War II”. References to “the Great War” peak between 1915 and 1941. But although its frequency drops thereafter, interest in the underlying events had not disappeared; instead, they are referred to as “World War I” (Fig. 1C).

These examples highlight two central factors that contribute to culturomic trends. Cultural change guides the concepts we discuss (such as “slavery”). Linguistic change, which, of course, has cultural roots, affects the words we use for those concepts (“the Great War” versus “World War I”). In this paper, we examine both linguistic changes, such as changes in the lexicon and grammar, and cultural phenomena, such as how we remember people and events.

The full data set, which comprises over two billion culturomic trajectories, is available for download or exploration at www.culturomics.org and ngrams.googlelabs.com.

The size of the English lexicon. How many words are in the English language (9)?

We call a 1-gram “common” if its frequency is greater than one per billion. [This corresponds to the frequency of the words listed in leading dictionaries (7) (fig. S3)]. We compiled a list of all common 1-grams in 1900, 1950, and 2000, based on the frequency of each 1-gram in the preceding decade. These lists contained 1,117,997 common 1-grams in 1900, 1,102,920 in 1950, and 1,489,337 in 2000.

Not all common 1-grams are English words. Many fell into three nonword categories: (i) 1-grams with nonalphabetic characters (“%”, “3.14159”), (ii) misspellings (“because”, “tabben”), and (iii) foreign words (“sensitivo”).

