

---

# Open-domain Cross-lingual Question Answering from Unstructured Documents

Günter Neumann

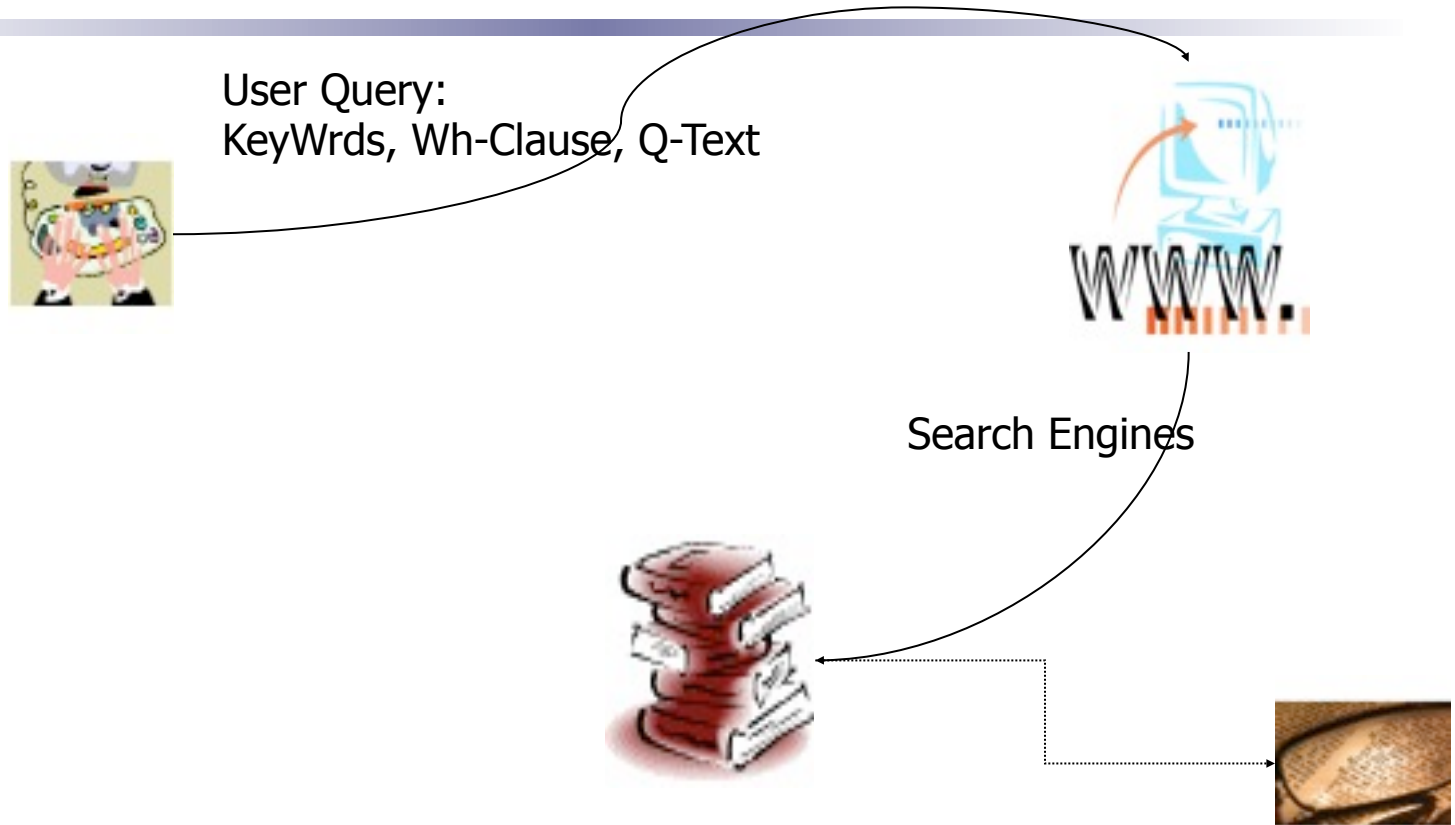
LT-lab, DFKI, Saarbrücken,

June, 2005

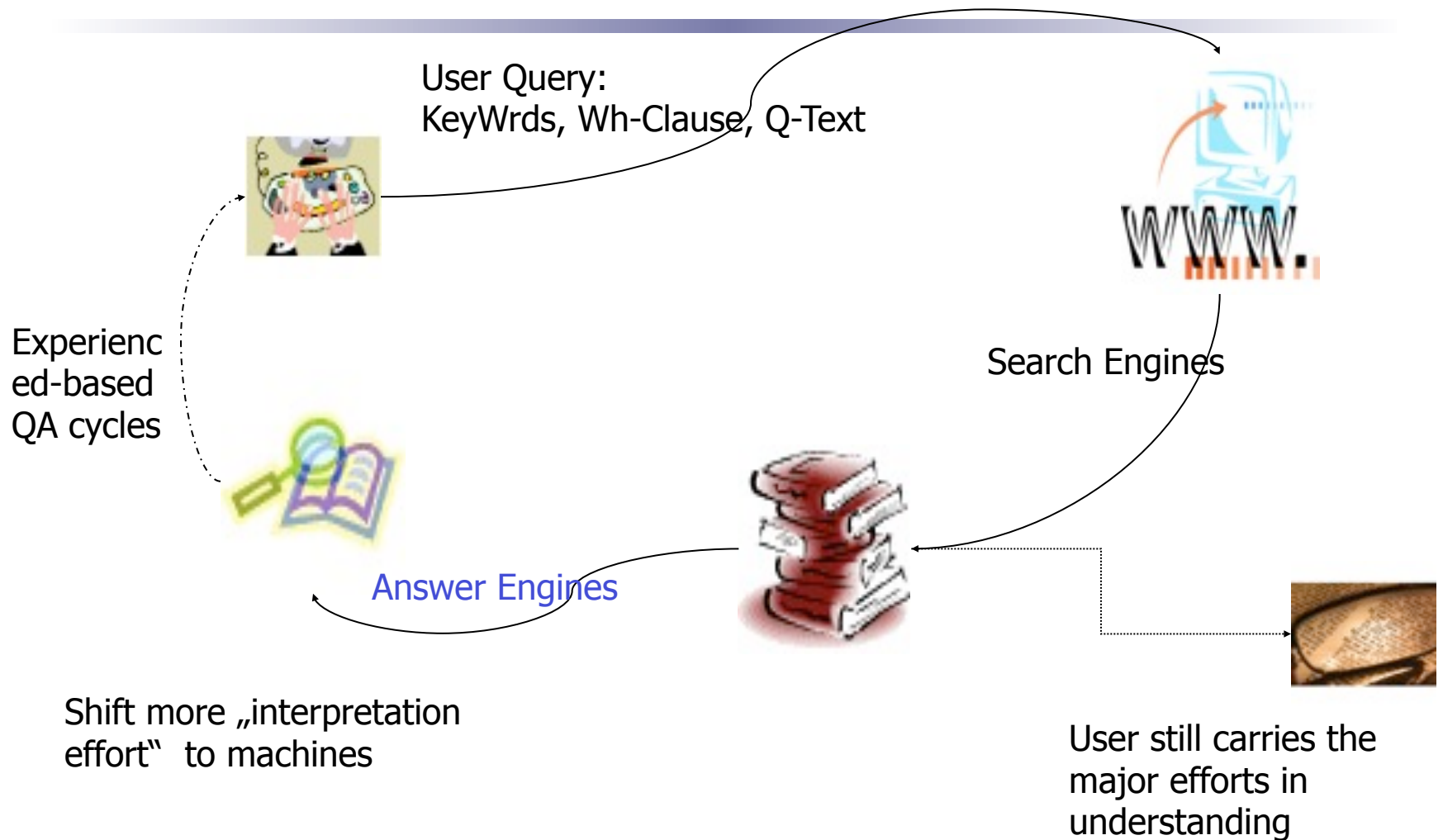
# Motivation: From Search Engines to Answer Engines

---

# Motivation: From Search Engines to Answer Engines

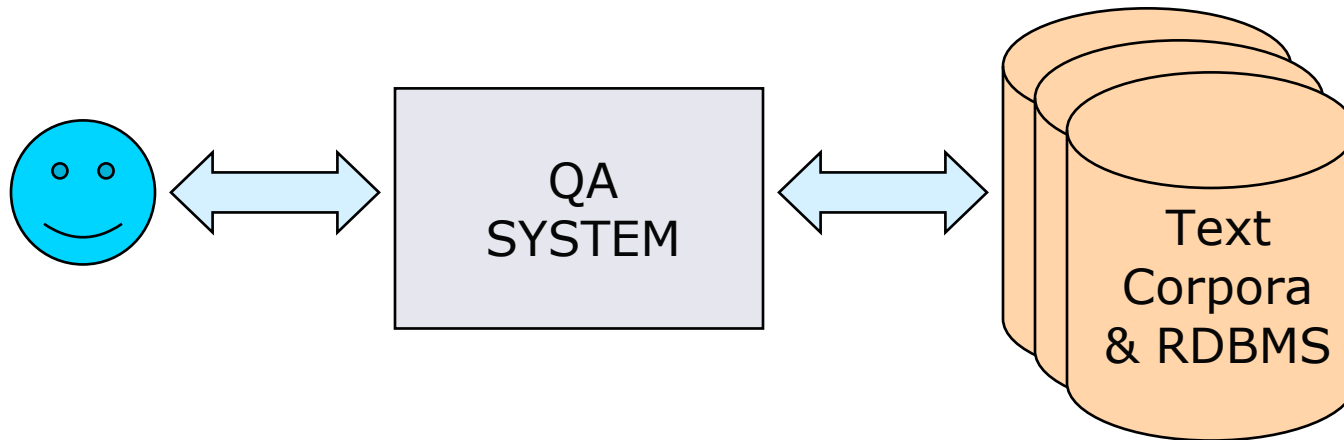


# Motivation: From Search Engines to Answer Engines



# Question Answering

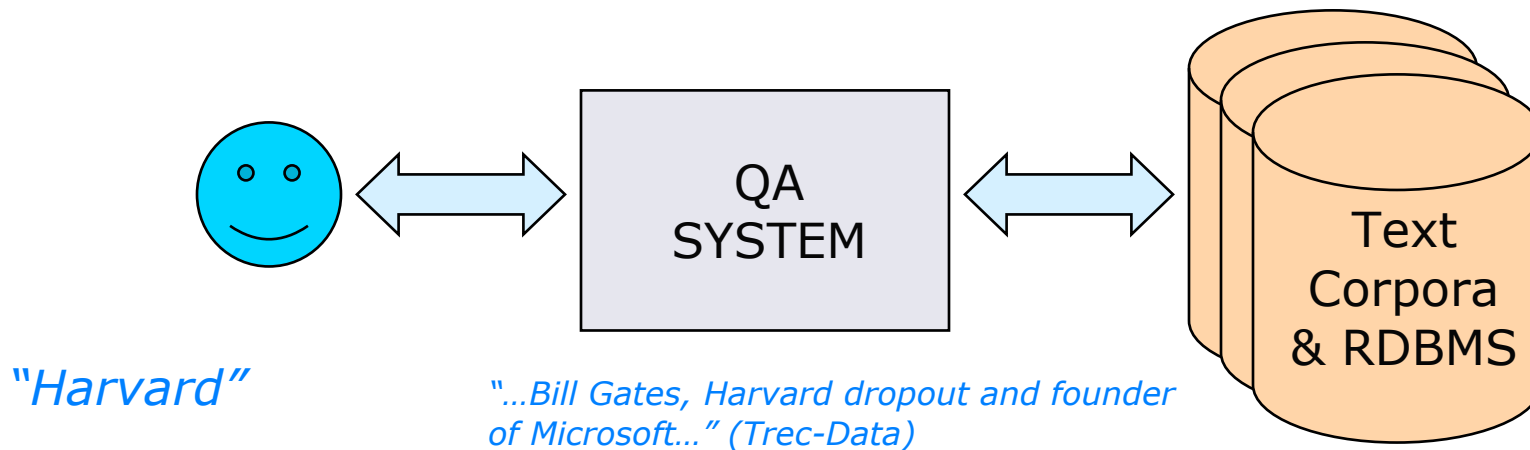
- Input: a question in NL; a set of text and database resources
- Output: a set of possible answers drawn from the resources



# Question Answering

- Input: a question in NL; a set of text and database resources
- Output: a set of possible answers drawn from the resources

*"Where did Bill Gates go to college?"*

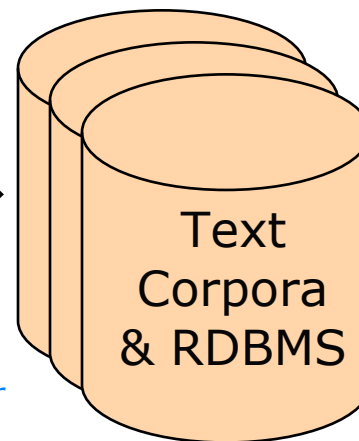
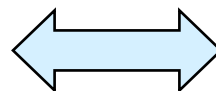
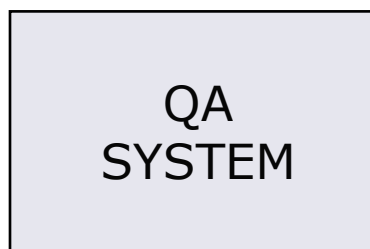
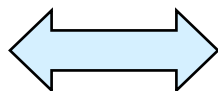


# Question Answering

- Input: a question in NL; a set of text and database resources
- Output: a set of possible answers drawn from the resources

*"Where did Bill Gates go to college?"*

*"What is the rainiest place on Earth?"*



*"Harvard"*

*"Mount Waialeale"*

*"...Bill Gates, Harvard dropout and founder of Microsoft..." (Trec-Data)*

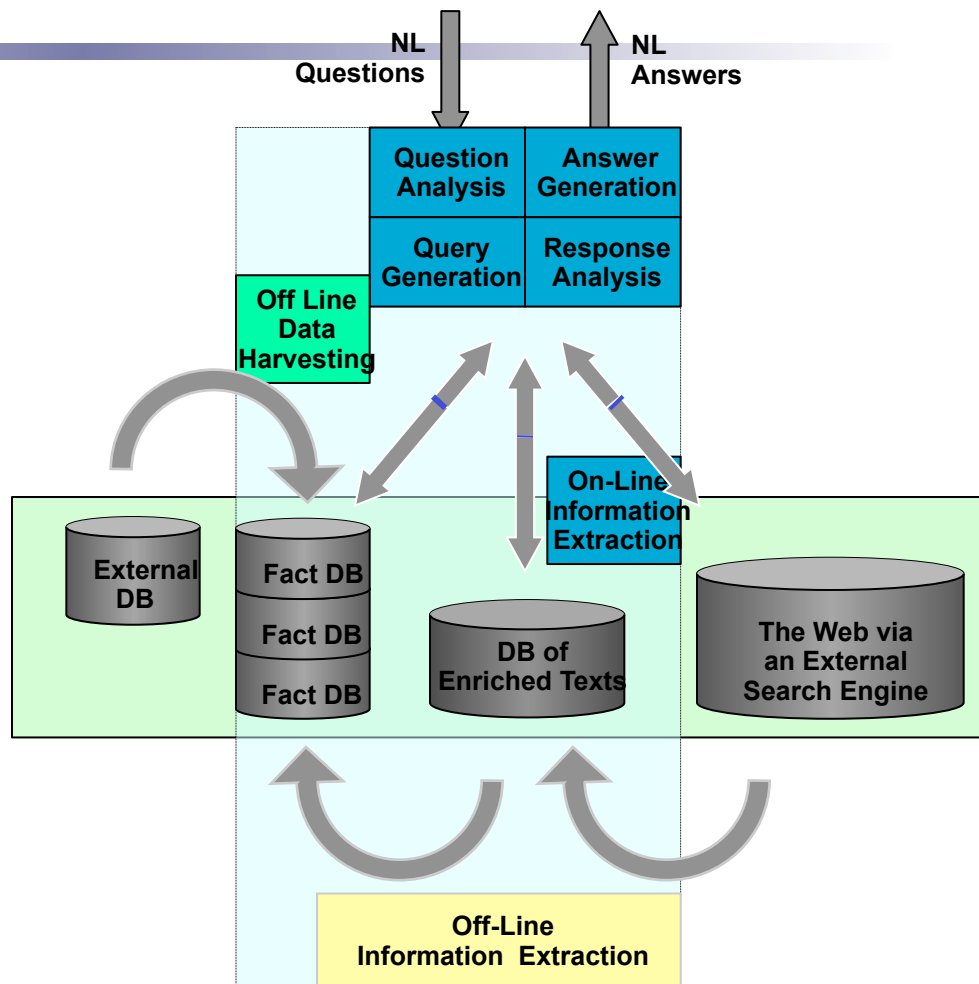
*"... In misty Seattle, Wash., last year, 32 inches of rain fell. Hong Kong gets about 80 inches a year, and even Pago Pago, noted for its prodigious showers, gets only about 196 inches annually. (The titleholder, according to the National Geographic Society, <sup>3</sup> is Mount Waialeale in Hawaii, where about 460 inches of rain falls each year.) ..."* (Trec-Data; but see Google-retrieved Web page.)

