

# GOBASE: an organelle genome database

Emmet A. O'Brien\*, Yue Zhang, Eric Wang, Veronique Marie, Wole Badejoko,  
B. Franz Lang and Gertraud Burger

Robert-Cedergren Center for Bioinformatics and Genomics, Département de Biochimie, Pavillon Roger-Gaudry,  
Université de Montréal, 2900 Edouard-Montpetit, Montreal QC, Canada H3T 1J4

Received September 11, 2008; Revised October 10, 2008; Accepted October 13, 2008

## ABSTRACT

The organelle genome database GOBASE, now in its 21st release (June 2008), contains all published mitochondrion-encoded sequences (~913 000) and chloroplast-encoded sequences (~250 000) from a wide range of eukaryotic taxa. For all sequences, information on related genes, exons, introns, gene products and taxonomy is available, as well as selected genome maps and RNA secondary structures. Recent major enhancements to database functionality include: (i) addition of an interface for RNA editing data, with substitutions, insertions and deletions displayed using multiple alignments; (ii) addition of medically relevant information, such as haplotypes, SNPs and associated disease states, to human mitochondrial sequence data; (iii) addition of fully reannotated genome sequences for *Escherichia coli* and *Nostoc* sp., for reference and comparison; and (iv) a number of interface enhancements, such as the availability of both genomic and gene-coding sequence downloads, and a more sophisticated literature reference search functionality with links to PubMed where available. Future projects include the transfer of GOBASE features to NCBI/GenBank, allowing long-term preservation of accumulated expert information. The GOBASE database can be found at <http://gobase.bcm.umontreal.ca/>. Queries about custom and large-scale data retrievals should be addressed to [gobase@bch.umontreal.ca](mailto:gobase@bch.umontreal.ca).

## INTRODUCTION

The amount of information available in generalist molecular sequence databases such as GenBank (1) continues to grow, and this information becomes more diverse and complex as we discover new biological phenomena. Therefore, there is an increasing need for expert databases

specializing in particular areas of molecular biology. Specialist databases provide expert curation of data, and access to that data in a flexible and well-integrated fashion serves a purpose complementary to generalist databases such as GenBank.

GOBASE is one such specialist database, which has been collecting, curating and publishing data concerning mitochondrial and chloroplast genomes since 1995 (2–5). Organelle genomes are of biological interest for a wide range of studies, such as molecular taxonomy, molecular mechanisms of trans-splicing and RNA editing, and non-Mendelian inherited metabolism-related disease in humans. GOBASE contains a number of different categories of data, such as nucleic acid and protein sequences, genetic maps, taxonomic data and RNA secondary structures. All gene and product names have been assigned from a locally maintained standard list, and this combines with a powerful and flexible interface to allow a wide range of complex searches. While initially GOBASE was designed primarily to address issues of comparative biology, such as the diversity of organelle genome structure in eukaryotes (e.g. 6,7), we have more recently added functionality specific to the human mitochondrial genome in GOBASE, such as searches by haplotype and disease state, which are of medical interest.

## DATA CONTENT

GOBASE release 21 (June 2008) contains 913 000 mitochondrial sequences including 737 000 genes, and 250 000 chloroplast-encoded sequences including 174 000 genes, derived mostly from GenBank releases up to 164. The large number of complete organelle genomes available makes GOBASE a valuable resource for phylogenomics, with 6300 complete mitochondrial genomes and 213 chloroplast genomes. This number has increased almost 4-fold since the previous report.

