

ORegAnno
open regulatory annotation

Canada's Michael Smith

www.bcgsc.ca

1. These authors contributed equally to this work; 2. Canada's Michael Smith Genome Sciences Centre; 3. University of Manchester

3. Implementation

Our understanding of gene regulation is currently limited by our ability to collectively synthesize and catalogue transcriptional regulatory elements stored in scientific literature. The Open Regulatory Annotation (OREgAnno) database is a dynamic collection of regulatory elements, including transcription factor binding sites, enhancers, and binding sites, and regulatory mutations (SNPs and haplotypes). OREgAnno is a web-based database that allows users to submit and retrieve regulatory elements and new annotations from users worldwide. Submissions to OREgAnno are immediately cross-referenced to Ensembl, dbSNP, Entrez Gene, the NCBI Taxonomy database, and other public databases. OREgAnno also includes a comprehensive catalog of regulatory regions, and 10,000 regulatory polymorphisms or haplotypes from 9 species. We are currently in the process of adding a large number of additional records from the literature and other public databases. OREgAnno is a community-driven, open-access community-based forum for annotation of cis-regulatory sequences. It is a first-of-its-kind, open-access structured encyclopedia of cis-regulatory elements and allows for the first time a comprehensive, open-access, and publicly accessible, curated, and verified gene identifiers (Essembl or Ensembl) ensure maximum compatibility with the existing literature. OREgAnno is a valuable resource for researchers studying cis-regulatory sequences represents a valuable resource for researchers investigating transcriptional regulatory or regulatory variation and provides an open-access system for continued development and expansion. OREgAnno is available directly through MySQL, <http://www.oreganno.org/>, or online at <http://www.oreganno.org/>.

Fig1. The ORegAnno resource consists of a (primarily) Java-based web application for the curation, storage and distribution of literature derived regulatory sequences. Entries are cross-referenced against a number of external databases (dbSNP, Ensembl, eVOC, Pubmed), visualized through Ensembl or UCSC browsers, and freely available to the public through direct database access (db01.bcgsc.ca), a perl API or XML.

Figure 1. The ORegAnno Resource

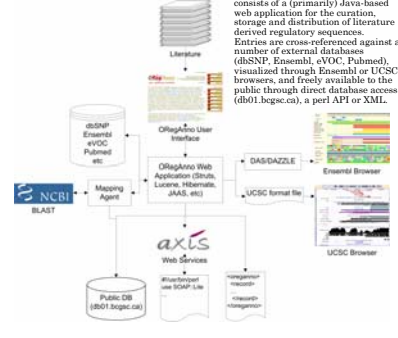


Figure 2. Database schema for mySQL

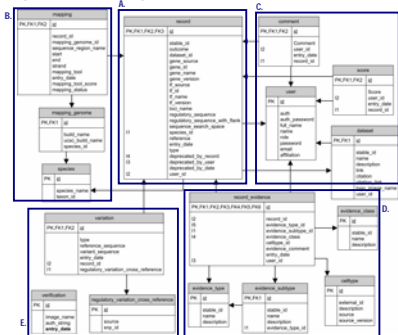


Fig 2. (A) Every evoRc/Anno record consists of a stable id, record type, species, reference, outcome, target gene, transcription factor (if known), sequence identifier. (B) The species and species name are used to derive genomic coordinates by BLAST; (C) Each record is associated with the user who entered it as well as the history of comments and scores it has received. If the record was acquired from an existing database it will be linked to that dataset's information; (D) If the record is a polymorphism or haplotype the variant sequence is also stored as well as any external links for that variant; (E) Each record will normally have some evidence for the function of the sequence from the literature or other sources. Evidence is stored as a list of links to literature, and subtypes (see table 2). If known, the cell type used for the experiments can also be stored using the eVOC cell type ontology[31].

3. Implementation

Figure 3. The ORegAnno User Interface



Fig 3. The ORegAnno user interface provides: (A) Login status; (B) Current contents of database (with link for detailed view); (C) Options to login/logout or create a new user (login only required for annotation); (D) Search engine (powered by Lucene) for basic or advanced searching; (E) Annotation forms for regulatory regions, binding sites, polymorphisms or haplotypes (can be used to create new annotations or to edit existing ones); (F) Tools to download an Ensembl or Entrez target gene; (G) Database downloads/access are available through xml dumps, direct mysql database access, or a perl API (using SOAP); (H) Help documentation provides walkthroughs, guidelines for annotation, and other useful information. (I) A donation page gives credit to major contributors and links to a complete ORegAnno user list.

6. Visualizations

Figure 5. Genome browser views for ORegAnno records.



Fig 5. (A) Ensembl and (B) UCSC views allow the user to visualize any ORegAnno sequence in its genomic context.

4. Evidence

Table 1. Sample of Evidence types and subtypes

Evidence type	Evidence subtype
Electrophoretic Mobility Shift Assay (EMSA)	Direct gel shift
	Supershift
	Gel shift competition
Reporter Gene Assay	Transient transfection luciferase assay
	Chromosomal acetyltransferase (CAT) Assay
	In-vivo GFP Expression Assay
	Dual luciferase reporter gene assay
	In-vivo LacZ Expression Assay
Protein Binding Assay	Chromatin immunoprecipitation (ChIP)
	DNAse Footprinting Assay
	Yeast 1-hybrid assay

Table 1. Each ORegAnno record is associated with one or more pieces of evidence. Oreganno currently contains 9 types and 30 subtypes of evidence. A user with administrator status can add new evidence types and subtypes as needed.

8. Acknowledgments

We would like to acknowledge the Wasserman lab (<http://www.cisreg.ca/ykwon/>), and James Fickett (<http://www.chi.umn.edu/MTIR/HomePage.html>) for generously making their regulatory element catalogues publicly available. We thank the OreAnno users for their continuing efforts to improve this resource through manual curation and record validation.

funding | We gratefully acknowledge funding from Genome Canada, Genome British Columbia and the BC Cancer Foundation. SBM was supported by the Natural Sciences and Engineering Research Council (NSERC) and the Michael Smith Foundation for Health Research (MSFHR). OLC was supported by the Canadian Institutes of Health Research (CIHR), NSERC and MSFHR. We thank Dr. Michael MCGS, Michael Smith Foundation for Health Research.

references | 1. Kelso *et al.* 2003; 2. Bergman *et al.* 2005; 3. Ho Sui *et al.* 2005; 4. Wasserman and Fickett. 1998; 5. Ponomarenko *et al.* 2001.