FACHHOCHSCHULE STUTTGART
HOCHSCHULE DER MEDIEN

# Automated Retrieval of Information in the Internet by Using Thesauri and Gazetteers as Knowledge Sources

Prof. Dr. Wolf-Fritz Riekert
Fachhochschule Stuttgart – Hochschule der Medien (HdM)
University of Applied Sciences Stuttgart – School of Media

mailto:riekert@hdm-stuttgart.de
http://v.hdm-stuttgart.de/~riekert

---

## CONTEXT

FACHHOCHSCHULE STUTTGART
HOCHSCHULE DER MEDIEN

Concept and prototype development under commission of the German Federal Environment Agency, Berlin (1996-1998) by:

- Research Institute of Applied Knowledge Processing (FAW), Ulm (Riekert, Wiest, Fuchs, Klingler)
- Condat AG, Berlin (Nouhuys, formerly: CAdMAp GmbH)

The concepts described served as an input into the following German environmental information systems:

- German Environmental Information Network (GEIN)
- Geographical Information System Environment (GISU)

Software development and maintenance now:

- Ernst Basler + Partner (GISU)
- Sema Group (GEIN)

---

## INFORMATION RESOURCES IN THE INTERNET

FACHHOCHSCHULE STUTTGART
HOCHSCHULE DER MEDIEN

Categories of information resources:

- multimedia documents
- data
- application services

Supply exploding

- Problem: orientation ("lost in hyperspace")
- powerful search tools required

---

## SEARCH ENGINES

FACHHOCHSCHULE STUTTGART
HOCHSCHULE DER MEDIEN

Search engines are based on a full text index which intentionally covers the whole Web

- Retrieval via Web browser (string search)
- Index maintained by "robots" "crawling" along hyperlinks
- No additional efforts required from information suppliers
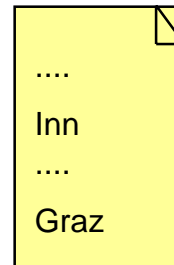
But:

- Search terms are interpreted only textually
- No semantic interpretation
- Full text index can only be used for textual resources

## EXAMPLE

Query:

"*Accommodation, Styria*"

Search Engine

....

Inn

....

Graz

---

## METAINFORMATION SYSTEMS

Metainformation systems support semantic criteria for indexing and retrieval:

- Thematic references (e.g., "Accommodation")
- Spatial references (e.g., "Styria")
- Temporal references (e.g., "July 11-12, 2002")

Indexing (i.e., entering the metainformation) is done manually by the system administrator or information suppliers:

- Higher information quality (compared to search engines)
- Higher workload imposed on system administrator or information suppliers

---

## EXAMPLES OF METAINFORMATION SYSTEMS

Examples of metainformation systems (taken from the environmental domain):

- GEIN: German Environmental Information System (Germany)
- GISU (Meta Component): Geographic Information System Environment (German Federal Environment Agency)
- UDK: Environmental Data Catalogue (Germany, Austria)
- CDS: Environmental Catalogue of Data Sources (European Environmental Agency)
- NGSC: National Geospatial Clearinghouse (USA)

---

## METAINFORMATION SYSTEMS: EXAMPLE GEIN (WWW.GEIN.DE)

GEIN Thesaurus Search - Netscape

File   Edit   View   Go   Communicator   Help

**Topic**
please add search words.
type any word :
how to match?                many selected terms ("or")

**Area**
please add search words.
type any area name :
how to match?                many selected area names ("or")

**Time**
please add a date
type any date :
how to match?                single date or period

you may search with the selected terms now english :          search now

**The Portal of German Environmental Information**

Document: Done

## SPECIFICATION AND PROCESSING OF SEMANTIC CRITERIA

Requirements

- Vocabulary for the specification of thematic, spatial and temporal references of information resources
- Techniques for the automated processing of thematic, spatial and temporal references

