

KNOWLEDGE-BASED LINGUISTIC ANNOTATION OF DIGITAL CULTURAL HERITAGE COLLECTION

Tuukka Ruotsalo, Lora Aroyo and Guus Schreiber

Speaker: Chenhua
Date: 24th Feb 2010

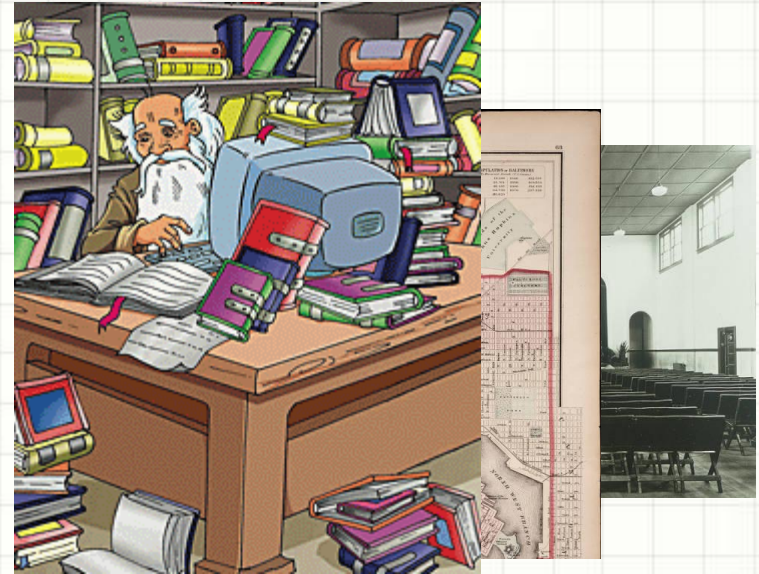
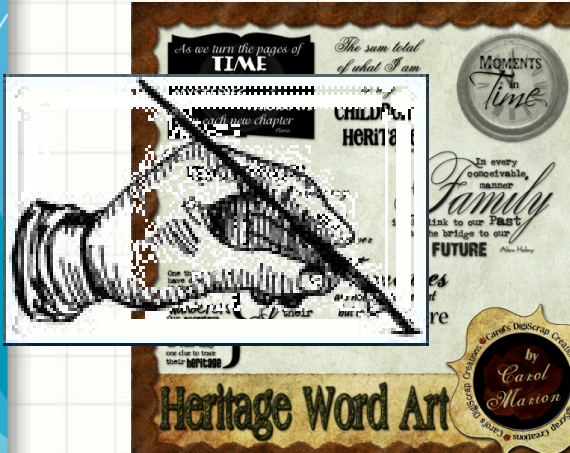
Outline

- Introduction
- Motivation
- Methodology
- Experimental Results
- Conclusion

Introduction

- Paris was painted in 1888.
- In Paris, Van Gogh painted the work in 1888.

Motivation



Better run ...



