# Summarizing Biological Literature with BioSumm

## Elena Baralis, Alessandro Fiori
### Politecnico di Torino

The availability of increasingly wider text repositories requires effective techniques to manage the huge mass of unstructured information there contained (e.g., navigate, analyze and represent it in the most suitable way). Particularly, in the biological and biomedical domain a huge amount of information is daily generated and contributed by a vast research community spread all over the world. Repositories like PubMed Central, the U.S. National Institutes of Health (NIH) free digital archive of biomedical and life sciences journal literature, nowadays contain billions of documents.

## System Architecture

✔ Fully modular
✔ Allows integrating plugins addressed to a specific task (e.g., clustering, web search).
✔ Change Domain dictionary to tailor the grading function according to the application domain

**Online search & local repository.** Keyword search of scientific papers on:
- *Google Scholar,*
- *PubMed Central (PMC)*
- *PubMed*

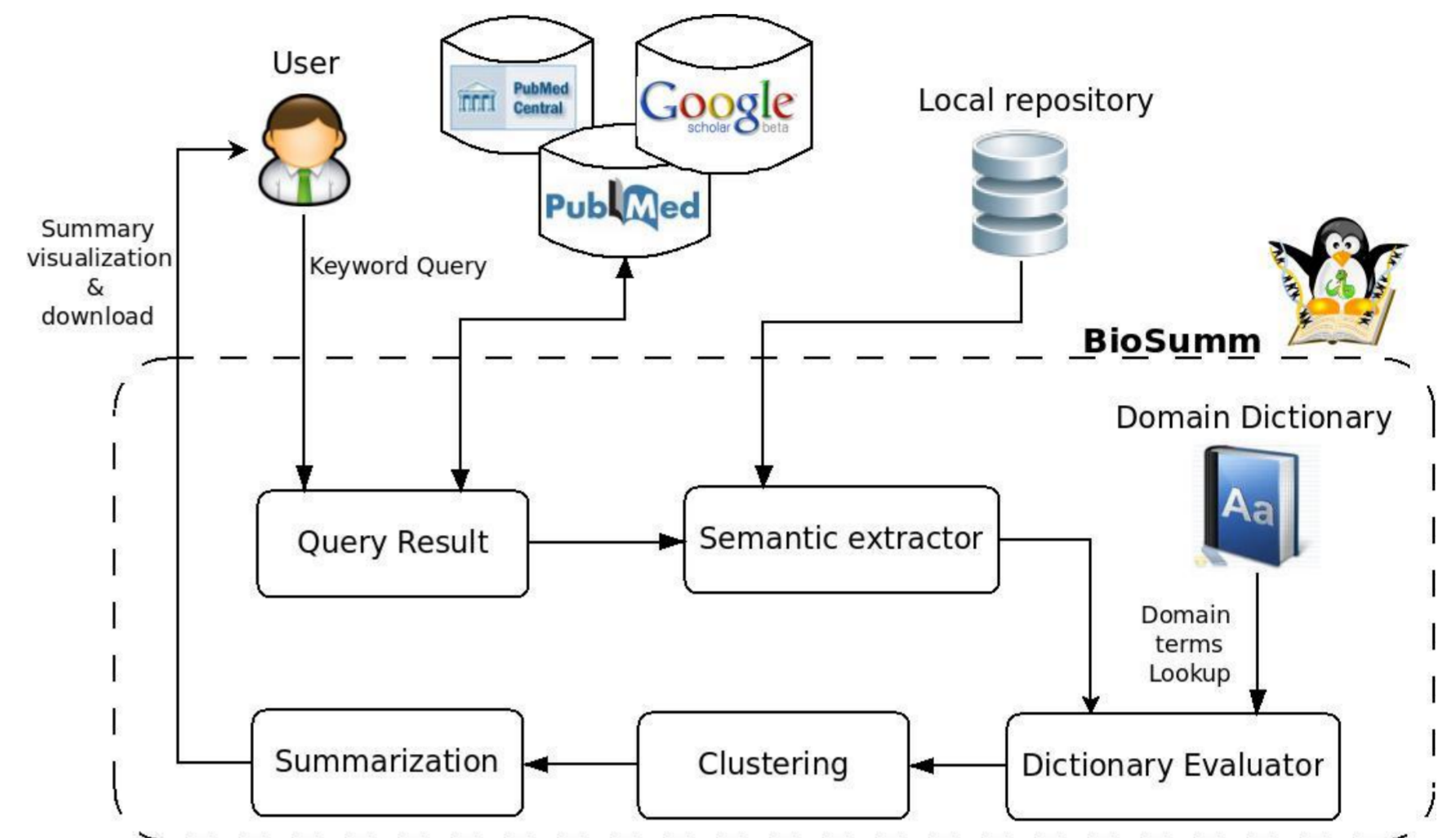Or load locally stored documents in pdf and xml formats.

