# Crosslingual
# Information Retrieval
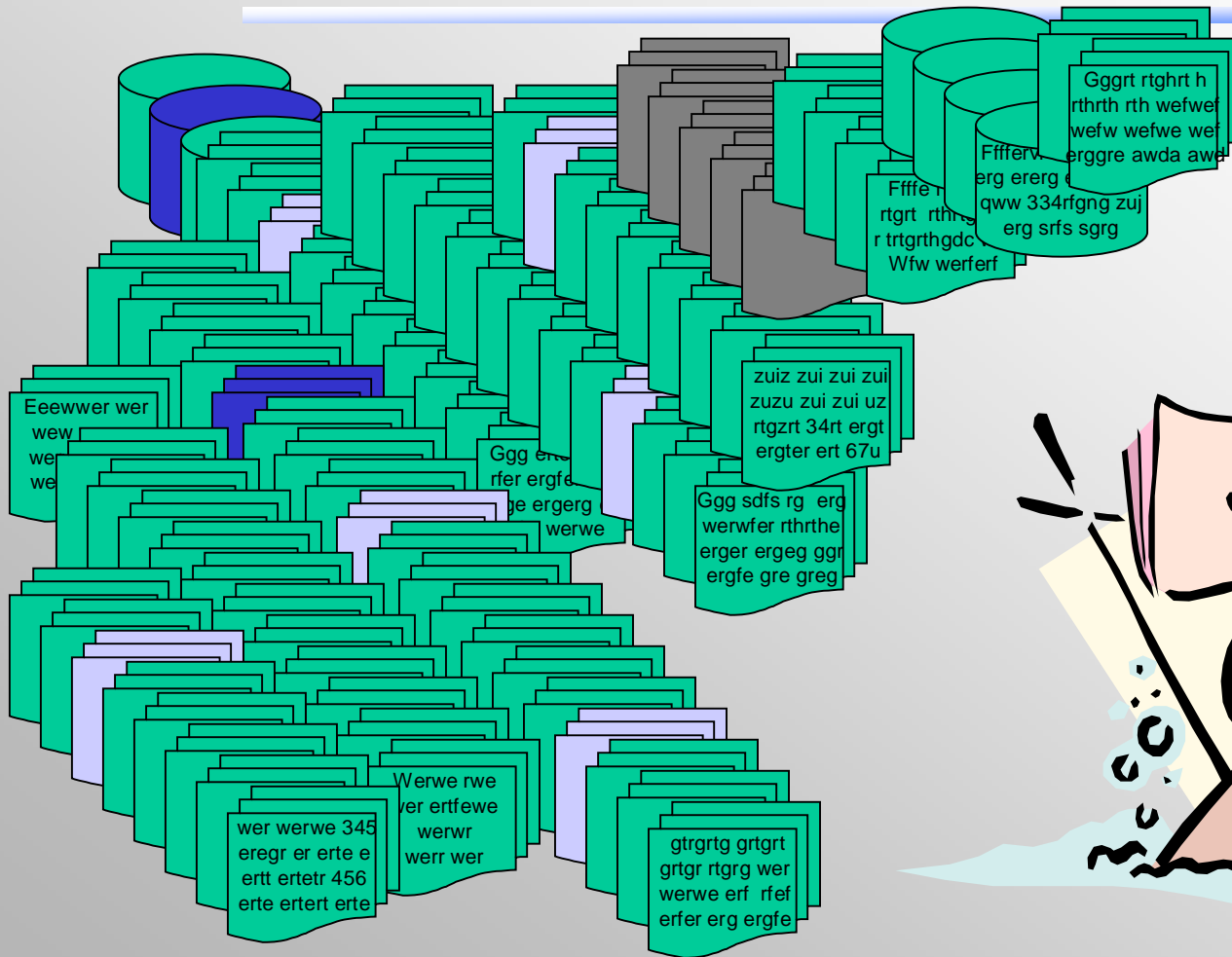# (CLIR)

Gr. Thurmair
Pisa 2004

# Outline

- Current retrieval
  - Likely improvements
- Why this does not easily work

- Design of a CLIR
  - Text translation, index translation, query translation

- Architecture of QTE based system
  - Analysis
  - Index
  - Query translation
  - Search & ranking
  - Document retranslation

I From IR to CLIR

# Looking for trouble?

# Architecture of an Indexing System

- Example: Lucene Open Source IR System
  - www.apache.org/lucene
  - Current Version = V1.4

- Software Modules:
  - Document
    - Document maintenance (set of fields)
  - Analysis
    - Tokeniser, Stop list, Index term Generator
  - Index
    - Writer: create / add documents
    - Reader: access index data
  - Search
    - Query Parser
    - Queries (Phrases / Boolean / …
    - Searcher (Hits, Highlighting …)

