

Wenn die Quellen überfließen. Spitzweg und Big Data



**BIG
DATA**

*Manfred Thaller,
Universität zu Köln*

*Graz, 27. Februar
2015*

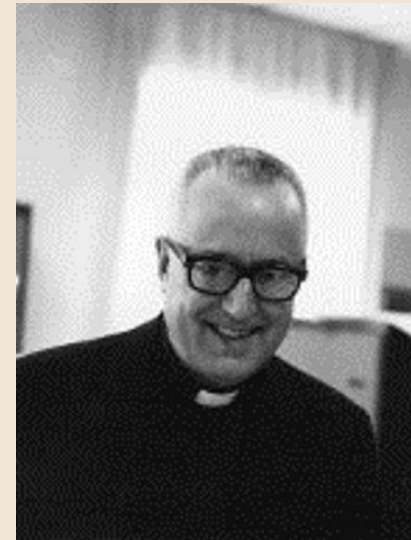
DHd2015

... Gelehrte

Donit on melioth
melior tanta rancid
a nat a pypus quib
ois fillo agnus in
opt un meli non
pugnat ut ell p
1101100011110010010
Historisch
1010111010101011111
Kulturwissenschaftliche
1010010100100101010
Informationsverarbeitung



Spitzweg, Gelehrter Mönch am Schreibtisch

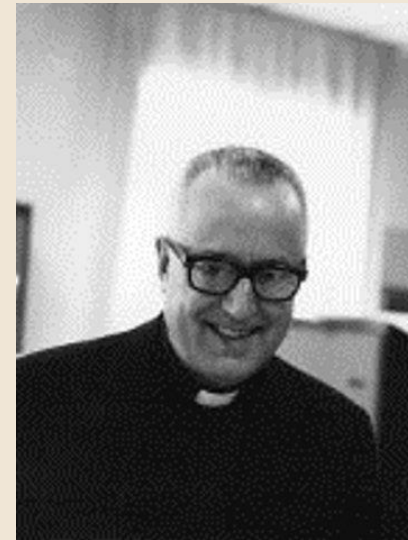


Padre Roberto Busa, 1913 -2011

Science Fiction Autoren und Gelehrte



E. E. „Doc“ Smith, 1890 - 1966



Padre Roberto Busa, 1913 -2011

Häresie 1

Es ist für Geisteswissenschaftlerinnen *nicht*
zwingend notwendig, fantasielos und intellektuell
ängstlich zu sein.

„Big Data“ und Gelehrte



Computer: neue Prosopographien

Arnold J. Toynbee, 1889-1975

„Big Data“ und Gelehrte



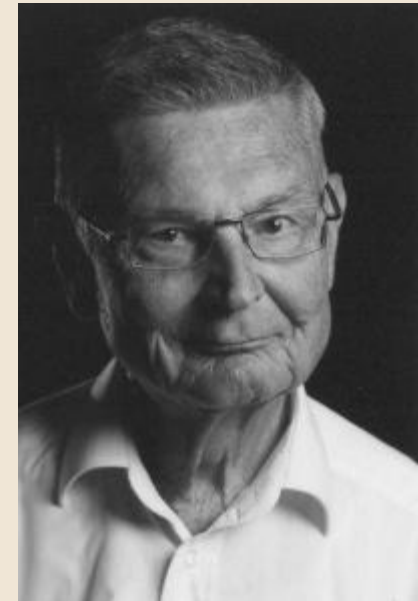
[Computer:] neue Prosopographien **Prosopographia Regnorum Orbis Latini**
Karl Ferdinand Werner, 1924 - 2008

„Big Data“ und Gelehrte

Donit on meloflu
meloflu rāna rāna
ā nū a pūpū rānū
ōs fello agūm
op nū melū nū
pūgnū nū ellē pū

11011000111110010010
10101110101010111111
1010010100100101010

Historisch
Kulturwissenschaftliche
Informationsverarbeitung



Computer: (?) neue Prosopographien

**Memorialüberlieferung, 1974 ff.
[Karl Schmid,] Joachim Wollasch**

Häresie 2

Auch methodisch sehr konservative Grundhaltungen sind mit dem Wunsch vereinbar, grundsätzlich Neues zu wagen.

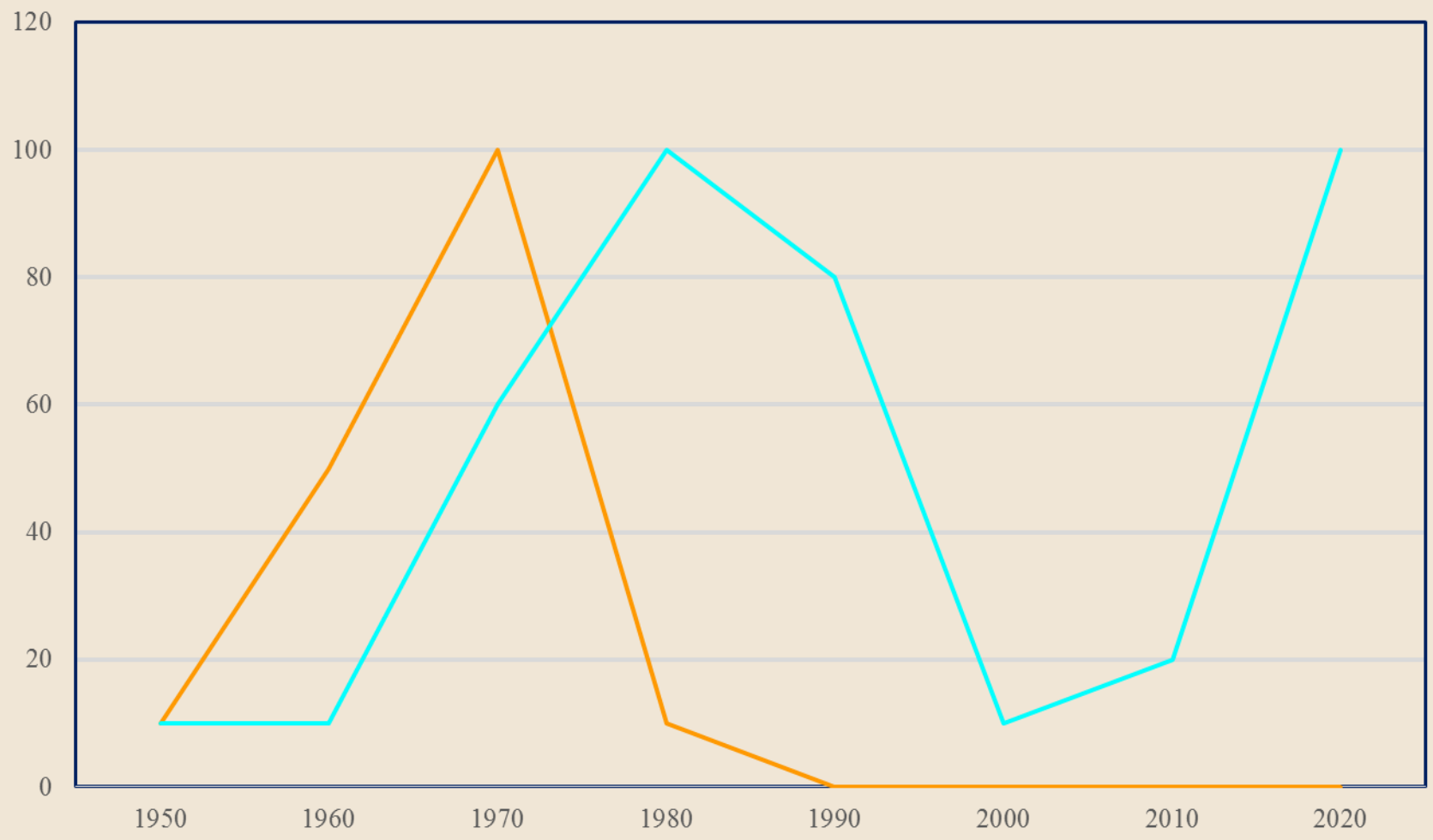
Das kann auch *ohne* IT gründlich schief gehen.