Research Question

**Is there a smart way to
annotate such massive collection?**



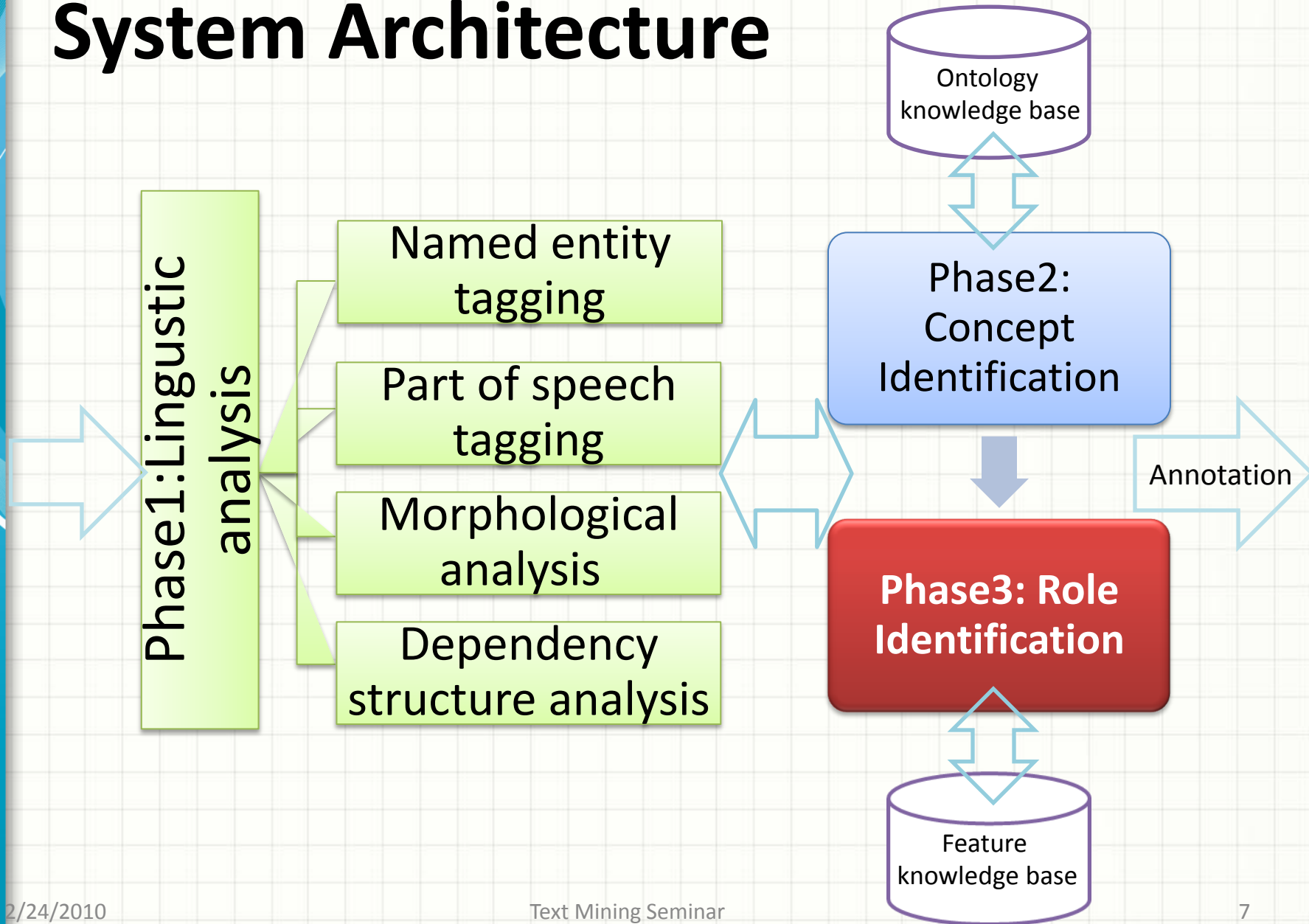
Methodology

- Background knowledge
 - Structured vocabulary
 - Enhance performance of retrieval
- Automatic annotation
 - Concept identification
e.g. Paris as a city
 - Role identification
e.g. Paris as a subject matter



The queen is portrayed in her full regalia. Everything here emphasizes her royal status: the crown, her ermine robes and the canopy. This official portrait was made in the studio of the painter Frans Pourbus II. It is a copy of the portrait in the Louvre in Paris. The queen is Marie de Médicis (1573–1642), a member of the renowned Italian family and wife of the King of France, Henry IV (1553–1610), whom she married in 1600. Henry IV was an ally of the Dutch Republic in the struggle against Spain. In 1638 Marie de Médicis, now a widow, visited Amsterdam, where she was received with great ceremony. For the occasion Joachim van Sandrart painted a militia painting in which she is also portrayed.

