# Language Technology and the Semantic Web

## Dr. Günter Neumann

http://www.dfki.de/~neumann
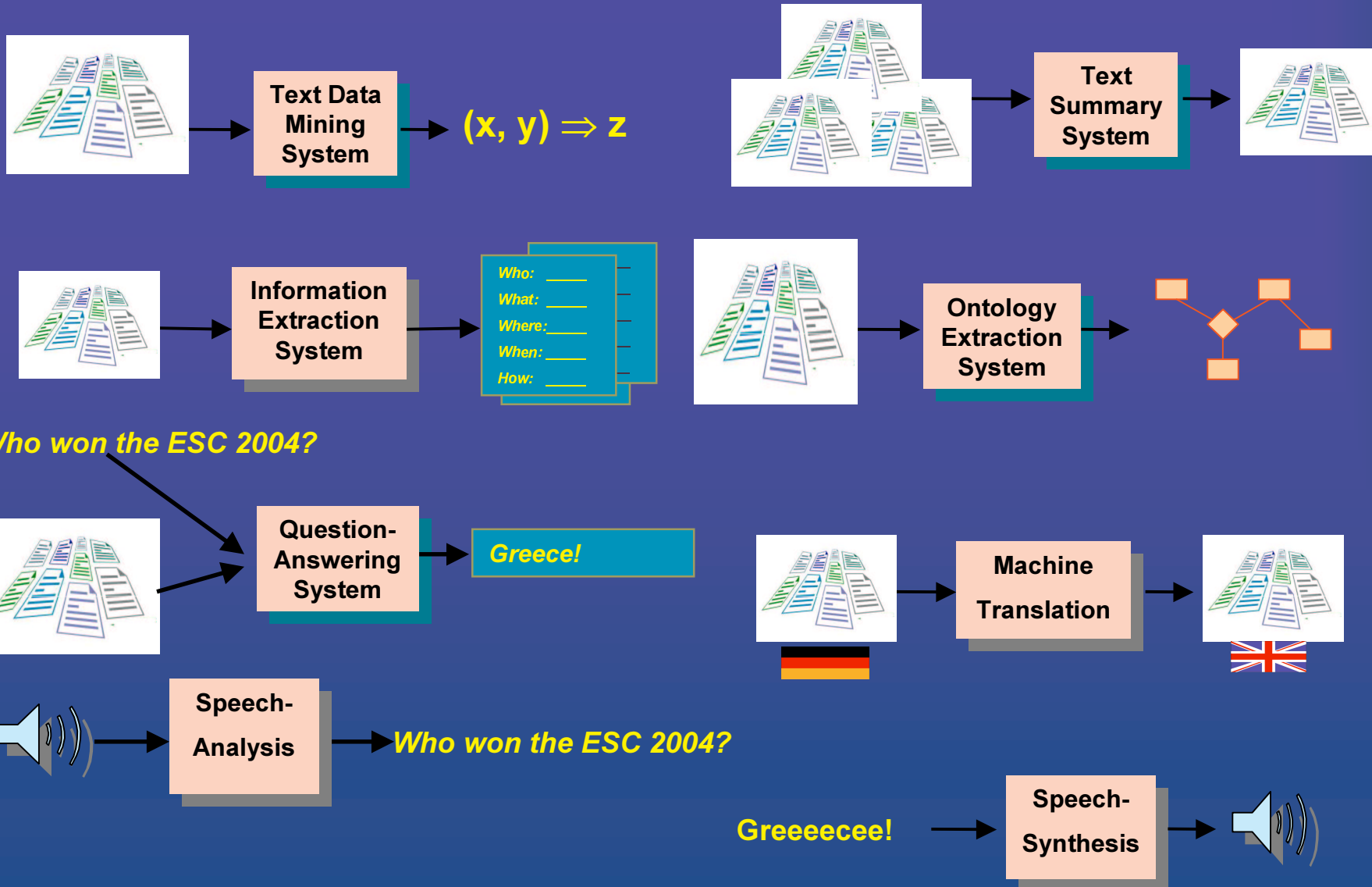
Language Technology-Lab

DFKI, Saarbrücken

# Overview

- Language Technology
- Semantic Web
- Information Extraction
- Information Access

