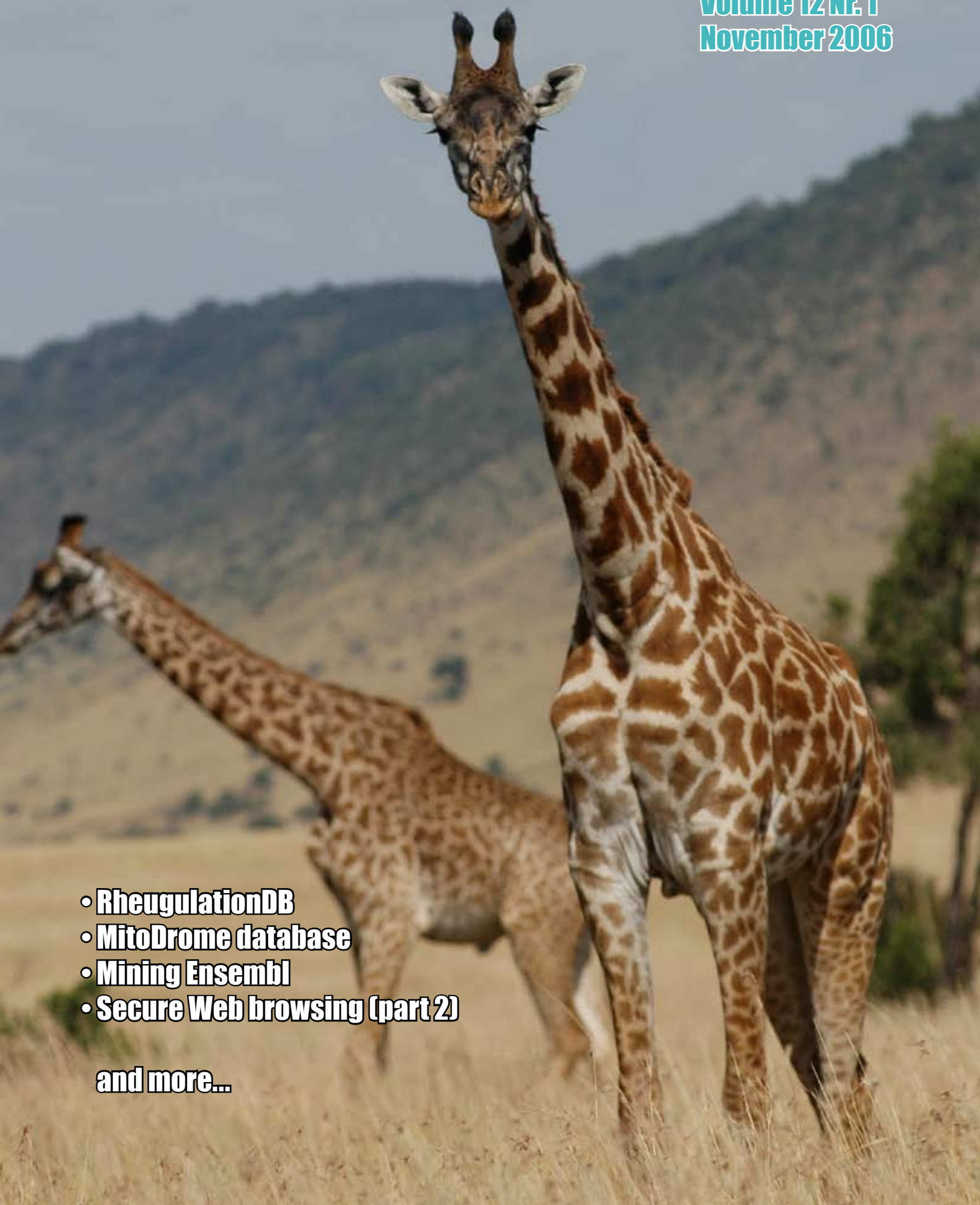


# EMBnet.news

Volume 12 Nr. 1  
November 2006



- **RheugulationDB**
- **MitoDrome database**
- **Mining Ensembl**
- **Secure Web browsing (part 2)**

**and more...**

# Editorial

EMBnet news is the newsletter of the world's largest Bioinformatics community. Bioinformaticians and regular users refer to it as a source of practical information on a wide range of subjects, justifying the ever increasing number of downloads. A new life cycle for EMBnet news is beginning. The editorial board is placing an extra effort to produce this publication more dynamically and with more timely news. Several options are on the table. It is your chance, as a reader, to influence these choices by sending us suggestions and valuable feedback. Please do so by sending e-mail to [emb-pr@embnet.org](mailto:emb-pr@embnet.org).

The editorial board: Erik Bongcam-Rudloff, Domenica D'Elia, Pedro Fernandes, Kimmo Mattila, Lubos Klucar, and Gonçalo Guimaraes Pereira.



Protein Spotlight (ISSN 1424-4721) is a periodical electronic review from the SWISS-PROT group of the Swiss Institute of Bioinformatics (SIB). It is published on a monthly basis and consists of articles focused on particular proteins of interest. Each issue is available, free of charge, in HTML or PDF format at

<http://www.expasy.org/spotlight>

We provide the EMBnet community with a printed version of issue 64. Please let us know if you like this inclusion.

Cover picture: Reticulated Giraffe (*Giraffa camelopardalis reticulata*), Masai Mara reserve, Kenya, August 2006 [© Erik Bongcam-Rudloff]

# Contents

Editorial	2
RheugulationDB	3
Course report	5
The MitoDrome database	6
Mining Ensembl	12
methBLAST server	15
Secure Web browsing (2)	17
Single handed node manag.	19
Protein spotlight 64	21
Node information	23

## Editorial Board:

Erik Bongcam-Rudloff, The Linnaeus Centre for Bioinformatics, SLU/UU, SE  
 Email: [erik.bongcam@bmc.uu.se](mailto:erik.bongcam@bmc.uu.se)  
 Tel: +46-18-4716696  
 Fax: +46-18-4714525

Domenica D'Elia, Institute for Biomedical Technologies - CNR, Bari, IT  
 Email: [domenica.delia@ba.itb.cnr.it](mailto:domenica.delia@ba.itb.cnr.it)  
 Tel: +39-80-5929674  
 Fax: +39-80-5929690

Pedro Fernandes, Instituto Gulbenkian, PT  
 Email: [pfern@igc.gulbenkian.pt](mailto:pfern@igc.gulbenkian.pt)  
 Tel: +315-214407912  
 Fax: +315-214407970

Gonçalo Guimaraes Pereira, UNICAMP, BR  
 Email: [goncalo@unicamp.br](mailto:goncalo@unicamp.br)  
 Tel: +55-19-37886237/6238  
 Fax: +55-19-37886235

Lubos Klucar, Institute of Molecular Biology, SAS Bratislava, SK  
 Email: [klucar@embnet.sk](mailto:klucar@embnet.sk)  
 Tel: +421-2-59307413  
 Fax: +421-2-59307416

Kimmo Mattila, CSC, Espoo, FI  
 Email: [kimmo.mattila@csc.fi](mailto:kimmo.mattila@csc.fi)  
 Tel: +358-9-4572708  
 Fax: +358-9-4572302

## RheugulationDB

an integrated resource for  
Rheumatoid Arthritis research



**José Teles**

MSc Student in Bioinformatics.  
Instituto de Medicina Molecular,  
Lisboa, PT, Instituto Gulbenkian de  
Ciência, Oeiras, PT (a\_teles@yahoo.com)

<http://www.rheugulationdb.com>

Rheumatoid Arthritis (RA) is a chronic inflammatory disease, with a strong autoimmune basis and preferential involvement of joint synovial tissue, affecting approximately 1% of the world population. Several interconnected pathways are involved in RA physiopathology, and specific molecules may be key players at particular spatial and temporal contexts of the disease. Due to this multifactorial nature, reductionist approaches have been successful in revealing some of these molecules but failed to provide an integrated

picture of the molecular events underlying RA onset and progression, generating a growing need for extensive integration of all the available information. Such an approach could provide us with a more global, biologically accurate picture of the molecular bases underlying RA immunological imbalance, with practical implications on disease diagnosis, prognosis and therapy.

The Rheumatology Research Unit (Instituto de Medicina Molecular - Lisbon) with the support of Abbott Immunology has been aiming to fill this gap, by providing a dynamic platform of updated information on the molecular complexity of the disease, in an open-access on-line resource suited for all rheumatology professionals and students: RheugulationDB. The population of the database was achieved by means of an extensive and thorough screening of the available literature and biomedical databases, with the resulting pool of information being stored in a MySQL relational database. Until this moment, more than 320 genes and 18 cell types with relevance in the molecular context of RA have been annotated and classified according to their presence in characteristic molecular processes of disease, based on more than 140 fundamental references. All these contents can be browsed

The screenshot shows the RheugulationDB website interface. At the top, there is a navigation bar with buttons for 'RHEUMATOID ARTHRITIS', 'BROWSE', 'TOOLS', 'LINKS', and 'CONTACT US'. Below the navigation bar, the search results for the gene TNF are displayed. The results include the official symbol (TNF), description (tumor necrosis factor (TNF superfamily, member 2)), aliases (TNFSF2, DIF), and previous symbols (TNFA). A central table provides links to various databases and resources:

<p>GENE INFORMATION</p> <p><a href="#">ENTREZ GENE</a> <a href="#">ENSEMBL</a> <a href="#">KEGG</a></p>	<p>VARIABILITY AND DISEASE ASSOCIATION</p> <p><a href="#">dbSNP</a> <a href="#">HAPMAP</a> <a href="#">HGMD</a></p>
<p>PROTEIN INFORMATION</p> <p><a href="#">UNIPROT</a> <a href="#">PFAM</a></p>	<p>PHARMA</p> <p><a href="#">PHARMGKB</a></p>
<p>ONTOLOGIES</p> <p><a href="#">GO</a></p>	<p>REFERENCES</p> <p><a href="#">PUBMED</a> <a href="#">OMIM</a></p>

Figure 1. Searching for a gene by name shows the outreach of the system.

very intuitively on-line by two means: graphical representations of the molecular processes, with direct links for the specific information on each gene involved; and querying of the databases for more specific information concerning a particular gene, cell type or molecular pathway. For each gene, we provide links to several external databases, allowing the user to increase the specificity of analysis for particular aspects, such as genes (Entrez Gene, Ensembl, KEGG), proteins (Uniprot, PFAM) or variability and disease association (dbSNP, HapMap, HGMD). Also, the links for the abstracts of all the references used in the literature review are provided, for easy access to particular papers. Worthy of remark is the fact that all graphical notations used in this project are the result of a parallel project – Rheugulation graphics – in which we aim at developing a standard representation for each of the molecules and cell types displayed on RheugulationDB. This will soon be made available for the use of the community, and we expect to further improve this notation, and address other utilizations, such as the dynamic creation of molecular pathways on demand.

