# OCAS: Ontology-Based Corpus and Annotation Scheme

Alexander Grothkast, University of Kaiserslautern
Benjamin Adrian, DFKI
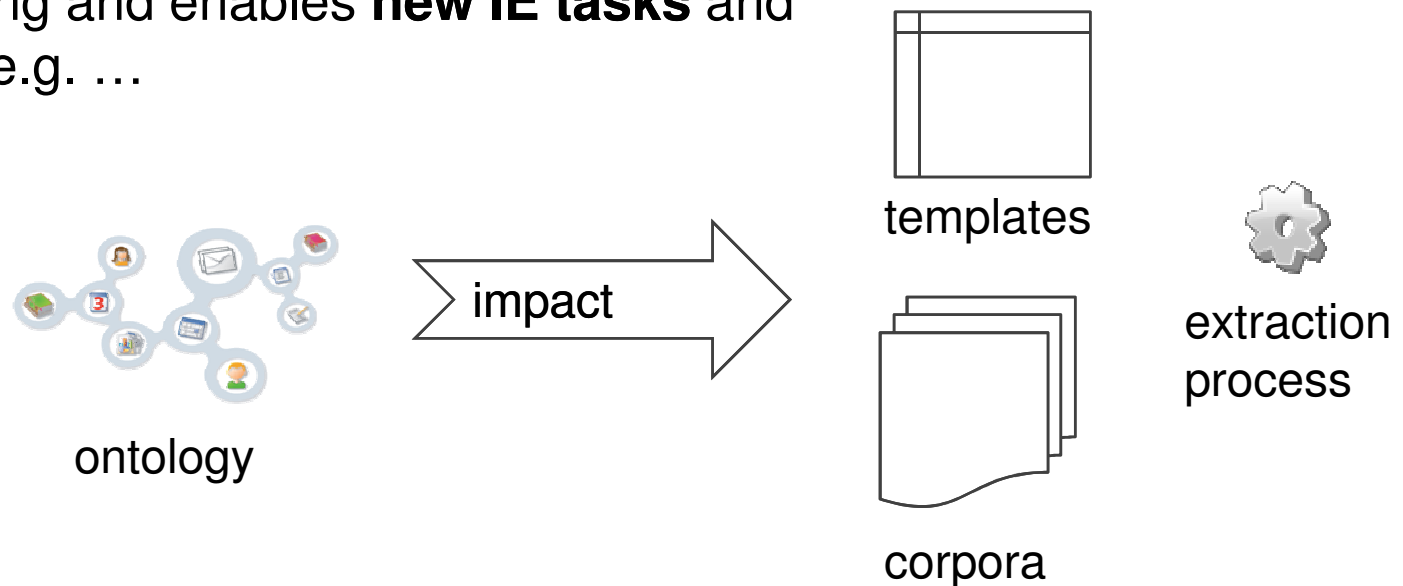Kinga Schuhmacher, DFKI
Andreas Dengel, DFKI

# Outline

1. Evaluation of Ontology-based Information Extraction

2. The Process of an Ontology-based Corpus and Annotation Scheme

3. A Review about the OCAS 2008 corpus

4. Summary

5. Outlook

# Evaluating OBIE systems

OBIES: In IE, **Ontologies** are used for enhancing results by the use of **symbolic domain knowledge**.

Compared to traditional IE, Ontology-based IE changes existing and enables **new IE tasks** and architectures, e.g. …