„Big Data“ und Gelehrte



<http://republicofletters.stanford.edu/>

Massenprosopographien



— Massenprosopographie manuell — Massenprosopographie IT 10

... Gelehrte

Donit on meloflu
meloflu rāna rāna
ā nā a pūpū rāmbi
ōs fello aggras
op nūc nūc nūc
pūgnā nūc ellē pūgnā

11011000111110010010
10101110101010111111
1010010100100101010

Historisch
Kulturwissenschaftliche
Informationsverarbeitung



Sir Anthony Kenny

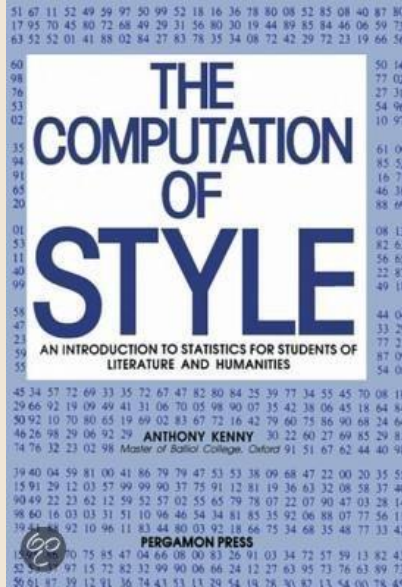
... Gelehrte



Sir Anthony Kenny

Geschichte der abendländischen Philosophie, 2014

... Gelehrte



Sir Anthony Kenny

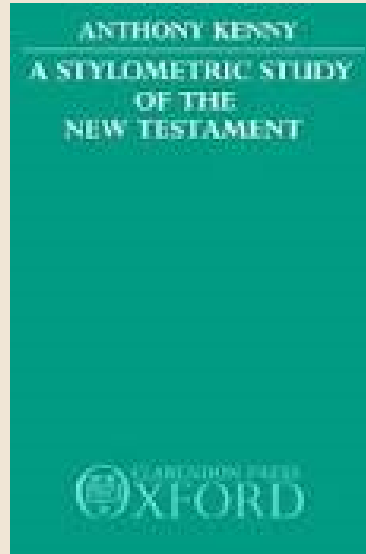
The computation of style, 1981

... Gelehrte

Donit on meliothi...
melior tanta rancia...
n'na a pypus...
ov fello...
op nix...
pugnat...

11011000111110010010
10101110101010111111
1010010100100101010

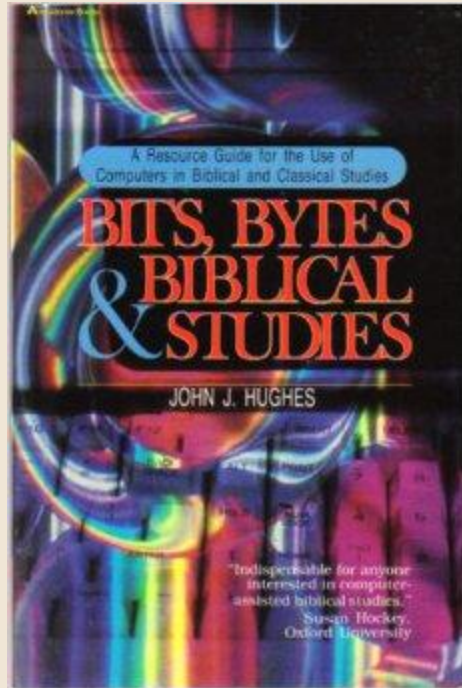
Historisch
Kulturwissenschaftliche
Informationsverarbeitung



Sir Anthony Kenny

A Stylometric Study of the New Testament, 1986

Gelehrte im Kontext



John J. Hughes: Bits, Bytes and Biblical Studies: A Resource Guide for the Use of Computers in Biblical and Classical Studies, 1987

Lexikonformat, 643 pages

Ein Blick zurück: 1993

Neben den Vorteilen enormer Speicherkapazitäten historischer Fakten und Daten rückt nun immer stärker auch die Möglichkeit präziser Dokumentationen von bildlichen Dokumenten durch EDV-Verarbeitungen in den Vordergrund der moderner Geschichtsforschung – und natürlich die rasche Verfügbarkeit und Überprüfung dieser Dokumente durch die internationale Kommunikationsvernetzung.

Ein Blick zurück: 1993

Digitale
Historisch
Kulturwissenschaftliche
Informationsverarbeitung

Geschichtsforschung per Computer

Die modernen Historiker arbeiten mit weltweiten Netzwerken

Walter Müller

Graz – Bis in die frühen Sechziger Jahre blieb die historische Wissenschaft beschränkt auf einige wenige monolithische Blöcke wie Wirtschafts-, Rechts- oder zuletzt Sozialgeschichte. In den folgenden Jahren verästelten sich die Geschichtswissenschaften in rund 40 bis 50 unterschiedliche Disziplinen. Ein Teil der Wissenschaftler begann, verstärkt elektronische Datenverarbeitung in ihren Forschungen, die immer

vernetzter wurden, einzusetzen. Dabei entwickelten sich die Historiker von reinen Computeranwendern zu Programmentwicklern.

Neben den Vorteilen enormer Speicherkapazitäten historischer Fakten und Daten rückt nun immer stärker auch die Möglichkeit präziser Dokumentationen von bildlichen Dokumenten durch EDV-Verarbeitungen in den Vordergrund der moderner Geschichtsforschung – und natürlich die rasche Verfügbarkeit und Überprüfung dieser Dokumente durch die internationale Kommunikationsvernetzung.