# What is the problem in search?

- Automatic Indexing: **From strings to concepts!**
  - Machines process *strings* – humans look for concepts

- Take into account variablitiy and flexibility of human language
  - Reduce variance / normalise to the relevant concepts of a text

- Finding relevant terms is a linguistic job
  - Linguistics = science of language

- Indexing is a linguistic task, among others with the support of statistic components

# What is an index term?

XML Topic Maps: XTL V1.0 (2001)

This specification provides a model and grammar for representing the structure of information resources used to define topics, and the associations (relationships) between topics. Names, resources, and relationships are said to be *characteristics* of abstract subjects, which are called *topics*. Topics have their characteristics within *scopes*: i.e. the limited contexts ... One or more interrelated documents employing this grammar is called a "topic map."

...

A **topic** is a resource within the computer that stands in for (or "reifies") some real-world *subject*. Examples of such subjects might be the play *Hamlet*, the playwrighter **William Shakespeare**, or the "authorship" relationship.

=> is *Shakespeare* a relevant term?

(38 occurrences, one of top 50 Terms in the corpus)

# Current Index Term creation techniques

- Tokenisation

  – All index terms are single words!

- Stemming

  – Intends to neutralise morphology variants

- Context

  – Intends to capture syntactic structure

- Clustering

  – Intends to neutralise semantic variants

| | |
|---|---|
| *author* | |
| *Author* | *ing* |
| *author* | *ities* |
| *author* | *s* |
| *Author* | *s* |
| *author* | *ship* |

*a topic maps several elements into a commmon cluster ...*

*topic element subject map XTM name association scope ID document resource occurrence class role indicator reference child variant relationship specification ... Shakespeare ...*

# Morphology: Stemmer

- **Stemming**: Reduction of Base forms to stems

  (often: Reduction of <u>Text</u>forms to stems)

- Technique: Cutting the **Affixes**
  - Problem: Cut prefixes as well?
  - Problem: Allomorphs

- Types of mistakes
  - Improper Linking of terms (Precision ↓)
  - Improper splitting of terms (Recall ↓)

- Effect
  - increases Recall
  - Increases noise if unprecise
  - Sometims necessary in CLIR (en verb - fr noun)

| | |
|---|---|
| Verarbeitung | -> arbeit |
| Funktionalität | -> funkt |
| Integrierbarkeit | -> integr |
| überreagieren | -> reag |
| Überreaktion | -> reak? |
| Verdichtung | -> dicht |
| Versorgung | -> sorg |
| Regionalisierung | -> reg |
| Regierung | -> reg |
| Regen, regnerisch | -> reg |

- Indexed ‚terms' in the corpus are:

| | |
|---|---|
| map | topic |
| Map | Topic |
| mapping | topic(s |
| Mapping | topicMap |
| maps | topicMaps.Org |
| map's | topicref |
| Maps | topics |
| | topic's |
| | Topics |

- How must the query look like?

Voorhees / Harmann 2000 (TREC8-Overview):

Standard Retrieval (ad-hoc task):

- No improvement in recent years
- Example: Precision SMART-System:
  between 0.15 and 0.4

„the results of the ad-hoc task have plateaued in recent years"

„the absolute performance of the task is less than ideal"

# From IR to CLIR

Multilingual Extension:
- Translate the query
- Translate the result

This is not enough!

documents

Indexing:

Tokenisation
Index Term detection

Query

**translate**

**translate**

Result

Document
Base

Index

# From IR to CLIR

- Current system structure is insufficient for good IR
  - Add more intelligence to the analysis
- Current structure also insufficient for CLIR
  - On index term creation
    - Index terms must be translated!
      - Single words, stemmed words … don't work
  - On index maintenance
    - Index must be language-sensitive …
  - On search
    - Search must be translated
    - Documents must be re-translated

# From IR to CLIR: more **Analysis Intelligence!**

- Failure of standard procedures
  - **Truncation / Stemming**
    - Dictionaries do not contain translations for stems or text forms
    - *deal$* or *author$* are not translatable
    => Lemmatiser, Decomposer required

  - **Context** and context operators
    - Most terms to be translated are **Multiword terms**.
      Indizes contain only single words.
    - Translation of  *money* NEAR *laundering*?
      > De: no context operation as just 1 word (*Geldwäsche*)
      > Es: wrong sequence (*lavado de dinero*)
    => Dependency analysis requireed

  - Ranking / **Clustering**
    - Always relates to monolingual term universe
    - parallel Korpora if available do not really require CLIR

II CLIR Design

## monolingual, multilingual, cross-lingual

| | | |
|---|---|---|
| Query (en) | →monolingual→ | Documents (en) |

| | | |
|---|---|---|
| Query (en) | → | Documents (en) |
| | multilingual | |
| Query (en) | → | Documents (de) |

| | | |
|---|---|---|
| Query (en) | ⤫ | Documents (en) |
| Query (en) | crosslingual | Documents (de) |

Query language                          Document language

# Crosslingual Information Retrieval

Drogenhandel
in Italien?