Furthermore we provide to less specialized users, such as students of the biomedical field or RA patients searching for more specific information, with a considerable amount of basic information regarding disease etiology, physiopathology, clinical aspects and management, which may prove useful not only as a general informative vehicle, but also allowing the less experienced user to grasp the true potential of integrating all this information in a dynamical information search environment. Finally, a set of utilitarian tools are being developed, in order to assist clinicians and investigators on their basic daily routine. The most immediate example is the Disease Activity Score (DAS) 28 calculator, which will be available on-line soon, allowing the easy calculation of this fundamental parameter of disease activity, as well as the subsequent printing and/or pdf download for patient medical record purposes.

After this first phase of implementation, the most valuable feature of RheugulationDB is the possibility of dynamically browsing a wealth of information on the molecular events of RA, in an integrated fashion, resulting from the very own

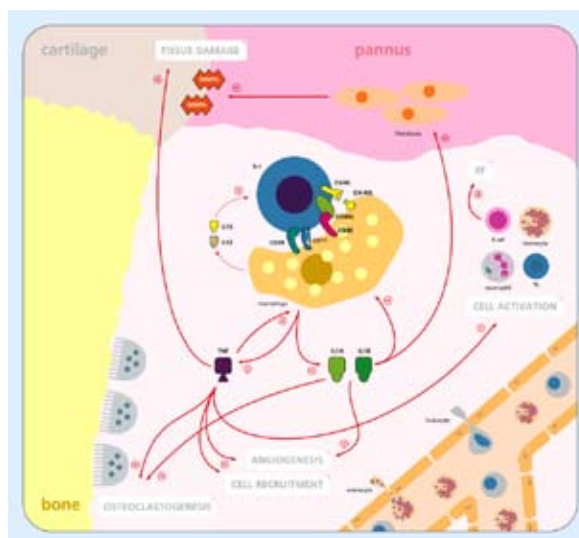


Figure 2. An example entry: PIVOTAL ROLE OF TNF AND IL-1 - Orchestration of inflammatory events in RA.

design and architecture of the database (Figure 1). This feature will be further explored in order to increase the complexity of results that can be obtained by querying the database in the future, thus eliciting new prospects in fields like candidate-gene studies or pharmacogenomic analysis and drug design. We expect RheugulationDB to be a valuable tool in the study of RA, useful for investigators and clinicians on their daily routine, potentially widening the scope of present analysis and contributing for a more accurate understanding of all the molecular intervenients on RA as well as their dynamic interplay. From this global analysis a more precise and oriented investigation will be possible for the events influencing disease susceptibility, prognosis and therapy response. The database has been - and will be - undergoing intensive updating, and as in any project of this sort, the interaction with the users is fundamental, in order to iterate new functionalities, and correct some aspects of the current implementation which may prove inadequate. We welcome all possible suggestions, comments and criticisms that can improve this resource, for the benefit of all those who can use on in their daily practice.

## Course Report: Bogotá Nov. 2005

The Swiss and the Colombian EMBnet nodes joined to organise a course entitled «**Microarray data analysis**». The main goal was to offer a complete overview of the methods and tools for the analysis of high throughput gene expression data (Microarray, Affymetrix chips, SAGE, MPSS) using the open source statistical package R with BioConductor and other modules. The schedule was as follows:

### Day 1

General introduction to the technologies, scope and experimental design.

Web tools and databases, software used

**Practicals:** R tutorial

### Day 2

Comparison of the techniques (microarray-SAGE/MPSS-PCR, Affy-Agilent)

Scanning, normalization, quality control, annotations of arrays

**Practicals:** normalization, annotation, good/bad quality of slides

### Day 3

Exploratory data analysis, unsupervised clustering and discrimination part 1

**Practicals:** data visualisation, clustering, class discovery

### Day 4

Exploratory data analysis, supervised clustering and discrimination part 2

**Practicals:** data interpretation, Gene Ontology, pathways

### Day 5

Short paper presentations, critical discussion. Summary.

This program was prepared and given by Drs. Thierry Sengstag, Eugenia Migliavacca and Pierre Farmer from the group of Mauro Delorenzi (Swiss Institute of Bioinformatics and National Centre for Competence in Research).

For the first time, several lectures were given via Internet using the Marratech video-meeting system controlled by the EMBnet.org server. Apart from a few network problems on Friday (day 5) the quality of the sound and the images was absolutely perfect. Students and professors were



Figure 1. The Marratech system.

impressed by this e-learning tool (Figure 1).

This system is very convenient, much more ecological and cheaper than a flight to Colombia. But, of course this doesn't replace the pleasure to meet people for a nice Swiss cheese fondue organised by the Swiss team!

About 20 students from Colombia and other Latino-american countries (Costa Rica, Peru, Venezuela) attended the course. We would like to thank the Ibero-american network for funding the travel expenses of those students. On the final day, an official ceremony allowed the students to receive their diploma (Figure 2).

### Links

Web site of the course & file repository

<http://www.co.embnet.org/microarreglos/>  
<ftp://ftp.ch.embnet.org/pub/MAcourseBogota05>

Emiliano Barreto & Laurent Falquet



Figure 2. The diploma ceremony.

# The MitoDrome database and recent development



**Domenica D'Elia, A. Turi, F. Licciulli, D. Catalano & Cecilia Saccone**

Italian EMBnet node, Institute for Biomedical Technologies (ITB), National Research Council, Via Amendola, 122/D, 70126 Bari, Italy

<http://www2.ba.itb.cnr.it/MitoDrome/>

## INTRODUCTION

MitoDrome is a web-based database which provides genomic annotations about nuclear genes of *D. melanogaster* encoding for mitochondrial proteins. The main concept of MitoDrome is to use the *Drosophila* as a model to provide data that can be useful for the investigation of issues related to the cell energy metabolism and mitochondrial biogenesis.

Data in the database are derived from *in silico* analysis relying on Human versus *D. melanogaster* sequence comparison. Here we give a brief rationale of the database and update you on improvements carried out in this last year on both new MitoDrome genome annotations and the web interface. For a more detailed description we suggest you read the article recently published (see reference 1).

## SHORT HISTORY AND NEW DEVELOPMENTS

When the MitoDrome database was launched in 2003 (2), it provided a complete set of nucleus-encoded mitochondrial genes of *D. melanogaster* whose identification was the result of a comparative study carried out using human protein sequences from the SwissProt database, versus the *Drosophila* genomic sequence, ESTs and cDNA data available in the FlyBase database. The new features implemented in 2006 are the

annotation and comparative analysis of OXPHOS genes in the genomes of two completely sequenced species of diptera, *D. pseudoobscura* and *A. gambiae*, the availability of data through a user-friendly web interface and the implementation of a flexible query system and a sequence export tool.

OXPHOS system annotations stored in MitoDrome have been powered by human intervention in the validation and integration of results deriving from a composite comparative study. These are: analyses of sequence homology, structural orthology (conservation of the number of introns, their location in the coding sequence and conservation of the reading frame with respect to the flanking exons), spatial organization of genes in the general structure of the genome and the presence of duplicates (see reference 3). This effort was prompted by the opportunity that the comparison of species at different levels of divergence offers to gain new and interesting insights into the conservative and/or differential evolution of genes and gene regulatory sequences.

## CONTENT

Currently MitoDrome consists of two major sections, namely the "DROME LIST", which contains the annotations of 285 *D. melanogaster* nuclear genes covering all mitochondrial components, and the "OXPHOS system" section. In this section MitoDrome annotates 78 OXPHOS orthologous genes in each species investigated and 47 duplicates (20 in *D. melanogaster*, 19 in *D. pseudoobscura* and eight in *A. gambiae*). For each of the identified genes the following data are available: gene sequence and structure, gene splice-site variants, structure and sequence of expressed transcript, predicted protein sequence, comparison (Clusters) of the gene structure (including splice-site variants, when present), protein sequences alignment (including the human counterpart), phylogenetic tree of orthologous and paralogous-OXPHOS genes in the three species (*D. melanogaster*, *D. pseudoobscura* and *A. gambiae*) and annotation of 43 P-insertion mutant alleles of OXPHOS *D. melanogaster* genes grouped on the basis of their functional classification.

## HOW TO ACCESS DATA AND EXPORT SEQUENCES

The new version of the MitoDrome database is available at the following address: <http://www2.ba.itb.cnr.it/MitoDrome/>. Every data set annotated in MitoDrome can be directly accessed using the related textual hyperlink on the database home page.

### *D. melanogaster* gene collection

*D. melanogaster* genes annotation can be browsed and exported in different ways and formats by accessing the "DROME LIST" page from the database home page. Here you can access the complete list of the *D. melanogaster* genes (as shown in Figure 1) or part of it choosing to search genes on the basis of their functional classification (result shown in Figure 2).

Information on genes is available for consultation in a flat file format. Gene entries and sequences can be extracted through the SRS system implemented on our server clicking on the entry name of the gene of interest listed in the gene table. The section "Documents & Downloads" make it possible to save both the complete list of *D. melanogaster* genes along with related FlyBase

ID and protein description and the flat file of the whole gene collection.