150 Historiker aus 23 Ländern trafen sich dieser Tage in Graz zur 8. Jahrestagung der internationalen „Association for History and Computing“ (AHC), um Erfahrungen dieser neuen computerunterstützten Forschung auszutauschen. Erstmals wurden auch Wissenschaftler aus dem südosteuropäischen Raum verstärkt ins wissenschaftliche Gespräch eingebunden.

Leonid Borodkin von der Moscow State University erinnerte daran, daß die osteuropäischen Universitäten nach wie vor unter einem Mangel an elektronischen Ressourcen leiden. Westeuropäische Wissenschaftler versuchen daher seit einiger Zeit, östliche Universitäten mit EDV-Anlagen auszustatten. Erst kürzlich sei ein Rechenpool für Geisteswissenschaftliche Institute in Estland von norwegischen Universitäten finanziert worden, ergänzte AHC-Chef Manfred Thaller.

Es sei im Osten gegenwärtig klimatisch eine günstige Zeit für die Geisteswissenschaften, so Thaller. In Polen etwa stehe momentan eine Sprach- und EDV- Ausbildung weit höher im Kurs als die gängigen betriebswirtschaftlichen Studien.

Eine der zentralen Aufgaben der künftigen EDV-unterstützten Geschichtswissenschaften: Die Weiterentwicklung international einheitlicher Kommunikationsebenen und Parameter.

ARD THEMEN
puter-
ologien

Ein Blick zurück: 1993



Ein Blick zurück: 1993

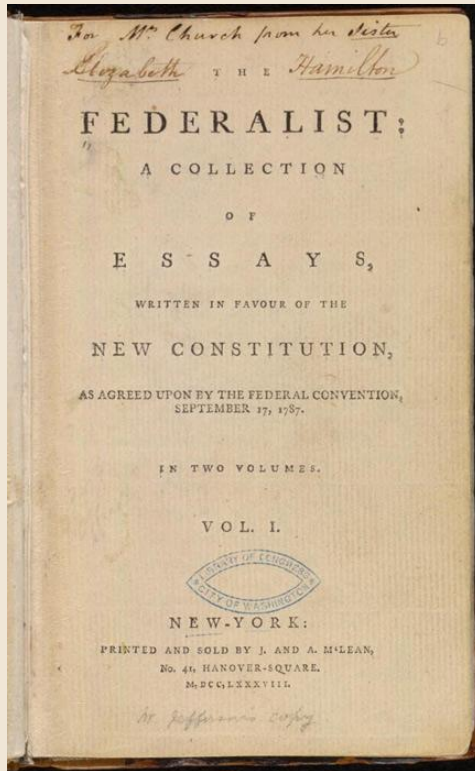
1011000111110010010
Historisch
Kulturwissenschaftliche
Informationsverarbeitung



... Gelehrte

Donit on melioth...
melior tanta rancid...
ā nū a pūmū cūbū...
ōis fello agūmū...
op nūc mēlū nōp...
pūgnā nō ellē pū...

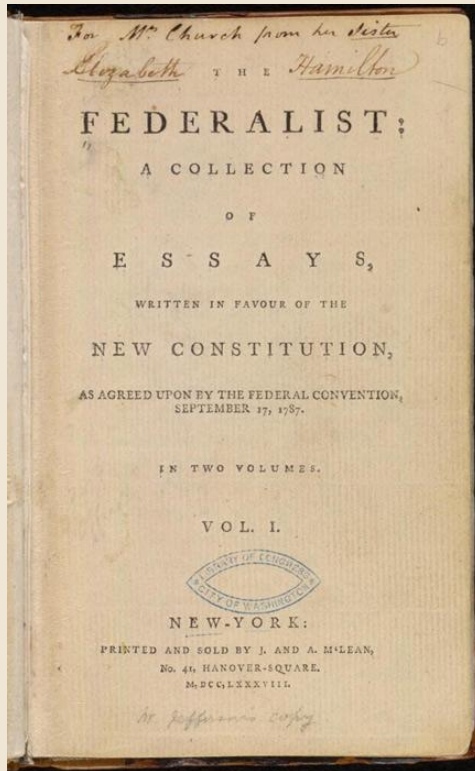
11011000111110010010
Historisch
1010111010101011111
Kulturwissenschaftliche
1010010100100101010
Informationsverarbeitung



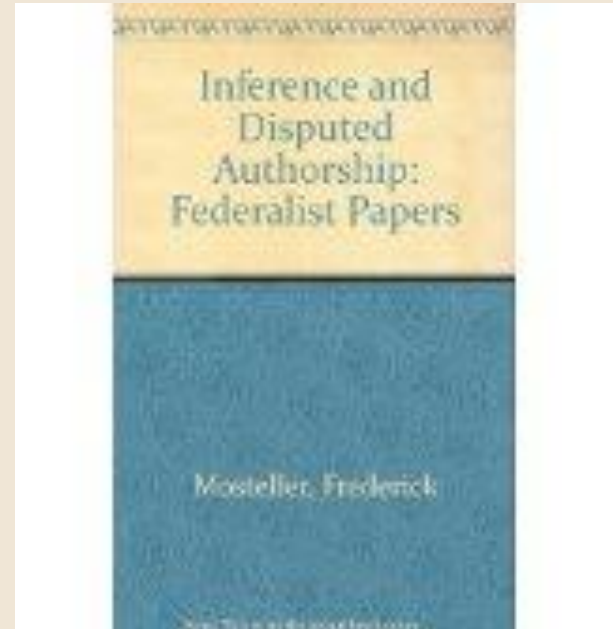
Frederick Williams & Frederick Mosteller, 1941

The Federalist Papers, 1786

Anfänge



The Federalist Papers, 1786



**Frederick Mosteller and David L. Wallace,
Inference and Disputed Authorship: *The Federalist*,
1964**

Anfänge

Philip J. Stone: The General Inquirer, 1961



IBM 7090, 1960er

2015: <http://www.wjh.harvard.edu/~inquirer/>

Status quo

Summary

- There is 1 document in the corpus with a total of 6,939 words and 1,999 unique words.
- Most frequent words in the corpus: **military** (83), **soldiers** (56), **soldier** (49), **men** (26), **experience** (22), More...