Drug import
to Germany?

Dok 1:
Winston Harton has been identified to be the drug
dealer who sold 350 gr of cocaine at Milan train station May
1996.

Dok 2:
Europol Edu has confirmed that the number of drugs
smuggled in via scheduled flights from Columbia is
increasing every year.

Dok 3:
An der ostdeutschen Grenze ist es wiederholt zu Schießereien
zwischen rivalisierenden Banden im Drogenkrieg gekommen.
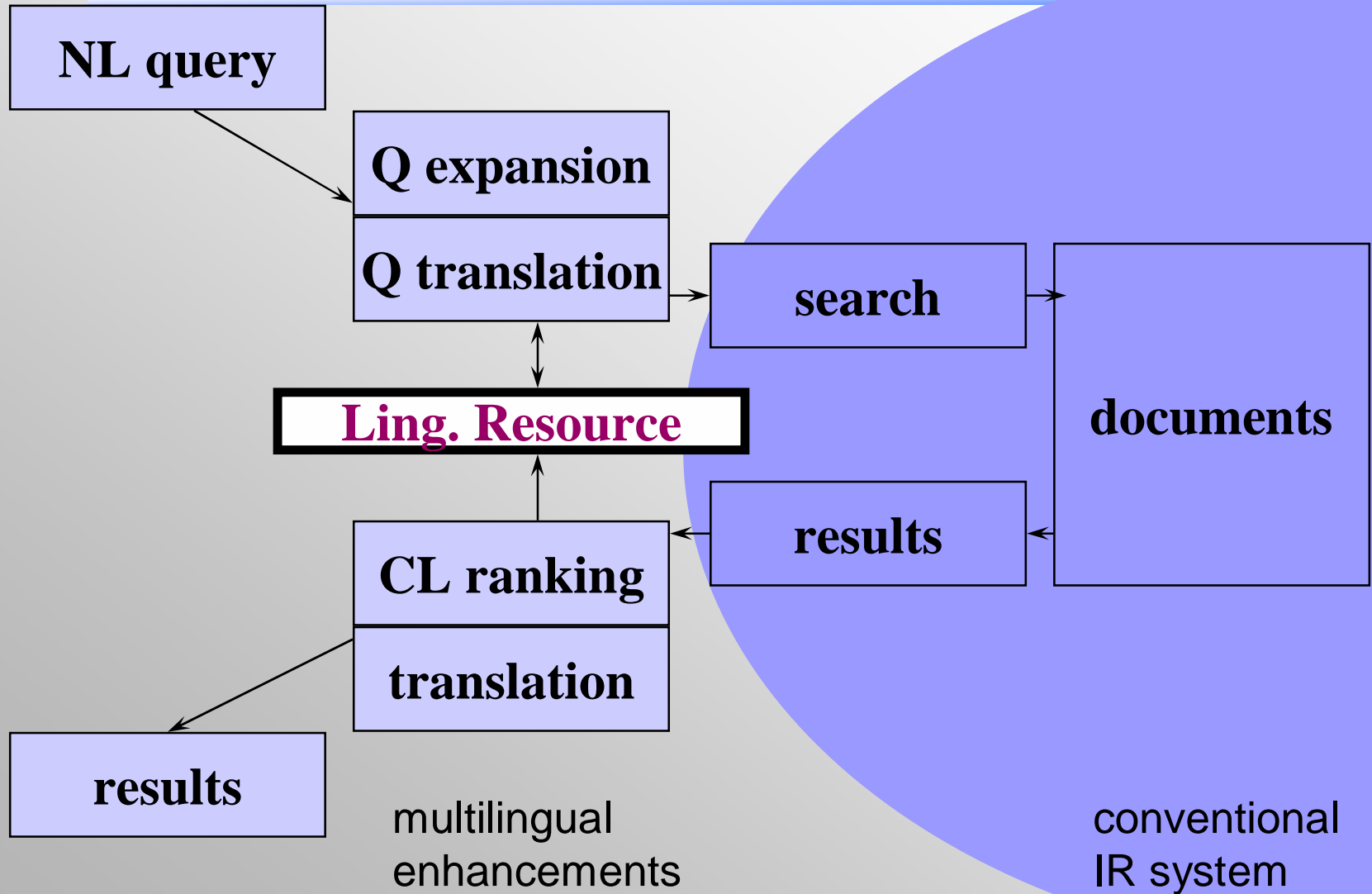2 Personen wurden getötet, mehrere verletzt.