### OXPHOS gene collection

The OXPHOS gene collection can be browsed in two different ways, by precompiled queries and through the search page.

Hyperlinks on the left hand side of the home page (OXPHOS System section) can be used to browse the database for all the OXPHOS genes annotated or for only those of one of the organisms studied (*A. gambiae*, *D. melanogaster*, *D. pseudoobscura*). The query report is shown in Figure 3.

Clicking on the column headers *Organism*, *Entry Name*, *CHR* and *Gene* you can alphabetically order the result of your search for species, MitoDrome entry name, chromosome location or gene name respectively. If you are interested in storing the result of your search you can use the "Save search" option up on the left hand side of the report. This produces an xls file that you can save on your computer and manage for your own purposes. If you are interested in extracting sequences, a clickable button in the upper part of the page gives access to the sequence

**MitoDrome**

HOME SEARCH CLUSTER BLAST DROME LIST CONTACT

OXPHOS System

- All Entries
- A. gambiae*
- D. pseudoobscura*
- D. melanogaster*
- Clusters
- P-Insertion mutant *D. melanogaster*

*D. melanogaster* gene set

[Go to the complete list](#)  
[www2.ba.itb.cnr.it/oxphos/classification](#)

The same data have also been organised in a flat-file format and can be retrieved using the Mitochondrial databases section at our [SRS](#) server

Documents & Downloads:

MitoDrome ID	ENTRY NAME	SW. PROT (HUMAN)	PROTEIN NAME	AA(N) HUM/DRO	M/Sim (%)	Drosophila GENE	FlyBase ID	MAP POSITION	Mutant allele
MITDROME10361	KBL	Q25600	2-AMINO-3-KETOBUTYRATE COENZYME A LIOASE	419417	66/78	CG10361	FBm0036208	68E1	
MITDROME07433	GABT	P80404	4-AMINO-BUTYRATE AMINOTRANSFERASE	500486	52/71	CG7433	FBm0036227	76D8-E1	
MITDROME03017	HEM0	P22527	5-AMINO-LEVULINIC ACID SYNTHASE	587539	59/74	<i>Alt.</i> CG3017	FBm0030764	60B8	
	HEM1	P13126		640539	47/73				
MITDROME10922	THIL	P24752	ACETYL-COA ACETYLTRANSFERASE	427410	65/60	CG10922	FBm0029949	7C1	
MITDROME04600	THIM	R42765	3-KETOACYL-COA THIOLEASE	397598	59/74	<i>RP2</i> CG4600	FBm0040064	30E4	

DAPEG Genetic Unit University of Bari, Italy

Figure 1. *D. melanogaster* gene collection search page and genes list.

## MitoDrome

HOME

SEARCH

CLUSTER

BLAST

DROME

*D. melanogaster* genes Functional Classification

## Oxidative phosphorylation

- [Complex I: NADH ubiquinone oxidoreductase](#)
- [Complex II: Succinate dehydrogenase](#)
- [Complex III: Ubiquinol-cytochrome C oxidoreductase complex](#)
- [Complex IV: Cytochrome c oxidase](#)
- [Complex V: F0/F1 ATP synthase](#)
- Others

## Carbohydrat

- Tric

- Pyr

- Oth

## Amino acid

## Metabolism

- Glyc

- Fat

- Fat

- Oth

## Nucleotide

## Sulfur mat

## Complex I: NADH ubiquinone oxidoreductase

MitoDrome ID	ENTRY NAME	SW-PROT (HUMAN)	PROTEIN NAME	AA%O HUMDRO	M-Sim (%)	Drosophila GENE	FlyBase ID	MAP POSITION
MTDROME08680	NUMM	Q23380	13 KDA-A SUBUNIT	124/126	93/98	CG8680	FBm0031684	2506
MTDROME06463	NUFM	Q16718	13 KDA-B SUBUNIT	115/124	66/61	CG6463	FBm0036100	6782
MTDROME11455	NIFM	Q43920	15 KDA SUBUNIT	105/101	25/97	CG11455	FBm0031228	21B1.2
MTDROME12203	NUYM	Q43181	18 KDA SUBUNIT	175/183	44/62	CG12203	FBm0031101	1827
MTDROME03683	NUPM	E51920	19 KDA SUBUNIT	171/175	44/52	CG3683	FBm0031504	60E1
MTDROME09172	NUKM	Q23251	20 KDA SUBUNIT	213/221	60/73	CG9172	FBm0030718	15A.5
MTDROME02014	NUKM	Q23251	20 KDA SUBUNIT	213/212	64/74	CG2014	FBm0030662	92B1
MTDROME03944	NUMM	Q00212	23 KDA SUBUNIT	210/217	69/75	MD2, CG3944	FBm0012567	92A.5

Figure 2. *D. Melanogaster*: OXPHOS Complex I genes list.

export page, which allows the export of any of the annotated sequences or sub-sequences such as gene, exon, intron sequences as well as transcript and UTRs or protein sequences for the selected entries. To see the content of each one of the entries listed in the query report you have to click on the eye icon in the *View* column. This gives access to the gene entry. An example of a MitoDrome entry is shown in Figure 4.

This picture reproduces only the upper part of the entry, by scrolling down it, you can read information about the protein product as well as use the link to related external resources such as the GO annotation and the SwissProt AC if available. You can get information about CDS, mRNA, UTRs and gene sequences also including data on their structural organization. The UTR regions are shown in blue, codifying regions are displayed in red, upstream and downstream gene sequences are displayed in black as are intron regions. If you are interested in extracting sequences a click-

able button in the Reference Sequence section of the entry gives direct access to the sequence export page. Another feature worth mentioning is that each entry is linked with the entry reporting information about orthologous and paralogous genes identified in the genome of the other organisms annotated in the database (Cluster). You can get information about it by clicking on the "associated cluster" button and access the cluster entry.

## Clusters

Information on Cluster of orthologous and paralogous genes in the three species of diptera compared with their human counterpart and each other is collected in the database section marked "Clusters". You can access this information either using the hyperlink on the left hand side of the home page or by clicking on the CLUSTER page available from the home page navigation tool bar.



The result is a table listing the Clusters stored in the database. Even this list can be alphabetically ordered by clicking on the column headings and the cluster entry accessed by clicking on the eye icon in the View column. The Cluster entry is organized in four different sections. The first section gives general information about the Cluster. A clickable button, "Protein sequence alignment", allows access to the Multialign output file of the Cluster sequences alignment. The second section provides a visual comparison of the gene maps of orthologs and paralogs contained in the Cluster, with the conserved exon/intron junctions indicated by vertical dashed lines. The "TREE" section shows the evolutionary tree derived from the analysis of the alignment of the relevant protein sequences using the ProtML program of the Phylogenetic Analysis Package MOLPHY (<http://ftp.cse.sc.edu/bioinformatics/molphy/>). At the bottom you can find a clickable list of Cluster members (orthologs and duplicates), which gives you access to the entry of each of them, and a clickable button (*extract sequence*) for the export of Cluster sequences and sub-sequences.

## P-Insertion mutant alleles in *D. melanogaster*

This database section provides a list of mutant insertion alleles of *D. melanogaster* genes putatively involved in the OXPHOS system biogenesis and function. It has been compiled mostly using information available from FlyBase and from the BDGP P-Element Gene Disruption Project.

You can access this database section from the textual hyperlink in the left hand side of the MitoDrome home page. As a result the system gives you a table which summarizes information about the number of *D. melanogaster* disrupted genes annotated in the database, classified on the basis of their functional activity and hyperlinked to the relevant sub-group table (see Figure 5).

## SEARCH page and BLAST tool

You can use the database SEARCH page to browse the database making use of different search criteria compared to those used in the precompiled queries and/or to combine different search criteria for a more targeted query. Search fields such as ENTRY NAME, GENE NAME

Select	Organism	Entry Name	CHR	Gene	Protein	Cluster	View
<input checked="" type="checkbox"/>	Anopheles gambiae str. PEST	AG_EG15300	2L	agEG15300	NADH-ubiquinone oxidoreductase 13 kDa-B subunit (NUFM) (By similarity)		
<input checked="" type="checkbox"/>	Anopheles gambiae str. PEST	AG_EG14117	3R	agEG14117	NADH-ubiquinone oxidoreductase 13 kDa-A subunit (NUMM) (By similarity)		
<input checked="" type="checkbox"/>	Anopheles gambiae str. PEST	AG_EG13302	3R	agEG13302	NADH-ubiquinone oxidoreductase 15 kDa subunit (NIPM) (By similarity)		
<input checked="" type="checkbox"/>	Anopheles gambiae str. PEST	AG_EG18085	2L	agEG18085	NADH-UBIQUINONE OXIDOREDUCTASE 18 kDa SUBUNIT (NUYM) (By similarity)		
<input checked="" type="checkbox"/>	Anopheles gambiae str. PEST	AG_EG15210	2L	agEG15210	Ubiquinol-cytochrome C reductase complex 7.2 kDa protein (UCR7) (By similarity)		
<input checked="" type="checkbox"/>	Anopheles gambiae str. PEST	AG_EG19249	2L	agEG19249	NADH-UBIQUINONE OXIDOREDUCTASE 19 kDa SUBUNIT (NUPM) (By similarity)		
<input checked="" type="checkbox"/>	Anopheles gambiae str. PEST	AG_EG16939	X	agEG16939	NADH-ubiquinone oxidoreductase 20 kDa subunit (NUKM) (By similarity)		
<input checked="" type="checkbox"/>	Anopheles gambiae str. PEST	AG_EG12298	2R	agEG12298	NADH-ubiquinone oxidoreductase 20 kDa subunit (NUKM) (By similarity)		
<input checked="" type="checkbox"/>	Anopheles gambiae str. PEST	AG_EG9698	2R	agEG9698	NADH-ubiquinone oxidoreductase 23 kDa subunit (NUKM) (By similarity)		

Figure 3. Query report .

and CHROMOSOME let you browse the database also using a list of search terms. The BLAST tool, hosted on our node server, can be used to browse the database using any type of sequence against the OXPPOS sequence collection. User registration is necessary but it is free of charge and the results of the BLAST job are sent to you by an e-mail.

