



18.02.2016 – Workshop Göttingen / Hildesheim

GEORG ECKERT
INSTITUT
Leibniz-Institut für internationale
Schulbuchforschung

 **DIPF**
Bildungsforschung
und Bildungsinformation

 TECHNISCHE
UNIVERSITÄT
DARMSTADT

 Stiftung Universität Hildesheim
2003

Welt der Kinder

Children and their World

User Centered Development for the DH – Experiences from the project

Ben Heuwing, Christa Womser-Hacker, Thomas Mandl
Institut für Informationswissenschaft und Sprachtechnologie, Universität Hildesheim

heuwing | womser | mandl@uni-hildesheim.de



Children and their World

Assumption: Text books reflect contemporary world interpretation of children:

General school attendance started

Textbooks as main information source for young adults

No travelling for most people

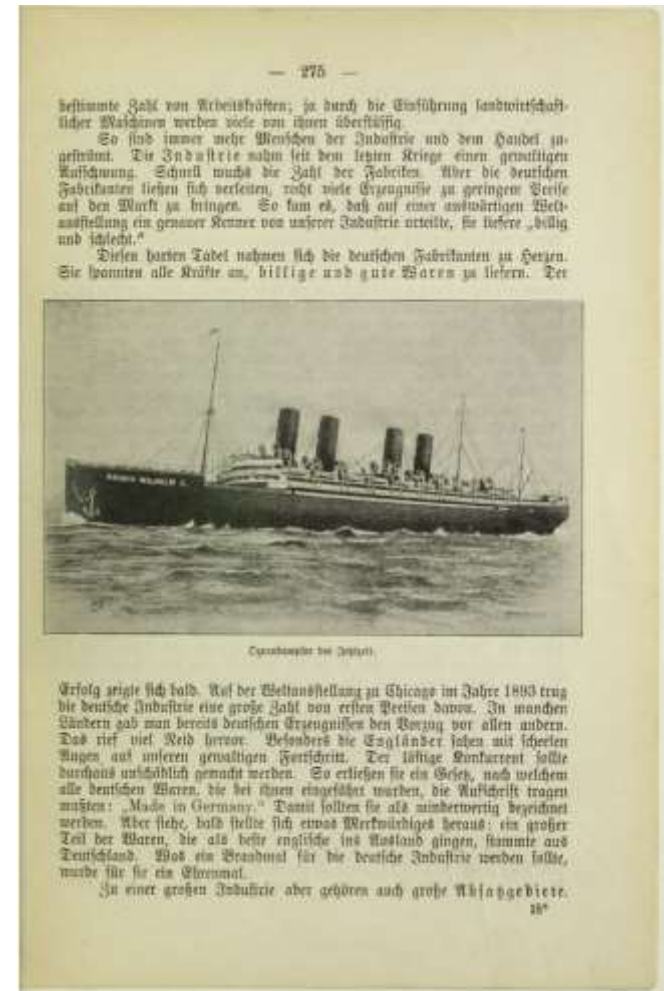
Historians want to explore German text books from the period between 1850 and 1918

Accelerated production of knowledge

Simultaneous processes of Globalization & Nationalization

Combination

of established hermeneutic methodology with innovative DH-methods



Funded by the Leibniz Association

Georg Eckert Institute for International Textbook Research

Prof. Dr. Simone Lässig, Dr. Robert Strötgen, Dr. Andreas Weiß, Maik Fiedler

Ubiquitous Knowledge Processing Lab (UKP) TU Darmstadt, and German Institute for International Educational Research (DIPF)

Prof. Dr. Iryna Gurevych, Dr. Richard Eckart de Castilho, Carsten Schnober

Institute for Information Science and Language Technology in Hildesheim (IWIST)

Prof. Dr. Christa Womser-Hacker, Prof. Dr. Thomas Mandl, Ben Heuwing

Other Partners

Institute of Popular Culture Studies, University of Zurich

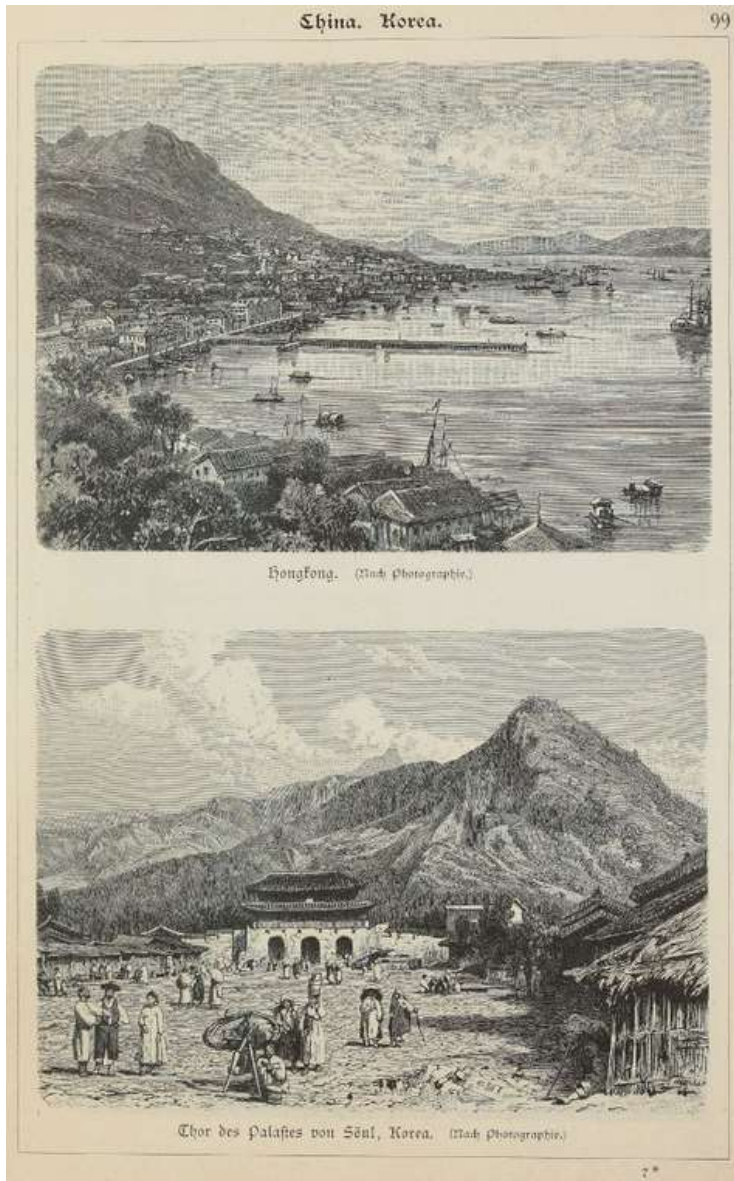
Bavarian State Library

Göttingen Centre for Digital Humanities

Universitätsbibliothek Technische Universität Braunschweig



Research Questions: Example



Perceptions of alterity and race in the context of colonialism

- How are foreign countries described?
- How is the role of Germany in the World presented?
- What is the role of citizens of foreign countries in Germany?
- How widespread were notions of race that were later pursued by national socialism?