# Hybrid QA Architecture

## Hypothesis

real-life QA systems will perform best if they can

- *combine* the virtues of domain-specialized QA with open-domain QA
- *utilize* general knowledge about frequent types and
- *access* semi-structured knowledge bases



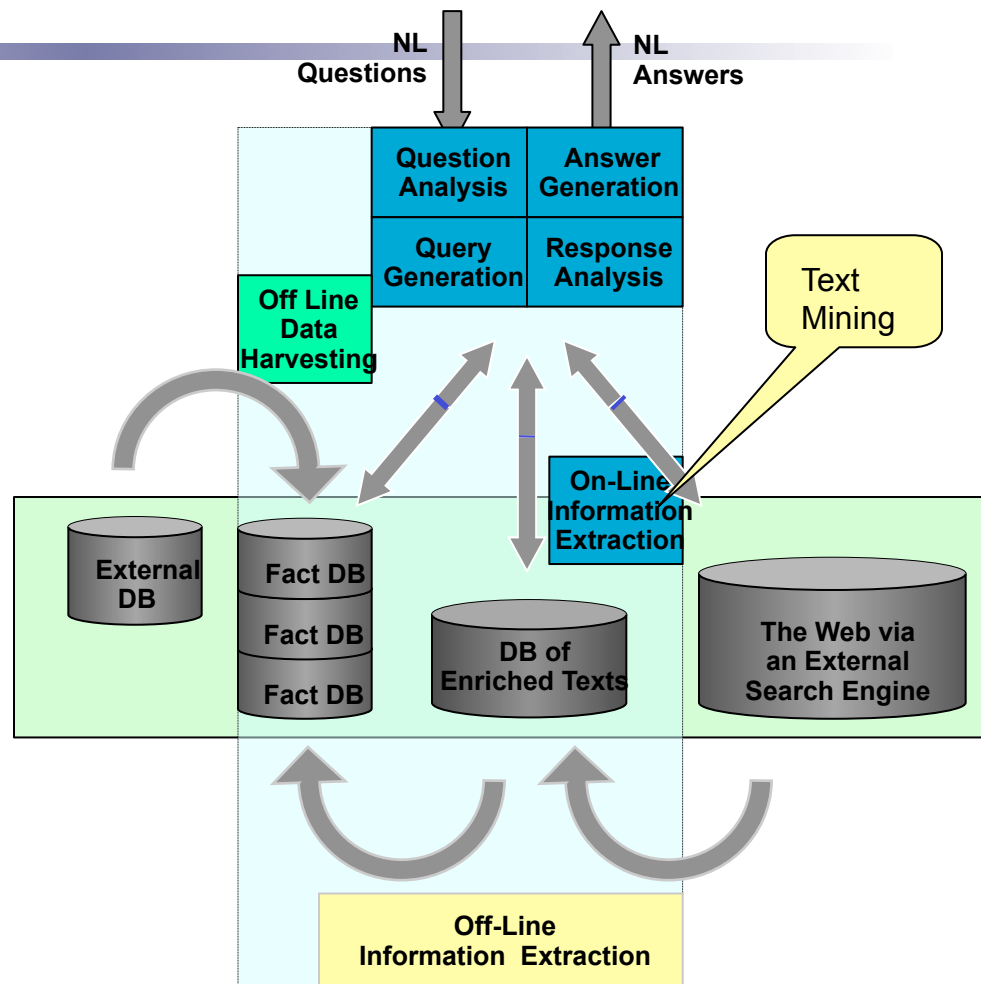


# Hybrid QA Architecture

## Hypothesis

real-life QA systems will perform best if they can

- *combine* the virtues of domain-specialized QA with open-domain QA
- *utilize* general knowledge about frequent types and
- *access* semi-structured knowledge bases



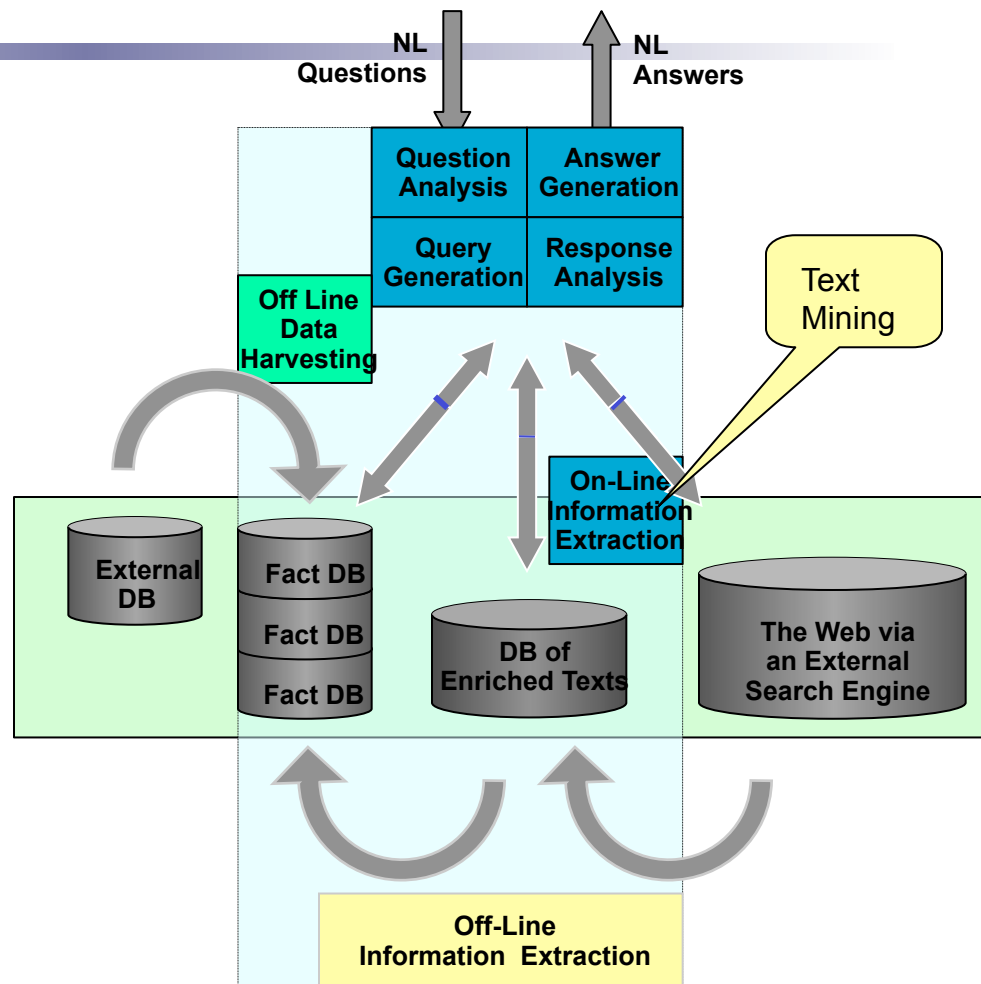
# Hybrid QA Architecture

## Hypothesis

real-life QA systems will perform best if they can

- *combine* the virtues of domain-specialized QA with open-domain QA
- *utilize* general knowledge about frequent types and
- *access* semi-structured knowledge bases

Advertisement:  
DFKI project Quetal  
2003-2005



# Design Issues

- Foster bottom-up system development
  - Data-driven, robustness, scalability
  - From shallow to deep NLP
- Large-scale answer processing
  - Coarse-grained uniform representation of query/documents
  - Text zooming
    - » From paragraphs to sentences to phrases
  - Ranking scheme for answer selection
- Common basis for
  - Online Web pages
  - Large textual sources

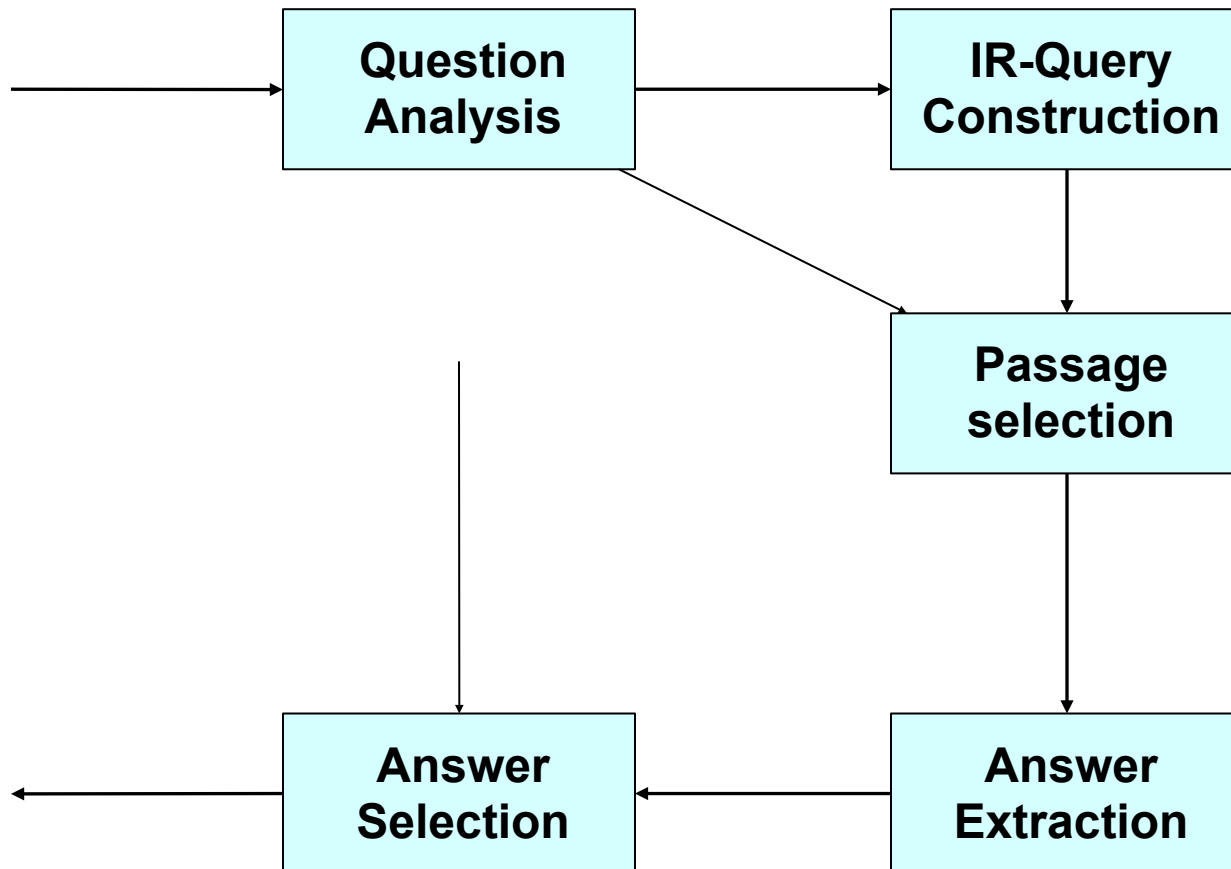


# Open-Domain Question Answering

---

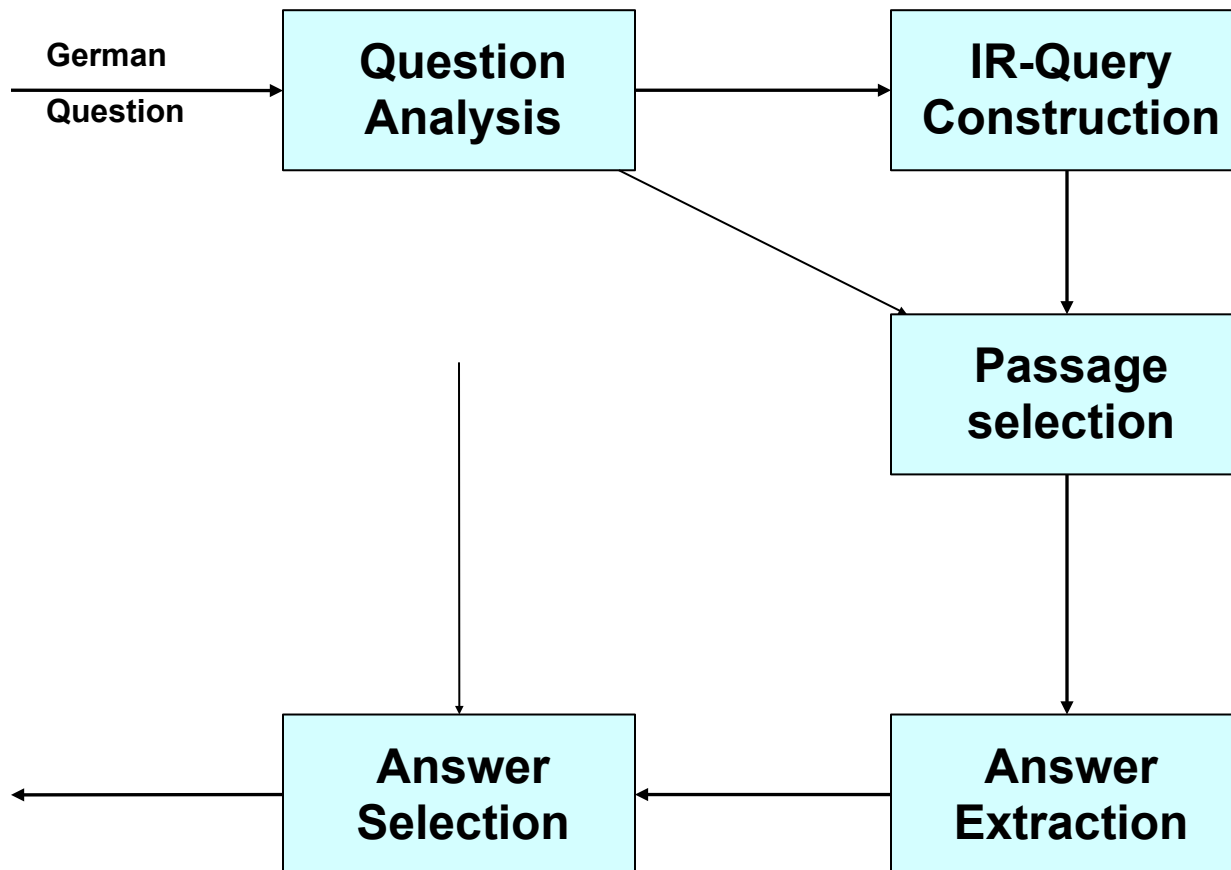
- **Open domain**
  - No restriction for the domain and type of question
  - No restriction on document source
- **Combines**
  - Information retrieval
  - Information extraction
  - Text mining
  - Computational Linguistics
- **Cross-lingual ODQA**
  - Express query in language X
  - Answer from documents in language Y

# Open-Domain Question-Answering



# Open-Domain Question-Answering

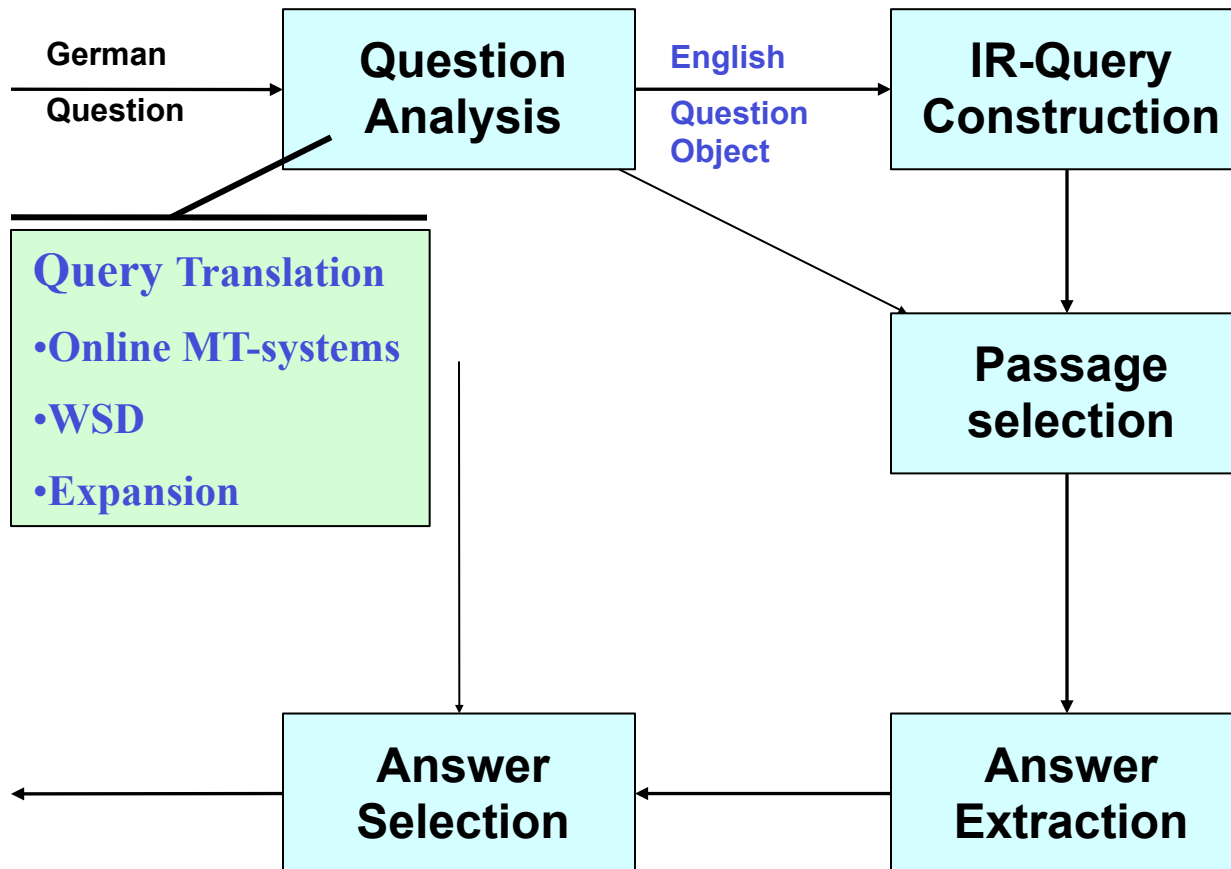
*“Mit wem ist David Beckham verheiratet?”*



