

# *Proteus:*

**a concept browsing interface  
towards conventional  
Information Retrieval Systems**

**Oreste Signore**

*CNUCE - Institute of CNR - via S. Maria, 36 - 56126 Pisa (Italy)*  
Phone: +39 (50) 593201 - E.mail: oreste@ICNUCEVM.CNUCE.CNR.IT

**Alfredo Maria Garibaldi\* - Maurizio Greco\***

\*Presently: *TECSIEL - via Enriques, 28 - 00146 Roma*  
Phone: +39 (6) 505091

**3rd International Conference on  
Database and Expert Systems Applications**

**DEXA'92**

**September 2-4, 1992**

*Valencia*

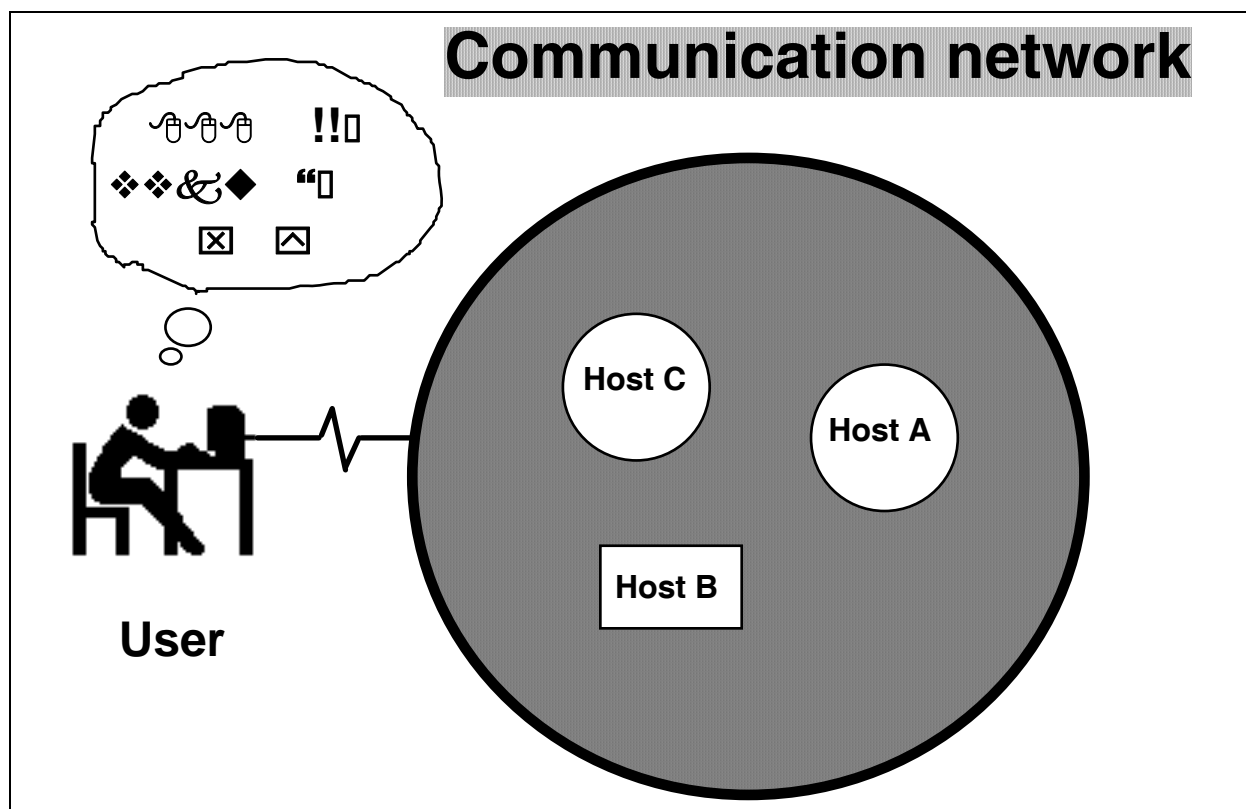
# Contents

- **The access to the on-line information**
- **The user interface problem**
- **Proteus**
  - **The architecture**
  - **The “conceptual document”**
  - **The formulation of the query**
  - **The browsing on the thesauri (the graphical display)**
  - **The query**
  - **A sample**
- **Conclusions and future work**

# The myth of the on-line information

- O** I intend to have access to an on-line data bank...
- m** No problem, *all you need is a modem and a NUI*  
That's all! Just connect to the data bank of interest.