The SEARCH page along with the BLAST tool are accessible from the MitoDrome home page by clicking on the SEARCH and BLAST buttons, respectively, on the home page navigation tool bar.

## ONGOING WORK AND FUTURE DEVELOPMENTS

The MitoDrome database has been designed to host a whole range of functional products related to mitochondria and to expand the annotation of mitochondrial related genes to an unlimited number of species. Ongoing current work, by the collaborating research group of Prof. Corrado Caggese in the Department of Genetics and Microbiology of Bari University, is focused on annotation of the OXPPOS genes in the ge-

nomics draft sequences of 9 more *Drosophila* species (*D. simulans*, *D. yakuba*, *D. erecta*, *D. ananassae*, *D. willistoni*, *D. persimilis*, *D. mojavensis*, *D. virilis* and *D. grinschawi*) available at the DroSpeGe BLAST server (<http://bugbane.bio.indiana.edu:7151/blast/>). Annotation of the OXPPOS genes of *B. mori* and *A. mellifera* is also under way as genomic sequences from these organisms are made available from the on-going project implemented by the collaborating NGHRI-funded Genome Sequencing Centres (National Human Genome Research Institute: <http://www.genome.gov/10000905>). As soon as the refined data is available it will be annotated into MitoDrome and made public available through its web-interface.

## REMARKS

The MitoDrome database has been developed in collaboration with Prof. Corrado Caggese, Department of Genetics and Microbiology of Bari University, Italy. Any comment, suggestion or proposals for collaboration is welcome. The MitoDrome home page provides a link to the CONTACT page. Here you can find information about

**MitoDrome** user: guest login

HOME SEARCH CLUSTER BLAST DRONE LIST CONTACT

Entry View

**MitoDrome ID: AG\_EG18985**

Organism: *Anopheles gambiae* str. PEST

Gene Name: agEG18985

Product: NADH-UBIQUINONE OXIDOREDUCTASE 18 KDA SUBUNIT (NUYM) (By similarity)

associated cluster

**GENE MAP**

Chromosome: 2L

Celera Gene: agCG52482

Citogenetic Map: 27A

Gene Map

pre-mRNA

agEG18985

189 177 174

100 b

■ = CDS ■ = UTRs

5' → 3'

Comment

This record is derived from the following: gb|AAAB01008807.1  
 gb|EM628145.5, mRNA sequence gb|EM644998.5, mRNA sequence The following contributed to reference sequence development: bases 1..3070 - AAAB01008807.1 7095454..7092385

**PRODUCT**

Figure 4. Gene entry.

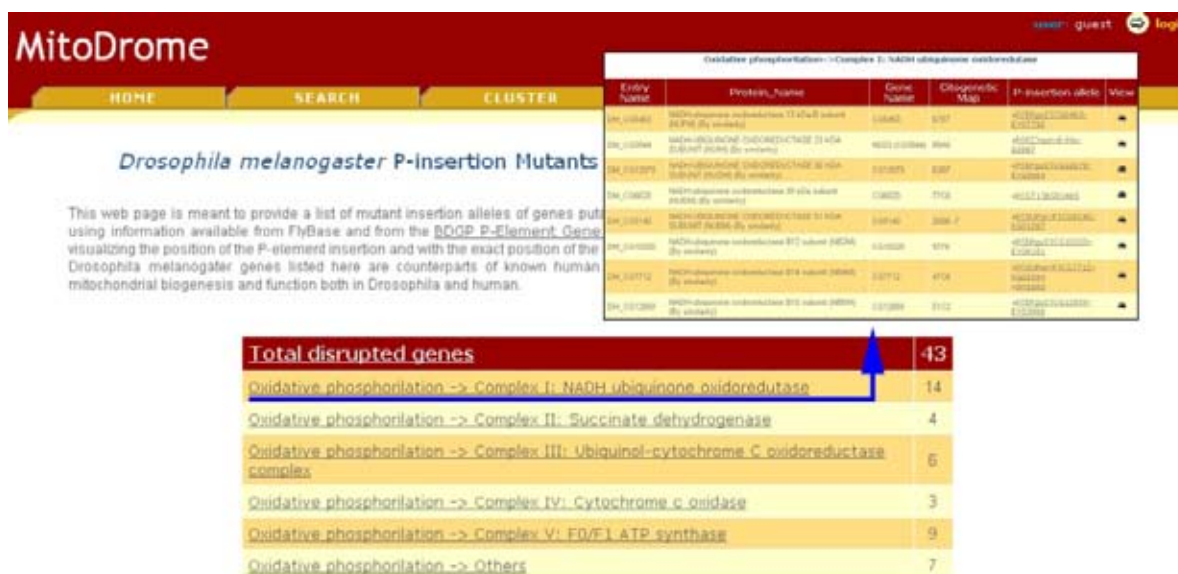


Figure 5. *D. melanogaster* OXPHOS P-insertion mutant alleles list with the detail of a search for those belonging to the OXPHOS Complex I sub-group. Hyperlinks connect each gene with the gene map picture, visualizing the exact position of the P-element as determined by a BLAST-N analysis, and with the FlyBase or BDGP P-Element Gene Disruption Project web sites.

people working on MitoDrome, their role and use phone/e-mail address for contact.

## REFERENCES

- 1) D'Elia D, Catalano D, Licciulli F, Turi A, Tripoli G, Porcelli D, Saccone C, Caggese C. The MitoDrome database annotates and compares the OXPHOS nuclear genes of *Drosophila melanogaster*, *Drosophila pseudoobscura* and *Anopheles gambiae*. *Mitochondrion*, 2006;6(5):252-257.
- 2) Sardiello M, Licciulli F, Catalano D, Attimonelli M, Caggese C: MitoDrome: a database of *Drosophila melanogaster* nuclear genes encoding proteins targeted to the mitochondrion. *Nucleic Acids Res* 2003, 31(1):322-4.
- 3) Tripoli G, D'Elia D, Barsanti P, Caggese C: Comparison of the oxidative phosphorylation (OXPHOS) nuclear genes in the genomes of *Drosophila melanogaster*, *Drosophila pseudoobscura* and *Anopheles gambiae*. *Genome Biol* 2005, 6(2):R11.

## Announcements

The BioMed Central launch a new Journal:

Algorithms for Molecular Biology

It is freely available at

<http://www.almob.org/info/about/>

Algorithms for Molecular Biology is an open access, peer-reviewed journal that encompasses all aspects of algorithms for molecular biology and genomics. Areas of interest include but are not limited to: algorithms for RNA and protein structure analysis, gene prediction and genome analysis, comparative sequence analysis and alignment, phylogeny, gene expression, machine learning, and combinatorial algorithms.

## Mining Ensembl



**Lisa Mullan**

EMBL Outstation - Hinxton,  
European Bioinformatics,  
Institute, Wellcome, Trust  
Genome Campus, Hinxton,  
Cambridge, CB10 1SD,  
United Kingdom

<http://www.ensembl.org>

### About Ensembl

The information available as data in the Ensembl interactive website is also accessible through a database search. BioMart is a database search interface allowing generation of complex queries within the data resource. This can be used to answer complex questions about a specific dataset and the results compared with other datasets.

All data and database and search infrastructure for Ensembl and BioMart is freely available to the entire research community. This practical uses the archived information but can be duplicated on the most recent release.

### Aim

To identify all currently known human protein-coding disease genes that have an essential splice site, more than two transcripts and a homologue in mouse. Both the transcribed 5' and 3' UTRs should be available.

Access <http://www.ensembl.org> and select **v35 Nov 2005** on the left-hand yellow column. Now select **Data mining [BioMart]** from the left hand side of the purple archive interface. Select

Ensembl 35

Homo sapiens genes (NCBI35)

and hit **next ▶**. This set contains 34,294 transcripts. Check and select the following:

GENE:

Disease genes

Transcript count  $\geq 3$

Entries with a 5' UTR only

Entries with a 3' UTR only

Gene type

Status

MULTI SPECIES COMPARISONS:

Homologous Mouse Genes

SNP:

Has essential splice site

and hit **next ▶** to retrieve 18 entries.

**Ensembl Transcript ID** is checked by default. Also check the following:

GENE:

Ensembl Gene ID

Description

Disease OMIM ID

Transcript count

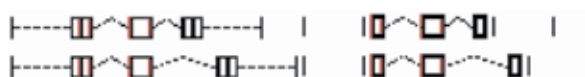
UniProt/Swiss-Prot ID

Disease description ID

and hit **export**. Seventeen genes have been identified that match the criteria set out, and each link accesses the relevant page in Ensembl. The HTML table data is static and data can be exported in a variety of file formats for further processing. Return to the previous page and **Se-**

**lect the Attribute Page:**

From the altered options, select **Peptide**,



check the following **Header Information:**

Chromosome

Coding Start

Ensembl Gene ID

Coding End (Chr bp)

and uncheck anything else. Hit **export** to retrieve the peptide sequences of the different splice transcripts for each of the 18 genes in fasta format.

## GLOSSARY

### Essential Splice Site

This refers to a single nucleotide polymorphism (SNP) present in the first two bases (GT), or the final two bases (AG) of an intron.

### Homologue

Homology indicates an evolutionary relationship between sequences. Divergent, but related sequences in the same genome are called paralogues, and those across species are known orthologues. Orthologues are functional equivalents.

### OMIM

On-line Mendelian Inheritance in Man – an online encyclopaedia the genotypes and resulting diseases and syndromes in Humans.

### Protein-coding genes Genes

containing exons that will eventually be translated proteins. There are 22,218 protein coding gene predictions in Human genome of Ensembl Release 35. They are distinguished from other gene types in the database, such as pseudogenes those coding for a variety of non-coding RNA molecules.

### Transcript

The sequences of exons that have been transcribed from original gene and will generally be translated to a protein. mature transcript is displayed in TransView and has no intervening sequences. Transcripts displayed in other views show intervening sequences as context information only.