templates
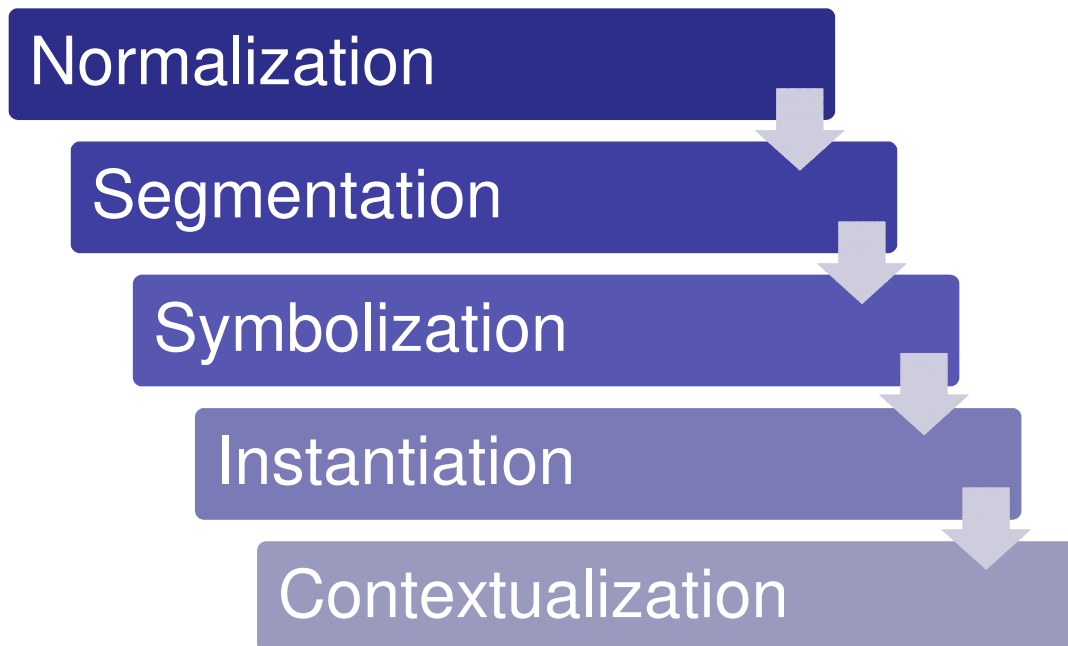
impact

ontology

extraction process

corpora

Traditional IE evaluation methods do not suffice. (Maynard, AHM 2005)
How do we evaluate an ontology-based IE system?

# Generic architecture of an OBIE system

A comparison of different OBIE systems needs a shared view on architectures.

The **OBIE layer cake** abstracts from a common OBIE process.

