

Exploring and analyzing linguistic variation

Elke Teich
Saarbrücken
Germany

Linguistic variation

Methods

text-based vs. variationist
data-based vs. data-driven
macro vs. micro-analysis

Description

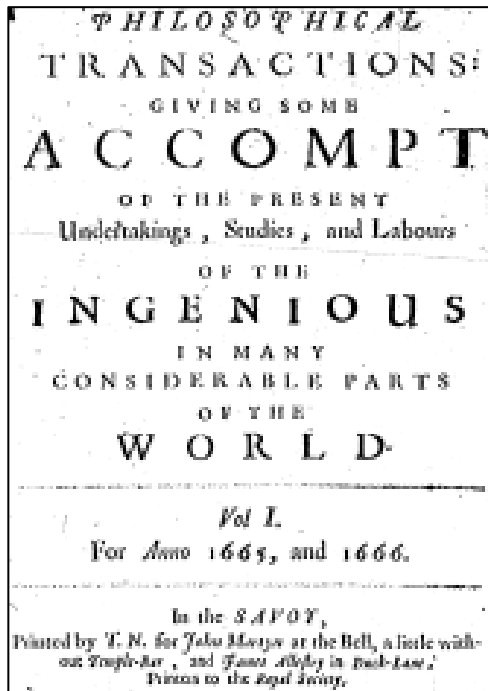
variables: time, age,
gender, context of
situation/register
linguistic levels:
phonetic, lexical,
syntactic, semantic

Scientific English
(1650-1850)

Models

entrenchment
innovation
diffusion
typicality
productivity
...

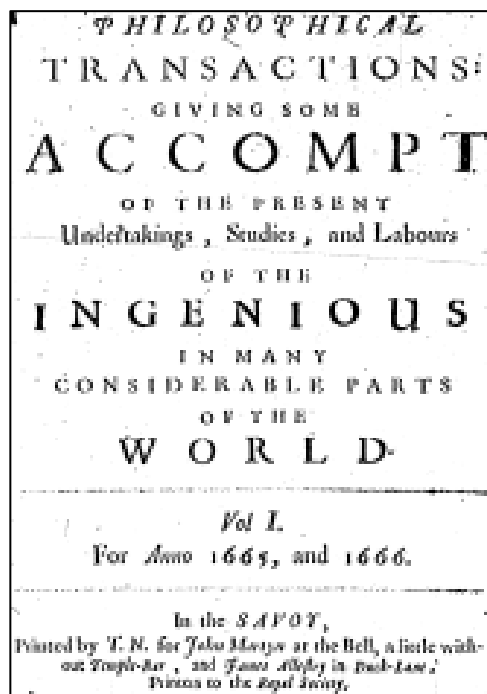
Historical example: Royal Society (article)



This I did with much solicitude further inquire into; whereupon I found not only one hollowness, but as often as I cut the Nerve asunder, the hollowness still continued therein, and I found in some places not only one cavity, but two or three cavities at once;

Coxe, Daniel. 1674. "A continuation of Dr. Daniel Coxe's Discourse, Touching the Identity of All Volatil Salts, and Vinous Spirits; Together with Two Surprising Experiments Concerning Vegetable Salts, Perfectly Resembling the Shape of the Plants, Whence They Had Been Obtained". *Philosophical Transactions (1665-1678)* 9. The Royal Society: 169–82.

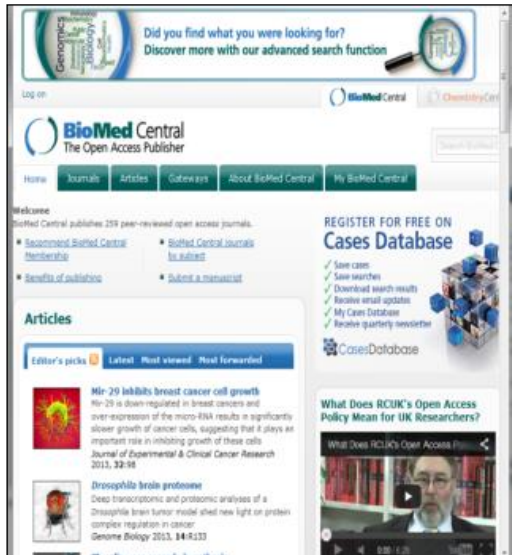
Historical example: Royal Society (article)



This I did with much solicitude further inquire into; whereupon I found not only one hollowness, but as often as I cut the Nerve asunder, the hollowness still continued therein, and I found in some places not only one cavity, but two or three cavities at once;

Coxe, Daniel. 1674. "A continuation of Dr. Daniel Coxe's Discourse, Touching the Identity of All Volatil Salts, and Vinous Spirits; Together with Two Surprizing Experiments Concerning Vegetable Salts, Perfectly Resembling the Shape of the Plants, Whence They Had Been Obtained". *Philosophical Transactions (1665-1678)* 9. The Royal Society: 169–82.

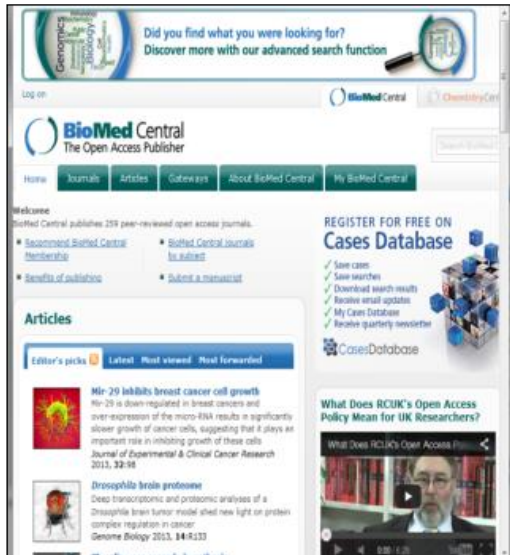
Contemporary example: Biomed Central (abstract)



We report the discovery of a novel downstream target of BCR-ABL signalling, PRL-3 (PTP4A3), an oncogenic tyrosine phosphatase. Analysis of CML cancer cell lines and CML patient samples reveals the upregulation of PRL-3.

<http://www.biomedcentral.com/>

Contemporary example: Biomed Central (abstract)



We report the discovery of a novel downstream target of BCR-ABL signalling, PRL-3 (PTP4A3), an oncogenic tyrosine phosphatase. Analysis of CML cancer cell lines and CML patient samples reveals the upregulation of PRL-3.

