# Showing Figures and Captions in the BioText Journal Search Engine

**Marti A. Hearst\*, Michael A. Wooldridge †, Jerry Ye\* and Anna Divoli\***

*School of Information, University of California, Berkeley, CA 94720, †California Digital Library, Oakland, CA 94612

## Summary

The BioText Search Engine is a freely available web-based application that provides biologists with new ways to access the scientific literature.

One novel feature is the ability to search and browse article figures and their captions.
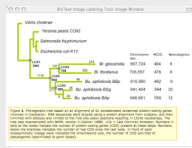
## Motivation

All prominent literature search systems are based on information from abstracts alone when:

• Researchers in the area of text-mining have started to investigate approaches for full-text analysis.

• PubMed Central has made available a large collection of full-text articles (Open Access), overcoming licensing restrictions from the publishers.

• Figure captions and figures can be especially useful for locating experimental results.
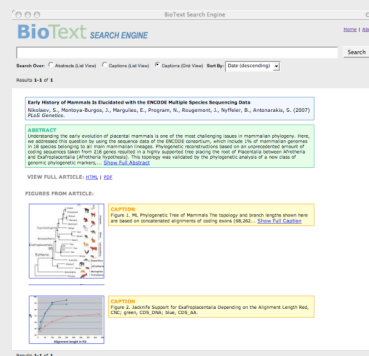
## The Views

The BioText Search Engine allows users to search in **"Abstracts - list view", "Captions - list view"** and **"Captions - grid view".** The interface is carefully designed according to usability principles and techniques.

Clicking on any figure opens a new window with a **large version** of the figure accompanied with its caption.



All views lead to the **"Endgame view"**. "Endgame view" displays a summary of a paper that an abstract or caption originates. This summary comprises of the *TITLE, CITATION, ABSTRACT*, all *FIGURES and corresponding CAPTIONS* from the paper in a list and HTML & PDF links to the paper.



## Abstracts - list view
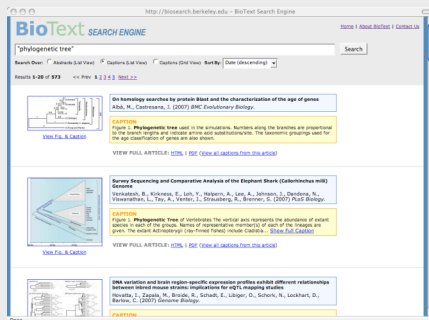
Searches in *TITLES, ABSTRACTS, AUTHOR NAMES*

Returns in a list *TITLE, CITATION, ABSTRACT*, thumbnails of *FIGURES* from the paper, HTML & PDF links to the paper and link *to the "Endgame view"*
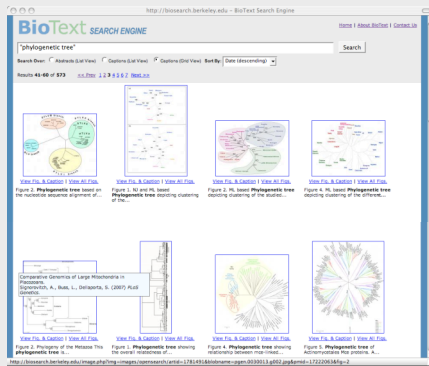


## Captions - list view

Searches in *CAPTIONS*

Returns in a list *TITLE, CITATION, CAPTION*, corresponding *FIGURE*, HTML & PDF links to the paper and link *to the "Endgame view"*



## Captions - grid view

Searches in *CAPTIONS*

Returns corresponding *FIGURES* in a grid, short *CAPTION* excerpt, *CITATION* in tooltip and link *to the "Endgame view"*



## Technical details

The current system indexes all **Open Access** articles available at PubMed Central. This collection consists of more than 150 journals, 20,000 articles, and 80,000 figures (new articles are downloaded and indexed daily).

The **figures** are **stored locally**, in order to be able to present thumbnails quickly.

The **Lucene** open source search engine issued to index, retrieve, and rank the text (using the **default statistical ranking**).

For tokenization, the standard analysis settings for Lucene are used: words are split at punctuation characters and hyphens, unless there is an umber in the token, and uses lower-casing, simple stemming, and stop word removal.

**Publication date** is stored as a separate field and can also be used to **sort** the results.

The interface is **web-based** and is implemented in python and PHP. Logs and other information are stored using MySQL.

Available at: **http://biosearch.berkeley.edu**

## Future work

The search engine is a work in progress. More functionality will be added over time. We plan to:

• Provide **full-text search**. Since the usability of different ranking functions for biology articles is still not well-understood, we plan to do usability testing, research how different sections should be weighted differently for different query types and investigate how best to show excerpts or summaries from full text before supporting this feature.

• Augment the caption search by indexing the **parts of the full text that refer to the caption**.

• Provide search over **table captions**.

• Incorporate topical features such as **genes/proteins** and **organisms**.

• For the grid view, we plan to provide **grouping according to categories** that are of interest to biologists, such as "sequence alignments" and "phylogenetic trees". (We are building a classifier for figures and their captions. We have developed an image annotation interface and are soliciting help with hand-labeling mated caption classifier.)

Additional future developments on the BioText Search Engine will depend on **feedback** and **requests** we receive from users, and the results of extensive **usability testing**.

### References
Marti A. Hearst, Anna Divoli, Harendra Guturu, Alex Ksikes, Preslav Nakov, Michael A. Wooldridge and Jerry Ye (2007) *"BioText Search Engine: beyond abstract search"* Bioinformatics; doi:10.1093/bioinformatics/btm301

Marti A. Hearst, Anna Divoli, Jerry Ye and Michael A. Wooldridge (2007) *"Exploring the efficacy of caption search for bioscience journal search interfaces"* ACL 2007 Workshop on BioNLP, Prague, Czech Republic

**BioText** SEARCH ENGINE

**http://biosearch.berkeley.edu/**

UC Berkeley School of Information