Normalization

Segmentation

Symbolization

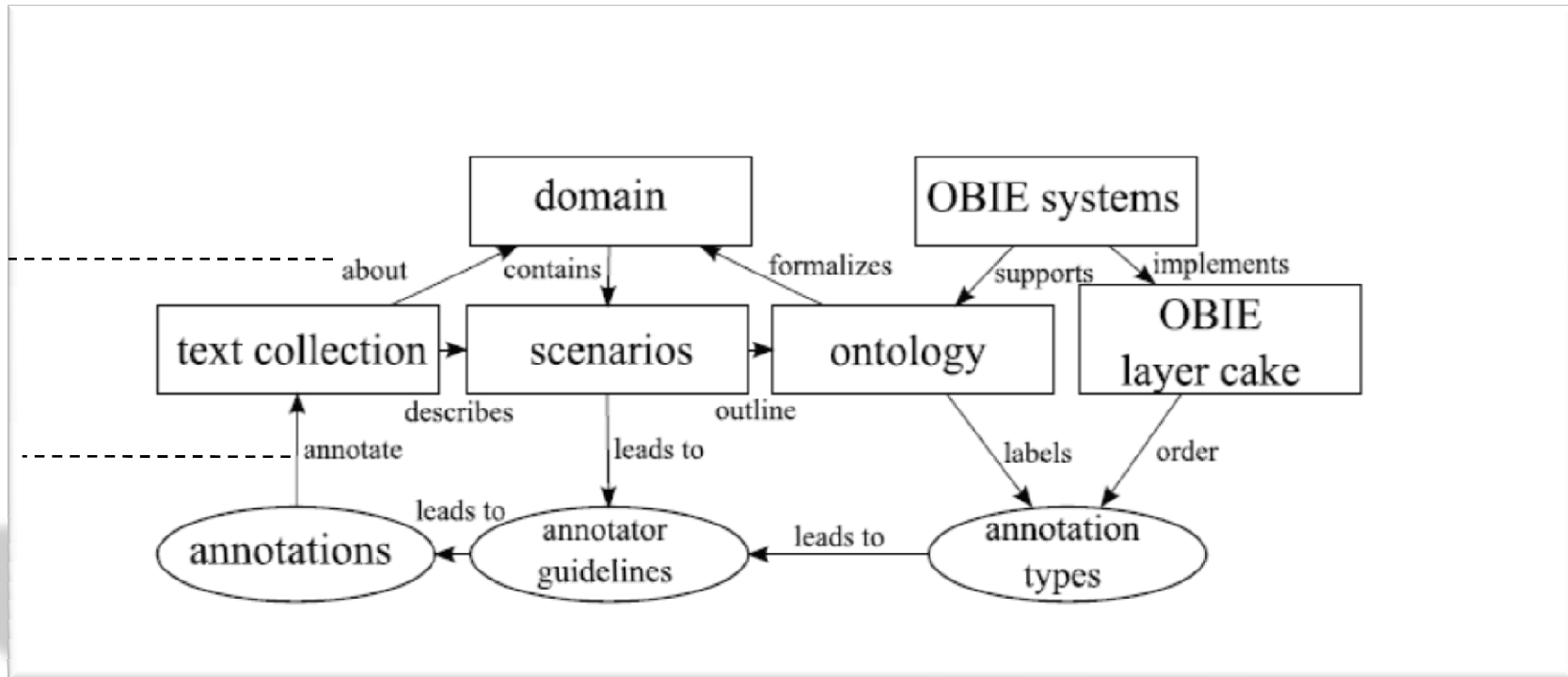Instantiation

Contextualization

For evaluating a certain OBIE systems with OCAS, its tasks have to be mapped to the OBIE layer cake.
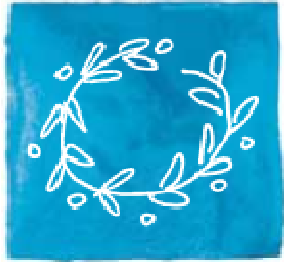
# Ontology-Based Corpus and Annotation Scheme

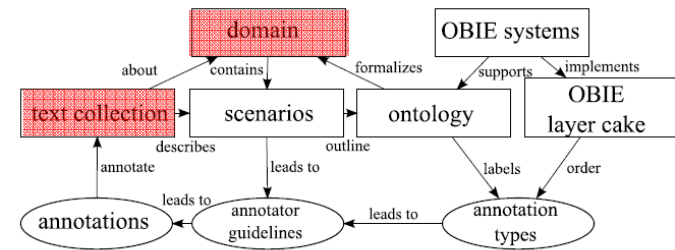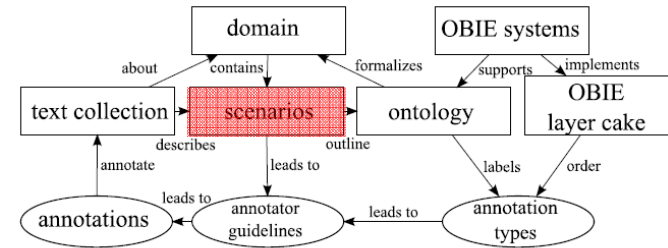The process for creating high quality OBIE corpora:

# Domain and Text Collection



News articles (from ABC, BBC) about the Olympic summer games 2004 comprise several benefits …
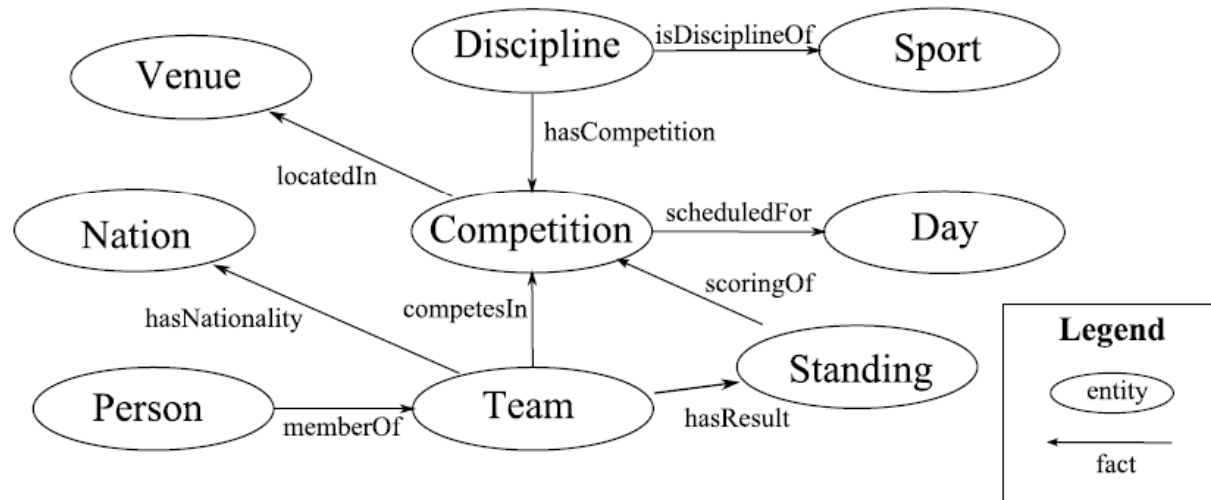
- **Closeness**, limited to a few, but strictly defined concepts
- **Compactness**, concepts of domain model are highly coherent
- **Richness**, concepts can be covered completely by a text corpus
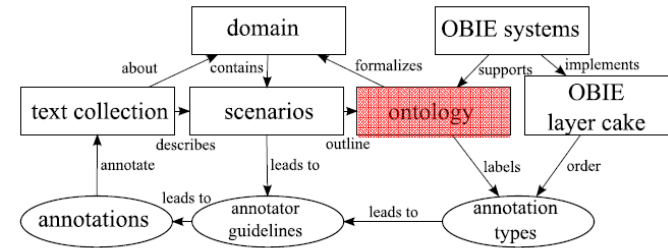
# Scenario Maps



Formal representation frame of a text.