Multilingual Query          Multilingual Text Base

# Architecture

**Problem:**

index (de) ↔ text base (de)

?? ← index (de)

search request (en) →

## 1. Translate **Text**:

| index (de) | ↔ | text base (de) |

| search request (en) | ↔ | index (en) | ← | text base (en) |

## 2. Translate **Index**:

| index (de) | ↔ | text base (de) |

| search request (en) | ↔ | index (en) |

## 3. Translate **Search**:

| search request (de) | ↔ | index (de) | ↔ | text base (de) |

| search request (de) |

# Multilingual Indexing

Options:

- Translate **Texts**
  - Case 1: into all relevant languages
    - Transforms crosslingual problem into monolingual one
  - Case 2: into Pivot language
    - Sub-case of above

- Translate **index terms**

- Translate **query**
  - Case 1: Query into n *mono*lingual document pools
  - Case 2: Query into a *multi*lingual document pool

# Translation of texts 2

- Problem solved by multiplying the texts
  - Make texts available in all languages
  - multilingual (= several monolingual) retrieval

- Feasibility:
  - Required in some applications
    - Patents, multilingual states (EG, Belgium, …)
  - Impossible in other areas (Internet)

- Evaluation:
  - From costly to impossible
  - Results depend on translation quality
    - translation dictionary updates invalidate search on existing document pool (->retranslate everything)

Drogenhandel
in Italien?

| | |
|---|---|
| Drogen | drug |
| Drogenhändler | drug dealer |
| verkaufen | sell |
| Kokain | cocaine |
| Mailand | Milan |
| schmuggeln | smuggle |
| Drogenkrieg | drug war |
| ostdeutsch | East German |
| Schießerei | shooting |

Drug import
to Germany?

Dok 1:
Winston Harton has been identified to be the drug dealer who sold 350 gr of cocaine at Milan train station May 1996.

Dok 2:
Europol Edu has confirmed that the number of drugs smuggled in via scheduled flights from Columbia is increasing every year.

Dok 3:
An der ostdeutschen Grenze ist es wiederholt zu Schießereien zwischen rivalisierenden Banden im Drogenkrieg gekommen. 2 Personen wurden getötet, mehrere verletzt.

# Translation of index 2

- Idea:
  - multilingual Index
    - Analyse query in query language, translate terms
    - Search with all document language index terms
  - (Problem of retranslation of the hits)

- Feasibility:
  - Not feasible
    - Ambiguity of index terms
    - Multiword terms not in index
    - Context dependency of translations
  - => Organise the index as a special resource!

> *Fehler:   mistake, fault, error, bug*
> *nuclear: Kern~, zentral, nuklear*
> *power:   Macht, Kraft, Strom*
> *plant:    Pflanze, Unternehmen*

> *nuclear power plant:*
> *zentrale Kraftpflanze?*

# Translation of query 1

Drogenhandel
in Italien?

Drogen — drug
Drogenhändler — drug dealer
Drogenkrieg — drug war
Kokain — cocaine
Mailand — Milan
ostdeutsch — East German
Schießerei — shooting
schmuggeln — smuggle
verkaufen — sell

Drug import
to Germany?

cocaine
Drogenkrieg
drug
drug dealer
Milan
ostdeutsch
Schießerei
sell
smuggle

Dok 1:
Winston Harton has been identified to be the drug dealer who sold 350 gr of cocaine at Milan train station May 1996.

Dok 2:
Europol Edu has confirmed that the number of drugs smuggled in via scheduled flights from Columbia is increasing every year.

Dok 3:
An der ostdeutschen Grenze ist es wiederholt zu Schießereien zwischen rivalisierenden Banden im Drogenkrieg gekommen. 2 Personen wurden getötet, mehrere verletzt.

# Translation of query 2

- Approach: Translation query
  - Analyse and translate the query terms
  - Search in (monolingual) Backend-System

- Evaluation
  - Backend database stays unchanged
  - Translation changes do not affect document base
  - Cross-lingual component as system frontend
    - contains multilingual linguistic resource
    - Which is also usable for re-translation
    - And can be maintained independently
  - Crosslinguality is tansparent for the users
  - Fine-tuning between frontend and backend required

# Potential CLIR architecture



NL query

Q expansion

Q translation

Ling. Resource

CL ranking

translation

results

search

documents

results

multilingual enhancements

conventional IR system

## III CLIR components

Document Analysis

Index creation

Query processing

    Translation

    Expansion

Ranking

Result re-translation

# Document Analysis 1

- **Language Identification**
  - Needed to load the proper analysis resources:
    - Tokenisers, stop lists, grammars… depend on language
- **Text Classification**
  - Used as a content filter: Focus on relevant domains
  - N-grams or terms as features
- **Named Entity** recognition

  Search for *Bush$_{-Person}$*, not for ‚*bush*'

  Search for *Verona$_{-Place}$*, not for *Verona$_{-Person}$*

- **Conceptual Indexing**
  - Find relevant / ‚good' *and translatable* index terms

# Document Analysis 2

- Resulting index terms must be **translatable**
  - **Base forms** vs. Text forms vs. Stemming
  - **Concepts** vs. Strings
    - Word senses / Homonyms
    - Concepts are often multiwords:
  - **Multiword** index terms?