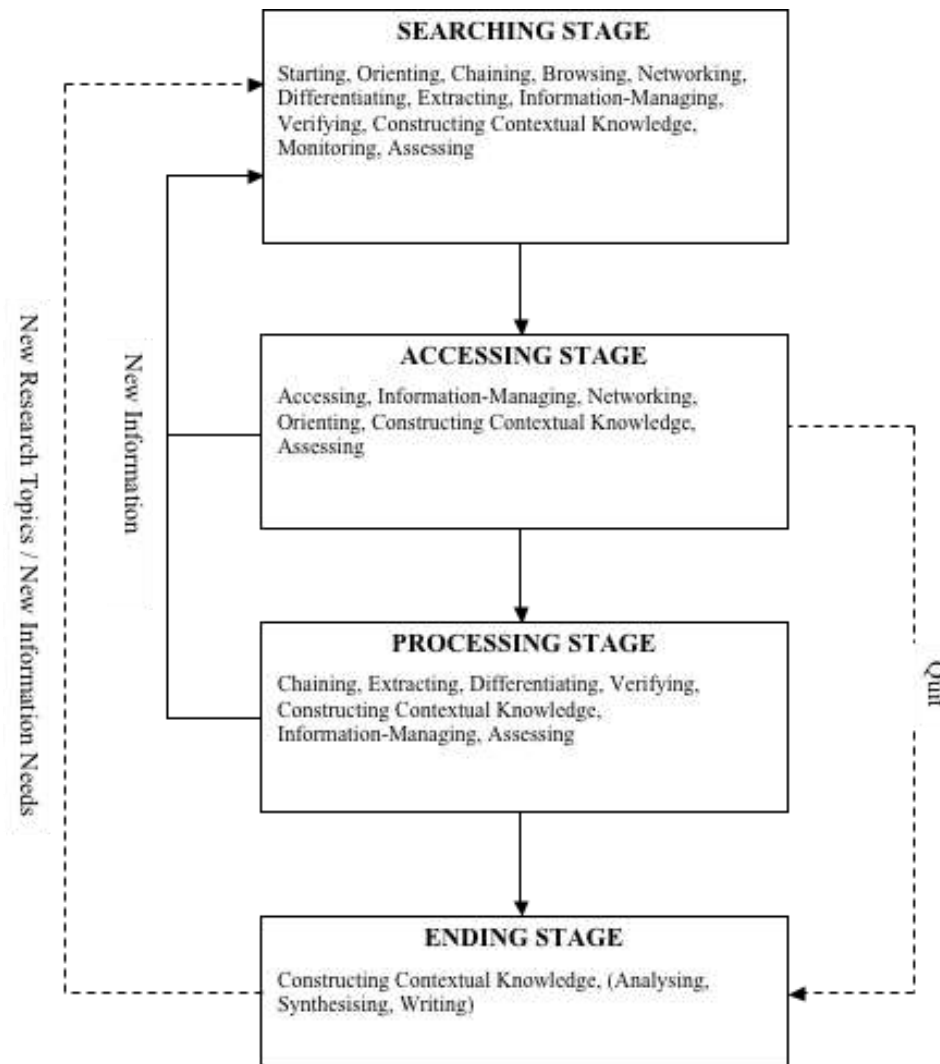
Information Seeking Behavior in the DH
Requirements Analysis in the DH
Interactive Analysis of Large Text Corpora

Research Interests of Information Science

Information Seeking of Historians

- Starting
- Orienting**
- Chaining
- Browsing
- Differentiating
- Verifying
- Extracting
- Accessing
- Monitoring
- Networking
- Information-managing
- Constructing contextual knowledge**
- Assessing**

Analysis? Corpus Level?
Distant Reading?



Rhee, Hea Lim (2012): Modelling historians' information-seeking behaviour with an interdisciplinary and comparative approach

Corpus



Data: ~6,000 digitized textbooks comprising ~900,000 pages

(~3,500 books, ~600,000 pages available yet)

Cannot be penetrated by hermeneutic methods alone

OCR

OCR-Errors: Estimated accuracy rate: 90-95% (characters)

Spelling Variations in the original document and across corpus

Metadata

Author, year & place of publication, publisher

Project internal classification of: *type of school, religious denomination (catholic/protestant), grade*

Subcorpus of childrens' books ~500

to be added

Pre-processing

OCR post-processing

Goal: identify erroneous words, replace with correct spelling

Rule-based: static mapping for recurring OCR mistakes

Statistical approaches: learn recurrent errors

Pre-processing for Retrieval and Topic Modeling:

Stop-word filtering, lemmatization

Filtering: capitalized words, sentiment words, named entities

Open Problems

Normalization of orthography?

Changes in meaning?

reiches des mittelalters erdteile kaiser als römische zum friedrichs großen inhalt
den bis germanen krieg gesehen allgemeine unter deutschland übersicht iii und
reich erdkunde ersten bei große brandenburgisch für zeittafel deutschlands zeit
völkerwanderung friedrich abschnitt karl das dem deutsche ein vorwort
dritter periode auflage neue von erster der gründung europa erste chr
geschichte 1648 revolution mit vom gegenwart aus reformation über anhang
zur zeitraum könig brandenburg durch römer auf nach inhaltsübersicht
inhaltsverzeichnis wilhelm zeitalter die römischen zweiter mittelalter
deutschen

Topic Modeling

Generative Model based on co-occurrences of terms in documents

Latent Dirichlet Allocation (LDA) - Blei 2003

A topic comprises all visible vocabulary words with weights assigned

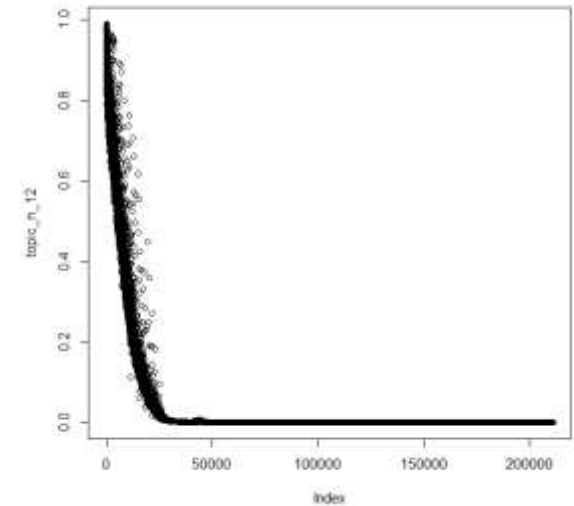
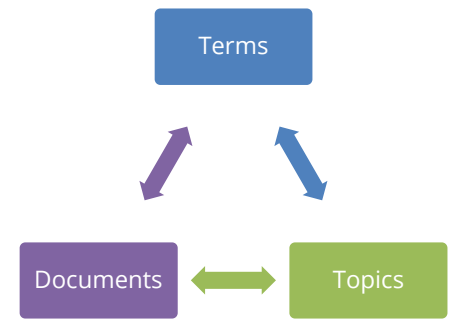
Intuition: words with largest weights are representative for a topic t

Application:

1. Generate k topics \rightarrow k topic weights for each word in vocabulary
2. Each document is a mixture of topics
For an (unseen) document, infer the proportion of each topic in that document (based on the model)

Topic number k : 50, 100, 150, 200, ...