Words in the Entire Corpus	Count	Trend
military	83	
soldiers	56	
soldier	49	
men	26	
experience	22	
forces	22	
life	22	
period	21	

Word Trends

Relative Frequency

Keywords in Context

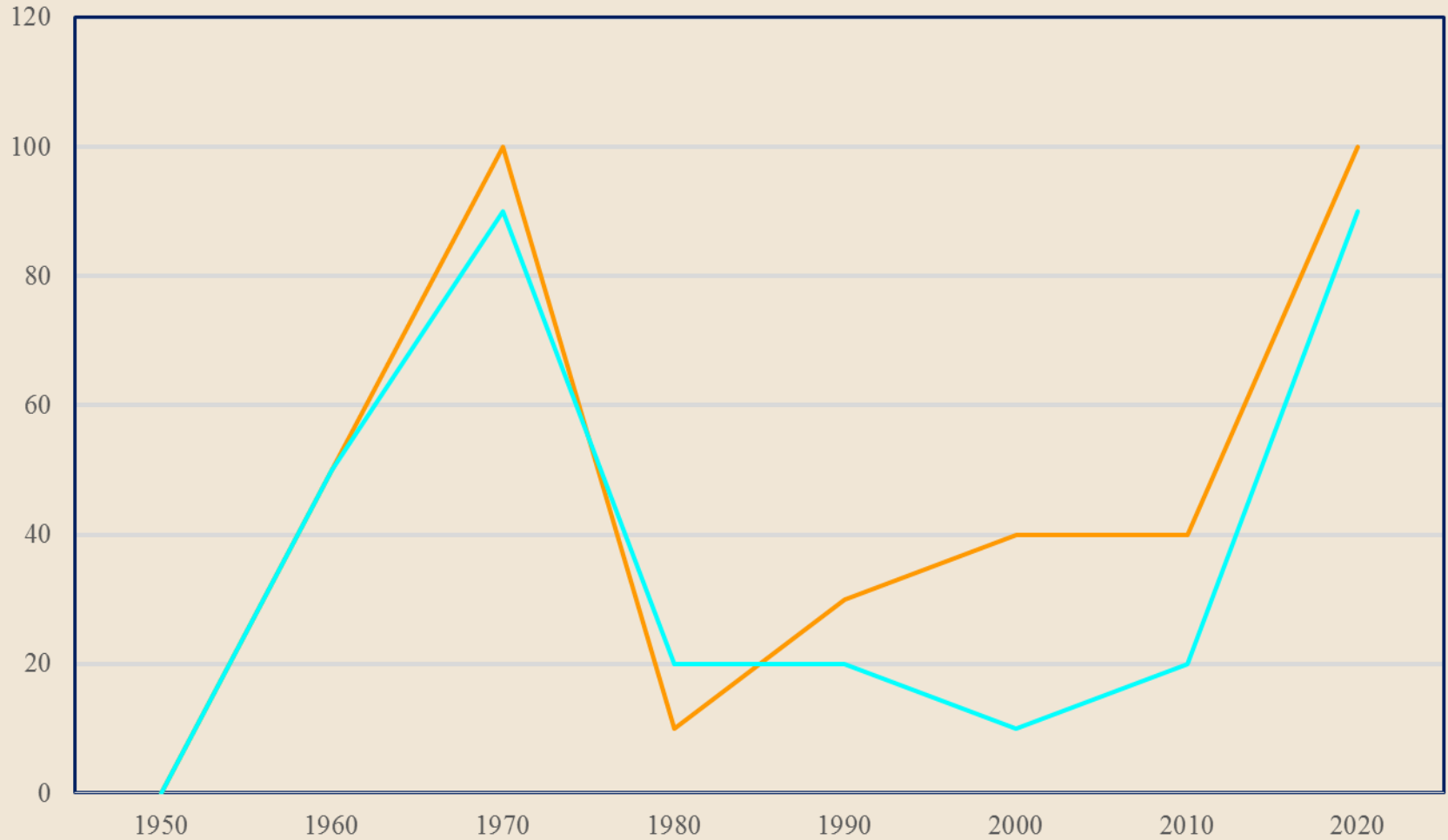
Type	Count	Relative	Trend
identity	11	15.85	↔

Stéfan Sinclair & Geoffrey
Rockwell
<http://hermeneuti.ca/>

Verläufe

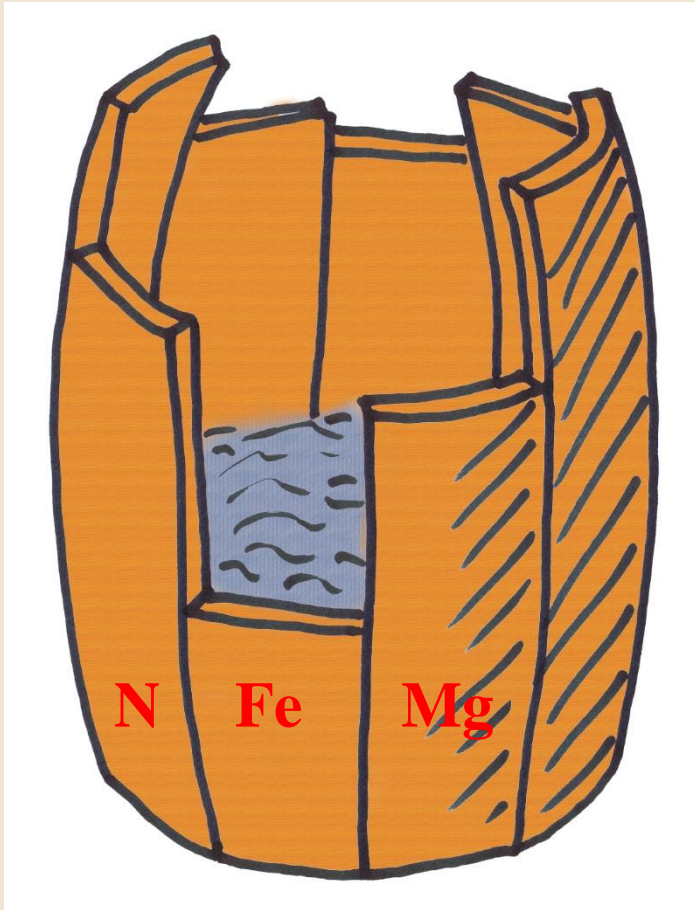
*Don't or melibria
melibria rāna rāna
ā nā a pūpū rāna
ōs fillo agnō rāna
pūpū nō ēllē pūpū*

