

# RepWeb: A Web-Based Search Tool for Repeat-Related Literatures

Taeha Woo<sup>1,2,†</sup>, Younguk Kim<sup>1,†</sup>, Jekeun Kwon<sup>1</sup> and Jungmin Seo<sup>1,2\*</sup>

<sup>1</sup>Korean BioInformation Center, Korea Research Institute of Bioscience and Biotechnology, Daejeon 305-806, Korea  
<sup>2</sup>Genome Research Center, Korea Research Institute of Bioscience and Biotechnology, Daejeon 305-806, Korea

## Abstract

Repetitive sequences such as SINE, LINE, and LTR elements form a major part of eukaryotic genomes. A literature search tool that summarizes the information contained within repeat elements would provide biologists in the field of genomics with a useful tool for analyzing genomic sequence features. We developed a java program designed to make literature access easier by using two search engines simultaneously. RepWeb is a web-based search system that provides a user friendly interface for searching the reference data and journals for information related to repeat elements by using the search engines, Google Scholar and PubMed, simultaneously. It provides an interface that displays the repeat element-related biological information, and includes useful functions such as the production of a repeat tree, clickable links to PubMed and Google Scholar, exporting, and sorting a field into date, author, journal and title.

**Availability:** RepWeb is freely available from the following web addresses: <http://www.repeatome.org>  
<http://www.repweb.org>  
<http://bioportal.kobic.re.kr:8080/RepWeb>

**Keywords:** repeat element, PubMed, Google Scholar, E-Utilities

## Introduction

Repetitive DNA sequences include simple repeats and transposon-derived interspersed repeats such as SINE, LINE, and LTR transposons (Jurka *et al.*, 2000). These sequences comprise a large part of the genomes of higher organisms. Despite the value of repetitive sequences, their

role in many organisms is largely unknown, and the search tools which allow users to explore literature related to repeat terms are not well developed.

PubMed is one of the most important and powerful bibliographical information resources for biologists, and includes links to full text articles and other related resources (<http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?DB=pubmed>). However, PubMed is not optimized to search for specific interests such as repeat elements.

The search engine, Google Scholar, provides a simple way to search for scholarly articles on a broad range of topics, within a large database. The Google Scholar service is based primarily on the full-text archives of some of the largest scholarly publishers, combined with some open access A/I databases, preprint and reprint repositories, and pages from presumably academic WEB sites. It provides not only a link to the cited and citing articles, but also calculates and reports the citation scores of those articles (Jacso, 2005). Its purpose is to locate scholarly literature across all disciplines in many formats and to offer the best scholarly search experience for users (Dean, 2005).

Repbase is a database of prototypic sequences of repetitive DNA from different eukaryotic species. Repeat library information can be acquired from Repbase, which is one of several databases for repeat genomic elements. However, Repbase does not provide information pertaining to journal articles that contain repeat terminology.

The Entrez Programming Utilities (E-Utilities) are comprised of seven server-side programs that are enlisted to provide a stable interface for the Entrez query and database utility at the National Center for Biotechnology Information (NCBI). The E-Utilities are based on a preset URL syntax designed to translate a set of input parameters into appropriate values that can be used with NCBI software to search for and recover the requested information. The E-Utilities make up the ordered interface to the Entrez system, and include 23 databases covering diverse sets of biomedical data, which include gene records, nucleotide and protein sequences, three-dimensional molecular structures, and biomedical literature (<http://www.ncbi.nlm.nih.gov/books/bv.fcgi?rid=coursework.chapter.eutils>).

In this paper we propose a web-based system, RepWeb, for searching reference information and journals related to repeat elements, which provides a user friendly interface. This is the first literature search web tool specifically designed to search for repeat elements reported to date.

<sup>†</sup>Both of these authors contributed equally to the work.