<http://www.biomedcentral.com/>

Linguistic development of scientific discourse

Assumptions

diversification, specialization

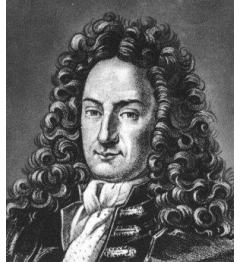
→ *denser encoding of information*

use of compact/reduced linguistic forms
(e.g. compounds, reduced relative clauses)

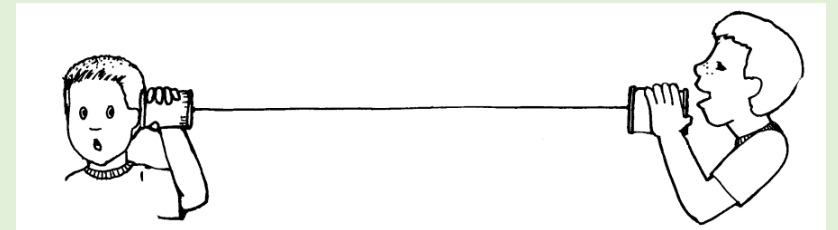
professionalization, institutionalization

→ *conventionalization*

use of fairly fixed vocabularies
(terminology, formulaic expressions)



optimal code



for communication



(cf. Aylett & Turk 2004; Jaeger 2010; Piantidosi, Tily & Gibson 2011)

- choice of a particular linguistic encoding depends on (predictability in) context

Information
Density (ID)



Surprisal

$$-\log_2 P(\textit{unit} \mid \text{Context})$$

- contextually determined predictability is appropriately indexed by Shannon's notion of information

Which features are involved in diachronic change in scientific writing?
How can ID/surprisal help capture this change?

Feature detection: typicality

- Discover representative and distinctive features

relative
entropy

Analyses

- grammatical typicality

Feature inspection: productivity

- Observe ambient context

average
surprisal

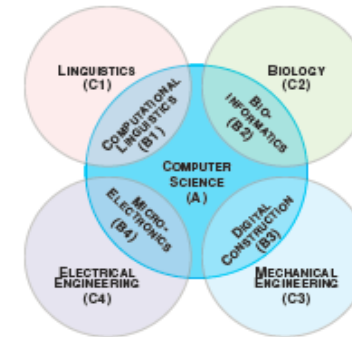
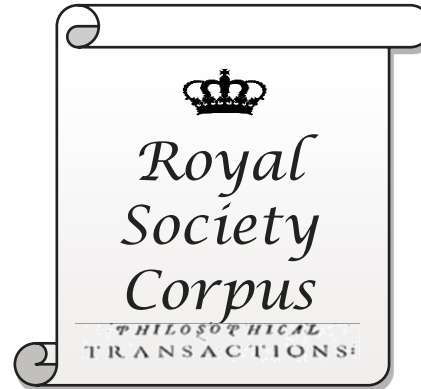
- lexical productivity

Data

SCIENTIFIC



1506–1700



1970s & 2000s

EME (1500–1700)

LME (1700–1900)

ME (1900–present)

MIXED

HELSINKI CORPUS
1500–1700



1641–1700

ARCHER
1600–1999



1640–1740

CLMET

1700-1920

BROWN FAMILY

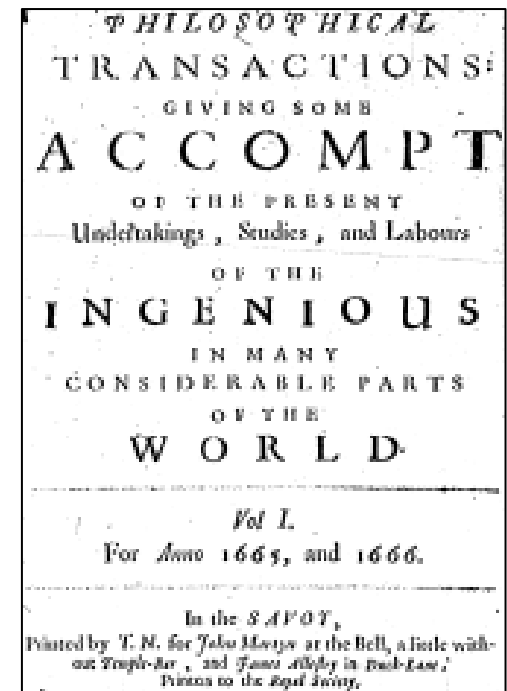
1931, 1961, 1991



Royal Society Corpus (RSC)

book reviews
articles
miscellaneous
obituaries

Journal\Type	brv	fla	mis	nws	total
Philosophical Transactions (1665-1678)	124	641	154	–	919
Philosophical Transactions (1683-1775)	154	3903	338	–	4395
Philosophical Transactions of the Royal Society of London	–	2531	283	–	2814
Abstracts of Papers Printed ...	–	1316	15	–	1331
Abstracts of Papers Communicated ...	–	429	5	–	434
Proceedings of the Royal Society of London	–	1476	38	14	1528
total	278	10296	833	14	11421



size: ca. 35 million tokens, source: XML (JSTOR)

1-, 10-, 50-year time periods

(Kermes et al. 2016)

available from:



Feature detection

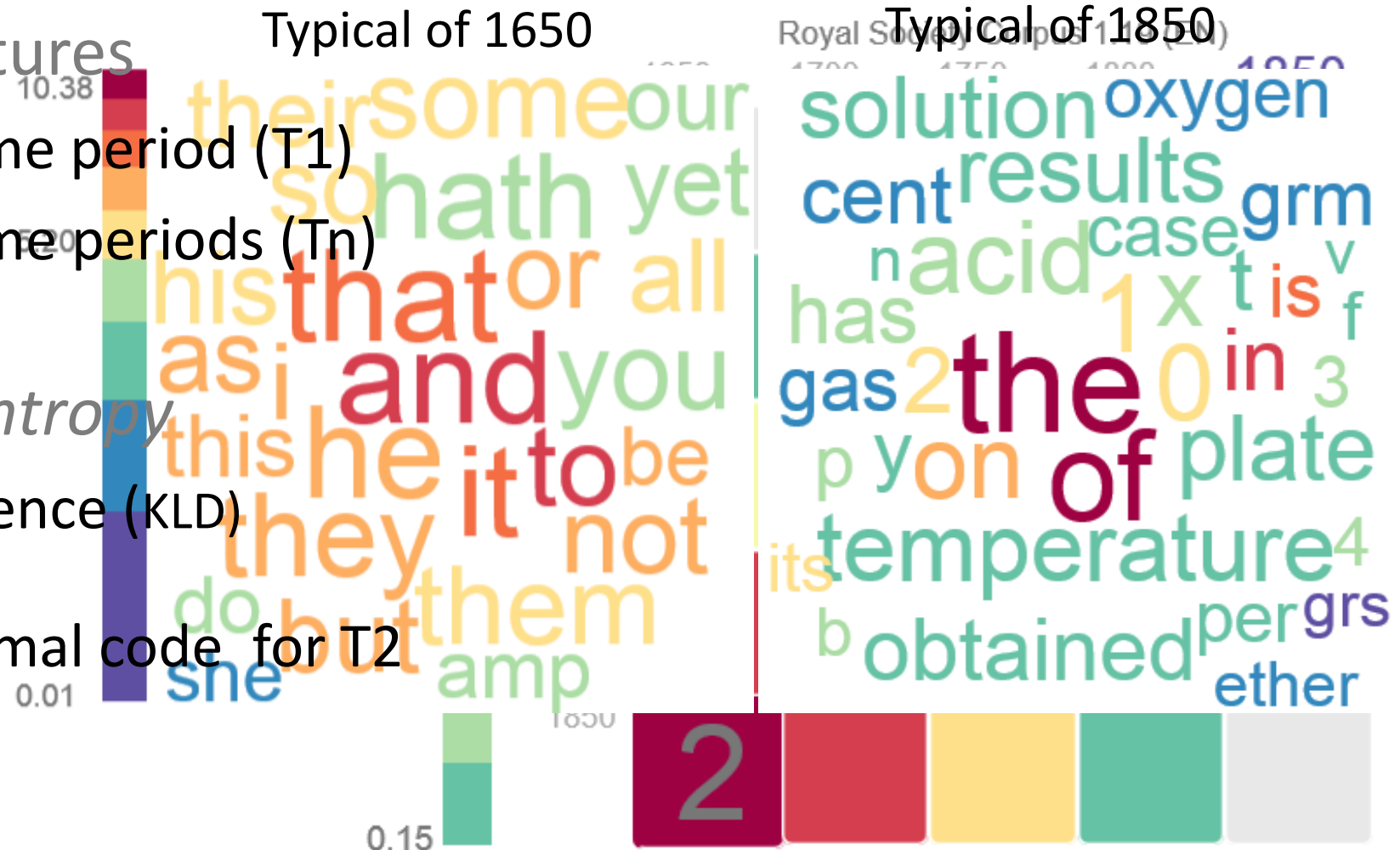
Typical linguistic features

- *representative* of a time period (T1)
- *distinctive* to other time periods (Tn)

Measure: *Relative Entropy*
Kullback-Leibler Divergence (KLD)

→ encode T1 with optimal code for T2

(Fankhauser et al. 2014)



Feature inspection

Typical features in context

In a given context, a unit has
low predictability → high surprise
high predictability → low surprise

Measure: *Average Surprisal (AvS)*

$$AvS(unit) = \frac{1}{|unit|} \sum_i -\log_2 p(unit|context_i)$$

(Genzel & Charniak 2002; Degaetano-Ortlieb et al. to appear)

*I Cannot enough wonder
at the strange agreement
of the thoughts of that acute
French Gentleman, Monsieur Auzont,
in the Hypothesis of the Comets
motion, with mine;*

Analyses

Typical lexico-grammatical patterns over time

A1: Grammatical typicality

A2: Lexical productivity

A1: Approach

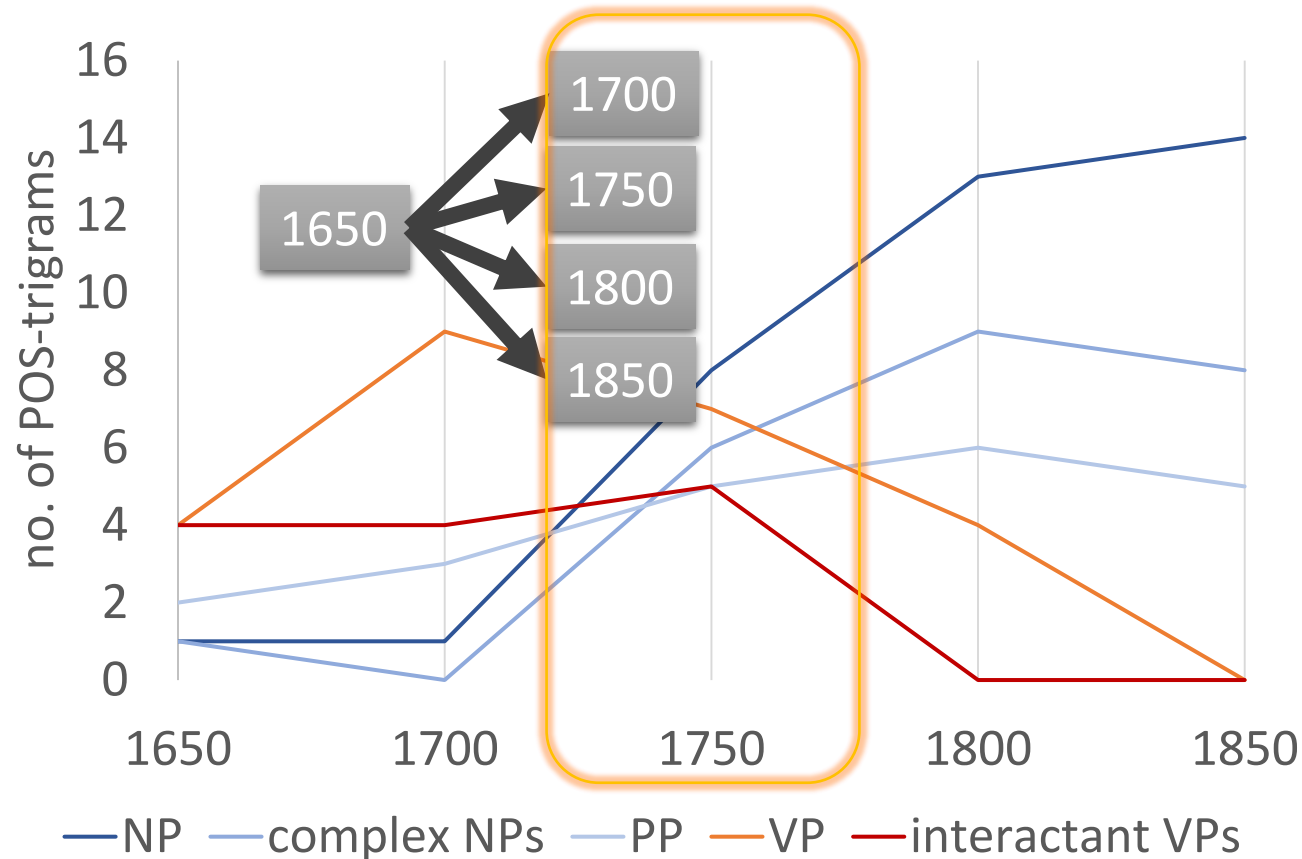
Approximate grammatical patterns with POS-trigrams

$$D_{KL}(T1||T2) = - \sum_i p(\text{trigram}_i|T1) \log_2 \frac{p(\text{trigram}_i|T1)}{p(\text{trigram}_i|T2)}$$

- Extract all POS-trigrams from corpus (e.g. Det-Adj-N such as *a few minutes*)
- Exclude sentence markers, symbols and foreign words
- Group into phrase types (NP, VP etc)

A1: Grammatical typicality

Comparison Typical phrase types over time



verbal (interactional)
→
nominal (informational)