1101100011110010010
Historisch
1010111010101011111
Kulturwissenschaftliche
11010010100100101010
Informationsverarbeitung



— Authorship analysis — Sentiment analysis

Versuch einer Verallgemeinerung

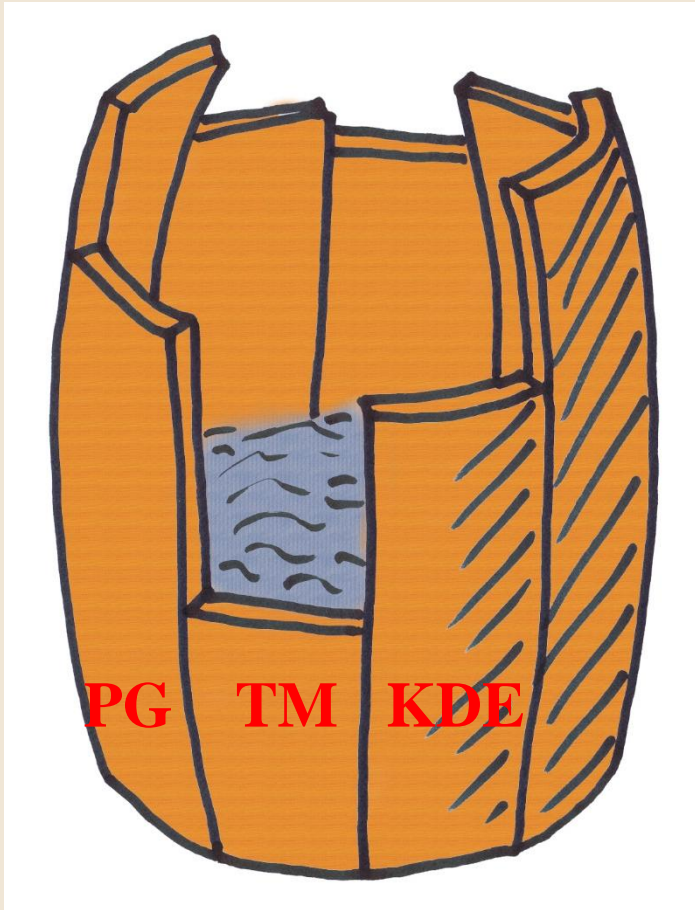


N = Stickstoff

Fe = Eisen

Mg = Magnesium

Versuch einer Verallgemeinerung



**PG =
Prozessorgeschwindigkeit**

TM = Theoretische Modelle

**KDE = Aufwand für die
Datenerfassung**

Big data ...



Eine Gasturbine von General Electric produziert an einem Tag 588 Gigabyte Daten.

Süddeutsche Zeitung, 41 / 19. Februar 2015, p. 21, Max Hägler: Transparenz ist zu wenig

Big data ...



Big data == $n * FK * VK$

FK = Filterkomplexität

VK = Verarbeitungskomplexität

Filterkomplexität I



(a)

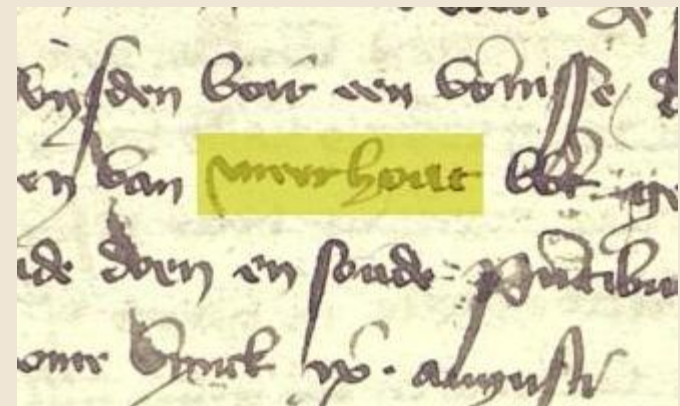
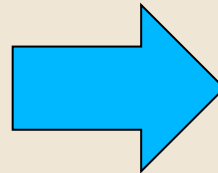
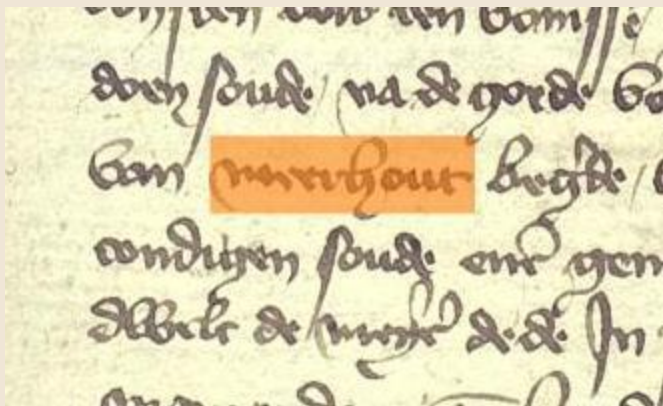


(b)

H. Mara - http://www.academia.edu/1217620/GigaMesh_and_Gilgamesh_-_3D_Multiscale_Integral_Invariant_Cuneiform_Character_Extraction

Filterkomplexität I

Meerhout in the *scheppenregister* of the city of Leuven:



L. Schomaker, MONK: <http://monk.target.rug.nl/>

Filterkomplexität I

You are transcribing page 1 of chapter 40 of **Reichsgericht** collection.

4.

Kläger zurück bezahle. Auf nach Abschluß dieser Vereinbarung

Kläger zurück bezahlen erst nach Abschluß dieser Vereinbarung habe von

gegenüber dem Beklagten...
dieser Aufforderung ist eine spezielle Leistungspflicht nicht auf
dem fact jedoch dem Kläger die über den zurückzuführen, in der

Filterkomplexität I

↓

HOME ARCHIV TRAINING TUTORIUM RESSOURCEN MyADFONTES ?

ARCHIV > [Archivaufgabe 1](#) >

Bestände und Findmittel – Aufgabe

10. Urbarium de Anno 1594



10. Urbarium de Anno 1594. Inzucht
articul. t. 2. dus n. 41. an Jan.
langb. von Jost Rodel aus n. 40. 7. 2.
Jhm. GH. 42.

11. Inzucht von Langhofen Urbarij
von 1594. N. In Abzinsung
Urbaris Jahressumme ist in 1594
Jahre von andern Titul Platten
den Jasso Jass, da ist an Jahr 1534.
Jhm. 1597. GH 42

EINGABE PRÜFEN

TIPPS ANZEIGEN TRANSKRIPTION ANZEIGEN

TUTORIUM

ZURÜCK ÜBERSICHT WEITER

Filterkomplexität II

Donit on melioth
melior rāna rāna
ā nā a pūpū rānā
sū fello rānā
op nū meli nū
pūpū nū sū rānā

11011000111110010010
1010111010101011111
1010010100100101010

Historisch
Kulturwissenschaftliche
Informationsverarbeitung



Verarbeitungskomplexität

Donne on melo...
melora rāna rāna...
ā nū a pūpū...
ōs fello...
op nūc...
pūgnā...