*... and so ...*



- p** Data banks reside on hosts, and are managed by  
*Information Retrieval Systems*

# **Accessing on-line data bank services: some drawbacks**

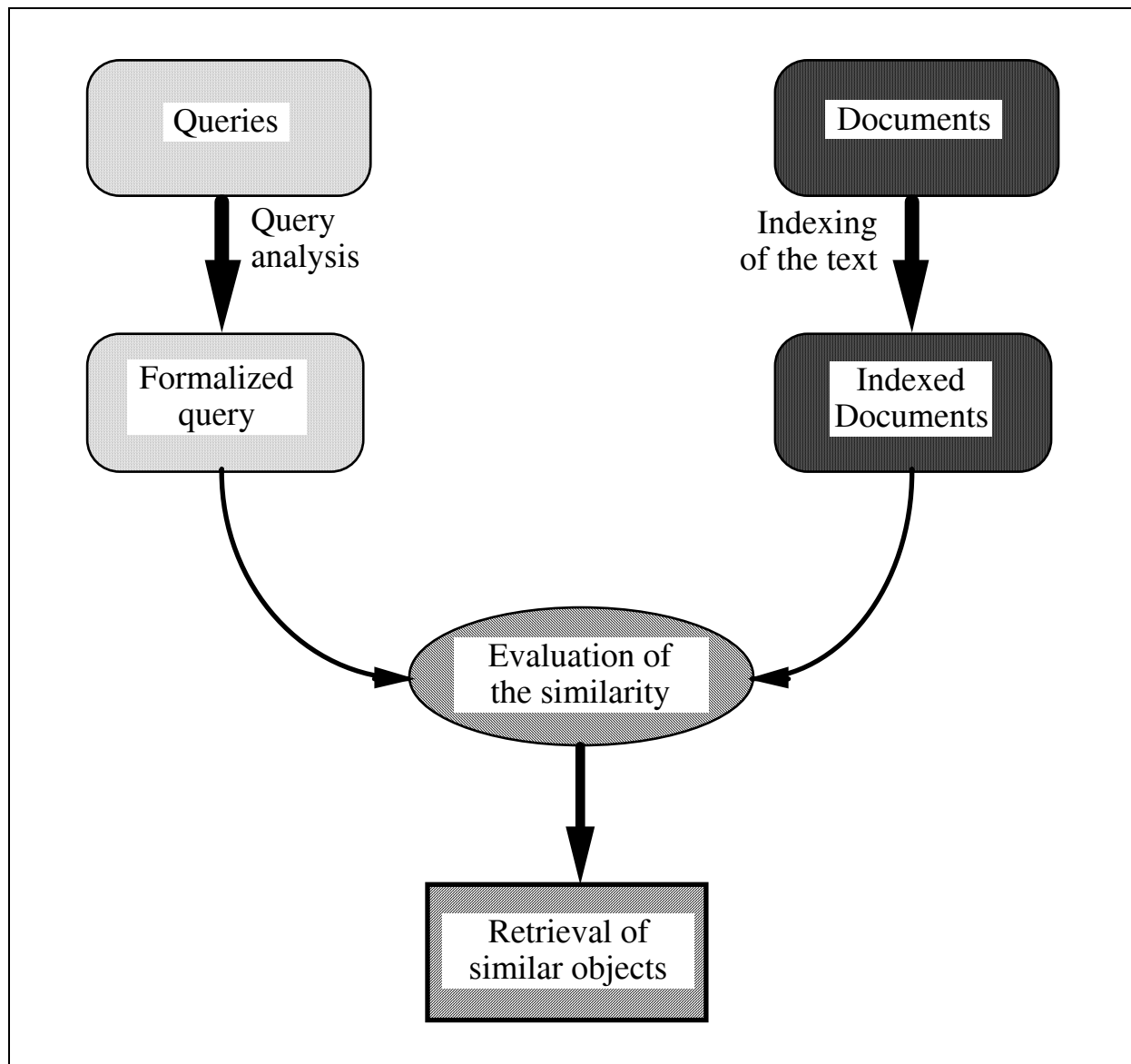
**The growth of the information market has been limited by several factors:**