Larger k increases topic specificity: more intuitive topics, but more redundancy and less clear inference



Sample Topics

Top 10 Words

Capital letters (200 topics):

Sieg, Eroberung, Zug, Tod, Herrschaft, Krieg, Iii, Kampf, Unterwerfung, Aufstand
Alexander, Perser, Babylon, Cyrus, Syrien, Israel, Euphrat, Aegypten, Asien, Darius
Reich, deutsch, Reiche, deutsche, Grenze, groß, Mark, Nation, Herrscher, Nachfolger

Named entities (200 topics):

Athen, Athener, Sparta, Athens, Griechenland, griechischen, Salamis, athenischen, Xerxes,
Aristides

Joseph, Isaak, Abraham, Juden, Jakob, Gott, Kanaan, Israel, Abrahams, Gottes

Westfalen, Bielefeld, Minden, Paderborn, Dortmund, Herford, Arnsberg, Soest, Hagen, Hamm

Sentiment words (50 topics).

Recht, Schuld, hart, streng, Klage, verbieten, bestrafen, schwer, Weise, schuldig

Tod, Erbe, Besitz, Anspruch, Streit, Macht, Recht, anerkennen, Vertrag, Konstanz

führen, Macht, erklären, mächtig, unterstützen, unabhängig, vollständig, glänzend,
Unterstützung, kräftig

Large Corpus

Rich Metadata

Assessments on Topics

Manual Annotations (?)

Lexical resources: Sentiment lexica (?)

Disambiguation of Historic Place Names, Person Names (x)

External Ontologies (DNB)

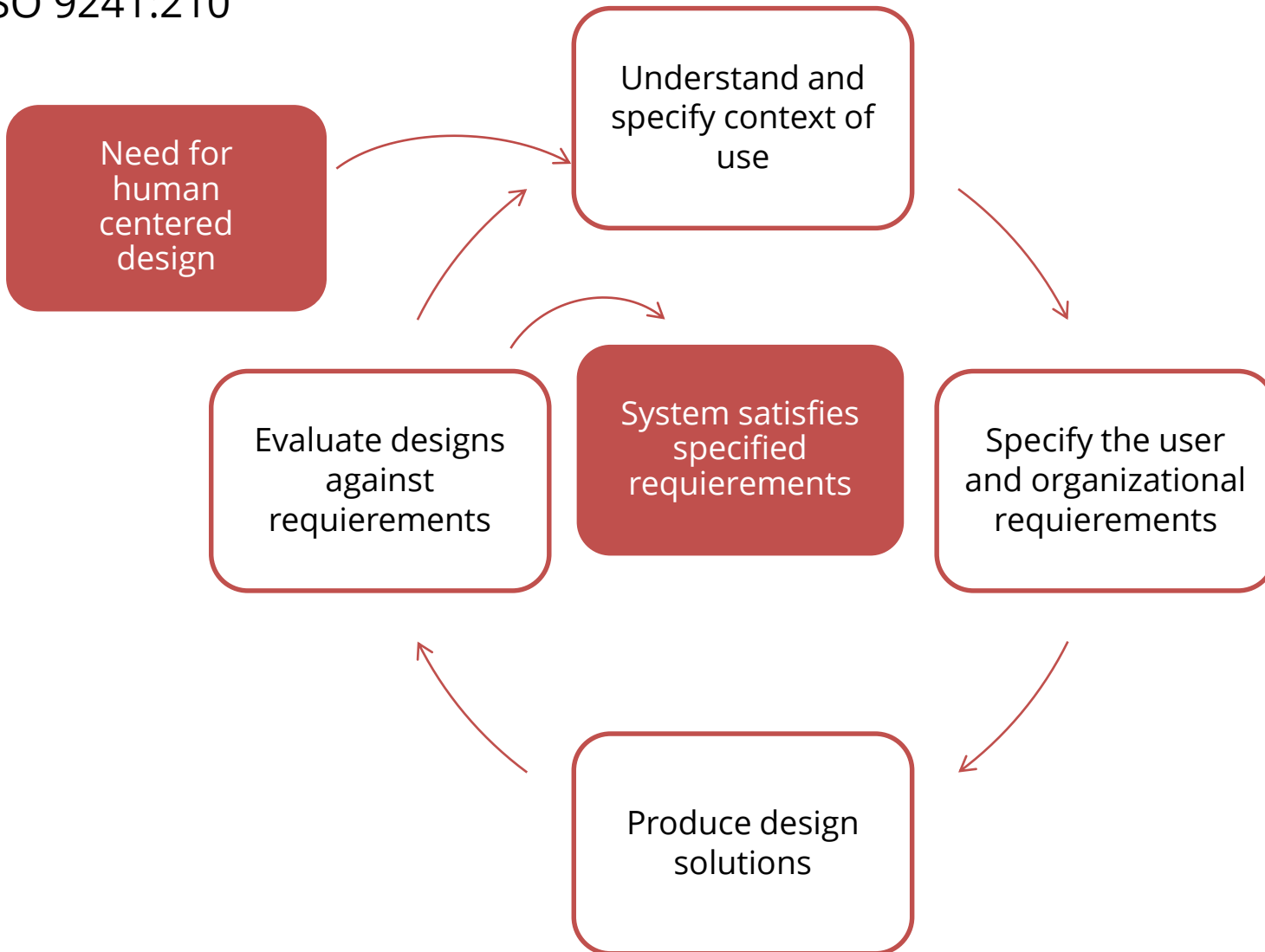
Summary: Resources

**How can these resources be employed
for a given research questions?**

Co-Development of research methodology and technology

Human-centred design process + Participation

ISO 9241:210



**Understand and specify
context of use**

Understand and specify context of use

Interviews

Contextual Interviews with 5 Historians

Including two from project WdK

Analysis of information representations used for research

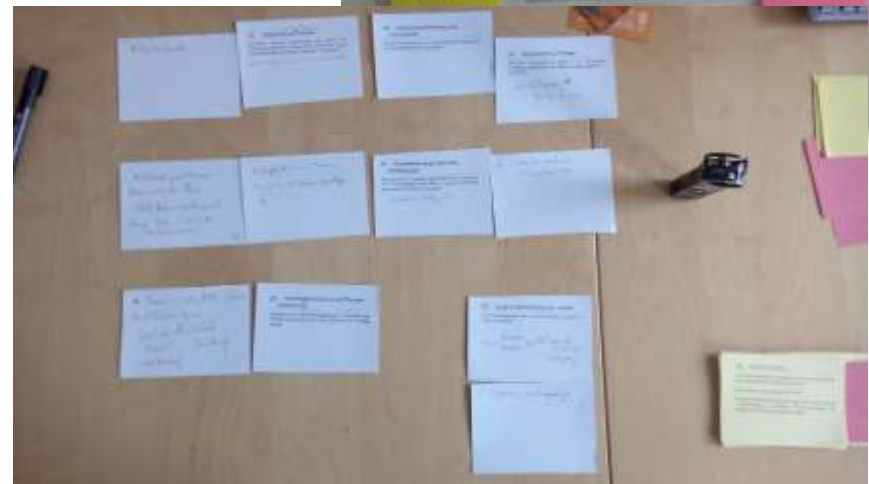
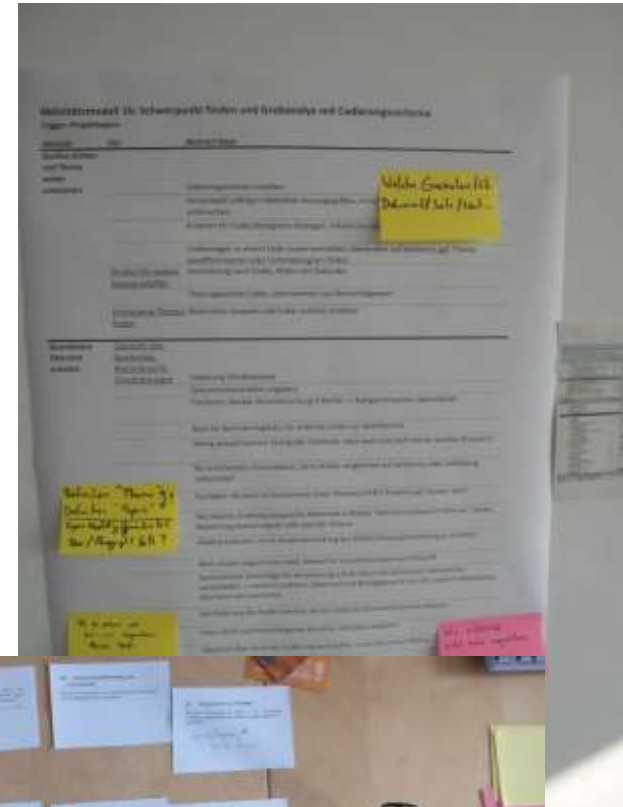
Internal Workshop to discussion results