\*Corresponding author: E-mail [jmseo@kribb.re.kr](mailto:jmse@kribb.re.kr),  
Tel +82-42-879-8133, Fax +82-42-879-8139  
Accepted 3 April 2007

## Features of RepWeb

For the terms and classification of repeats, text and repeat information were obtained from the RepeatMasker library of Repbase (Ver. 11.10, 11-13-2006) for each repeat element (Jurka *et al.*, 2000) (Fig. 1). Repbase is a database of prototypic sequences representing repetitive DNA from different eukaryotic species. RepWeb relies on repeat dictionary, query terms, and word filters provided by Google scholar and E-Utilities (Fig. 2). We can provide various filters (article filters, author filters, subject area filters, etc.) for these queries using E-Utilities and Google Scholar. For example, if a keyword has "Alu" in the repeat name search box, its search terms will be as follows in

PubMed : "Alu" and (repeat [tw] or "repetitive element" [tw] or "repetitive sequences" [tw] or "Repetitive Sequences, Nucleic Acid" [MH] or "Interspersed Re petitive Sequences" [MH]). In Google Scholar, they will be : intitle:alu + (repeat OR "repetitive element" OR "repetitive sequences" or "Repetitive Sequences, Nucleic Acid" OR "Interspersed Repetitive Sequences").

These results are automatically parsed and assigned to each search engine when submitted. RepWeb also has a function to allow users to download each search result text file and its link to the Repbase database in the GIRI. The Web crawler and parser are implemented by Python, and Google Scholar is used to acquire more information from the WEB. The user can search using a specific search

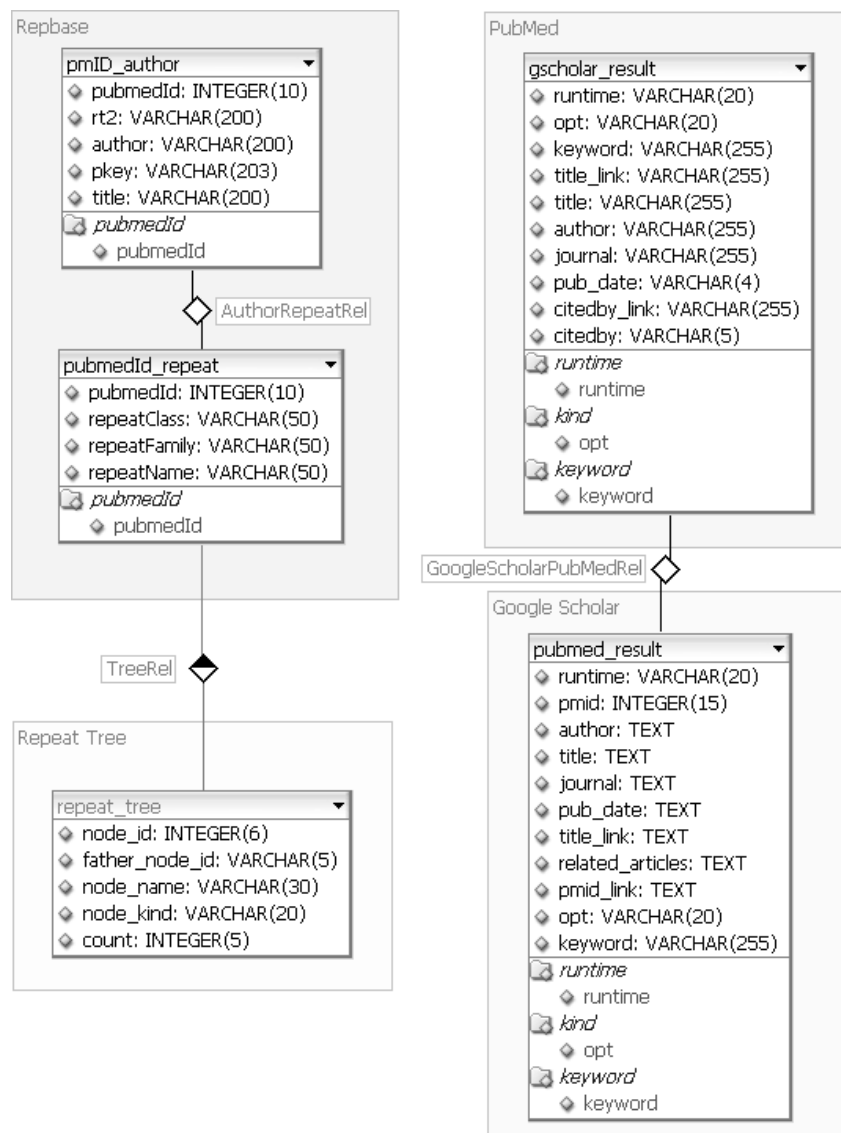


Fig. 1. Database Schema in RepWeb.