# Human Language Technology

- *Human Language Technology LT* – covers
  - The design and implementation of algorithms, data and electronic devices for processing of natural language (text and speech), and
  - Their integration into real-world applications and products

- Language Technology defines the engineering part of computational linguistic

# LT-methods cover many areas

Text Data Mining System

$(x, y) \Rightarrow z$

Text Summary System

Information Extraction System

Who: _____
What: _____
Where: _____
When: _____
How: _____

Ontology Extraction System

**Who won the ESC 2004?**

Question-Answering System

*Greece!*

Machine Translation

Speech-Analysis

**Who won the ESC 2004?**

**Greeeecee!**

Speech-Synthesis

Multi/cross-linguality is of great importance in all these areas!

# LT as embedded part of applications

- **Human-Machine Communication**

- **Data-oriented Knowledge Acquisition**

### Integration

- Modularity
- Multi-media
- Software-Engineering standards

### High Performance

- Real-time
- Robustness
- Scalability
- Adaptation
- Evaluation

# Language Technology

- **LT-Methods**
  - Named Entity-Recognition
  - PoS/Sem-Tagging
  - Controlled Languages
  - Integration of shallow & deep NLP („text zooming")
  - Reference-resolution
  - NL-oriented ontologies

- **Core technology**
  - Efficient data structures
  - Weighted finite state automata
  - Machine learning
  - Statistical inference

- Already a successful technology transfer
  - Industry (Microsoft, IBM, Siemens, Telekom, ...) & Spin-offs, competence centers, ...
  - Speech-systems, MT, Editors, Text-Mining, Knowledge-Mining Content-Management, ...

- Newest Technology Hype: the Semantic Web
  - What role does it play for LT?

# The Semantic Web (SW)

- Tim Berners-Lee, 1998:
    - "This document is a plan for achieving a set of connected applications for data on the Web in such a way as to form a consistent logical web of data (semantic web)."
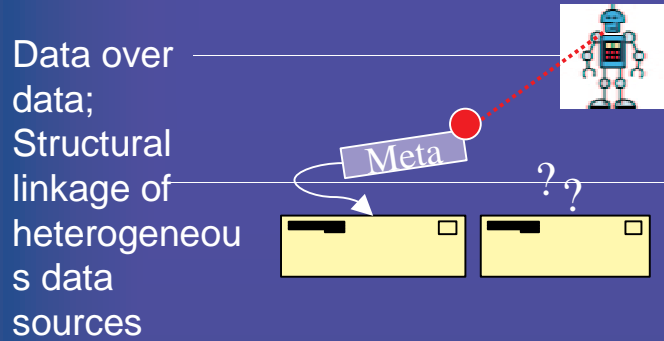


- Tim Berners-Lee et al., 2001
    - "… an extension of the current web in which information is given well-defined meaning, better enabling computers and people to work in cooperation."
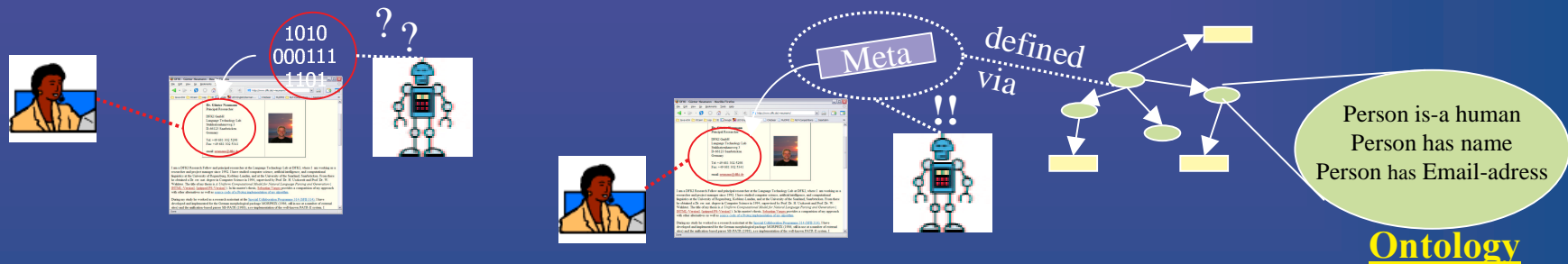
# SW – illustrated

**2** Add meta-data

The existing web will further emerge, so that computers can understand content on-line, to better help humans to organize, search, and exchange information.