11011000111110010010
Historisch
10101110101010111111
Kulturwissenschaftliche
1010010100100101010
Informationsverarbeitung

NC STATE UNIVERSITY

VIRTUAL PAUL'S CROSS PROJECT

A DIGITAL RE-CREATION OF JOHN DONNE'S GUNPOWDER DAY SERMON

LONDON 1622

Overview Churchyard Acoustics Preacher Occasion Sermon Support

<http://vpcp.chass.ncsu.edu/> --- „John Donne“

Verarbeitungskomplexität

Donne an meliorum
meliorum tanta tanta
a tua a pupis quibus
siv filio agitur in
opere meo meliorum
pugnat ut esse patet

11011000111110010010
Historisch
10101110101010111111
Kulturwissenschaftliche
1010010100100101010
Informationsverarbeitung



2500 Position 7

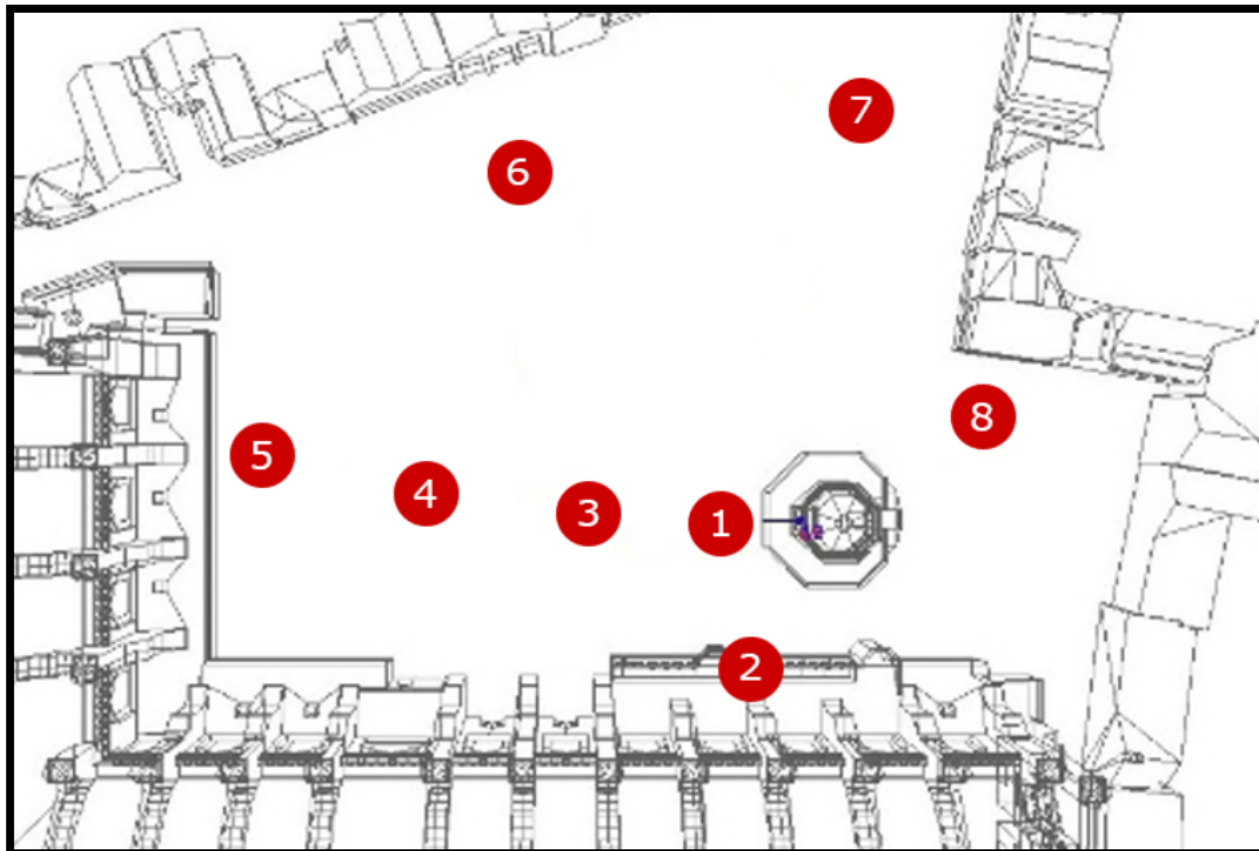
Choose a Location:

500 People	2500 people
Location 1	Location 1
Location 2	Location 2
Location 3	Location 3
Location 4	Location 4
Location 5	Location 5
Location 6	Location 6
Location 7	Location 7
Location 8	Location 8

1200 people	5000 people
Location 1	Location 1
Location 2	Location 2
Location 3	Location 3
Location 4	Location 4
Location 5	Location 5
Location 6	Location 6
Location 7	Location 7
Location 8	Location 8

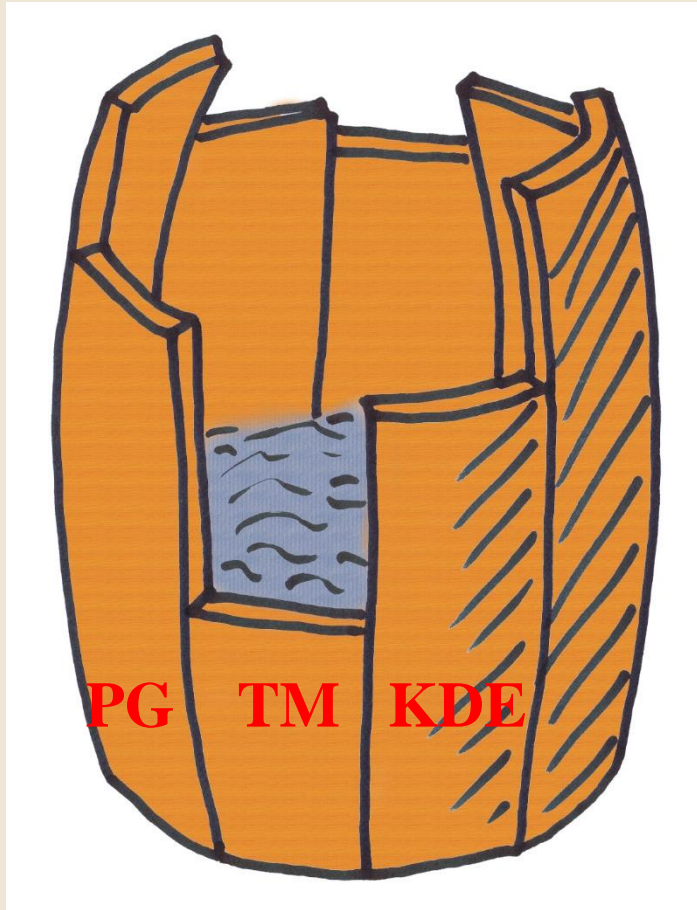
<http://vpcp.chass.ncsu.edu/> --- „John Donne“

Verarbeitungskomplexität



Choose a Location:	
500 People	2500 people
Location 1	Location 1
Location 2	Location 2
Location 3	Location 3
Location 4	Location 4
Location 5	Location 5
Location 6	Location 6
Location 7	Location 7
Location 8	Location 8
1200 people	5000 people
Location 1	Location 1
Location 2	Location 2
Location 3	Location 3
Location 4	Location 4
Location 5	Location 5
Location 6	Location 6
Location 7	Location 7
Location 8	Location 8

Versuch einer Verallgemeinerung



**PG =
Prozessorgeschwindigkeit**

TM = Theoretische Modelle

**KDE = Aufwand für die
Datenerfassung**

Big data ...

