# Multilingual System for Web-Information: The State of The Art

## Feiyu Xu
## DFKI, LT-Lab
## Germany

# Multilingual Information System

❑ Motivation

❑ Strategies

❑ MIETTA System

# Motivation

- ❏ More and more web information are encoded in other languages than English, for example, Chinese 7%

  - ✦ English is loosing its dominance

- ❏ Improve the inter-culture and business communication

  - ✦ Quick access to information is important for every day life

- ❏ Text REtrieval Conference (TREC) (http://trec.nist.gov/)

  - ✦ English, Spanish, Chinese, etc.

- ❏ Information Retrieval for Asian Language Conference (IRAL)
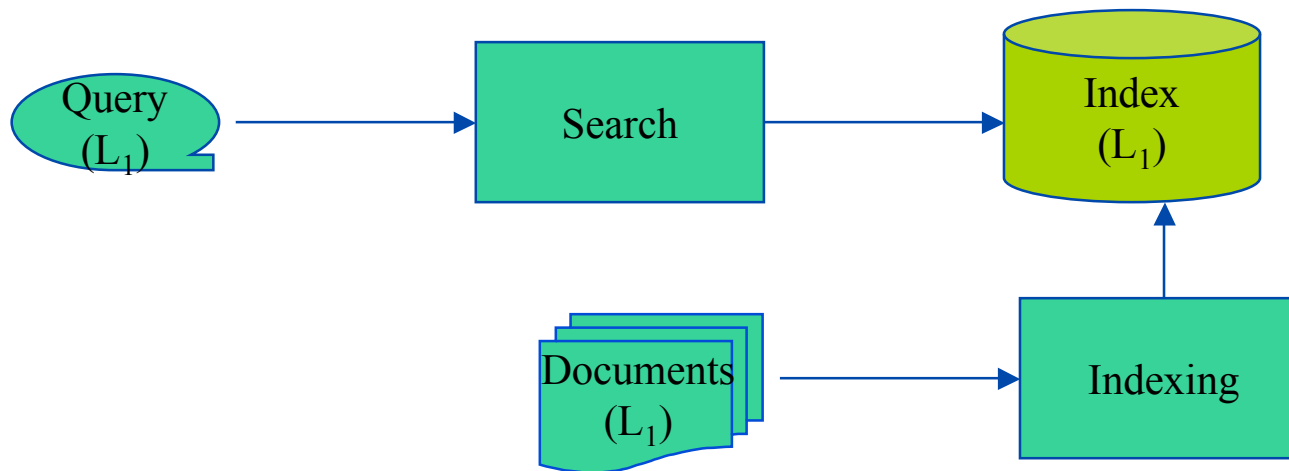
- ❏ European ESPRIT consortium (French, Belgian, German)

# What is Information Retrieval (http://www.lt-world.org)

- **Synonyms**: document retrieval

- **Definition**: Information Retrieval is the process of locating information that fits a user's requirements, where the requirements are usually expressed as a search query. The fit of the retrieved information with the information need is referred to as "relevance" …

- http://www.lt-world.org/HLT_Survey/ltw-chapter7-2.pdf

# What is Monolingual Information Retrieval?

❑ Query and information to be looked for are encoded in a same language

# What is Multilingual Information Retrieval?

❑ An extension of the general information retrieval problem

❑ Finding information, e.g., web documents which are not encoded in the same language as the query is encoded in

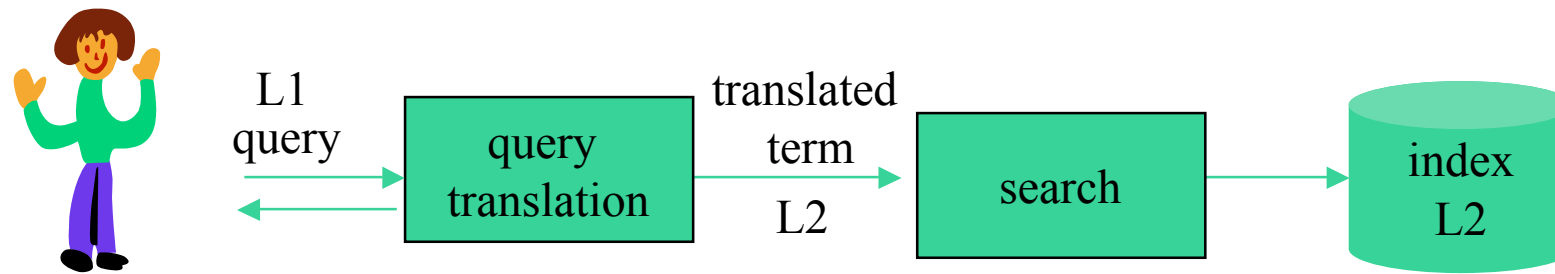❑ Similar terms: "crosslingual information retrieval" and "translingual information retrieval"