Data over data; Structural linkage of heterogeneous data sources

Meta

? ?

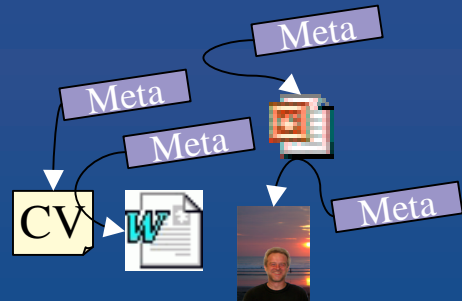**3** Ontologies associate meaning to meta-data

SW exists of meta-data and links to global ontolgoies, which define the meaning of terms.

An ontology serves as a structural vocabulary for the interpretation of domain-specific terms.

**4** Strukturiertes Web von Daten

1010 000111 1101

? ?

Meta

*defined via*

Person is-a human
Person has name
Person has Email-adress

**Ontology**

**5** The SW does not only consider Web-pages
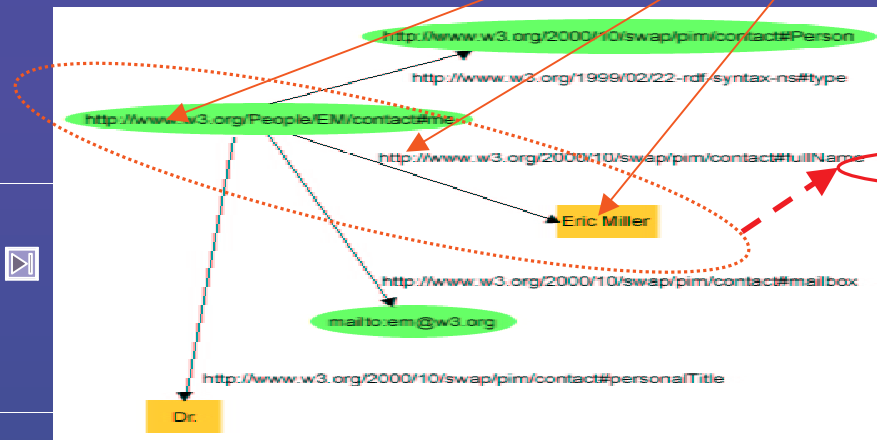
Meta
Meta
Meta
Meta

CV

**6** How will I use the SW?

• Intelligent information search;
• Automatic support for the management of my personal information on the SW

# RDF and OWL: Modeling data on the SW

**1  RDF: Resource Description Framework**

RDF is language for the representation of meta-data over web resources.
RDF-statements are triples of the form **(Subj, Pred, Obj).**



**2  XML & N3 sind alternative RDF-Syntaxen**

**XML schematically: <Subj> <Pred> Obj </Pred> </Subj>**

**N3:**
@prefix rdf: http://www.w3.org/1999/02/22-rdf-syntax-ns# .
@prefix contact: http://www.w3.org/2000/10/swap/pim/contact# .
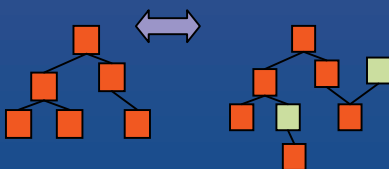@prefix EM: http://www.w3.org/People/EM/contact# .

EM:me rdf:type                contact:person .
EM:me contact:full-name       "Eric Miller" .
EM:me contact:personaiTitie "Dr." .
EM:me contact:mailbox  rdf:resoure "mailto:em@w3.org .
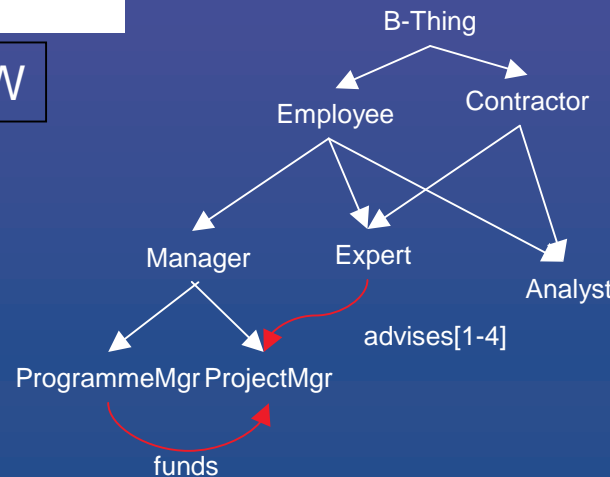
**3  OWL: Web Ontology Language**



•some RDF-statements have a fix interpretation (is-a, =, inverseOf, card, ...)

