# CryptoDB: a *Cryptosporidium* bioinformatics resource update

**Mark Heiges, Haiming Wang, Edward Robinson, Cristina Aurrecoechea, Xin Gao[1,2], Nivedita Kaluskar[1], Philippa Rhodes[1], Sammy Wang[1], Cong-Zhou He[1], Yanqi Su[1], John Miller[1], Eileen Kraemer[1] and Jessica C. Kissinger[2,]***

Center for Tropical and Emerging Global Diseases, [1]Department of Computer Science and [2]Department of Genetics, University of Georgia, Athens, GA, USA

## ABSTRACT

**The database, CryptoDB (http://CryptoDB.org), is a community bioinformatics resource for the AIDS-related apicomplexan-parasite, *Cryptosporidium*. CryptoDB integrates whole genome sequence and annotation with expressed sequence tag and genome survey sequence data and provides supplemental bioinformatics analyses and data-mining tools. A simple, yet comprehensive web interface is available for mining and visualizing the data. CryptoDB is allied with the databases PlasmoDB and ToxoDB via ApiDB, an NIH/NIAID-fundedBioinformatics Resource Center. Recent updates to CryptoDB include the deposition of annotated genome sequences for *Cryptosporidium parvum* and *Cryptosporidium hominis*, migration to a relational database (GUS), a new query and visualization interface and the introduction of Web services.**

## INTRODUCTION

The Apicomplexan parasite *Cryptosporidium* is a global causative agent of severe and chronic diarrheal disease in humans and other animals. As no reliable chemo- or immuno-therapy is currently available, infections can be life threatening for people with a compromised immune system, such as AIDS patients. The pathogen is typically spread via contaminated drinking water and is resistant to water chlorination and filtration (1). Because of the water safety threat to public health, *Cryptosporidium* is ranked as a Category B Biodefense Pathogen by the National Institutes of Health. Bioinformatics analysis plays an important role in understanding the biology of and identification of potential drug targets in this medically important parasite. To aid the research community in this line of inquiry, the online database CryptoDB continues to update and expand its role of warehousing and interfacing *Cryptosporidium* genome sequence, annotation, sequence analysis and other *Cryptosporidium*-related information.

## UPDATED DATASET

Version 3.0 of CryptoDB was released in April 2005 and contains the published genome sequence and annotation for *Cryptosporidium hominis*, strain TU502 (2) and *Cryptosporidium parvum*, strain IOWA (3). The database houses copies of assembled genome contigs and gene annotations deposited in GenBank (4) by the sequence generators. The *C.parvum* genome sequence is represented by 18 contigs ranging in size from 17 kb to 1.2 Mb in length and annotated with 3885 total genes. The *C.hominis* genome sequence is represented in 1422 contigs ranging in size from a few hundred to 90 thousand base pairs in length and annotated with 3956 total genes. *C.parvum* chromosome 6 has been independently sequenced and annotated (5) and is represented in the database.

In addition to the data provided by genome sequencing efforts, ∼6 Mb of genome survey sequence (GSS) and expressed sequence tag (EST) (6) data are incorporated. The ESTs are clustered into RNA transcripts and aligned to the genome using the methodology applied at ApiEST-DB (7).

Gene annotations provided by the genome sequencing centers are augmented with supplemental analyses. Pre-computed BLASTX analyses of *Cryptosporidium* contigs versus the GenBank non-redundant protein database and EST alignments to contig sequences offer supporting evidence for gene predictions. Potential syntenic relationships of the *C.hominis* and *C.parvum* contig sequences are calculated and graphically displayed. Protein feature predictions of signal peptides and transmembrane domains are provided. Open reading frames >50 and >100 amino acids in length have been calculated for all nucleic acid sequences in all six reading frames.

*To whom correspondence should be addressed. Tel: +1 706 542 6562; Fax: +1 706 542 3910; Email: jkissing@uga.edu

All sequence datasets are available for bulk download in FASTA format. Programmatic access to selected resources is provided via Web service interfaces.

## IMPROVED DATABASE AND WEB INTERFACE

CryptoDB 3.0 is backed by a relational database utilizing the Genomics Unified Schema (GUS, GUSdb.org) (8) and Oracle 10g. Migration to a relational database architecture marks a major improvement over previous releases because of the new services and resources that can now be offered.