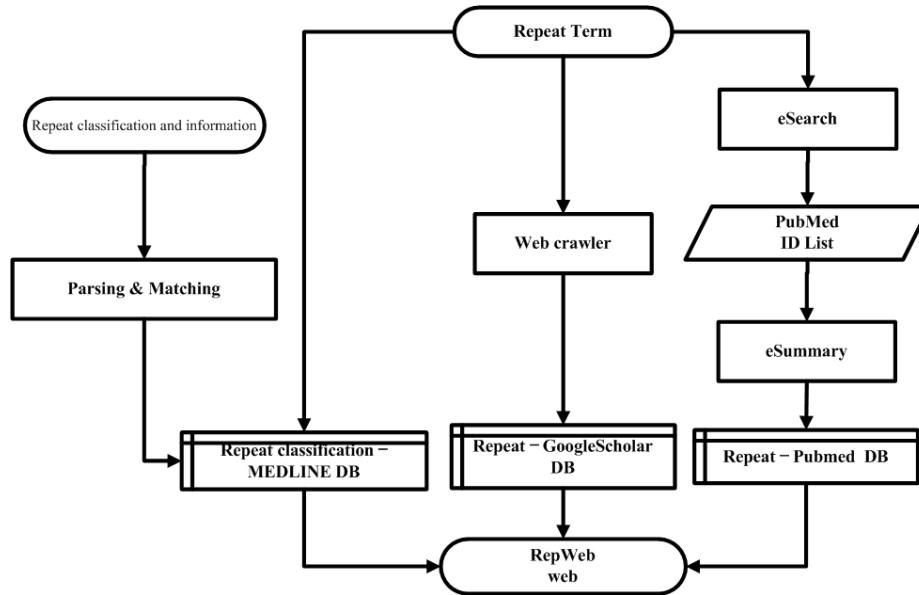


Fig. 2. Flowchart for searching a repeat term in RepWeb.

**A.**

**C.**

PHID	PubMed entry ID of searched journal
AUTHOR	Author of searched journal
REPEAT	Class, family and name about inputted keyword
JOURNAL	Journal name of searched journal
TITLE	Title of searched journal
PUBLICATION DATE	publication date of searched journal
CITED BY	In Google Scholar, cited journal information
RELATED ARTICLES	In PubMed, related journal information

STATISTICS			
SEARCH KEYWORD	Alu		
NO. OF RESULT	Repbase	Google Scholar	PubMed
	9	99	100

AUTHOR ▼▲	JOURNAL ▼▲	PUBLICATION DATE ▼▲	RELATED ARTICLES
Varzani A, Stephan W, Stepanov V, Raicu F, Cojocaru R, Roschin Y, Glavce C, Dergachev V, Spirdonova M, Schmidt HD, Weiss E	J Hum Genet	2007;52(4):309-16	Related Articles 17387576
Population history of the Dniester-Carpathians: evidence from <b>Alu</b> markers.			
Ennaffaa H, Amor MB, Yacoub-Loueslati B, El-Khil HK, Gonzalez-Perez E, Moral P, Mica-Neyra H, Elgaied A	Ann Hum Biol	2006 Sep-Dec;33(5-6):634-40	Related Articles 17381051
<b>Alu</b> polymorphisms in Jerba Island population (Tunisia): comparative study in Arab and Berber groups.			
Kumar PP, Mehta S, Purbey PK, Notani D, Jayani RS, Purohit HD, Raju DV, Ravi DS, Bhande RR, Mitra D, Galande S	J Virol	2007 Mar 21;	Related Articles 17376900
SATB1-binding sequences and <b>Alu</b> -like motifs define a unique chromatin context in the vicinity of HIV-1 integration sites.			
Udaka T, Okamoto N, Aramaki M, Tori C, Kosaki R, Hosokai N, Hayakawa T, Takahata N, Takahashi T, Kosaki K	Am J Med Genet A	2007 Apr 1;143(7):721-6	Related Articles 17334995
An <b>Alu</b> retrotransposon-mediated deletion of CHD7 in a patient with CHARGE syndrome.			
Athanasiadis G, Estebani E, Via M, Duggoujon JM, Moschonas N, Chaabani H, Moral P	Eur J Hum Genet	2007 Feb 28;	Related Articles 17327877
The X chromosome <b>Alu</b> insertions as a tool for human population genetics: data from European and African human groups.			
Yunusbayev B, Kutuev I, Khusanov R, Guseinov G, Khunudinova E	Hum Biol	2005 Aug;78(4):465-76	Related Articles 17278621
Genetic structure of Dagestan populations: a study of 11 <b>Alu</b> insertion polymorphisms.			
Pichler I, Mueller JC, Stefanov SA, De Grandi A, Volgato CB, Pinggera GK, Mayr A, Cognigni M, Pioner F, Heisinger T, Pramstaller PP	Hum Biol	2006 Aug;78(4):441-64	Related Articles 17278620
Genetic structure in contemporary south Tyrolean isolated populations revealed by analysis of Y-chromosome, mtDNA, and <b>Alu</b>			