In OCAS, we used two scenario maps:
I.   Single person
II.  Team

# Ontology
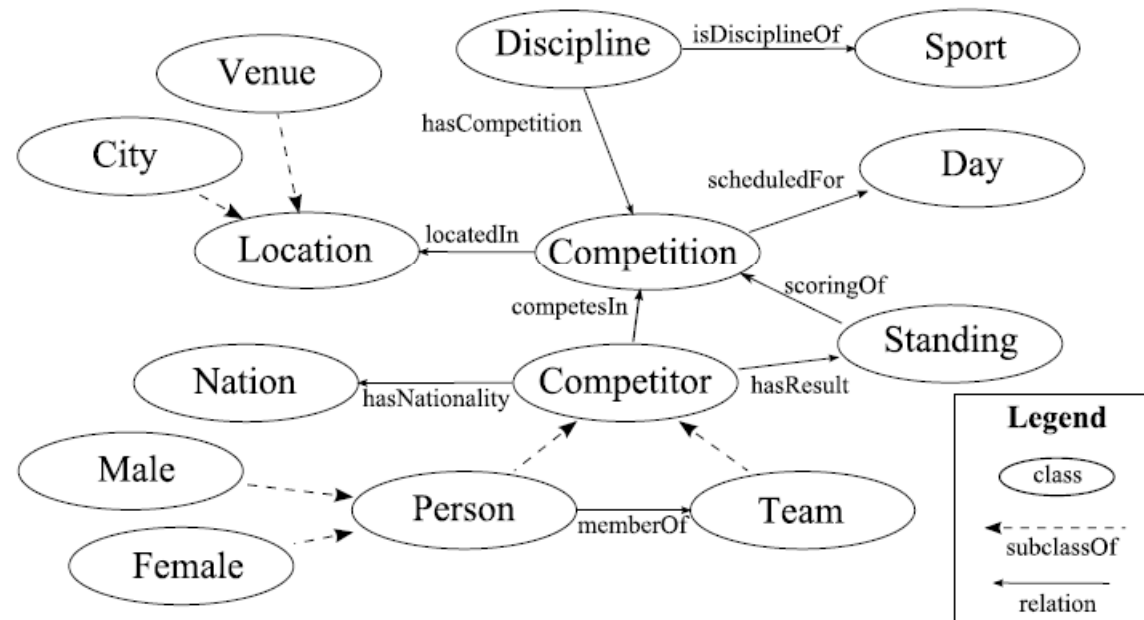


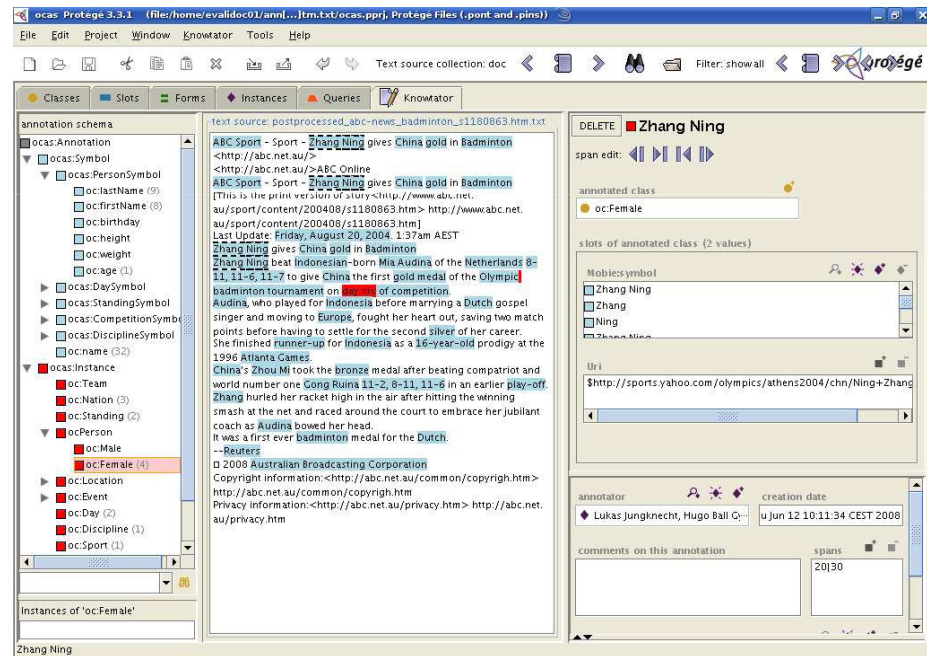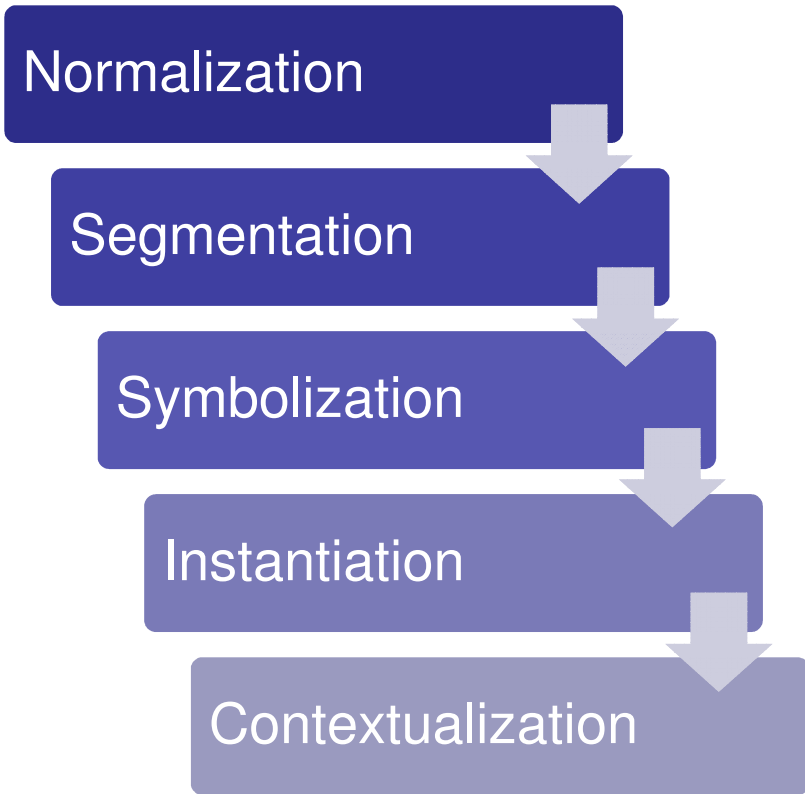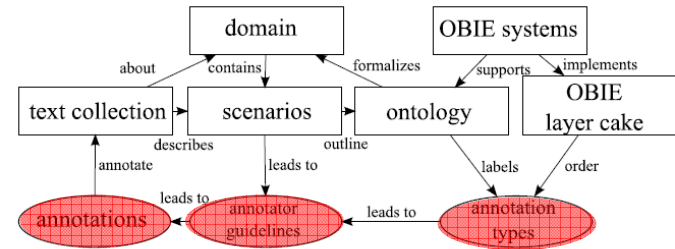Ontology was evolved from aggregated scenario maps.



