

# Big Data in den Digital Humanities

Dr. Jochen Tiepmar

Abteilung Automatische Sprachverarbeitung, Universität Leipzig

# Fragen

- Was bedeutet Digital Humanities & Big Data?
- Digital Humanities sind keine klassische Big Data Anwendung. Wo sind die Unterschiede?
- Welche Besonderheiten der Digital Humanities sind dabei zu berücksichtigen?
- Was sind praktische DH-Anwendungen für Big Data

# Digital Humanities (DH)

- Supporting research in humanities with tools from computer science  
Datenbank-Spektrum, Springer Verlag, 2015.  
<http://dx.doi.org/10.1007/s13222-014-0177-7>
- Schnittstelle zwischen Informatik und Geisteswissenschaften
- Kombiniert Methoden von  
Archeologie, Kunst, Kultur & Sozialwissenschaften, Geschichte, Linguistik,  
Literatur, Philosophie, Musik, Politik,...
- Mit Tools aus  
Datenvisualisierung, Data Mining, Mapping , Information Retrieval, Statistik,  
Text Mining,...
- Wir betrachten hier nur textorientierte DH

# DH Anfänge

- Index Thomisticus:
  - Konkordanz der Schriften des Thomas von Aquin
  - Texte mit 11 Millionen laufenden Wortformen
  - entspräche auf Lochkarten gespeichert einem Papiergewicht von mehr als 100 Tonnen
  - Erstes Projekt im Bereich Humanities Computing
  - Kooperation mit IBM, kontinuierliche Nutzung neuer Speichertechnologien (Magnetband, Festplatten, CD)
  - Erstellung von 56 gedruckten Bänden mit 70.000 Seiten

# DH Anfänge



*Roberto  
Busa Archive*

# Big Data

- Volume      ->      Umfang der Daten
  - Velocity    ->      Verarbeitungs/Berechnungsgeschwindigkeit
  - Veracity    ->      Verlässlichkeit der Informationen/Fehlerrate
  - Variety     ->      Heterogenität von Daten/Tools/,...
- 
- Je nach Anwendung:
    - Visualization, Value, Volatility, Vulnerability, Validity, Variability, (Quality)

# Big Data

- In klassischen Big Data-Anwendungsbereichen liegt der Fokus nicht gleichmäßig auf allen Vs
- Häufig werden Verbesserungen in Volume und Velocity als die einzigen Aufgaben der Informatik angesehen
- Kernanwendungen: Sensornetzwerke und Umweltdaten, Simulationen in Wissenschaft und Industrie