More recently (5), we have added bacterial genome sequences for reference purposes. As of release 21 GOBASE includes three complete bacterial genomes: *Escherichia coli* K12; the alpha-proteobacterium

\*To whom correspondence should be addressed. Tel: +1 514 343 6111; Fax: +1 514 343 2210; Email: [eobrien@bch.umontreal.ca](mailto:eobrien@bch.umontreal.ca)





a)

Retrieve		Clear	
Select a gene or genes to see mutations in, OR select a range of positions on the human mitochondrial genome.			
	Gene Name	Start	end
<input type="checkbox"/>	trnF(gaa)	578	648
<input type="checkbox"/>	rns	649	1602
<input type="checkbox"/>	trnV(uac)	1603	1671
<input type="checkbox"/>	rnl	1672	3229
<input type="checkbox"/>	trnL(uaa)	3230	3304
<input checked="" type="checkbox"/>	nad1	3307	4263
<input type="checkbox"/>	trnI(gau)	4263	4331

b)

Mutation Position	Mutation	Mutation Type	Disease Associated	Disease id	Gene-Name
3497	C->T	substitution	Leber Hereditary Optic Neuropathy	9	nad1

7 8901234567 8901234567 8901234567 8901234567 8901234567 8901234567 8901234567 8901234567 8901234567 8901234567  
 Refseq 3477 CAAAGAGCCCTAAAACCCG[CACATCTACCATCACCCCTACATCACCGCCCGACCTTAGCTCTCACC 3546  
 CAAAGAGCCCTAAAACCCG[CACATCTACCATCACCCCTACATCACCGCCCGACCTTAGCTCTCACC

Sequence(s) with Mutation

There are 22 sequences with this mutation at position 3497.

- |                          |                          |                          |                          |
|--------------------------|--------------------------|--------------------------|--------------------------|
| <a href="#">9449636</a>  | <a href="#">9554055</a>  | <a href="#">9554123</a>  | <a href="#">9554144</a>  |
| <a href="#">9950266</a>  | <a href="#">10099587</a> | <a href="#">10100022</a> | <a href="#">10101102</a> |
| <a href="#">10101582</a> | <a href="#">10101807</a> | <a href="#">10102167</a> | <a href="#">10102587</a> |
| <a href="#">10102827</a> | <a href="#">10102872</a> | <a href="#">10103622</a> | <a href="#">10103682</a> |
| <a href="#">10103982</a> | <a href="#">10104222</a> | <a href="#">10104297</a> | <a href="#">10104417</a> |
| <a href="#">10105707</a> | <a href="#">10105857</a> |                          |                          |

Figure 3. (a) Human mutation query page, allowing the user to select the gene(s) of interest and specify the range of positions on the sequence to search for mutations. (b) Result page showing details for an individual mutation.

sequence or gene-coding regions, selectable via buttons from the Gene query page. There are a small number of unusual cases, such as trans-spliced genes, where there is no straightforward correspondence between a single gene and a contiguous linear region of the source sequence record. The GOBASE database structure has now been modified to address these cases transparently. Sequences of complex gene-coding regions are assembled in advance, stored and made available in query results through the same interface as conventional linear genes.

All sequences retrieved from GOBASE now come with detailed literature references derived from the source GenBank records. Journal, author and title are provided, and a direct link to the appropriate PubMed entry if one exists.

Because of practical constraints, any given query in GOBASE returns at most 5000 results. Users wishing to execute custom queries retrieving larger amounts of data are invited to contact the GOBASE team at gobase@bch.umontreal.ca so that the query can be run directly on the database via SQL.

IMPLEMENTATION

The GOBASE database is implemented in version 7.4.1 of the PostgreSQL relational database management system with a web interface written in v4.3.8 of the PHP scripting language. The graphics on the gene pages are generated using the GD module for Perl/PHP, version 2.0.25. Perl (5.8.0) scripts are used to download data from GenBank and process it into GOBASE. All procedures are executed on PCs with two 2.4 GHz or 2.8 GHz Intel Xeon CPUs.

FUTURE PLANS

Specialized databases with all their valuable information are prone to disappearance (15), mostly because of funding constraints, unless transferred to sustainable public databases. We are therefore collaborating with scientists at NCBI to establish a database based on the content of GOBASE as an auxiliary to GenBank. This database will focus on the additional data that expert curation at GOBASE has generated, notably the curated gene and

product names and synonyms and RNA secondary structure data, thus providing a permanent repository for two decades of curation of organelle genome data.