Big data without big theory is big bullshit!

„Sie hören ja nichts anderes mehr: *big data*. Aber big data without big theory is big bullshit!“ Ohne Theorie, ohne Verständnis des Menschen und seiner Psyche komme man nicht weiter. Management sei eine *liberal art*, eine Geisteswissenschaft.

Süddeutsche Zeitung, 261 / 13. November 2014, p. 19, „Nahaufnahme“:
Interview mit Richard Straub, ehemaliger Chief Learning Officer (Leiter der Personalentwicklung) bei IBM.

Big data ...

What Do You Do with a Million Books?

Greg Crane, 2006

Ja, abgesehen von den 1000, die klassische Texte enthalten, was eigentlich?

Was um alles in der Welt machen wir dann mit 20.000 Archivkilometern?

Eine Definition von eScience

Naturwissenschaft in einer Arbeitsumgebung, in der

- (1) der Zugang zur für eine Forschungsfrage benötigten Information,
- (2) die Analyse dieser Information
- (3) und die Publikation der gewonnenen Ergebnisse

gleichermaßen gut durch die (verteilte) Informationstechnologie unterstützt wird.

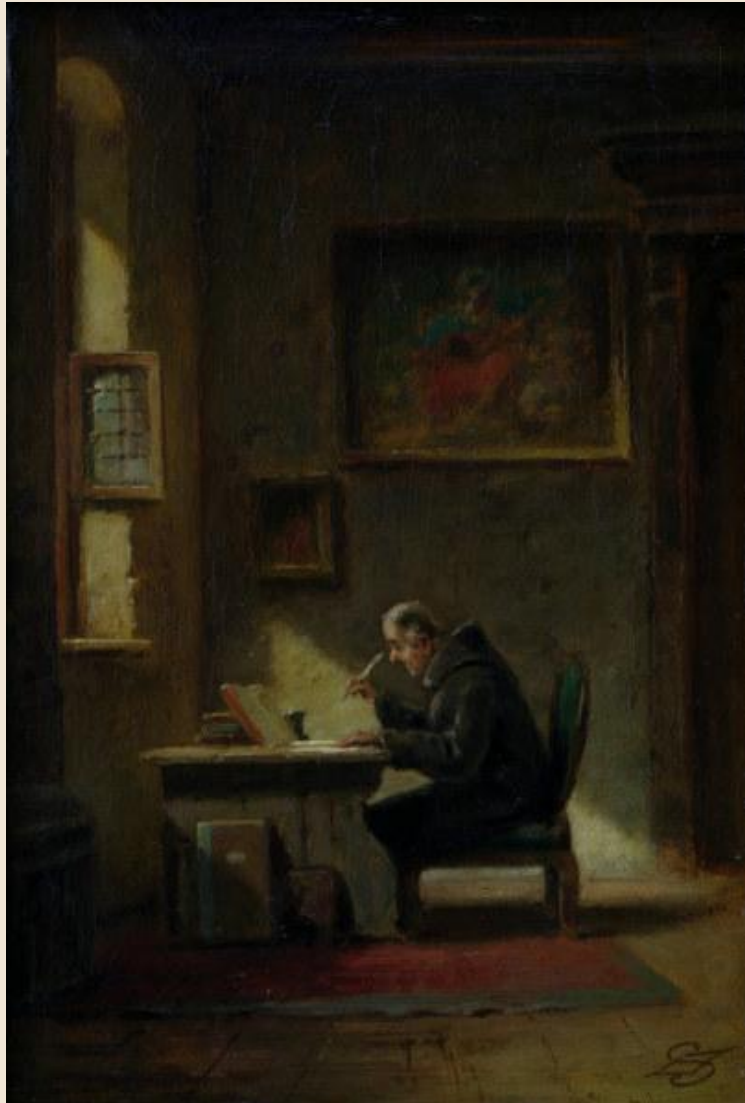
Eine Definition von eHumanities

Geisteswissenschaft in einer Arbeitsumgebung, in der

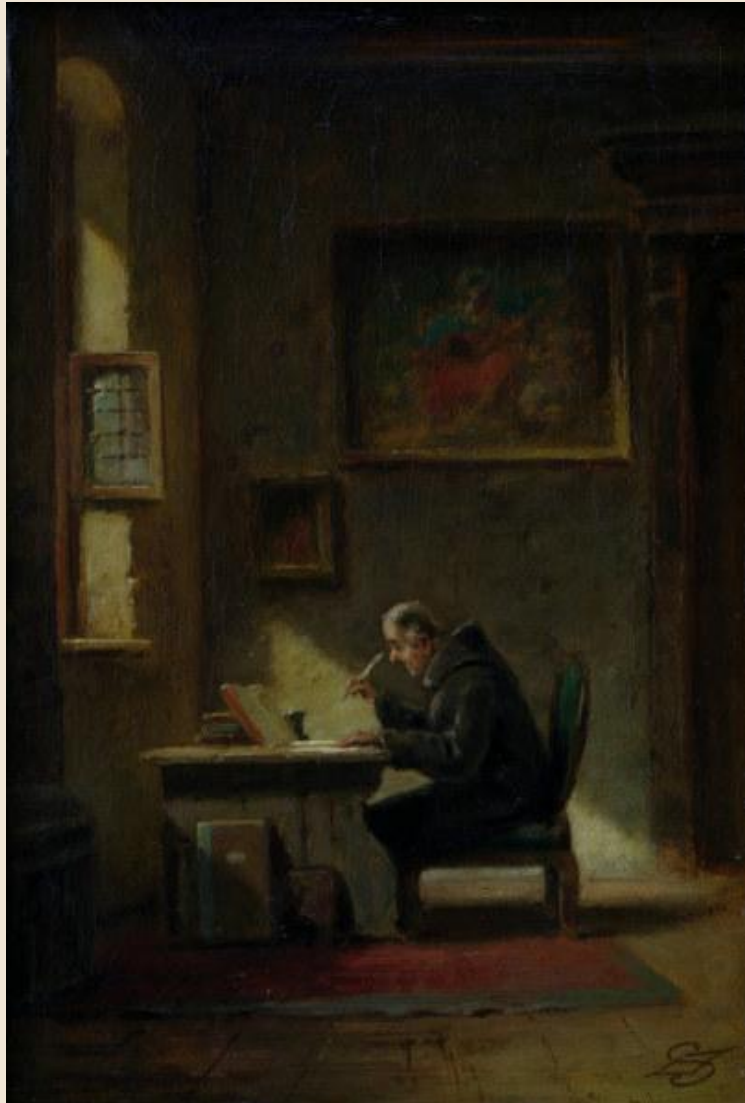
- (1) der Zugang zur für eine Forschungsfrage benötigten Information,
- (2) die Analyse dieser Information
- (3) und die Publikation der gewonnenen Ergebnisse

gleichermaßen gut durch die (verteilte) Informationstechnologie unterstützt wird.

