






Babbage Institute for Knowledge- and Information Technologic

**efficient indexing
for
advanced document
management**



Vervenne D.
BIKIT, Gent

Some milestones...

- Clay tablets 3500 BC 
- Printing on paper: Gutenberg 1454 
- Typewriter by Remington, 1874 
- windows-PC at Xerox Parc, 1973 
- Wwweb at Cern/Ncsa, 1993 

Vervenne ©


**What is indexing ?
How can we index ?
Thesaurus-based indexing ?**

document indexing: what ?

Indexing is the process of creating **meta-data** about the document(s)

these meta-data represent knowledge about the documents, in order to classify and retrieve them (cf search engine).


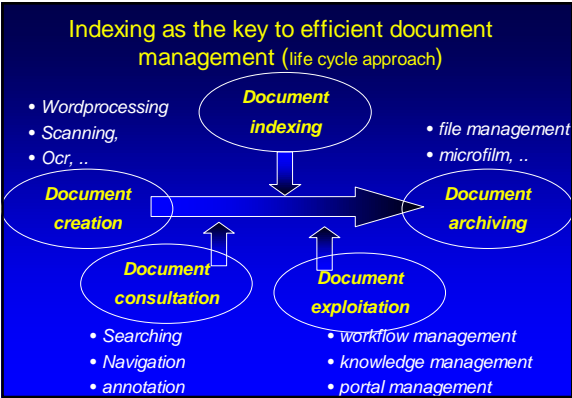
Labels on files represent meta-data.



document indexing: history ?

When printing was invented, documents had no title.

Titles have been added to books for indexing purpose in order to **classify** them by librarians.

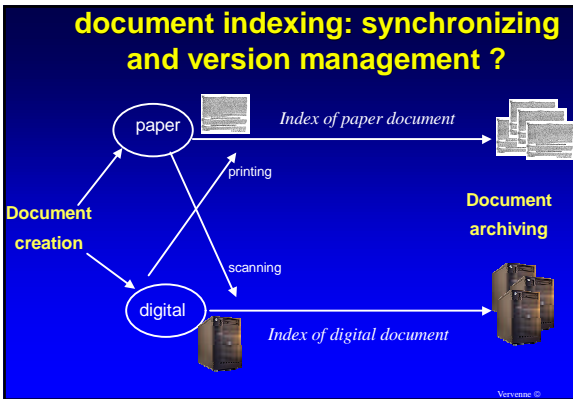
Babbage Institute for Knowledge- and Information Technologe

document indexing: problem ?

How to classify documents such that they are retrievable in the most efficient way ?

- using authority files with reference terms or codes ?
- linking all relevant semantic associations ?

Verveen ©



What is indexing ?

How can we index ?

- manual
- automatic

thesaurus-based indexing ?

Verveen ©

1: manual document indexing

Authors of documents can add meta-data, e.g.:

- give a **relevant** name to the file
- add **keywords** in pre-defined fields
- add headers or footers with relevant meta-data

(specialized) **Indexers** can add meta-data, e.g.:

- fill in external file forms
- fill in separate databases with links to the doc.

Verveen ©

manual document indexing: exampl.

Filenames can contain meta-data about their content:

e.g.:

- *heldere_fazantensoep_met_sherry.txt*
- *final-report-Octopus-version3.pdf*

some names have an internal logic (and are mostly cryptic for non-experts)

e.g.:

- *1996sep10_eindtermen.doc*
- *MOC-051299-DVF123a.doc*

Verveen ©

(manual) document indexing

Indexing is traditionally a (formal) activity done by the expert **librarian**: s/he wanted to classify books and journals in order to retrieve them easily (by paper forms)

Cohen M. BCS-78435
An introduction to Logic
 1972
 Routledge & Kegan,
 London
 ISBN 0 7100 1197 9 pp465

Verveen ©

Babbage Institute for Knowledge- and Information Technologie

Manual Indexing info in the document

E.g., through customized fields in the word-processor

Manual Indexing info in the document

E.g., Default Properties field of MsWord

field indexing with manual support of thesaurus

2: Automatic document indexing

How can we extract meta-data automatically ?

- Statistical approach ?
- Neural network approach ?
- Thesaurus-based approach ?

statistical indexing

The **frequency** of all individual words in the document is counted and stored in an alphabetic index.

This approach is used for 'full text' search e.g. Internet search.

Our cat

Our cat is very lazy. She sleeps all the time. The mouse and the dog, living next door, are not afraid of our cat. The result is a very friendly neighborhood.