Fig. 3. Web interface of RepWeb. (A) Search interface of RepWeb. The user is able to search the repeat term with four inputs. (B) Repeat Tree. This tree is classified according to repeat class, family and name. In addition, users can search repeat information for the selected repeat. (C) Result of searching in Repbase with a table view.

term which will return hundreds of hits. After the initial search, Web Crawler delivers the results and parses the information. The final results are presented as a table on the web (Fig. 3C).

The web server, RepWeb, has been developed to

search journals related to the repeats in Repbase, Google Scholar and PubMed. Through the web interface, researchers are able to retrieve journals pertaining to the queried keyword within three resources. To acquire the journal-related Repbase, RepWeb constructs a local

relational database which integrates two public databases; repeat library data and PubMed. RepWeb is a web interface based on a Tomcat server with dynamic content generated by a Java Server Page (JSP) and AJAX (Asynchronous JavaScript and XML). By applying these technologies, RepWeb provides users with a dynamic, responsive, and faster browser. The interface is presented as a repeat tree menu that allows users to search easily while selecting a particular repeat (Fig. 2B).

## Results

RepWeb is a Java application that provides a user-friendly interface for searching the reference information and journals related to repetitive sequences within the genomes of various organisms. RepWeb has a distinctive advantage, in that it retrieves the results from two search engine sites simultaneously. It relies on two databases, using the search term results from Google Scholar and PubMed, and has a Repbase database which matches the Repbase library with information from PubMed. The application is freely available from <http://www.repeatome.org> or <http://www.repweb.org> or directly from <http://biportal.kobic.re.kr:8080/RepWeb>.

## Future Directions

We will update the search function of RepWeb by sequence, and apply it to a gene pathway using our pipeline and web database system. RepWeb will continue to incorporate newly established repeat elements. The role of RepWeb in providing literature information for biologists will continue to adapt in line with new developments and the needs of researchers in this field.

## Acknowledgements

We are grateful to anonymous RepWeb reviewers for their comments. This work was supported by the Korean Ministry of Science and technology under grant number M10508040002-06N0804-00210 and M10407010001-06N0701-00110. This work was also supported by the KRIBB Research Initiative Program.

## References

- Dean, G. (2005). A look at Google Scholar, PubMed, and Scirus: comparisons and recommendations., In *the Journal of the Canadian Health Libraries Association*, pp. 85-89.
- Entrez Utilities. [http://eutils.ncbi.nlm.nih.gov/entrez/query/static/eutils\\_help.html](http://eutils.ncbi.nlm.nih.gov/entrez/query/static/eutils_help.html).
- Google Scholar. <http://schoar.google.com/scholar/about.html>.
- Jacso, P. (2005). Comparison and analysis of the citedness scores in web of science And Google Scholar., In *Lecture Notes in Computer Science*, Fox, E.A., Neuhold, E.J., Premsmit, P., Wuwongse, V. ed. (Springer-Verlag, Berlin), pp. 360-369.
- Jurka, J., Kapitonov, V.V., Pavlicek, A., Klonowski, P., Kohany, O., Walichiewicz, J. (2000). Repbase update: a database and an electronic journal of repetitive elements. *Trends Genet.* 16, 418-420.
- PubMed. <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?DB=pub Med>.
- RepBase homepage. Genetic Information Research Institute. <http://www.girinst.org/>