



THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

### BioMart Central Portal

**Citation for published version:**

Guberman, JM, Ai, J, Arnaiz, O, Baran, J, Blake, A, Baldock, R, Chelala, C, Croft, D, Cros, A, Cutts, RJ, Di Génova, A, Forbes, S, Fujisawa, T, Gadaleta, E, Goodstein, DM, Gundem, G, Haggarty, B, Haider, S, Hall, M, Harris, T, Haw, R, Hu, S, Hubbard, S, Hsu, J, Iyer, V, Jones, P, Katayama, T, Kinsella, R, Kong, L, Lawson, D, Liang, Y, Lopez-Bigas, N, Luo, J, Lush, M, Mason, J, Moreews, F, Ndegwa, N, Oakley, D, Perez-Llamas, C, Primig, M, Rivkin, E, Rosanoff, S, Shepherd, R, Simon, R, Skarnes, B, Smedley, D, Sperling, L, Spooner, W, Stevenson, P, Stone, K, Teague, J, Wang, J, Wang, J, Whitty, B, Wong, DT, Wong-Erasmus, M, Yao, L, Youens-Clark, K, Yung, C, Zhang, J & Kasprzyk, A 2011, 'BioMart Central Portal: an open database network for the biological community', *Database: The Journal of Biological Databases and Curation*, vol. 2011, pp. bar041. <https://doi.org/10.1093/database/bar041>

**Digital Object Identifier (DOI):**

[10.1093/database/bar041](https://doi.org/10.1093/database/bar041)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Publisher's PDF, also known as Version of record

**Published In:**

Database: The Journal of Biological Databases and Curation

**Publisher Rights Statement:**

This article is Open Access

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



## Original article

# BioMart Central Portal: an open database network for the biological community

Jonathan M. Guberman<sup>1</sup>, J. Ai<sup>2</sup>, O. Arnaiz<sup>3</sup>, Joachim Baran<sup>1</sup>, Andrew Blake<sup>4</sup>, Richard Baldock<sup>5</sup>, Claude Chelala<sup>6</sup>, David Croft<sup>7</sup>, Anthony Cros<sup>1</sup>, Rosalind J. Cutts<sup>6</sup>, A. Di Génova<sup>8</sup>, Simon Forbes<sup>9</sup>, T. Fujisawa<sup>10</sup>, E. Gadaleta<sup>6</sup>, D. M. Goodstein<sup>11,12</sup>, Gunes Gundem<sup>13</sup>, Bernard Haggarty<sup>5</sup>, Syed Haider<sup>14</sup>, Matthew Hall<sup>15</sup>, Todd Harris<sup>16</sup>, Robin Haw<sup>1</sup>, S. Hu<sup>2</sup>, Simon Hubbard<sup>17</sup>, Jack Hsu<sup>1</sup>, Vivek Iyer<sup>18</sup>, Philip Jones<sup>7</sup>, Toshiaki Katayama<sup>19</sup>, R. Kinsella<sup>7</sup>, Lei Kong<sup>20</sup>, Daniel Lawson<sup>21</sup>, Yong Liang<sup>1</sup>, Nuria Lopez-Bigas<sup>13</sup>, J. Luo<sup>7</sup>, Michael Lush<sup>22</sup>, Jeremy Mason<sup>15</sup>, Francois Moreews<sup>23</sup>, Nelson Ndegwa<sup>7</sup>, Darren Oakley<sup>18</sup>, Christian Perez-Llamas<sup>13</sup>, Michael Primig<sup>24</sup>, Elena Rivkin<sup>1</sup>, S. Rosanoff<sup>7</sup>, Rebecca Shepherd<sup>9</sup>, Reinhard Simon<sup>25</sup>, B. Skarnes<sup>18</sup>, Damian Smedley<sup>7</sup>, Linda Sperling<sup>3</sup>, William Spooner<sup>16,26</sup>, Peter Stevenson<sup>5</sup>, Kevin Stone<sup>15</sup>, J. Teague<sup>9</sup>, Jun Wang<sup>20</sup>, Jianxin Wang<sup>1</sup>, Brett Whitty<sup>1</sup>, D. T. Wong<sup>2</sup>, Marie Wong-Erasmus<sup>1</sup>, L. Yao<sup>1</sup>, Ken Youens-Clark<sup>16</sup>, Christina Yung<sup>1</sup>, Junjun Zhang<sup>1</sup> and Arek Kasprzyk<sup>1,\*</sup>