•**Sharing** of information between individuals from multiple documents ⇒ Web of data from heterogeneous sources
•Semantic of OWL as basis for  inference mechanism over these data structures.

**4  Relevante Aspekte für das SW**

standardization, Web-globalization, distribution of resources

**5  Ontology Mapping**
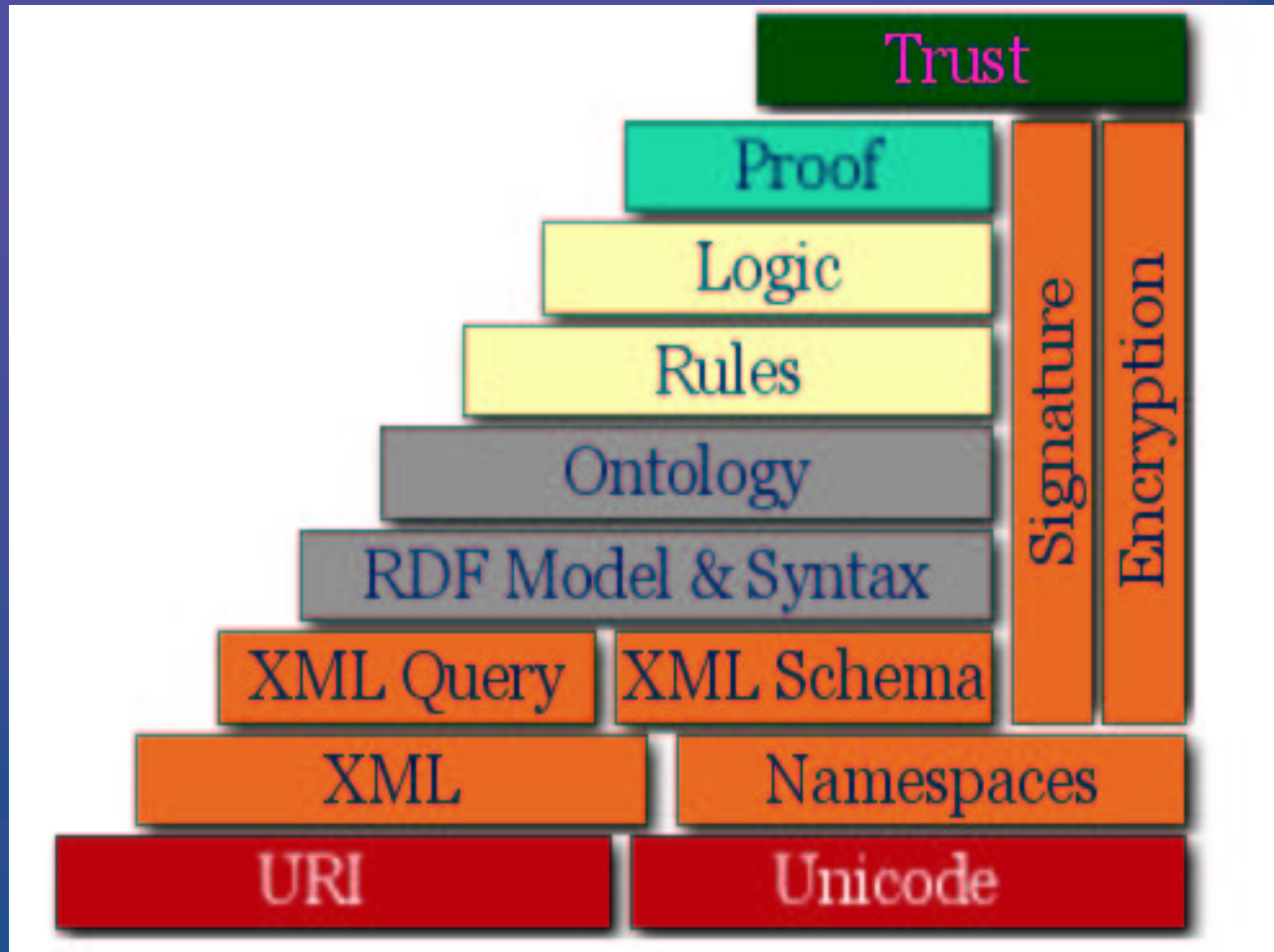


Mapping between distributed, local ontologies

# The SW-pyramid

Basic research

Current focus of
major efforts

Established standards

# Relevance of LT for SW

1 During the transi[tion] WWW to the SW, L[T] technology

**Intelligent Information Extraction**

2 As long as the hum[an] [is in the] Loop", NL will remain to [be the] core Human-SW communic[a]tion device.

**Intelligent Information Access**

3 Humans will also in the future exchange knowledge via NL documents: Semantically annotated documents as Human-SW interface