Full index

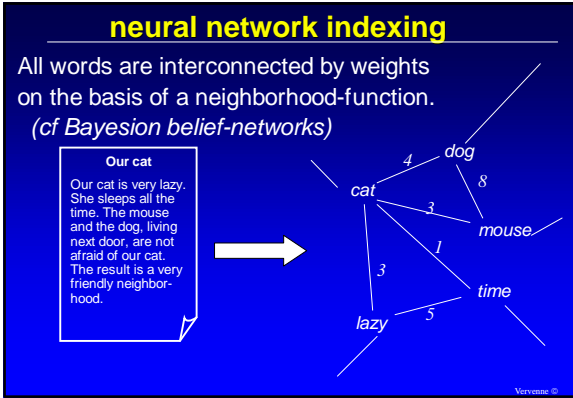
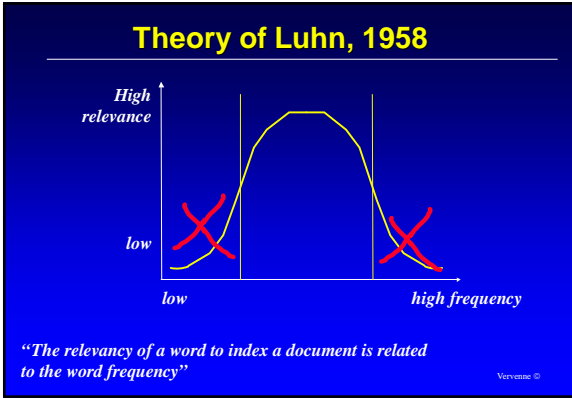
Afraid 1
cat 3
dog 1
is 3
lazy 1
living 1
neighborhood 1
next 1
our 3
result 1
she 1
sleeps 1
time 1
the 4

Stopword list

Manual Indexing info in the document

Meta Tags in HTML for Web Search Engines

Babbage Institute for Knowledge- and Information Technologe



Neural networks: usage

Neural networks are efficient for very large collections of documents ...

... But the indexing process operates as a black box ...

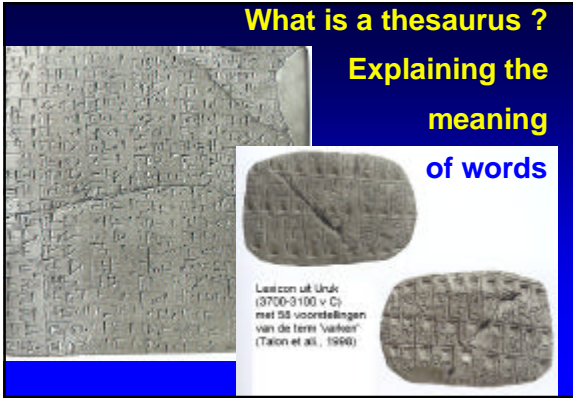
Cf. Autonomy

What is indexing ?
How can we index ?
Thesaurus-based indexing ?

What is a thesaurus ?

“The thesaurus (...) is a vocabulary of controlled indexing language, formally organized so that a priori relationships between concepts are made explicit, to be used in information retrieval systems, whether these are databases, or printed indexes or catalogues”

Aitchison et al., 1997



Babbage Institute for Knowledge- and Information Technologe

What is a thesaurus ? : more than synonyms

Webster's Dictionary refers to Greek word 'thesauros' (Gr. θησαυρος) to be translated as "treasure room" :

L., treasure, store, collection, Gk thesauros: A book containing a store of words or of information about a particular field or set of concepts.

Webster Third New International Dictionary, 1971, Vol. III, p. 2374.
(MS-Word does not support a thesaurus !!)

First thesaurus-idea from John Wilkins, 1668
'*Essay towards a Real Character and a Philosophical Language*'.

Encyclopedia Britannica Macropaedia, Vol. 5, 1974, p. 719.

What is a thesaurus ? : 2 ISO norms

International Organisation for Standardization.
ISO 5964: Guidelines for the Establishment and Development of Multilingual Thesauri. 2nd ed. - Geneva: ISO, 1985.

International Organisation for Standardization.
ISO 2788: Guidelines for the Establishment and Development of Monolingual Thesauri. 2nd ed. Geneva: ISO, 1986.

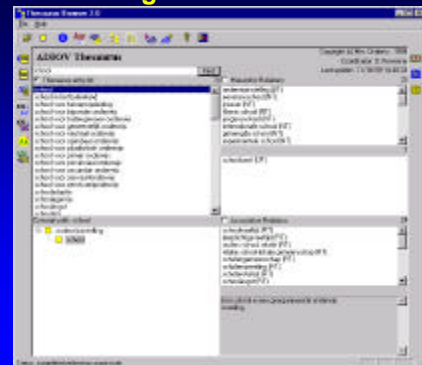
What is a thesaurus ? ..ISO-approach

A semantic network of words with standardised relations:

- **hierarchic relations** (broader/narrow terms)
e.g., *Europe --> broader term --> Belgium*
- **equivalent relations** (use/use for relations)
e.g., *VRT -->use for --> Vlaamse Radio Televisie*
- **related terms** (rt relations)
e.g. *stripverhaal --> rt --> ontspanning*
- **scope notes** (example)
e.g. "*an AKB is a monitor module for automatic adjustment of the phosphor radiation*"

IKEM Thesaurus Manager Module

ISO-compatible database format (mono-lingual)



A lexicon is Not a thesaurus !!