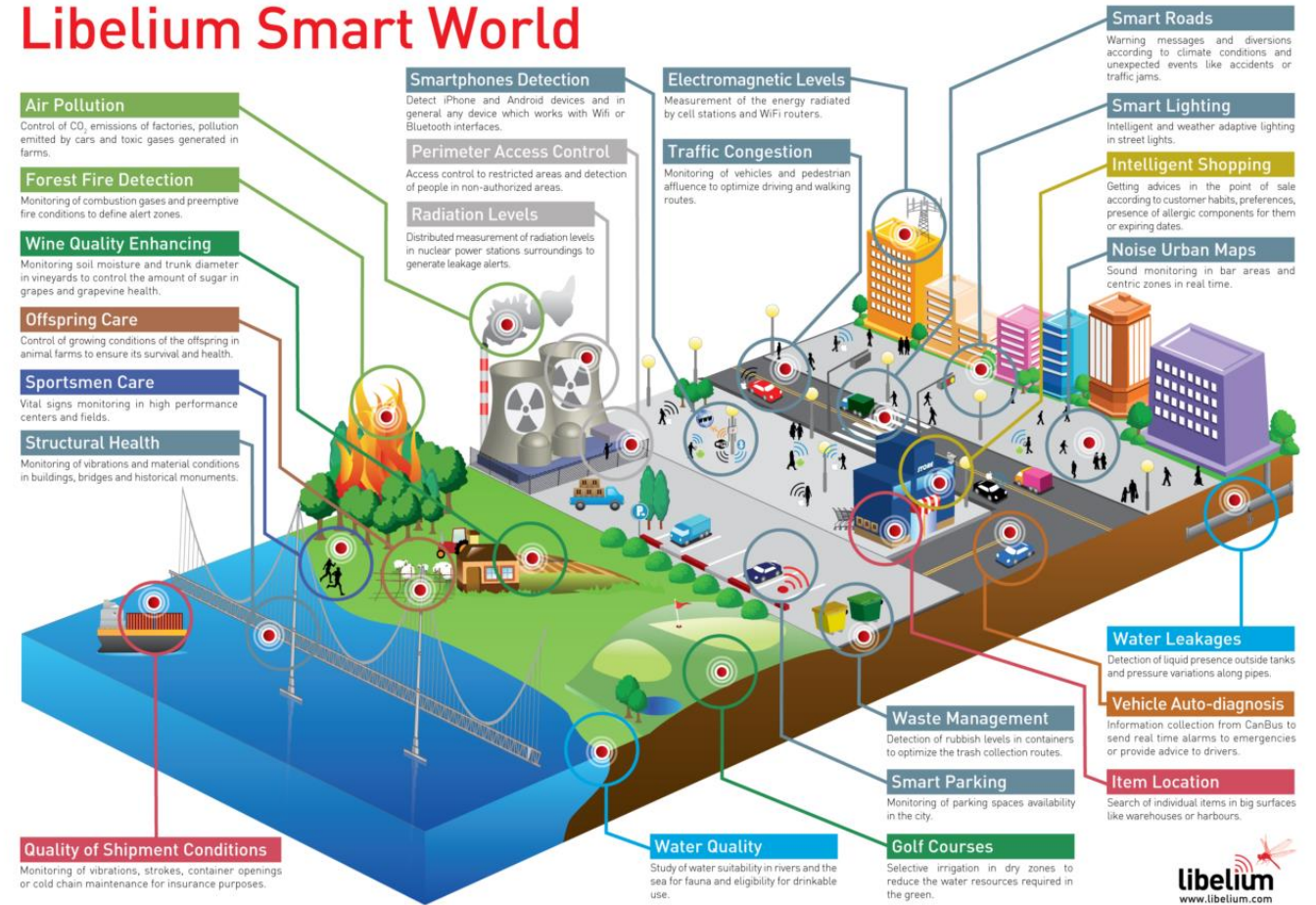
# Big Data

- Ein paar Anwendungsbeispiele



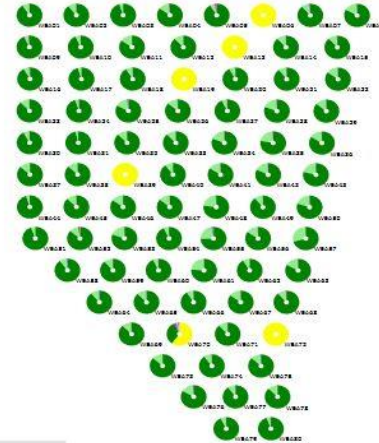
# Big Data

## Libelium Smart World

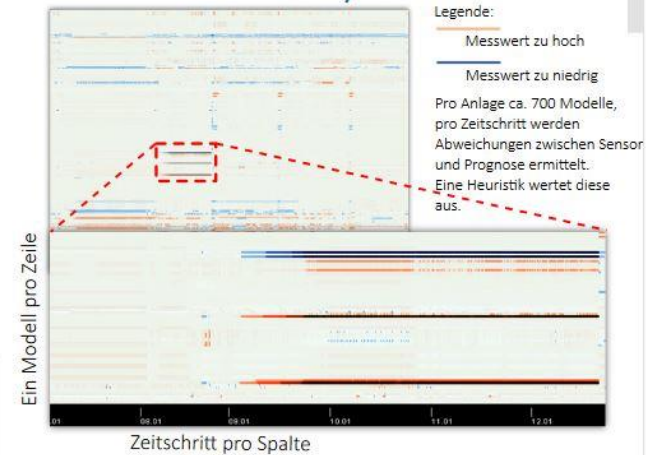


# Big Data

## Vollautomatische Zustandsüberwachung von Windparks



## Datenstromanalyse



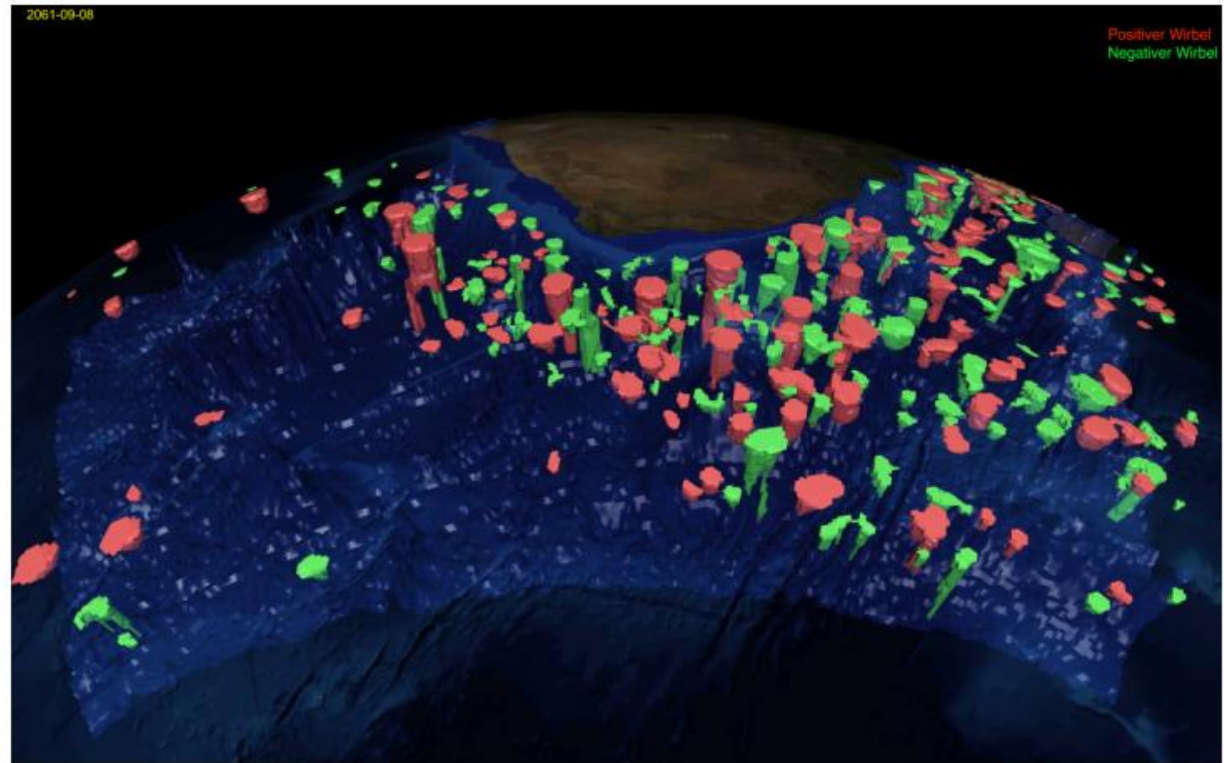
## Ergebnisse und Mehrwerte

- Erkennung von Sensordefekten, Trendänderungen, Druckverluste, Beschädigungen an Baugruppen, Leckagen, Ausrichtungsfehler
- Ermöglicht prädiktive Wartung
- Frühzeitige Erkennung von schleichenden Ertragsverlusten: Fehlausrichtung einer WEA im Offshorebereich ca. 40k € Schaden pro Monat