A1: Grammatical typicality

1750s – a period of transition

1750 vs.	phrase type	example
1650/1700	NP	<i>the freezing point, the same time</i>
	NP (complex)	<i>the effects of, degree of heat</i>
	PP	<i>of nitrous air, an account of</i>
	AdjP (complex)	<i>small quantity of, great number of</i>
<hr/>		
1800/1850	VP	<i>repeated the experiment, found the sum</i>
	VP (interact.)	<i>I then took, I did not</i>
	VP (modal)	<i>could not perceive/find</i>
	NP (complex)	<i>part of it, end of it</i>
	Cl (interact.)	<i>as I found</i>

new: informational

- nominal patterns
- *specialized* terminology

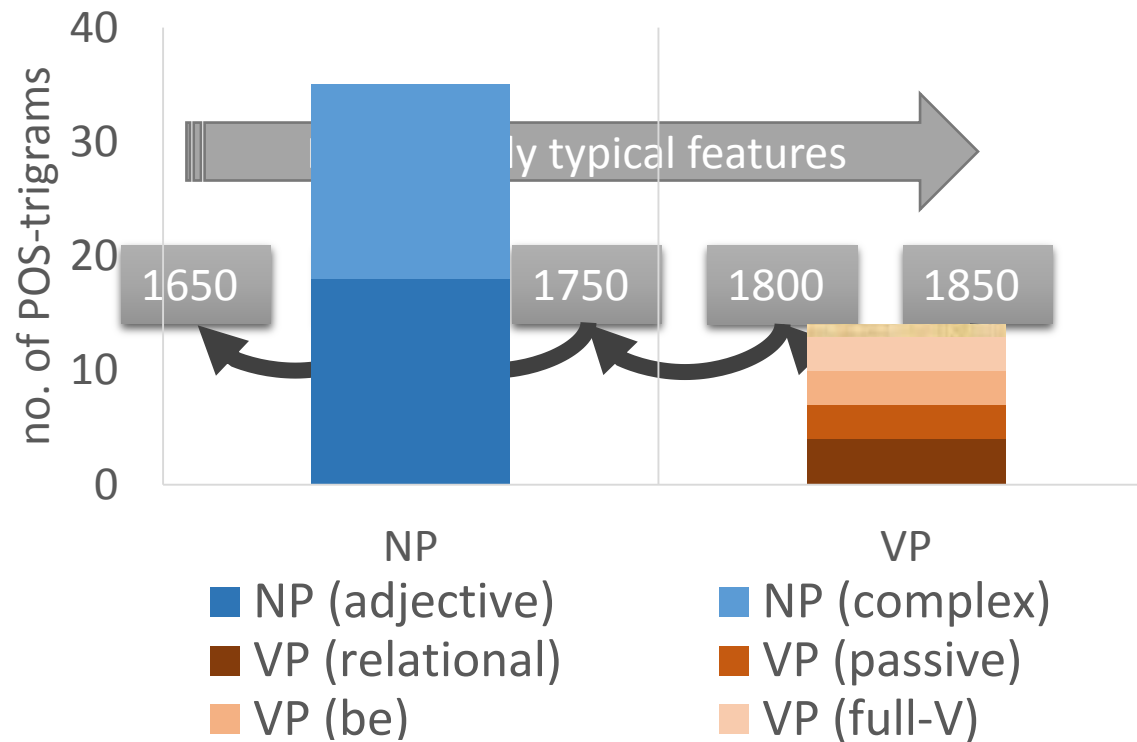
old: interactional

- verbal patterns (interactant, past tense, modals), *author-as-agent*
- nominal patterns referring to *general concepts*

(1751: Royal Society starts review process → *professionalization*)

A1: Grammatical typicality

Comparing typical phrase types



NPs related to terminology

NP adjective (*carbonic/muriatic acid gas*)

NP complex (*center of gravity*)

VPs related to “scientific” style

VP relational (*star is blue, r is odd*)

VP passive (*effect is produced, account is given*)

VP full-V (*the author concludes/states/considers*)

VP red-rel (*effect produced by, line drawn from*)

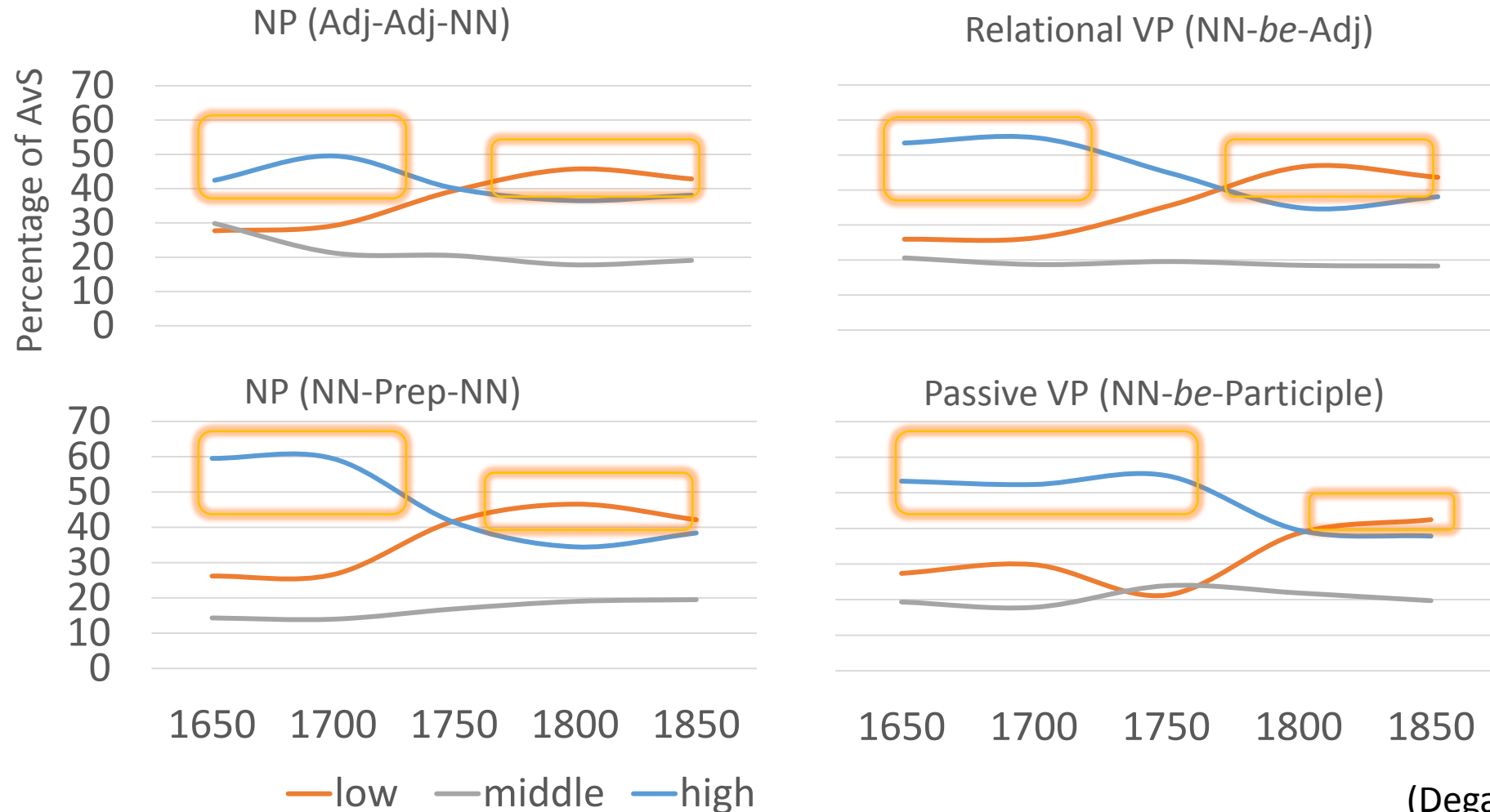
A2: Lexical productivity

- Inspect increasingly typical POS-trigrams in ambient context
- Do they become more or less productive?