Approach

- Thesaurus to support specification and processing of thematic references
- analogously: „Gazetteer" to support specification and processing of spatial references
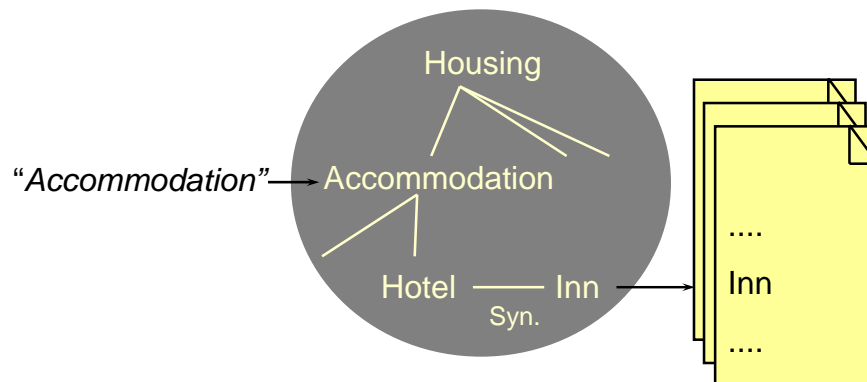- Handling of temporal references: relatively easy, not an issue of this talk

## THESAURUS

A Thesaurus is a structured collection of terms with the following properties:

- Terms provide a controlled vocabulary for the specification of thematic references,
- Terms can be used for both indexing and retrieval.
- Terms are more than simple keywords.
- Terms form a semantic network established by:
  - ⇨ synonym relationship (inn - hotel)
  - ⇨ generalization hierarchy of broader / narrower terms (accommodation - hotel)
  - ⇨ linkage via related terms (accommodation - tourism)

## THESAURUS-SUPPORTED QUERY PROCESSING

Query          Thesaurus          Information Resources



"Accommodation" → Accommodation

Housing

Hotel —— Inn
Syn.

Inn

....

....

## BLACK BOX SEARCH PROBLEM: A THESAURUS CAN HELP

Problem: Information resources are searched for by using a form in most metainformation systems ("black box search")

- It is not clear which level of detail is required while specifying a query
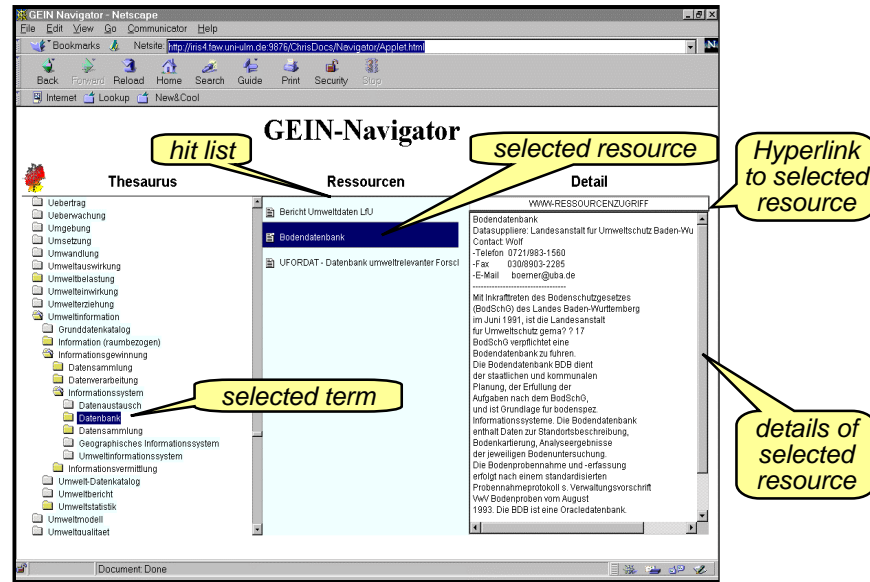  - ⇨ Many casual users dislike form-based search interfaces

Requirement: Hierarchical directories to access the information resources

- However: Manual maintenance of hierarchical directories very time-consuming

Solution: Use a thesaurus for the automated generation of a hierarchical directory

