## Slide 1

eAQUA

ASV

**An infrastructure for eHumanities**
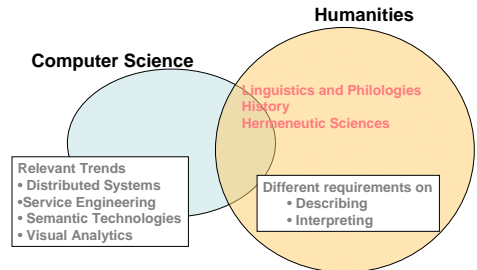
*Workshop on Historical texts*
Boston, 2010/01/13

Gerhard Heyer
(using slides by Marco Büchler, Volker Boehlke, and Charlotte Schubert)
Automatische Sprachverarbeitung
Computer Science Department
University of Leipzig

## Slide 2

eAQUA

- Computer Science and Humanities
- eAQUA
- eHumanities, Digital Humanities, and Humanities
- eHumanities and eScience Infrastructures

Gerhard Heyer

2

## Slide 3

eAQUA — Computer Science and Humanities

**Humanities**

**Computer Science**

Linguistics and Philologies
History
Hermeneutic Sciences

**Relevant Trends**
• Distributed Systems
• Service Engineering
• Semantic Technologies
• Visual Analytics

**Different requirements on**
• Describing
• Interpreting

3

Gerhard Heyer

## Slide 4

eAQUA — Computer Science and its applications

| 1940-1960 | Scientific Computing |
|---|---|
| 70ies | Data bases, digitizing business processes (Wirtschaftsinformatik) |
| 80ies | Digitizing electro mechanical applications, SGML |
| 90ies | Digitizing analogue media, linking distributed ressources: http, HTML, XML |
| since 2000 | Internet based services, knowledge management |

4

## Slide 5
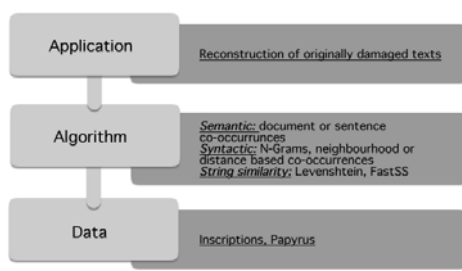
eAQUA — Computer Science and its applications

- Replacement of analogue by *digital* media and processes
- Increasing impact of digital media and proces-sing models on traditional work flows based on analogue media
- Digitizing media and work flows creates new methods and applications
- *eHumanities:* Digitizing media and work flows in the Humanities
- BMBF project eAQUA a good example

5

Gerhard Heyer

## Slide 6

eAQUA — eAQUA Methodology

**Application** — Reconstruction of originally damaged texts

**Algorithm** — *Semantic:* document or sentence co-occurrences
*Syntactic:* N-Grams, neighbourhood or distance based co-occurrences
*String similarity:* Levenshtein, FastSS

**Data** — Inscriptions, Papyrus

Gerhard Heyer

6

## Slide 7

- **Atthidographers** – Classification of Atthidographers
- **Plato** - Aftermath of Plato's writings
- **Plautine metric** – Metrical analysis of Latin comedies
- **CAMENA** - Knowledge network of Latin CAMENA corpus
- **Inscriptions** - Extraction of significant templates for different kinds of inscriptions like release documents (slave trading)
- **Papyri** - Classification of papyri (e. g. slave trading contracts) and text completion of fragmentary texts
- **Mental maps** - Building of century based mental maps and highlighting differences

Course Text Mining for Classical Studies (TM4CS): http://www.eaqua.net/e_3.html

## Slide 8

### Challenges from Comuter Science point of view

- **Integrating textual ressources**
  Conversion of different formats and standards

- **Software engineering issues**
  Getting classicists to understand the intent and effects of text mining algorithms and clearly define their functional requirements

- **Semantic technologies that address the long tail**
  Modifying and amending statistical and pattern based methods to effectively deal with rare events

## Slide (Example subproject Atthidographs - left)

Searching for „Ἀτθίδος"

Wort: *Ἀτθίδος ( 19575 )*
Anzahl: *276*
Häufigkeitsklasse: *14*
Normalisiert gleiche Formen: *Ἀτθίδος (276); Ἀτθίδος (2); ατθίδος (1);*

Wörter mit gleicher Grundform: *Ἀτθίδος (276); Ἀτθίς (82); Ἀτθίδα (73); Ἀτθίς (47); Ἀτθίς (32); Ἀτθίδων (26); Ἀτθίδας (120); Ἀτθί (4); Ἀτθίων (2); Ἀτθίδος (2); Ἀτθίων (2); Ἀτθίδες (2); Ἀτθίθ (1); Ἀτθίθ (1); Ἀτθί (1); Ἀτθίθ (1); Ἀτθίθ (1);*