<sup>1</sup>Ontario Institute for Cancer Research, Toronto, M5G 0A3, Canada, <sup>2</sup>School of Dentistry and Dental Research Institute, University of California Los Angeles (UCLA), Los Angeles, 90095-1668, USA, <sup>3</sup>Centre de Génétique Moléculaire UPR3404, Centre National de la Recherche Scientifique (CNRS), F-75794 Paris cedex 16 Paris, France, <sup>4</sup>MRC Harwell, Harwell Science and Innovation Campus, Oxfordshire, OX11 0RD, UK, <sup>5</sup>MRC Human Genetics Unit, Institute of Genetics and Molecular Medicine, Western General Hospital, Edinburgh, EH4 2XU, UK, <sup>6</sup>Centre for Molecular Oncology and Imaging, Barts Cancer Institute, Queen Mary University of London, E1 2AD, UK, <sup>7</sup>EMBL-European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge, CB10 1SD, UK, <sup>8</sup>Center for Mathematical Modeling and Center for Genome Regulation, University of Chile, Blanco Enclada, Chile, <sup>9</sup>Cancer Genome Project, Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Cambridge, CB10 1SA, UK, <sup>10</sup>Kasuz DNA Research Institute, Chiba, 292-0818, Japan, <sup>11</sup>Department of Energy, Joint Genome Institute, Walnut Creek, USA, <sup>12</sup>Genomics Division, Lawrence Berkeley National Laboratory, Berkeley, 94720, USA, <sup>13</sup>Research Unit on Biomedical Informatics, Department of Experimental and Health Sciences, Universitat Pompeu Fabra, Barcelona, 08003, Spain, <sup>14</sup>Computer Laboratory, University of Cambridge, Cambridge, CB3 0FD, UK, <sup>15</sup>Mouse Genomic Informatics Group, The Jackson Laboratory, Bar Harbor, 04609, USA, <sup>16</sup>Cold Spring Harbor Laboratory, Cold Spring Harbor, 11724, USA, <sup>17</sup>Faculty of Life Sciences, The University of Manchester, Manchester, M13 9PL, UK, <sup>18</sup>Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, CB10 1SD, UK, <sup>19</sup>University of Tokyo, Tokyo, Japan, <sup>20</sup>University of Tokyo 7-3-1, Hongo, Bunkyo-ku, Tokyo, Japan 113-0033, <sup>21</sup>Center for Bioinformatics, State Key Laboratory of Protein and Plant Gene Research, College of Life Sciences, Peking University, Beijing, 100871, P.R. China, <sup>22</sup>VectorBase, European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, CB10 1SD, UK, <sup>23</sup>HUGO Gene Nomenclature Committee (HGNC), European Bioinformatics Institute (EMBL-EBI) Wellcome Trust Genome Campus, Hinxton, CB10 1SD, UK, <sup>24</sup>Sigenae, Institut National de la Recherche Agronomique (INRA), St-Gilles, 75338, France, <sup>25</sup>Inserm U625, University of Rennes, Rennes, 35043, France, <sup>26</sup>International Potato Center (CIP), Lima, 1558, Peru and <sup>26</sup>Eagle Genomics Ltd., Babraham Research Campus, Cambridge, CB22 3AT, UK

\*Corresponding author: Tel: +647 258 4321; Fax: 647 258 4321; Email: arek.kasprzyk@gmail.com

Submitted 16 May 2011; Revised 8 August 2011; Accepted 10 August 2011

BioMart Central Portal is a first of its kind, community-driven effort to provide unified access to dozens of biological databases spanning genomics, proteomics, model organisms, cancer data, ontology information and more. Anybody can contribute an independently maintained resource to the Central Portal, allowing it to be exposed to and shared with the research community, and linking it with the other resources in the portal. Users can take advantage of the common interface to quickly utilize different sources without learning a new system for each. The system also simplifies cross-database searches that might otherwise require several complicated steps. Several integrated tools streamline common tasks, such as converting between ID formats and retrieving sequences. The combination of a wide variety of databases, an easy-to-use interface, robust programmatic access and the array of tools make Central Portal a one-stop shop for biological data querying. Here, we describe the structure of Central Portal and show example queries to demonstrate its capabilities.