Fig. 2. The ontology for the OCAS2008 test corpus.

The instance base was populated with data from Wikipedia and Yahoo! sports.

# Annotation Types



**Normalization**

**Segmentation**

**Symbolization**

**Instantiation**

**Contextualization**



Knowtator, Ogren, P.V., Proc. of
North American Chapter of ACL 2006

# Annotation Types



...

...

Datatype properties

Explicit, implicit instances
and object properties

Explicit, implicit facts and
classified datatype properties



Knowtator, Ogren, P.V., Proc. of
North American Chapter of ACL 2006

# Reviewing OCAS

## Facts

- Selected **121 rich articles** from 5,000 candidates. The selected articles contain a total of 31,102 words.

- **6 persons** annotated during **8 days**.

- Annotation process took a total of **176 person hours**.

- Results in **5.66 hours** for annotating **1,000 words**.

## Planned Quality Assurance

**Completeness** - documents have to be annotated completely with respect to scenario maps.

**Consistency** between annotations, e. g. each annotated instance must also be marked as symbol.

**Correctness** of annotations, i.e. every annotated instance in the corpus respects domain ontology's instance set and uses the same URIs.

# Summary

**OCAS comprises:**

a **generic OBIE architecture** called *OBIE layer cake for comparing OBIE systems* by similar subtasks,
a document **corpus of 121 news articles** (ABC, BBC) with 31,000 words about a closed domain (Olympic Summer Games 2004),
a compact **domain ontology** about the Olympic Summer Games 2004 including more than 40,000 instances from Wikipedia and Yahoo!,
two **annotation scenarios** that extend traditional template-based evaluations,
an **annotation set** that contains typed annotations according to the ontology and the *OBIE layer cake,*
annotations that concern **symbols**, **instances**, explicitly occurring **facts** and implicit induced facts
**human created annotations** according to predefined specifications.