Example: GEIN Navigator (prototype developed at FAW Ulm)

# PROTOTYPICAL GENERATION OF A HIERARCHICAL DIRECTORY

hit list

selected resource

Hyperlink to selected resource

selected term

details of selected resource

---

# A PROCEDURE TO GENERATE A HIERARCHICAL DIRECTORY

- Create a "weeded" thesaurus consisting of all relevant terms, i.e.:
  - ⇨ take all terms used as an index for existing information resources,
  - ⇨ add recursively all broader terms,
  - ⇨ disregard all other terms
- Display thesaurus in a hierarchical presentation (Windows Explorer-like), starting from "toplevel terms"
- Special highlighting indicates which terms
  - ⇨ directly lead to hits,
  - ⇨ possess narrower terms leading to hits
- Provide navigation paths to the metainformation records and from there to the original information resources

---

# METAINFORMATION SYSTEMS VS. SEARCH ENGINES

Metainformation system:

- Easy retrieval by using semantical criteria
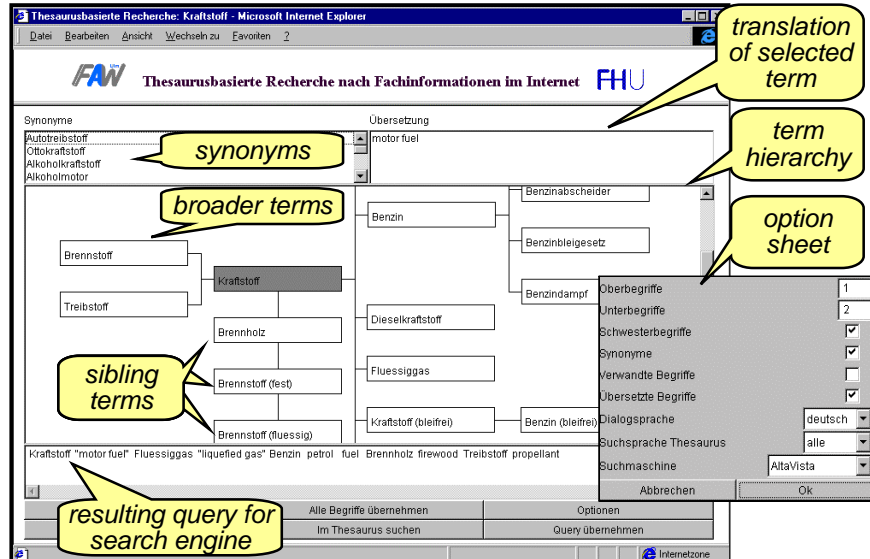- But: Indexing very expensive for administrators or information suppliers

Search engine:

- Indexing very easy, no work imposed on suppliers
- But: only textual processing of search criteria

Synthesis:

- Combination of the advantages of search engines and metainformation systems: Thesaurus-based preprocessor for search engines

---

# COMBINE THE ADVANTAGES

| | Indexing inexpensive | Semantic processing of search terms |
|---|---|---|
| search engine | ✘ | — |
| metainformation system | — | ✘ |
| search engine with thesaurus-based preprocessor | ✘ | ✘ |

## THESAURUS-BASED PREPROCESSOR FOR SEARCH ENGINES

*translation of selected term*

*term hierarchy*

*option sheet*

*synonyms*

*broader terms*

*sibling terms*

*resulting query for search engine*

---

## SPATIAL REFERENCES IN TRADITIONAL SYSTEMS

Problem: Spatial references in traditional systems are handled very poorly (if they are handled at all):

- Rigid vocabulary
  - ⇨ Usually only one single spatial reference system supported (coordinates only, names only)
- No intelligence
  - ⇨ It cannot be recognized if one region encloses another

Solution: Specification of spatial references through geographic objects (geoobjects)