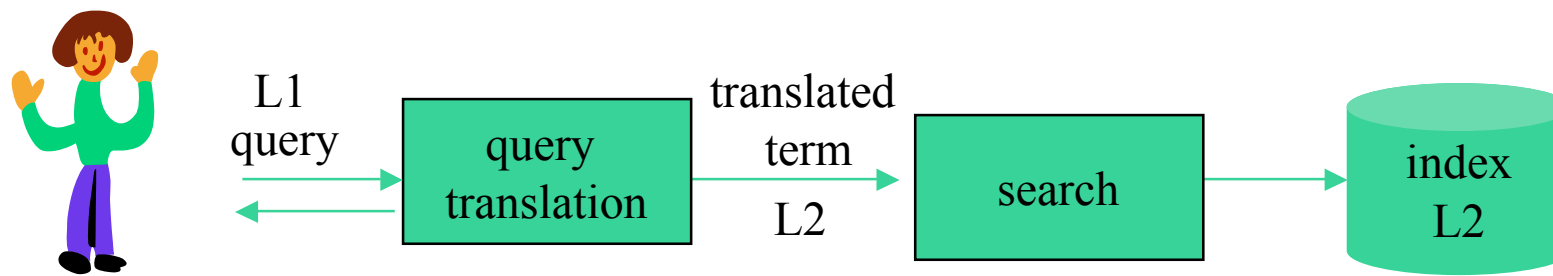# Different Multilingual Information Retrieval Strategies Supported by Language Technologies

❑ Online query translation

  ✸ Help user to formulate his query in a foreign language

❑ Online document translation

  ✸ Translate the found document into the query language

❑ Offline document translation

  ✸ Make web documents multilingual available

❑ Combination of information extraction and multilingual generation

  ✸ Make database information multilingual available and allow the free text retrieval of database information

Source: Feiyu Xu, 2002

# Query Translation

❑ Help user to formulate their query in another language



L1
query → query translation → translated term L2 → search → index L2

❑ The primary problem is that short queries provide less context for word sense disambiguation, and inaccurate translations lead to bad recall and precision

❑ How can the user access the content of the found document?

L1
query

translated
term

query
translation

L2

search

index
L2

# MULINEX System

Source: Feiyu Xu, 2002

# MULINEX System

Source: Feiyu Xu, 2002

# MULINEX System

# Online document translation

❑ Translating the found the documents into query language, for example, google

Query L1

Machine Translation (from $L_i$ to $L_1$)

Search

Index $(L_1,\ldots,L_n)$

Documents $(L_1,\ldots,L_n)$

Indexing

Query L1

Search

Machine Translation (from $L_i$ to $L_1$)

Index $(L_1,\ldots,L_n)$

Documents $(L_1,\ldots,L_n)$

Indexing

DFKIt

# Online document translation (Google)

**Source: Feiyu Xu, 2002**

# Offline Document Translation

❑ Automatic offline translation

- ✷ Source text is translated into target languages
- ✷ Index is constructed from translation
- ✷ Search term in one language yields original and translated documents

query

L1

search

index
L1

indexing

original
documents
L2

document
translation

translated
documents
L1

❑ A higher translation and retrieval performance, since the full original document provides more context for disambiguation. The word sense disambiguration problem is less severe than query translation

❑ The main limitation is the duplication of the indices, and the translated documents also need to be stored

❑ The offline translation is practically not viable due to big cost of computation and storage for the general search engines like Alta-Vista, Yahoo, etc.