### UniProt

The Universal Protein Resource. A protein sequence database split into two sections. The SwissProt section contains curated sequences, the TrEMBL section contains automated annotation only.

### UTR (UnTranslated Region)

The transcribed portion of the gene is not converted to the final protein. UTRs are annotated to the extent supported by EMBL mRNA records. As there is no guarantee that UTR sequences EMBL records are complete there is similarly no guarantee that Ensembl genome analysis and annotation pipeline has enough biological evidence to predict complete UTR regions.

## About Ensembl

Ensembl — a joint project between the European Bioinformatics Institute (EMBL-EBI) and the Wellcome Trust Sanger Institute — provides a framework for working with the genomes of higher animals (metazoans). It presents, via an interactive website, the human genome together with other genomes that are important for addressing complex questions in medical research and molecular biology.

Genome information is taken directly from the relevant sequencing source and is fed into the Ensembl automated annotation pipeline. The Ensembl analysis and annotation pipeline is based on a rule set of heuristics a human annotator would use.). All Ensembl gene predictions are based on experimental evidence, which is imported via manually curated UniProt/Swiss-Prot, partially manually curated NCBI RefSeq, automatically annotated UniProt/TrEMBL records.

## Contact the Ensembl helpdesk at

[http://www.ensembl.org/Homo\\_sapiens/helpview](http://www.ensembl.org/Homo_sapiens/helpview)

---

## Announcement



### Vital-IT Transnational access Programme: an opportunity for european researchers



#### Summary

The Vital-IT high performance facility of the Swiss Institute of Bioinformatics, through funding provided by the EU 6th Framework Programme, is pleased to invite proposals for cost-free use of its Integrated Computational Genomics Resource from individuals, institutions and companies from any of the EU Member and Associated States (except Switzerland).

#### Introduction

Vital-IT is an innovative life science informatics facility providing computational resources, consultancy and training to connect fundamental and applied research. It is a collaboration between the Swiss Institute of Bioinformatics (SIB), the Universities of Lausanne, Geneva and Basel, the Ludwig Institute for Cancer Research, the Swiss Federal Institute of Technology, Lausanne (EPFL), Hewlett Packard Company and Intel Corporation.

The Vital-IT facility comprises three clusters of state-of-the-art 64-bit servers, a large shared memory (64 GB) quad-processor Itanium® server, three dedicated database (SQL) servers, and over 10 TB of shared storage with appropriate backup.

The Integrated Computational Genomics Resource mirrors all the internationally-important genome, transcriptome and proteome sequence repositories, as well as providing access to a range of locally developed specialized databases.

A team of five programmers, with experience in software engineering, parallelization and optimization, in database applications, and in the design of algorithms for life science applications, supports Vital-IT's activities

#### Visiting Developer Programme

Visiting Developers may stay for a period from one week to two months at Vital-IT, with a likely average of one month. Their activities may include the development of new software for HPC applications in life science, parallelization and optimization of existing software for the specific hardware architecture of Vital-IT, and large data analysis projects making use of the rich database collection offered by this facility. They will be provided with office space, free access to all hardware and software resources required by the project, and technical assistance from Vital-IT staff.

Visiting developers will be selected primarily on the basis of a project proposal that will be evaluated by a review panel including external experts in high performance computing, bioinformatics and genomics, including representatives from industry.

#### Remote Access Programme

Remote Access to the HPC infrastructure and computational genomics environment will be provided via a new user-oriented program. Successful applicants will be provided with a user account on Vital-IT and adequate CPU and disk storage quota to carry out the proposed project. This programme is primarily intended for projects that depend on database and software resources that were developed at Vital-IT and cannot easily be ported to another HPC centre.

Requests for remote access to Vital-IT will also be evaluated by the review panel mentioned above. Remote users of the Vital-IT platform should prepare and submit their jobs according to detailed guidelines (available at <http://www.vital-it.ch/vitalit-tech-support.htm>) similar to those applied to existing users. These jobs should be able to run without requiring further assistance from Vital-IT personnel, and will partly be selected according to this criterion.

#### Training Programme

Users, in particular those with limited technical experience, will also be able to attend courses on the technical aspects of the infrastructure, to learn how to take full advantage of it. The training will mainly target new users from European countries. The courses, which will be offered periodically at the Vital-IT facility, will typically extend to 1 or 2 days, and will be open to graduate students, post-doctoral fellows, and more senior researchers.

#### How to Apply

Proposals for both the Visiting Developer and Remote Access Programmes should be submitted using the Application Form (available at <http://www.vital-it.ch/ICGR-Application.htm>). Applications will be reviewed regularly and every effort made to give applicants a timely decision.

#### Website for Further Information

<http://www.vital-it.ch>

# The CMGG methBLAST server hosted by BEN

## Introduction



**Guy Bottu**

Belgian EMBnet Node (BEN), ULB Campus de la Plaine, blv. du Triomphe, 1050 Brussels, Belgium



**David Coornaert**

Eukaryotes have enzymes that are specific for CG sequences in DNA and convert the cytosine into a 5-methylcytosine. Because m5c has an increased transition rate, CG is rarefied except in some regions ("CpG islands"), which are usually located in the promoter region of genes. It is assumed that methylation represses the gene. In cancerous cells as well global hypomethylation as hypermethylation in the promoter region of certain genes has been reported, hence the growing interest in the measurement of CpG methylation levels.

There exist several methods to assess CpG methylation; most involve a treatment of the DNA with sodium bisulphite, which transforms unmethylated cytosine into uracyl, followed by the quantification of modified and unmodified sequences

<http://medgen.ugent.be/methBLAST/>

Figure 1. The methBLAST input page. In the top part you can type in a primer pair or alternatively a series of primers concatenated and separated by a short poly-N track. In the middle part you can choose the genome (human/mouse/rat) and eventually further limit the sequences to be searched by means of the result of a keyword search, which is performed in the NCBI Entrez server.

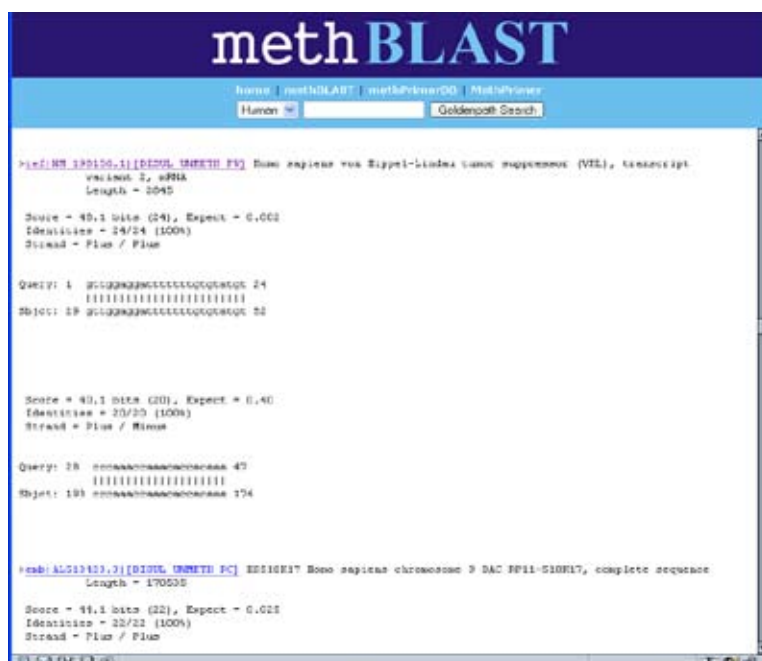


Figure 2. Part of a methBLAST output. The pair of PCR primers has given a "hit" with a sequence corresponding to sulphite modified unmethylated RefSeq:NM\_198156.

using either specific PCR primers, separation by electrophoresis, labelling with radioactive dCTP/dTTP, restriction enzymes or sequencing. These methods involve the use of PCR primers that amplify exclusively methylated and unmethylated DNA or amplify both.

### The CMGG server for methylated DNA specific primers

It is common practise to test whether prospective primers are unique for the DNA region of interest by a BLAST search against a complete genome. The classic genome searching facilities at NCBI, EBI, etc. are however not adequate because the primers are designed to anneal to bisulphite modified DNA, for which the sequence differs significantly from the original DNA. To address this problem the research team of Jo Vandesompele at the Center for Medical Genetics of Ghent University Hospital (Belgium) has set up a server with a database of validated primers (methPrimerDB) and a BLAST server (methBLAST) where you can search against sequence sets that correspond to an in silico bisulphite modified genome. Currently they support only human, mouse and rat. For each genomic set there are sequences labelled BISUL\_METH\_FW, BISUL\_METH\_RC, BISUL\_UNMETH\_FW and BISUL\_UNMETH\_RC. These sets

correspond to unmethylated DNA, where all the C's have been replaced by T's, and methylated DNA, where a C just before a G has been preserved; furthermore, since after bisulphite treatment the two strands are not complementary anymore, as well the forward as the reverse strand is considered.

### The collaboration between CMGG and BEN

Facing a lack of local computing power, the CMGG team contacted the BEN team and a collaboration was started up (at the side of BEN mainly taken care of by D. Coornaert). A system was set up where the methBLAST Web server delegates the BLAST searches to a computer of BEN, on which the databanks are updated with each release of GenBank.

### Reference

Pattyn F., Hoebeeck J., Robbrecht P., Michels E., De Paepe A., Bottu G., Coornaert D., Herzog R., Speleman F. and Vandesompele J. methBLAST and methPrimerDB: web-tools for PCR based methylation analyses of cancer-related genes. BMC bioinformatics, in press.



## Secure Web browsing: How to safely surf the web? (part 2)



**George Magklaras**

Senior Computer Systems Engineer, The Biotechnology Centre of Oslo, University of Oslo

<http://www.biotek.uio.no/>

### Exploring the security settings of Mozilla Firefox

The 'Mozilla Firefox' project represents the most famous Open Source alternative to IE. It's cross platform ability, the well thought interface and the security conscious behaviour of the development team with fast security bug fixes and lack of controversial technologies such as ActiveX make it an ideal web browsing tool. Its browsing settings are simpler than the IE ones. So, let's start exploring them. After launching Firefox, go to the 'Edit' top toolbar menu and start from the 'Privacy' settings on the left as shown in Figure 4.