- **technological problems:**
  - reliability and speed of the transmission networks
  - TTY based dialog
  - accounting problems
  - high costs (transmission and access)
  
- **User interface problems:**
  - complex login procedures
  - differences between query languages
  - data bank structure
  - contents of the fields

**Some solutions proposed at EEC level:**

- **Common Command Language (CCL, also proposed as an ISO standard)**
- **special purpose workstations**

# Functional description of an IRS



p **Stored documents are represented by a set of:**

**index terms**

# The user interface

p In the Information retrieval approach it is necessary to:

- identify the content of the documents **to be stored**
- identify the content of the user's query
- **begin a search based on a specific topic**

p Users need support to:

- define or refine **the topics which interest them**
- **improve their knowledge of the subject**
- **clearly express the concepts to be found out**
- refine or enlarge the concept **they were starting from**

## ?! A menu based interface?

<b>Enter the required values</b>		
<b>Author</b>	=	_____ OR _____ OR _____
<b>Year</b>	=	_____
<b>Argument</b>	=	_____ OR _____ OR _____

**No, thanks!**

*(I am able to understand the boolean logic, what I'm looking for is the content and the organisation of the data bank)*

# Thesauri and concept browsing interface

- p Typical problems:
  - lack of precision of the utilized terms
  - incorrect identification of the element to be found
  
- p An approach based upon knowledge:
  - assures a consistency in the representation between indexers and users
  - provides a semantic system of concepts accessible to indexers and users
  - creates a mechanism of communication that fits the key concepts.
  
- p The user may be supported by the:  
*explicit representation of the relationships of synonymy, preference and hierarchy between terms.*  
That is building a

**thesaurus**

- p ***A concept browsing interface:***
  - allows the user to identify the descriptors
  - creates a communication mechanism to fit the key concepts
  - takes advantage from the availability of a high quality thesaurus.
  
- Some prototypes:
  - CANSEARCH (Pollitt, 1983)
  - CALIBAN (Frei and Jauslin, 1983)
  - COALSORT (Monarch and Carbonnell, 1987)

# Proteus

## p **Aims:**

- easier interaction with conventional systems
- interaction with a “conceptual document”, independent from the physical implementation
- formulation of the query in ISO standard language
- tools for the identification of the relevant concepts

[

**more effective queries**

## p **Tools:**

- **Dictionary manager**
- **Graphic browser on structured thesauri**  
*(graphical representation of the concepts and of their hierarchy, synonymy and preference relationships)*
- **Management of special fields**
- **Windows for the insertion of the query values**