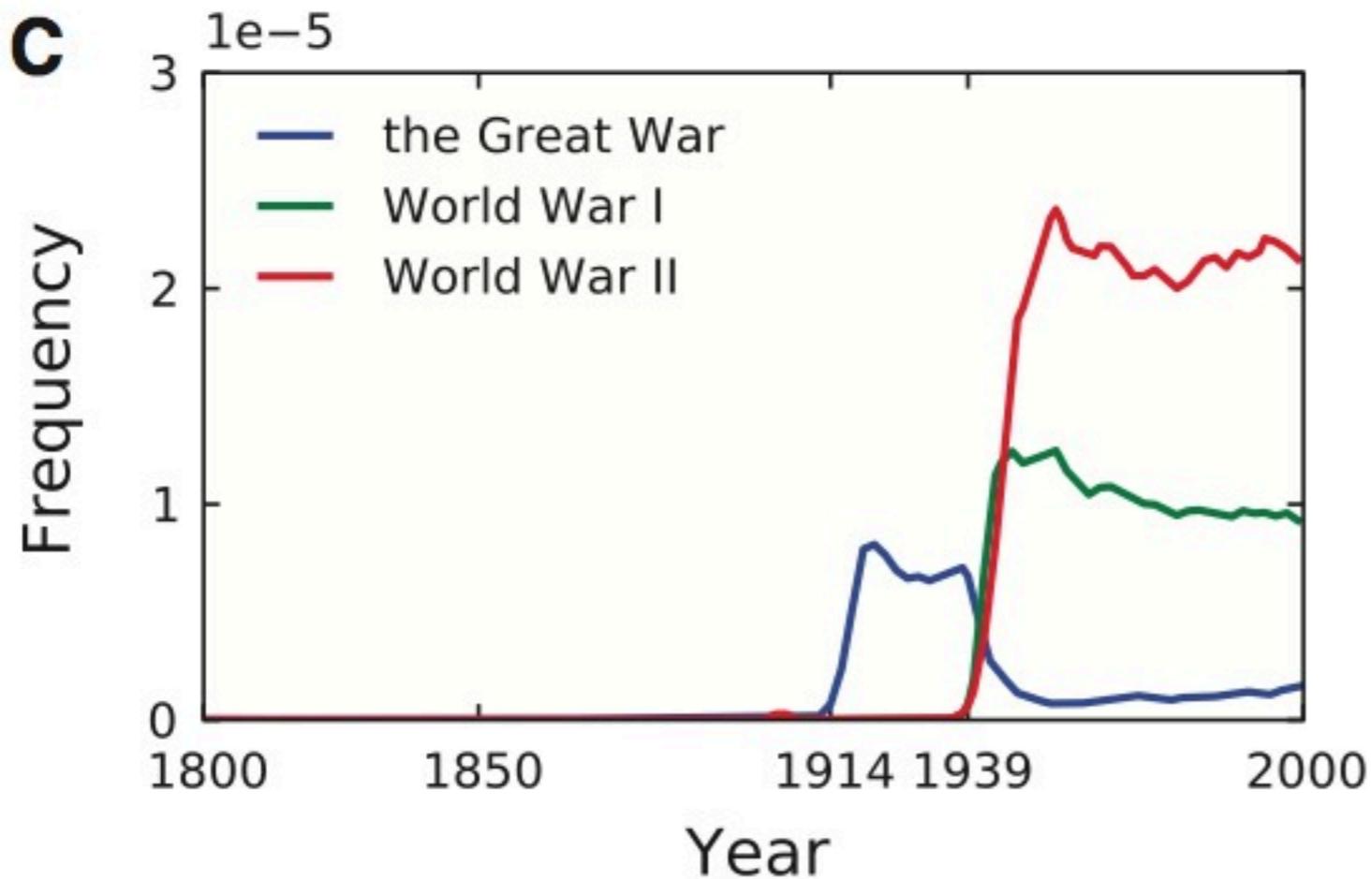
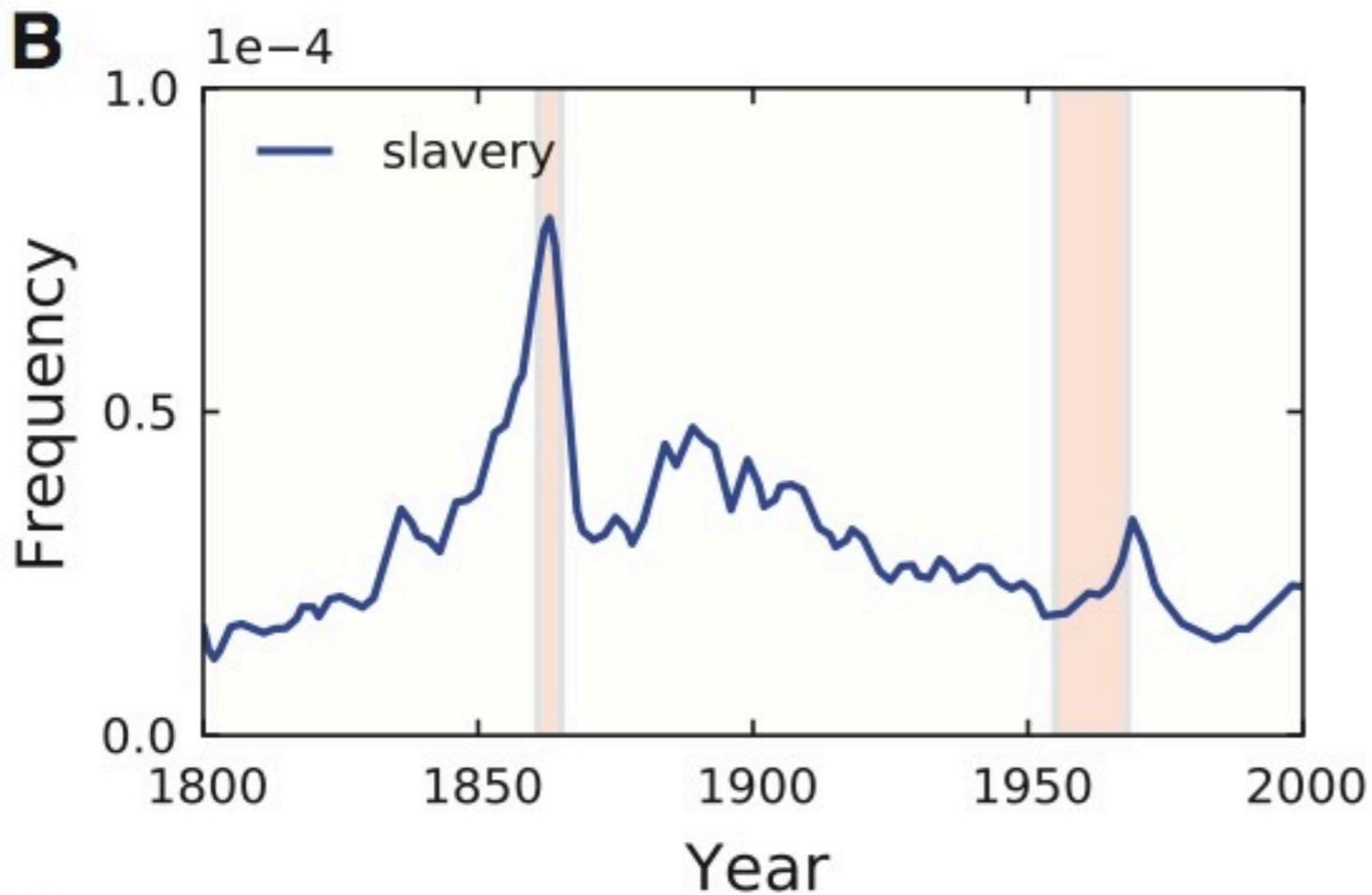
To estimate the number of English words, we manually annotated random samples from the lists of common 1-grams (7) and determined what fraction were members of the above nonword categories. The result ranged from 51% of all common 1-grams in 1900 to 33% in 2000.