# Big Data



## Darstellung der Wasservolumen



Scheuermann

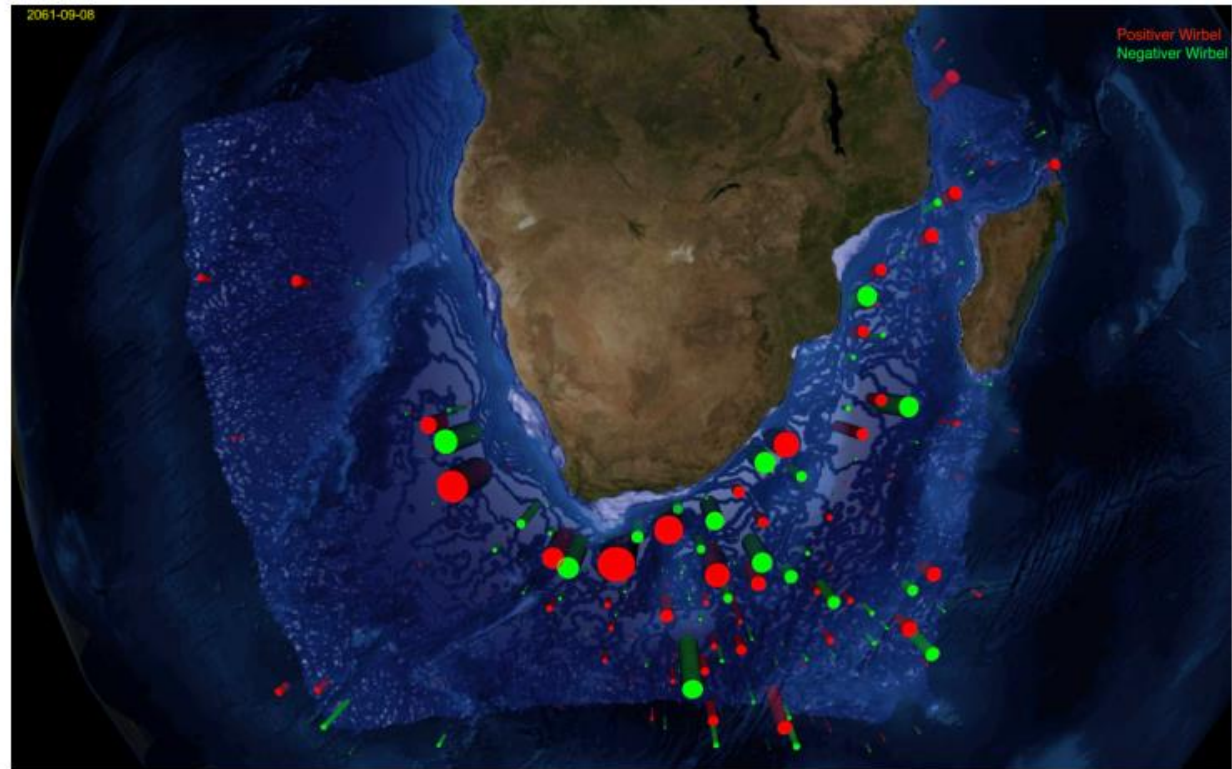
UNIVERSITÄT LEIPZIG



# Big Data



## Darstellung der Wasservolumen



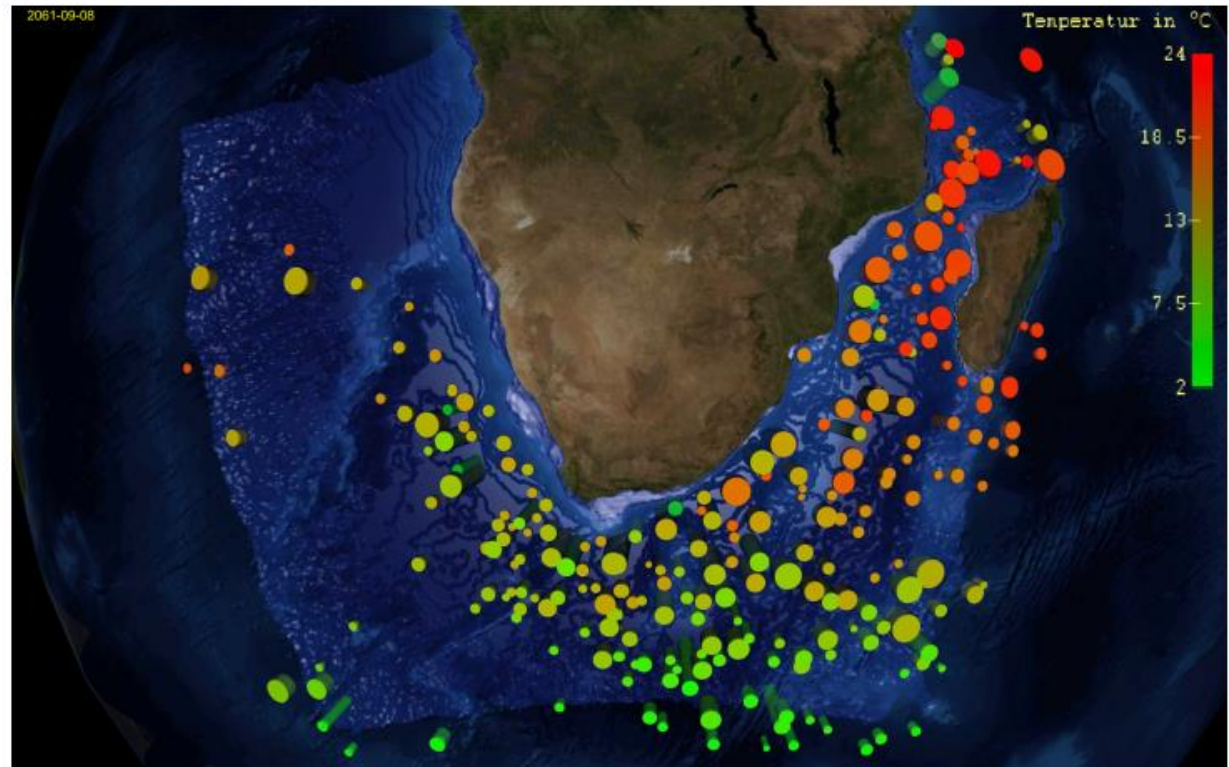
Scheuermann

UNIVERSITÄT LEIPZIG

# Big Data



## Temperatur der Wirbel



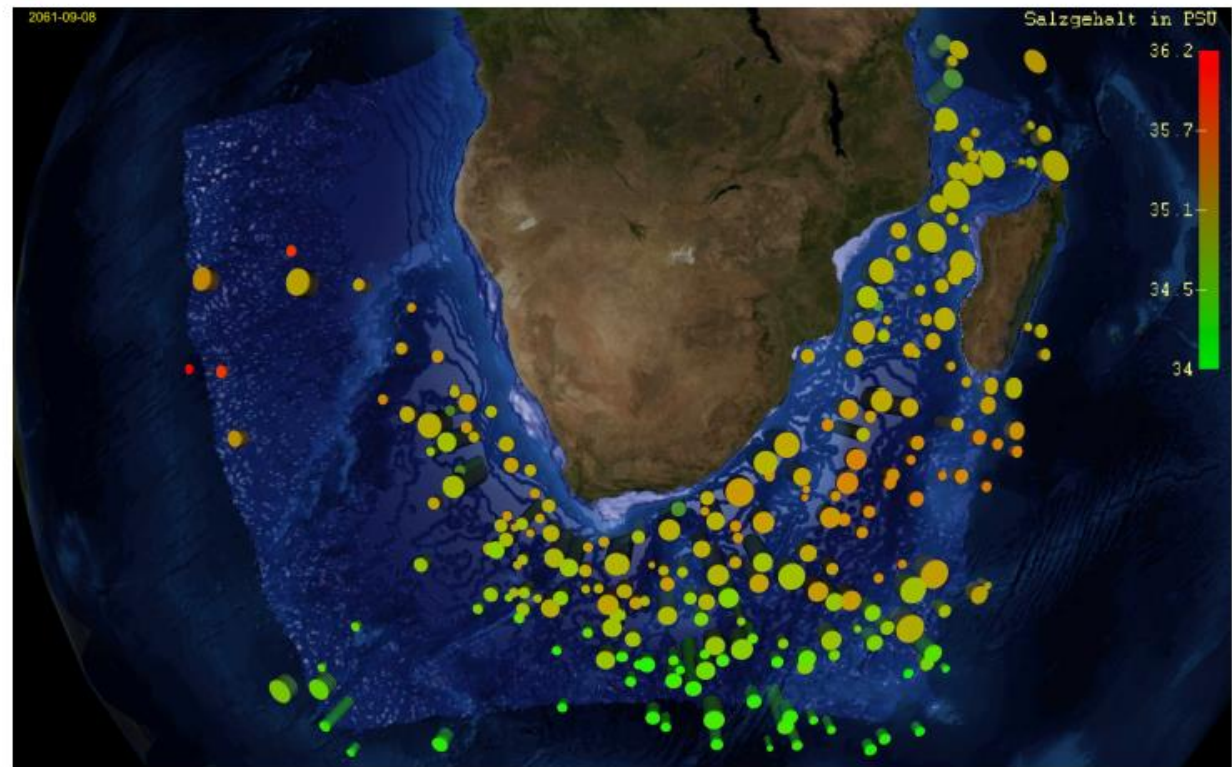
Scheuermann

UNIVERSITÄT LEIPZIG

# Big Data



## Salzgehalt der Wirbel



Scheuermann

UNIVERSITÄT LEIPZIG

# Big Data



## Ausblick: Datenmengen

### **Bisherige Arbeit:**

- 1 Jahr (365 Zeitschritte), kleiner Ozeanteil (75 GB)

### **Aktuelle Arbeit:**

- Analyse über längeren Zeitraum (200 Jahre)
- Analyse über dem ganzen Globus (ca. 50-fache Fläche)
- Mehr Variablen (Biomasse, Gasgehalt)

⇒ 10.000-20.000 fache Datenmenge, etwa 1 Petabyte

⇒ Verteilte Analyse als Notwendigkeit

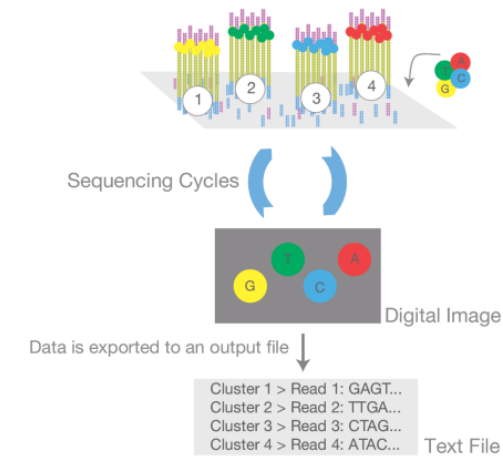
### **Zukünftige Arbeit:**

⇒ In-situ Analyse als Ziel (mehr Zeitschritte, nochmals Faktor 24)

# Big Data

## Next Generation Sequencing

### C. Sequencing



Sequencing reagents, including fluorescently labeled nucleotides, are added and the first base is incorporated. The flow cell is imaged and the emission from each cluster is recorded. The emission wavelength and intensity are used to identify the base. This cycle is repeated “n” times to create a read length of “n” bases.