The 'History', 'Saved Form Information', 'Saved Passwords' and 'Download Manager History' features are designed to control how traceable are your personal web actions. 'History' is a log of the

URLs of all the web sites you have visited for a certain number of days. The 'Saved Passwords' feature is also an important security setting. If you visit a web site where you are required to enter a password to access information, your browser will ask you whether you would like to store that password locally, so that you will not have to type it repeatedly. This is convenient, but not a good choice, when you work from public or unsafe machines. Thus, if you are concerned about leaving saved passwords in unsafe machines, then you can do one of three things listed below, by expanding the 'Saved Passwords' field:

i) You can make sure that the 'Remember Passwords' option is not ticked and clear all the passwords currently stored.

ii) You can select which passwords can be stored by reviewing the saved password list and refusing your browser to store passwords for some of the websites.

iii) You can set up a master password to lock the access to the rest of your passwords. This will not allow someone to view the clear-text form of your passwords unless they know this password.

Option i) is the safest one but the most inconvenient when it comes to scenarios where the user handles a large number of sites with different passwords. Number ii) is error prone and number iii) is by far the best solution that balances convenience and security. Figure 5 illustrates what you need to do, in order to set up a master password.



Figure 4. Firefox 'Preferences'-'Privacy' Menu.

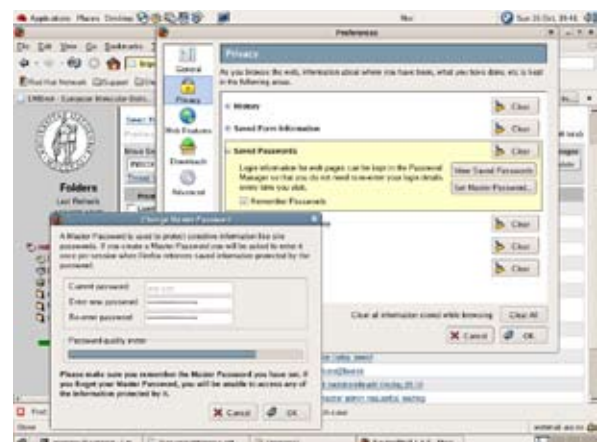


Figure 5. Setting a Firefox 'master password'.

From the expanded 'Saved Passwords' setup menu, click on the 'Set Master Password...' option and you will get to the 'Change Master Password' dialogue box. If you have never set a master password before, the 'Current password' field will be blank, so go ahead and enter (and verify) a new master password. Notice that as you type the new master password, the 'Password quality meter' bar will give you an indication of how easy is to guess your chosen password. This password is important, so don't select 'password' or other common dictionary words as your master password. If your password is too easy to guess there is not really much value in setting one up, so make sure you combine numbers, capital, small case letters and symbols, to make the bar go further to the right and increase the quality of your password. After setting the master password, verify that when you try to review the saved passwords in clear text, you are asked for the master password. In addition, make sure you do not forget the master password and take care of backing your Firefox profile, so disk failures won't make you loose these settings.

Moving on to the Firefox 'Cookies' settings (still on the 'Privacy' settings area), you have mechanisms to control which web sites are allowed to install cookies in your machine. Expand the 'Cookies' settings and see if the option 'Allow sites to set cookies' is ticked. If you do not want sites to install cookies, then you should uncheck this option. Since this may impact the functionality of some web sites, the strategy I follow is to allow only selected web sites to allow cookies, as shown in Figure 6.

I unchecked the 'Allow sites to set cookies' option and then clicked the 'Exceptions' button. Then, in the 'Exceptions' window that pops up, I type the URLs of the sites that are allowed or never allowed to get their cookies in my computer, clicking on the 'Allow' or 'Block' button respectively. An alternative approach would be to leave the 'Allow sites to set cookies' and use the Exceptions button to explicitly define the sites that should be banned (blocked) from storing cookies in my hard drive. Thus, if you have unchecked the 'Allow sites to set cookies' option, it makes sense to define only sites that should be allowed. On the contrary, if you wish to use cookies, the Exceptions should list only the sites you wish to block. Figure 7 illustrates both allowed and block sites

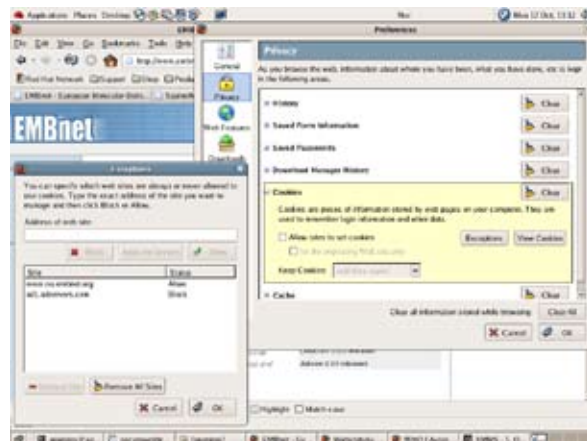


Figure 6. Selective Cookies handling in Firefox.

with the 'Allow sites to set cookies' option just to emphasize the different buttons of the 'Exceptions' window.

After reviewing privacy related features, the 'Web Features' options are concerned with browsing safety. From minor issues, such as the option to block the annoying pop-up windows with adverts to important ones, such as the permission of web sites to install software or execute Java and Javascript code. I normally dislike the option of having all web sites capable of installing programs such as embedded file and program viewers. So, I allow only certain trusted web sites to install software. Bear in mind that the execution of Java and Javascript code via the browser is not normally dangerous, as long as the versions of your browser (for Javascript) and your Java Virtual Machine are up-to-date and free of security bugs. However, if either of these components contains an exploitable bug and you visit a rogue web site, there is no browser setting that can prevent the exposure of your machine. If you disable these features you will almost certainly loose functionality on certain web sites. So, leave them on, but ask your system administrator to verify that you are on the latest and the greatest in terms of the Java and Firefox versions you are running, if you cannot do it yourself. A last tip is to go in the 'Advanced' settings field and make sure that the Security sub-field has the 'Use SSL 2.0', 'Use SSL 3.0' and 'Use TLS 1.0' options checked, to ensure compliance with the appropriate Transport Layer security standards.

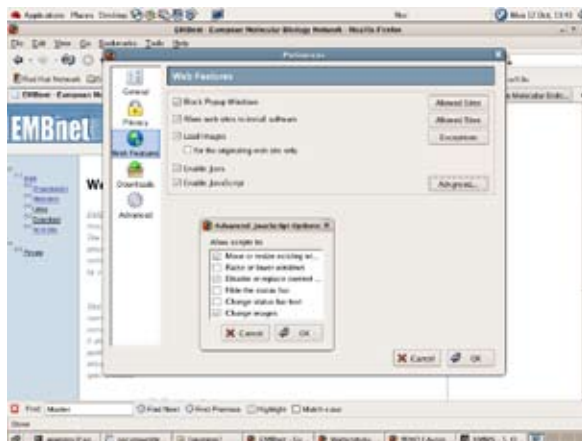


Figure 7. Web Features options in Firefox.

## Conclusions

There is a lot more into security than just the issues we have discussed in this article. However, the knowledge of these basic issues should put you in charge of basic web browsing safety issues. As an epilogue, here is a summary of some basic things you should have in mind to keep your computer safe, in addition to the issues we have discussed:

- 1) Always attend to the operating system updates (or patches) to keep your computer safe.
- 2) Always have your anti-virus signatures updated. This should be done on a daily basis. Laptop users that spend considerable amounts of time off an Internet connection should take note of the urgency to update their virus signatures, the next time they connect to the Net.
- 3) Your applications also need updates. In particular, you should always patch your web, mail and document browsers.
- 4) For your day-to-day computing activities, avoid using an account with elevated privileges, but an ordinary one. For example, Windows users should not use an Administrator account and UNIX-based/Linux users should not use the 'root' account. In that way, you can contain the effects of a system breach.

That's it. Happy (and safe) web surfing!

## Single handed node management



**Pedro Fernandes**

EMBnet node manager for Portugal, Instituto Gulbenkian de Ciência, Oeiras, PT  
([pfern@igc.gulbenkian.pt](mailto:pfern@igc.gulbenkian.pt))

The Portuguese node of the EMBnet has been running at the Instituto Gulbenkian de Ciencia (IGC), since 1992, with a very small group of people (between two and four) managing it. It has maintained a continued service with installed software and maintained databases, together with a continued educational programme, in which more than 700 people have received entry level training in Bioinformatics (one or two courses per year) and specific training in selected subjects (<http://gtpb.igc.gulbenkian.pt/bi-courses>). For three years, a substantial part of



Figure 1. Our new cluster system

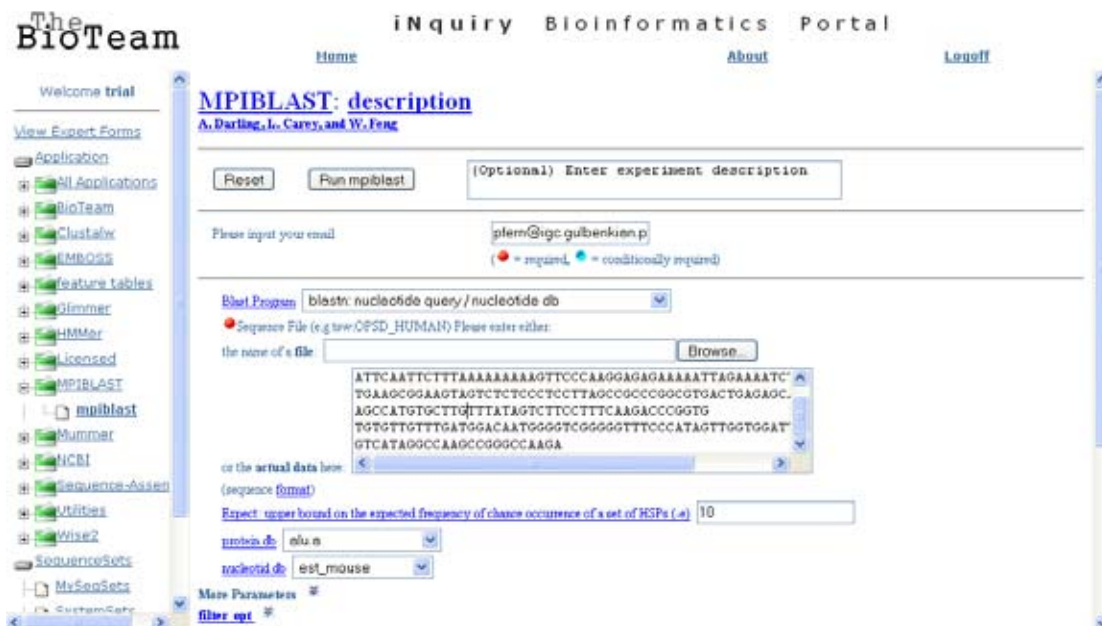


Figure 2. BioTeam iNquiry: mpiblast user interface.

the first Portuguese MSc course in Bioinformatics was also taught, and this represents about 2100 hours of intensive, professional training, of which the first 30 MSc graduates are a vivid result.