## ACKNOWLEDGEMENTS

The authors would like to thank Ilene Mizrachi, Susan Schaefer, Tatiana Tatusova and Jim Ostell at NCBI; Chris Cesaire, Ousman Diallo, and Olivier Tremblay-Savard for contributions to the development of the RNA editing functionality in GOBASE, and Allan Sun for systems administration.

## FUNDING

This project was funded by grants MOP-15331 and MOP-84453 from the Canadian Institute for Health Research (CIHR, Genetics Institute). Funding for open access charge: CIHR.

*Conflict of interest statement.* None declared.

## REFERENCES

1. Benson,D.A., Karsch-Mizrachi,I., Lipman,D.J., Ostell,J. and Wheeler,D.L. (2008) GenBank. *Nucleic Acids Res.*, **36**, D25–D30.
2. Korab-Laskowska,M., Rioux,P., Brossard,N., Littlejohn,T.G., Gray,M.W., Lang,B.F. and Burger,G. (1998) The Organelle Genome Database Project (GOBASE). *Nucleic Acids Res.*, **26**, 138–144.
3. Shimko,N., Liu,L., Lang,B.F. and Burger,G. (2001) GOBASE: the organelle genome database. *Nucleic Acids Res.*, **29**, 128–132.
4. O'Brien,E.A., Badidi,E., Barbasiewicz,A., deSousa,C., Lang,B.F. and Burger,G. (2003) GOBASE – a database of mitochondrial and chloroplast information. *Nucleic Acids Res.*, **31**, 176–178.
5. O'Brien,E.A., Zhang,Y., Yang,L., Wang,E., Marie,V., Lang,B.F. and Burger,G. (2006) GOBASE – a database of organelle and bacterial genome information. *Nucleic Acids Res.*, **34**, D697–D699.
6. Lang,B.F., Gray,M.W. and Burger,G. (1999) Mitochondrial genome evolution and the origin of eukaryotes. *Annu. Rev. Genetics.*, **33**, 351–397.
7. Burger,G., Gray,M.W. and Lang,B.F. (2003) Mitochondrial genomes: anything goes. *Trends Genet.*, **19**, 709–716.
8. Koski,L.B., Gray,M.W., Lang,B.F. and Burger,G. (2005) AutoFACT: an automatic functional annotation and classification tool. *BMC Bioinform.*, **6**, 151.
9. Covello,P.S. and Gray,M.W. (1989) RNA editing in plant mitochondria. *Nature*, **341**, 662–666.
10. Benne,R., Van den Burg,J., Brakenhoff,J.P., Sloof,P., Van Boom,J.H. and Tromp,M.C. (1986) Major transcript of the frameshifted coxII gene from trypanosome mitochondria contains four nucleotides that are not encoded in the DNA. *Cell*, **46**, 819–826.
11. Hoch,B., Maier,R.M., Appel,K., Igloi,G.L. and Kössel,H. (1991) Editing of a chloroplast mRNA by creation of an initiation codon. *Nature*, **353**, 178–180.
12. Attimonelli,M., Acceturro,M., Santamaria,M., Lascaro,D., Scioscia,G., Pappad,G., Russo,L., Zanchetta,L. and Tommaseo-Ponzetta,M. (2005) HmtDB, a human mitochondrial genomic resource based on variability studies supporting population genetics and biomedical research. *BMC Bioinform.*, **1**, S4.
13. Wheeler,D.L., Barrett,T., Benson,D.A., Bryant,S.H., Canese,K., Chetvernin,V., Church,D.M., Dicuccio,M., Edgar,R., Federhen,S. *et al.* (2008) Database resources of the National Center for Biotechnology Information. *Nucleic Acids Res.*, **36**, D13–D21.
14. Ruiz-Pesini,E., Lott,M.T., Procaccio,V., Poole,J.C., Brandon,M.C., Mishmar,D., Yi,C., Kreuziger,J., Baldi,P. and Wallace,D.C. (2007) An enhanced MITOMAP with a global mtDNA mutational phylogeny. *Nucleic Acids Res.*, **35**, D823–D828.
15. Merali,Z. and Giles,G. (2005) Databases in peril. *Nature*, **23**, 1010–1011.