The CryptoDB web interface provides a set of forms through which users can easily query the annotation and pre-computed analysis data (Figure 1A). Queries for contig sequence, gene and protein features are possible and can be restricted to either or both of the hosted species genomes. At the gene level, users can conduct text searches of gene product descriptions, search for genes by RNA type (mRNA, rRNA, snoRNA and tRNA) and find genes having alignments to *C.parvum* ESTs. For protein features, users can select genes predicted by SignalP (9) to encode a signal peptide or predicted by TMHMM (10) to contain transmembrane domains. Users may also retrieve a specific gene by locus tag or a contig
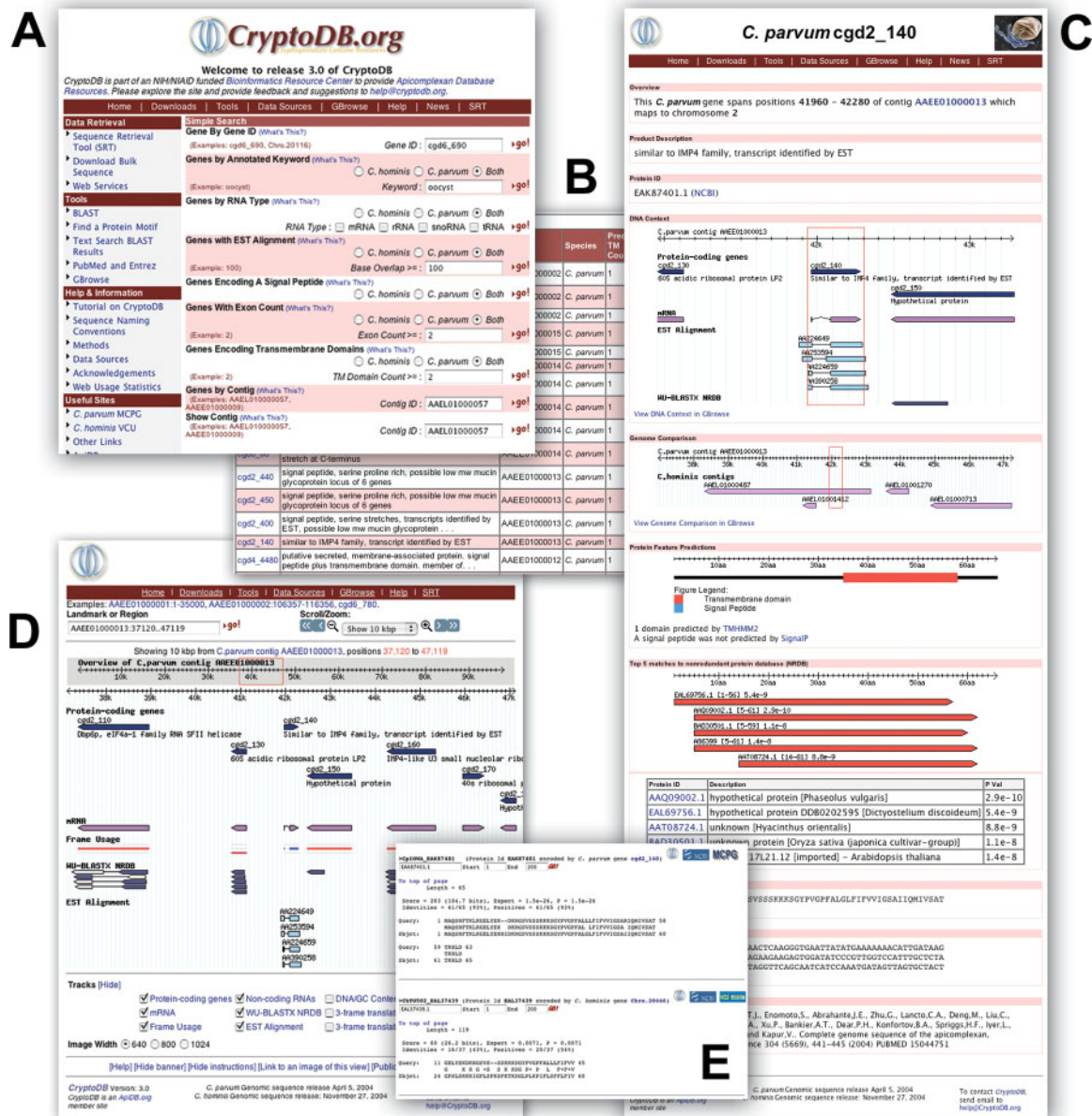


**Figure 1.** Database functionality. (**A**) Searches are initiated via queries provided on the web site's front page. (**B**) The results are returned as a summary table with links to detailed record pages. (**C**) Detailed record page with summary of all available data/information for this gene or contig sequence. (**D**) GBrowse of genomic region of interest provides a graphical view of annotations and similarity analyses. (**E**) Search results, such as BLAST reports, are linked to detailed records and to a sequence retrieval utility.

sequence by accession number. *Ad hoc* data selections not obtainable via the provided queries may be requested by email to help@cryptodb.org.