Using this technique, we estimated the number of words in the English lexicon as 544,000 in 1900, 597,000 in 1950, and 1,022,000 in 2000. The lexicon is enjoying a period of enormous growth: The addition of ~8500 words/year has increased the size of the language by over 70% during the past 50 years (Fig. 2A).

Notably, we found more words than appear in any dictionary. For instance, the 2002 *Webster’s Third New International Dictionary* (W3), which keeps track of the contemporary American lexicon, lists approximately 348,000 single-word wordforms (10); the *American Heritage Dictionary of the English Language, Fourth Edition* (AHD4) lists 136,361 (11). (Both contain additional multiform entries.) Part of this gap is because dictionaries often

¹Program for Evolutionary Dynamics, Harvard University, Cambridge, MA 02138, USA. ²Cultural Observatory, Harvard University, Cambridge, MA 02138, USA. ³Institute for Quantitative Social Sciences, Harvard University, Cambridge, MA 02138, USA. ⁴Department of Psychology, Harvard University, Cambridge, MA 02138, USA. ⁵Department of Systems Biology, Harvard Medical School, Boston, MA 02115, USA. ⁶Laboratory of Large, Harvard University, Cambridge, MA 02138, USA. ⁷Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology (MIT), Cambridge, MA 02139, USA. ⁸Harvard Medical School, Boston, MA 02115, USA. ⁹Harvard College, Cambridge, MA 02138, USA. ¹⁰Google, Mountain View, CA 94043, USA. ¹¹Knightsbridge, London, UK. ¹²Department of Organismic and Evolutionary Biology, Harvard University, Cambridge, MA 02138, USA. ¹³Department of Mathematics, Harvard University, Cambridge, MA 02138, USA. ¹⁴Broad Institute of Harvard and MIT, Harvard University, Cambridge, MA 02138, USA. ¹⁵School of Engineering and Applied Sciences, Harvard University, Cambridge, MA 02138, USA. ¹⁶Harvard Society of Fellows, Harvard University, Cambridge, MA 02138, USA.

¹⁷These authors contributed equally to this work. [To whom correspondence should be addressed, E-mail: jb.michel@gmail.com (J.-B.M.); eraz@mit.edu (E.L.A.)]



Tikje **simplistisch** (?),
 maar hoe beter doen?