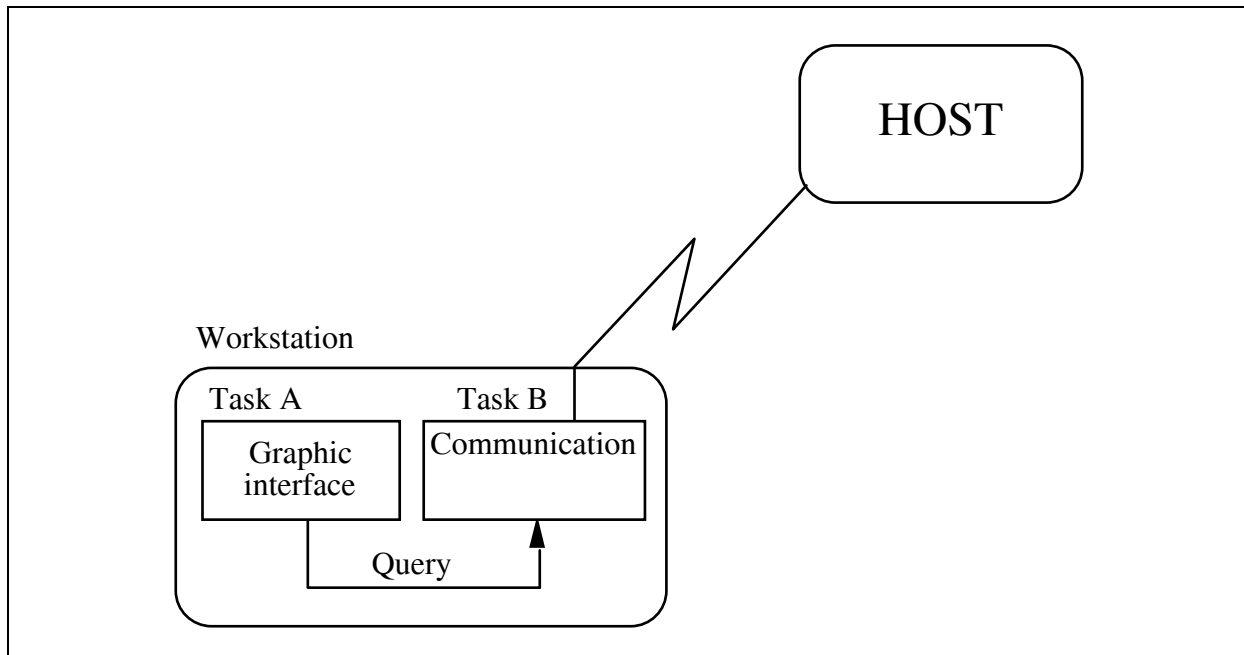
# General architecture

p **The basic idea:**

- prepare the query in the ISO standard syntax
- translate it into the query language of the system available on the host the user is connected to

[

the interface and query formulation support functions are separated from the data bank query phase

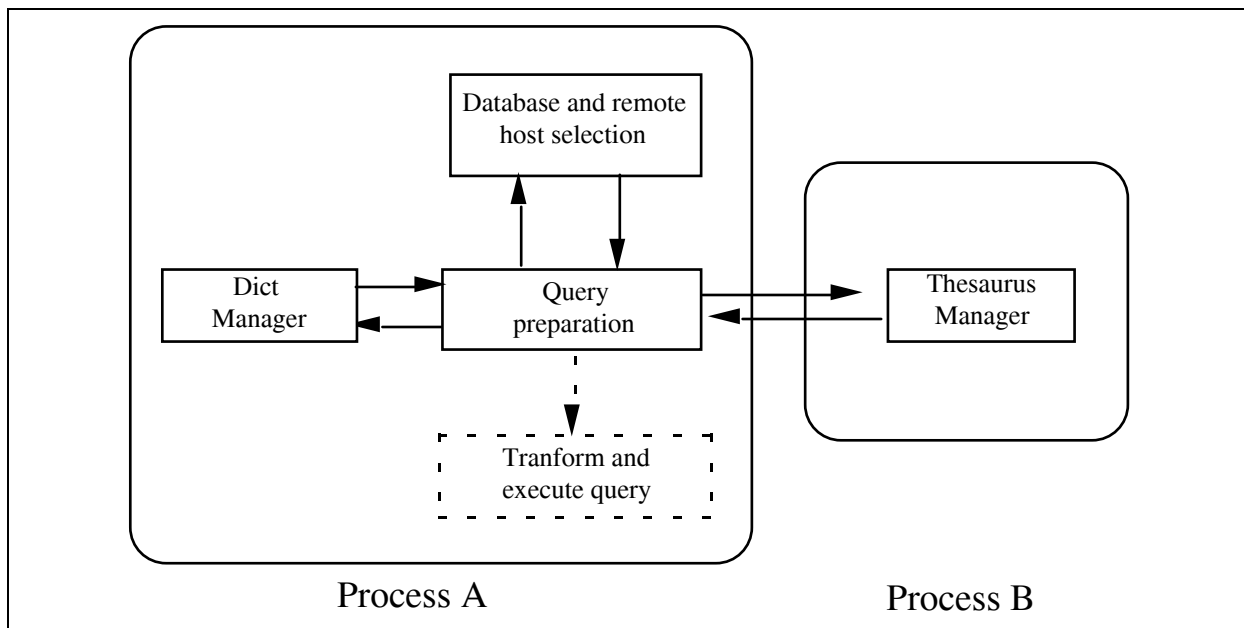


**Two tasks**

- **Adaptable to every Information Retrieval System**
- **Available on widespread personal computers**  
(PS/2, running OS/2)

# The software architecture

- Two different programs (A and B)
- Their executing instances are the process A and the process B.
- For each thesaurus a thesaurus manager has been coded, but they have the same basic structure.