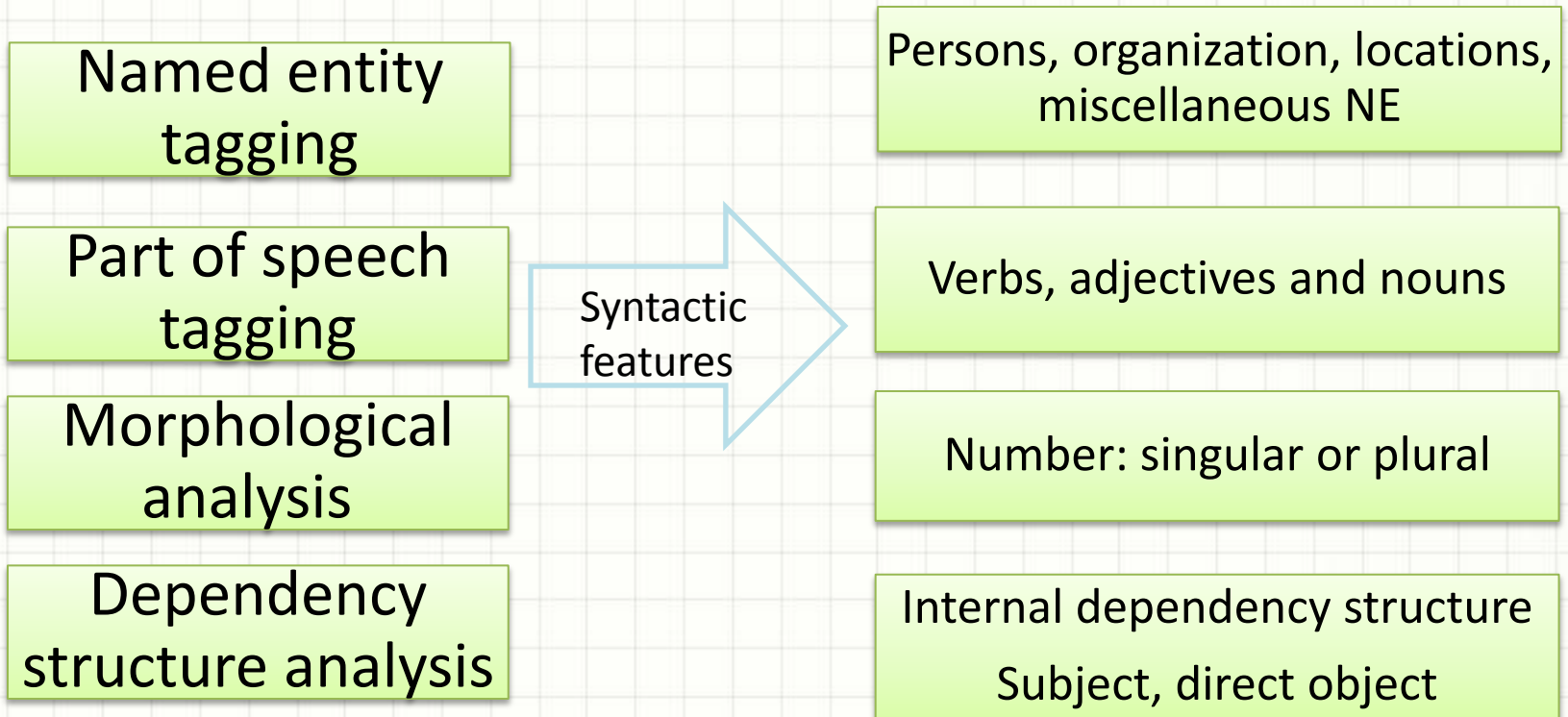
System Architecture



Knowledge Base

- Art and Architecture Thesaurus (AAT)
- Getty Thesaurus of Geographic (TGN)
- Union List of Artist Names (ULAN)
- WordNet
- etc.

