

# Introduction to NLP methods: *Overview of Possible Methodologies used in Text Retrieval*

---

Marco Büchler  
Göttingen Centre for Digital Humanities

Lyon, France  
June 2<sup>nd</sup> 2014

# A fundamental question

---

**How can the computer really support to identify lines of transmissions (text re-use) on big data?**

# A “fundamental answer”

---

## Naive method:

- Compare every text chunk (like sentence) with each other.
- TLG:  $5,500,000 * 5,500,000 = \mathbf{3.025e13}$  comparisons
- Assumption: Comparison rate of **1000 sentences/sec.**
- This process would run about **3.025e10** seconds or more than **959 years.**

# Does Big Data Analysis make any sense?

---

## Serendipity Effect

# Methodology

---

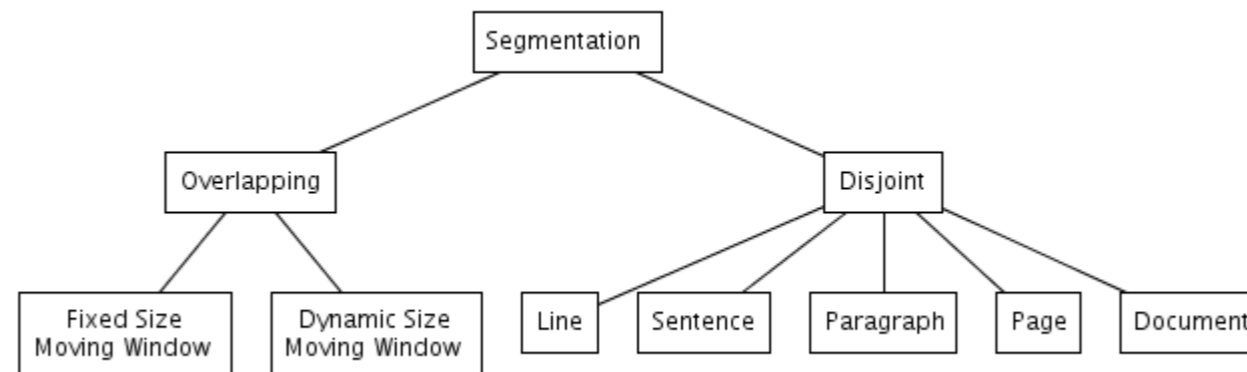
**7-Level-Architecture** to deal with *Data Diversity*:

- a. Segmentation:** Defining *Re-use Units*
- b. Preprocessing:** Cleaning of *Re-use Units*
- c. Featuring:** Creating of *Fingerprints*
- d. Selection:** Selection of *Features* from *Fingerprint*
- e. Linking:** Linking of *Re-use Units* given common Features
- f. Scoring:** Scoring of Re-use Overlaps of pairwise linked Re-use Units
- g. Postprocessing:** optional Post-Processing of the Re-use Graph

**Implemented in TRACER software:** more than a million permutations of implementations of different levels are recently possible