## *Two processes*

(A can activate B and both can continue in parallel)

- Multitasking implemented both at the process level and at the thread level.  
Each process is made by two execution threads:
  - a *first thread* receives the user input and passes the time consuming tasks to the second thread;
  - a *second thread* executes the time consuming tasks (as a series of queries to the support database) avoiding locking the system waiting for the completion of them.
- Processes are synchronized by semaphores, threads are synchronized by an exchange of messages.

# The “conceptual document”

Document fields

DB Name ArtObjects

Author  
Object  
Subject  
State of conservation  
Location  
Historical notes

End Selection

Exec

Cancel

## p Four different types of field:

**Text** *free text field*

**Cod** *the field takes values that are a coded representation of the perceived values*  
(splitting in several subfields, codified values, etc.)

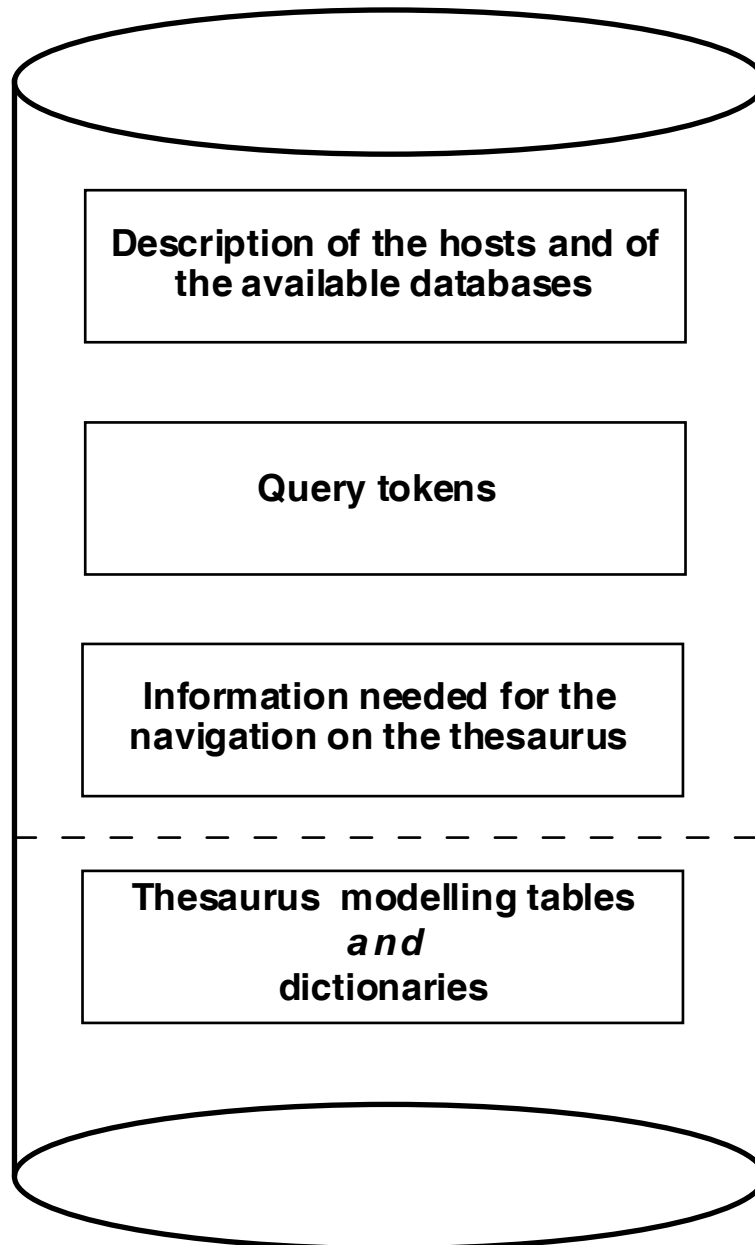
**Dict** *the field can take values from a definite list of values*

**Thes** *the field values are semantically arranged according to their relationships of preference, equivalence or hierarchy*