Time has come for change again. The team is now reduced to myself and the IGC no longer participates in the management of the MSc, now taken autonomously by the FCUL in Lisboa. Moreover, an old application for funding has been approved and a new machine has been purchased with 90% of public funding, specifically to continue the node activity, as it was considered worth while supporting by the referees of the project.

The new system is now in place, running final pre-production tests. It is a cluster system, based on the IBM e-server technology; the central machine is a dual Xeon and there are 60 high density nodes, based on dual PPC dual core machines arranged in five blade centres; the high availability storage access node is controlled by two dual Xeon machines, with crossed fiber controllers, linked to a 4300 array with 14 disks totaling 4.2TB. The whole set runs Suse SLE9.

The scarce manpower to run it called for an installation that provides the most wanted Bioinformatics methods ready to go in parallelized form. We chose to provide the BioTeam iNquiry suite to

start with. This means that, off the shelf we can offer our users, taking advantage of the parallel architecture, the following software: EMBOSS, hmmer, clustalw, wise, seqio, tgicl, mpiblast, glimmer and ncbi-blast. The BioTeam also provides a self-updating data service that populates our disks with fresh data overnight. This arrangement makes it possible for a single handed node manager to dedicate his time to other things, such as planning and running courses, supporting users and providing advice. The system is rather versatile in providing expansions, as new applications can be added for access through the iNquiry portal by writing rather simple XML scripts. Of course, the system also accepts SSH based requests, and this is bound to be useful in routine mass generated requests, such as the ones required by IGC's new Genotyping Unit, or distributed simulation jobs issued by IGC's Theoretical Biology groups.

Production level operation with this system is starting in full form in September 2006. With this we trust we will continue pleasing our population of about 200 registered national users, plus the local community, of approximately the same size.

Running a system like this is feasible and less of a hurdle compared to the past.

## A grey matter

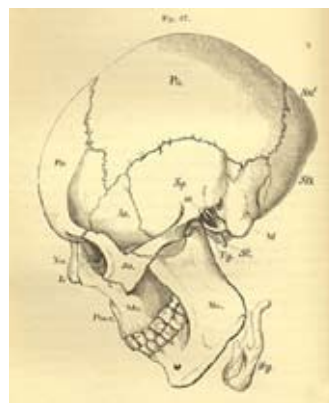
Vivienne Baillie Gerritsen

The human brain has been a hot issue for centuries. Physical anthropology flourished in the 19<sup>th</sup> century and with it the science of craniometry compounded by a growing belief in biological determinism. Intelligence – that intangible quality – was quantified and said to be dependent on brain size. Criminality was based on facial features and cranial particularities. And the notion of racism became a bodily measurement. Thankfully, the 20<sup>th</sup> century offered the necessary wherewithal to tone down all these beliefs thanks to an ever-growing knowledge of the molecular processes going on inside the human body. Intelligence is no longer quantifiable and cannot be defined according to the size of a human brain. Criminality has nothing to do with someone's looks and population genetics have demonstrated that the notion of racism has no real meaning. Despite all this, it is clear that modern humans would not be where they are, were it not for the size of their brain, and its grey matter. And we now know of a number of proteins that are involved in such a process, one of which is a protein known as microcephalin.

The human brain is big. It is bigger than our fellow primates' brains and it is bigger than our ancestors'. And its bigness is no doubt due to forces driven by increasing social complexity rather than mere adaptation to ever-changing geographical or meteorological surroundings. Comparatively, the relative size of our head is not really all that different from a rat's. One of the great differences however is the arrangement of our cerebral cortex – or 'grey matter' – under our protective skull. To cut a long story short, the cerebral cortex houses our neurons and their fibers. Over the millennia, it has become so large that – for spatial reasons – it has had to fold itself into a complex convoluted structure to fit into our skull. Such a structure does not exist in a rat's brain where the cerebral cortex is a smooth structure with no pleats. So what are the molecular processes at the heart of such largeness and convolution?

Diseases which affect the size of the brain can shed some light. Primary microcephaly is one. Primary microcephaly is an affection in which babies are born with a brain far smaller than the average newborn's brain. It remains small throughout their life and those affected show signs of cognitive shortcomings. Interestingly, there is no other kind of affection; any other function orchestrated by the brain is fulfilled normally. This finding has prompted scientists to qualify the condition as atavistic, i.e. people

affected with primary microcephaly have a brain which can be compared to that of our ancestors'.



**Fig.1** The human skull from T. H. Huxley's *Lectures on the Elements of Comparative Anatomy*. (London. John Churchill & sons, 1865)

Microcephalin is clearly involved in defining brain size. Though it is far from the only one. People with primary microcephaly bear a truncated form of the protein. Truncating part of a protein can be quite ruthless – as is the case with pruned microcephalin. The protein loses its function altogether and as a result the brain is not fully developed at birth. It is for this reason that some believe that such a mutation may zap the human brain back in time; to a time before

the emergence of modern humans. But how could a protein have such a drastic effect on our brain?

Two events must have a marked effect on brain size: neuron proliferation and neuron apoptosis. Could microcephalin have something to do with either of these? It seems so. Microcephalin has what is known as BRCT domains. These are domains which are found in proteins belonging to the whole of the animal kingdom and are known to be involved both in protein-protein and DNA interactions. As a result, proteins with BRCTs are likely to be involved in DNA damage and repair mechanisms, and consequently in cell-cycle control leading to cell proliferation or apoptosis. Microcephalin has been shown to have a role in DNA damage response and probably goes about it by regulating a number of other proteins directly involved in DNA repair.

It is not unique to humans. Microcephalin orthologs are found in all other mammals, and some argue that others are also likely to be found in all chordates. What is particularly interesting is that microcephalin has evolved rapidly within the primate lineages. What is more, there is a particular variant of human

microcephalin that seems to have been – and may well still be – under the grasp of positive selection for the last 30'000 years only. This is well after the emergence of modern man, and once modern man had left Africa to discover the rest of the world. Indeed, this particular haplotype is, for the great majority, found in Eurasians. Some have even been so far as to suggest that the Eurasian microcephalin could be concomitant with novel forms of 'art' and symbolism...

So? Microcephalin is most likely involved in brain size but it cannot bear on its own the burden of modern human brain size and what such a notion implies. As a number of researchers have pointed out, though it is largely present in the fetal brain, it is also found in other tissues where it must have some other function. Likewise, there are no doubt many other proteins which also have a say in brain size. Making a direct link between microcephalin and brain size is therefore perhaps a little hasty. However, what is sure is that the protein is involved in DNA repair mechanisms. Certain types of cancer are due to DNA damage and microcephalin may well have a tumor suppressor function – something researchers will be looking into...

### Cross-references to Swiss-Prot

Microcephalin, *Homo sapiens* (Human): Q8NEM0

### References

1. Evans P.D., Gilbert S.L., Mekel-Bobrov N., Vallender E.J., Anderson J.R., Vaez-Azizi L.M., Tishkoff S.A., Hudson R.R., Lahn B.T.  
Microcephalin, a gene regulating brain size, continues to evolve adaptively in humans  
*Science* 309:1717-1720(2005)  
PMID: 16151009
2. Jackson A.P., Ponting C.  
Evolution of primary microcephaly genes and the enlargement of primate brains  
*Curr. Opin. Genet. Dev.* 15:241-248(2005)  
PMID: 15917198
3. Gould S.J.  
'The mismeasure of man'  
W.W. Norton & Company, New York, London, 1981, chap. 3 ('Measuring heads')

## National Nodes

### Argentina

Oscar Grau  
IBBM, Facultad de Cs. Exactas, Universidad Nacional de La Plata  
Email: grau@biol.unlp.edu.ar  
Tel: +54-221-4259223 Fax: +54-221-4259223  
<http://www.ar.embnet.org>

### Australia

Sonia Cattley  
RMC Gunn Building B19, University of Sydney, NSW, 2006  
Email: scattley@angis.org.au  
Tel: +61-2-9531 2948  
<http://www.au.embnet.org>

### Austria

Martin Grabner  
Vienna Bio Center, University of Vienna  
Email: martin.grabner@univie.ac.at  
Tel: +43-1-4277/14141  
<http://www.at.embnet.org>

### Belgium

Robert Herzog, Marc Colet  
BEN ULB Campus Plaine CP 257  
Email: rherzog@ulb.ac.be mcolet@ulb.ac.be  
Tel: +32 2 6505146 Fax: +32 2 6505124  
<http://www.be.embnet.org>

### Brazil

Gonçalo Guimaraes Pereira  
Laboratório de Genômica e Expressão - IB  
UNICAMP-CP 6109  
13083-970 Campinas-SP, BRASIL  
Tel: 0055-19-37886237/6238  
Fax: 0055-19-37886235  
Email: goncalo@unicamp.br  
<http://www.br.embnet.org>