High AvS values indicate higher productivity → type variation
i.e. grammatical patterns attract new lexical items (and spread to new contexts)

Low AvS values indicate lower productivity → conventionalized use
i.e. grammatical patterns are used in same contexts

A2: Lexical productivity



From 1650 to 1750:
quite unpredictable
→ high variation

From 1750 onwards:
more predictable
→ conventionalized use

(Degaetano-Ortlieb & Teich 2016)

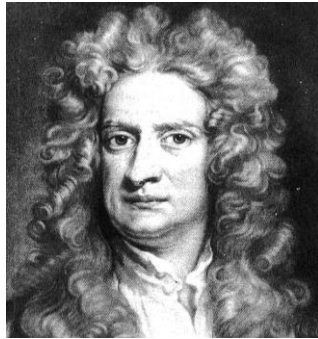
A2: Lexical productivity

Lexical realizations of Adj-Adj-NN

period	examples	freq. (pM)
1650	<i>dark brown colour</i>	7 (2.70)
	<i>next foregoing tract</i>	6 (2.32)
	<i>cold fair weather</i>	4 (1.55)
1750	<i>obedient/obliged humble servant</i>	135 (23.25)
	<i>heavy/light inflammable air</i>	110 (17.47)
	<i>diluted vitriolic acid</i>	29 (6.96)
1850	<i>concentrated sulphuric acid</i>	104 (8.93)
	<i>carbonic acid gas</i>	64 (5.50)
	<i>complete differential coefficient</i>	49 (4.21)

- from general to specific concepts
- towards terminological patterns that increase in frequency

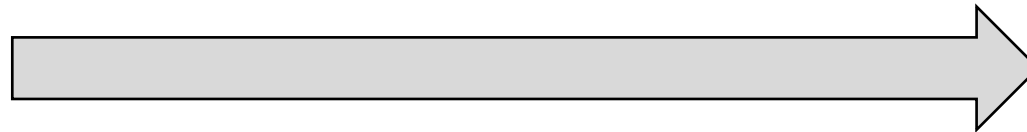
Results (so far)



1750



*interactant
verbal patterns*

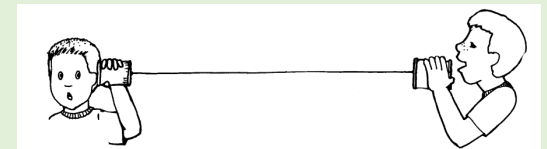


*complex nominal patterns
relational verbal patterns*

→ increasing use of denser forms of encoding

→ conventionalization

optimal code



for communication

Approach (so far)

- Relative entropy
 - good for comparing whole corpora
 - reveals typical features (beyond simple frequency)
 - indicator of degree of difference
- Average Surprisal
 - good for inspecting ambient context
 - reveals differences in ID of alternative encodings
 - indicator of productivity

Macro-
Analysis:
Contrast
(Linguistics)

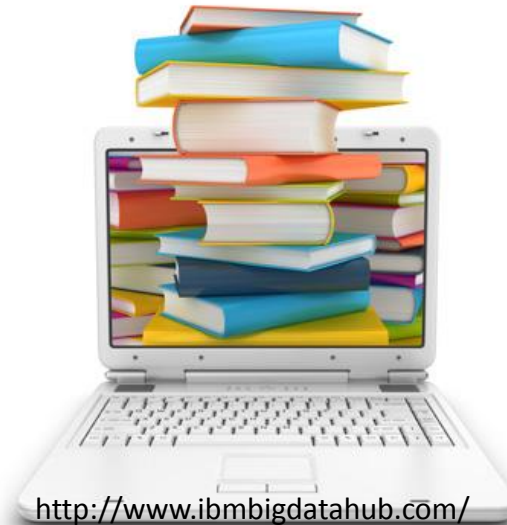
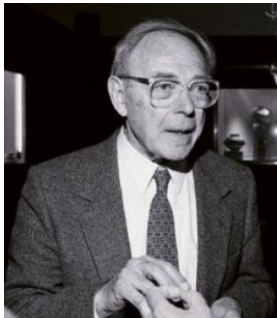
Micro-
Analysis:
Context
(Philology)

A variationist's (*my*) dream come true



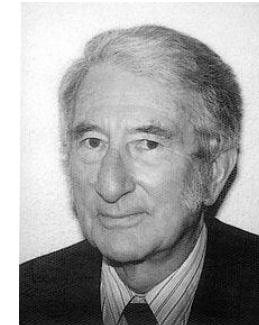
Linguistics (Science)

- valid generalizations
- modeling linguistic processes
- (psychology)



Linguistics (Philology)

- textual scrutiny
- modeling linguistic experience in context (culture, time)
- (sociology)