4 NL-generation of information in form of NL-Text, e.g., heterogeneous resources, dynamically created reports, newspapers, …

CV

# Information Extraction (IE)

**Template:**

ManagementSuccession

*PersonIn:* _____

*PersonOut:* _____

*Position:* _____

*Organisation:* _____

*TimeIn:* _____

*TimeOut:* _____

Text classification

Linguistic processing
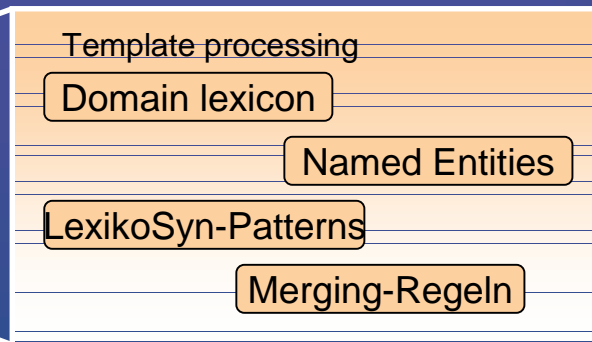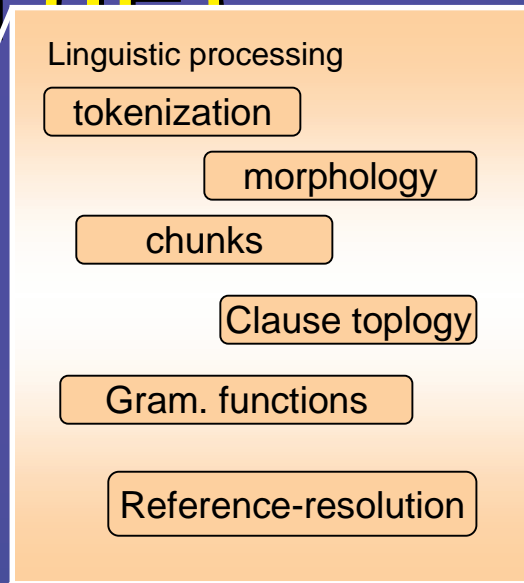
Template processing

Linguistic processing

- tokenization
- morphology
- chunks
- Clause toplogy
- Gram. functions
- Reference-resolution

Template processing

- Domain lexicon
- Named Entities
- LexikoSyn-Patterns
- Merging-Regeln

documents

**Dr. Hermann Wirth**, bisheriger **Leiter** der **Musikhochschule München**, verabschiedete sich heute aus dem Amt. Der 65jährige tritt seinen wohlverdienten Ruhestand an. Als seine Nachfolgerin wurde **Sabine Klinger** benannt. Ebenfalls neu besetzt wurde die Stelle des Musikdirektors. Annelie Häfner folgt Christian Meindl nach.