- Geoobjects are more than names or coordinates
- They possess both names and coordinates
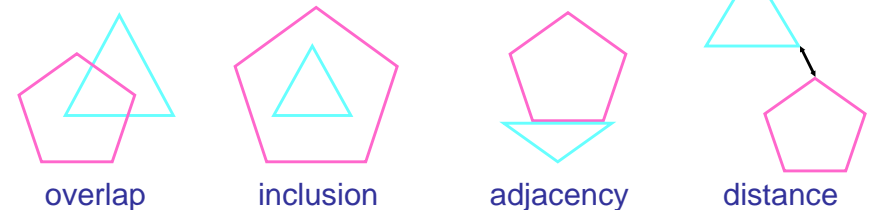- Geometrical and topological relationships can be computed

---

## GAZETTEER

A Gazetteer is a structured collection of geographic data objects (geoobjects)

- Geoobjects can be used to specify spatial references
- Spatial reference: n:m-relationship between information resources and geoobjects in the gazetteer
- Geoobjects may possess the following properties:
  - ⇨ name (e.g., "Styria")
  - ⇨ geometry (e.g., coordinates describing a polygon)
  - ⇨ type (e.g., "state")
  - ⇨ unique identifier (e.g., administrative code)
  - ⇨ optional: hierarchy (e.g., administrative hierarchy)
  - ⇨ optional: synonymous names, translated names

---

## TOPOLOGICAL AND GEOMETRICAL RELATIONSHIPS

Topological and geometrical relationships, e.g.,

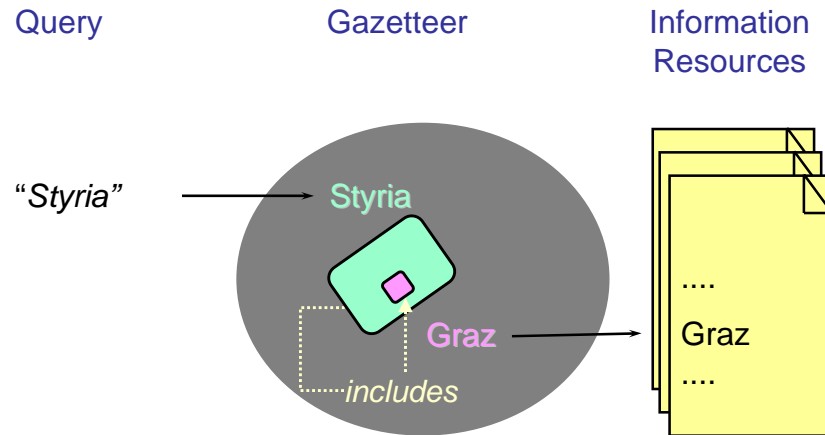

overlap        inclusion        adjacency        distance

can be derived from geometry (i.e., coordinates)

With the help of these relationships,
a flexible geographic vocabulary can be used for

- indexing and
- retrieval purposes

## GAZETTEER-BASED RETRIEVAL

Query          Gazetteer          Information
                                  Resources

"*Styria*" → Styria

Graz → Graz

*includes*

....
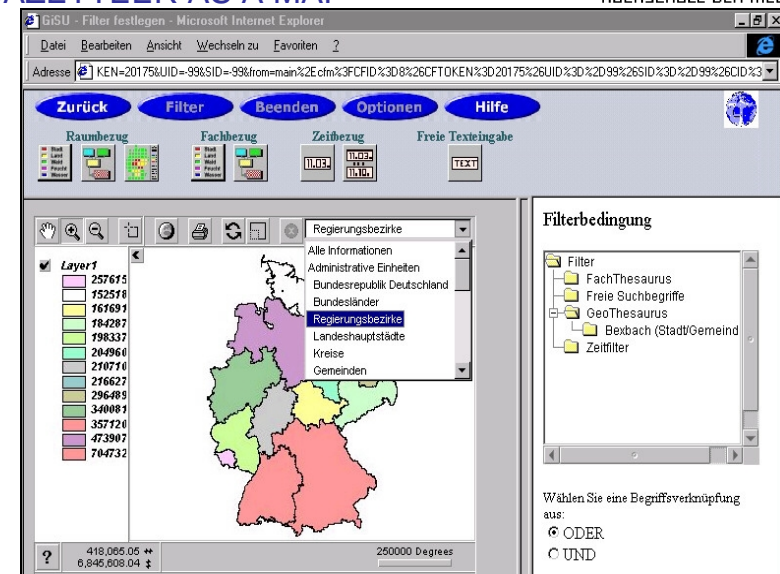Graz
....