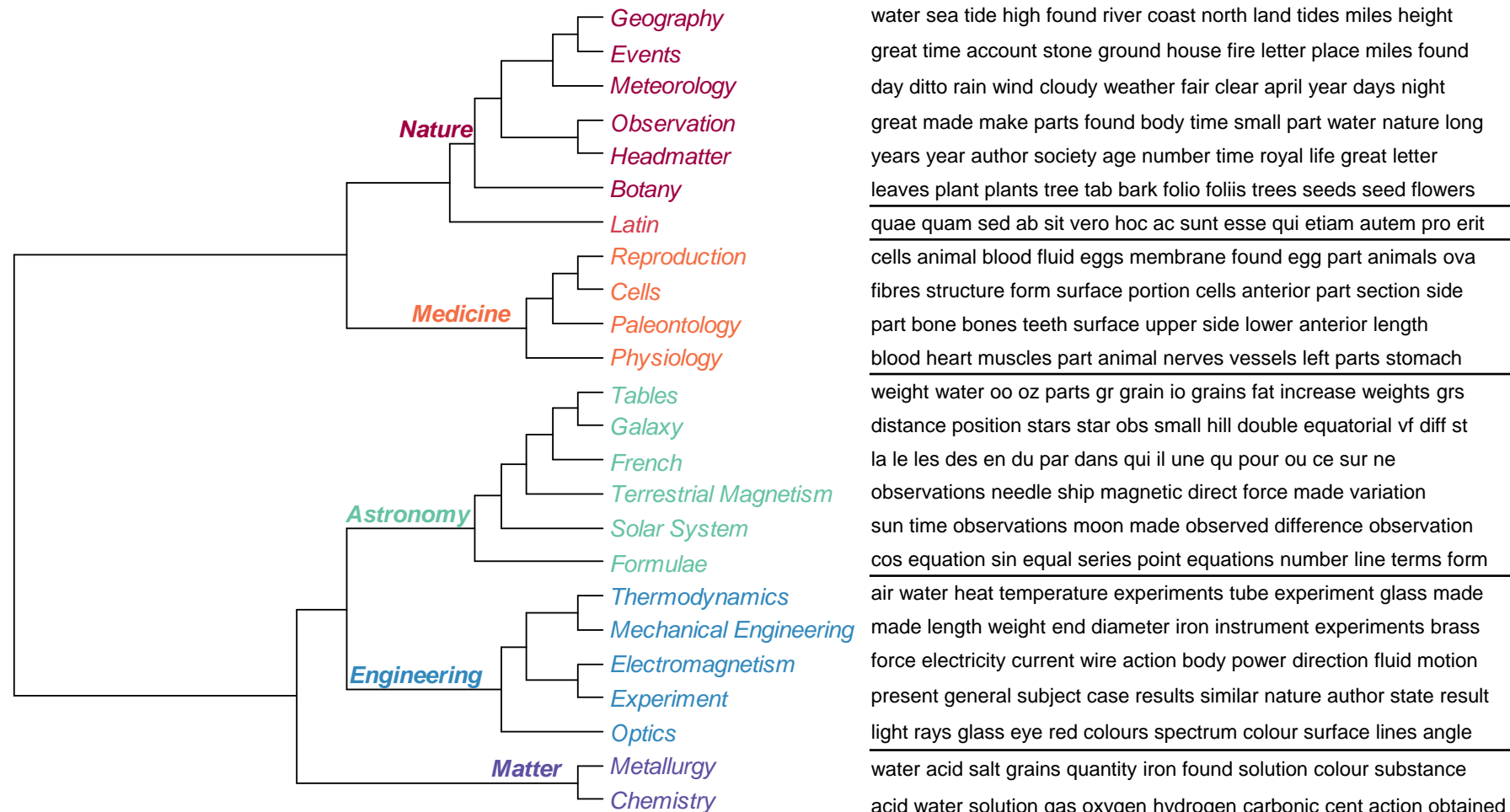


Informatics

- modeling information processes
- robust, reliable systems
- (engineering)



Current & future work

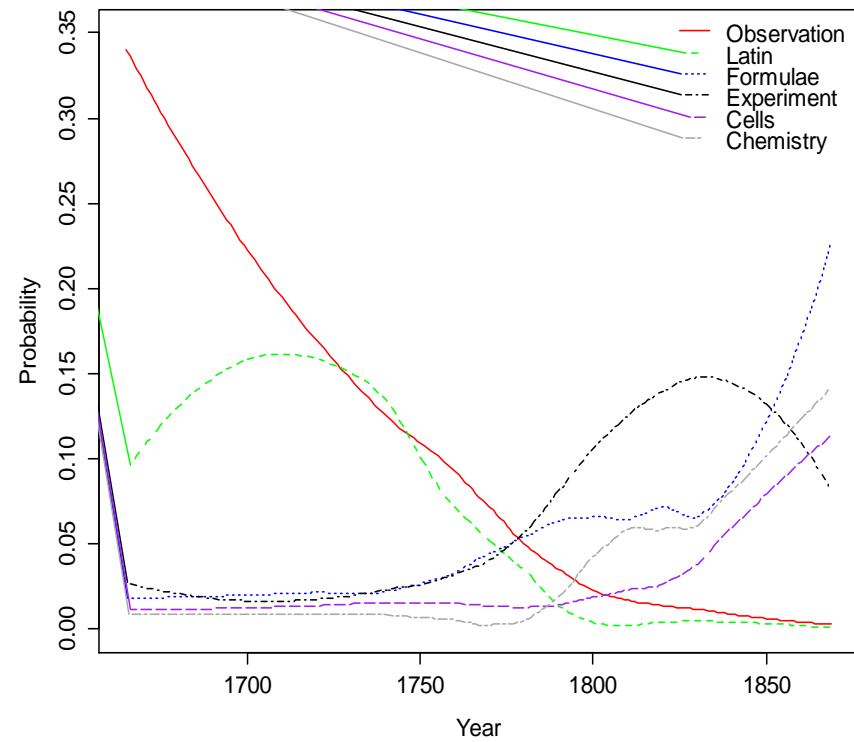


(Fankhauser et al. 2016)

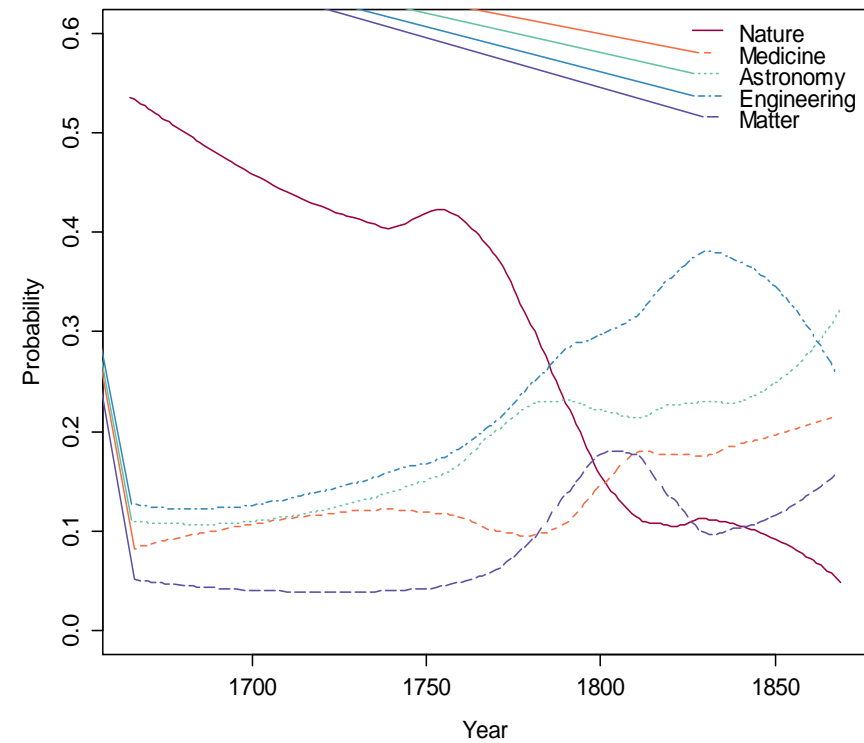
- Factorize Document-Word distribution $P(w_i|d)$ into Topic-Word distributions $P(w_i|z_k)$ and Document-Topic distributions $P(z_k|d)$

$$P(w_i|d) = \sum_k P(w_i|z_k)P(z_k|d).$$

- Dimensionality Reduction: Represent documents by means of 20-100 topics instead of > 100.000 types (different words)
- Here: 24 Topics on documents with stopword exclusion



Selected Topics



Topic Groups

- Observation

- First half: Few, dominant topic (groups) – skewed distribution
- Second half: many topic (groups) – even distribution