### Chile

Juan A. Asenjo  
Centre for Biochemical Engineering and Biotechnology (ClByB), University of Chile  
Beauchef 861, Santiago, Chile  
Tel: +56 2 6715140  
Fax: +56 2 6991084  
Email: juasenjo@ing.uchile.cl  
<http://www.embnet.cl>

### China

Jingchu Luo  
Centre of Bioinformatics  
Peking University  
Beijing 100871, China  
Tel: 86-10-6275-7281  
Fax: 86-10-6275-9001  
Email: luojc@pku.edu.cn  
<http://www.cn.embnet.org>

### Colombia

Emiliano Barreto Hernández  
Instituto de Biotecnología  
Universidad Nacional de Colombia  
Edificio Manuel Ancizar  
Bogota - Colombia  
Tel: +571 3165027 Fax: +571 3165415  
Email : ebarreto@ibun.unal.edu.co  
<http://www.co.embnet.org>

### Costa Rica

Allan Orozco  
University of Costa Rica (UCR), School of Medicine,  
Department of Pharmacology and ClinicToxicology  
San Jose, America Central  
Costa Rica  
Email: allanorozco@gmail.com  
Tel: +506 2074489  
<http://www.dftc.ucr.ac.cr/>

### Cuba

Ricardo Bringas  
Centro de Ingeniería Genética y Biotecnología,  
La Habana, Cuba  
Email: bringas@cigb.edu.cu  
Tel: +53 7 218200  
<http://www.cu.embnet.org>

### Finland

Kimmo Mattila  
CSC, Espoo  
Email: kimmo.mattila@csc.fi  
Tel: +358 9 4572708  
Fax: +358 9 4572302  
<http://www.fi.embnet.org>

### France

Jean-Marc Plaza  
INFOBIOGEN, Evry  
Email: plaza@infobiogen.fr  
Tel: +33 1 60 87 37 11 Fax: +33 1 60 87 37 96  
<http://www.fr.embnet.org>

### Hungary

Endre Barta  
Agricultural Biotechnology Center  
Szent-Gyorgyi A. ut 4. Godollo,  
Email: barta@abc.hu  
Tel: +36 30-2101795  
<http://www.hu.embnet.org>

### India

Akash Ranjan  
Laboratory of Computational Biology & Bioinformatics  
facility, Centre for DNA Fingerprinting and Diagnostics  
(CDFD), Hyderabad  
Email: akash@cdfd.org.in  
Tel: +91 40 7155607 / 7151344 ext:1206  
Fax : +9140 7155479  
<http://www.in.embnet.org>

### Israel

Leon Esterman  
INN (Israeli National Node) Weizmann Institute of Science  
Department of Biological Services, Biological Computing  
Unit, Rehovot  
Email: Leon.Esterman@weizmann.ac.il  
Tel: +972- 8-934 3456  
<http://www.il.embnet.org>

### Italy

Cecilia Saccone  
CNR - Institute of Biomedical Technologies  
Bioinformatics and Genomic Group  
Via Amendola 168/5 - 70126 Bari (Italy)  
Email: saccone@area.ba.cnr.it  
Tel. +39-80-5482100 - Fax. +39-80-5482607  
<http://www.it.embnet.org>

### Mexico

Cesar Bonavides  
Nodo Nacional EMBnet, Centro de Investigación sobre  
Fijación de Nitrógeno, Cuernavaca, Morelos  
Email: embnetmx@cifn.unam.mx  
Tel: +52 (7) 3 132063  
<http://embnet.cifn.unam.mx>

### The Netherlands

Jack A.M. Leunissen  
Dept. of Genome Informatics  
Wageningen UR, Dreijenlaan 3  
6703 HA Wageningen, NL  
Email: Jack.Leunissen@wur.nl  
Tel: +31 317 484074  
<http://www.nl.embnet.org>

### Norway

George Magklaras  
The Norwegian EMBnet Node  
The Biotechnology Centre of Oslo  
Email: admin@embnet.uio.no  
Tel: +47 22 84 0535  
<http://www.no.embnet.org>

### Pakistan

Raheel Qamar  
Department of Biosciences, COMSATS Institute of  
Information Technology, Park Road, Chak Shahzaad  
Campus, Chak Shahzaad  
Islamabad, Pakistan  
Email: Raheel\_qamar@comsats.edu.pk  
Tel: +0092-333-5119494  
[http://www.ciiit.edu.pk/Departments\\_&\\_Faculties/Link=Detail&f=Departments%5F%26%5FFaculties&SMID=10](http://www.ciiit.edu.pk/Departments_&_Faculties/Link=Detail&f=Departments%5F%26%5FFaculties&SMID=10)

### Poland

Piotr Zielenkiwicz  
Institute of Biochemistry and Biophysics  
Polish Academy of Sciences Warszawa  
Email: piotr@pl.embnet.org  
Tel: +48-22 86584703  
<http://www.pl.embnet.org>

### Portugal

Pedro Fernandes  
Instituto Gulbenkian de Ciencia  
Unidade de Bioinformatica  
2781-901 OEIRAS  
Email: pfern@igc.gulbenkian.pt  
Tel: +351 214407912 Fax: +351 2144079070  
<http://www.pt.embnet.org>

### Russia

Sergei Spirin  
Biocomputing Group, Belozersky Institute Moscow  
Email: sas@belozersky.msu.ru  
Tel: +7-095-9395414  
<http://www.genebee.msu.ru>

### Slovakia

Lubos Klucar  
Institute of Molecular Biology SAS Bratislava  
Email: klucar@embnet.sk  
Tel: +421 2 5930 7413  
<http://www.sk.embnet.org>

### South Africa

Ruediger Braeuning  
SANBI, University of the Western Cape, Bellville  
Email: ruediger@sanbi.ac.za  
Tel: +27 (0)21 9593645  
<http://www.za.embnet.org>

### Spain

José M. Carazo, José R. Valverde  
EMBnet/CNB, Centro Nacional de Biotecnología, Madrid  
Email: carazo@es.embnet.org,  
jrvalverde@es.embnet.org  
Tel: +34 915 854 505 Fax: +34 915 854 506  
<http://www.es.embnet.org>

### Sweden

Nils-Einar Eriksson, Erik Bongcam-Rudloff  
Uppsala Biomedical Centre, Computing Department,  
Uppsala, Sweden  
Email: nils-einar.eriksson@bmc.uu.se  
erik.bongcam@bmc.uu.se  
Tel: +46-(0)18-4714017, +46-(0)18-4714525  
<http://www.embnet.se>

### Switzerland

Laurent Falquet  
Swiss Institute of Bioinformatics, Génopode-UNIL, CH-1015  
Lausanne Email: Laurent.Falquet@isb-sib.ch  
Tel: +4121 692 4078 Fax: +4121 692 4065  
<http://www.ch.embnet.org>



## Specialist Nodes

### EBI

Rodrigo López  
EBI Embl Outstation, Wellcome trust Genome Campus,  
Hinxton Hall, Hinxton, Cambridge, United Kingdom  
Email: rls@ebi.ac.uk  
Phone: +44 (0)1223 494423  
<http://www.ebi.ac.uk>

### ETI

P.O. Box 94766  
NL-1090 GT Amsterdam, The Netherlands  
Email: wouter@eti.uva.nl  
Phone: +31-20-5257239  
Fax: +31-20-5257238  
<http://www.eti.uva.nl>

### ICGEB

Sándor Pongor  
International Centre for Genetic Engineering and  
Biotechnology  
AREA Science Park, Trieste, ITALY  
Email: pongor@icgeb.trieste.it  
Phone: +39 040 3757300  
<http://www.icgeb.trieste.it>

### IHCP

William Moens  
Institute of Health and Consumer Protection  
Via E. Fermi 1 - 21020 Ispra (Varese), Italy  
Email: william.moens@jrc.it  
Phone: +390332786481  
<http://ihcp.jrc.cec.eu.int/>

### ILRI/BECA

Etienne deVilliers  
International Livestock Research Institute  
PO Box 30709, Nairobi 00100, Kenya  
Email: e.villiers@cgiar.org  
Phone: +254 20 4223000  
[www.becabioinfo.org](http://www.becabioinfo.org)

### LION Bioscience

Thure Etzold  
LION Bioscience AG, Heidelberg, Germany  
Email: Thure.Etzold@uk.lionbioscience.com  
Phone: +44 1223 224700  
<http://www.lionbioscience.com>

### MIPS

H. Werner Mewes  
Email: mewes@mips.embnet.org  
Phone: +49-89-8578 2656  
Fax: +49-89-8578 2655  
<http://www.mips.biochem.mpg.de>

### UMBER

Terri Attwood  
School of Biological Sciences, The University of Manchester,  
Oxford Road, Manchester M13 9PT, UK  
Email: attwood@bioinf.man.ac.uk  
Phone: +44 (0)61 275 5766  
Fax: +44 (0) 61 275 5082  
<http://www.bioinf.man.ac.uk/dbbrowser>



EMBnet.news  
ISSN 1023-4144

Dear reader,

If you have any comments or suggestions regarding this newsletter we would be very glad to hear from you. If you have a tip you feel we can print then please let us know. Please send your contributions to one of the editors. You may also submit material by e-mail.

Past issues of EMBnet.news are available as PostScript or PDF files. You can get them from the EMBnet organisation Web site:

<http://www.embnet.org/download/embnetnews>

### Publisher:

EMBnet Executive Board  
c/o Erik Bongcam-Rudloff  
Uppsala Biomedical Centre  
The Linnaeus Centre for Bioinformatics, SLU/UU  
Box 570 S-751 23 Uppsala, Sweden  
Email: erik.bongcam@bmc.uu.se  
Tel: +46-18-4716696

Submission deadline for the next issues:

December 10, 2006

EMBnet.news is an official publication of the EMBnet organisation  
[www.embnet.org](http://www.embnet.org)