**Database URL:** <http://central.biomart.org>.

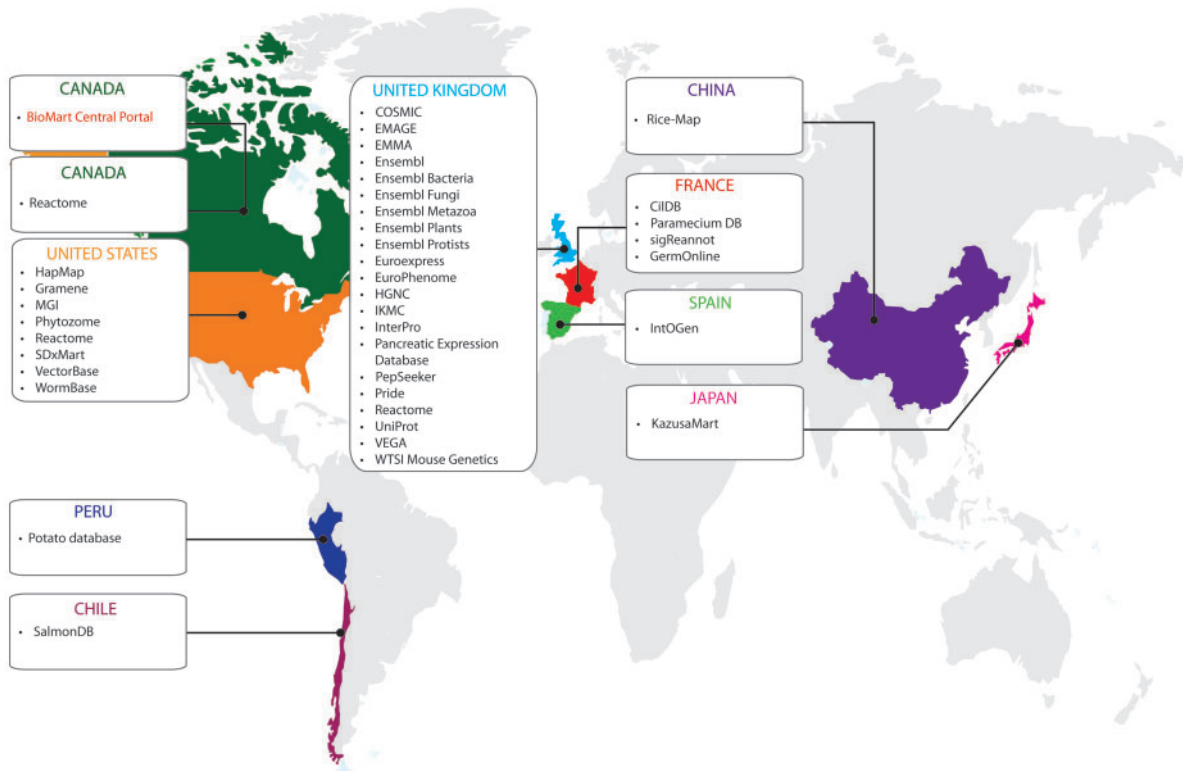


Figure 1. Databases available on the BioMart Central Portal and their host countries (April 2011).

## Project description

### Introduction

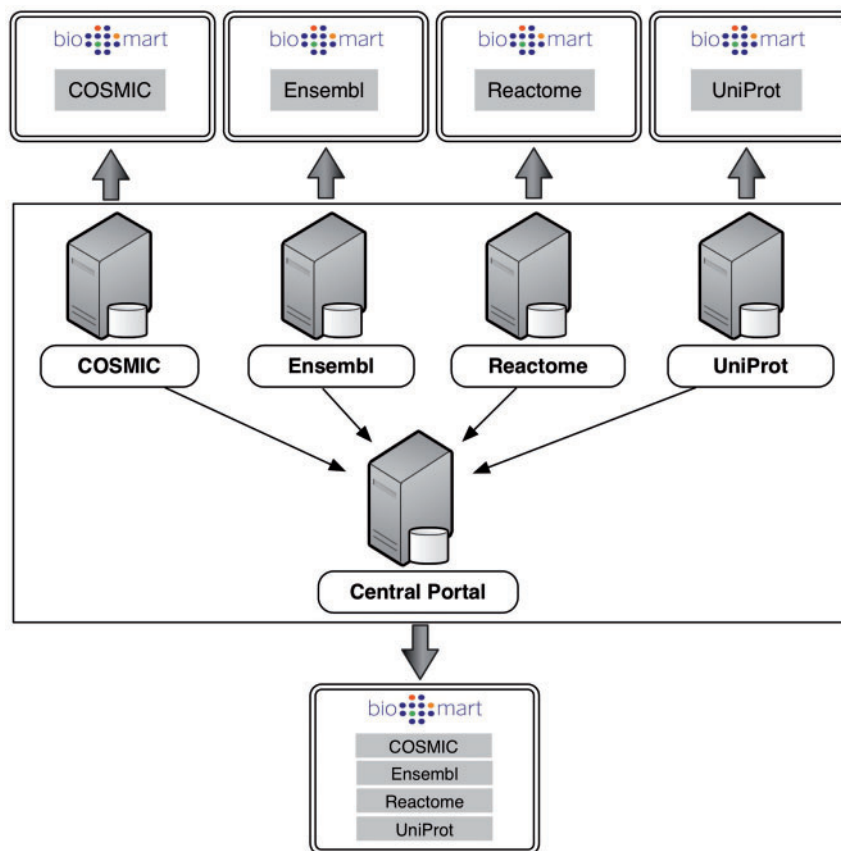
BioMart is a free, open-source, federated database system (1–3). It is cross-platform and supports many popular relational database managements systems, including MySQL, Oracle, PostgreSQL, SQL Server and DB2. The software is data-agnostic, and can therefore be easily adapted to existing data sets. It is expandable and customizable through a plug-in system, and is open-source so the community can participate in deeper development. Furthermore, BioMart can seamlessly connect geographically disparate databases, facilitating collaboration between different groups. These features have catalyzed the creation of BioMart Central Portal, a first of its kind community-supported effort to create a single access point integrating many different, independently administered biological databases (Figure 1).