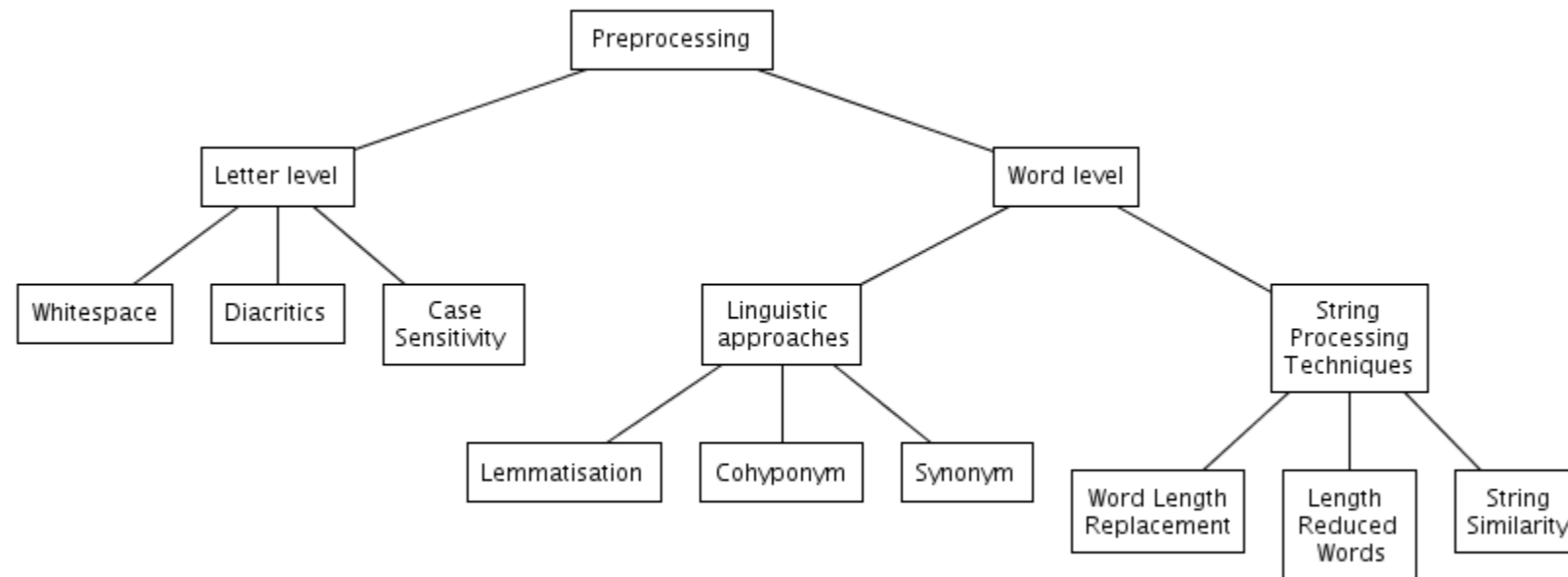
# Segmentation

---



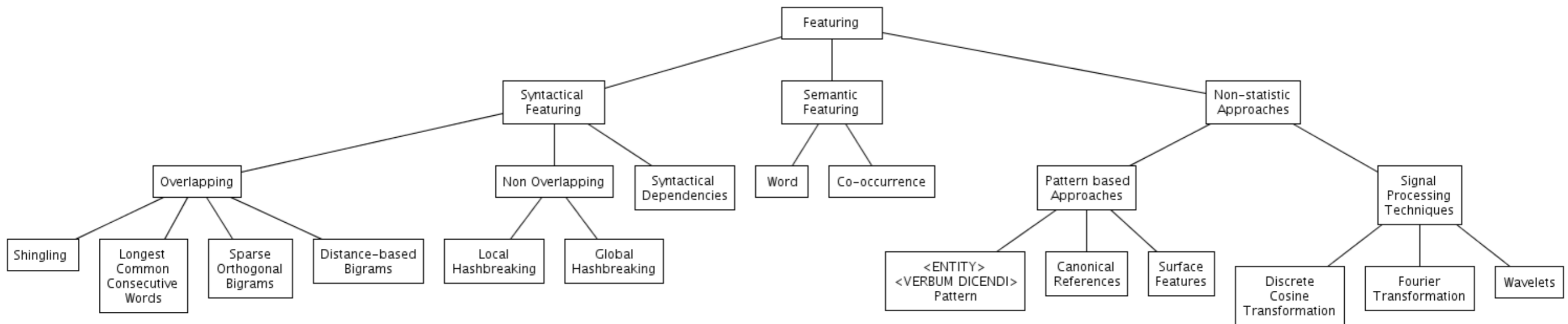
# Preprocessing

---



# Featuring

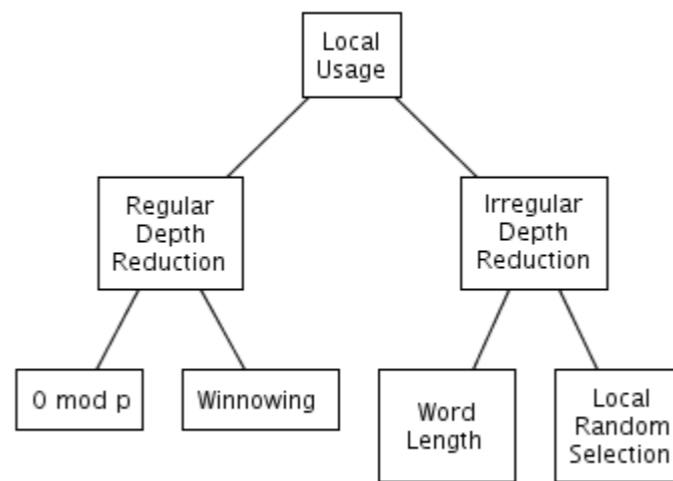