ManagementSuccession

*PersonIn:* Klinger

*PersonOut:* Wirth

*Position:* Leiter

*Organisation:* Musikhochschule München

*TimeIn:* _____

*TimeOut:* 3.4.2002

# IE for semantic annotation

Identification of IE-sub-tasks:
- basic entities (e.g., proper names)
- binary relations between entities
- n-ary relations/events

**Machine learning!**

Automatic Content Extraction (ACE)

- Spezification of an IE-core-ontology
- Annotation-specification & -tools
- Templates as specializations of the IE-core-ontology (also multi-templates)

IE as core for semantic annotation
- identification
- discovery
- validation
- evaluation

of semantic relationships & as basis for the automatic creation of meta data

# IE for semantic annotation

domain
ontology

Domain lexicon

IE-core
ontology

IE-core
system

inference
engine

NL-oriented
ontology

{ <t1, rel?, t2> }
<NP, VG, NP>
<NE, ?, NP>
<NE, ofPP, NE>

LT as basis of
• concept identification
• determination of plausible
structural relation candidates

# Example for entities & their mentions

[COLOGNE, [Germany]] (AP) _ [A [Chilean] exile] has filed a complaint against [former [Chilean] dictator Gen. Augusto Pinochet] accusing [him] of responsibility for [her] arrest and torture in [Chile] in 1973, [prosecutors] said Tuesday.
[The woman, [a Chilean] who has since gained [German] citizenship], accused [Pinochet] of depriving [her] of personal liberty and causing bodily harm during [her] arrest and torture.

Person
Organization
Geopolitical Entity

# LT-challenges

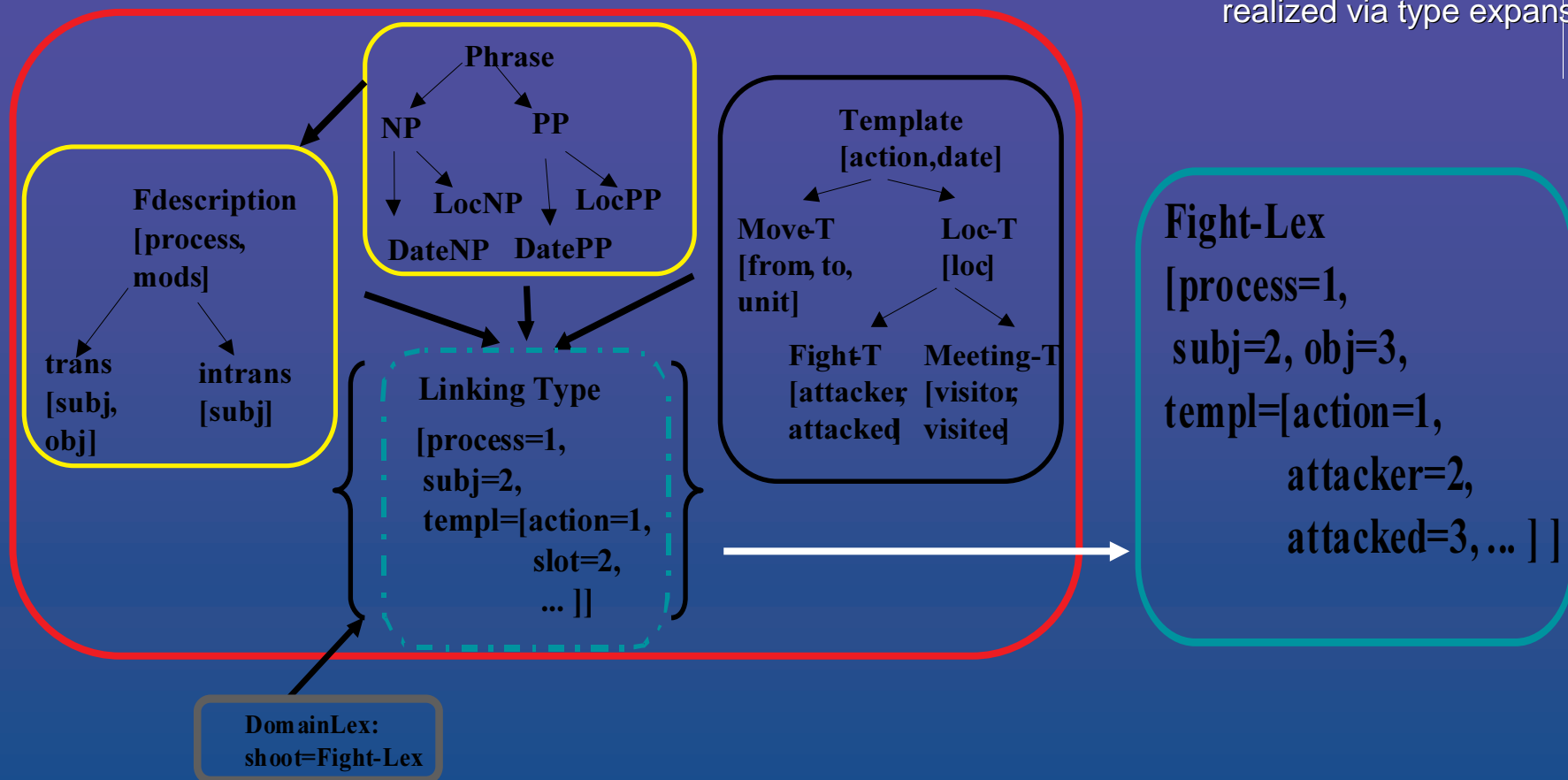Identification of verbalizations/mentioning of concepts/instances