Ultraconserved words point to deep language ancestry across Eurasia

Mark Pagel^{1,2,3,4}, Quentin D. Atkinson⁵, Andrea S. Calude⁶, and Andrew Meade⁷

¹School of Biological Sciences, University of Reading, Reading, Berkshire RG6 2AA, United Kingdom; ²Santa Fe Institute, Santa Fe, NM 87501; ³School of Psychology, University of Auckland, Auckland 1142, New Zealand; ⁴Linguistics Programme, University of Waikato, Hamilton 3240, New Zealand

Edited by Colin Renfrew, University of Cambridge, Cambridge, United Kingdom, and approved April 15, 2012 (received for review October 21, 2011)

The search for ever deeper relationships among the world's languages is bedeviled by the fact that most words evolve too rapidly to preserve evidence of their ancestry beyond 5,000 to 8,000 y. On the other hand, quantitative modeling indicates that some "ultraconserved" words exist that might be used to find evidence for deep linguistic relationships beyond that time barrier. Here we use a statistical model, which takes into account the frequency with which words are used in common everyday speech, to predict the existence of a set of such highly conserved words among seven language families of Eurasia postulated to form a linguistic superfamily that evolved from a common ancestor around 15,000 y ago. We derive a dated phylogenetic tree of this proposed superfamily with a time-depth of ~14,450 y, implying that some frequently used words have been retained in related forms since the end of the last ice age. Words used more than once per 1,000 in everyday speech were 7- to 10-times more likely to show deep ancestry on this tree. Our results suggest a remarkable fidelity in the transmission of some words and give theoretical justification to the search for features of language that might be preserved across wide spans of time and geography.

cultural evolution | phylogeny | historical linguistics

The English word *brother* and the French *frère* are related to the Sanskrit *bhrāta* and the Latin *frater*, suggesting that words as more sounds can remain associated with the same meaning for millennia. But how far back in time can traces of a word's genealogical history persist, and can we predict which words are likely to show deep ancestry?

These questions are central to ancient times and to efforts to identify linguistic world's languages (1-5). Evidence for 30 such as Amerind (6), linking most of the New World, and Nostratic (7-9) and Fin the major language families of Eurasia-identification of putative "cognate" words in history), the sound and meaning correspondence are expected to last long enough to indicate that they derive from a common ancestor. Such evidence is often criticized for it words are thought to suffer from too much erosion to allow secure identification of 5,000 to 8,000 y (11, 12), and second, even cognates can be identified, proponents have been unable to provide statistical evidence that they are not the result of chance between unrelated languages (13). Both objections can be overcome if a class of words exists whose members' soundness are expected to last long enough to indicate that they derive from a common ancestor. Such ultraconserved words are thought to be preserved across wide spans of time and geography.