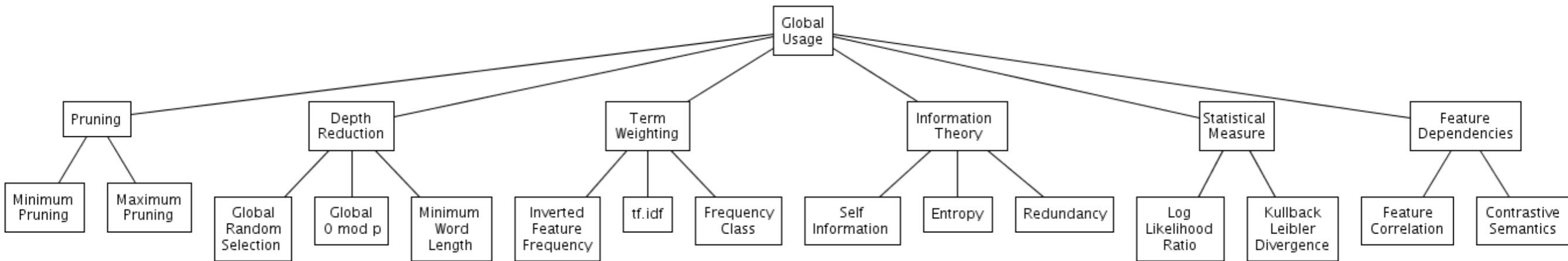
---





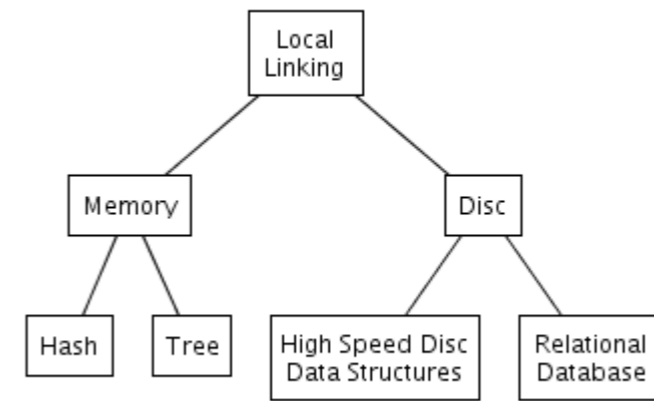
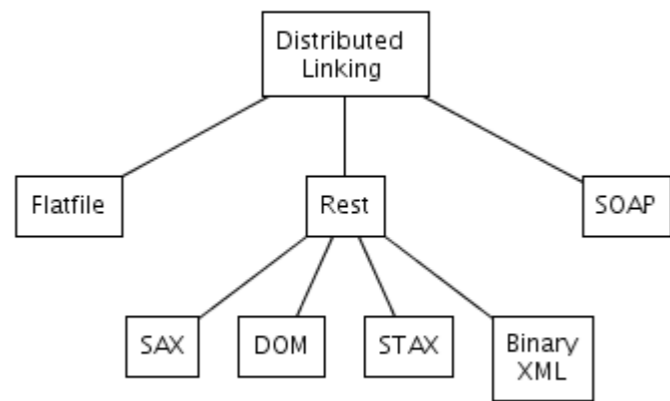
# Selection

