

#1. Find concepts in scenario

The assistant began showing the scenario to Dennis. It was suggested the use of **Digital Library** concept to begin searching.

#2. Disambiguate given a target (concept)

The assistant showed the following options according to wikipedia, and ask him to choose only one option.

```
Wikidata item search
Number of results:      10

Results:
1      digital library (Q212805) - library in which collections are stored in electronic media formats and accessible via computers
2      Digital Library for Dutch Literature (Q2451336) - website about Dutch language and Dutch literature
3      Digital Library Federation (Q5275906) - organization
4      Digital Library of Mathematical Functions (Q24534) - online project meant to be a collection of special functions
5      Digital Library of Slovenia (Q3435281) - digital library supported by National and University library of Slovenia
6      Digital Library Perspectives (Q53952043)
7      Digital library of the Lombardy (Q28531719)
8      Digital Library from the Meiji Era (Q11638378)
9      Digital Library of India (Q56480935)
10     Digital Library of the History of Friesland (Q2133797) - digital library
```

figure 2. different contexts where digital library is used.

Dennis chose option 1.

The assistant processed option1 text in wikifier.org

Wikifier results

library in which collections are stored in electronic media formats and accessible via computers

figure 3. finding concepts for the option selected

#3. Process a paragraph-article from wikipedia The assistant took the main article related to “digital library” in wikipedia and looked for its concepts in wikifier.org

Wikifier results

A digital library, digital repository, or digital collection, is an online database of digital objects that can include text, still images, audio, video, or other digital media formats. Objects can consist of digitized content like print or photographs, as well as originally produced digital content like word processor files or social media posts. In addition to storing content, digital libraries provide means for organizing, searching, and retrieving the content contained in the collection. Digital libraries can vary immensely in size and scope, and can be maintained by individuals or organizations.[1] The digital content may be stored locally, or accessed remotely via computer networks. These information retrieval systems are able to exchange information with each other through interoperability and sustainability. [2]

figure 4. finding concepts in a related wikipedia article

#4. Show concepts according elicitor selections

The assistant showed to Dennis the concepts collected (from fig.3 and fig4.)

library, electronic media, computers, digital library, digital repository, digital, online, database, still images, digital media, digital content, word processor, social media, digital libraries, digital content, computer networks, information retrieval, systems, information, interoperability, sustainability

The assistant asked Dennis to select 5 keywords

Dennis selected: **Information, digital content, database, online, digital libraries**

#5. Create a corpus using keywords selected by elicitor

The assistant asked to create a query using “digital library” in the query:

Dennis chose: "digital library" + "digital content"

For the query, the results in Github were only 2 projects. The assistant suggested to improve his query.

Dennis chose: "digital library" + information

GitHub results were 11. The assistant asked to improve the query one more time:

Dennis chose: "digital library" + database

GitHub results were 20. The assistant asked to choose one of his queries based on results.

Dennis chose **"digital library" + database**

The assistant created the corpus in <http://corpus-retrieval.herokuapp.com/> for the query

The tool returned 388 projects

#6. Find keywords in Github

The assistant looked the context of 1st keyword: "digital library" in the corpus.

It was found 417 texts

the assistant looked the texts with the 2nd keyword: “database”

it was found 57 texts / 417

after pruning duplicated texts, assistant got 56 unique texts with both keywords

#7. Processing texts to find GitHub keywords (top30)

"invenio" "framework" "project" "system" "patron" "redshift" "access" "californica"
"clemson" "contentdm" "mysql" "initiative" "model" "admin" "search" "columbia"
"repository" "management" "journal" "book" "software" "user" "research" "information"
"omeka" "download" "knowledge" "queue" "blaster" "modem"

The assistant showed the concepts saying that apparently these are words well related (in order of appearance) to the concept of Digital Library and Database

The assistant asked Dennis to select 5 keywords:

Dennis asked: **framework, system, access, repository, information**

#8. Create a new query

The assistant showed the selections made by Dennis and asked to create a final query

1g: "digital libraries", "information", "digital content", "database", "online"

2g: framework, system , access, repository, information

Dennis created: "digital library" repository

#9. the ten first projects from the corpus created

1. <https://github.com/vphill/pyoaiharvester>
2. <https://github.com/snkim/AutomaticKeyphraseExtraction>
3. <https://github.com/PerseusDL/canonical-greekLit>
4. <https://github.com/PerseusDL/canonical>
5. <https://github.com/SeerLabs/CiteSeerX>
6. https://github.com/PerseusDL/treebank_data
7. <https://github.com/fcrepo4/fcrepo4>
8. <https://github.com/code4lib/ruby-oai>
9. <https://github.com/PerseusDL/lexica>
10. https://github.com/PerseusDL/catalog_data