# Facts Sheet - MIETTA

❑ Title: MIETTA -Multilingual Information Extraction for Tourism
and Travel Assistance

❑ Funding: EU Language Engineering Sector of TAP (HLT-IST)

❑ Technical Partners: DFKI, Celi, University of Helsinki, Polito,
Unidata

❑ User Partners: Commune DI Rome, City of Turku,
Staatskanzlei of the Saarland

# Objectives

❑ Multilingual internet portal and specialised information system for tourist information

Five languages: English, Finnish, French, German, Italian

Three regions: Rome, Saarland and Turku

❑ Integrated access to heterogeneous data sources and make it fully transparent to end users whether they are searching in

  ★ WWW documents or

  ★ Databases

# Offline Document Translation in MIETTA

❑ Use document translation as the main strategy. The reason is that it allows direct access to the content, it provides better performance within a restricted domain

❑ Use LOGOS for document translation, which covers the following directions:

✦ German⇒ English, French, Italian

← English⇒ French, German, Italian, Spanish

❑ The final document collection in MIETTA after the document translation yielded an almost fully covered multilingual setup.

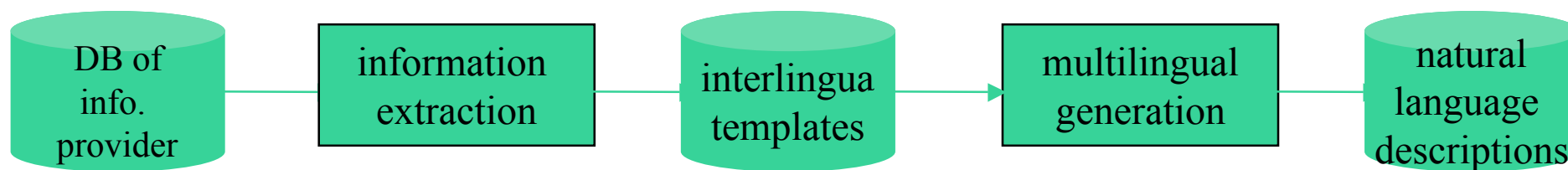# Information Extraction and Multilingual Generation

❑ Motivation

✦ Make the database content more structured and multilingual accessible.

✦ Apply the same free text retrieval method to the generated descriptions as to the web documents

| DB of info. provider | → | information extraction | → | interlingua templates | → | multilingual generation | → | natural language descriptions |

DB of info. provider → information extraction → interlingua templates → multilingual generation → natural language descriptions

# Information Extraction

❑ The objective of information extraction is twofold:

  ✦ To extract the domain relevant information (templates) from the unstructured data so that the user can access more facts and more accurately

  ✦ To normalise the extracted data in a language independent format to facilitate the multilingual generation

❑ Three steps for template extraction in MIETTA

  ✦ Natural language shallow processing: named entities, np, vp

  ✦ Normalisation: converting information into a language independent format

  ✦ Template filling: mapping the extracted information into template slots by employing specific template filler rules

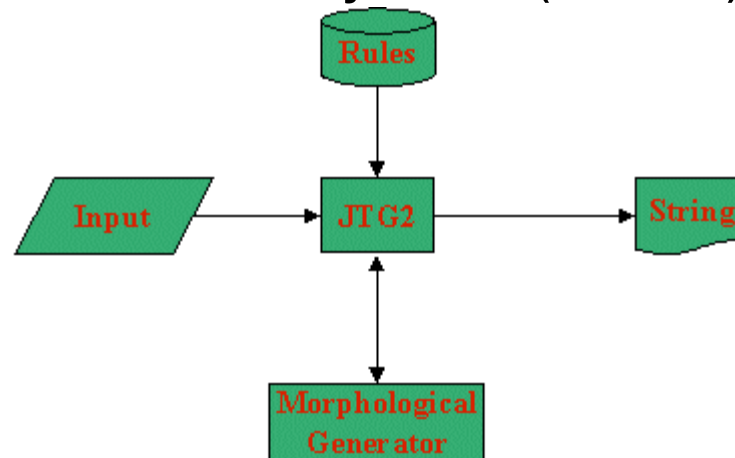# Example of IE

## German text from an event calendar in Saarland

St. Ingbert: -Sanfte Gymnastik für Seniorinnen und Senioren, montags von 10 bis 11 Uhr im Clubraum, Kirchengasse 11.