## p Others (picture, sound, etc.) can be added

## The supporting database

- p **Three groups of tables (absolutely general)**



- p **The thesaurus modelling tables represent the semantic relationships between the terms**

# The query formulation phase

p The user can get a description of any document field

p Field values can be entered in the appropriate window

The diagram shows a window titled "Query box". Inside the window, the label "Author" is positioned to the left of a rectangular text input field. Below the input field, there are three buttons: "OK", "Save and continue", and "Cancel".

p Free text field values are accepted as entered

p Coded field values are translated by a field dependent module (*user coded*)

p Dictionary controlled fields are automatically checked or the values are chosen from a multiple selection window

p Thesaurus controlled fields are automatically checked.  
If the user needs to refine the concept he/she is interested to, a graphic browsing on the thesaurus can take place.

## Two sample cases

p **ICONCLASS** (Iconographic Classification System):

- a classification system set up by the Leiden University (Holland) between 1944 and 1984
- used for the description of the iconography of the western art (about 30.000 subjects)
- keywords are attached to the different subjects, and the code supplies a good linguistic independence

*ICONCLASS is tree structured*

p **GEODOC**:

- a multilingual thesaurus agreed by the "Commission on Geological Documentation" (COGEODOC) of the IUGS (International Union of Geological Sciences) as a common terminological thesaurus for Earth Sciences, to be adopted by different documentation centres.
- More than 5.000 terms.

*GEODOC is a multi-tree thesaurus*

## ICONCLASS: a typical hierarchy

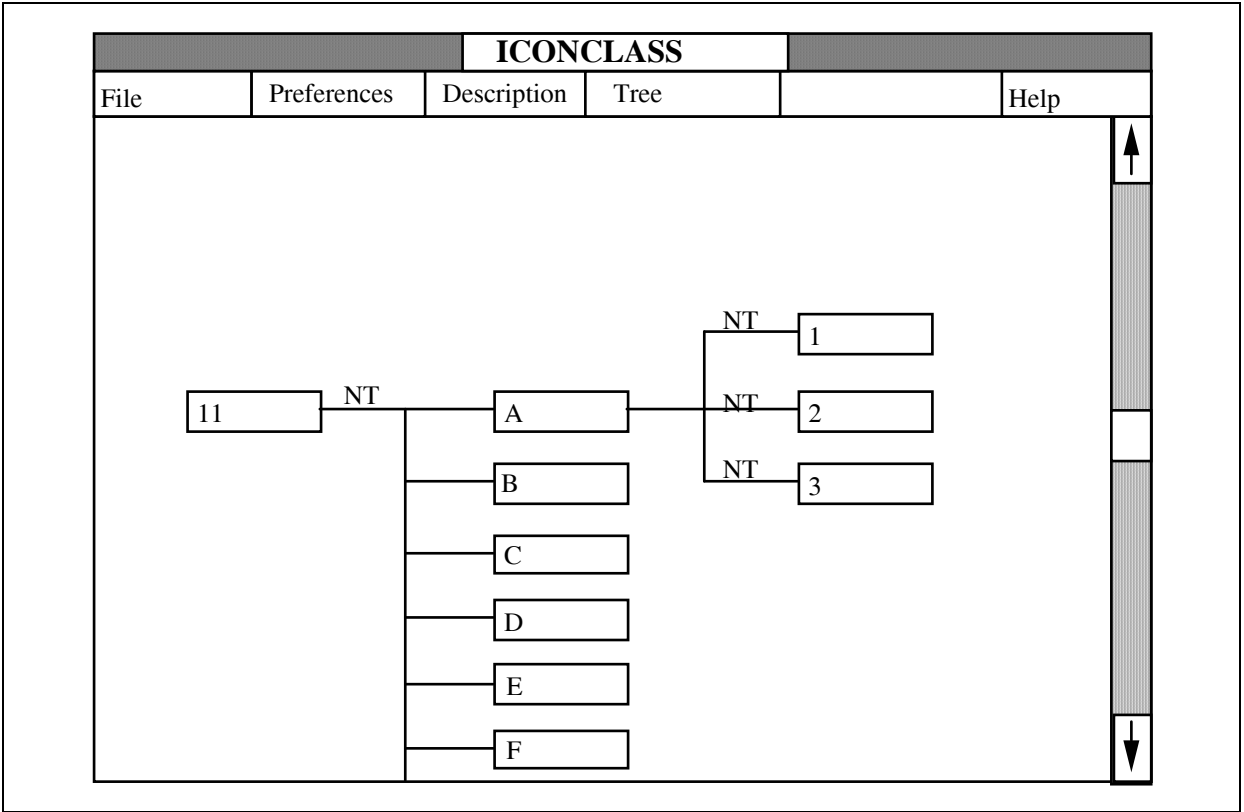
<b>Code</b>	<b>Subject</b>
<b>1</b>	<b>Religion and magic</b>
<b>10</b>	<b>(Symbolic) representations in relation to creation, cosmos, cosmogony, universe, and life (in the broadest sense)</b>
<b>11</b>	<b>Christian religion</b>
<b>11 A</b>	<b>Deity, God (in general) in relation to Christian religion</b>
<b>11 A 1</b>	<b>God the creator</b>
<b>11 A 11</b>	<b>God measuring the Universe (with compasses)</b>
<b>11 A 2</b>	<b>Divine Nature</b>
<b>11 A 21</b>	<b>Divinity</b>
<b>11 A 22</b>	<b>Symbols in relation to Divine Nature</b>
<b>11 A 22 1</b>	<b>Cyrcele symbolizing God's perfectness</b>
<b>11 A 23</b>	<b>God's perfections</b>
<b>11 A 3</b>	<b>God's wrath</b>
<b>11 A 31</b>	<b>'Flagello di Dio'</b>
<b>11 B</b>	<b>Father, Son and Holy Gost in relation to Trinity</b>
...	...
<b>11 C</b>	<b>God the Father</b>
..	...
<b>11 D</b>	<b>Christ</b>
...	...
<b>11 E</b>	<b>The Holy Gost</b>
...	...
<b>11 F</b>	<b>The Virgin Mary</b>
...	...