- Measure for topical diversification

- Shannon-Entropy on topic distributions per year (*ent*)

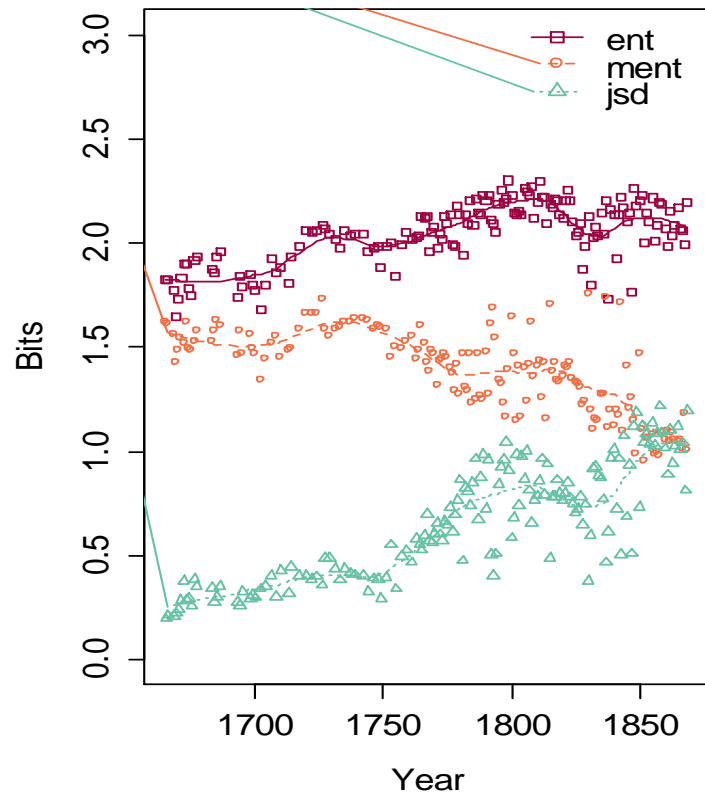
$$H(P_y) = - \sum_k P(z_k|y) \log_2 P(z_k|y)$$

- Skewed -> low entropy, Even -> high entropy

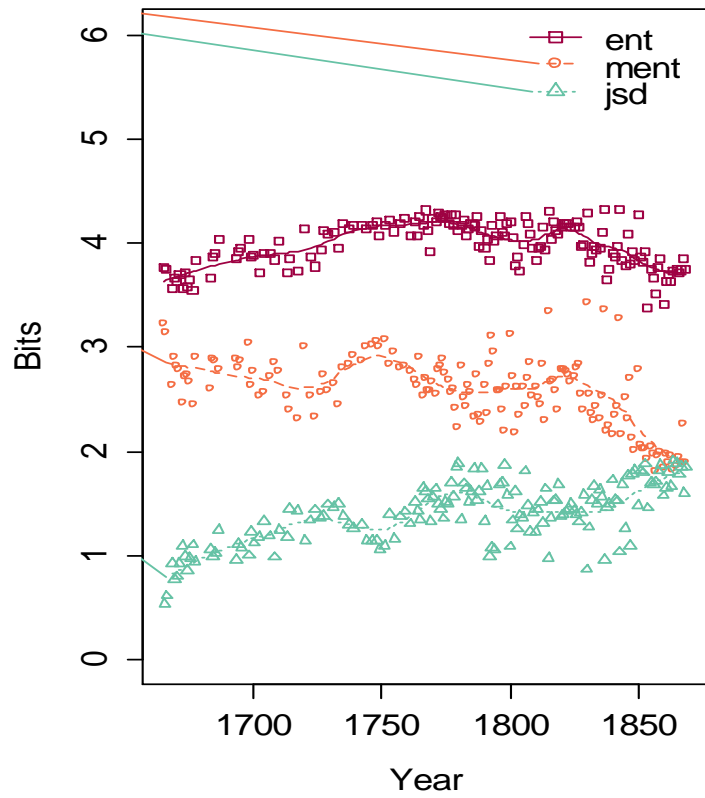
- Measure for topical specialization

- Mean Shannon-Entropy of *individual* document-topic distributions per year (*ment*)