- Additional Info: **Normalisation** aspects
  - Capitalisation / Diacritics – foreign characters
  - Mapping of transliteration systems
    - *Yeltsin – Jelzin, Abou Moussa – Abbu Mussah*
  - Phonetic search support for names

# Document Analysis 3

Language Identification

Topic Identification

Conceptual Indexing

Information Extraction

Name Search

Analysis

XML Analysis Data Structure

Text text text text text text text text text text text text text text text text text

**Language:** .....
**Topic:** .....
**Concepts:** .....
**Persons:** .....

Backend System

**XML Structure for analysis result (backend independent)**
**Enriching documents with intelligence**

# Document Analysis 4: Example

# Index organisation 1

- Case 1: one document pool *per language*
  - Needs a Language Identifier at indexing time
  - monolingual index per Pool (de, es, fr, ...)
  - Multiple queries: *one* query *per pool* (de -> de, en -> en, ...)
- Case 2: a *common* document pool
  - => multilingual index!
  - Problems with interferences between languages
  - Problems with Stop lists
  - One query: *all* queries go to one pool
    - de -> deenfr, en -> deenfr, ...
    - (Language Identifier at runtime)

| | |
|---|---|
| *de-en:* | *kind post* |
| | *not as* |
| *de-fr:* | *Seine ...* |
| *en-fr:* | *car ...* |

# Index organisation 2

- Preselection of documents into monolanguage pools
  - Better Indexing
    - Both for conventional (Stoplisten) and linguistic indexing
  - Simpler search
    - Better Performance: less processing at runtime
    - No inter-language clashes
  - But: Multiple search required

- Option depends on Backend System capabilities
  - Possibility of multiple data pools
  - Language as document attribute
  - …

# Query **Analysis**

**Before it can be translated
the query must be <u>analysed</u>!**

- Find key concepts and translatable units

  *„drug and heroin dealers in the former German Democratic Republic"*

  – Resolve gapping (*drug dealer / heroin dealer*)
  – Remove untrabslatable (function) words after analysis
  – Identify multiword units
    - Drug dealer / heroin dealer / German Democratic Republic
    - Ingle wod analysis fails …
  – Remove noise words (*former*) (?)
  => Precise linguistic means are required

# Query Translation

- Success of CLIR depends of the **quality of translation**:

    *Nuclear -> zentral  AND power -> Macht AND plant -> Pflanze* ☹
    *Kernkraftwerk* ☺

- **Specialised resource** for the relevant domain
    - No general purpose dictionary helps
    - Even domain specific glossary can be suboptimal

        Glossary has *Droge -> drug* but text has *narcotics*
    - Resource must match the document base

- Problem of **unknown terms**
    - No hits if unknown parts are left in the query, except
    - Proper Names: translate / transliterate …
        - *Persia -> Persien*   (must fit the indexing strategy)

- Problem of target language multiwords (flexer)

# Options for Query Translation

## Approaches

- Translation by **parallel corpora**
- Translation by **MT**
- Translation using bi-/multilingual **dictionaries**
- Translation using **concept hierarchies**

## Problem

- Terminology must match the texts!!

  (Aventinus: coverage ~ 16%, TQPro: coverage ~ 30%)

- Missing translations

Termbank has:    *drug -> Medizin*
Texts have   :            *-> Drogen*

search result: lausy

# Query translation by corpora

- Procedure
  - provide parallel corpora (Bible …)
  - Alignment on sentence level
  - Identify corresponding terms (bilingual term alignment)
    - (improvement: multiword terms, lemmatisers)
  - Use those terms for search

- Evaluation
  - Translations match the texts well
  - but: parallel corpora often not available
    - (if so: MLIR possible)
    - New data: gaps in translation

# Query translation with MT

- Problems with MT translations
  - short Queries have **few context** for disambiguation
    - Main tproblem: multiword terms

  - MT-dictionaries must fit the **domain**
    - Usually not the case: Babelfish etc. with focus on GV
    - *several* MT-Systems do't help (Eurospider Trec8)
    - But: tuned MT came out well in MUMIS (medical IR)

  - 1:n translations must be **disambiguated**
    - How? (By subject area tag?)