*English: St. Ingbert: -Gentle Gymnastic for seniors, every Monday from 10:00 to 11:00 am, in Club room, Kirchengasse 11*

Event:

```
Name:        gymnastic
Addressee:   seniors
time:
             start time:10
             end time: 11
             weekly: yes
             weekday: 1

location:    city name: St. Ingbert
             address: Club room
                      Kirchengasse 11
```

DFKI

# Multilingual Generation

❑ Template Generation system (JTG/2)



❑ Language independent input allows for easy extension of the generation component to other languages

**Source: Feiyu Xu, 2002**

# Example

Level1: Event
Level2: Theater
Level3:
Event-Name: Faust
StartDate: 21.10.99
PlaceName: Staatstheater
Address: Schillerplatz, 66111 Saarbrücken
Phone: 0681-32204

English:

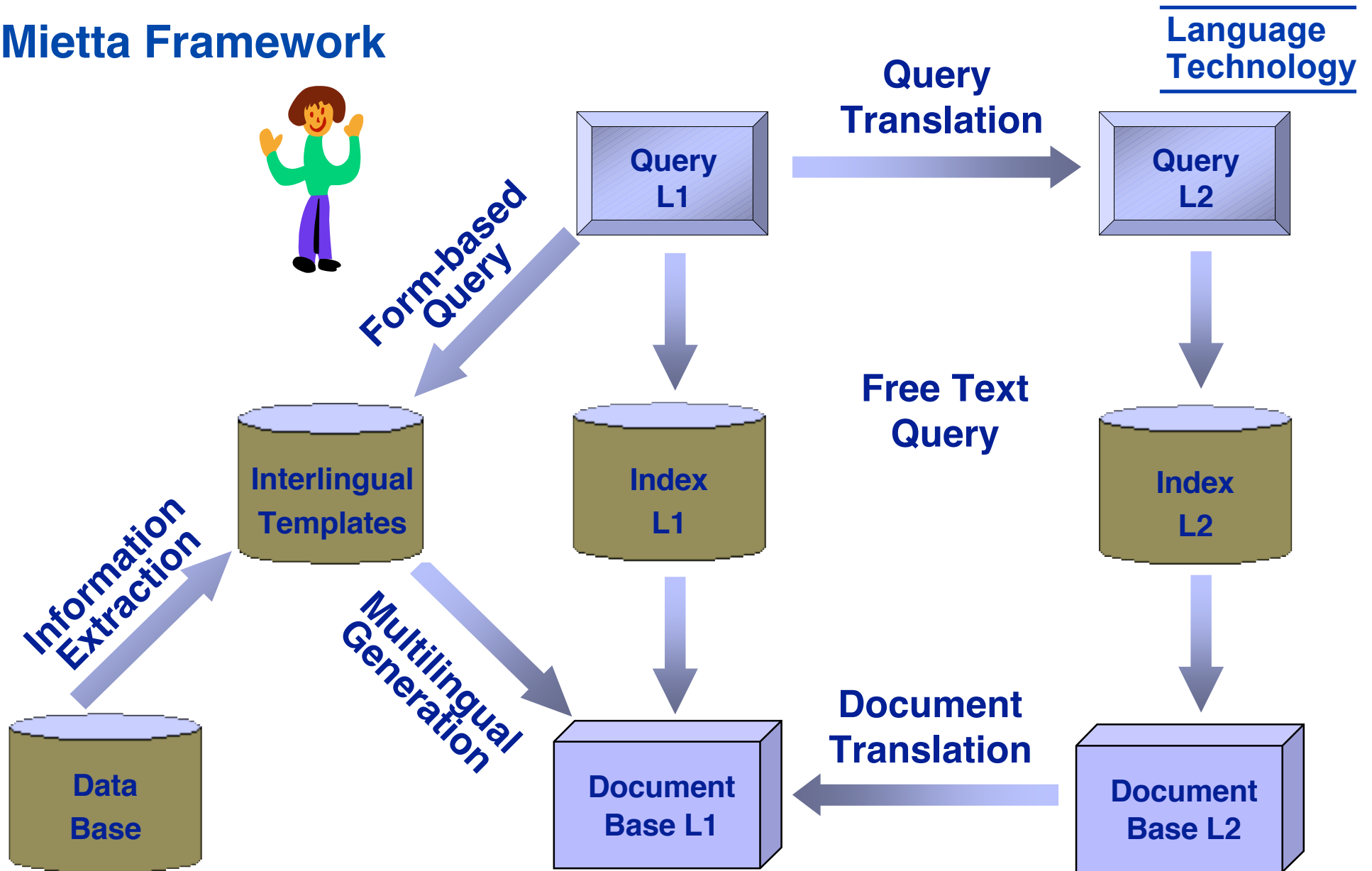The theater show Faust will take place at the Staatstheater in Schillerplatz 1, 66111 Saarbrücken (in the downtown area).