# Open-Domain Question-Answering

*“Mit wem ist David Beckham verheiratet?”*

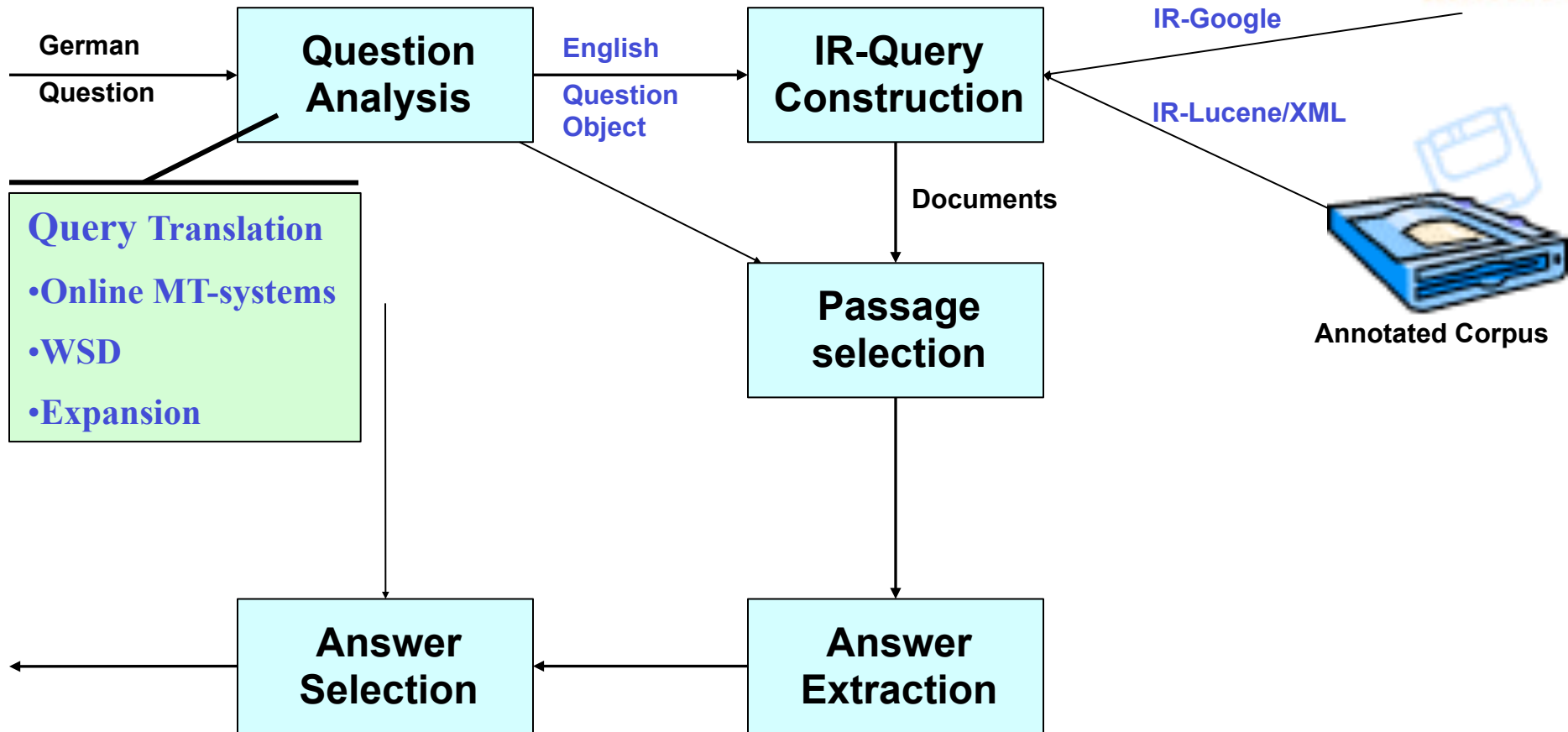
*{person:David Beckham, married, person:?}*



# Open-Domain Question-Answering

*“Mit wem ist David Beckham verheiratet?”*

*{person:David Beckham, married, person:?}*

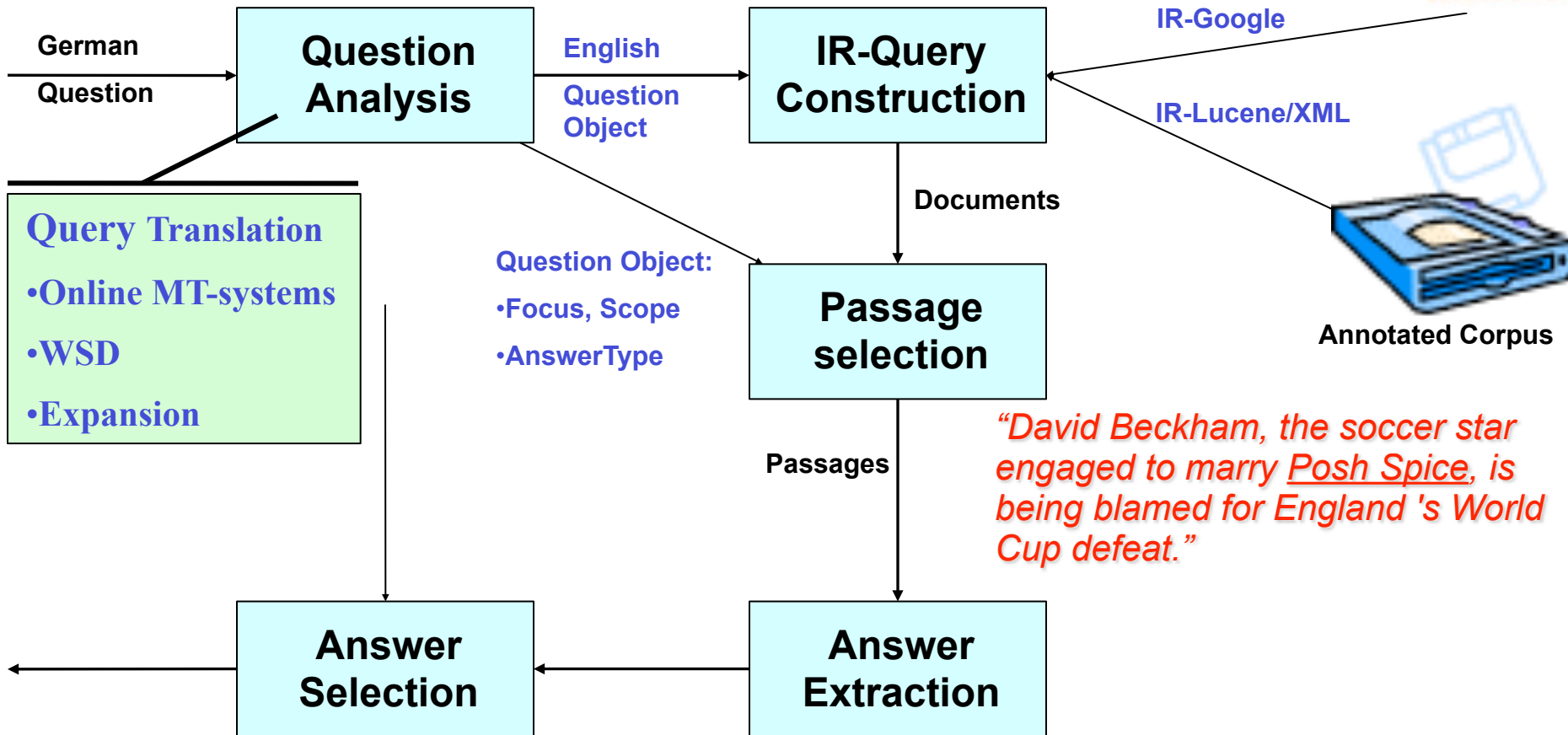




# Open-Domain Question-Answering

*“Mit wem ist David Beckham verheiratet?”*

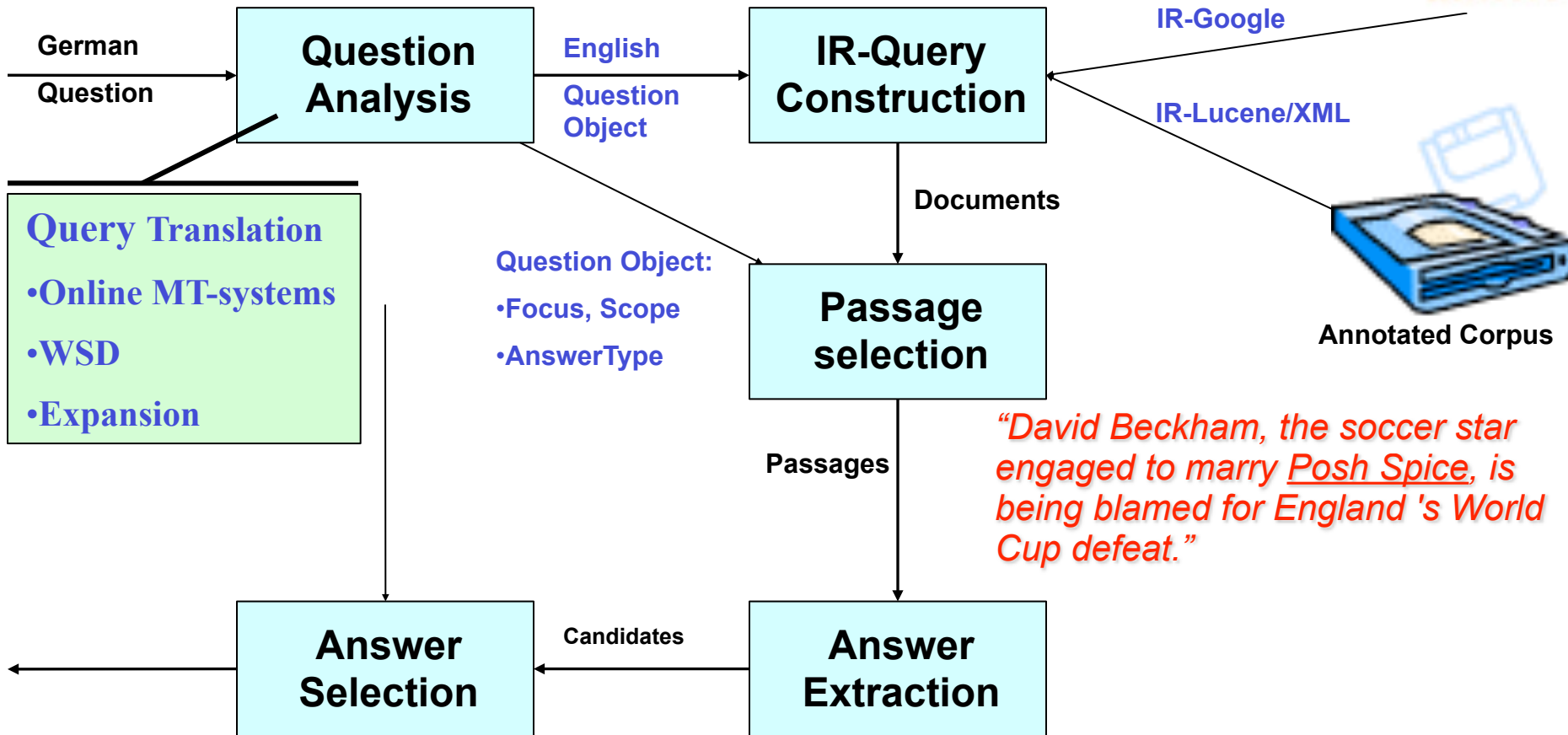
*{person:David Beckham, married, person:?}*



# Open-Domain Question-Answering

*“Mit wem ist David Beckham verheiratet?”*

*{person:David Beckham, married, person:?}*

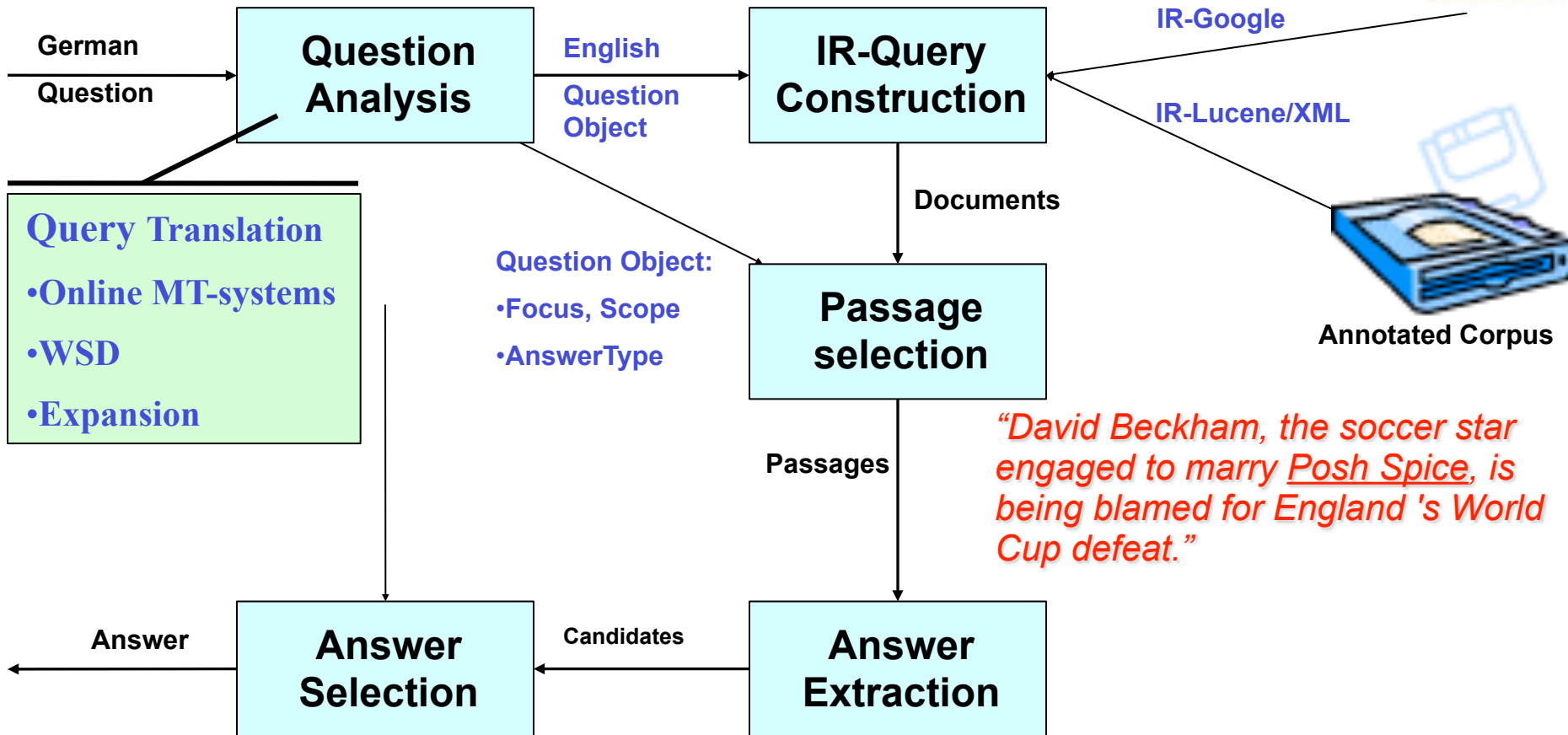


*{person:David Beckham, person:Posh<sup>7</sup>Spice}*

# Open-Domain Question-Answering

*“Mit wem ist David Beckham verheiratet?”*

*{person:David Beckham, married, person:?}*



*Posh Spice*

©, PD Dr. Günter Neumann, Lt-lab, DFKI

*{person:David Beckham, person:Posh<sup>7</sup> Spice}*

# Open-Domain Question Answering: Some Details

---

- Multi-dimensional annotation of unstructured text corpora
- Cross-lingual query processing
- Clef-2004 participation

# Why multi-dimensional annotation of un-structured text?

- The assumption is that a **structural analysis** of un-structured texts **towards** the type of information that can be the **focus of questions**, will support the retrieval of relevant small textual information units through **informative IR-queries**.
  - From candidate document retrieval to candidate answer retrieval.
- However, since we cannot foresee all the different user's interests/questions, a **challenging research question** is:
  - How detailed can the structural analysis be made without putting over a "straitjacket" of a particular view on the un-structured source?
- The assumptions here are:
  - Questions and answers are somewhat related ("questions influence the information geometry and hence, the information view and access", see also Rijsbergen, 2004)
  - There is a bias between off-line and on-line answer extraction.

# Some initial experiments

We have performed some experiments focusing on the relationship between the size of information units and answer containment (using the QA-test set from Clef-2003).

<b>#N</b> <b>Unit-Type</b>	<b>1</b>	<b>5</b>	<b>10</b>	<b>20</b>	<b>30</b>	<b>40</b>	<b>50</b>	<b>100</b>
<i>Sentences*</i>	37.9	58.2	65.8	69.6	70.8	72.1	74	75.9
<i>Sentences</i>	28.4	53.1	60.1	67	70.2	72.7	72.7	74.6
<i>Passages*</i>	39.8	63.2	68.3	73.4	74	75.3	76.5	77.8
<i>Passages</i>	31.6	60.7	67.7	71.5	74.6	77.2	77.2	80.3
<i>Documents*</i>	47.4	69.6	76.5	80.3	81.6	82.9	82.9	83.5
<i>Documents</i>	46.2	68.3	77.8	82.2	82.2	83.5	84.1	85.4

As a result we hypothesized that it is reasonable to use NE-annotated sentences as major retrieval units for the IR-engine

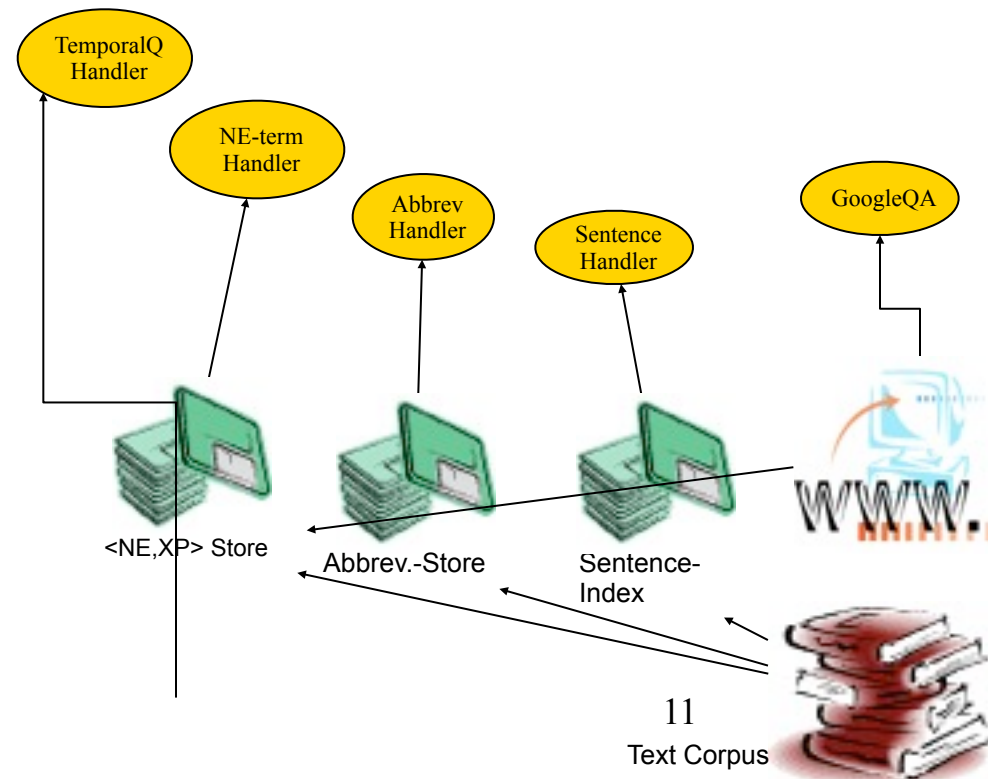
⇒

Simplified answer extraction process & no need of special passage extraction methods

Precision of retrieval for different unit types and top N units retrieved, namely documents, passages, sentences – and their NE-annotated correspondents (marked by \*).