www.pnas.org/cgi/doi/10.1073/pnas.1119407109

RESEARCH ARTICLE

Quantitative Analysis of Culture Using Millions of Digitized Books

Jean-Baptiste Michel^{1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35,36,37,38,39,40,41,42,43,44,45,46,47,48,49,50,51,52,53,54,55,56,57,58,59,60,61,62,63,64,65,66,67,68,69,70,71,72,73,74,75,76,77,78,79,80,81,82,83,84,85,86,87,88,89,90,91,92,93,94,95,96,97,98,99,100,101,102,103,104,105,106,107,108,109,110,111,112,113,114,115,116,117,118,119,120,121,122,123,124,125,126,127,128,129,130,131,132,133,134,135,136,137,138,139,140,141,142,143,144,145,146,147,148,149,150,151,152,153,154,155,156,157,158,159,160,161,162,163,164,165,166,167,168,169,170,171,172,173,174,175,176,177,178,179,180,181,182,183,184,185,186,187,188,189,190,191,192,193,194,195,196,197,198,199,200,201,202,203,204,205,206,207,208,209,210,211,212,213,214,215,216,217,218,219,220,221,222,223,224,225,226,227,228,229,230,231,232,233,234,235,236,237,238,239,240,241,242,243,244,245,246,247,248,249,250,251,252,253,254,255,256,257,258,259,260,261,262,263,264,265,266,267,268,269,270,271,272,273,274,275,276,277,278,279,280,281,282,283,284,285,286,287,288,289,290,291,292,293,294,295,296,297,298,299,300,301,302,303,304,305,306,307,308,309,310,311,312,313,314,315,316,317,318,319,320,321,322,323,324,325,326,327,328,329,330,331,332,333,334,335,336,337,338,339,340,341,342,343,344,345,346,347,348,349,350,351,352,353,354,355,356,357,358,359,360,361,362,363,364,365,366,367,368,369,370,371,372,373,374,375,376,377,378,379,380,381,382,383,384,385,386,387,388,389,390,391,392,393,394,395,396,397,398,399,400,401,402,403,404,405,406,407,408,409,410,411,412,413,414,415,416,417,418,419,420,421,422,423,424,425,426,427,428,429,430,431,432,433,434,435,436,437,438,439,440,441,442,443,444,445,446,447,448,449,450,451,452,453,454,455,456,457,458,459,460,461,462,463,464,465,466,467,468,469,470,471,472,473,474,475,476,477,478,479,480,481,482,483,484,485,486,487,488,489,490,491,492,493,494,495,496,497,498,499,500,501,502,503,504,505,506,507,508,509,510,511,512,513,514,515,516,517,518,519,520,521,522,523,524,525,526,527,528,529,530,531,532,533,534,535,536,537,538,539,540,541,542,543,544,545,546,547,548,549,550,551,552,553,554,555,556,557,558,559,560,561,562,563,564,565,566,567,568,569,570,571,572,573,574,575,576,577,578,579,580,581,582,583,584,585,586,587,588,589,590,591,592,593,594,595,596,597,598,599,600,601,602,603,604,605,606,607,608,609,610,611,612,613,614,615,616,617,618,619,620,621,622,623,624,625,626,627,628,629,630,631,632,633,634,635,636,637,638,639,640,641,642,643,644,645,646,647,648,649,650,651,652,653,654,655,656,657,658,659,660,661,662,663,664,665,666,667,668,669,670,671,672,673,674,675,676,677,678,679,680,681,682,683,684,685,686,687,688,689,690,691,692,693,694,695,696,697,698,699,700,701,702,703,704,705,706,707,708,709,710,711,712,713,714,715,716,717,718,719,720,721,722,723,724,725,726,727,728,729,730,731,732,733,734,735,736,737,738,739,740,741,742,743,744,745,746,747,748,749,750,751,752,753,754,755,756,757,758,759,760,761,762,763,764,765,766,767,768,769,770,771,772,773,774,775,776,777,778,779,780,781,782,783,784,785,786,787,788,789,790,791,792,793,794,795,796,797,798,799,800,801,802,803,804,805,806,807,808,809,810,811,812,813,814,815,816,817,818,819,820,821,822,823,824,825,826,827,828,829,830,831,832,833,834,835,836,837,838,839,840,841,842,843,844,845,846,847,848,849,850,851,852,853,854,855,856,857,858,859,860,861,862,863,864,865,866,867,868,869,870,871,872,873,874,875,876,877,878,879,880,881,882,883,884,885,886,887,888,889,890,891,892,893,894,895,896,897,898,899,900,901,902,903,904,905,906,907,908,909,910,911,912,913,914,915,916,917,918,919,920,921,922,923,924,925,926,927,928,929,930,931,932,933,934,935,936,937,938,939,940,941,942,943,944,945,946,947,948,949,950,951,952,953,954,955,956,957,958,959,960,961,962,963,964,965,966,967,968,969,970,971,972,973,974,975,976,977,978,979,980,981,982,983,984,985,986,987,988,989,990,991,992,993,994,995,996,997,998,999,1000}

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 use therein of "ultraconserved" words, 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 zoology, the evolution of grammar, collective memory, the adoption of technology, the pursuit of fame, ownership, and historical epidemiology. Cultureless extends the boundaries of rigorous quantitative inquiry to a wide array of new phenomena spanning the social sciences and the humanities.

Reading small collections of carefully chosen works enables scholars to make pertinent inferences about trends in human thought. However, this approach rarely enables precise measurement of the underlying phenomena. Attempts to introduce quantitative methods into the study of culture (1-4) have been hampered by the lack of suitable data.

We report the creation of a corpus of 5,195,769 digitized books containing ~4% of all books ever published. Computational analysis of this corpus enables us to observe cultural trends and subject them to quantitative investigation. "Cultureless" extends the boundaries of scientific inquiry to a wide array of new phenomena.