Wenn die Quellen überfließen ...



Wenn die Quellen überfließen ...



Kritik

Konzentration auf Infrastruktur verhindert Diskussion, was in den Geisteswissenschaften mit IT erreicht werden soll.

(W. McCarty: <http://digitalhumanities.org/humanist/>)

Mechanistische Übernahme aus den Hard Sciences; zu starke Konzentration auf Werkzeuge trivialisiert die methodisch konzeptuelle Diskussion.

(P. Svensson

<http://digitalhumanities.org/dhq/vol/5/1/000090/000090.html> im Kontext von: [.../vol/3/3/000065/000065.html](http://digitalhumanities.org/dhq/vol/3/3/000065/000065.html) und [.../vol/4/1/000080/000080.html](http://digitalhumanities.org/dhq/vol/4/1/000080/000080.html))

Kritik & Angst

The digital humanities are at a rhetorical and institutional crossroads, having just enough critical mass of public attention to warrant a torrent of speculation.

Shawna Ross, “In Praise of Overstating the Case: A review of Franco Moretti, *Distant Reading*”, in: DDQ 8/1 (2014)

<http://www.digitalhumanities.org/dhq/vol/8/1/000171/000171.html>

Such worries have highlighted the conflict rhetoric characterizing DH meta-discourse — with its familiar debates about distant versus close reading, about tool builders versus tool users, about rejecting versus incorporating theory, about using DH skills to be helpful around the office versus being more protective of our time. What unites these debates is the common tension between, on the one hand, desiring the fertility and excitement created by the

Kritik & Angst

Such worries have highlighted the conflict rhetoric characterizing DH meta-discourse — with its familiar debates about distant versus close reading, about tool builders versus tool users, about rejecting versus incorporating theory, about using DH skills to be helpful around the office versus being more protective of our time.

... on the one hand, desiring the fertility and excitement created by the coexistence of a diversity of approaches and interactions among very differently-minded scholars and, on the other, fearing that collaboration-related compromises could devalue DH work, subsequently re-marginalizing DH workers.

Shawna Ross, “In Praise of Overstating the Case: A review of Franco Moretti, *Distant Reading*”, in: DDQ 8/1 (2014)

<http://www.digitalhumanities.org/dhq/vol/8/1/000171/000171.html>

Häresie 3

Mir ist völlig egal, was die „Digital Humanities“ sind.

Mich interessiert, welche epistemischen Möglichkeiten sich aus der Anwendung von Algorithmen auf geisteswissenschaftliche Quellen ergeben.

*Nennt das Ding
doch „Isabella“.*



Häresie 4

Jede geisteswissenschaftliche Disziplin gehört zu den Digital Humanities.

Digital Humanities sind jedoch ein intellektuelles Abenteuer und ein anspruchsvolles Unterfangen.

Die korrekte Bedienung eines Smartphones qualifiziert nicht als Digital Humanist.



Thallers Ungleichheiten

Die Digital Humanities sind nicht identisch

- (1) Mit digitalen Bibliotheken.
- (2) Mit der Korpuslinguistik.
- (3) Mit dem Textmining.
- (4) ...

Sie *nutzen* diese für epistemische Erweiterungen der Geisteswissenschaften souverän.

PS: Souverän genug um keine Angst vor der Informatik zu haben.

Häresie 5

Die modernen Informationssysteme haben die Bedingungen für die Geisteswissenschaften so gründlich verändert, dass wir den Mut haben müssen, die wissenschaftliche Praxis zu ändern.

Dazu brauchen wir eine ernsthaft geführte epistemische Debatte, in voller Kenntnis der Konzepte der Digital Humanities.

PS: Und keine darüber, wer im Sandkasten mitspielen darf.

Spitzwegs Welt

*D*om̄t̄ oñ̄ meliōt̄ia
melior̄ tān̄cā tān̄cō
ā nū̄ a p̄p̄uō cūb̄i
ōis filio aḡḡr̄is
op̄ nū̄ melī nōp̄
p̄p̄uō nō ell̄ p̄p̄uō

11011000111110010010
1010111010101011111
1010010100100101010

Historisch
Kulturwissenschaftliche
Informationsverarbeitung



Spitzwegs Welt

*D*omit om̄ melioſa
meliora cāra cāra
ā nū a pūna cūbi
sū fillo agna
op nū meli nū
pūna nū ell pū

11011000111110010010
1010111010101011111
1010010100100101010

Historisch
Kulturwissenschaftliche
Informationsverarbeitung



Unsere Welt

*Quia cum melioribus
meliora cetera ceteris
a tua a purius cibus
siv filio affinis in
op me meli non nulli
pugnata ut esse pat*

11001100011110010010
Historisch
1010111010101011111
Kulturwissenschaftliche
1010010100100101010
Informationsverarbeitung



Unsere Welt

Dies ist die Welt der Informationssysteme.

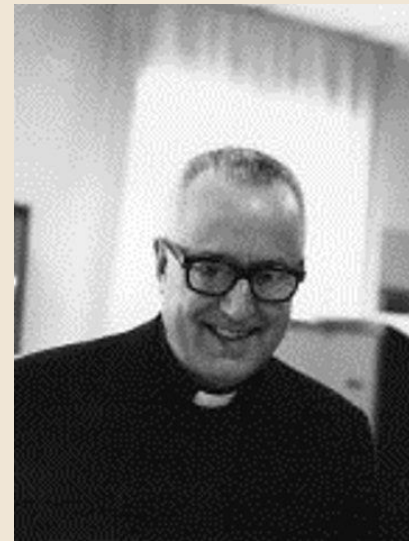
Sie wartet nicht auf uns.

Wenn die Geisteswissenschaften die IT nicht nutzen, tun es eben andere.

Bitte, bitte: Mut zur Fantasie ...



E. E. „Doc“ Smith, 1890 - 1966



Padre Roberto Busa, 1913 -2011

Danke!

manfred.thaller@uni-koeln.de