<b>11 G</b>	<b>The Angels</b>
...	...

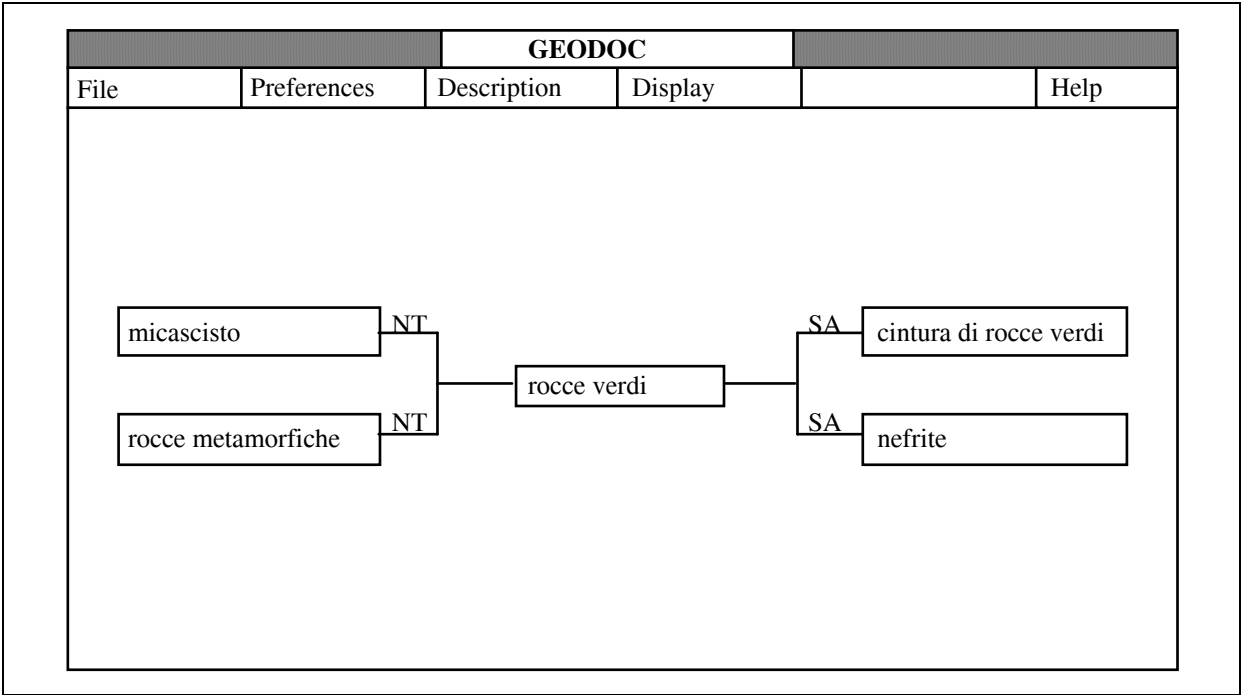
# The graphic display of the thesauri

- p The Thesaurus Manager exhibits a different behaviour depending on the specific thesaurus
  
- p The basic functionalities are the same:
  - the thesaurus is represented as a graph  
*(the nodes are the thesaurus terms, the arcs are the connecting relationships)*
  - the arcs can be labelled with the names of the relators
  - "point and click" interaction style
  - three levels are shown
  
- p Multi-tree thesauri are displayed as a "butterfly display":
  - the central level is always constituted by a single node
  - on the right the domain
  - on the left the codomain
  
- p The user:
  - identifies the kind of relationships he/she is interested in
  - selects a term (either directly, or from a list box)
  - gets the display of a neighbour of the term  
*(only the selected relationships are shown)*
  - can use the various options available from the menu bar
  - can set some default ("*double click*") actions
  - can move around the structure  
*(and ask for detailed explanation of the terms)*
  - can extend the tree towards the root or towards the leaves.
  
- p Any term can be selected in order to build a query to the database: all the terms which are selected in a single navigation will be OR'ed.
  
- p The colours emphasize the role of the nodes

# The two different displays



**The graphic display of the neighbouring of the code '11 A'**



**A “butterfly display” for the GEODOC thesaurus (italian version)**

# The query

- p The tokens which constitute the query are stored in table.  
The query:

(Author : Michelangelo OR Caravaggio)  
AND  
(Subject : '11 A' OR '11 D 2' OR '11 F')

is stored as:

Field	Value
Author	'Michelangelo' OR 'Caravaggio'
Subject	'11 A' OR '11 D 2' OR '11 F'

- p The resulting query may be edited  
(presently only the boolean logic can be modified)

Query :

1 AND 2

1 Author: Michelangelo OR Caravaggio

2 Subject : '11 A' OR '11 D 2' OR '11 F'