### D. Alignment & Data Analysis

Reads

```
ATGGCATTGCAATTTGACAT
TGGCATTGCAATTTG
AGATGGTATTG
GATGGCATTGCAA
GCATTGCAATTTGAC
ATGGCATTGCAATT
AGATGGCATTGCAATTTG
```

Reference Genome

```
AGATGGTATTGCAATTTGACAT
```

Reads are aligned to a reference sequence with bioinformatics software. After alignment, differences between the reference genome and the newly sequenced reads can be identified.



# Big Data

## The Data Deluge

- Genome size (human) 3 Gb
- Transcripts:
  - ~ 20,000 genes
  - $10^6 \dots 10^7$  RNA products (crude estimate)
- Sequencing run (Illumina HiSeq 2500)  $6 \times 10^{11}$  characters in  $6 \times 10^9$  reads
- larger experiments:  $\gg 1000$  samples  $\Rightarrow \sim 10^{15}$  characters, i.e., in the Petabyte range

# Hintergrund

- Canonical Text Services Projekt in Leipzig
  - Dissertationsthema: Implementation and Evaluation of the Canonical Text Service Protocol as Part of a Research Infrastructure in the Digital Humanities
  - Webservice für persistente Textzitation über Projektgrenzen hinweg
  - Verschiedenste frei verfügbare Korpora
    - Deutsches Textarchiv, Textgrid, Voices of the Holocaust, Grimms Märchen, Perseus, Parallel Bible Corpus,....
- <http://cts.informatik.uni-leipzig.de/>

# Hintergrund



- Scalable Data Solutions (ScaDS) in Leipzig
  - Zusammenarbeit verschiedener Forschungsgruppen in Dresden & Leipzig
    - Visualisierung, Datenbanken, NLP, DH
  - Ziel: Aufbau eines nationalen Kompetenzzentrums für Big Data
  - CTS eine der entwickelten Lösungen
- <https://www.scads.de/>

# Motivation

- Eine der Hauptfragen in ScaDS:
- „Warum ist dein Projekt Big Data?“

# Motivation

- Eine der Hauptfragen in ScaDS:
- „Warum ist dein Projekt Big Data?“
- Offensichtliche Antwort: „Weil ich mit Terabytes an Daten arbeite!“

# Motivation

- Eine der Hauptfragen in ScaDS:
- „Warum ist dein Projekt Big Data?“
- Offensichtliche Antwort: „Weil ich mit Terabytes an Daten arbeite!“
- Mein Problem: Selbst große Textkorpora nur ein paar Gigabyte groß

# Beispielanwendung CTS Index

Datensatz	Dokumente	Statsche Referenzen (CTS URNs)	HDD Größe Index/Daten (MB)	Sprache (n)
Textgrid	91133	7284050	2730/2621	Deu
Deutsches Text Archiv	8190	21911559	7841/12887	Deu
Perseus	1693	273126	490/65	Grc, Lat, Farsi, Eng
Parallel Bible Corpus	831	8818067	3330/4972	Ca 800
TED Talk Subttle Transripts	52987	15292408	2642/2507	105

Textkorpora idR thematisch/zeitlich/... zusammenhängend

Analysen idR nur über wenige Dokumentenparametrisierungen

Trendanalyse benötigt bspw. nur Dokumente in 1 Sprache

Eher dezentrale als zentrale Datenmenge

→ Abgesehen von (eher unüblichen) zentralen Archiven ist Volume für Textdaten (Primärdaten) kein Problem

# The Big Vs

- Volume      ->      Umfang der Daten
  - Velocity    ->      Verarbeitungs/Berechnungsgeschwindigkeit
  - Veracity    ->      Verlässlichkeit der Informationen/Fehlerrate
  - Variety     ->      Heterogenität von Daten/Tools/,...
- 
- Je nach Anwendung:
    - Visualization, Value, Volatility, Vulnerability, Validity, Variability, (Quality)



# 4 Big V in DH

- Volume

- Primärdaten vs Sekundärdaten
    - Text vs. Annotation/Metadata
  - Primärdaten generell unproblematisch bezogen auf Volume

- Sekundärdaten können beliebig aufgebläht werden
    - Volume –Bezug projektabhängig

Daten im ASV-Wortschatz	Primärdaten (Bytes) (Sätze)	Sekundärdaten (Bytes) (alles andere)
deu_mixed_2011	37, 270, 576, 048	517, 020, 294, 364
deu_news_2011	3, 672, 898, 564	59, 421, 534, 187
deu_newscrawl_2011	3, 735, 178, 336	222, 879, 231, 073

# 4 Big V in DH

- Velocity
  - Verarbeitungsgeschwindigkeit
  - Streamingbasierte Workflows
    - (bspw. real time Twitter Analyse)
    - Frameworks wie Apache Flink, Spark, ...
  - Zeitintensive Berechnungen (bspw. Zitatanalyse, Topic Modeling)
    - Parallelisierung
    - Computer Cluster

# Velocity - Parallelisierung

- Auf mehreren Ebenen durchführbar
- Dokumente, Kapitel, Sätze,...
- Über Daten oder Workflows oder Quellcode
- Erlaubt Beschleunigung über Computer-Cluster oder Threads
- Parallelisierung erzeugt immer Overhead (bspw. Verwaltungsaufwand)
- Beschleunigungsfaktor sollte etwa dem Parallelisierungsfaktor entsprechen