  - MT-Systems don't have **Ontologies** for Q-Expansion

  => (better way: Generate MT-dictionaries out of CLIR-Resource)

# Query Translation using dictionaries

- **Ocnventional MRD's**
  - Look up the Article, Search with all Words?

    much Noise (Univ. New Mexico)

  - Select the reading: Search with „best" reading?

    often incorrect (Twenty-One in TREC)

  - Formal use often difficult
    - Intended for human lookup
    - Much implicit information (examples …

**=> Termbanks / specialised resources**
  - Selection by using subject area tags
  - Expansion by using specialised links
    (PAROLE / SIMPLE / **NORNA**)

# Example: Dictionary Article (LKG)

## Urteil

n **1.** judg(e)ment, *der Geschworenen*: verdict (*a.* fig.) (*Straf~, Strafmaß*) sentence (*Scheidungs~*) decree: **das ~ verkünden** pronounce judg(e)ment (*od.* sentence); -> **fällen 2.** (*Ansicht*) judg(e)ment, opinion: **sich ein ~ bilden** form a judg(e)ment (**über** *acc.* about, on); **darüber kann ich mir kein ~ erlauben!** I am no judge (of that)!

## Ton

m **1.** tone, sound (*a. Film, TV, Ggs. Bild*); F **er hat keinen ~ gesagt** he didn't say a word; **2.** (*Redeweise*) tone: **ich verbiete mir diesen ~!** I don't take that tone with me!; **3.** *ling.* accent, stress: **den ~ legen auf** (acc) *a.fig.* emphasize **4.** (*Farb~*) tone (*a. phot.*), shade: **ein rötlicher ~** a tinge of red.

- Univ NewMex takes all targets
- Twenty-One takes first: *judg(e)ment* (0 hits) und *tone*

# Translation with concept nets

- Technique:
  - Vreate a conceptional hierarchy („ontology") where nodes combine clusters of monolingual entries
  - Look up concept for term
  - Look up target terms for concept

- Advantage
  - Integration of Expansion and translation

- Problem
  - Find right concept link (reading)
    - word sense disambiguation: difficult for short Queries
    - Subject area? Often not available, e.g. WordNet

# Example concept net



enclosure
drug
   hard drug
   medicine; medication; medicament; medicina
   anaesthetic; anesthetic; anesthetic agent; ana
   addictive drug
   abortifacient; abortifacient drug; abortive; abor
   tobacco; baccy
   substitute drug
   narcotic; narcotic drug
      cocaine type drug
      cocaine; cocain
        crack; crack cocaine; Hong Kong rock
        cocaine base; free-base
      opiate
        morphine; morphia
          morphine base
          heroin; diacetyl morphine; diacetylr
            heroin base
          benzylmorphine
        opium
        ▶ raw opium
        benzethidine
        codeine
          acetyldihydrocodeine
        acetylmethadol; LAAM
      methadone; methadon
   psychoactive drug; mind-altering drug
   sedative; sedative drug; tranquilizer; tranquilliz
   soporific; hypnotic
   soft drug
   dance drug

# Translation: Disambiguation Strategies

- „natural Selection"
  - *plant -> Pflanze / Werk*,
    *power plant -> Machtpflanze* in no text

- Subject area selection
  - typical MT-Technology: „take terms from IT"

- Specialisation of the linguistic resource
  - WordNet: activate only a segment (IT),
    eliminate potentially wrong readings

- Corpus-Selection
  - bilingual Reference Corpus
  - Take translation with is in the best fitting cluster

# Query Expansion

- **Conceptual expansion**
  - Searcher's and writer's formulations need not match
  - Query Expansion: enrich the query by useful alternatives
    Synonyms, narrower or similar terms:
    *vehicle? -> car, truck, BMW, …*
  - A multilingual concept net helps to do this
    - But: Which concept to pick in case of homonyms

- Expansion on the **document language** side
  If source language synonyms have the same translation,
  expansion is useless

  | Persia -> Persien | Persia -> Persien |
  | Iran -> Persien ☹ | SYN Iran ☺ |

- **Morphological expansion**
  If index has full forms: create all full forms

# Expansion techniques

- Document / corpus based
- Concept hierarchy based
  - How deep to go?
- Interactive?
- Expansion is language-specific!