$$H_{mean}(P_y) := 1/ny \sum_{d_j \in y} H(P_{d_i})$$



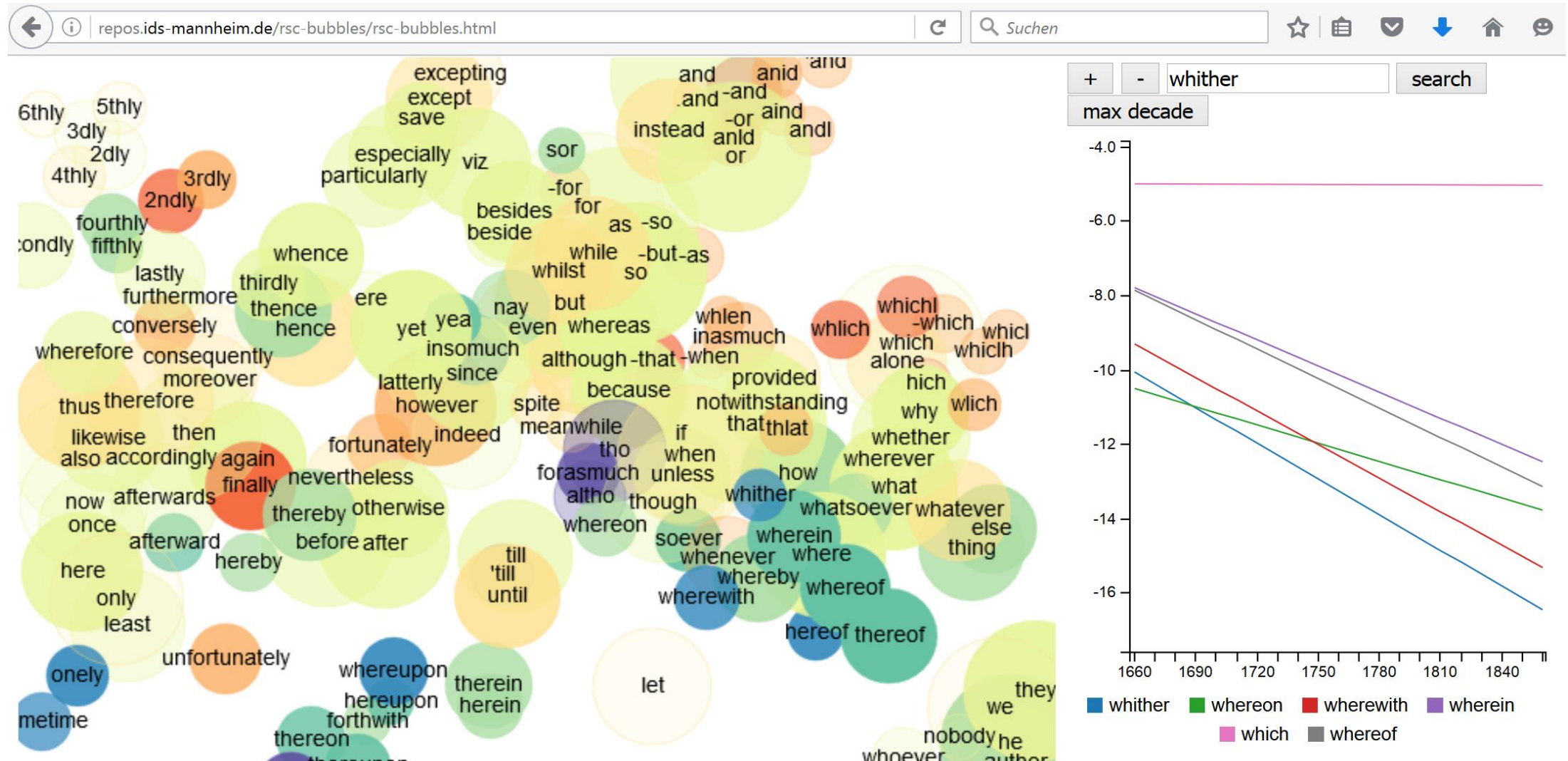
Topic Groups



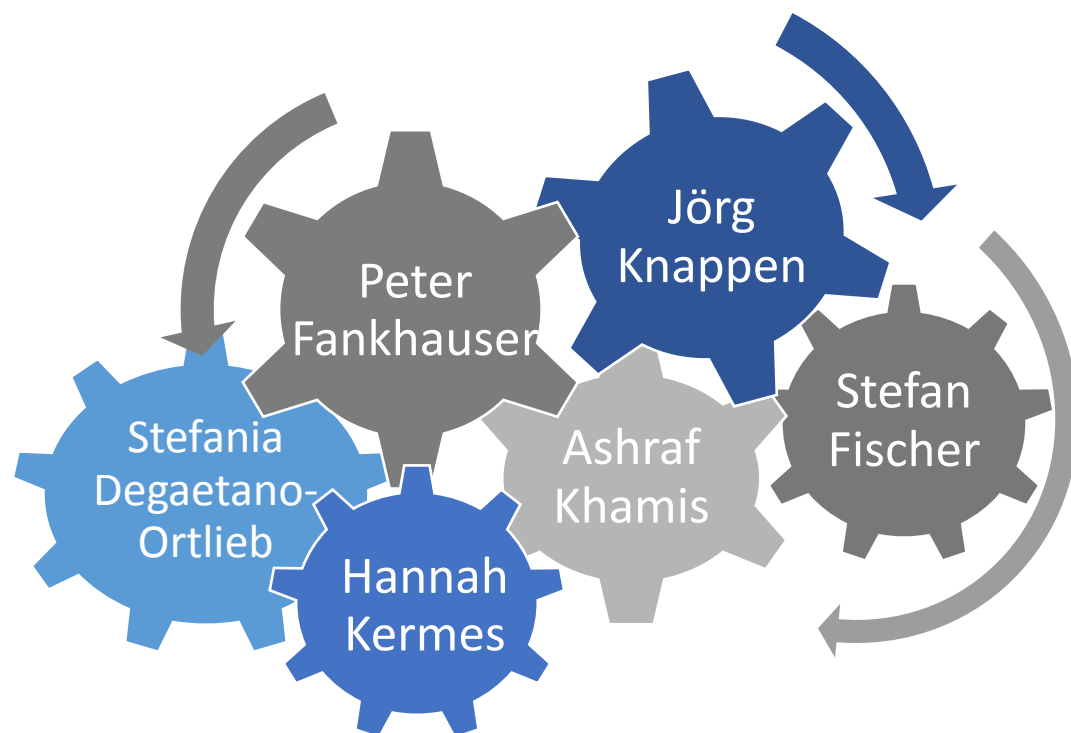
Topics

-
- Topical evolution of science
 - Diversification: more evenly distributed topics overall: *ent* increases
 - Specialization: topically more specific documents: *ment* decreases
 - Balance: Overall Complexity (*ent*) vs. Individual Complexity (*ment*):
 $ent - ment = (\text{generalised}) \text{ Jensen Shannon Divergence } jsd \text{ increases}$
 - Complexity balance as a general principle?
 - Other distributions/models (e.g. word/pos ngrams)
 - Other contextual dimensions (e.g. register)
 - Example (Brown family): Registers become more diverse, documents more specific

Word embeddings (word2vec)



Thanks to team and sponsors



References

- Matthew Aylett and Alice Turk. The Smooth Signal Redundancy Hypothesis: A Functional Explanation for Relationships between Redundancy, Prosodic Prominence, and Duration in Spontaneous Speech *Language and Speech* 47:31-56, March 2004 .
- Matthew W. Crocker, Vera Demberg, and Elke Teich. Information Density and Linguistic Encoding (IDeaL). *KI - Künstliche Intelligenz*, 30(1):77-81, 2016.
- Stefania Degaetano-Ortlieb and Elke Teich. Information-based modeling of diachronic linguistic change: from typicality to productivity. In *Proceedings of Language Technologies for the Socio-Economic Sciences and Humanities* (LATECH'16), ACL, Berlin, Germany, 2016.
- Stefania Degaetano-Ortlieb, Hannah Kermes, Ashraf Khamis and Elke Teich. An Information-Theoretic Approach to Modeling Diachronic Change in Scientific English. In Carla Suhr, Terttu Nevalainen and Irma Taavitsainen (eds), *Selected Papers from Varieng - From Data to Evidence* (d2e) 2015, Helsinki, to appear.
- Peter Fankhauser, Jörg Knappen, and Elke Teich. Exploring and Visualizing Variation in Language Resources. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation* (LREC), Reykjavik, 2014.
- Peter Fankhauser, Jörg Knappen and Elke Teich. Topical Diversification over Time in the Royal Society Corpus, *Proceedings of DH 2016*, Krakow, Poland, 2016.
- Dmitriy Genzel and Eugene Charniak. Entropy rate constancy in text. In *Proceedings of ACL*, Philadelphia, 2002.

References

- John Hale. A probabilistic Earley parser as a psycholinguistic model. In *Proceedings of the 2nd Conference of the North American Chapter of the Association for Computational Linguistics*, volume 2, pages 159–166, Pittsburgh, PA, 2001.
- Florian T. Jaeger. Redundancy and reduction: Speakers manage syntactic information density. *Cognitive psychology* 61(1):23-62, 2010.
- Hannah Kermes, Stefania Degaetano-Ortlieb, Ashraf Khamis, Jörg Knappen and Elke Teich. The Royal Society Corpus: From Uncharted Data to Corpus. In *Proceedings of the 10th International Conference on Language Resources and Evaluation (LREC)*, Portorož, Slovenia, 2016.
- Roger Levy. 2008. Expectation-based syntactic comprehension. *Cognition* 106(3):1126–1177.
- Steven T. Piantadosi, Harry Tily and Edward Gibson. Word lengths are optimized for efficient communication, *Proceedings of the National Academy of Sciences* 108(9):3526-3530, 2011,.