# Velocity - Paralleles Wortzählen (Flink)

```
public class WordCountExample {
    public static void main(String[] args) throws Exception {
        final ExecutionEnvironment env = ExecutionEnvironment.getExecutionEnvironment();

        DataSet<String> text = env.fromElements(
            "Who's there?",
            "I think I hear them. Stand, ho! Who's there?");

        DataSet<Tuple2<String, Integer>> wordCounts = text
            .flatMap(new LineSplitter())
            .groupBy(0)
            .sum(1);

        wordCounts.print();

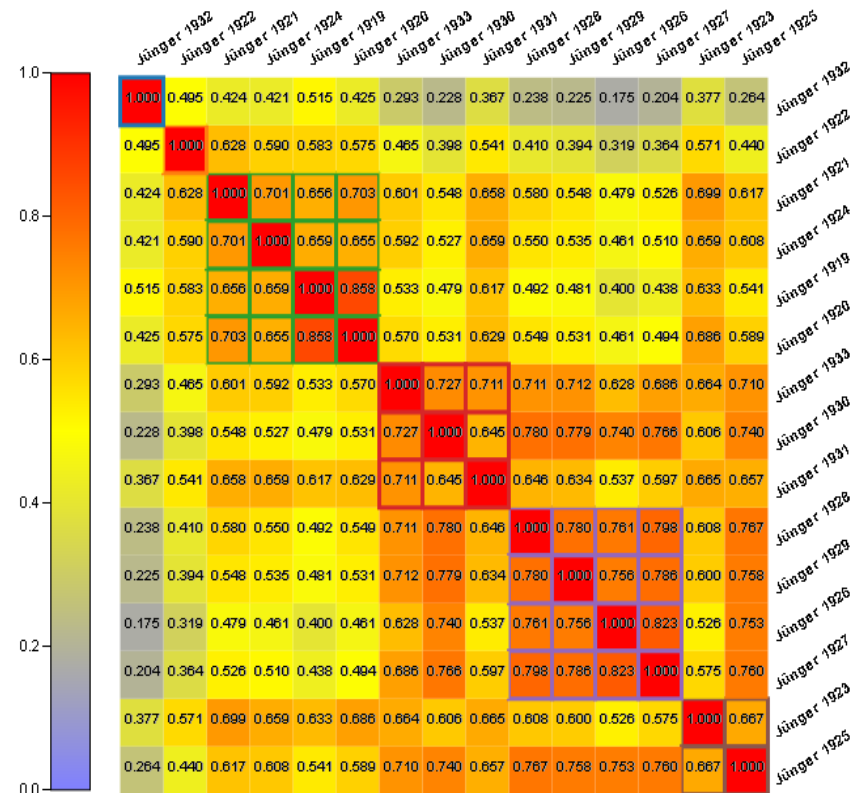
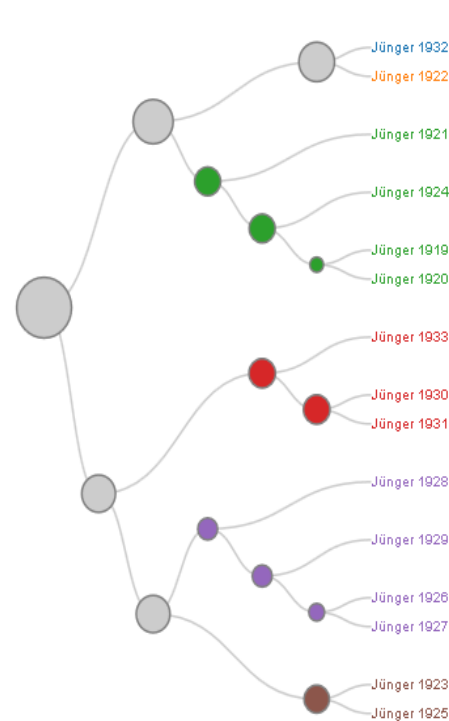
        env.execute("Word Count Example");
    }

    public static class LineSplitter implements FlatMapFunction<String, Tuple2<String, Integer>> {
        @Override
        public void flatMap(String line, Collector<Tuple2<String, Integer>> out) {
            for (String word : line.split(" ")) {
                out.collect(new Tuple2<String, Integer>(word, 1));
            }
        }
    }
}
```