# Example: CL Search in Sensus

# Search and ranking

- Search requires **multiple queries**
  - Select proper index per language
    - English query to Englich index
      German query to German index, …
  - Do textbase/language-specific (statistical) expansion
  - Come back with several document sets
    - E.g.: 30 hits in French, 12 in German, 25 in English

- **Ranking**
  - Is a hit in French better than a hit in English?
  - Current ranking approaches are monolingual
    - Based on a common set of terms (found/nonfound):
      What if the term sets are completely disjoint?
  - Alternatives:
    - Rank the languages internally, present them separated
    - Find a common basis of comparison (e.g. translation)
    - Create a sequence of queries with released constraints

# Result translation

- Foreign language search results must be brought back into the interface language! (NORNA)



- Options:
  - Machine translation
  - Term replacement
- Target: understand the content / judge relevancy

# Example: Term Substitution Output

# Complete System

Language identification

**Intelligent Indexing:**
Concepts, NamedEntities, Classification
Normalisations, Phonetics

Query

Translation

Lex

Translation

ML Ranking

Index

Document Base

# Consequence for the resource

Engineering  aspects when implementing a
multilingual information processing system

- – Focus on cross-lingual information retrieval
- – Features of such a system containing multilingual
  extensions

- – Requirements for **resources** to support such a system
- – Maintenance and build-up of such resources

# Resources for Crosslingual Retrieval

CLIR are the most complex ones for resources!

1. **Document Analysis** requires linguistic analysis (TREC)
   - lemmatiser, head-modifier-parser, word sense disambiguation

2. **Query Expansion** + Feedback  improves search results significantly (TREC)

3. **Crosslinguality** requires translation of queries and hit documents (CLEF)

All these steps need **linguistic resources**

Resources must be **maintained**
   - Tools to build them
   - Tools for maintenance: Single point of administration
   - Tools for exchange: Easy exchange

# Resources for Document Analysis

**Monolingual** lexical resource, covering

Minimum: **Morphosyntax** for index term creation
- POS and inflection
- (additional requirement for CLIR: multiword terms!)

+: information on **argument structures** etc.
- Dependency trees / head-modifier analysis

+: Some **semantic and collocation** information
- Named Entity recognition
- Word sense disambiguation

# Resources for Query Expansion

- Expansion needs **relations** of a **ConceptNet**
  - Strict hierarchical relations (automatic expansion)
    - Synonyms, narrower terms
  - Additional relations (interactive expansion)
    - Part-of, see-also, ...

- ConceptNet must be **descriptive**
  - Contains all searchable terms of a corpus
    - Unlike thesaurus (which is normative)
  - => To be built for each single application
    - ConceptNet Building Tools are important!

# Query and Document Translation

- Query Translation needs **terminology data**
  - Concept-oriented term banks
  - 50-70% multiword terms

- Document Retranslation needs **MT lexicons**
  - Monolingual and bilingual (transfer) information
  - Reading vs. Concept organisation
  - Transfers: tests&actions, statistical clusters

- Translation resources must be **consistent**
  => Exchange requirements between system components

- Translation resources must **fit the text base**
  => Building tools

# Resources: Requirements

CLIR needs a resource like the following:

- Organised as a **Multilingual ConceptNet**
  - (language-specific wordnets do not really help in CLIR)

- Nodes: **Concepts**
  - Concepts are abstract, defined by definitions
  - Concepts have **multilingual expressions** („concept")
    - -> **Query translation**
  - Each monolingual concept has its terms
  - Terms have (monolingual) **linguistic descriptions**
    - -> **Indexing / Analysis**

- Links: **Relations**
  - Strict hierarchies (is_a), additional relations
    - -> **Query expansion**

# ConceptNet vs. Dictionary

- ConceptNet is not a dictionary
  - **Organisation**
    - Dictionary is organised in **lemmata**
      - All concepts with same lemma form a dictionary *article*
    - ConceptNet is organised in **concepts**
      - Lemmas for different concepts are different entries
  - **Relations**
    - Dictionary does <u>not</u> explicit hierarchical relations
    - Concepts use relations to build **hierarchies** & networks
- But: It uses dictionary information
  - (basic linguistic information, subject area)

# ConceptNet vs. Thesaurus

- ConceptNet is not a thesaurus / ontology
  - Thesaurus is based on **canonical** terms
    - Thesaurus terms *represent* concepts
      other words are not to be used
  - ConceptNet is based on **used** terms
    - Many terms can represent a concept
    - Analysis base for a ConceptNet is a *corpus*
      all terms in the corpus are in the ConceptNet

- But: Both model **relations** between concepts
  - ConceptNets have more terms (synonyms)
  - ConceptNets have more relation types (EuroWordNet)

# ConceptNet vs. Terminology Entry

- Terminology entry is canonical, not descriptive

- Assumes 1:1 relationship between languages

- Does not provide:
  - Explicit hierarchical structure
  - Detailed linguistic annotations
  - Idiosyncratic relations between languages
    - Esp.: transfer tests and actions
      (can differ from en > fr and da > de within one term!)