Gene pages and contig sequence pages provide a detailed view of annotation and analysis for a given record in the database. Gene pages contain a text overview of the gene, including the coordinate position on its contig and product description when available (Figure 1C). GBrowse (11) has been utilized to provide a graphic display of annotated gene features and the data mapped to the genome, BLAST hits, ESTs, etc. (Figure 1D).

The web interface includes a mechanism to allow users to readily download the sequences and other attributes associated with their query result set. A query history permits users to track their searches and combine them into more complex queries across data types (e.g. 'list all genes on chromosome 3 that contain transmembrane domains').

## ANALYSIS AND RESEARCH TOOLS

Several tools for data mining augment the published annotations and pre-computed analyses. Users may BLAST their own sequences against the genomic contig, annotated protein, GSS and EST sequence databases (Figure 1E). A motif search tool finds protein sequences with PROSITE (12) or user-defined amino acid patterns. Keywords from *Cryptosporidium* genomic sequences versus GenBank NRDB BLASTX results are indexed and searchable. In each case, the results contain links back to detailed gene, protein or contig record pages or to external databases (e.g. GenBank) as appropriate (Figure 1B).

To facilitate tracking of the latest literature, PubCrawler (13) is used to poll NCBI's PubMed and GenBank each week day for new *Cryptosporidium*-related updates.

## AFFILIATIONS

CryptoDB is a member of ApiDB.org, an NIH/NIAID funded Bioinformatics Resource Center (BRC) for Biodefense and Emerging or Re-Emerging Infectious Diseases (www.niaid. nih.gov/dmid/genomes/brc/default.htm). Other ApiDB members include the genome databases for *Plasmodium* (PlasmoDB) (14,15) and *Toxoplasma* (ToxoDB) (16) and the Apicomplexan EST database, ApiEST-DB (7). CryptoDB and other member databases are linking to ApiDB in a coordinated effort to promote comparative studies and ease of access across these apicomplexan genomes.

## WEB SERVICES AND DATABASE NEWS

To facilitate database integration with ApiDB, other NIAID BRC's and programmatic access of CryptoDB by others, web services for CryptoDB have been implemented. Web services are pieces of software that can communicate across the Internet to build distributed applications. They can do this regardless of the software used for their implementation as long as they use a common protocol, SOAP (17). CryptoDB uses SOAP and provides published WSDL files and sample client software in Java (using Axis) and PERL (using SOAP::Lite). Currently, one service that retrieves FASTA sequence files is active. Additional services and infrastructure (18) are planned.

To facilitate the dissemination of news and updates related to CryptoDB, we have established a Really Simple Syndication news feed (RSS) that is displayed on the home page of CryptoDB and ApiDB and can be read by any RSS news aggregator.

## FUTURE PLANS

CryptoDB is fully funded and staffed with biologists and software developers with close ties to software developers for GUS, ToxoDB and PlasmoDB. This fertile ground will support many opportunities for frequent database updates and expansions with new data types, analyses, data-mining tools and visual displays. Gene ontology terms and protein feature signatures from InterProScan (19) analyses will be included with gene records. Improvements to visualization of genome-wide synteny are planned. SRI International's Pathway Tools software (20) is being added to facilitate analyses of metabolic pathways in both annotated *Cryptosporidium* genome sequences. Future releases of CryptoDB will publish this information, as 'CryptoCyc' for querying and visualization in a graphical display.

To facilitate data sharing, CryptoDB has the capacity to activate a Distributed Annotation Server (DAS) (21) via a DAS-GUS adapter if needed.

CryptoDB will continue to work with the ApiDB consortium to further integrate its resources with other apicomplexan genome sites. Database federating technologies, web services and portal designs are currently being implemented toward this end. Data exchange and interoperability with other NIAID Bioinformatics Resource Centers will be a continued effort.

## REFERENCES

1. Fayer,R. (1997) *Cryptosporidium and Cryptosporidiosis.* CRC Press, Inc., Boca Raton, FL.
2. Xu,P., Widmer,G., Wang,Y., Ozaki,L.S., Alves,J.M., Serrano,M.G., Puiu,D., Manque,P., Akiyoshi,D., Mackey,A.J. *et al.* (2004) The genome of *Cryptosporidium hominis. Nature,* **431**, 1107–1112.

3. Abrahamsen,M.S., Templeton,T.J., Enomoto,S., Abrahante,J.E., Zhu,G., Lancto,C.A., Deng,M., Liu,C., Widmer,G., Tzipori,S. *et al.* (2004) Complete genome sequence of the apicomplexan, *Cryptosporidium parvum*. *Science*, **304**, 441–445.