Hier werden die 1en pro Wort gruppiert und zusammenaddiert

Hier wird jedes Wort mit einer 1 gesammelt

# Velocity - Vergleichende Korpusanalyse



More frequent in Jünger 1920

Word

Feuer

standen

zurück

zusammen

kurzen

Schatten

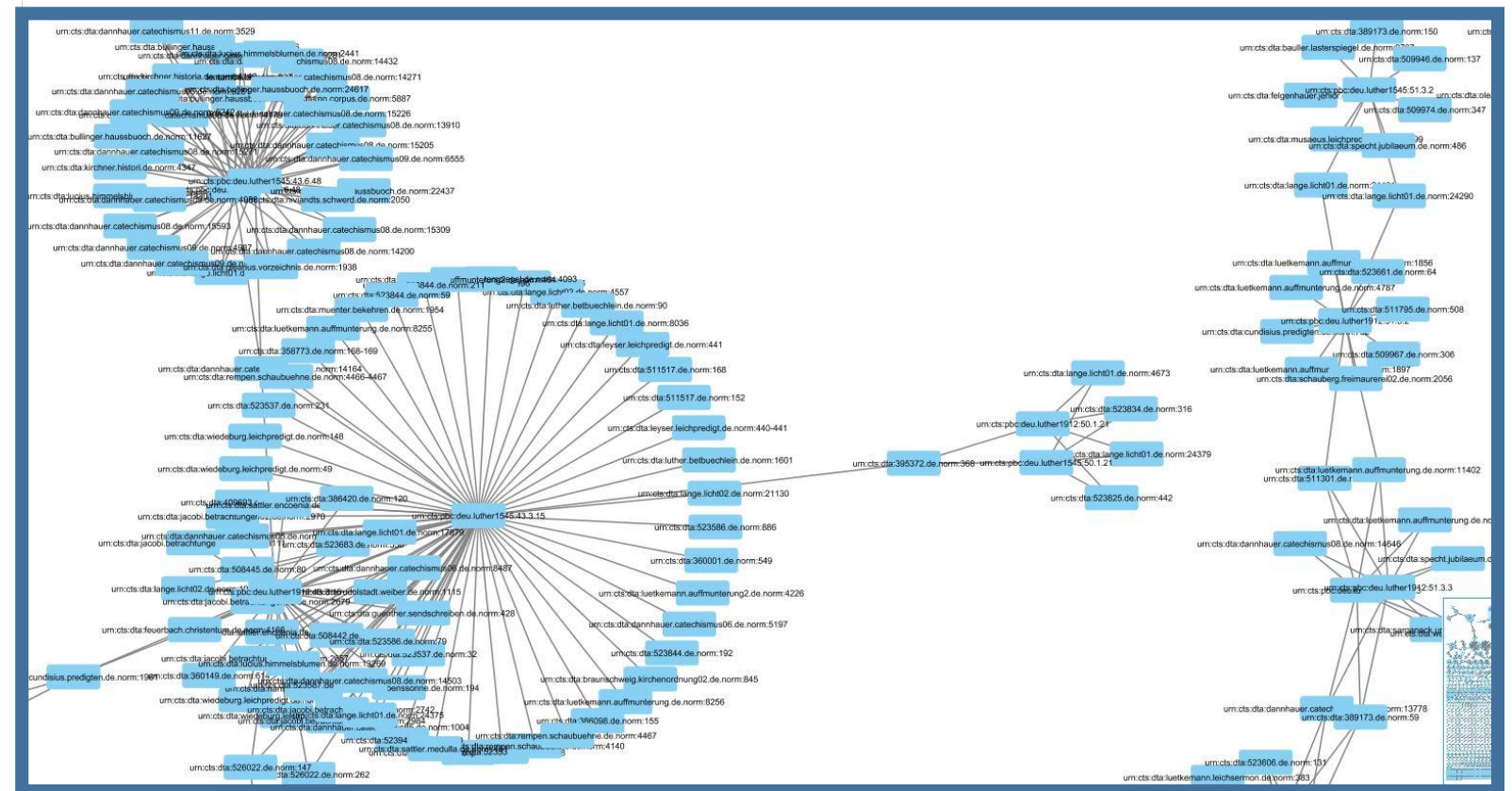
Kein

springt

größeres

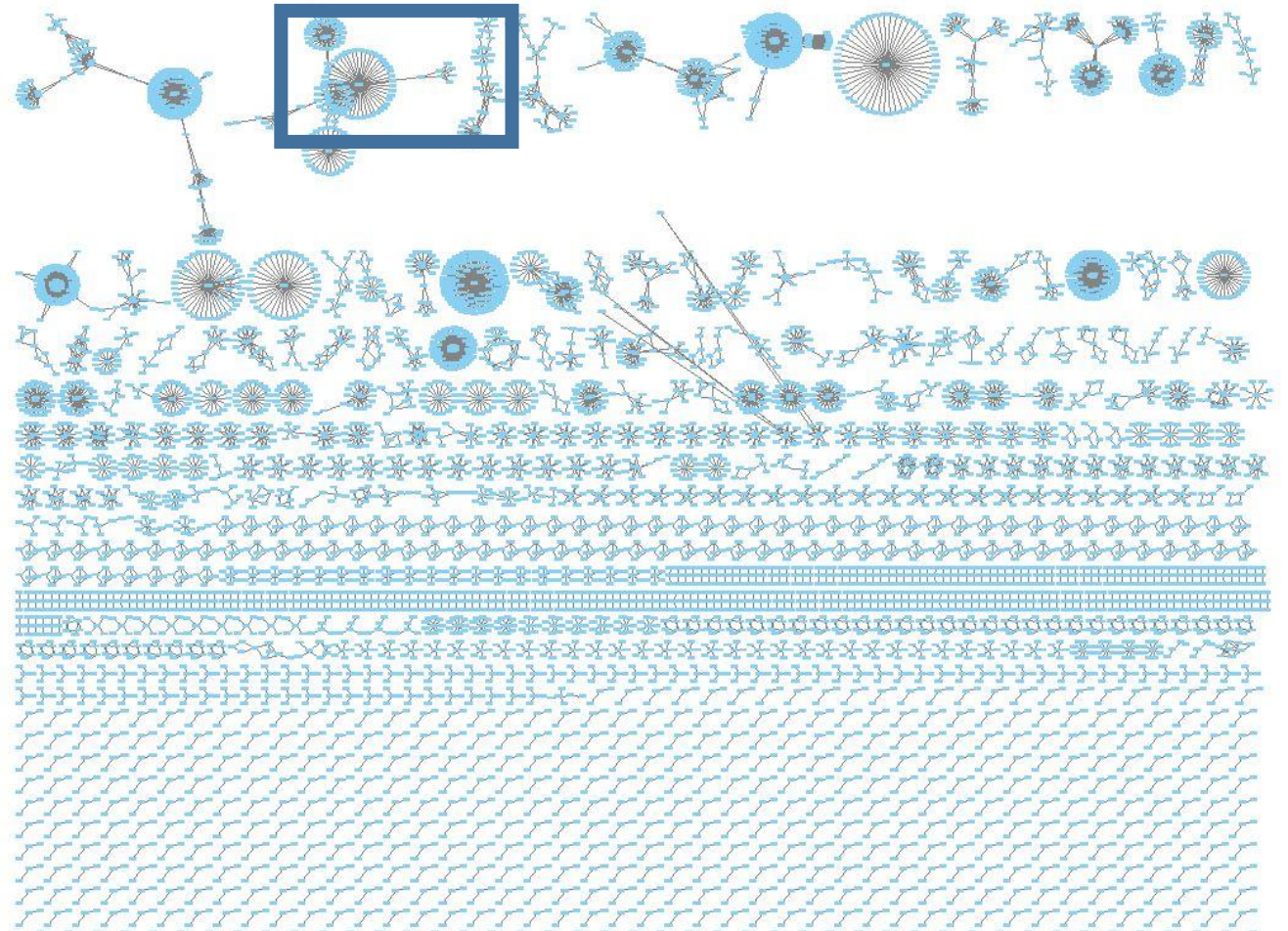
überragende

# Velocity - Text Reuse/Zitationsanalyse 1





# Velocity - Text Reuse/Zitationsanalyse 1



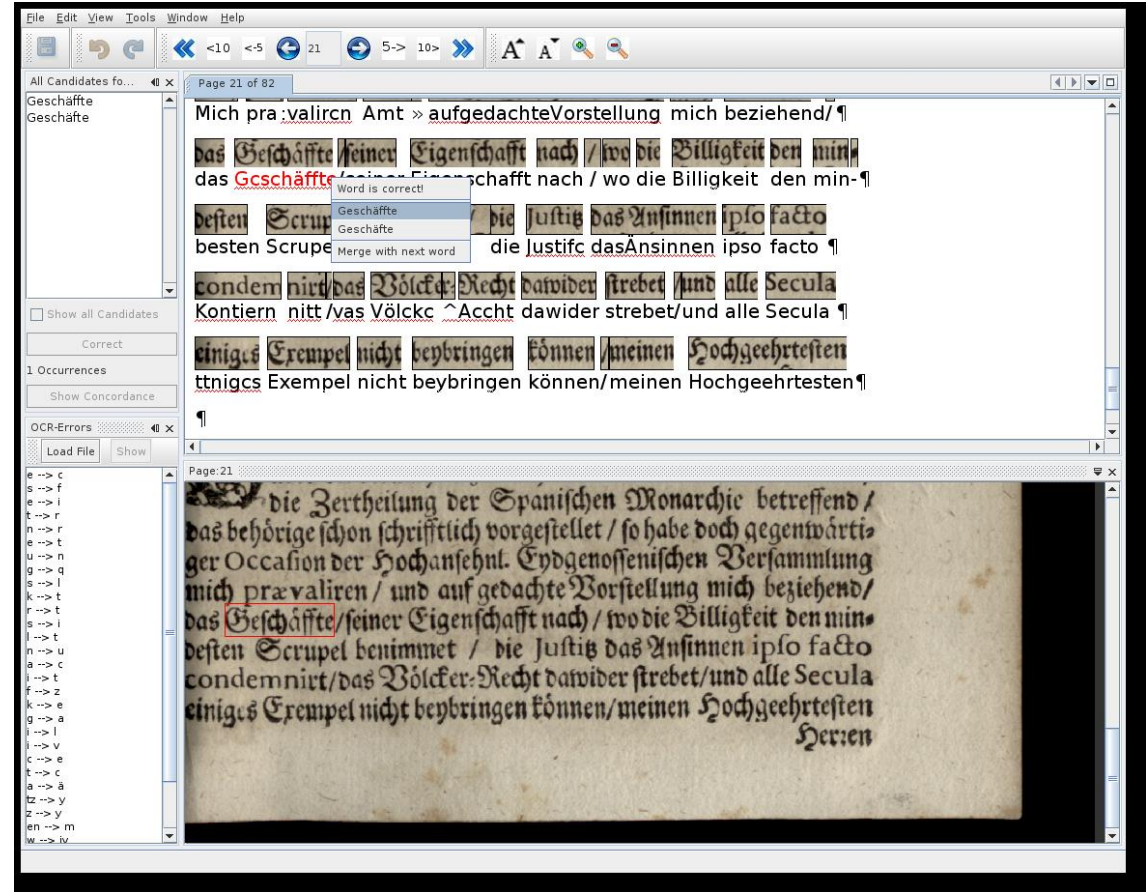
# 4 Big V in DH

- Veracity
  - Verlässlichkeit der Daten/Fehlerrate
  - OCR Korrektur automatisch/manuell
  - Schrifterkennung in Bildern
    - Machine Learning, Neural Networks
  - Normalisierung von Text
    - Punctuation, Rechtschreibfehler vs Korpuspezifische Abweichungen (zeitl. Kontext)