The scheduled date is Thursday, October 21, 1999. Phone: 06 81-32204

Finnish:

Teatteriesitys Faust järjestetään Staatstheaterissa, osoitteessa Schillerplatz 1, 66111 Saarbrücken (keskustan alueella). Tapahtuman päivämäärä on 21. lokakuuta 1999. Puhelin: 06 81-32204.
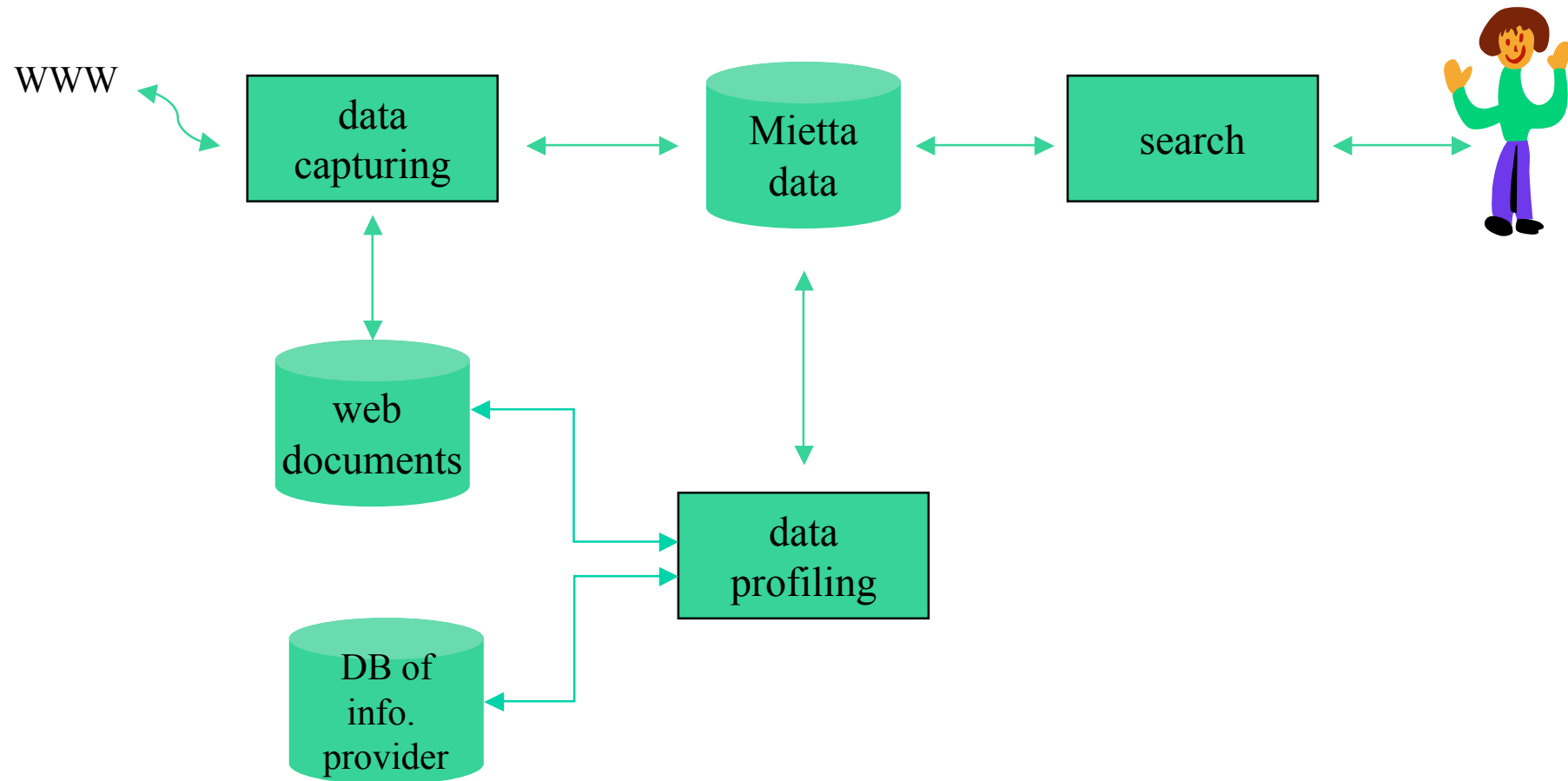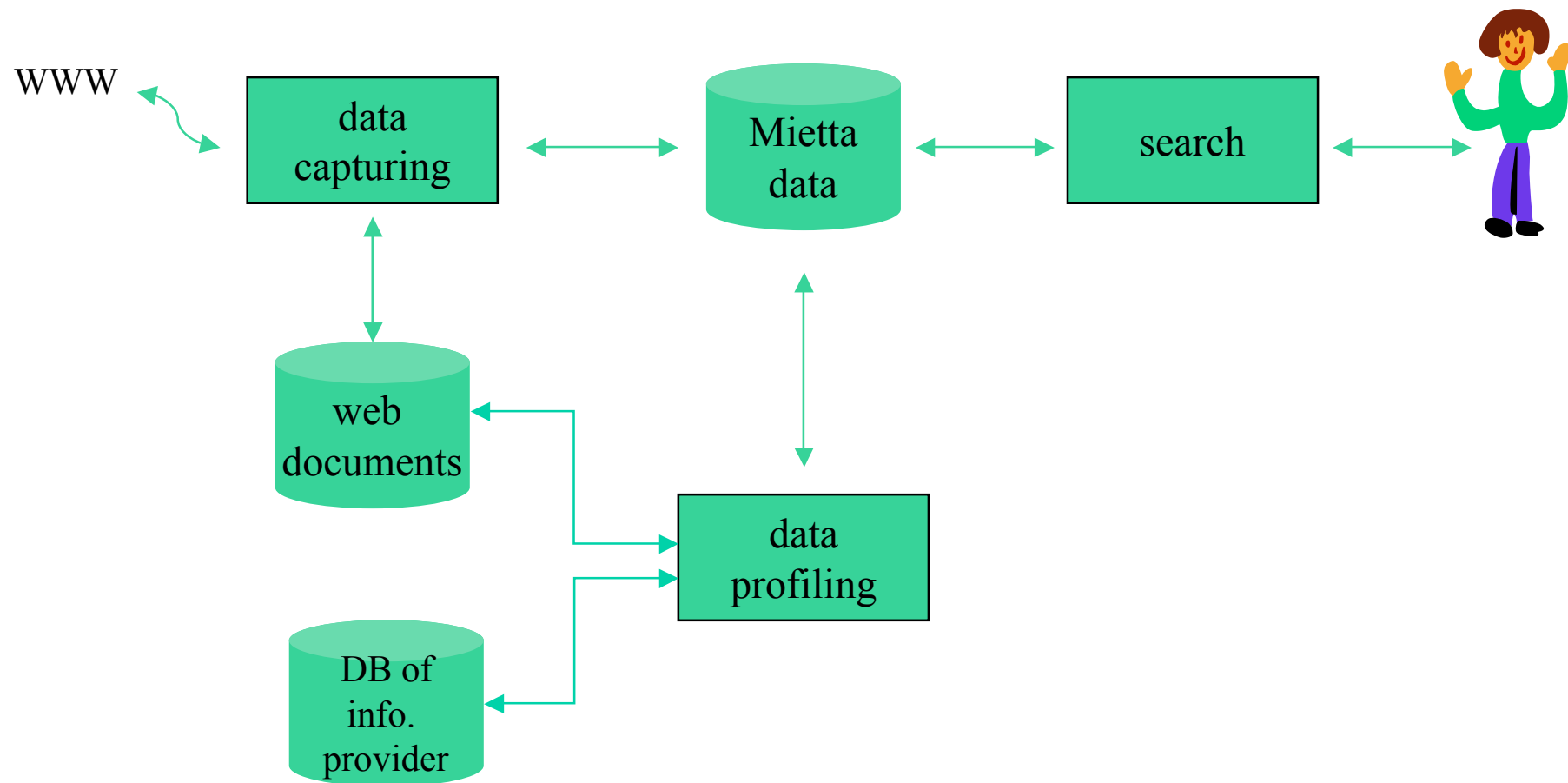
# Mietta Framework