---



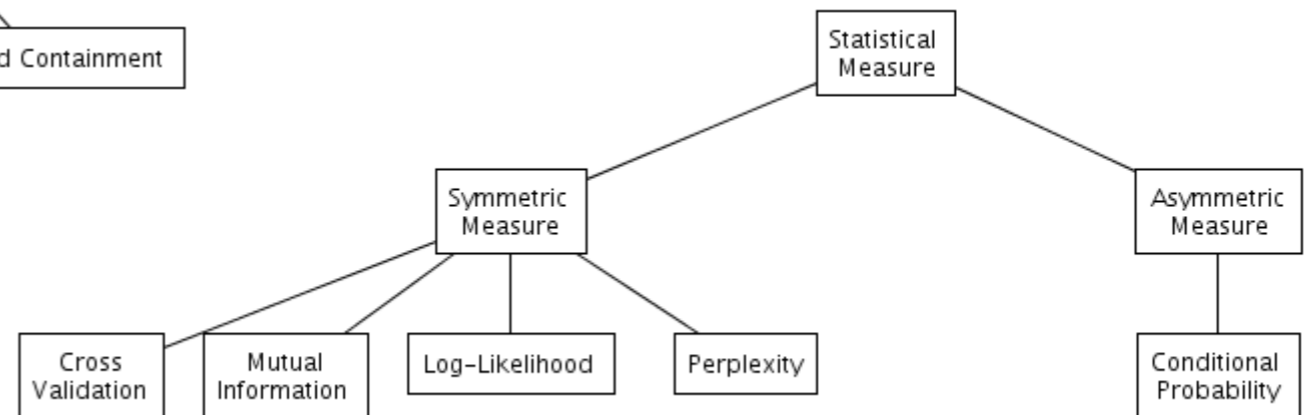
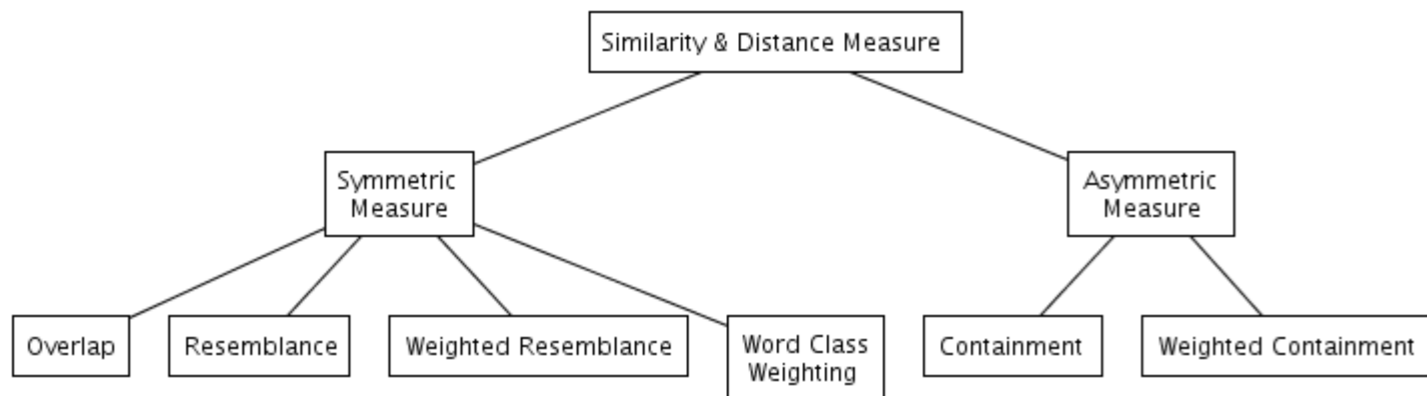
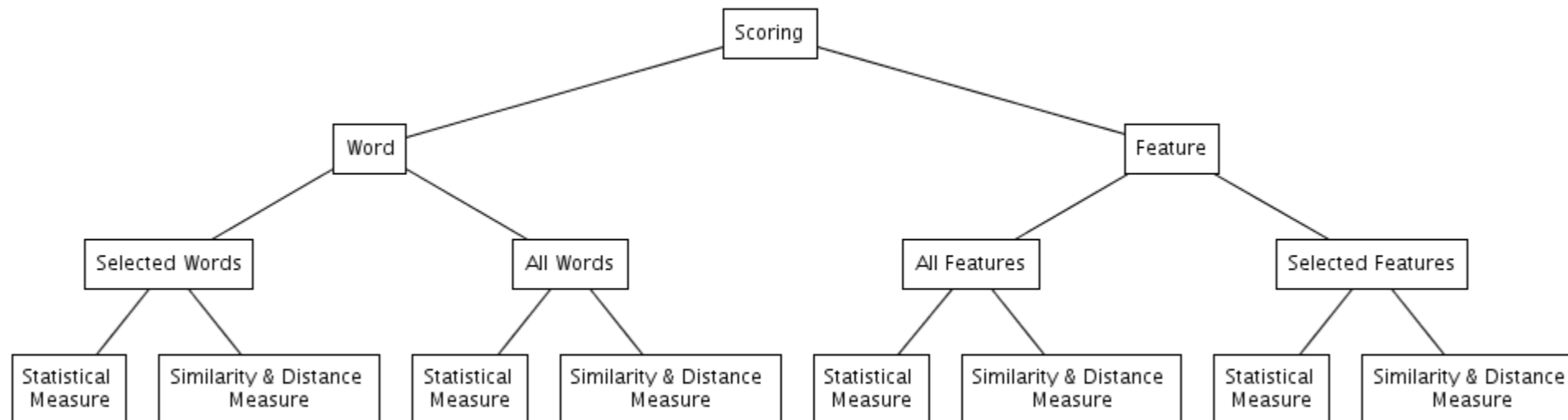
# Linking