Alphabetic filter queries Exact substring queries Fuzzy queries

Language choice →

Lexicon entry term →

Translations of the entry →

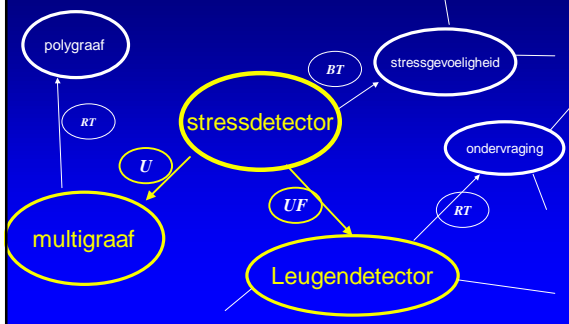
BIKIT

A thesaurus in MS.Word is Not a thesaurus !!

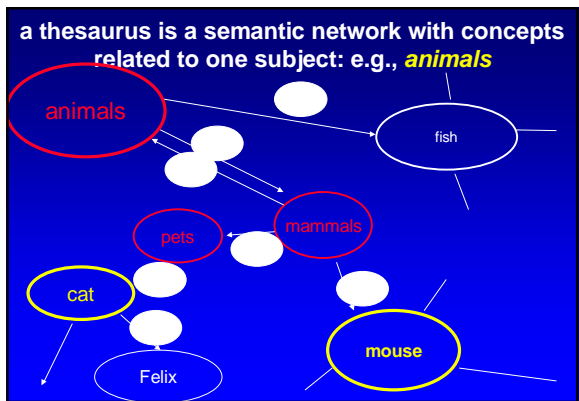
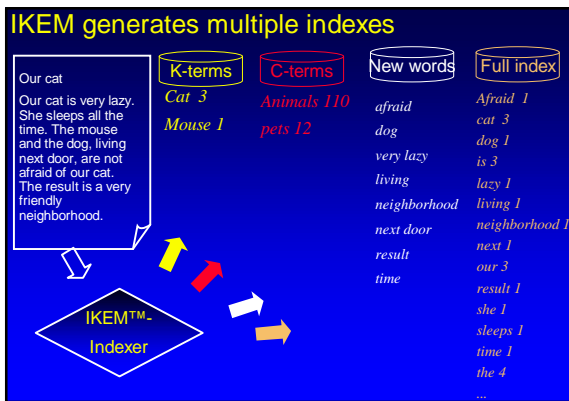
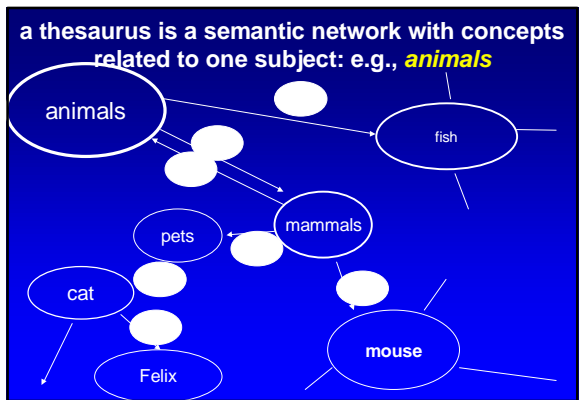
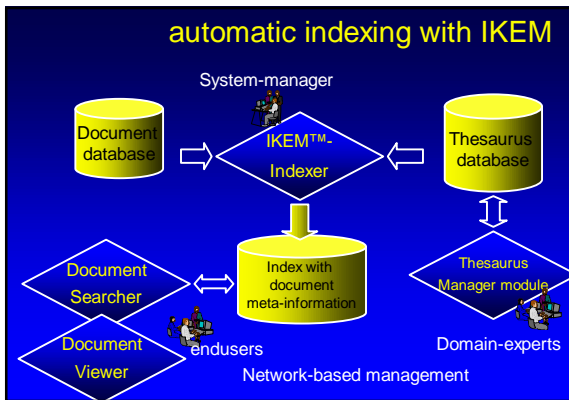
BIKIT

Babbage Institute for Knowledge- and Information Technologie

Relations between terms can be user-dependent:
vb. Wat is de voorkeurterm ?



IKEM:
Example of thesaurus-based indexing toolkit



Babbage Institute for Knowledge- and Information Technologic

What is *linking* about ?
 D. Vervenne
 This document describes the facilities to link documents that are not complete. Since we consider electronic documents as **collections** of conceptual **subdocuments**, it is relevant to process all subdocuments in appropriate ways. We therefore need facilities that address parsing **dependencies** that can be interpreted as **links**.