---

## GAZETTEER: APPLICATIONS

- Browser for geographical names
- Specification of spatial references on a cartographic interface
- Automated extension of queries: information resources in a certain geographic or topological neighborhood can be included into the scope of the query
- Easy transition between textual and geometrical representation of spatial references
- Text analysis for the automated spatial indexing of textual information resources
- Graphical display of spatial references as "footprints", e.g. to visualize a search result

---

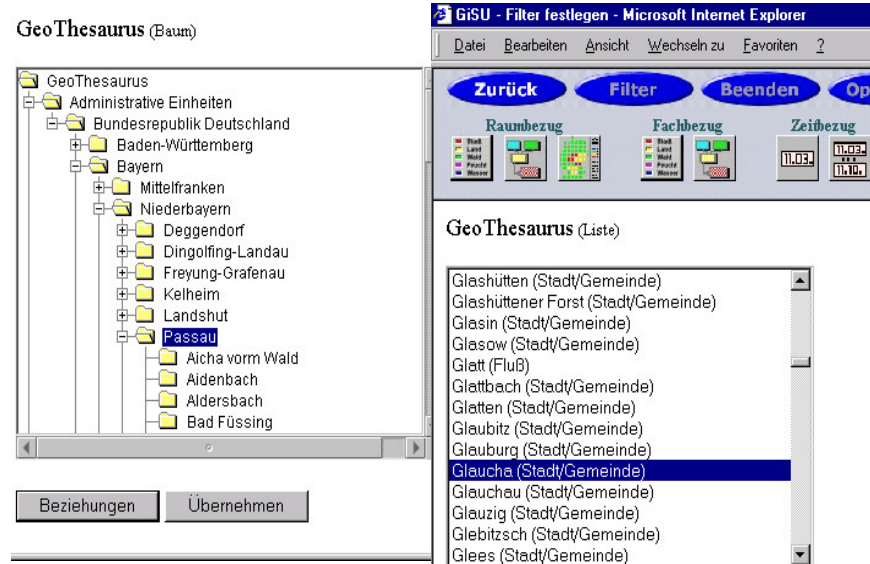## THE GERMAN "GEOTHESAURUS": AN EXAMPLE FOR A GAZETTEER

The German Federal Environment Agency developed a Gazetteer known as "Geothesaurus"

- It contains about 100 000 administrational, topographical and environmental entities in the form of geoobjects
- The geometries are rastered in a 3×3 km² grid
  - ⇨ The whole geothesaurus can be represented in a relational database (no "geographic information system" required)
- Application in two German environmental metainformation systems:
  - ⇨ GEIN (German Environmental Information Network)
  - ⇨ GISU (Geographic Information System Environment)

---

## EXAMPLE GISU: DISPLAYING THE GAZETTEER AS A MAP

GeoThesaurus (Baum)

GeoThesaurus (Liste)

GiSU - Filter festlegen - Microsoft Internet Explorer

- Metainformation systems and search engines can be enhanced considerably by thesauri and gazetteers
- New attractive user interfaces: maps, directories, network graphics instead of blackbox search
- Flexible vocabulary for the specification of thematic and spatial references
  - ⇨ Automated reformulation, extension, and translation of terms and geographic locations
- Thesauri and gazetteers are knowledge structures which are relatively stable and application-independent
  - ⇨ maintenance relatively inexpensive
  - ⇨ reusable in multiple applications
- Investments in thesauri and gazetteers pay