# Veracity - Schrifterkennung 1

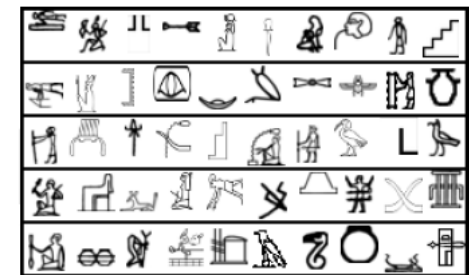
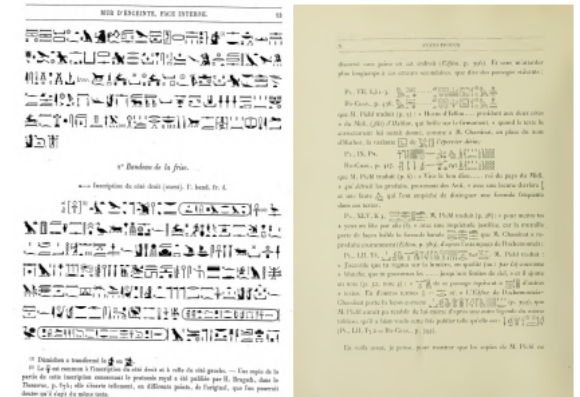
Interaktive OCR-Plattform  
*PoCoTo*



# Veracity - Schrifterkennung 2

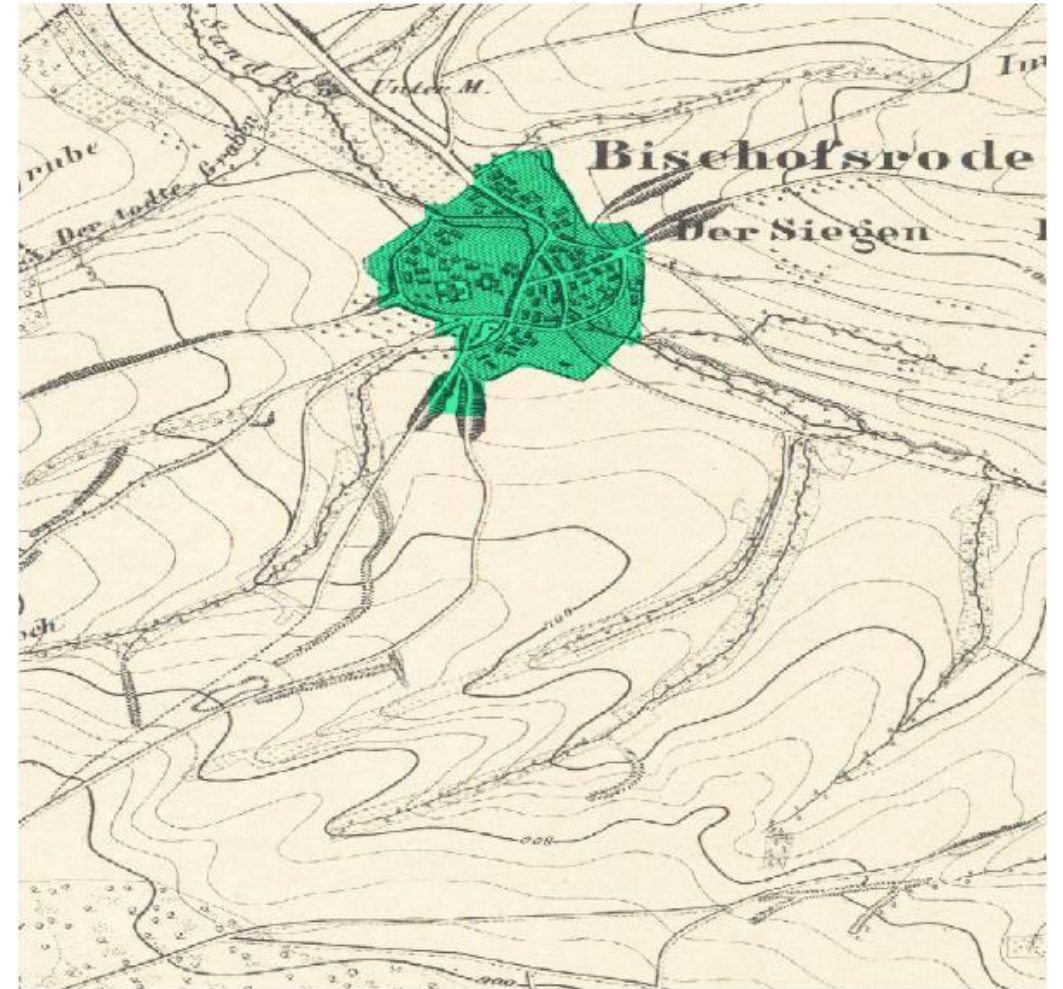
# Erkennung ägyptischer Hieroglyphen

- Ziel: Erkennung der Hieroglyphen in digitalisierten Büchern der frühen Ägyptologie
- 6832 unterschiedliche Hieroglyphen
- Test der Objekterkennung an Sequenzen mit jeweils 10 zufällig ausgewählten Hieroglyphen



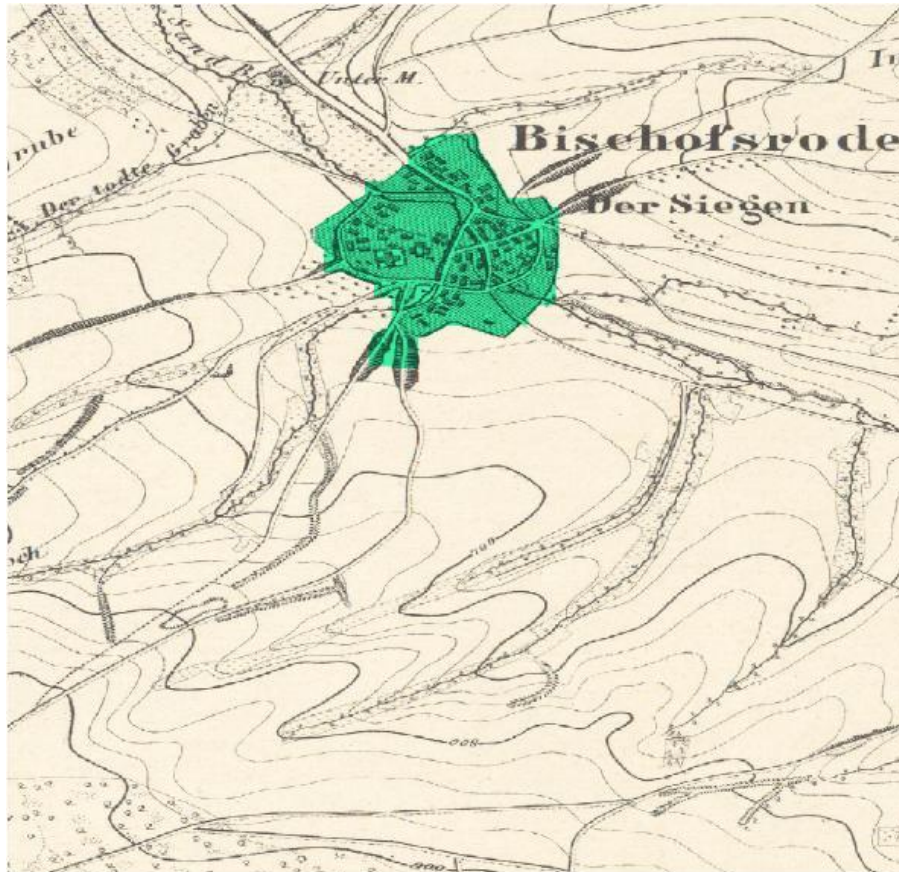
# Veracity - Ortsnamenerkennung

- Ebenfalls OCR:
- Kartenbeschriftung erkennen





# Veracity - Mustererkennung



# 4 Big V in DH

- Variety
  - Heterogenität & Interoperabilität
  - Dateiformate
    - TEI/XML, Docbook, ...
  - Datenzugang
    - Online vs Offline
    - Zip-Pakete, Git, Projektspezifische Lösungen
  - Tools
    - Input/Import-Formate