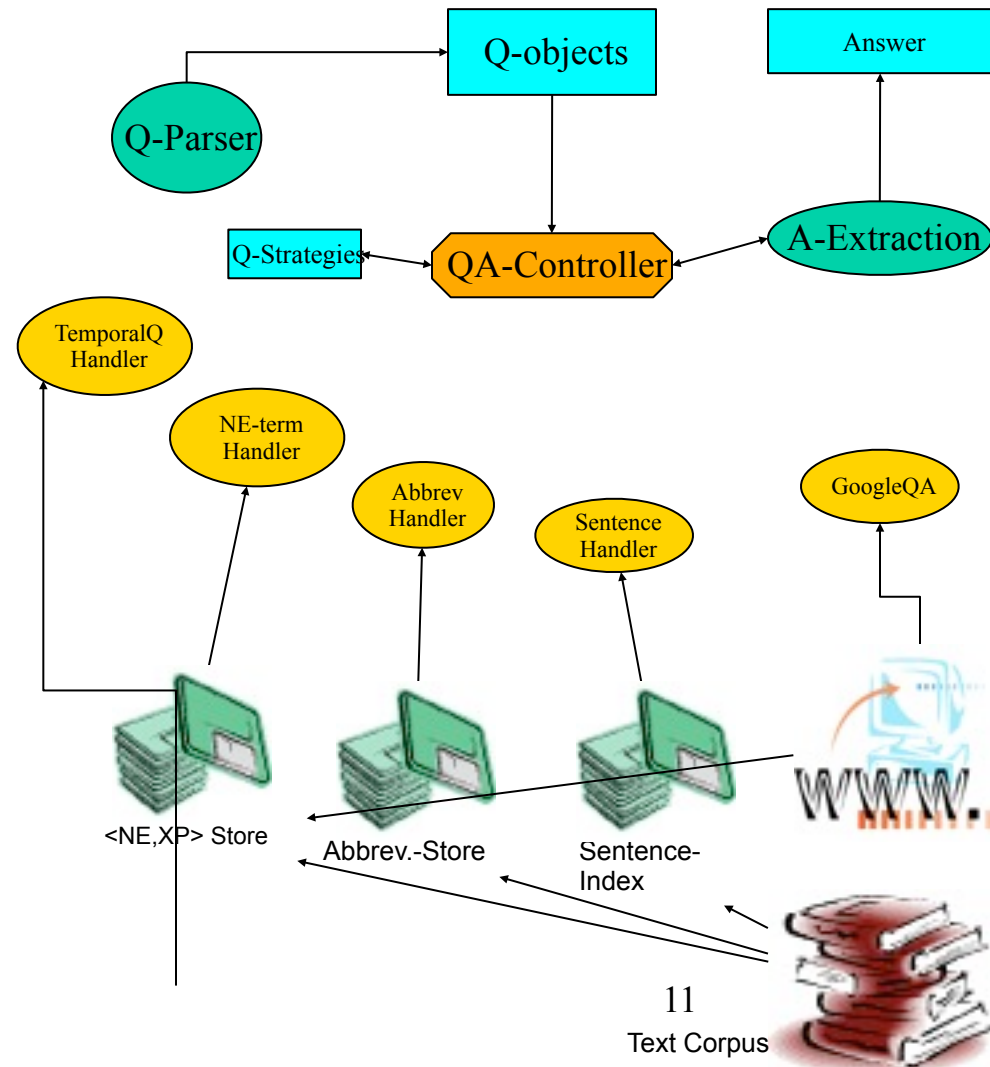
# Open-domain Question Answering: Multi-dimensional annotation

- Idea: off-line annotation of the data collection, which support
  - Query-specific indexing (Q-strategies), and
  - Answer extraction
- Sentence-level pre-processing proved valuable
  - Sentences-boundary
  - Named Entity + Co-reference
  - Abbreviations
  - NE-anchored tuples



# Open-domain Question Answering: Multi-dimensional annotation

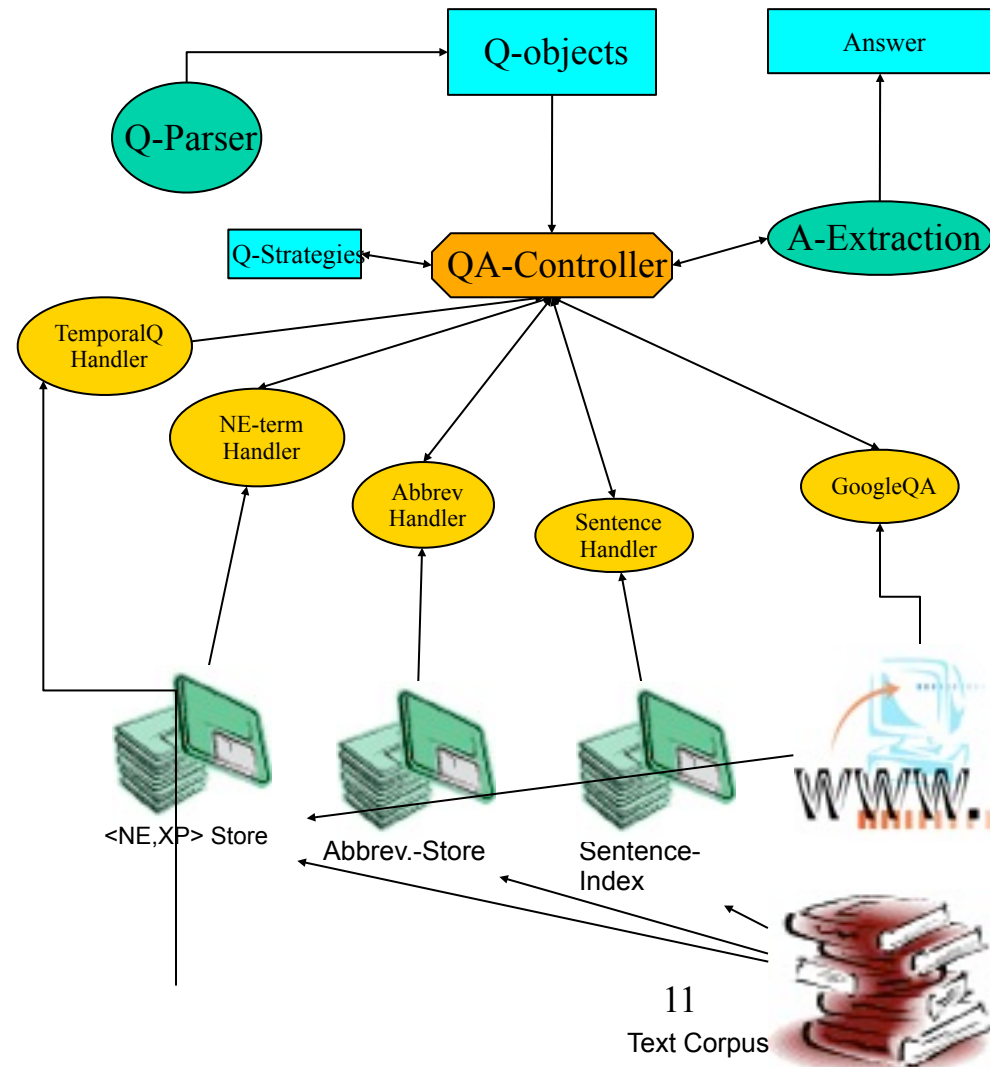
- Idea: off-line annotation of the data collection, which support
  - Query-specific indexing (Q-strategies), and
  - Answer extraction
- Sentence-level pre-processing proved valuable
  - Sentences-boundary
  - Named Entity + Co-reference
  - Abbreviations
  - NE-anchored tuples





# Open-domain Question Answering: Multi-dimensional annotation

- Idea: off-line annotation of the data collection, which support
  - Query-specific indexing (Q-strategies), and
  - Answer extraction
- Sentence-level pre-processing proved valuable
  - Sentences-boundary
  - Named Entity + Co-reference
  - Abbreviations
  - NE-anchored tuples



# Why Query Analysis for Open-Domain?

---

- In principle possible:
  - No query analysis necessary
  - Send to IR (e.g, Lucene or Google) NL-Query as it is
- However
  - Would not allow to control IR through specific search queries
  - Answer extraction & selection not possible

# Query Analysis

---

- Hence, we want to compute an internal query object, to support
  - Construction of “meaningful” IR-search queries, e.g.,
    - » “What is a battery?” -> “define:battery”
    - » “Wer hat das Schloss Charlottenburg erbaut?” -> “Schloss Charlottenburg”+ erbauen erbaute erbaut
    - » “Name a scientist who won the Nobel Price in physics.” ->  
neType:PERSON text:scientist text:won text:”Nobel Price”+ text:physics+
  - Control/guide answer selection
    - » Type of question (e.g., definition, completion)
    - » Expected answer type (e.g., PERSON, LOCATION, ...)

# Why is it Difficult?

---

- Open-domain:
  - Should not rely on sophisticated use of knowledge bases
  - Should also work, even if input is non-well-formed and in-complete
  - Nevertheless should help to direct search engine and answer processing

# Robust Interpretation of NL Queries

- German syntax (SMES):
  - Topological parsing
  - Local Subgrammars for Wh-phrases
- Re-representation
  - Distributed Representation for Dependency Structure
- Query analysis
  - Major information
    - » Q-type (description/definition/...)
    - » A-type (Person/Location/...)
    - » Scope (further constraints for A-type)
  - Q-type determination using Wh-meta terms
    - » “What **type** of bridge is the Golden Gate Bridge?”
  - Corpus-driven approach for Wh-domain terms
    - » “What is the **capital** of Somalia?”
- Determines control-information for QA-strategy selection

# Examples (and more)

---

„Was für eine Art Tier ist der Hund? "

↙  
<IOOBJ msg='quest' q-type='C-HYPONYM' q-weight='1.0'>

<A-TYPE>ANIMAL</A-TYPE>

<SCOPE>hund</SCOPE>

...

“In welcher Stadt lebte Picasso?”

↙  
<IOOBJ msg='quest' q-type='C-COMPLETION' q-weight='1.0'>

<A-TYPE>LOCATION</A-TYPE>

<SCOPE>stadt</SCOPE>

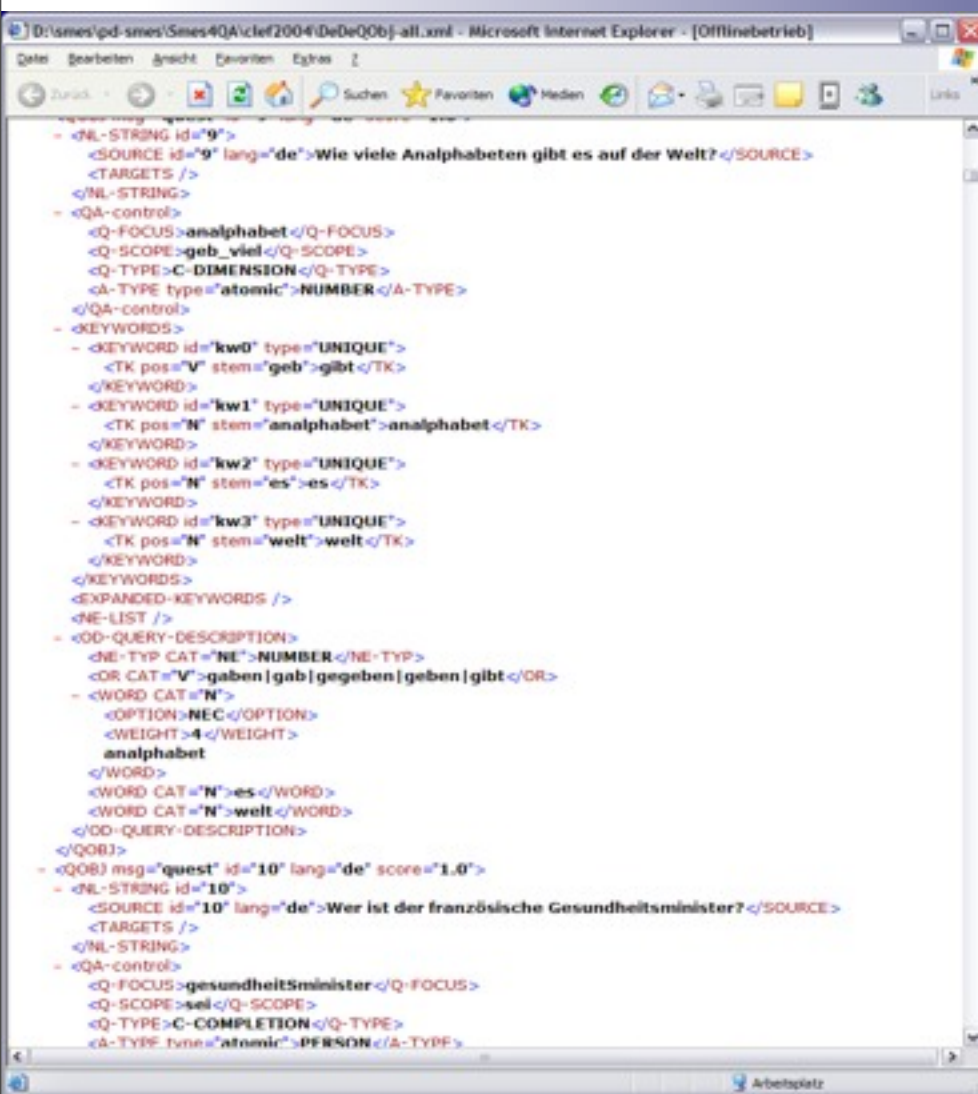
# Open-Domain NL-Query Analysis

The screenshot shows a Microsoft Internet Explorer window displaying XML analysis results for two German queries. The first query is "Wie viele Analphabeten gibt es auf der Welt?" and the second is "Wer ist der französische Gesundheitsminister?". The XML output includes metadata like source and language, QA-control information, keywords, and expanded keywords.

```
<NL-STRING id="9">
  <SOURCE id="9" lang="de">Wie viele Analphabeten gibt es auf der Welt?</SOURCE>
  <TARGETS />
</NL-STRING>
<QA-control>
  <Q-FOCUS>analphabet</Q-FOCUS>
  <Q-SCOPE>geb_viel</Q-SCOPE>
  <Q-TYPE>C-DIMENSION</Q-TYPE>
  <A-TYPE type="atomic">NUMBER</A-TYPE>
</QA-control>
<KEYWORDS>
  <KEYWORD id="kw0" type="UNIQUE">
    <TK pos="V" stem="geb">gibt</TK>
  </KEYWORD>
  <KEYWORD id="kw1" type="UNIQUE">
    <TK pos="N" stem="analphabet">analphabet</TK>
  </KEYWORD>
  <KEYWORD id="kw2" type="UNIQUE">
    <TK pos="N" stem="es">es</TK>
  </KEYWORD>
  <KEYWORD id="kw3" type="UNIQUE">
    <TK pos="N" stem="welt">welt</TK>
  </KEYWORD>
</KEYWORDS>
<EXPANDED-KEYWORDS />
<NE-LIST />
<OD-QUERY-DESCRIPTION>
  <NE-TYP CAT="NE">NUMBER</NE-TYP>
  <OR CAT="V">gaben|gab|gegeben|geben|gibt</OR>
  <WORD CAT="N">
    <OPTION>REC</OPTION>
    <WEIGHT>4</WEIGHT>
    analphabet
  </WORD>
  <WORD CAT="N">es</WORD>
  <WORD CAT="N">welt</WORD>
</OD-QUERY-DESCRIPTION>
</QOR>
<QOR msg="quest" id="10" lang="de" score="1.0">
  <NL-STRING id="10">
    <SOURCE id="10" lang="de">Wer ist der französische Gesundheitsminister?</SOURCE>
    <TARGETS />
  </NL-STRING>
  <QA-control>
    <Q-FOCUS>gesundheitsminister</Q-FOCUS>
    <Q-SCOPE>wel</Q-SCOPE>
    <Q-TYPE>C-COMPLETION</Q-TYPE>
    <A-TYPE type="atomic">PERSON</A-TYPE>
  </QA-control>
</QOR>
```

# Open-Domain NL-Query Analysis

“How many illiterates are there in the world?”



```
D:\smes\pd-smes\smes4QA\clef2004\DeDeQ0b]-all.xml - Microsoft Internet Explorer - [Offlinebetrieb]
Datei Bearbeiten Ansicht Favoriten Extras Z
Zurück Suchen Favoriten Medien
- <NL-STRING id="9">
  <SOURCE id="9" lang="de">Wie viele Analphabeten gibt es auf der Welt?</SOURCE>
  <TARGETS />
  </NL-STRING>
- <QA-control>
  <Q-FOCUS>analphabet</Q-FOCUS>
  <Q-SCOPE>geb_viel</Q-SCOPE>
  <Q-TYPE>C-DIMENSION</Q-TYPE>
  <A-TYPE type="atomic">NUMBER</A-TYPE>
  </QA-control>
- <KEYWORDS>
  <KEYWORD id="kw0" type="UNIQUE">
    <TK pos="V" stem="geb">gibt</TK>
  </KEYWORD>
  <KEYWORD id="kw1" type="UNIQUE">
    <TK pos="N" stem="analphabet">analphabet</TK>
  </KEYWORD>
  <KEYWORD id="kw2" type="UNIQUE">
    <TK pos="N" stem="es">es</TK>
  </KEYWORD>
  <KEYWORD id="kw3" type="UNIQUE">
    <TK pos="N" stem="welt">welt</TK>
  </KEYWORD>
  </KEYWORDS>
  <EXPANDED-KEYWORDS />
  <NE-LIST />
- <OD-QUERY-DESCRIPTION>
  <NE-TYP CAT="NE">NUMBER</NE-TYP>
  <OR CAT="V">gaben|gab|gegeben|geben|gibt</OR>
  <WORD CAT="N">
    <OPTION>REC</OPTION>
    <WEIGHT>4</WEIGHT>
    analphabet
  </WORD>
  <WORD CAT="N">es</WORD>
  <WORD CAT="N">welt</WORD>
  </OD-QUERY-DESCRIPTION>
  </QOR>
- <QOR> msg="quest" id="10" lang="de" score="1.0"
- <NL-STRING id="10">
  <SOURCE id="10" lang="de">Wer ist der französische Gesundheitsminister?</SOURCE>
  <TARGETS />
  </NL-STRING>
- <QA-control>
  <Q-FOCUS>gesundheitsminister</Q-FOCUS>
  <Q-SCOPE>wel</Q-SCOPE>
  <Q-TYPE>C-COMPLETION</Q-TYPE>
  <A-TYPE type="atomic">PERSON</A-TYPE>
```



# Open-Domain NL-Query Analysis

“How many illiterates are there in the world?”