# Outlook

Verify quality of OCAS 2008

Evaluate iDocument with OCAS 2008

Evaluate GATE's OBIE facilities with OCAS 2008
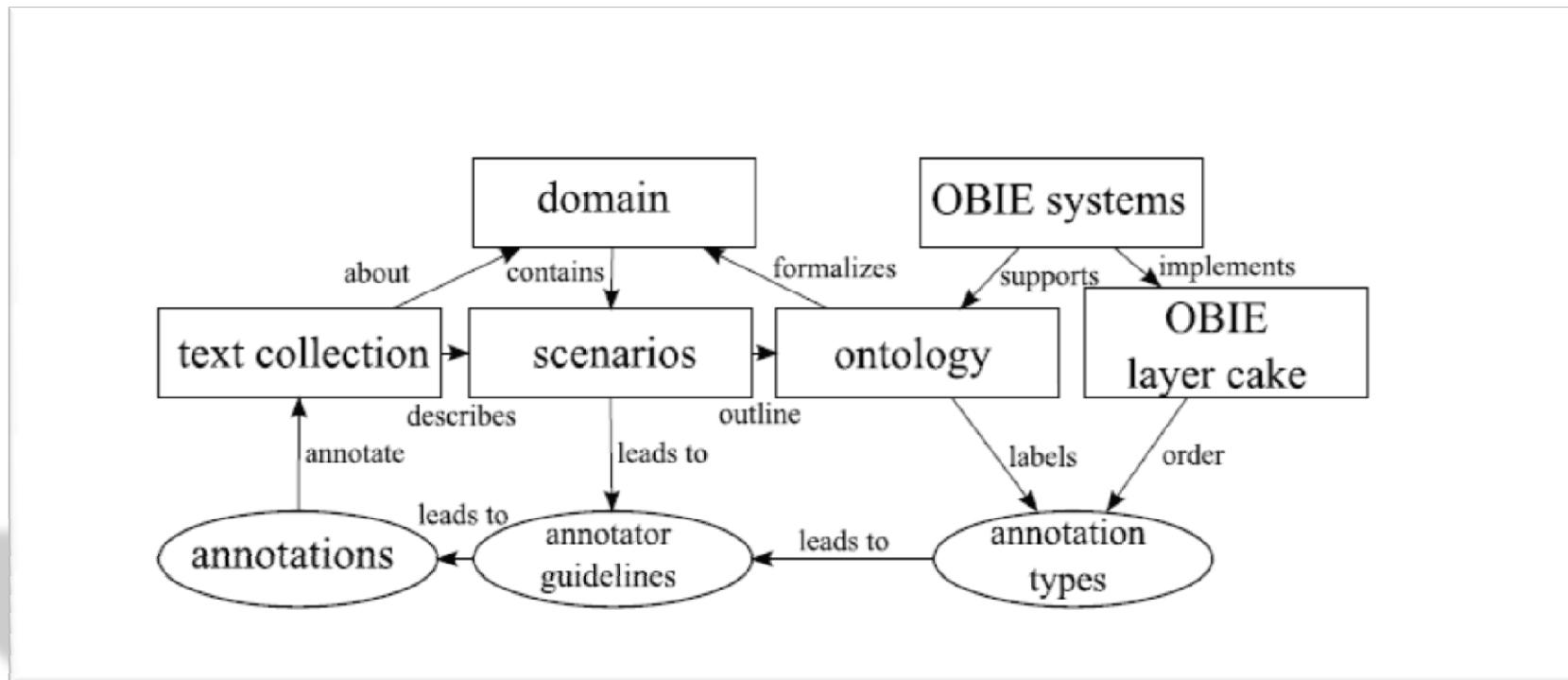
Publish OCAS
(RDF dump of ontology, annotated text corpus, test results)

Other ideas?

# Thanks for attention

Any questions, tips, comments or remarks?



More information on OCAS will appear soon at http://idocument.opendfki.de