Linguistic Analysis



Concept Identification

- Define(chunking) and map meaningful units to concepts in structured vocabularies
- Perform differently for nouns, verbs and NE's



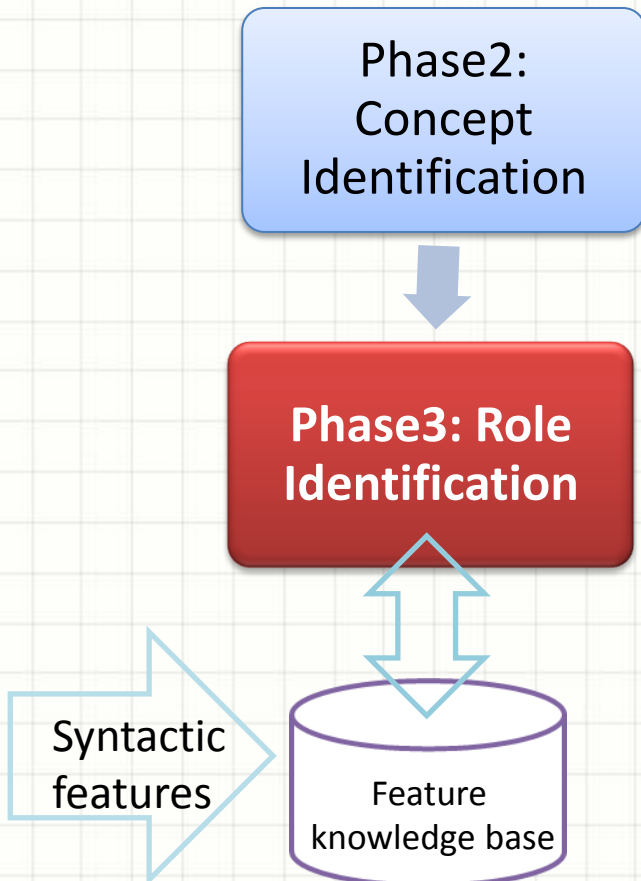
Mapping chunks, NE's, bi-words to KB

Examples for matching NEs:
NE tagged with persons
ULAN → others → WordNet

Syntactic features

Phase2: Concept Identification

Role Identification



- Difference between concept and role identification
 - “Rembrandt” is an instance of concept “person”, independent of context
 - “Rembrandt” can take various role , e.g, creator or subject of artworks, dependent of context
- How to do role identification task?
 - SVM
 - Based on features:
 - syntactic and semantic
 - E.g. PoS tag, Voice of a sentence verb, PoS path parsing constituent to verb or predicate

Evaluation

- Using a collection of natural language descriptions of artworks.
 - ARIA collection from Rijksmuseum Amsterdam
 - 250 artworks randomly selected
 - Typical descriptions on “what, who, where, when and which people or culture related to the artworks
- Using 3 structured vocabularies (Knowledge Base)
 - AAT, TGN,ULAN and WordNet
- Using an artwork annotation schema
 - Visual Resources Association(VRA) specialized on artwork

Evaluation (Cont.)

- Comparison

- Compare it to the performance of human annotators (overall accuracy?)
- Compare it to the performance of base-line method (Benefit from knowledge base?)

- Comparison Criteria

- F1 measure: $F1 = \frac{1}{\alpha Precision} + \frac{1}{(1-\alpha) Recall}$ ($\alpha = 0.5$)
- The higher F1 value, the better.