**Document structure extractor.** Parsing and extraction of available components (e.g., title, authors, journal, abstract, body, keywords).

**Clustering.** Reduce the heterogeneity of the retrieved documents.
➤ Document collection represented as a matrix tf-idf
➤ *Bisecting K-means* clustering method is performed

**Dictionary evaluator.** Compute the semantic weights of grading function according to the terms in the domain dictionary. The Domain dictionary contains genes and proteins names from *BioGrid* database.

**Summarization.** Based on a traditional statistic summarizer, it biases sentence selection using the information contained in the *Domain Dictionary.*
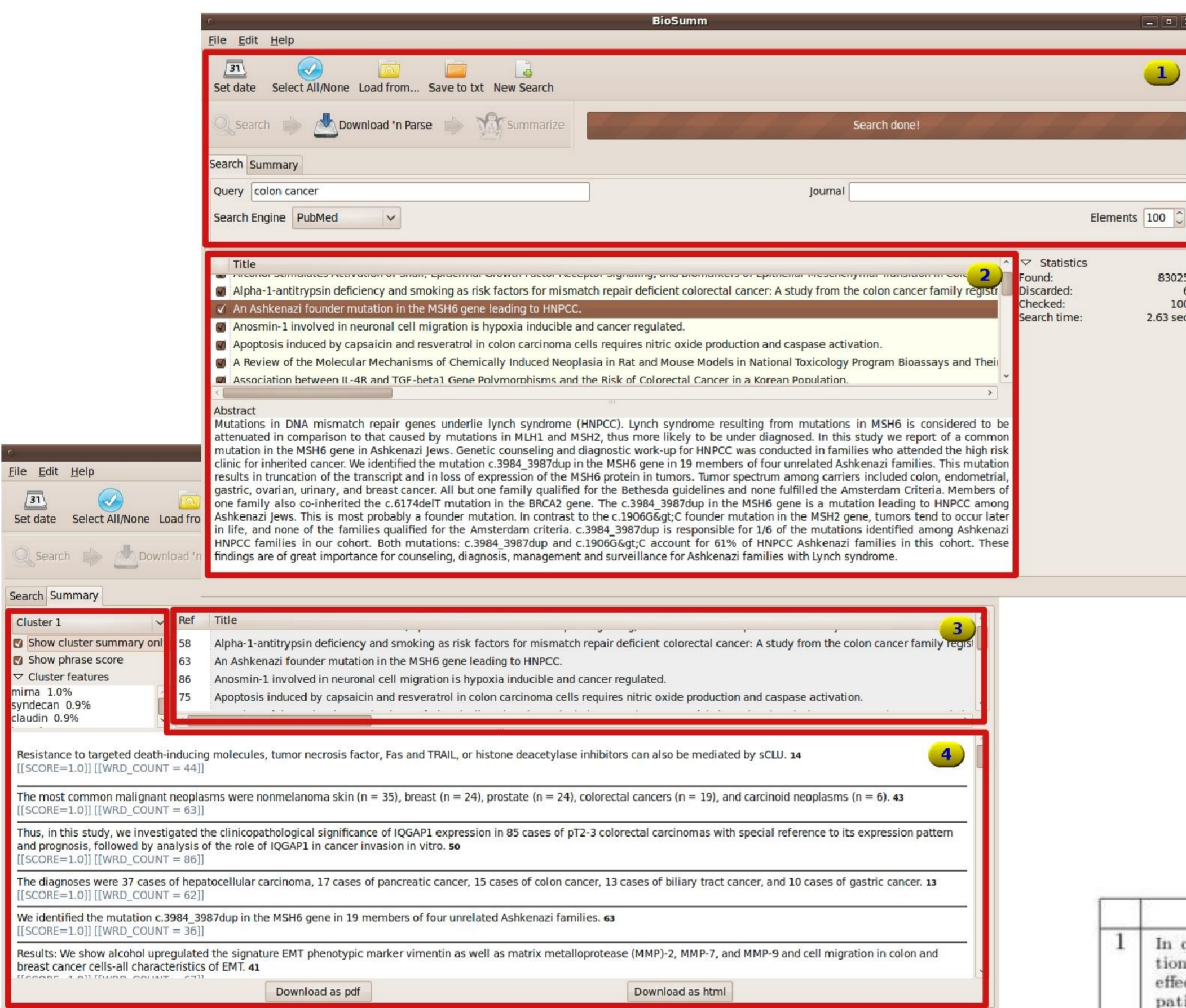