The corpus has emerged from Google's effort to digitize books. Most books were drawn from over 40 university libraries around the world. Each page was scanned with custom equipment (7), and the text was digitized by means of optical character recognition (OCR). Additional volumes, both physical and digital, were contributed

¹Program for Evolutionary Dynamics, Harvard University, Cambridge, MA 02138, USA; ²Cultural Observatory, Harvard University, Cambridge, MA 02138, USA; ³Institute for Quantitative Social Sciences, Harvard University, Cambridge, MA 02138, USA; ⁴Department of Psychology, Harvard University, Cambridge, MA 02138, USA; ⁵Department of Systemic Biology, Harvard Medical School, Boston, MA 02115, USA; ⁶University of Laquila, Italy; ⁷Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁸Department of Physics, Harvard University, Cambridge, MA 02138, USA; ⁹Department of Computer Science, Harvard University, Cambridge, MA 02138, USA; ¹⁰Department of Statistics, Harvard University, Cambridge, MA 02138, USA; ¹¹Department of Applied Mathematics, Harvard University, Cambridge, MA 02138, USA; ¹²Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ¹³Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ¹⁴Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ¹⁵Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ¹⁶Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ¹⁷Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ¹⁸Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ¹⁹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ²⁰Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ²¹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ²²Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ²³Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ²⁴Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ²⁵Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ²⁶Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ²⁷Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ²⁸Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ²⁹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ³⁰Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ³¹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ³²Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ³³Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ³⁴Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ³⁵Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ³⁶Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ³⁷Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ³⁸Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ³⁹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁴⁰Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁴¹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁴²Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁴³Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁴⁴Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁴⁵Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁴⁶Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁴⁷Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁴⁸Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁴⁹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁵⁰Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁵¹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁵²Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁵³Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁵⁴Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁵⁵Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁵⁶Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁵⁷Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁵⁸Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁵⁹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁶⁰Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁶¹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁶²Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁶³Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁶⁴Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁶⁵Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁶⁶Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁶⁷Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁶⁸Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁶⁹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁷⁰Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁷¹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁷²Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁷³Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁷⁴Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁷⁵Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁷⁶Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁷⁷Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁷⁸Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁷⁹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁸⁰Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁸¹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁸²Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁸³Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁸⁴Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁸⁵Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁸⁶Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁸⁷Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁸⁸Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁸⁹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁹⁰Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁹¹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁹²Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁹³Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁹⁴Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁹⁵Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁹⁶Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁹⁷Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁹⁸Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ⁹⁹Department of Mathematics, Harvard University, Cambridge, MA 02138, USA; ¹⁰⁰Department of Mathematics, Harvard University, Cambridge, MA 02138, USA

Ancient symbols, computational linguistics, and the reviewing practices of the general science journals

Richard Sproat¹
Center for Spoken Language Understanding

1. Introduction

Few archaeological finds are as evocative as artifacts inscribed with symbols. Wherever an archaeologist finds a potsherd or a seal impression that seems to have symbols scratched or impressed on the surface, it is natural to want to "read" the symbols. As if the symbols come from an undeciphered or previously unknown symbol system it is common to ask what language the symbols supposedly represent and whether that system can be deciphered.

One of the first questions that really should be asked is whether the symbols are a writing system, as linguists usually define it, is a symbol system that uses a language. Familiar examples are alphabets such as the Latin, Greek, and Cyrillic alphabets, alphasyllabaries such as Devanagari or Tamil, syllabaries such as the Japanese Kana, and morphosyllabic systems like Chinese characters. But if a symbol system does not encode language, as in the case of heraldry, mathematics (used to represent dance), and boy scout merit badges are not writing systems that represent things, but do not function as part of a writing system.

Whether a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this, but these have always been used under the assumption that the symbols are writing, and the techniques are used to uncover patterns that might aid in the decipherment. Patterns of symbol distribution might indicate a symbol system is not linguistic; for example, odd repetition patterns that a symbol system is unlikely to be writing. But until recently, the only statistical techniques could be used to determine that a symbol system is writing or not is a difficult question to answer definitively in the affirmative if one can develop a verifiable system of language or languages. Statistical techniques have been used to help with this,

Discussie

- Goed voor **visibiliteit** vakgebied
- **Kruisbestuiving** exacte wetenschappen
- **Afgunst** speelt mee?
- Maar ook beschermen eigen vakgebied:
 - Retoriek Big Data wordt 'misbruikt'...
 - Methodologie unaniem afgewezen...
 - Onderzoekers zijn allesbehalve geesteswetenschappers...
- **Hebben de Geesteswetenschappen een eigen 'toptijdschrift' nodig?**