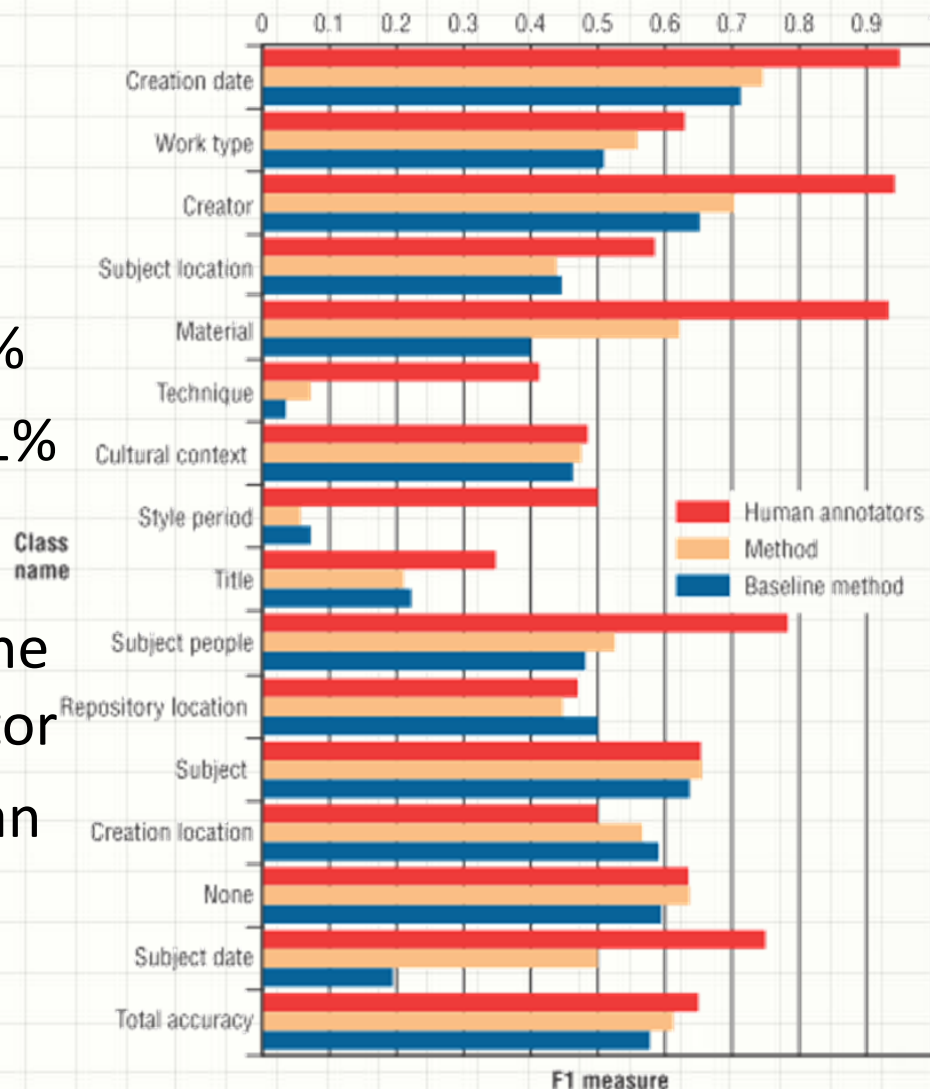
Experimental Results

- **Accuracy**

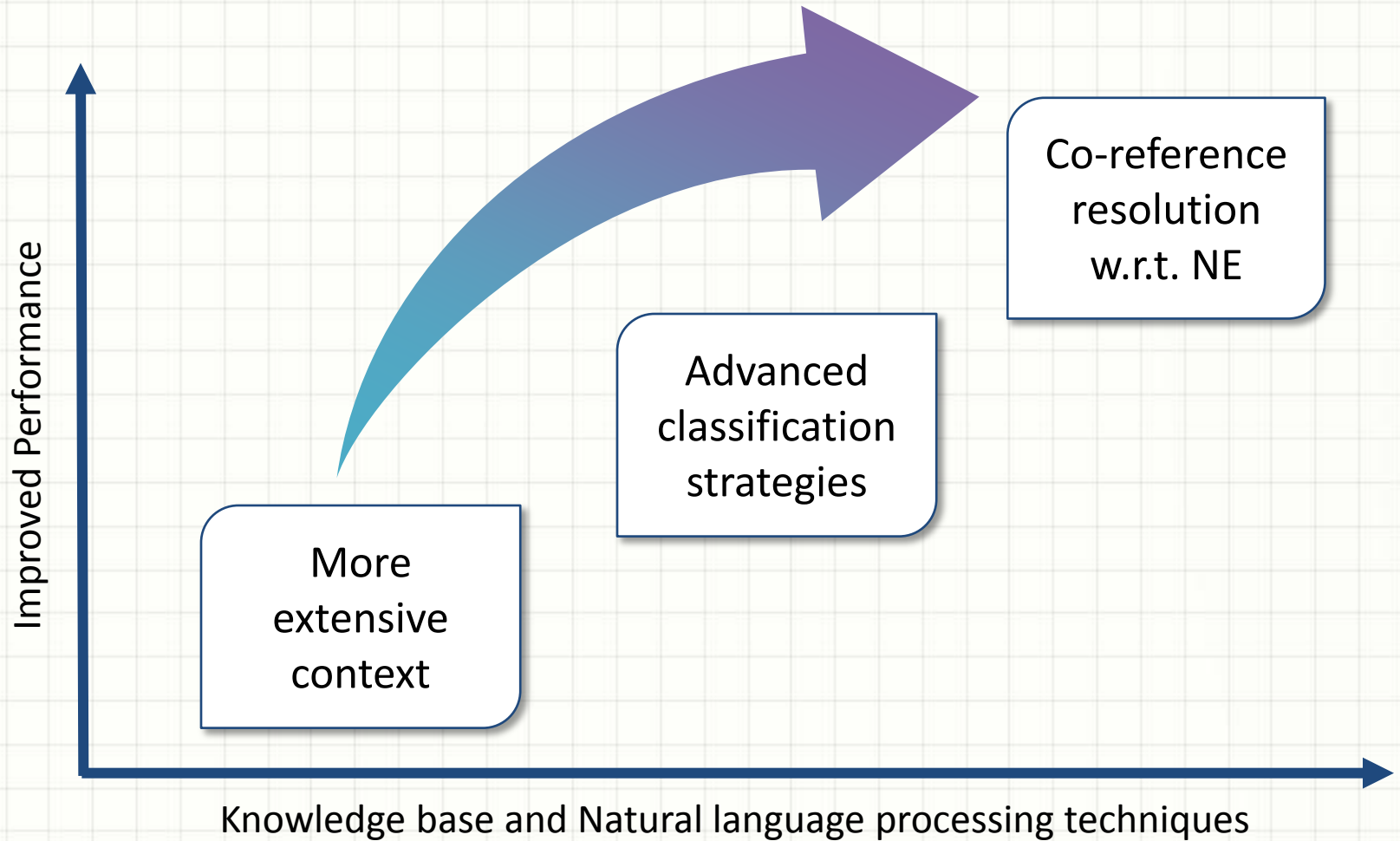
- 61.2%
- Baseline method: 57.8%
- Human Annotator: 65.1%

- **Discussion**

- Performance close to the level of human annotator
- Performance better than baseline method

