

# Automatische Extraktion von Fachterminologie aus kunsthistorischen Volltexten

Juliane Bredack

SKILL 2014

Studierendenkonferenz Informatik

25. September 2014

# Überblick

- Einführung
  - Was sind Mehrwortgruppen?
  - Lingo und die Komponenten
  - Das Reallexikon zur Deutschen Kunstgeschichte
- Praktisches Vorgehen
- Ergebnisse
  - Extraktion von Mehrwortgruppen
- Fazit

# Bedeutung von Mehrwortgruppen

- Erweiterung der Suchmöglichkeiten im Information Retrieval
- Bedeutungsdifferenzierung auf der Ebene eines einzelnen Wortes
  - Bank → Deutsche Bank
- Aussagekraft und Deutlichkeit
- Zusammenhängende lexikalische Einheit → besitzen als Ganzes eine Bedeutung

# Lingo – das Indexierungssystem

- System zur linguistisch und statistisch basierten automatischen Indexierung
- Wörterbuchgestützte Grundform-identifizierung bzw. –reduktion, Kompositazerlegung, algorithmische und lexikalische Mehrwortgruppenerkennung, Synonymrelationierung

# Programmmodul sequencer

- Algorithmische Identifizierung und Extraktion noch unbekannter Mehrwortgruppen

*Automatische Extraktion fachterminologischer Mehrwortgruppen mit Lingo.*

```
<Automatische=[(automatisch/a)]>  
<Extraktion=[(extraktion/s)]>  
<fachterminologischer=[(fachterminologisch)/a]>  
<Mehrwortgruppen=[(mehrwortgruppe/s)]>  
<mit=[(mit/c)]>  
<Lingo=[(lingo/s)]>  
:./PUNC
```

# Programmmodul sequencer

Bildung von Wortmustern, basierend auf Wortklassen →  
Identifizierung von Wortgruppen

*Automatische Extraktion fachterminologischer  
Mehrwortgruppen mit Lingo.*

*Automatische [a] Extraktion [s] fachterminologischer [a]  
Mehrwortgruppen [s]*

AS automatische extraktion

AS fachterminologische  
mehrwortgruppe

# Reallexikon zur Deutschen Kunstgeschichte - RDK

- Nachschlagewerk zur Realienkunde seit 1937
- Hohe Zahl an kunsthistorischen Sachbegriffen und Eigennamen (Personennamen, Geografika), spezifische Fachinformationen



Quelle: Flammenschwert RDK, Spalten 701, 702, Band 9

# Mehrwortgruppen im Reallexikon zur Deutschen Kunstgeschichte

- Kombinationen von Sachbegriffen  
*„Darbringung im Tempel“*
- Sachbegriffe und Personennamen  
*„Verkündigung an Maria“*
- Sachbegriffe und Geografika  
*„Flucht nach Ägypten“*
- Funktionswörter (Präpositionen, Konjunktionen, Artikel)



# Mehrwortgruppen mit Funktionswörtern im Reallexikon zur Deutschen Kunstgeschichte

- Syntaktisch abgeschlossene MWGs durch die Verbindung von Fachtermen und Funktionswörtern  
*„Szene aus dem Leben Jesu“*
- Inhaltliche Spezifizierung durch komplexere MWG-Konstruktionen  
*Leben, Szene → „dargestellte Szene aus dem Leben Christi“*
- Bedeutungsdifferenzierung einzelner Begriffe  
*Zeit → „Kunst der altchristlichen Zeit“*

# Praktisches Vorgehen – Erstellung von Wörterbüchern

- Modifikation vorhandener Wörterbücher, die kunsthistorisches Fachvokabular enthalten (Sachbegriffe, Personennamen, Geografika)
- Erstellung eines Wörterbuchs mit Funktionswörtern (Artikel, Konjunktionen, Präpositionen)
- Allgemeines Rechtschreibwörterbuch

# Praktisches Vorgehen – Vergabe neuer Wortklassen

<b>Wortklasse</b>	<b>Bedeutung</b>
E	Kunsthistorische Fachbegriffe
A	Adjektiv
K	Komposita
C	Präposition
R	Artikel
U	Konjunktion

# Praktisches Vorgehen – Definition von Extraktionsmustern

- Bildung der Extraktionsmuster:
  - Erstellung von Kriterien zur Musterbildung
  - Folgen typischen Erscheinungsformen zusammengesetzter Terme im Deutschen
  - gezielte Platzierung von Funktionswörtern:
    - Artikel vor Substantiven (RE) oder Adjektiv-Substantiv Verbindungen (RAE)
    - Funktionswörter werden nur innerhalb eines Wortmusters platziert, beginnen und enden nie mit einem Funktionswort

ERAE „*Symbol des eucharistischen Christus*“

AERAE „*berühmte Komposition der italienischen Malerei*“

# Ergebnisse der Mehrwortgruppenextraktion

- Erstellung von Kriterien zur Bewertung fachterminologischer Mehrwortgruppen zur qualitativen Beurteilung der Extraktionsmuster und daraus resultierenden Mehrwortgruppen:
  - Semantische Abgeschlossenheit
  - Spezifizierung von Personennamen und Geografika
- Indexierung von ca. 8.000 doppelspaltigen Lexikonseiten
- Vier- und fünfteilige Erkennungsmuster liefern beste Ergebnisse

## Substantiv-Funktionswort-Adjektiv-Substantiv

ERKE                      *„darstellung der himmelserscheinung maria“*

## Adjektiv-Substantiv-Funktionswort-Substantiv

AKUE                      *„weißes leinentuch als altarbekleidung“*

## Adjektiv-Substantiv-Funktionswort-Adjektiv-Substantiv

AERAE                      *„epische dichtung der höfischen zeit“*

## Wortmuster, die Personennamen oder Geografika spezifizieren

Personennamen:

ECEE                      *„altar von georg raphael“*

Geografika:

KCEE                      *„wandgemälde in der capella greca“*

# Fazit – Erfolgreiche Extraktion von Mehrwortgruppen mit Hilfe algorithmischer Verfahren

- Fachterminologische Mehrwortgruppen mit kunsthistorischem Bezug werden extrahiert
- Syntaktisch komplex strukturierte Mehrwortkonstruktionen durch Funktionswörter im Extraktionsalgorithmus
- Analyse umfangreicher Dokumentkollektionen hinsichtlich enthaltener Mehrwortgruppen möglich
- Modifizierung der Ergebnisse durch Korrektur der Wörterbücher und/oder Wortklassen möglich
- Allerdings: Analyse sehr zeitintensiv, Fachwissen zur Kunstgeschichte notwendig

# Kontakt

Juliane Bredack

[julianebredack@web.de](mailto:julianebredack@web.de)