For administrators, participation in Central Portal offers several benefits. Central Portal can provide an instantly available and automatically updated source of annotations for other projects, as is done in the International Cancer Genome Consortium Data Portal (4). Being part of the community can also expose a database to a wide user base. Furthermore, because the BioMart software allows

administrators to easily create their own plug-ins, joining the community allows administrators to take advantage of the tools that others have created, thereby enhancing their own databases. Central Portal passes queries directly to the individual member servers, so administrators retain full control of their databases and their data (Figure 2).

For users, Central Portal offers a central repository for a vast array of biological data. BioMart can interoperate with other web sites, because results can be configured to link to outside resources; examples in Central Portal include KEGG pathway information (5–7) and Pancreatic Expression Database entries (8). The intuitive interface is consistent across all databases, so users familiar with one source can immediately transfer their skills to another data source. Since Central Portal is constantly updated, users are immediately exposed to new resources as they become available. In addition to the web-based interface, Central Portal also offers a wide variety of other access methods for more advanced querying, including application programming interfaces (APIs) for Java, SPARQL, REST and SOAP.

Moreover, both users and administrators benefit from the value gained by having individual databases connected in a central access point. By allowing data sets to be linked together, resources can be combined in novel ways,



**Figure 2.** Each individual server hosts its own instance of BioMart retrieving data from its own local database backend. Central Portal offers a unified access point to all of these databases, distributing queries to the appropriate servers.

potentially revealing unexpected connections or suggesting new avenues of inquiry. The strength of the Central Portal comes from the fact that it is created and supported by a large community, and, as a whole, it is greater than the sum of its parts.

### Interface

When viewing the Central Portal home page, users are presented with the main querying section, which is divided into three subsections: Identifier Search, Tools and Database Search (Figure 3).

The Identifier Search (Figure 3A) allows users to input gene identifiers in a number of formats (e.g. Gene name, Ensembl IDs, RefSeq IDs, etc.) and search for it across all of the member databases in the Portal. The result of the search links to a report page for the identifier, which summarizes key information about the search term taken from several sources (Figure 4). With this function users can quickly find information about a single identifier, and perhaps even locate resources that they did not realize were applicable to the target of their query.

The Tools section (Figure 3B) contains links to various data analysis tools in four categories: Gene retrieval, Variant retrieval, Sequence retrieval and ID Converter. The first two sections allow quick access to some of the largest and most popular databases contained in Central Portal. The third section, Sequence retrieval, allows easy querying of genomic and protein sequences in any of several formats (Figure 5). The fourth section, the ID Converter tool, allows users to enter or upload a list of identifiers in any format supported by a BioMart database, and retrieve the same list converted to any other supported format.

In the Database Search section (Figure 3C), users can access the individual member databases for querying through the BioMart interface. To make finding the relevant database easier, users can choose to browse databases by the type of information contained therein (Search by type) or by the organism with which the database is concerned (Search by organism). Browse by type is further subdivided into several categories such as Genome [e.g. Ensembl databases (9)], Gene annotation [e.g. HGNC (10)], Protein sequence and structure [e.g. InterPro (11)], Interactions and pathways [e.g. Reactome (12)], Gene

expression [e.g. EMAGE (13)], Cancer [e.g. COSMIC (14)] and Model organism databases [e.g. Gramene (15)], Search by organism is subdivided into categories for bacteria, plants, protists, invertebrates and vertebrates. After choosing a data set, users can construct queries using the basic BioMart concepts of attributes, which indicate what information should be returned, and filters, which restrict the database entries that are retrieved.