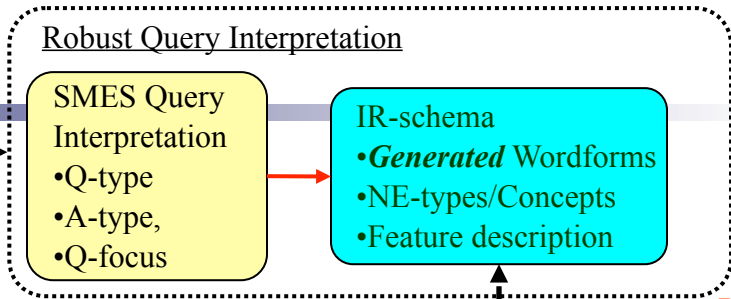
```

<NL-STRING id="9">
<SOURCE id="9" lang="de">Wie viele Analphabeten gibt es auf der Welt?</SOURCE>
<TARGETS />
<NL-STRING>
<QA-control>
<Q-FOCUS>analphabet</Q-FOCUS>
<Q-SCOPE>geb_viel</Q-SCOPE>
<Q-TYPE>C-DIMENSION</Q-TYPE>
<A-TYPE type="atomic">NUMBER</A-TYPE>
</QA-control>
<KEYWORDS>
<KEYWORD id="kw0" type="UNIQUE">
<TK pos="V" stem="geb">gibt</TK>
</KEYWORD>
<KEYWORD id="kw1" type="UNIQUE">
<TK pos="N" stem="analphabet">analphabet</TK>
</KEYWORD>
<KEYWORD id="kw2" type="UNIQUE">
<TK pos="N" stem="es">es</TK>
</KEYWORD>
<KEYWORD id="kw3" type="UNIQUE">
<TK pos="N" stem="welt">welt</TK>
</KEYWORD>
</KEYWORDS>
<EXPANDED-KEYWORDS />
<NE-LIST />
<OD-QUERY-DESCRIPTION>
<NE-TYP CAT="NE">NUMBER</NE-TYP>
<OR CAT="V">gaben|gab|gegeben|geben|gibt</OR>
<WORD CAT="N">
<OPTION>NEC</OPTION>
<WEIGHT>4</WEIGHT>
analphabet
</WORD>
<WORD CAT="N">es</WORD>
<WORD CAT="N">welt</WORD>
</OD-QUERY-DESCRIPTION>
</QOB>
<QOB msg="quest" id="10" lang="de" score="1.0">
<NL-STRING id="10">
<SOURCE id="10" lang="de">Wer ist der französische Gesundheitsminister?</SOURCE>
<TARGETS />
<NL-STRING>
<QA-control>
<Q-FOCUS>gesundheitsminister</Q-FOCUS>
<Q-SCOPE>wel</Q-SCOPE>
<Q-TYPE>C-COMPLETION</Q-TYPE>
<A-TYPE type="atomic">PERSON</A-TYPE>

```

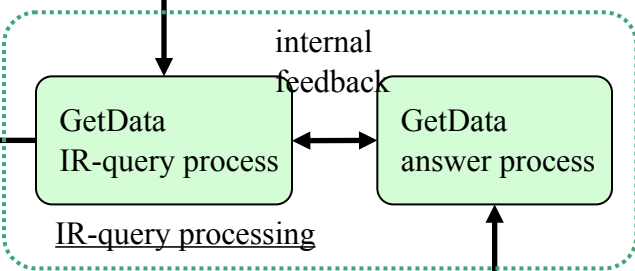
NL query

SMES Robust Parser  
•MorphoSyntax  
•Full Dependency trees

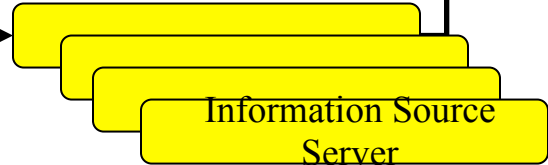


IR-schema

(external feedback)



GetData IR-Query



# Open-Domain NL-Query Analysis

“How many illiterates are there in the world?”

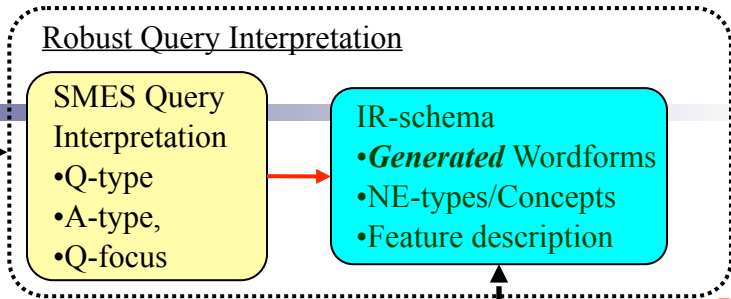
```

<NL-STRING id="9">
<SOURCE id="9" lang="de">Wie viele Analphabeten gibt es auf der Welt?</SOURCE>
<TARGETS />
<NL-STRING />
<QA-control>
<Q-FOCUS>analphabet</Q-FOCUS>
<Q-SCOPE>geb_viel</Q-SCOPE>
<Q-TYPE>C-DIMENSION</Q-TYPE>
<A-TYPE type="atomic">NUMBER</A-TYPE>
</QA-control>
<KEYWORDS />
<KEYWORD id="kw0" type="UNIQUE">
[[W=V, OR], ["gaben", "gab", "gegeben", "geben", "gibt"]],
[[W=N, NEC, 4], ["analphabet"]],
[[W=N], ["welt"]], [NE=Number]
</KEYWORD />
</NL-STRING />
<NE-LIST />
<OD-QUERY-DESCRIPTION>
<NE-TYP CAT="NE">NUMBER</NE-TYP>
<OR CAT="V">gaben|gab|gegeben|geben|gibt</OR>
<WORD CAT="N">
<OPTION>NEC</OPTION>
<WEIGHT>4</WEIGHT>
analphabet
</WORD>
<WORD CAT="N">es</WORD>
<WORD CAT="N">welt</WORD>
</OD-QUERY-DESCRIPTION>
</QOB>
<QOB msg="quest" id="10" lang="de" score="1.0">
<NL-STRING id="10">
<SOURCE id="10" lang="de">Wer ist der französische Gesundheitsminister?</SOURCE>
<TARGETS />
<NL-STRING />
<QA-control>
<Q-FOCUS>gesundheitsminister</Q-FOCUS>
<Q-SCOPE>sel</Q-SCOPE>
<Q-TYPE>C-COMPLETION</Q-TYPE>
<A-TYPE type="atomic">PERSON</A-TYPE>

```

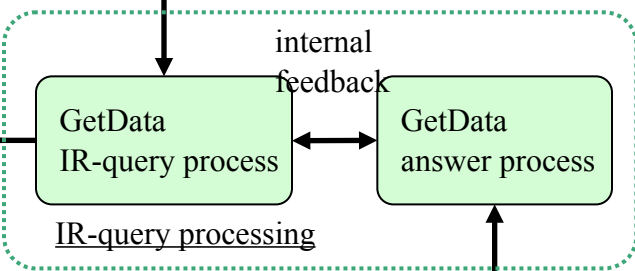
NL query

SMES Robust Parser  
 •MorphoSyntax  
 •Full Dependency trees

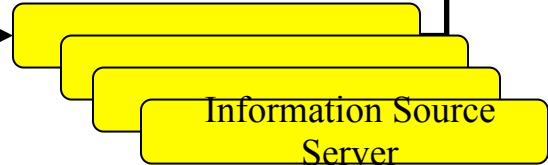


IR-schema

(external feedback)



GetData IR-Query



# Open-Domain NL-Query Analysis

“How many illiterates are there in the world?”

*NL-generated:* Word forms (morpho-syntactic paraphrases)

[[W=V, OR], ["gaben", "gab", "gegeben", "geben", "gibt"]],  
 [[W=N, NEC, 4], ["analphabet"]],  
 [[W=N], ["welt"]], [NE=Number]

NL query

SMES Robust Parser  
 •MorphoSyntax  
 •Full Dependency trees

Robust Query Interpretation

SMES Query Interpretation  
 •Q-type  
 •A-type,  
 •Q-focus

IR-schema  
 •Generated Wordforms  
 •NE-types/Concepts  
 •Feature description

IR-schema

(external feedback)

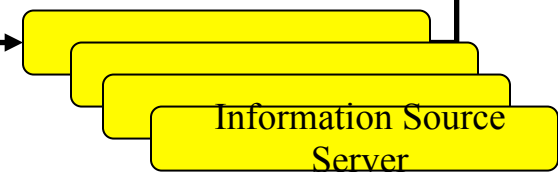
internal feedback

GetData IR-query process

GetData answer process

IR-query processing

GetData IR-Query



Information Source Server

# Open-Domain NL-Query Analysis

“How many illiterates are there in the world?”

*NL-generated:* Word forms (morpho-syntactic paraphrases)

[[W=V, OR], ["gaben", "gab", "gegeben", "geben", "gibt"]],  
 [[W=N, NEC, 4], ["analphabet"]],  
 [[W=N], ["welt"]], [NE=Number]

NL query

SMES Robust Parser  
 •MorphoSyntax  
 •Full Dependency trees

Robust Query Interpretation

SMES Query Interpretation  
 •Q-type  
 •A-type,  
 •Q-focus

IR-schema  
 •Generated Wordforms  
 •NE-types/Concepts  
 •Feature description

IR-schema

(external feedback)

internal feedback

GetData IR-query process

GetData answer process

IR-query processing

“neTypes:Number AND (gaben OR gab OR gegeben OR geben OR gibt) AND analphabet^4 AND welt”

Information Source Server

# Open-Domain NL-Query Analysis

“How many illiterates are there in the world?”

*NL-generated:* Word forms (morpho-syntactic paraphrases)

[[W=V, OR], ["gaben", "gab", "gegeben", "geben", "gibt"]],  
 [[W=N, NEC, 4], ["analphabet"]],  
 [[W=N], ["welt"]], [NE=Number]

NL query

SMES Robust Parser  
 •MorphoSyntax  
 •Full Dependency trees

Robust Query Interpretation

SMES Query Interpretation  
 •Q-type  
 •A-type,  
 •Q-focus

IR-schema  
 •Generated Wordforms  
 •NE-types/Concepts  
 •Feature description

IR-schema

(external feedback)

internal feedback

GetData IR-query process

GetData answer process

IR-query processing

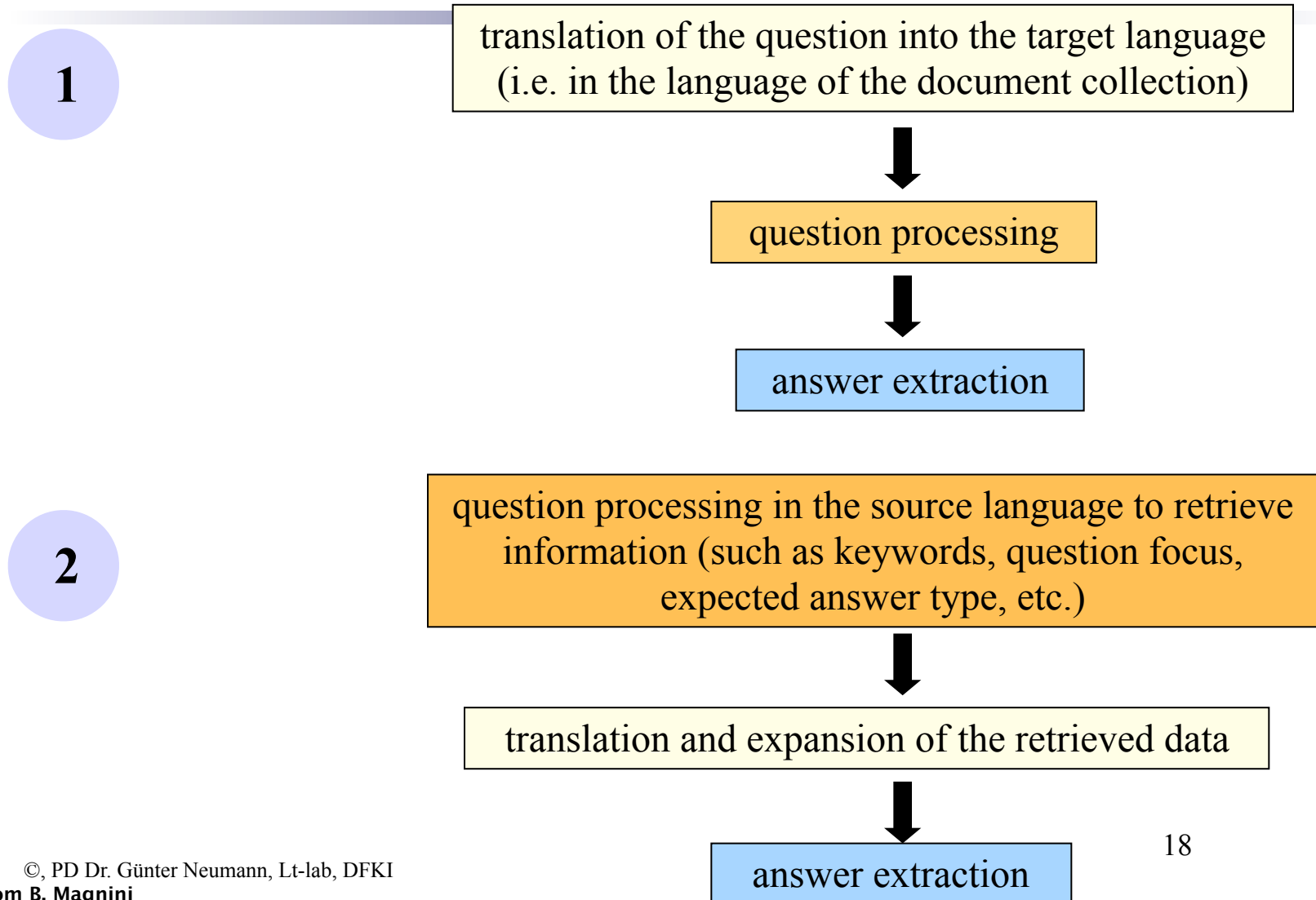
“neTypes:Number AND (gaben OR gab OR gegeben OR geben OR gibt) AND analphabet^4 AND welt”

“neTypes:Number OR (gaben OR gab OR gegeben OR geben OR gibt) OR +analphabet^4 OR welt”

Information Source Server

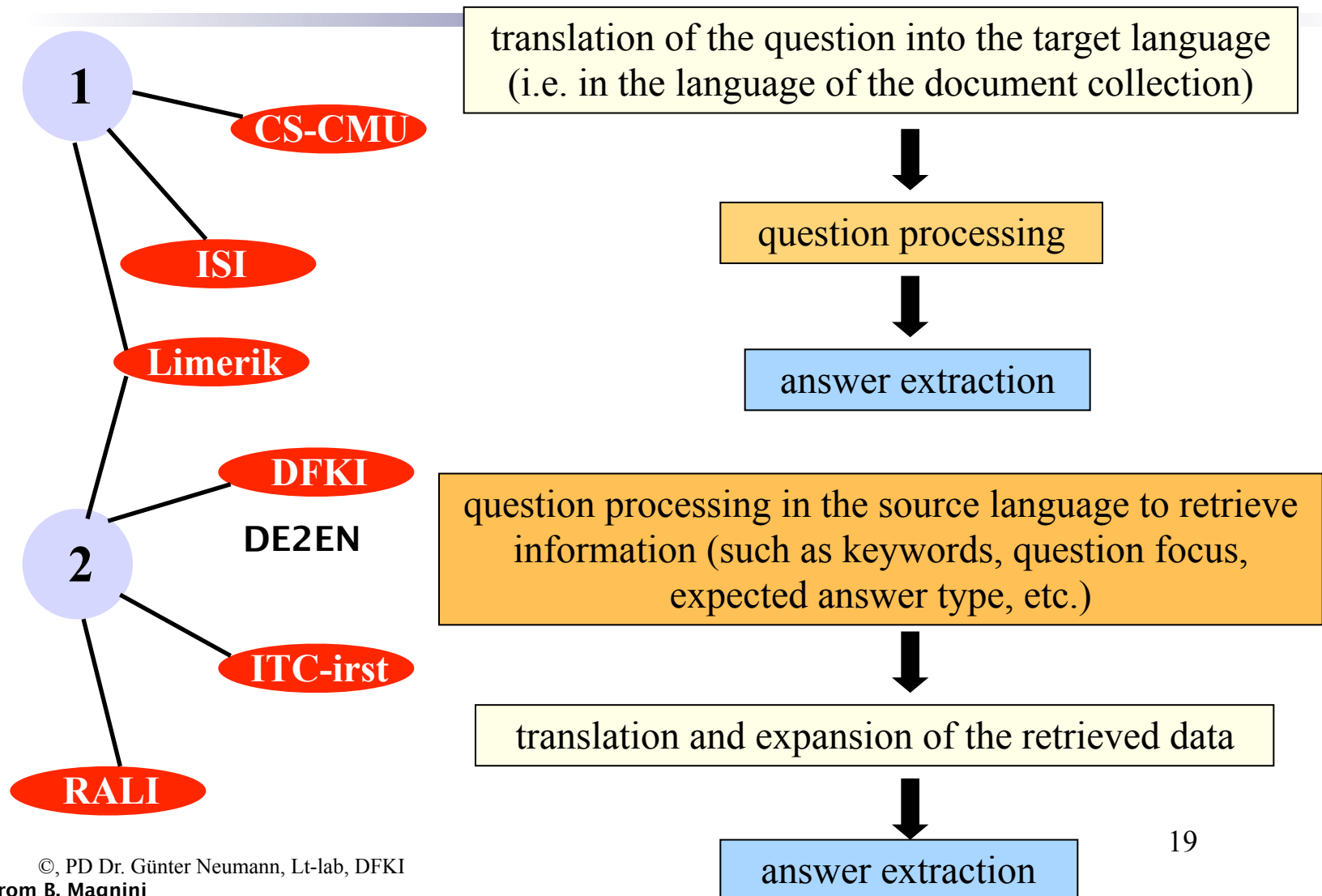
# Approaches in CL QA

Two main different approaches used in Cross-Language QA systems:



# Approaches in CL QA

Two main different approaches used in Cross-Language QA systems:



# Query Translation & Expansion

---

- First idea:
  - Only use EuroWordNet
  - Defines a word-based translation via synset offsets
- Experience
  - EuroWordNet too sparse on German side
  - Nevertheless introduced too much ambiguity
  - NE-translation is crucial
- So far, not very much of help



# Query Translation & Expansion

---

- First idea:
  - Only use EuroWordNet
  - Defines a word-based translation via synset offsets
- Experience
  - EuroWordNet too sparse on German side
  - Nevertheless introduced too much ambiguity
  - NE-translation is crucial
- So far, not very much of help
- Second idea:

# Query Translation & Expansion

---

- First idea:
  - Only use EuroWordNet
  - Defines a word-based translation via synset offsets
- Experience
  - EuroWordNet too sparse on German side
  - Nevertheless introduced too much ambiguity
  - NE-translation is crucial
- So far, not very much of help
- Second idea:
  - Use EuroWordNet

# Query Translation & Expansion

---

- First idea:
  - Only use EuroWordNet
  - Defines a word-based translation via synset offsets
- Experience
  - EuroWordNet too sparse on German side
  - Nevertheless introduced too much ambiguity
  - NE-translation is crucial
- So far, not very much of help
- Second idea:
  - Use EuroWordNet
  - Use external MT-services

# Query Translation & Expansion

---

- First idea:
  - Only use EuroWordNet
  - Defines a word-based translation via synset offsets
- Experience
  - EuroWordNet too sparse on German side
  - Nevertheless introduced too much ambiguity
  - NE-translation is crucial
- So far, not very much of help
- Second idea:
  - Use EuroWordNet
  - Use external MT-services
  - Overlap-mechanism for query expansion

# Query Translation & Expansion

---

- First idea:
  - Only use EuroWordNet
  - Defines a word-based translation via synset offsets
- Experience
  - EuroWordNet too sparse on German side
  - Nevertheless introduced too much ambiguity
  - NE-translation is crucial
- So far, not very much of help
- Second idea:
  - Use EuroWordNet
  - Use external MT-services
  - Overlap-mechanism for query expansion
- Crosslingual because

# Query Translation & Expansion

- First idea:
  - Only use EuroWordNet
  - Defines a word-based translation via synset offsets
- Experience
  - EuroWordNet too sparse on German side
  - Nevertheless introduced too much ambiguity
  - NE-translation is crucial
- So far, not very much of help
- Second idea:
  - Use EuroWordNet
  - Use external MT-services
  - Overlap-mechanism for query expansion
- Crosslingual because
  - Q-type & A-type from DE-Question Analysis

# Query Translation & Expansion

- First idea:
  - Only use EuroWordNet
  - Defines a word-based translation via synset offsets
- Experience
  - EuroWordNet too sparse on German side
  - Nevertheless introduced too much ambiguity
  - NE-translation is crucial
- So far, not very much of help
- Second idea:
  - Use EuroWordNet
  - Use **external** MT-services
  - Overlap-mechanism for query expansion
- Crosslingual because
  - Q-type & A-type from DE-Question Analysis
  - Synsets from EuroWN direct query expansion (**online alignment**)

# Example

## 1. Translation services for Word Sense Disambiguation (WSD)

---

Wo wurde das Militärflugzeug Strike Eagles 1990 **eingesetzt**?

## 2. Query expansion using EuroWordNet



# Example

## 1. Translation services for Word Sense Disambiguation (WSD)

---

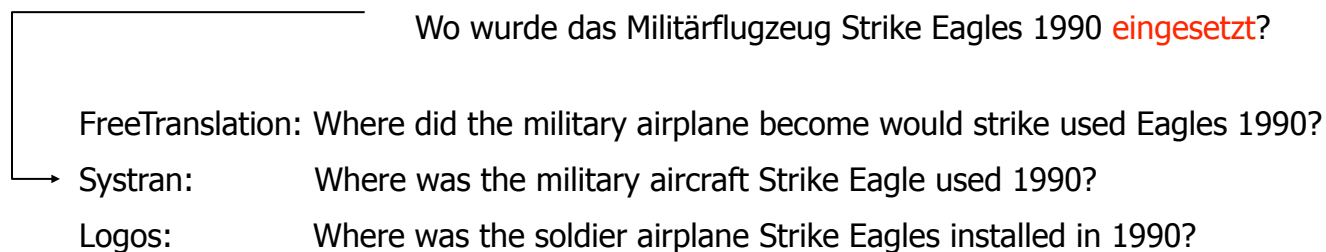


Wo wurde das Militärflugzeug Strike Eagles 1990 eingesetzt?

## 2. Query expansion using EuroWordNet

# Example

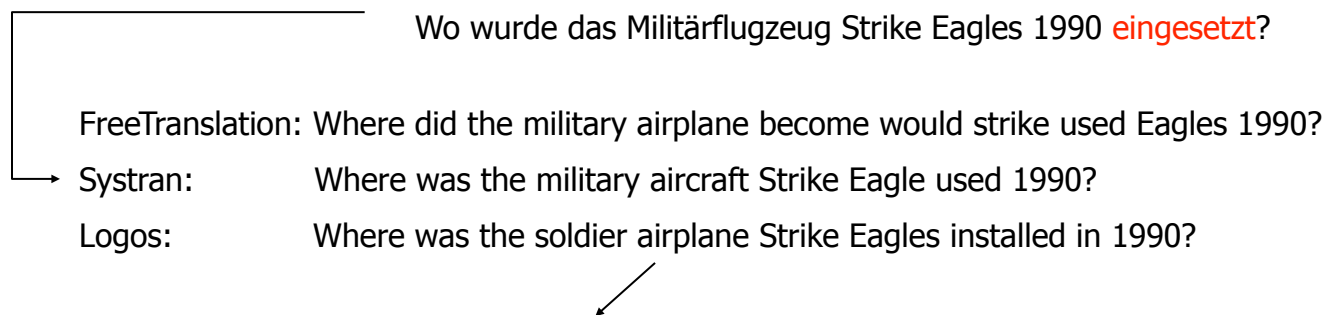
## 1. Translation services for Word Sense Disambiguation (WSD)



## 2. Query expansion using EuroWordNet

# Example

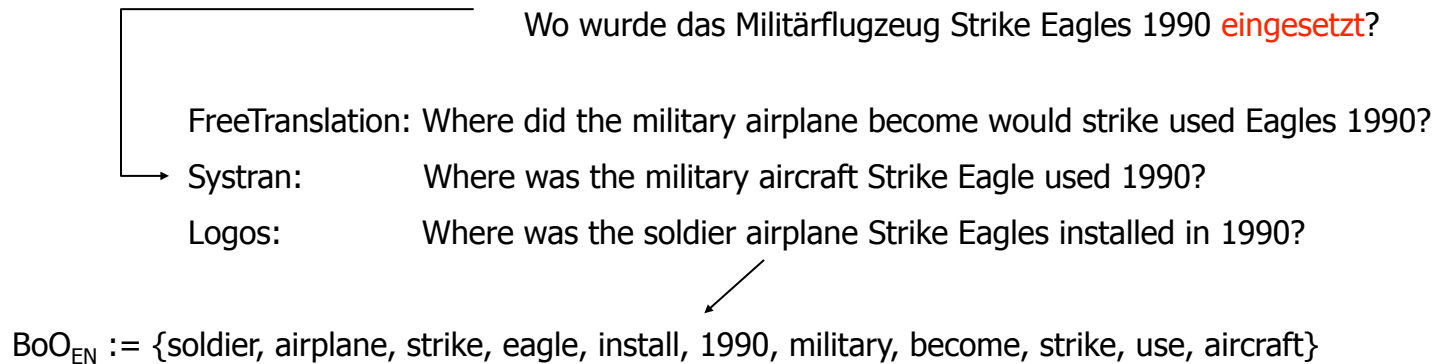
## 1. Translation services for Word Sense Disambiguation (WSD)



## 2. Query expansion using EuroWordNet

# Example

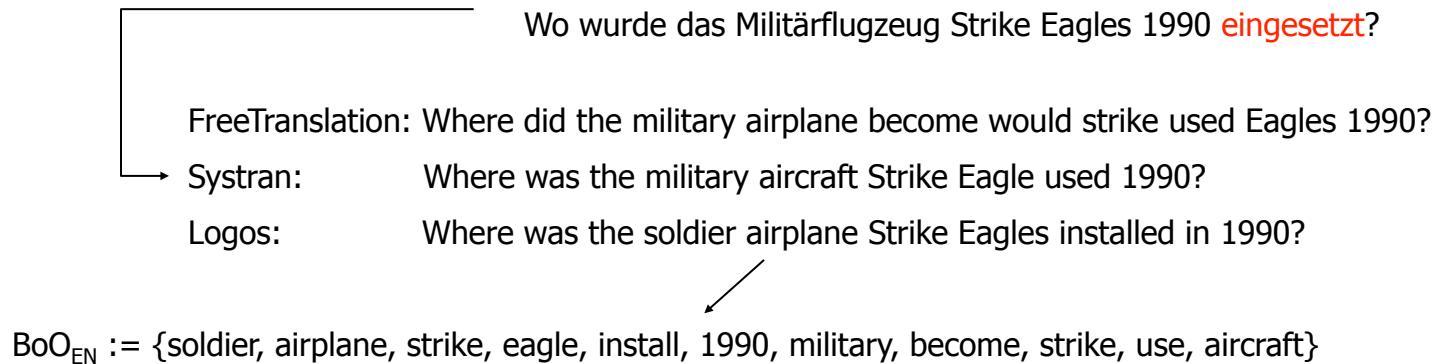
## 1. Translation services for Word Sense Disambiguation (WSD)



## 2. Query expansion using EuroWordNet

# Example

## 1. Translation services for Word Sense Disambiguation (WSD)

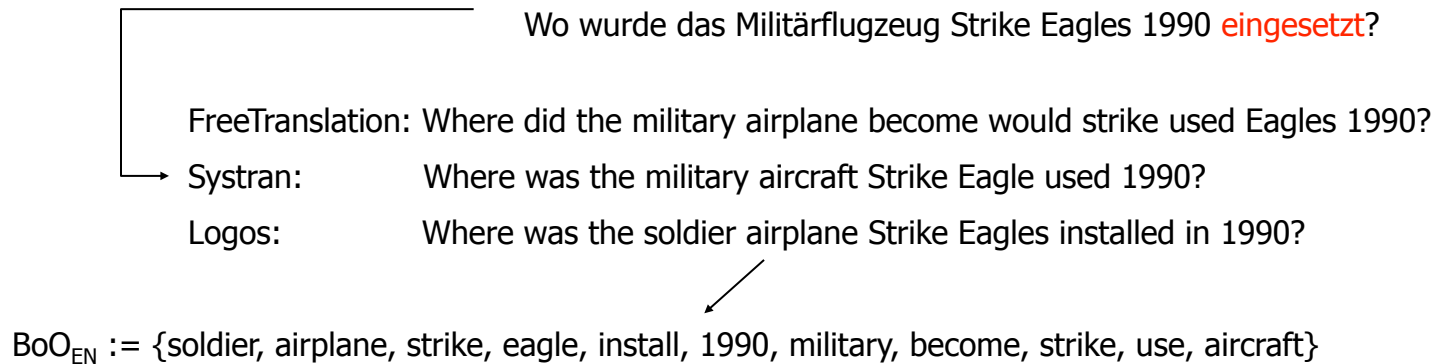


## 2. Query expansion using EuroWordNet

$\forall x \in \text{BoO}_{\text{EN}}$ : lookup(EuroWN);  
If  $x$  is unambiguous: extend  $\text{BoO}_{\text{EN}}$   
Else  $\forall \text{readings}(x)$ :  
get its aligned German readings &  
Look them up in  $\text{BoO}_{\text{GN}}$   
If successfully then add English terms to  
 $\text{BoO}_{\text{EN}}$

# Example

## 1. Translation services for Word Sense Disambiguation (WSD)



## 2. Query expansion using EuroWordNet

$\forall x \in \text{BoO}_{\text{EN}}$ : lookup(EuroWN);  
If  $x$  is unambiguous: extend  $\text{BoO}_{\text{EN}}$   
Else  $\forall \text{readings}(x)$ :  
get its aligned German readings &  
Look them up in  $\text{BoO}_{\text{GN}}$   
If successfully then add English terms to  
 $\text{BoO}_{\text{EN}}$

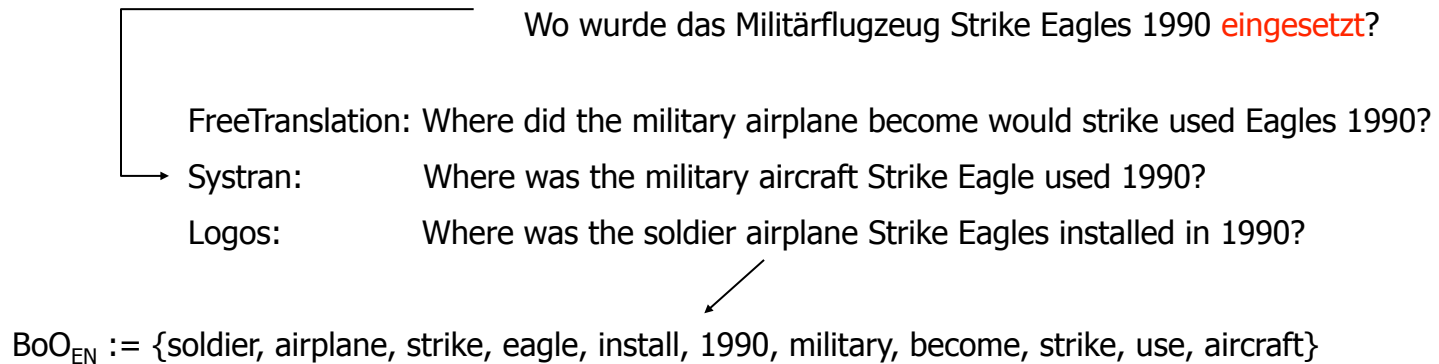
Reading-697925

EN: {handle, use, wield}

DE: {handhaben, hantieren}

# Example

## 1. Translation services for Word Sense Disambiguation (WSD)



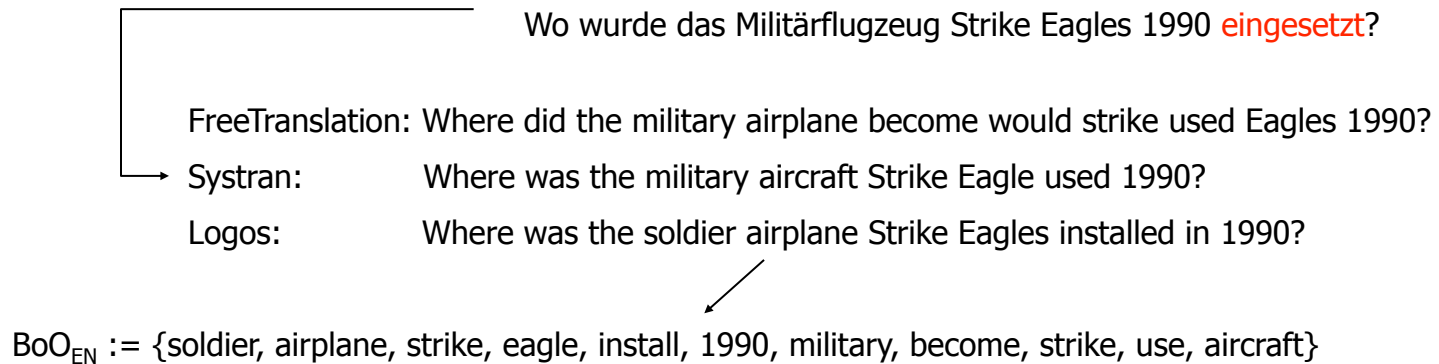
## 2. Query expansion using EuroWordNet

$\forall x \in \text{BoO}_{\text{EN}}$ : lookup(EuroWN);  
If  $x$  is unambiguous: extend  $\text{BoO}_{\text{EN}}$   
Else  $\forall \text{readings}(x)$ :  
get its aligned German readings &  
Look them up in  $\text{BoO}_{\text{GN}}$   
If successfully then add English terms to  
 $\text{BoO}_{\text{EN}}$

~~Reading-697925  
EN: {handle, use, wield}  
DE: {handhaben, hantieren}~~

# Example

## 1. Translation services for Word Sense Disambiguation (WSD)



## 2. Query expansion using EuroWordNet

$\forall x \in \text{BoO}_{\text{EN}}$ : lookup(EuroWN);  
If  $x$  is unambiguous: extend  $\text{BoO}_{\text{EN}}$   
Else  $\forall \text{readings}(x)$ :  
get its aligned German readings &  
Look them up in  $\text{BoO}_{\text{GN}}$   
If successfully then add English terms to  
 $\text{BoO}_{\text{EN}}$

~~Reading-697925~~

~~EN: {handle, use, wield}~~

~~DE: {handhaben, hantieren}~~

Reading-1453934:

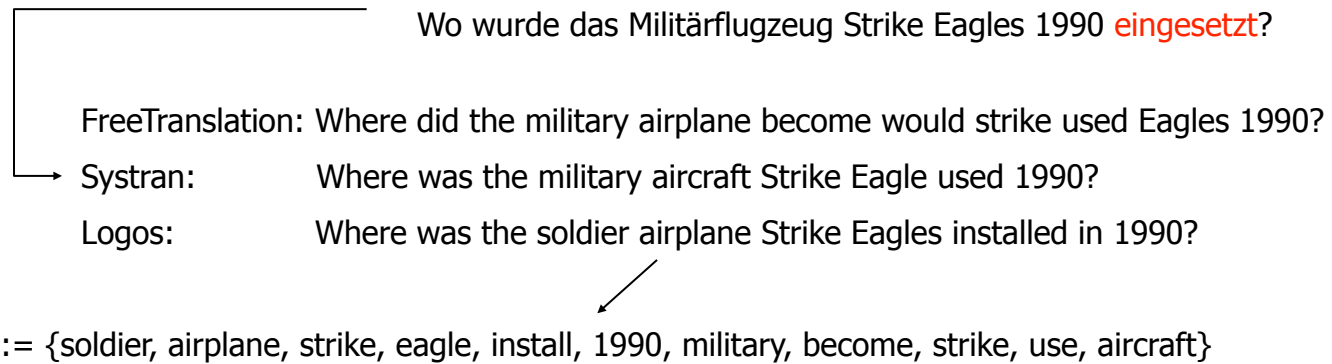
EN: {behave toward, use}

DE: not aligned



# Example

## 1. Translation services for Word Sense Disambiguation (WSD)



## 2. Query expansion using EuroWordNet

$\forall x \in \text{BoO}_{\text{EN}}$ : lookup(EuroWN);  
If  $x$  is unambiguous: extend  $\text{BoO}_{\text{EN}}$   
Else  $\forall \text{readings}(x)$ :  
get its aligned German readings &  
Look them up in  $\text{BoO}_{\text{GN}}$   
If successfully then add English terms to  
 $\text{BoO}_{\text{EN}}$