Weighting the concept distances that are used to represent the implicit knowledge from a document (pat. pend.)
 Copyright '97

Searching documents ?

The thesaurus supports the **Document Searcher**

Callouts: Hierarchic browsing, Search results, All thesaurus-terms, Full text search, Keyword search, Concept search, Update search

The document viewer

supports the XML-based document-structure recognition

Structure of document

Thesaurus construction And Maintenance models

Creation and updates of thesauri

Options available in IKEM™-TDocs

- 1) Import of hierarchic relations {BT,NT,U,UF,RT} using a tabbed ascii text format
- 2) Import of all relations {BT,NT,U,UF,RT} and Scope Note using a table format (like Excel)
- 3) (r) Import of hierarchic relations {BT,NT} from a tagged document collection (Autor's tagging)
 * this method is currently under development

BIKIT Copyright '97

Method #1 : error checking

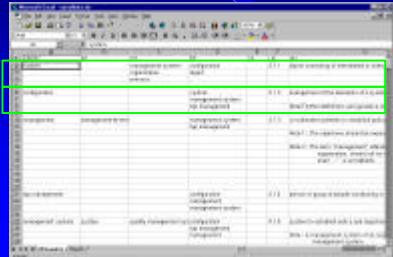
BIKIT

Babbage Institute for Knowledge- and Information Technologie

Method #2 : filling in a relations table

- A table is used with the following column headings: | ENTRY | BT | NT | RT | UF | LB | SN |
- Relations are entered in table rows, in this case in Excel

relations of "system"
relations of "configuration"
...



BIKIT

Method #3 : using tagged documents

- Step 1 Author's classification terms are attached to documents as tags (preferably in XML format)
- Step 2 Relevant expressions are extracted from document contents and saved in the update table
- Step 3 The update table is processed and relevant terms are attached to the authors classification terms in the thesaurus

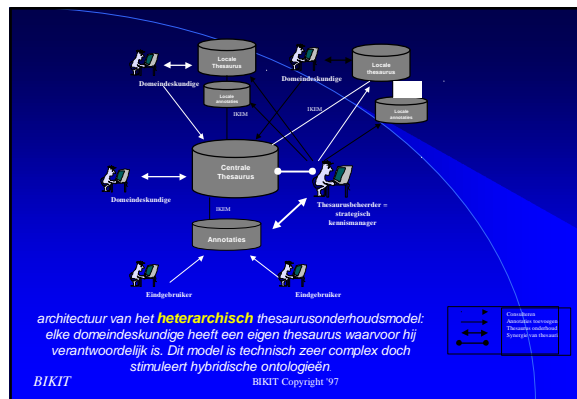
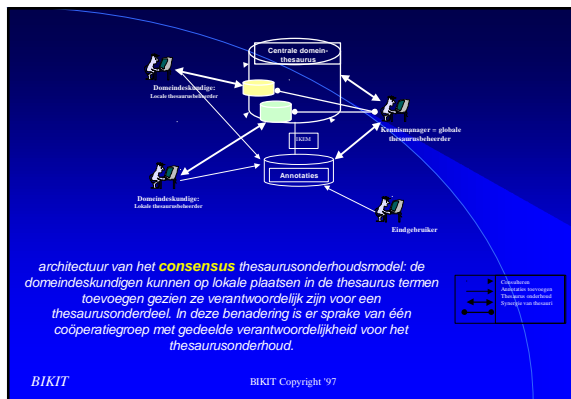
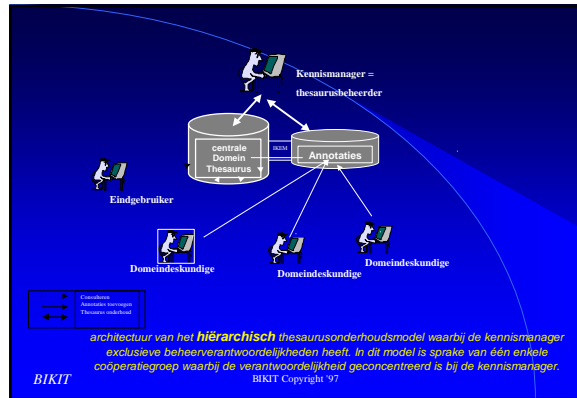
BIKIT

BIKIT Copyright '97

How to maintain a thesaurus ?

3 workflow models:

- hierarchic workflow
- consensus workflow
- heterarchic workflow



Babbage Institute for Knowledge- and Information Technologe

Conclusion ?

Indexing is a strategic process with important implications for many document management processes:

new trend is to start the indexing process during the creation process and to integrate with content managem.