4. Xu,P., Widmer,G., Wang,Y., Ozaki,L.S., Alves,J.M., Serrano,M.G., Puiu,D., Manque,P., Akiyoshi,D., Mackey,A.J. *et al.* (2004) Corrigendum: the genome of *Cryptosporidium hominis*. *Nature*, **432**, 415.

5. Bankier,A.T., Spriggs,H.F., Fartmann,B., Konfortov,B.A., Madera,M., Vogel,C., Teichmann,S.A., Ivens,A. and Dear,P.H. (2003) Integrated mapping, chromosomal sequencing and sequence analysis of *Cryptosporidium parvum*. *Genome Res.*, **13**, 1787–1799.

6. Strong,W.B. and Nelson,R.G. (2000) Preliminary profile of the *Cryptosporidium parvum* genome: an expressed sequence tag and genome survey sequence analysis. *Mol. Biochem. Parasitol.*, **107**, 1–32.

7. Li,L., Crabtree,J., Fischer,S., Pinney,D., Stoeckert,C.J.Jr, Sibley,L.D. and Roos,D.S. (2004) ApiEST-DB: analyzing clustered EST data of the apicomplexan parasites. *Nucleic Acids Res.*, **32**, D326–D328.

8. Davidson,S.B., Crabtree,J., Brunk,B., Schug,J., Tannen,V., Overton,G.C. and Stoeckert,J.C.J. (2001) K2/Kleisli and GUS: experiments in integrated access to genomic data sources. *IBM systems Journal*, **40**, 512–531.

9. Bendtsen,J.D., Nielsen,H., von Heijne,G. and Brunak,S. (2004) Improved prediction of signal peptides: SignalP 3.0. *J. Mol. Biol.*, **340**, 783–795.

10. Krogh,A., Larsson,B., von Heijne,G. and Sonnhammer,E.L. (2001) Predicting transmembrane protein topology with a hidden Markov model: application to complete genomes. *J. Mol. Biol.*, **305**, 567–580.

11. Stein,L.D., Mungall,C., Shu,S., Caudy,M., Mangone,M., Day,A., Nickerson,E., Stajich,J.E., Harris,T.W., Arva,A. *et al.* (2002) The generic genome browser: a building block for a model organism system database. *Genome Res.*, **12**, 1599–1610.

12. Sigrist,C.J., Cerutti,L., Hulo,N., Gattiker,A., Falquet,L., Pagni,M., Bairoch,A. and Bucher,P. (2002) PROSITE: a documented database using patterns and profiles as motif descriptors. *Brief. Bioinform.*, **3**, 265–274.

13. Hokamp,K. and Wolfe,K.H. (2004) PubCrawler: keeping up comfortably with PubMed and GenBank. *Nucleic Acids Res.*, **32**, W16–W19.

14. Bahl,A., Brunk,B., Crabtree,J., Fraunholz,M.J., Gajria,B., Grant,G.R., Ginsburg,H., Gupta,D., Kissinger,J.C., Labo,P. *et al.* (2003) PlasmoDB: the *Plasmodium* genome resource. A database integrating experimental and computational data. *Nucleic Acids Res.*, **31**, 212–215.

15. Kissinger,J.C., Brunk,B.P., Crabtree,J., Fraunholz,M.J., Gajria,B., Milgram,A.J., Pearson,D.S., Schug,J., Bahl,A., Diskin,S.J. *et al.* (2002) The *Plasmodium* genome database. *Nature*, **419**, 490–492.

16. Kissinger,J.C., Gajria,B., Li,L., Paulsen,I.T. and Roos,D.S. (2003) ToxoDB: accessing the *Toxoplasma gondii* genome. *Nucleic Acids Res.*, **31**, 234–236.

17. Stein,L. (2002) Creating a bioinformatics nation. *Nature*, **417**, 119–120.

18. Sivashanmugam,K., Miller,J.A., Sheth,A.P. and Verma,K. (2004) Framework for semantic web process composition. *International Journal of Electronic Commerce*, **9**, 71–106.

19. Quevillon,E., Silventoinen,V., Pillai,S., Harte,N., Mulder,N., Apweiler,R. and Lopez,R. (2005) InterProScan: protein domains identifier. *Nucleic Acids Res.*, **33**, W116–W120.

20. Karp,P.D., Paley,S. and Romero,P. (2002) The Pathway Tools software. *Bioinformatics*, **18** (Suppl. 1), S225–S232.

21. Dowell,R.D., Jokerst,R.M., Day,A., Eddy,S.R. and Stein,L. (2001) The distributed annotation system. *BMC Bioinformatics*, **2**, 7.