- But: Term metamodel similar to Concept nodes
  - Cf. TBX standard – not MILE compatible ☺

# ConceptNet vs. WordNet

- Similar: based on a kind of SynSet representation
- Definition of concept / synset node

- But:
  - WordNet is not multilingual
    - One hierarchy (per language)
    - EuroWordNet – ILI: ?
  - No linguistc annotations per term
  - No relationships between synset terms
    - (On transfer level: tests & actions / preferences)
    - On crossreference level: abbreviations, headwords, …

- WordNet is not a <u>Word</u>Net

# Are ConceptNets domain-independent?

- No. Multilingual ConceptNet is **domain specific**

  – Concept hierarchy is purpose-dependent

  – Users will do searches in their domain
    - Search terms sometimes are user-specific
    - Prepare resource based on user-corpus

  – General purpose terms are difficult to search & translate
    - 1:n translations introduce noise in the search
    - People tend to use specific terms for searching

- Link different ConceptNet instances
  – Use a common **top level ontology**

# Are ConceptNets language independent?

- **Yes!**
  - As far as CLIR applications are concerned
    - Language-independent definition of concepts
      - Definition by: defintion text, part-of-speech, subject-field
    - Language-independent definition of relations
      - *Floppy_Drives* are *part_of computers* in any language
- **No!**
  - As far as linguistic accuracy is concerned
    - Concepts are defined by (connotational) context
    - Concept hierarchies are sometimes language-specific
      - (legal system, educational system, ...)

=> Cover concepts in one ConceptNet
   Some concepts do not have terms for some languages

# Result: Linguistic Design

- Organise resources in **concepts**
  - Homonyms create noise in retrieval

- Base concepts on an ontology / hierarchical **net**
  - Start from a top level ontology (EuroWordNet)
  - Application specific hierarchy below the top level
  - Allow for additional relations (subset of EuroWordNet)

- Describe concepts with **terms (entries)**
  - Cover all relevant terms in the domain

- Enrich concepts with **multilingual** terms / entries
  - To be used as translations in query and text translation

- Describe terms with linguistic **annotations**
  - Needed for compilation into linguistic applications
  - Based on OLIF / EAGLES / MILE

# IV ConceptManager: Software Structure

Software Structure

Concept Manager Client
(Java)

Java rmi

Concept Manager Backend
(Java)

multi-editor
UTF8

SQL Database
(mySQL, SQLServer, Access)

# ConceptNet Manager: GUI

# Ontology Tree: Menu

# Active Concept: Concept Features

# Active Concept: Term Table

# Term Table: Linguistic Info

# Term Table: Cross Reference

# Term Table: Transfers

# Search – Result in Ontology Tree

# Advanced Query Dialogue

# Advanced Query Result List

# Concept Manager: Functionality

- **Concept level coding**
  - Add / remove concepts
  - Link concepts to each other

- **Term level coding**
  - Add / remove terms

- Linguistic level
  - Code linguistic features
  - Code translation information
  - Code crossreferences

- **Import entries**
  - OLIF import format
  - Unlinked folder

- **Delete nodes / trees**
  - Save subtrees in folder

- **Compile entries**
  - Export for query translation
  - Export for MT

- **Query the DB**
  - Simple queries
  - Extended Boolean queries

- This resource needs to be **administered**
  - Single point of administration!
  - => ConceptNet Administration Tool

- This resource must **feed the applications**
  - Resources for linguistic corpus analysis
  - Machine Translation dictionaries
    - Concept oriented vs. Reading oriented organisation
    - Bilingual directed vs. Multilingual
  - => Exchange of (parts of) resources
    - Relevance of Exchange format (OLIF)

- This resource must be **set up**!
  - Because it is application-specific
  - => Tools for corpus extraction and build-up

# Resource Build-Up

# More Intelligence ...

- Tools to create **linguistic annotations** for terms
  - Syntactic / semantic descriptions

- Tools to group terms into **concepts**
  - Monolingually (clustering, head-modifier-analysis)

- Tools to define **relations** between concepts
  - Hierarchical / non-hierarchical

- Tools to map **multilingual** concepts
  - (Extensions of current Bi-Extract)

(Business Constraint:
  - Create a ConceptNet as described with 10K concepts in 3 languages in 2 weeks time)

# Evaluation

- Only partial evaluations possible
- E.g.: IR using syntactic or semantic information:
  - POS Support
    - *Plays$_{-Noun}$ of Shakespeare* vs. *Child plays$_{-Verb}$*
    - No improvement due to performance of POS taggers in some language -> improve this first
  - Named Entity
    - Names like *Akrinski* in German texts are not mixed up
    - *Fisher$_{-person}$* increases precision but reduces recall if NE is not safe enough
  - First improve what you already know is a source of deterioration

**Thank you for your attention**

g.thurmair@linguatec.de
www.linguatec.de