Kookkurrenzähnliche Formen: *κατηγόμενος (0.2295); σηγκαρωσήσεις (0.2222); σπείφρ (0.2222); ὑπαρχάίαιν (0.2222); Ἀνδροτίων (0.22); Φιλόχορος (0.21); γυτής (0.1985); Ἑλλάνικος (0.19); Πραξίων (0.1849); Σπαράδα (0.1849); Σπάρανος (0.1803); Μεγαρικόν (0.1707); Ἑλετιατον (0.1571); ἀετή (0.15); FGeHia (0.14); Θαίτογνος (0.14); Σπήρων (0.1389); Δικηγίδας (0.1306); Ἀπίων (0.1304); ἱστοριαί (0.13); FHG (0.13); Ἔφρορος (0.13); ἰσμόν (0.13); Παντακλέσς (0.125); Τσαιηφαίον (0.123); Ἡρόδοτος (0.12); γλιτσίν (0.12); κσίματος (0.12); τῶμα (0.12); μελίσισης (0.12); Ἑλλάνικός (0.12); Φιλίων (0.12); γωφρσίσης (0.12); Μοννίγον (0.119); Ἱποθωντίδης (0.1146); Εἴδολος (0.11); κσκήσθσ (0.11); Φιρσκίνης (0.11); Κλείδημος (0.11); Πειρηγήσσας (0.1045); Χρσναλῶν (0.1019); τωίσσον (0.1); FGeH (0.1); ἱστορίαίον (0.1); ἱπίτιγ (0.1); Ἀθηγάιν (0.1); Ἀπολλίδσωρος (0.1); Θσυκιδίδης (0.1); FI (0.1); Ἀτθίθ (0.1); Πσιφίσσίν (0.1); μάντισς (0.1);*

## Slide (Example subproject Atthidographs - right)

Looking at co-occurrences to find „interesting" neighbours

Signifikante rechte Kookkumenzen für *'ατθίδος*

- νομάδας **is an interesting candidate**
- **Reference to quotation from the Atthis of Philochorus (3. / 2. century BC) in a normal word search hard to find**
- **Co-occurrence indicates that atthidographers may have invented a phase in Athenian history to show an evolutionary development from** *nomadism* **to settled life**

## Slide 11

**There also is an infrastructure aspect …**

## Slide 12

- **What we need is a transition from accidental to organized cooperation**
- *Centers of competence* **might mediate technical and organisational issues**
- **A platform of** *services* **allows all participants to share digital ressources and establisch a culture of** *best practices* **(cf. CLARIN)**

## Slide 13

- *Services* **can be implemented as webservices for sharing data and algorithms (SOAP)**
- **Every participant needs to provide his data and algorithms in a standardized form (WSDL)**
- **Webservices allow to**
  - **reuse data from aggregated results**
  - **generate specific annotated data irrespective of the location of the resource**



| Normalised words | Co-occurrences | Co-occurrences based similar words | Quotations |

## Slide 14

Ancient Greek philology

Ancient history

Textual reuse

Computer science

**Algorithms**

## Slide 15

- Started in April 2004
- Client implementations in Java, .NET, Perl, Delphi, Python and PHP are known
- Currently: 18 services installed
- Since Sept. 2006: 180M accesses
  - Service chaining
  - Behaviour of users
  - Relevant topics

Source: Marco Büchler and Gerhard Heyer: Leipzig Linguistic Services - A 4 Years Summary of Providing Linguistic Web Services.
In: Gerhard Heyer (Editor):Text Mining Services – Building and applying text mining based service infrastructures in research and industry.,
Leipzig, Germany, 2009

## Slide 16

process chain builder/validator

DSpin webservice — DSpin webservice — service repository

toolwrapper A — toolwrapper format 1 — toolwrapper B

resource A — tool b

## Slide 17

application layer — Webservice A, Webservice B, WebLicht, XML

configuration layer — WS registry management tool, XML, registry webservices

persistence layer — ISOcat registry, WS registry database

## Slide 18

## Slide 1 (top-left)

- Chaining consideren an instance of functional composition: c(b(a()))
- Basically the same problem to be solved like on the type checking level of a compile run: Check if all input needs of a certain function are satisfied: correct number and type of parameters (format is given by programming language).
- Our "parameters" (NLP-specific):
  - document is encoded using a certain format (TEI, …)
  - a certain concept/kind of information inside of a valid document (tokens, tags, …) => parameter
  - this concept/information is encoded using a certain datatype (utf8, tagset A, ...)

## Slide 2 (top-right)

## Slide 3 (bottom-left)

## Slide 4 (bottom-right)

- **Already done:**
  - Implementation of a basic registry for webservices (url, descriptions etc + metadata on format & input/ouput)
  - Several services of the DSpin prototype available
  - Implementation of a first version of the matchmaking algorithm
  - First successfull test in two different workflow/chaining tools (Tübingen, Leipzig)
- **Future developments:**
  - Implementation of a generic chain builder: automatically suggest a chain from startpoint (resource) to endpoint (tool, certain information, ...)
  - Open up the registry for metadata harvesting
  - Integrate other communities (like digital classics) and services