Prioritization of requirements

Finding common terminology and understanding of goals

Participation:

Cooperation and Co-Design



Specify user requirements

Proposed Workflow for analysis

1. Discover Dimensions of Analysis

- Navigate through corpus
- Find changes, pattern and trends for sub-collections
- Create document sets and select topics for in depth analysis
- Find relevant documents

2. Quantitative comparisons across dimensions

- Find differences between sub-sets (terms, topics)
- Direct comparisons: topic intensity, topic correlations
- Subjunctive interface

3. Qualitative Analysis & Comparison

- Identify, collect and compare relevant examples from texts

4. Evaluate and correlate results

- Be aware of ones own preconceptions
- Be aware of skewed distributions across dimensions in the collection

Entry Points for
Analysis &
Exploration
Serendipity

Analysis across
Multiple Contexts

Evidence from
Multiple
Perspectives

Detect and
Manage Bias in the
Collection

Create Design Solutions

Working Prototype

Welt der Kinder
Explorer - Beta
WdK_06_05

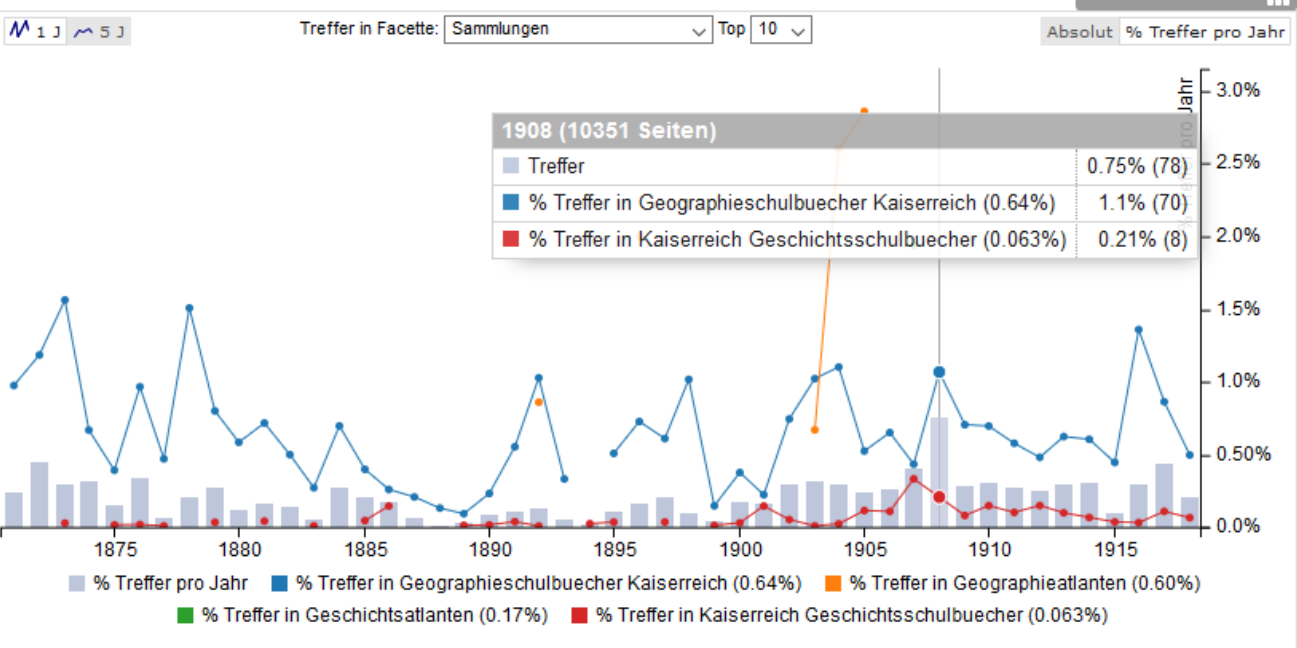
korea Suchen

Aktuelle Auswahl: [Anfrage:korea](#) [Erscheinungsjahr:\[1871 TO 1920\]](#)

2 Kategorien Seitenzahl | [Seiten % der Anfrage](#)

- Filter**
- Auflage**
 - k.A. oder Auflage 1 (457)
 - Auflage 2 und höher (622)
- Erscheinungsjahr**

1871 : 1920 OK



- Von 1850 bis 1920 (1079)
- [1850 - 1860](#) (0)
 - [1860 - 1870](#) (0)
 - [1870 - 1880](#) (154)
 - [1880 - 1890](#) (95)
 - [1890 - 1900](#) (94)
 - [1900 - 1910](#) (349)
 - [1910 - 1920](#) (387)
- [Alle](#)

- Sammlungen**
- Inhalt**
- Schultyp_Allg_wdk**
- Schultyp_wdk**
- Schulformen_opac**
- Bildungsstufe_opac**
- Erscheinungsort**
- Ort_erste_angabe_wdk**
- Regionen_opac**

Topic Model: TM Hauptwörter (50) TM

Topic: Topic auswählen

Seiten Gruppirt nach Werk Sortieren: Relevanz zu Anfragetermen Erscheinungsjahr

1079 Seiten 1 von 108 [weiter](#)

1. Bilder-Atlas zur Geographie der außereuropäischen Erdteile - S. 103
 1901 - Leipzig [u.a.] : Bibliogr. Inst.
 Autor: Geistbeck, Alois
 Auflage: 2
 Sammlungen: Geographieschulbuecher Kaiserreich

**Evaluate designs against
requirements**

Evaluating Topic Models

Assessment of Domain Experts (Historians)

Coherence of Topics

Value for research Questions

Dokumente zu den Topics: http://wdk.uk	Relevanz für Fragestellung			Überraschung	Kommentare
	Sehr relevant	relevant	Nicht relevant	Ist das Topic unerwartet?	
Topic					
T16: [Gebirge, Alpen, Klima, Höhe, Teil, Flüsse, Ebene, Norden, Süden, Boden]	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	
T15: [Heinrich, Kaiser, Otto, Friedrich, Papst, Italien, Deutschland, Sachsen,	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	

Task-Based Validation: Models & Interface

Cooperative, task-based exploration

Applying tools to real tasks

4 Hypotheses which are results of previous historical research

Extending Tools

Ad-hoc analysis in Excel et al.