# Variety - Forschungsinfrastrukturen

- Umgebungen (in DH) für
  - Evaluation
  - Auswahl
  - Wiederverwendung
  - Kombination
- von Tools, Daten und Metadaten

# Variety - CLARIN

- Zugang zu verschiedensten Daten & Metadaten via webservice
- Inhalte über gemeinsamen Metadatenstandart vereinheitlicht (CMDI)
- Federated Search
- Verkettete workflows via WebLicht
  - webapplications, keine locale Installation nötig
- Persistente Zitierbarkeit durch Persistent IDentifiers (PID)
- Aggregierter Inhalt aus verschiedenen Quellen als Ergebnis

# Variety - VLO



## Virtual Language Observatory

Explore the world of language resources and technology from different perspectives



[VLO](#) > [Faceted search](#) > [Search: "Leipzig"](#) > [Selections](#): Written Corpus ✕ German ✕

[Permalink](#) | [Report](#) | [Help](#)

### SEARCH

### SEARCH RESULTS

6 results Showing 1 to 6

[Wortschatz](#) Expand

Collected from newspaper texts, webcrawling, etc.: words (+frequency), cooccurrences (+graph), left/right neighbours, example sentences

Resources: | 1 other |

[deu\\_news\\_2008\\_100K](#) Expand

100.000 sentences of a German newspaper corpus based on material from 2008

Resources: | 1 text document | 1 other |

#### NARROW DOWN

Use the categories below to limit the search results to those matching the selected value(s).

- + LANGUAGE
  - German ✕**
- + COLLECTION
- RESOURCE TYPE
  - Written Corpus ✕**
- + COUNTRY
- + GENRE
- + FORMAT



# Variety - Federated Content Search (FCS)

The screenshot displays the CLARIN-D Federated Content Search (FCS) interface. At the top, the logo features a stylized network of nodes and lines. Below it, the text "CLARIN-D FEDERATED CONTENT SEARCH" is prominently displayed. A search bar contains the word "Leipzig", and a "Search" button is to its right. Navigation links for "About", "Search options", "Search results", and "Help" are provided. A toolbar includes options to "Clear", "Use WebLicht", "Export to Personal Workspace", and "Download". The search results are listed under the heading "deu\_news\_2010\_1M, ASV Leipzig". Each result entry consists of a snippet of text, a source identifier "Leipzig", and a brief description. Information icons are visible on the right side of each result row.

Search Results
My Sustainable World wird den in Leipzig aufgenommen, fachübergreifenden Dialog in den Bereichen Energy, Living und Mobility sowie Finanzierung nachhaltiger Konzepte fortführen.
Auch die zweite Mannschaft von Erzgebirge Aue konnte den Siegeszug von RB Leipzig nicht stoppen und unterlag zu Hause mit 0:1.
Michael Boris, der erst seit drei Tagen das Zepter als Trainer bei Schalke II schwingt, ging im Vorfeld mit seinen Jungs hart zu Werke: "Gestern in Leipzig haben wir kein Spiel gewonnen und nur drei Tore erzielt.
In Richtung Leipzig ist die Autobahn allerdings befahrbar.
So wandte sich kürzlich der Lehrer einer neunten Klasse aus Leipzig an das NDC.
Der Sprecher von Rasenballsport Leipzig, Hans-Georg Felder, sagte, nach der Umgestaltung werde zu erkennen sein, dass es sich um das Stadion von RB Leipzig handle.
Aktuell werden in Leipzig parallel zum BMW X1 der 3-Türer, das Coupé und das Cabrio der BMW 1er Rei
Die CDU im Landtag sprach von einem Affront gegen die Messestadt Leipzig.

# Variety - Federated Content Search (FCS)

The screenshot displays the CLARIN-D Federated Content Search (FCS) interface. At the top, the logo features a stylized network of blue nodes and lines, followed by the text "CLARIN-D FEDERATED CONTENT SEARCH". Below this is a search bar containing the text "Leipzig" and a blue "Search" button. Navigation links for "About", "Search options", "Search results", and "Help" are positioned below the search bar. A toolbar at the top of the results area includes a "Clear" button, a "Use WebLicht" dropdown, an "Export to Personal Workspace" button, and a "Download" dropdown. The results are organized into two sections, each with a header bar.

**TIGER, IMS, Universität Stuttgart**

Mit der Rheinbraun Verkaufsgesellschaft wurden drei Vertriebstöchter in Leipzig , Weimar und Potsdam gegründet , die eine flächendeckende Versorgung mit leichtem Heizöl , Diesel und Schmierstoffen gewährleisten sollen .

Inzwischen hat sich der Schleier über Leipzig und Zwickau , der Lausitz oder dem ' ' schwarzen Dreieck ' ' bei Zittau gelichtet .

**Dingler Online, Berlin-Brandenburg Academy of Sciences and Humanities**

Leipzig im Juli 1820.

Vor ungefähr 15 Jahren brachten es die Franzosen unter dem Namen "Chinesische Schminkblätter," (Rouge en feuille) das erstmal auf die Messe nach Leipzig von wo es sich allmählich auch dem östlichen und nördlichen Europa verbreitete.














# Variety - WebLicht

- Webservice zur Verkettung und Ausführung von Workflows
  - Workflowübergreifende Metabeschreibung
    - Kompatibilität zwischen Ausgabe und Eingabe verschiedener Elemente
    - CMDI

Input and Chain Selection

Run Tools

<b>Title [Plain Text]</b> Als ich in den Krieg zog, freute ich mich, wie wohl jeder junge Mensch in jenen Tagen, auf das, was vor uns ...	<b>SfS: To TCF Converter</b> Language: German Document Type: TCF TCF Version: 0.4 Text	<b>SfS: Tokenizer/Sentences -</b> newlinebounds <input type="text" value="false"/> Sentences Tokens	<b>SfS: POS Tagger - OpenNLP</b> Part of Speech: STTS Tagset	<b>Berlin: Tokens2Lexicon</b> Language: German Document Type: Lexicon Format TCF Version: 0.4 entries.type: types	<b>Berlin: dlexDB Types: Frequency</b> frequencies.type: absolute
	 	 	 	 	 

# Variety - Interoperabilität

- Austauschbarkeit und Wiederverwendbarkeit von Daten und Workflows
- Erlaubt „blindes“ Erstellen von Tools & Daten
- Erfordert übergeordnete, anwendungsunabhängige Regeln (Metaebene) für Interface & Access
- Reduziert Variety
- Für textbasierte DH: Canonical Text Services (nächste Woche)

# Zusammenfassung

- Datenvolumen ist nur 1 Aspekt von Big Data
  - Komplexe Probleme können auch mit kleineren Datensätzen aufwendig werden
- Big Data Herausforderungen können sein:
  - Texterkennung, Schrifterkennung, Fehlererkennung, Normalisierung
  - Zeitanforderungen, Durchsatz, in sich komplexe Berechnungen
  - Umgang mit Annotation, Metadata, Meta-metadata,...
  - Interoperabilität
- Vor allem in DH ist die Bestimmung von “Big Data-igkeit” einer bestimmten Fragestellung nicht trivial und universell übertragbar