~~Reading-697925~~

~~EN: {handle, use, wield}~~

~~DE: {handhaben, hantieren}~~

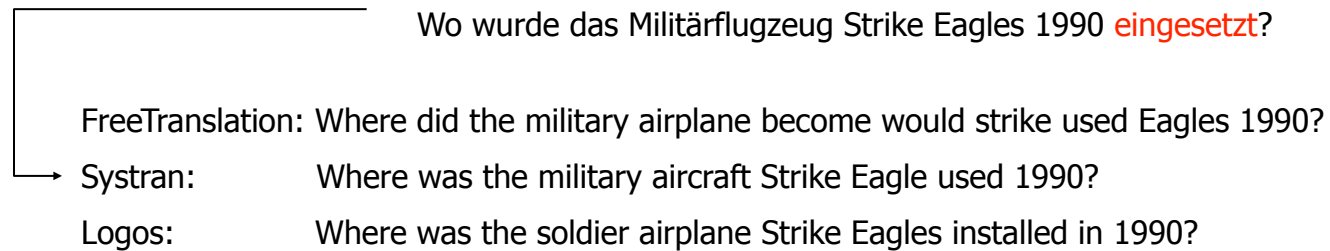
~~Reading-1453934:~~

~~EN: {behave toward, use}~~

~~DE: not aligned~~

# Example

## 1. Translation services for Word Sense Disambiguation (WSD)



$BoO_{EN} := \{\text{soldier, airplane, strike, eagle, install, 1990, military, become, strike, use, aircraft}\}$

## 2. Query expansion using EuroWordNet

$\forall x \in BoO_{EN}$ : lookup(EuroWN);  
If  $x$  is unambiguous: extend  $BoO_{EN}$   
Else  $\forall$  readings( $x$ ):  
get its aligned German readings &  
Look them up in  $BoO_{GN}$   
If successfully then add English terms to  
 $BoO_{EN}$

~~Reading-697925~~

~~EN: {handle, use, wield}~~

~~DE: {handhaben, hantieren}~~

~~Reading-1453934:~~

~~EN: {behave toward, use}~~

~~DE: not aligned~~

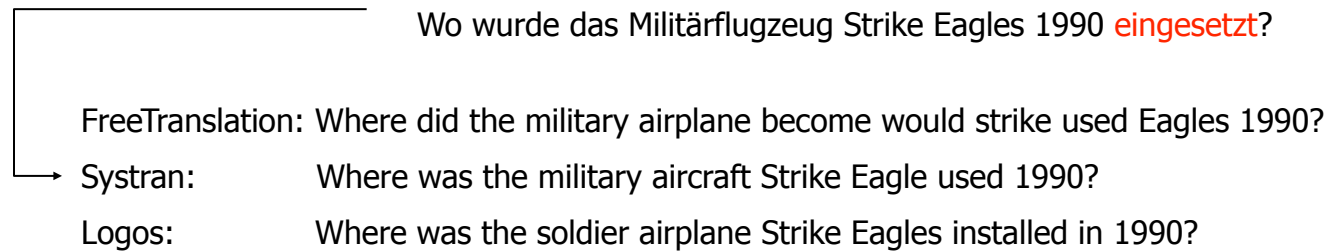
Reading-658243:

EN: {apply, employ, make use of, put to use, use, utilise, utilize}

DE: {anbringen, anwenden, bedienen, benutzen, **einsetzen**, ...}

# Example

## 1. Translation services for Word Sense Disambiguation (WSD)



BoO<sub>EN</sub> := {soldier, airplane, strike, eagle, install, 1990, military, become, strike, use, aircraft}

## 2. Query expansion using EuroWordNet

$\forall x \in \text{BoO}_{\text{EN}}$ : lookup(EuroWN);  
If  $x$  is unambiguous: extend BoO<sub>EN</sub>  
Else  $\forall \text{readings}(x)$ :  
get its aligned German readings &  
Look them up in BoO<sub>GN</sub>  
If successfully then add English terms to  
BoO<sub>EN</sub>

~~Reading-697925~~

~~EN: {handle, use, wield}~~

~~DE: {handhaben, hantieren}~~

~~Reading-1453934:~~

~~EN: {behave toward, use}~~

~~DE: not aligned~~

Reading-658243:

EN: {apply, employ, make use of, put to use, use, utilise, utilize}

DE: {anbringen, anwenden, bedienen, benutzen, **einsetzen**, ...}

# Cross-lingual Open-Domain Question Answering

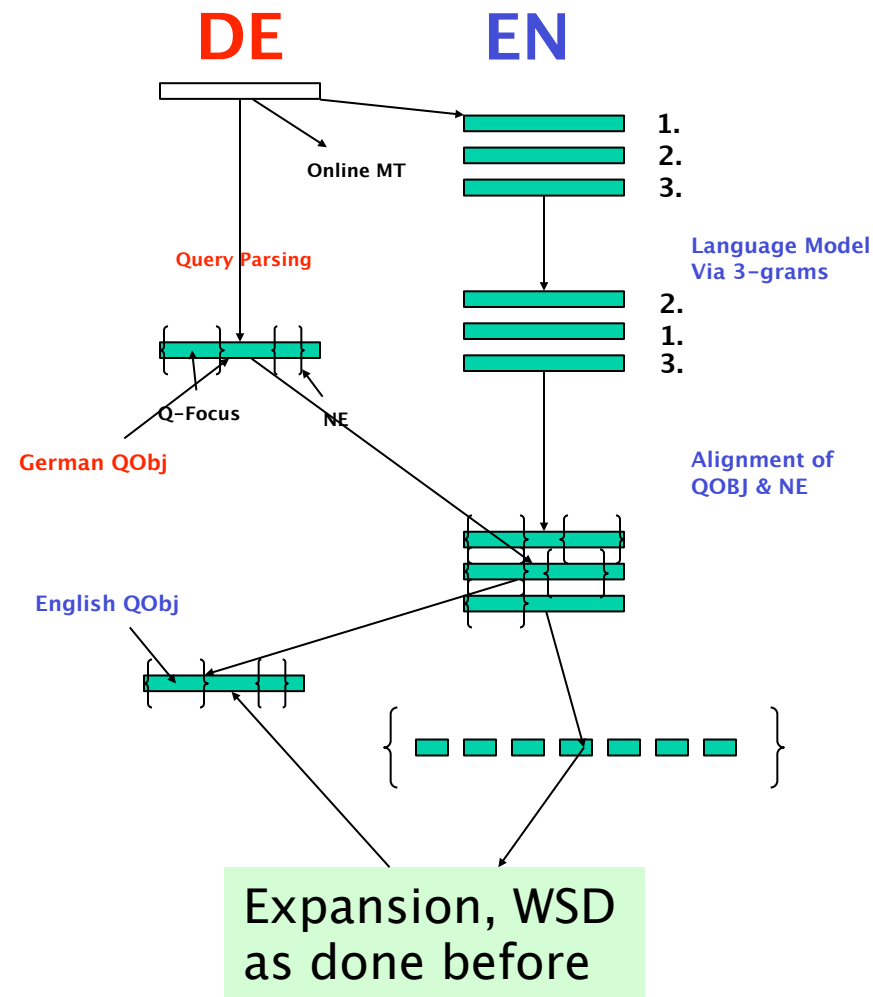
---

- How we have done it so far:
  - Use external MT-services
  - Overlap-mechanism &
  - EuroWordNet for query expansion
- Experience
  - External MT services
    - » used ako WSD
    - » as well as query expansion
  - Reduced degree of ambiguity
- Problems/Restrictions
  - Important information from German QueryObject (e.g., Q-focus, Q-scope) cannot be carried over
    - ⇒ Need of query parsing on English side
    - ⇒ restricted view on cross-linguality
  - Translations of NE
  - Translated strings are neutral wrt. corpus

# Hybrid NL-Query Translation for Cross-lingual ODQA

## Improvements

- **Language Model**
  - translations from the on-line MT systems are ranked according to a language model
  - 3-gram model of document corpus  
⇒ corpus-sensible ranking of translations
- **Alignment of Query-Information**
  - based on several filters (dictionary, PoS & string similarity)  
⇒ “transformation” of DE-QueryObject (Q-Focus) onto to EN-translation  
⇒ no need of parsing on English side
- **NE-specific alignment**
  - Not person names
  - but organizations, locations



# We wanted to use the Clef Forum for External Evaluation of Our ODQA approach: So, what is Clef?

---

- The **Cross-Language Evaluation Forum (CLEF)** supports global digital library applications by
  - developing an infrastructure for the **testing, tuning and evaluation of information retrieval systems** operating on European languages and
  - creating **test-suites of reusable data** which can be employed by system developers for benchmarking purposes.
- **QA@Clef 2003:**
  - Initial pilot evaluation for Cross-Lingual Open-Domain Question Answering (6 groups, among 3 from overseas)

# QA@Clef 2004 Track Setup – Task Definition

Given **200 questions** in a source language, find **one exact answer** per question in a collection of documents written in a target language, and provide a justification for each retrieved answer (i.e. the **docid** of the unique document that supports the answer).

S \ T	DE	EN	ES	FR	IT	NL	PT
BG							
DE							
EN							
ES							
FI							
FR							
IT							
NL							
PT							

6 monolingual and  
50 bilingual tasks.

**18 Teams**  
participated in 19  
tasks, submitting  
48 runs.

# Evaluation Exercise – Results (EN)

Results of the runs with English as target language.

Run Name	R	W	X	U	Overall Accuracy (%)	Accuracy over F (%)	Accuracy over D (%)	NIL Accuracy		CWS
								Precision	Recall	
bgas041bge n	26	168	5	1	13.00	11.67	25.00	0.13	0.40	0.056
dfki041deen	47	151	0	2	<b>23.50</b>	23.89	20.00	0.10	0.75	0.177
dltg041fren	38	155	7	0	19.00	17.78	30.00	0.17	0.55	-
dltg042fren	29	164	7	0	14.50	12.78	30.00	0.14	0.45	-
edin041deen	28	166	5	1	14.00	13.33	20.00	0.14	0.35	0.049
edin041fren	33	161	6	0	16.50	17.78	5.00	0.15	0.55	0.056
edin042deen	34	159	7	0	17.00	16.11	25.00	0.14	0.35	0.052
edin042fren	40	153	7	0	<b>20.00</b>	20.56	15.00	0.15	0.55	0.058
hels041fien	21	171	1	0	10.88	11.56	5.00	0.10	0.85	0.046
irst041iten	45	146	6	3	<b>22.50</b>	22.22	25.00	0.24	0.30	0.121
irst042iten	35	158	5	2	17.50	16.67	25.00	0.24	0.30	0.075
lire041fren	22	172	6	0	11.00	10.00	20.00	0.05	0.05	0.032
lire042fren	39	155	6	0	19.50	20.00	15.00	0.00	0.00	0.075



# QA@Clef 2005

---

- Temporal questions
- 23 groups
- DFKI
  - DE2DE, 2 runs
  - DE2EN
  - EN2DE, 2 runs

## QA@CLEF-2005

The 67 runs that have been submitted are divided into:

- 2 in the tasks with *BG* as target (both monolingual)
- 6 in the tasks with *DE* as target (3 monolingual and 3 cross-language)
- 12 in the tasks with *EN* as target (all cross-language)
- 18 in the tasks with *ES* as target (13 monolingual and 5 cross-language)
- 2 in the tasks with *FI* as target (both monolingual)
- 13 in the tasks with *FR* as target (10 monolingual and 3 cross-language)
- 6 in the tasks with *IT* as target (all monolingual)
- 3 in the tasks with *NL* as target (all monolingual)
- 5 in the tasks with *PT* as target (4 monolingual and 1 cross-language)

27

# Temporal Questions

---

## ■ Examples

- Who was world champion in soccer 1966?
- How often did German scientists win a Nobel Prize between 1950 and 1980?
- Who was US president during German's re-unification phase?

## ■ Some challenges

- Implicit temporal expressions
- Answer candidates contain explicit mentioning of dates, but query refers to interval

## ■ Our strategy

- Construct date-related underspecified IR-query
- Check consistency between date instances from answer candidates and query

# Web-based Query Analysis

- In principle the same approach as for Lucene-IR-query; Major differences
  - Google-specific IR-query construction
  - IR-query construction: Return N-best snippets
  - Snippet post-processing
    - » A-type + snippet-structure,
      - e.g, use <b> ... </b> for distance measure
    - » Construction of possible well-formed string
      - “Snippet grammar”
      - application of NL parsing tool
    - » Latent semantic kernels for answer extraction
- WebQA-based answer validation
  - Q-type/Q-focus + Query string (Q)  
+ Answer candidate (AC) =  
**direct\_answer\_string**
  - Google’s **Total Estimated Counts** (TEC) for re-ranking

©, PD Dr. Ginter Neumann, Leipzig, DFKI  
Ongoing: trustworthiness for answer candidates (utility theory; authority

# Web-based Query Analysis

- In principle the same approach as for Lucene-IR-query; Major differences
  - Google-specific IR-query construction
  - IR-query construction: Return N-best snippets
  - Snippet post-processing
    - » A-type + snippet-structure,
      - e.g, use <b> ... </b> for distance measure
    - » Construction of possible well-formed string
      - “Snippet grammar”
      - application of NL parsing tool
    - » Latent semantic kernels for answer extraction
- WebQA-based answer validation
  - Q-type/Q-focus + Query string (Q)  
+ Answer candidate (AC) =  
`direct_answer_string`
  - Google’s **Total Estimated Counts** (TEC) for re-ranking

**Q: What is the capital of Germany?**  
**AC: Berlin, Ney York**

©, PD Dr. Ginter Neumann, Leipzig, DFKI  
Ongoing: trustworthiness for answer candidates (utility theory; authority

# Web-based Query Analysis

- In principle the same approach as for Lucene-IR-query; Major differences
  - Google-specific IR-query construction
  - IR-query construction: Return N-best snippets
  - Snippet post-processing
    - » A-type + snippet-structure,
      - e.g, use <b> ... </b> for distance measure
    - » Construction of possible well-formed string
      - “Snippet grammar”
      - application of NL parsing tool
    - » Latent semantic kernels for answer extraction
- WebQA-based answer validation
  - Q-type/Q-focus + Query string (Q)  
+ Answer candidate (AC) =  
**direct\_answer\_string**
  - Google’s **Total Estimated Counts** (TEC) for re-ranking



**Q: What is the capital of Germany?**  
**AC: Berlin, Ney York**

©, PD Dr. Ginter Neumann, Leipzig, DFKI  
Ongoing: trustworthiness for answer candidates (utility theory; authority

# Web-based Query Analysis

- In principle the same approach as for Lucene-IR-query; Major differences
  - Google-specific IR-query construction
  - IR-query construction: Return N-best snippets
  - Snippet post-processing
    - » A-type + snippet-structure,
      - e.g, use <b> ... </b> for distance measure
    - » Construction of possible well-formed string
      - "Snippet grammar"
      - application of NL parsing tool
    - » Latent semantic kernels for answer extraction
- WebQA-based answer validation
  - Q-type/Q-focus + Query string (Q)  
+ Answer candidate (AC) =  
**direct\_answer\_string**
  - Google's **Total Estimated Counts** (TEC) for re-ranking

Q: What is the capital of Germany?  
AC: Berlin, Ney York

+ "Berlin"  
"is the canital of Germany" ~10  
TEC=331

©, PD Dr. Ginter Neumann, Leipzig, DFKI  
Ongoing: trustworthiness for answer candidates (utility theory; authority

# Web-based Query Analysis

- In principle the same approach as for Lucene-IR-query; Major differences
  - Google-specific IR-query construction
  - IR-query construction: Return N-best snippets
  - Snippet post-processing
    - » A-type + snippet-structure,
      - e.g, use <b> ... </b> for distance measure
    - » Construction of possible well-formed string
      - “Snippet grammar”
      - application of NL parsing tool
    - » Latent semantic kernels for answer extraction
- WebQA-based answer validation
  - Q-type/Q-focus + Query string (Q)  
+ Answer candidate (AC) =  
**direct\_answer\_string**
  - Google’s **Total Estimated Counts** (TEC) for re-ranking

Q: What is the capital of Germany?  
AC: Berlin, Ney York

+”Berlin”  
“is the canital of Germany”~10  
TEC=331

©, PD Dr. Ginter Neumann, Leipzig, DFKI  
Ongoing: trustworthiness for answer candidates (utility theory; authority

# Web-based Query Analysis

- In principle the same approach as for Lucene-IR-query; Major differences
  - Google-specific IR-query construction
  - IR-query construction: Return N-best snippets
  - Snippet post-processing
    - » A-type + snippet-structure,
      - e.g, use <b> ... </b> for distance measure
    - » Construction of possible well-formed string
      - "Snippet grammar"
      - application of NL parsing tool
    - » Latent semantic kernels for answer extraction

- WebQA-based answer validation

- Q-type/Q-focus + Query string (Q)  
+ Answer candidate (AC) =  
**direct\_answer\_string**
- Google's **Total Estimated Counts** (TEC) for re-ranking

Q: What is the capital of Germany?  
AC: Berlin, Ney York

+ "Berlin"  
"is the canital of Germany" ~ 10  
TEC = 331

+ "New York"  
"is the capital of Germany" ~ 10  
TEC = 75

29

©, PD Dr. Ginter Neumann, Leibniz-URK  
Ongoing: trustworthiness for answer candidates (utility theory; authority



# DFKI Related Future Work

---

- **SmartWeb**

- Mobile access to the Semantic Web
- ODQA for search in syntactic Web pages
- BMBF Verbund project (partners from research and industries, e.g., BMW, Siemens, T-Systems)

- **HyLab**

- QA for Personal Digital Memory
- BMBF funded DFKI project (start 2006)



# BMBF

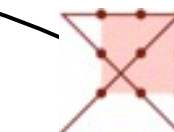
## Softwaresysteme

# SmartWeb



## AIFB

### DAIMLERCHRYSLER



European Media Lab

## IT-2006 und Futur-Programm

### BMBF: 13,7 Mio Euro (Dr. Reuse)

### Leiter: W. Wahlster

### Laufzeit: 2004-2007



Chair for  
Pattern Recognition  
FAU Erlangen-Nuremberg

## LMU IPSK

Ludwig-Maximilians-  
Universität München

## SIEMENS



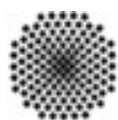
UNIVERSITÄT  
DES  
SAARLANDES

## ...T...Systems



SEMANTICS FOR THE WEB

IMS Institut für Maschinelle  
Sprachverarbeitung,  
Universität Stuttgart



Fraunhofer



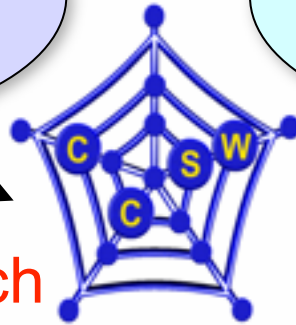
Institut  
Rechnerarchitektur  
und Softwaretechnik



# SmartWeb: Mobiler Breitbandzugang

Mobiler  
breitbandiger  
Internetzugang  
UMTS/WLAN

Mobiler  
Multimodaler Dialog  
Sprache, Gestik,  
Haptik, Mimik



Neu: Offener Themenbereich

Neu: Fragebeantwortung

## Semantisches Web

Automatisch  
annotierte  
Intranetseiten

Semantisch  
annotierte  
Webseiten

Klassische  
HTML-  
Webseiten

Sprachtechnologie,  
Informationsextraktion

---

DONE!