Interdisciplinary learnings

Coverage of the recent, national history of Germany and Prussia increases

because of a new orientation in Germany's official educational policy in the year 1890,

accompanied by a decrease of contents about ancient history

Example - Hypothesis

Teilnehmer 1

Festlegen: Untersuchungszeitraum 1880-1900, Topics

WdK_06_05

2 Kategorien

Seitenzahl | % der Ergebnisse

Filter

Auflage

Erscheinungsjahr

1880 : 1900 OK

Von 1850 bis 1920 (172726)

1850 - 1860 (0)

1860 - 1870 (0)

1870 - 1880 (0)

1880 - 1890 (76544)

1890 - 1900 (86239)

1900 - 1910 (9943)

1910 - 1920 (0)

Alle

Sammlungen

Inhalt

Schultyp_Allg_wdk

Schultyp_wdk

Schulformen_opac

Bildungsstufe_opac

Erscheinungsort

#	Topic	Mittel
	fuß]	
38	T37: [luther, lehre, dichter, schrift, kirche, werk, wittenberg, universität, lied, gedicht]	0,97%
17	T16: [handel, industrie, ackerbau, gewerbe, viehzucht, stadt, arbeiter, fabrik, arbeit, verkehr]	0,89%
31	T30: [friedrich, wilhelm, könig, kaiser, iii, karl, kurfürst, prinz, johann, preußen]	0,87%
32	T31: [wald, baum, tier, pferd, boden, feld, wiese, mensch, pflanze, rind]	0,86%
49	T48: [rom, römer, italien, rmer, krieg, stadt, spanien, volk, gallien, hannibal]	0,82%
2	T1: [gott, tempel, mensch, gtter, priester, zeus, held, erde, opfer, himmel]	0,78%
39	T38: [napoleon, schlacht, franzose, heer, sieg, juni, paris, preußen, mai, leipzig]	0,76%
47	T46: [stadt, hauptstadt, straße, eisenbahn, dorf, einw., handel, verkehr, einwohner, berlin]	0,71%
1	T0: [general, armee, heer, mann, truppe, schlacht, feind, franzose, festung, preußen]	0,66%
21	T20: [athen, athener, griechenland, sparta, stadt, spartaner, flotte, perser, insel, theben]	0,61%
35	T34: [insel, meer, küste, halbinsel, ozean, afrika, kap, land, europa, teil]	0,61%
23	T22: [papst, kirche, rom, kaiser, bischof, heinrich, gregor, erzbischof, papste, vii]	0,60%

Aktuelle Ausw



Topic: T

Seiten

172726 Seite

1. Kursus

Participant 1
Topics for 1880-1900

172726 Pages

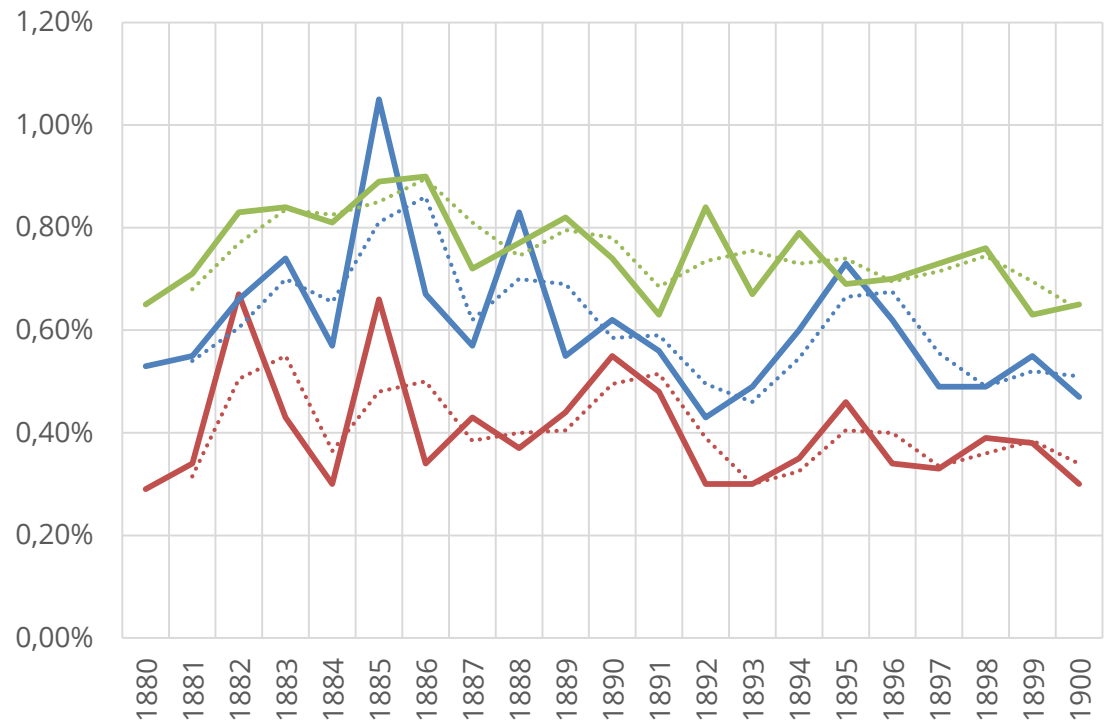
Antike:

T20: [athen, athener, griechenland, sparta, stadt, spartaner, flotte, perser, insel, theben]

T5: [rom, cäsar, senat, pompejus, jahr, sulla, antonius, marius, csar, provinz]

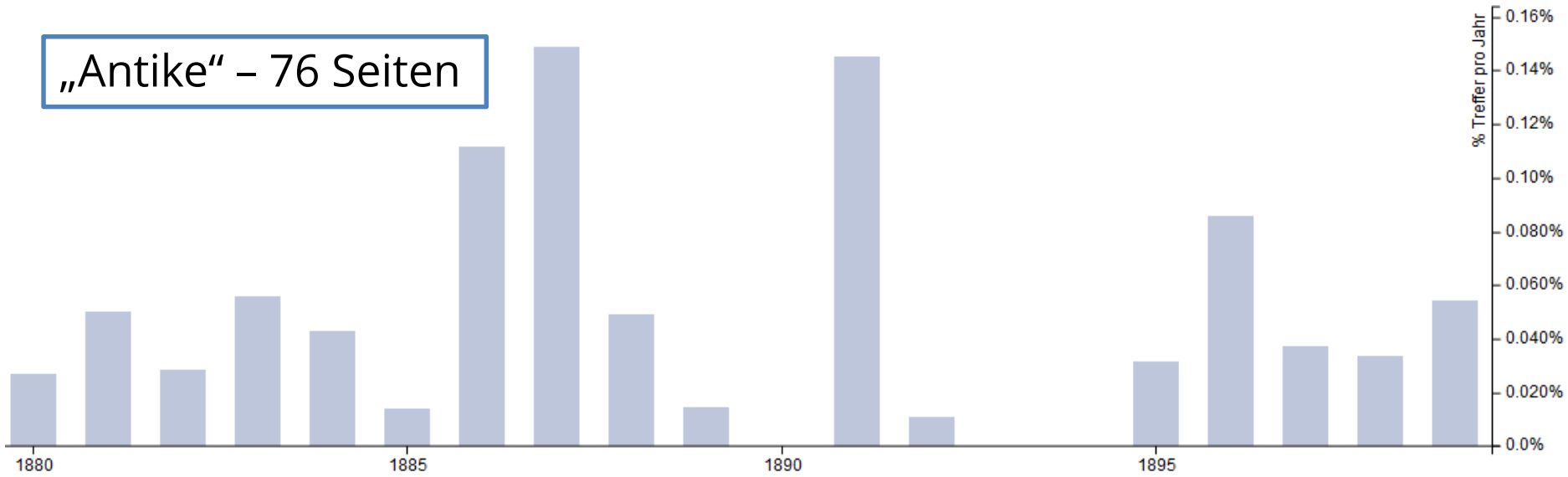
Französische Revolution & Befreiungskriege