- Linking of domain ontology and NL-oriented ontology (e.g., WordNet)
- Paraphrasing
- Metonymy ("Peking orgainzes the Olympic Games 2008.")
- Reference identification ("Chancellor Schröder, Schröder, the German chancellor, he, …")
- Analysis of sublanguages as basis for adaptive IE (cf. Grishman, 2001)

# Domain modeling in DFKI system SMES is realised using typed feature structures

○ Domain modeling via hierarchy of templates (black box), using the formalism TDL, which is also used to model hierarchies of linguistic objects ( yellow boxes).

○ The interface between domain knowledge and linguistic entities is specified via *linking types* (green box), which represent a close connection between concepts of the different layers, and which are accessible via the domain lexicon (brown & green box). Template-filling is then realized via type expansion.

**Phrase**

**NP**      **PP**

**LocNP**   **LocPP**

**DateNP   DatePP**

**Fdescription**
**[process,**
**mods]**

**trans**          **intrans**
**[subj,**        **[subj]**
**obj]**

**Linking Type**
**[process=1,**
**subj=2,**
**templ=[action=1,**
**slot=2,**
**... ]]**

**Template**
**[action,date]**

**Move-T**              **Loc-T**
**[from, to,**          **[loc]**
**unit]**

**Fight-T    Meeting-T**
**[attacker, [visitor,**
**attacked] visited]**

**Fight-Lex**
**[process=1,**
**subj=2, obj=3,**
**templ=[action=1,**
**attacker=2,**
**attacked=3, ... ] ]**

**DomainLex:**
**shoot=Fight-Lex**

# NL-annotations for the SW

Starting point: START multi-media QA system, by Boris Katz et al.

T-expression
<subject relation object>

Central issues
1. Sentence-based NL analysis
2. NL-annotations for multi-media information segments



*Bill surprised Hillary with his answer*
<<Bill surprise Hillary> with answer>
<answer related-to Bill>

Processing of huge text collections:
1. Extraction of relevant sentences from texts.
2. Syntax analysis
3. Annotation of the texts with syntax

NL-Question
*Whose answer surprised Hillary?*
<answer surprise Hillary>
<answer related-to **whom**>

7/2004, GN

# Haystack: the universal information client

http://haystack.lcs.mit.edu/

Motivation:
semantic annotation should be a side-effect of daily use of computer.

Idea:
Personalized information portal for all relevant services, like email, documents, calender, Web-pages, ...

Collection of all data uniformly via RDF-database

Programming language Adenine for the manipulation of frequent (i.e., as support for the implementation of specific service programs).

## Haystack RDF-database:

```
@prefix dc: http:77purl.org/dc/elements/1.1/
@prefix :  http://www.50states.com/data#

{ :State
     rdf:type    rdfs:Class ;
     rdfs:label  „State"
}
{ :bird
     rdf:type      rdf:Property ;
     rdfs:label     „State bird" ;
     rdfs:domain :State
}
{ :alabama
     rdf:type :State ;
     dc:title „Alabama" ;
     :bird „Yellowhammer" ;
     :flower „Camellia" ;
     :population „4447100" ;
     ...
}
```