Thank you for your attention!

---

# HyLaP

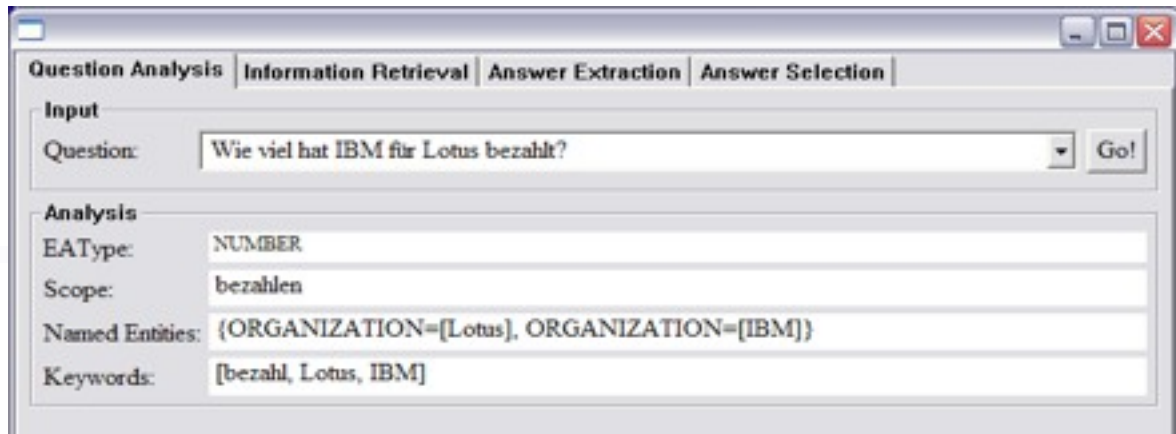
Hybrid Language Processing Technologies  
for a personal associative information  
access and management application

Günter Neumann & Hans Uszkoreit  
BMBF funded DFKI project, ~1.5 Euro  
Start: 2006

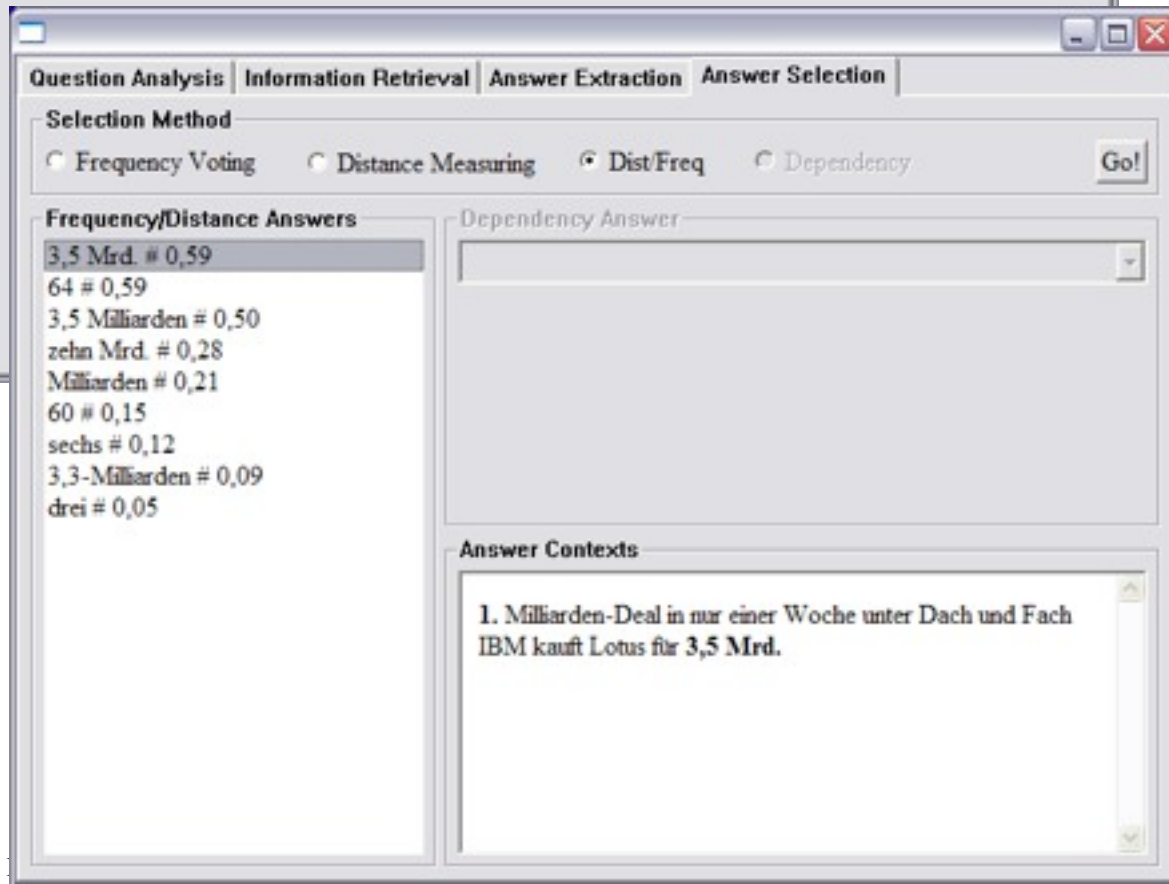
# Assumptions

---

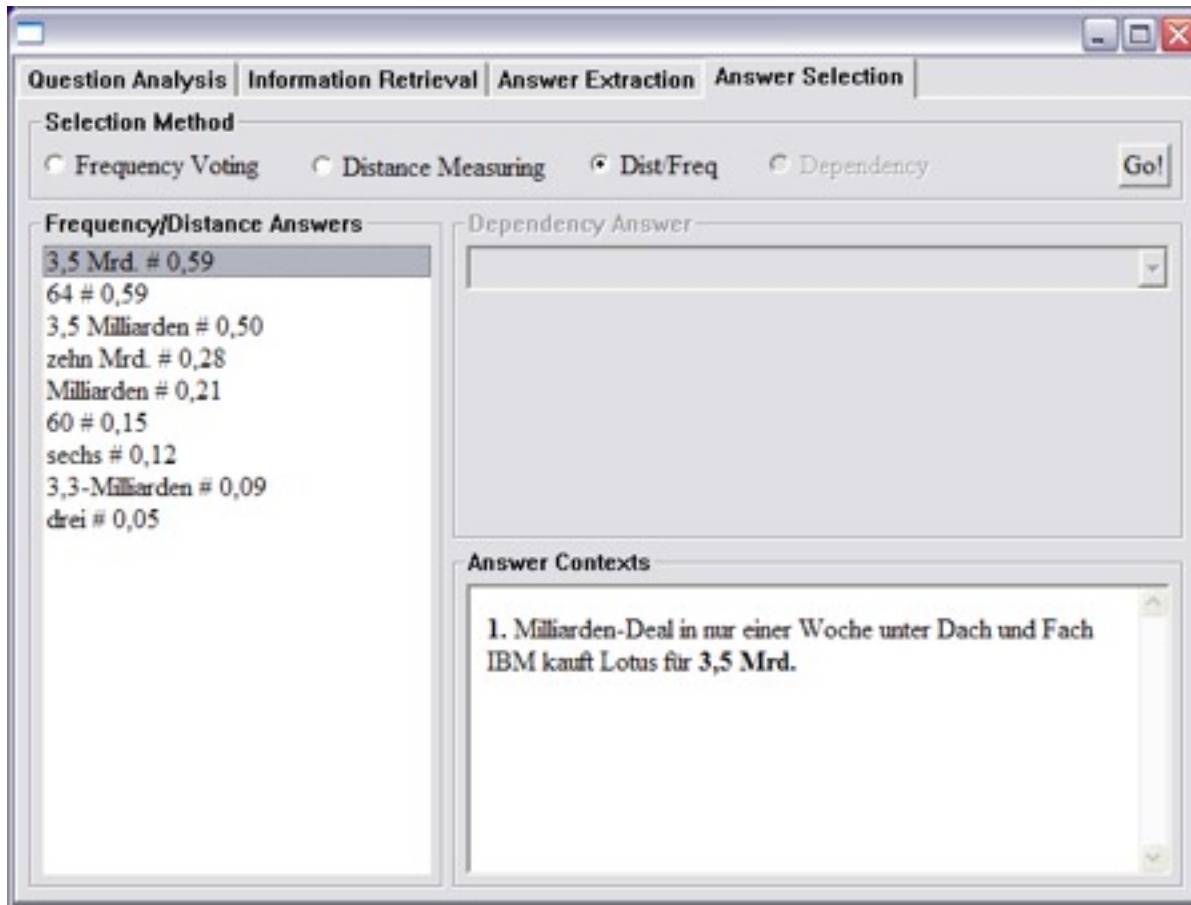
- Question answering is the most natural way of requesting information
- The management of personal digital content becomes a true challenge
- The exploitation of personal digital memory will change our lives
- Applications for authoring, browsing and commenting will converge
- Every document can become an interface to memory



Frageanalyse:  
 Zentral ist die Berechnung  
 Des Fragetyps (Faktenfrage,  
 Definitionsfrage) und des  
 Fragegegenstandes (Person,  
 Datum, Ort, ...).  
 Da hierbei linguistischbasierte  
 Verfahren eingesetzt müssen,  
 Besteht die Kunst darin, dies  
 Gleichermassen Präzise und  
 Robust machen zu können ohne  
 Dem Benutzer eine Einschränkung  
 Bei seinen Formulierungsbemühungen  
 Zu machen.



Ergebnisauswahl:  
 Wenn das System eine Menge von  
 Antwortkandidaten bestimmt hat,  
 müssen gemäß ihrer Güte sortiert  
 werden. Dadurch kann 'dem Benutzer  
 zB nur die beste Antwort geliefert  
 werden oder die N-besten, sodass der  
 Benutzer zB selber noch entscheiden  
 kann, welche er als die Beste Antwort  
 auffasst. Dies wird am besten dadurch  
 realisiert, dass zu jedem  
 Antwortkandidat auch der Satz, in  
 dem Antwort vorkommt angezeigt  
 wird. So kann der Benutzer stets auch  
 den Kontext und auch die Güte des  
 Systems selbst noch einmal  
 begutachten.

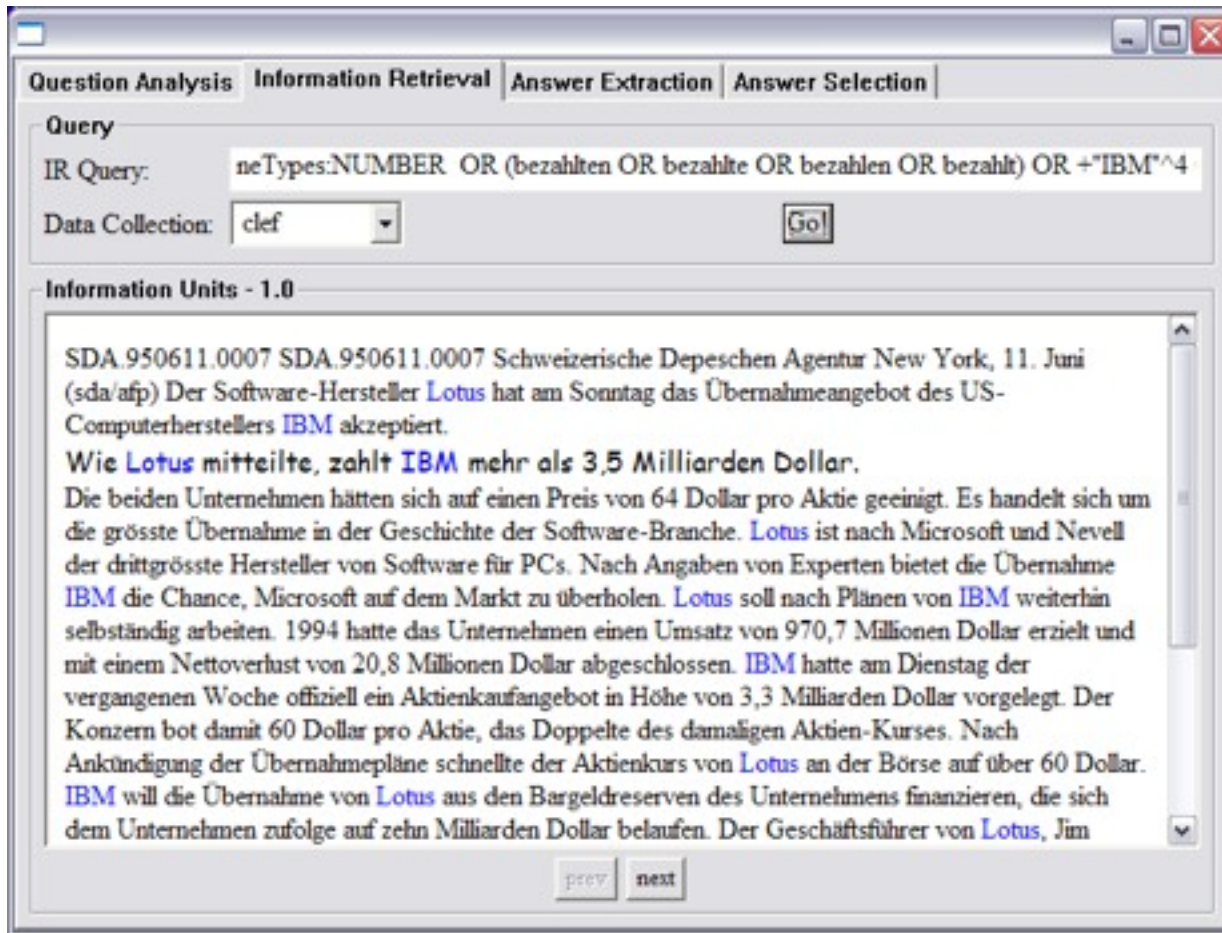


Ergebnisauswahl:  
 Wenn das System eine Menge von Antwortkandidaten bestimmt hat, müssen gemäß ihrer Güte sortiert werden. Dadurch kann dem Benutzer zB nur die beste Antwort geliefert werden oder die N-besten, sodass der Benutzer zB selber noch entscheiden kann, welche er als die Beste Antwort auffasst. Dies wird am besten dadurch realisiert, dass zu jedem Antwortkandidat auch der Satz, in dem Antwort vorkommt angezeigt wird. So kann der Benutzer stets auch den Kontext und auch die Güte des Systems selbst noch einmal begutachten.

Bei der Answer selection koennen verschiedene Verfahren eingesetzt werden, auch kombiniert:  
 -Frequenz  
 -Distanz zwischen Antwortkandidat und Fragewörtern  
 -Beides gemischt

Die Kunst hierbei ist es, solche Kriterien zu finden, die möglichst korrekte Antwortkandidaten von falschen Antworten unterscheiden. Jenachdem, ob das QA open-domain oder domain-spezifisch eingesetzt wird, können hjerzu auch Ontologien und spezielle Inferenzmechanismen eingesetzt werden. Desweiteren bietet sich hier die Möglichkeit, erfolgreiche Berechnungen zu speichern (episodic memory)





Dokumentenauswahl:  
Das Ergebnis der Frageverarbeitung muss so umgeformt werden, dass sie optimal zur Information retrieval eingesetzt werden kann. Hier benutzen wir sogar Verfahren zur linguistischen Generierung von Wörtern/Phrasen, um einen hohen recall zu erreichen. Die Kunst hier ist es, die Frageanalyse in eine optimale IR-Suchanfrage zu überführen, sodass bereits beim Retrieval nur sehr wenige, aber relevante Dokumente bestimmt werden.

Antwortextraktion:  
 Die durch die IR-Engine bestimmten Paragraphen, müssen für die Identifikation und Extraktion der Antwortkandidaten weiter analysiert werden, wobei neben IE-Techniken auch eine robuste NL-Verarbeitung eingesetzt wird. Zum Beispiel ist es hilfreich, die Beziehungen zwischen den einzelnen Wörtern genau zu bestimmen, damit zB genauer Abstandsmasse eingesetzt werden können. Auch ist es gerade für zusammengesetzte Wörter zu identifizieren.

Die Kunst ist zum einen der Einsatz sehr robuster und sehr schneller Sprachtechnologie, aber vor allem der Einsatz automatischer Methoden zur Skalierung relevanter Extraktionsparameter.

**Question Analysis** | Information Retrieval | Answer Extraction | Answer Selection

**Input**

Question:

**Analysis**

EAType:	PERSON
Scope:	seien
Named Entities:	{ORGANIZATION=[Dortmund]}
Keywords:	[Dortmund, sei, trainer]