T38: [napleon, schlacht, franzose, heer, sieg, juni, paris, preußen, mai, leipzig] Topicfilter

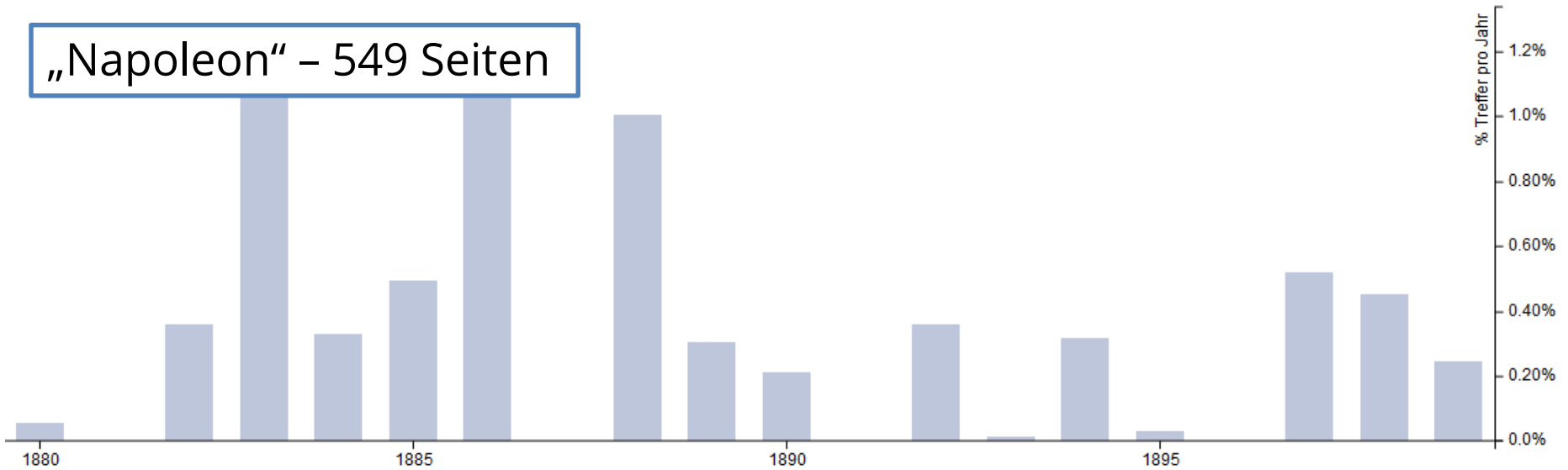


Participant 1
Term-Search on Textbooks for „Gymnasium“

„Antike“ – 76 Seiten



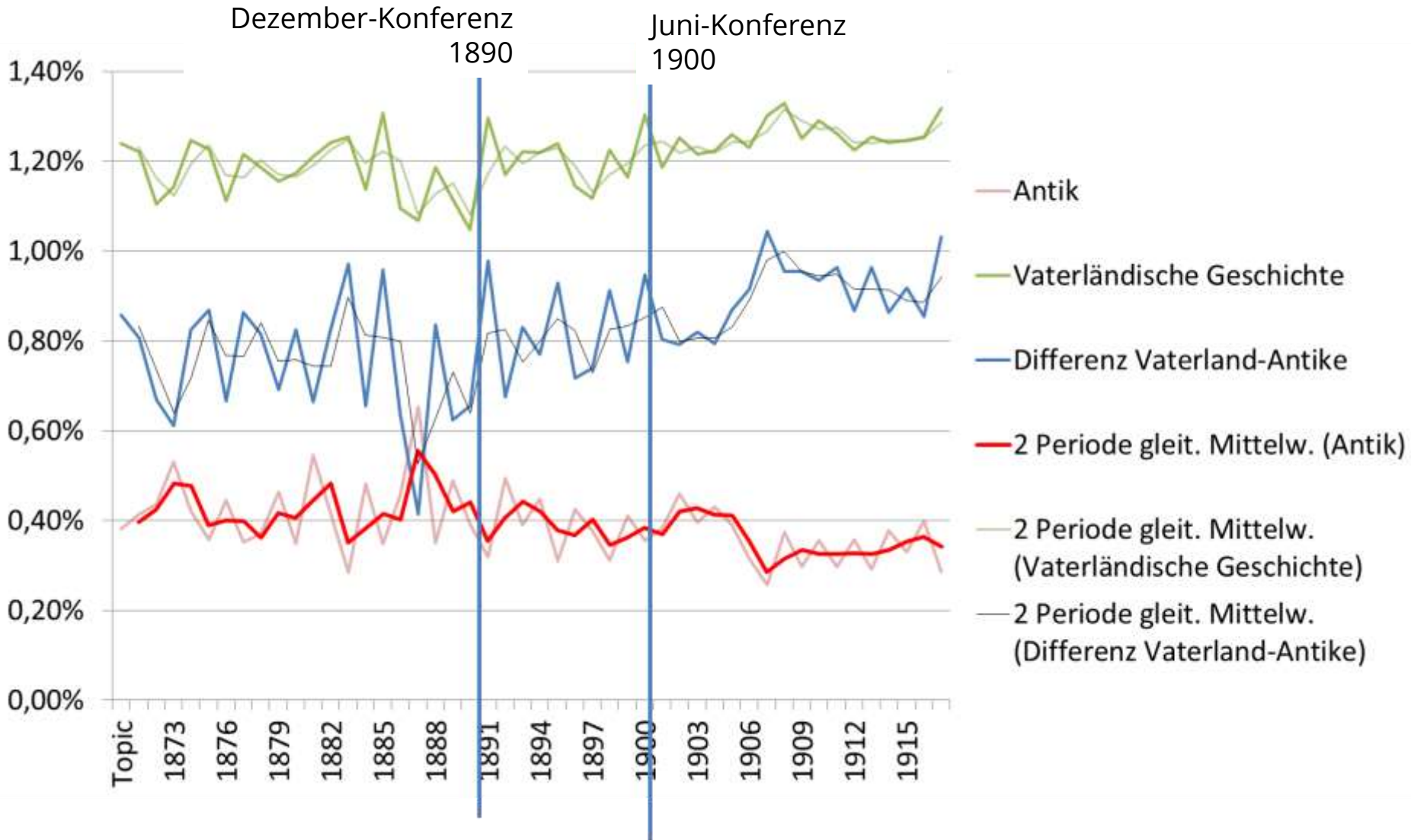
„Napoleon“ – 549 Seiten



Participant 2

Combination of Topics

Untersuchung des Verhältnisses für größeren Zeitraum

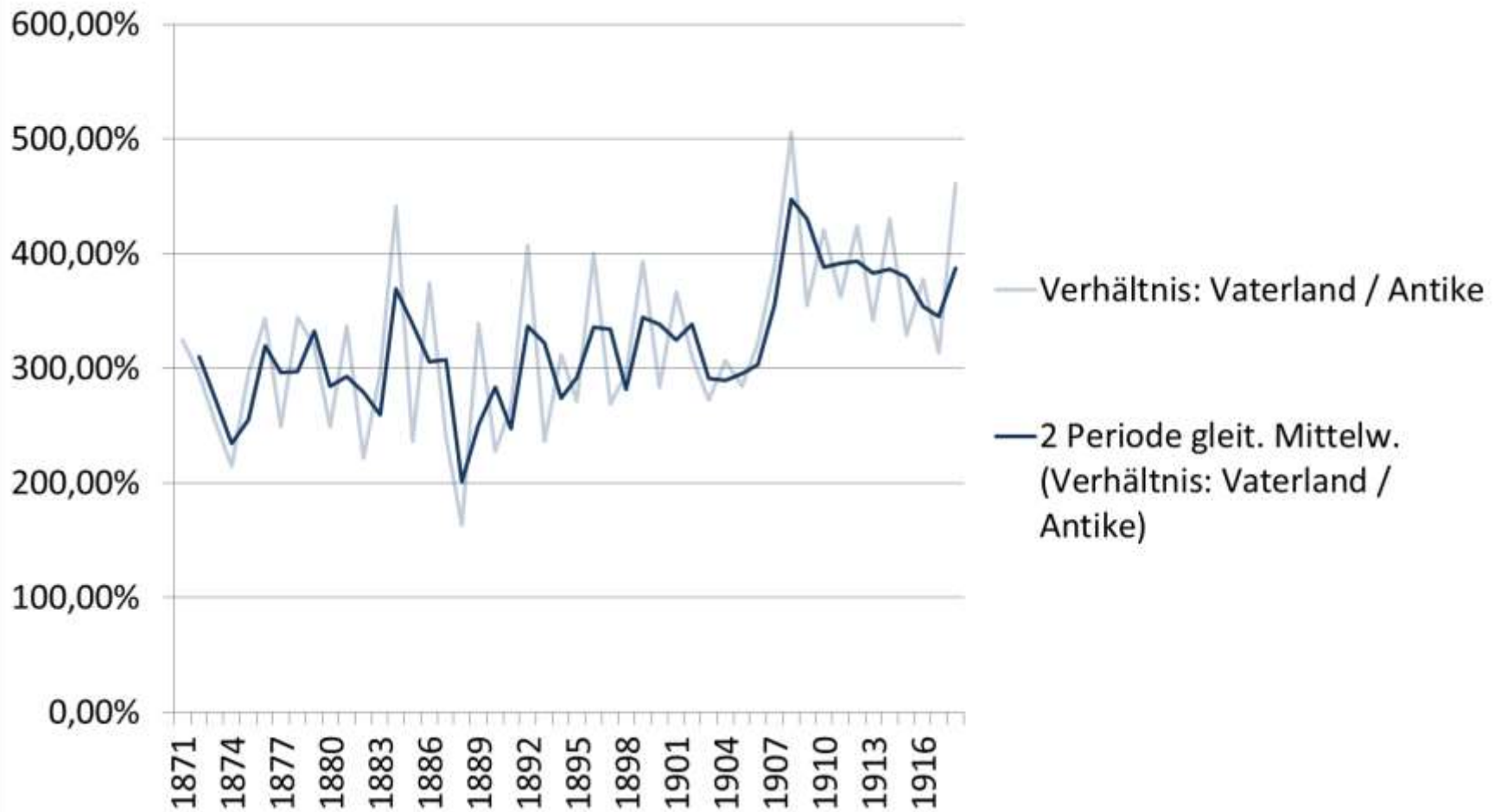


Teilnehmer 2

Kombination von Topics

Untersuchung des Verhältnisses für größeren Zeitraum

Verhältnis: Vaterland / Antike



Gesamtergebnis der Validierungsstudie

These	Tn	Beurteilung	Grundlage	Belastbarkeit	Transparenz
1: Kolonien - Erdkunde / Geschichte	1	Teilweise bestätigt	Topics	4	5
	2	bestätigt	Topics & Terme	3,5	4
2: Kolonien - Volksschulen / Höhere Schulen	2	bestätigt	Topics & Terme	4	4
3: Vaterländische Geschichte / Antike	1	Nicht eindeutig		3	4
	2	bestätigt (in der Tendenz)	Topics	4	3,5
4: Kriegsflotte Erdkunde / Geschichte	2	widerlegt	Terme	3	2

Results of the Validation Study

Tools and models appear to be valid

Importance of Transparency and Validity

Demonstrate validity using different tools and methods