## Natural language schema:

```
@prefix nl: http://www.ai.mit.edu/projects/infolab/start#

Add{ :stateAttribute
        rdf:type          nl:NaturalLanguageSchema ;
        nl:annotation @( :attribute „of" :state) ;
        nl:code          :stateAttributeCode
}
Add{ :attribute
        rdf:type         nl:Parameter ;
        nl:domain        rdf:Property ;
        nl:descrProp   rdf:label ;
}
Add{ :state
        rdf:type          :Parameter ;
        nl:domain         :State ;
        nl:descrProp    dc:title;
}

Method
:stateAttributeCode : state=state :attribute=attribute
        return (ask {state  attribute ?x })
```

**Ask{state=:alabama, attribute=:bird, ?x }**

**⇒ ?x= „Yellowhammer"**

**Antwort:** *Yellowhammer*

**:bird        ⇐ :attribute=„state bird"**
**:alabama  ⇐  :state=„Alabama"**

**Frage:** *What is the state bird of Alabama?*

# Example:
# Linking of t-expressions & RDF

@prefix nl: http://www.ai.mit.edu/projects/infolab/start#

Add{ :Person
    rdf:type        rdfs:Class ;
}

Add{ :homeAddress
    rdf:type        rdf:Property ;
    rdfs:domain   :Person ;

    nl:annotation  @(nl:subj „lives at" nl:obj) ;
    nl:annotation  @(nl:subj „'s home adress is" nl:obj) ;
    nl:annotation  @(nl:subj „'s apartment" nl:obj) ;

    nl:generation  @(nl:subj „'s home address is" nl:obj) ;

}

Remarks:

- NL-annotations as a means for controlling the paraphrasing potential of NL expressions

- Richer linguistic annotations are possible (e.g., fine-grained grammatical functions, agreement)

- Also relevant for user-oriented adaptation of service programs

# Natural language annotations for the SW

- NL used as meta-data
    - Readability of RDF
    - Supports transition from WWW to SW
    - NL-annotation specifies which kind of (NL)-question a meta-data is able to answer
      $\Rightarrow$ controlled question-answering systems

- Information access (IA) within SW
    - Development of programs, which help a user to locate, to collect, to compare and to link information

- NL is the most natural way for user to perform IA
    - SW should support in the same way IA using specialized languages/exchange formats & NL

# Relevance

- Approach is open for future extensions:
    - statistical-based models (add weight to the NL-annotations)
    - Machine Learning of NL-annotations on basis fo ontology-oriented IE (cf. Hovy et al. 2002)
- The current mechanism of NL-annotations is idiosyncratic, however at DFKI we plan the following:
    - Exploration of a linking mechanism between dependency structure and RDF/OWL
    - Foundation for novel template-based QA-strategies

# Example for the processing of complex questions

- Approach:
  - Select templates via Q-Type & Q-Focus:
    - Definition question, list-question
    - Person: born-where, born-when, business-what $\Rightarrow$ Ontology
  - Pro property P, select IR-Schema:
    - NL-based query-pattern
  - P might be:
    - From the set of known NE-types (person, location, date, …) $\Rightarrow$ answer-type
    - NL-Phrase, which "describes" P, in case no a-type can be determined
- Compute for each P für jede P one/several IR-Query-terms, e.g.,
  - NE-type:person & text:<query term>

„Wer ist Thomas Mann?"

Q-type=c-definiton,
focus=<Person, „Thomas Mann">

IR-Schemata:
<PERSON> "geboren in" <LOCATION>

"(neTypes:LOCATION AND +geboren
(text:\"Thomas Mann\" OR text:Mann))"

**Search engine**

# IE-based question answering

- Approach can also be used for template-based questions:
  - let $t \in T$, set of templates, which are known to the system – via IE-Ontology – e.g., "management-sucession-Template"
  - for all properties E of t, combine E with NL-schema
    - E.g., "Person-In" $\Rightarrow$ (<PERS> "is_successor_of" <PERS>)
- Answering of complex questions
  - As composition of the answering of – relative to the conceptual description – simple questions
  - Implementation of this approach as part of the DFKI project Quetal (prototype as part of DFKI's qa@clef-2004 system)
  - Interactive online IE through close integration of IE & IA

# Concluding remarks

- LT is a key technology for the construction of the Semantic Web
- Very high requirements on
  - Performance
  - Modularity & integration
  - scalability & on-demand availability
  - Domain & user adaptation
- Systematic evaluation of LT-methods
  - Driving power & revisions of futuer developments
- In the future, cognitive-based methods will be considered
  - as inspiration for more effectiv LT-methods