## Experimental results

BioSumm is neither a traditional summarizer nor a extractor of dictionary terms. It is designed to be a summarizer oriented to the biological domain. Thus, its summaries have both the expressive power of the traditional summaries and the domain specificity of documents produced by a dictionary entry extractor.

The difference with a traditional summarizer may be appreciated in the example. Six most graded sentences in BioSumm and in a traditional summary are reported. The results were produced by the experiments carried on the scientific journals freely available in *PubMed Central*. Specifically, they contain sentences belonging to a cluster of documents belonging to the Breast Cancer journal. The keywords of the cluster (the words describing its major topics) are *proband, Ashkenazi, Jewish* .

The comparison shows that BioSumm, although oriented on *biology*, is still able to cover all the major topics covered by a traditional summarizer. Moreover, its sentences are less generic and contains a lot of genes and proteins which are described in details and not only listed.

The results suggest that researchers that discover gene correlations by means of analysis tools (e.g., data mining tools) may exploit this framework to effectively support the biological validation of their results. Experimental results obtained by means of *ROUGE* are also reported.



## Graphical User Interface

**① Document search.** The user can set the parameters (e.g., keywords, journal, date of publication) to retrieve the documents from supported digital libraries.

**② Document browsing.** Management of retrieved documents to select the most relevant for summarization task. Documents may be individually excluded from the document collection on which both clustering and summarization are performed, thus further refining the analysis.

**③ Documents of cluster.** List of the documents belonging to a cluster identified by the clustering block. BibTex description is also available.

**④ Cluster summary.** Each sentence in the summary is scored by relevance with respect to the considered application domain. Furthermore, the cluster is described by means of the relevant keywords appearing in its documents.



Experimental comparison between BioSumm and general purpose summarizer (OTS)

| Dataset | BioSumm | | | OTS | | |
|---|---|---|---|---|---|---|
| | Precision | Recall | F-measure | Precision | Recall | F-measure |
| Breast Cancer | **0.08246** | **0.22553** | **0.11456** | 0.08026 | 0.21860 | 0.11141 |
| Arthritis Res | **0.09089** | **0.25362** | **0.12596** | 0.08844 | 0.24406 | 0.12197 |

ROUGE-2 evaluation on two different document collections

| Dataset | BioSumm | | | OTS | | |
|---|---|---|---|---|---|---|
| | Precision | Recall | F-measure | Precision | Recall | F-measure |
| Breast Cancer | **0.10038** | **0.28175** | **0.14053** | 0.09872 | 0.27599 | 0.1811 |
| Arthritis Res | **0.11095** | **0.31777** | **0.15498** | 0.10905 | 0.30888 | 0.15169 |

ROUGE-SU4 evaluation on two different document collections