Level of analysis

Topics in many cases are too specific or too general

→ Implement Topic-Combination as a tool

Planned Steps of Evaluation

User Studies, based on Word-Intrusion, Topic Intrusion

- Evaluate new approach to text modeling in comparison to prototype (Chang et al. 2009)

IR-Evaluation Studies

- TopWords from Topics (tf*idf) compared to ranking by Topic-Intensity derived from models

Additional interactive Evaluation Studies

- Focus on Usability
- Focus on subjective measures of Recall & Precision from Interactive Information Retrieval

Entry Points for Analysis & Exploration
Serendipity

Analysis across Multiple Contexts

Evidence from Multiple Perspectives

Detect and Manage Bias in the Collection



Detect Anomalies in Sub-Groups?

Scan for Pattern from external Ontologies?

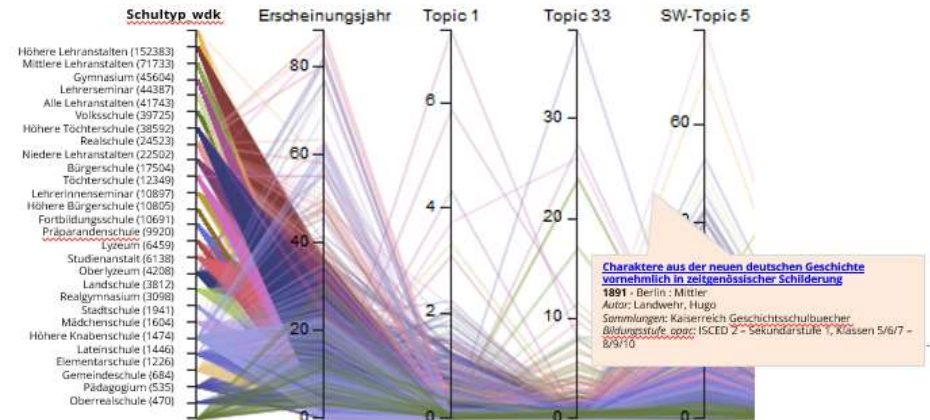


Visualizations grouped by work

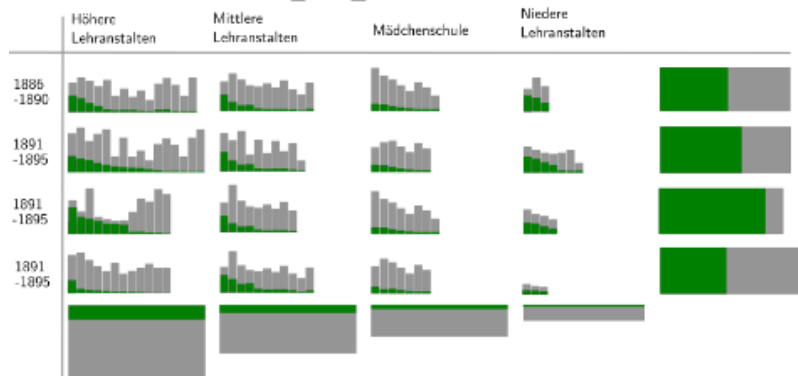
Multivariate Analysis

Include context information

Control for biases in the corpus across two variables



Seitenzahlen für T5: [rom, cäsar, sen ...] Filter: 2%
in Kategorie: Schultyp_Allg_wdk nach Erscheinungsjahr