Exec Cancel Back to selection

# Conclusions

- p Document retrieval requests the exact knowledge of concepts to be found
  
- p **Proteus makes easier to formulate queries on a document:**
  - conceptual structure of the document
  - automatic check against the available dictionaries and thesauri.
  
- p The thesaurus manager gives a snapshot of the structure of the thesaurus in a neighbour of the term of interest, and the user is allowed to navigate in the structure. The navigation on the thesaurus may be accomplished *in parallel* with the formulation of the query on other fields of the document.
  
- p The interface is totally independent from the databases, only the thesaurus manager has an intrinsic dependence on the thesaurus structure.
  
- p The interface agrees with the SAA standards
  
- p While the novice user can receive help in the formulation of the query, the expert user can use some shortcuts in order to speed up the process.

# Future work

## Some extensions are under study or development

- p **modification of the tokens of the query**
- p **implementation of the connection with the remote host according to the CPI standards, using APPC**
- p **the query will be saved for future use**  
*(establishment of a user profile)*
- p **multitasking navigation on several thesauri, using appropriate semaphores**  
*(but this enhancement seems not to be really useful for the user, who should follow several classification schemas in parallel)*
- p **use of the Thesaurus manager as a tool for the manipulation of structured thesauri**  
*(addition of terms, translation, moving of a set of terms, establishment of new relationships, etc.)*
- p **parametric representation of the thesaurus structure**
  - h **implementation of a single thesaurus manager, which will produce a display that will vary according to the structure of the thesaurus it is processing**
- p **implementation in a “client-server” architecture:**
  - h **a server will host the support database, and will implement the connection with an external host.**