---



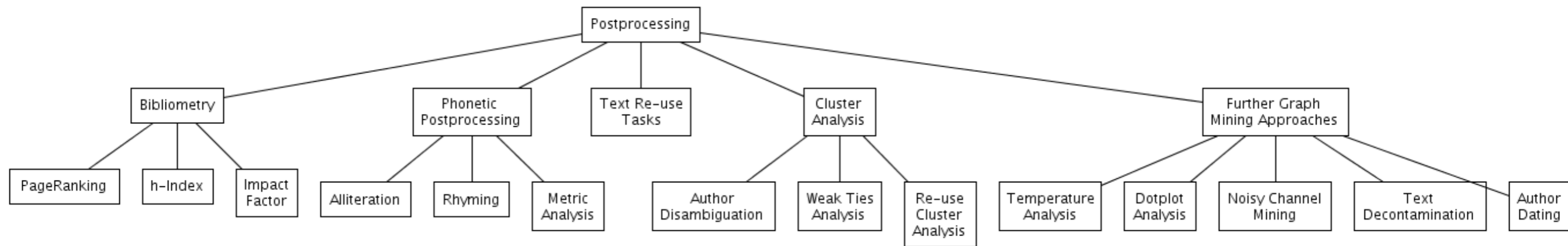
# Scoring

---



# (Postprocessing)

---



# Contacts

---

**For more details:**

<http://www.etraces.e-humanities.net>

<http://www.gcdh.de/en/>

***Google group for Historical Text Re-use:***

<http://groups.google.com/group/historical-text-re-use>

***Marco Büchler***

**Göttingen Centre for Digital Humanities**

**Georg August University Göttingen, Germany**

**[mbuechler@gcdh.de](mailto:mbuechler@gcdh.de)**