Incorporate External Knowledge Sources

Related Persons as Query-Suggestions

Analysis of narratives based on history

WDK Corpus-Explorer

Bismarck

Relevanz Jahr

Suchergebnisse

Die Völkerkunde im Unterricht an den höheren Schulen
1910 - Braunschweig : Graf - Woldemann, August (Autor)
Sammlung für Geographisch-historische Kaiserreich - Höhere Lehranstalten
Geschicht, wdk, Tuzgen - Bildungsziele, opac: Sonstige Lehrmittel, alle Lehrstufen

[12 wozu mit Unrecht hat man Pechel einen Vorwurf daraus gemacht], daß er von der braunbäuerischen Rasse der Äthiopier mehrere kleine Rassen abgetrennt hat, wie die Australier, Papuas und Dravida. Namentlich Forscher[?] haben später vorgeschlagen, die sämtlichen dunkelhäutigen Rassen als Teile zu der Rasse der Negroiden zusammenzufassen, doch sind ihnen hierin die meisten neueren Ethnographen nicht gefolgt. (Heterland) [?] sagt darüber: „Vielfach werden die selben genannten vier Rassen (Australische Rasse, papuatisch-melanesische Rasse, australische Rasse, dravidische Rasse) nur als verschiedene Typen einer einzigen negroiden östlichen Rasse aufgeführt. [...] Negroiden oder negroiden Rasse immer weniger Anhänger findet.“ Jedenfalls sind die Verbreitungsbereiche dieser sogenannten negroiden Rassen heute geographisch weit voneinander getrennt. [...]

Die Völkerkunde im Unterricht an den höheren Schulen - S. 14 (12)
1910 - Braunschweig : Graf - Woldemann, August (Autor)

Verbindungen anzeigen zu:

- [Bismarck, Otto Fürst von \(seit 1871\)\(1815-1898\)](#) ⓘ
Politiker, Reichskanzler
- [Bismarck, von](#) ⓘ
altmärkische Adelsfamilie

WDK Corpus-Explorer

Bismarck

Relevanz Jahr

Suchergebnisse

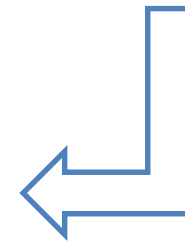
Die Völkerkunde im Unterricht an den höheren Schulen
1910 - Braunschweig : Graf - Woldemann, August (Autor)
Sammlung für Geographisch-historische Kaiserreich - Höhere Lehranstalten
Geschicht, wdk, Tuzgen - Bildungsziele, opac: Sonstige Lehrmittel, alle Lehrstufen

[12 wozu mit Unrecht hat man Pechel einen Vorwurf daraus gemacht], daß er von der braunbäuerischen Rasse der Äthiopier mehrere kleine Rassen abgetrennt hat, wie die Australier, Papuas und Dravida. Namentlich Forscher[?] haben später vorgeschlagen, die sämtlichen dunkelhäutigen Rassen als Teile zu der Rasse der Negroiden zusammenzufassen, doch sind ihnen hierin die meisten neueren Ethnographen nicht gefolgt. (Heterland) [?] sagt darüber: „Vielfach werden die selben genannten vier Rassen (Australische Rasse, papuatisch-melanesische Rasse, australische Rasse, dravidische Rasse) nur als verschiedene Typen einer einzigen negroiden östlichen Rasse aufgeführt. [...] Negroiden oder negroiden Rasse immer weniger Anhänger findet.“ Jedenfalls sind die Verbreitungsbereiche dieser sogenannten negroiden Rassen heute geographisch weit voneinander getrennt. [...]

Verbindungen zu:

- [Bismarck, Otto Fürst von \(seit 1871\)\(1815-1898\)](#) (10278)
Politiker, Reichskanzler ⓘ
- [Bismarck, August Wilhelm von \(1750-1783\)](#) (929/929)
preußischer Kriegsminister und Gesandter am dänischen Hof ⓘ
- [Bismarck, Levin Friedrich von \(1703 bis 1774\)](#) (0/0)
preußischer Justizminister ⓘ

(Treffer einzeln | Treffer gemeinsam)



Open questions...

Document ranking for topic models?

Relevant document with multiple top topics do not appear on top → not visible

More specialized modeling techniques necessary?

e.g. Interactive topic modeling? Dynamic and/or structural topic models?

Interpretation and trust in models and visualizations more important than perplexity/novelty
(Chuang et al. 2012)

Subjectivity and evaluative Language

Targeting topics, named entities, ...? Ideologies?

Statistical significance ...

... of differences in document sets?

... of correlations of topics in document sets?

Infrastructure needed for DH (Digital History) based on large corpora

Robust Correction of OCR

Robust Text-Normalization

cf. Efforts of BBAW

Robust Disambiguation of Entities (Events, Persons, Places, Dates) and Linking to external Knowledge Sources

Integration of manual Annotations at scale

Literature

- Blei, D. M. ; Ng, A. Y. ; Jordan, M. I. (2003): Latent dirichlet allocation. In: *Journal of Machine Learning Research* Bd. 3, S. 993–1022
- Chang, J. ; Boyd-Graber, J. L. ; Gerrish, S. ; Wang, C. ; Blei, D. M. (2009): Reading Tea Leaves: How Humans Interpret Topic Models. In: *NIPS*. vol. 22, pp. 288–296
- Chuang, J. ; Ramage, D. ; Manning, C. ; Heer, J. (2012): Interpretation and Trust: Designing Model-driven Visualizations for Text Analysis. In: *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems, CHI '12*. New York, NY, USA: ACM, S. 443–452
- Freund, L.; Toms, E. G; Waterhouse, J. (2005): Modeling the information behaviour of software engineers using a work-task framework. In: *Proceedings of the American Society for Information Science and Technology* Bd. 42, Nr. 1
- Rhee, Hea Lim (2012): Modelling historians' information-seeking behaviour with an interdisciplinary and comparative approach. In: *Information Research* Bd. 17, Nr. 4 - <http://InformationR.net/ir/17-4/paper544.html>
- Womser-Hacker, Christa (2013): Information Seeking Behaviour (ISB). In: Kühlen, Rainer; Semar, Wolfgang; Strauch, Dietmar (Hrsg.) *Grundlagen der praktischen Information und Dokumentation: Handbuch zur Einführung in die Informationswissenschaft und -praxis*. 6., völlig neu gefasste Ausgabe. De Gruyter Saur