Language Technology

Query Translation

Query L1 → Query L2

Form-based Query

Free Text Query

Information Extraction

Interlingual Templates

Index L1

Index L2

Multilingual Generation

Document Translation

Data Base

Document Base L1 ← Document Base L2

DFKIt

Language Technology

Query Translation

Query L1    Query L2

Form-based Query

Free Text Query

Information Extraction

Interlingual Templates

Index L1

Index L2

Multilingual Generation

Data Base

Document Base L1

Document Translation

Document Base L2

# The Overall MIETTA System

WWW

data capturing ⟷ Mietta data ⟷ search

web documents

DB of info. provider

data profiling

Source: Feiyu Xu, 2002

DFKIt

WWW

data capturing

Mietta data

search

web documents

data profiling

DB of info. provider

DFKIt

# Data Profiling

- ❑ Document translation, based on LOGOS machine translation system

- ❑  Information Extraction from database entries for template construction

- ❑  Multilingual generation from templates to obtain natural language descriptions

- ❑  Free text indexing

**DFKI***lt*

# TNO (ISM, VSM) Indexing Toolkit

❑ ISM: A lemma-based fuzzy index based on trigrams

❑ VSM: A vector space model index based on lemmatas



**Source: Feiyu Xu, 2002**

Mietta data

free text indexing

multilingual generated texts

web documents

translations

# Scalability of the Framework

❑ Adaptation to other domains

  ✶ Domain specific templates

  ✶ Domain Concept hierarchy

  ✶ Domain specific template filler rules

  ✶ Domain specific generation grammars

❑ Extension to other languages

  ✶ Natural language generation tool requires less effort for the development of a grammar rule set in a language

  ✶ Information extraction requires available language specific resources

  ✶ Document translation is dependent on the machine translation system

# Evaluation of the MIETTA System

❑ The standard relevance assessment model used in ad hoc and routing forums of TREC is difficult to apply to the complete MIETTA system because of

  ✴ Broad variety of search strategies

  ✴ Heterogenous data sources

❑ MIETTA is evaluated as technically "excellent" by EU

❑ Two projects are derived from MIETTA

  ✴ Natural science foundation of China Project in SJTU

  ✴ EU project of MIETTA to transfer the idea into product in XtraMind in Saarbrücken

# Conclusion: Innovative Technical Features

- ❑ Integration of different multilingual and crosslingual search technologies
- ❑ Combination of IE and multilingual generation
- ❑ Integration of DB and text document access
- ❑ Intelligent User Interface
- ❑ XML for advanced information management
- ❑ Localisation technologies for user interface and multilingual generation
- ❑ Highly suitable as a domain-specific information system and internet portal

Source: Feiyu Xu, 2002

# MIETTA Start Page: Choose Region

# Choose Language

# MIETTA Search Menu

Source: Feiyu Xu, 2002

# MIETTA Free Text Retrieval

Source: Feiyu Xu, 2002

# MIETTA Class-based Navigation

**Source: Feiyu Xu, 2002**

# MIETTA Class-based Navigation with Free Text

# MIETTA Form based Query

**Source: Feiyu Xu, 2002**

# Online Text Generation

| | |
|---|---|
| English | The theater **Staatstheater** is located in Schillerplatz 1, 66111 Saarbrücken (in the downtown area). Phone: 06 81-32204 . |
| Finnish | Teatteri **Staatstheater** sijaitsee osoitteessa Schillerplatz 1, 66111 Saarbrücken (keskustan alueella). Puhelin: 06 81-32204. |
| French | Le théâtre **Staatstheater** se trouve Schillerplatz 1, 66111 Saarbrücken (dans la zone du centre). Téléphone: 06 81-32204 . |
| German | Das Theater **Staatstheater** befindet sich in der Schillerplatz 1, 66111 Saarbrücken (im Stadtzentrum). Phone: 06 81-32204 . |
| Italian | Il teatro **Staatstheater** si trova in Schillerplatz 1, 66111 Saarbrücken (nella zona del centro). Telefono: 06 81-32204. |

# Result Presentation

❑ Result contains both database entries and documents

❑ All information is presented in uniform format

  ✶ Classified

  ✶ Ordered according to the relevance



Your search returned 3 webdocuments and 1 database entries.

Some of them belong to the following more specific classes. You can click on one of the classnames to get to these specific results.

| | | |
|---|---|---|
| Major Tourist Attractions | | 1 |
| Transport | 2 | |
| Events | 1 | |

| webdocuments that best match your search | | |
|---|---|---|
| LINE B | Transport | 63,64% |
| FASCIA BLU (BLUE SECTOR) | Transport | 30,71% |
| EASTER | Public and Town Festivals | 5,65% |

| database entries that best match your search | | |
|---|---|---|
| Colosseo (in the downtown area). Opening hours: 0900-1500. Phone: 067004261; Fax: 061234567. | Monuments | 100% |