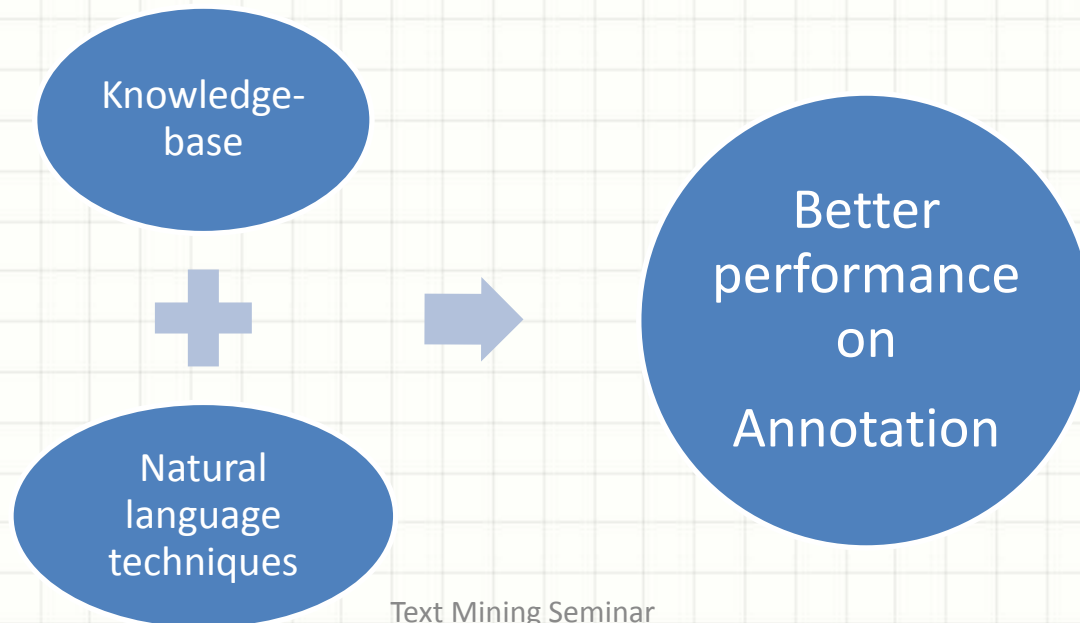


Further Discussions & Future Work



Summary

- Given a set of objects each accompanied by a text description, a set of structured vocabularies, a metadata schema, and a training set of annotations of the text descriptions, the method automatically produces annotations for the objects, and its performance is close to the level of human annotator.





THANKS!



APPENDIX

metadata

Table 1. Dublin Core Visual Resources Association metadata schema roles selected for annotation.

Role	Explanation	Value range*
Work type	Specific type of artwork being described.	AAT, WN
Title	Title or identifying phrase given to an artwork.	Literal
Material	Substance of which an artwork is composed.	AAT, WN
Technique	Production or manufacturing processes, techniques, and methods incorporated in the fabrication or alteration of the artwork.	AAT, WN
Creator	Names, appellations, or other identifiers assigned to an individual, group, or corporate body that has contributed to the design, creation, production, manufacture, or alteration of the artwork.	ULAN, Literal
Creation date	Date or range of dates associated with the creation, design, or production of the artwork.	Literal
Repository location	Geographic location and/or name of the repository locations entity whose boundaries include the artwork.	Literal
Creation location	Geographic location and/or name of the creation locations entity whose boundaries include the artwork.	TGN
Style period	A defined style, historical period, group, school, dynasty, movement, etc., whose characteristics are represented in the artwork.	AAT, WN
Cultural context	Name of the culture, people (ethnonym), or adjectival form of a country name from which an image originates, or the cultural context with which the artwork has been associated.	AAT, WN
Subject term	Terms or phrases that describe, identify, or interpret the artwork and what it depicts or expresses. These include generic terms that describe the work and the elements that it comprises.	AAT, WN
Subject people	Terms or phrases that describe, identify, or interpret particular people.	ULAN, WN, Literal
Subject location	Terms or phrases that describe, identify, or interpret geographic places.	TGN, WN
Subject date	Terms or phrases that describe, identify, or interpret time.	Literal

*Acronyms: AAT = Art and Architecture Thesaurus, TGN = Thesaurus of Geographic Names, ULAN = Union of Artist Names, WN = WordNet.

Feature knowledge base

Features for the constituent “regalia”:

1. Verb identifier: <http://www.w3.org/2006/03/wn/wn20/instances/synset-portray-verb-4>
2. Passive voice: true
3. Position before verb: false
4. Constituent identifier: <http://www.w3.org/2006/03/wn/wn20/instances/synset-regalia-noun-1>,
<http://e-culture.multimediant.nl/ns/getty/aat#300185696>
5. Constituent PoS: NN
6. Partial PoS path: IN
7. Partial Dependency Path: prep-in
8. Constituent Ontology Base: WN, AAT
9. Ontology root: <http://www.w3.org/2006/03/wn/wn20/instances/synset-artifact-noun-1>,
<http://e-culture.multimediant.nl/ns/getty/aat#30018711>,
<http://e-culture.multimediant.nl/ns/getty/aat#300264092>
10. Constituent Word Type: Noun