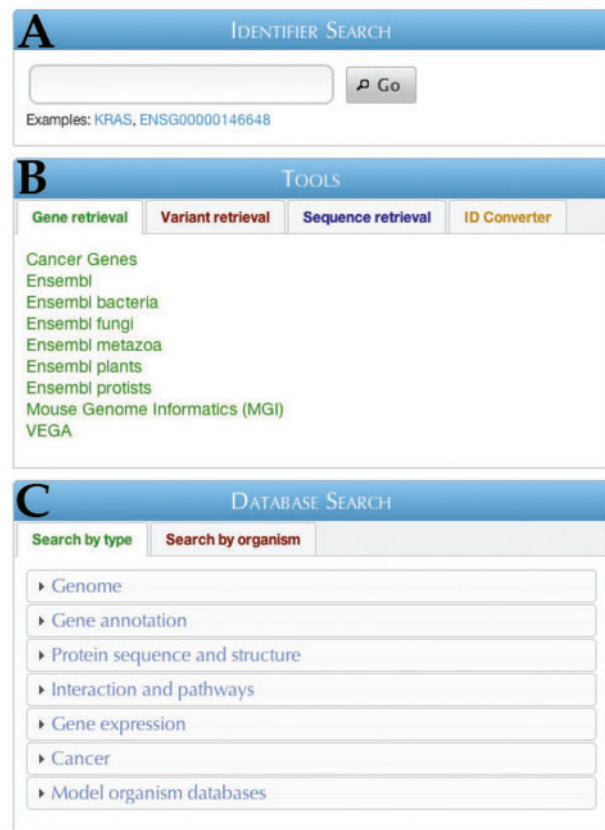
### Access methods

In addition to the graphical user interfaces, Central Portal also offers programmatic access to allow for automated querying. Several programming interfaces are available: an XML querying method that can be accessed via REST or SOAP requests, a full Java API and RDF querying via SPARQL. The syntax of any of the APIs is easy to use for programmers familiar with the basic BioMart concepts of attributes, filters and data sets. For example, to retrieve a list of filters for a given data set, a client could use the REST API and access the URL `/martservice/filters?datasets=datasetname`. Alternatively and equivalently, the client could use the Java API using the method `getFilters(datasetname)` to accomplish the same result. Because, there are a variety of APIs available, developers can choose the access method that makes the most sense for their specific applications and use cases.

To further ease the adoption of the APIs, the equivalent code of any query constructed in the web GUI can be retrieved in any of the API formats by clicking on the appropriate button on the query page; in this way, queries can be saved, modified and easily transferred from one format to another. It also provides a readily available graphical method of constructing complex API calls, which could be of use in certain tools or scripts.

## Data content

BioMart Central Portal contains a constantly growing list of data sources accessible by a wide variety of methods and tools. The following table reflects the contents of the portal as of May 2011:



**Figure 3.** The BioMart Central Portal home page. Three main entry points are available: (A) Identifier search, (B) Tools and (C) Database search.

Database	Location	Description	References
Cildb	CNRS, France	Database for eukaryotic cilia and centriolar structures, integrating orthology relationships for 33 species with high-throughput studies and OMIM	(16)
COSMIC	WTSI, UK	Somatic mutation information relating to human cancers	(14)
EMAGE	MRC HGU, UK	<i>In situ</i> gene expression data in the mouse embryo	(13)
EMMA	EBI, UK	Mouse mutant strain information	(17)
Ensembl	WTSI/EBI, UK	Genome databases for vertebrates and other eukaryotic species	(9)
Ensembl Bacteria	EBI, UK	Genome databases for bacteria	(9)
Ensembl Fungi	EBI, UK	Genome databases for fungi	(9)
Ensembl Metazoa	EBI, UK	Genome databases for metazoa	(9)

(continued)

Database	Location	Description	References
Ensembl Plants	EBI, UK	Genome databases for plants	(9)
Ensembl Protists	EBI, UK	Genome databases for protists	(9)
Eurexpress	MRC HGU, UK	Transcriptome atlas database for mouse embryo	(18)
EuroPhenome	MRC Harwell, UK	Mouse phenotyping data	(19)
GermOnline	Inserm, France	Cross-species microarray expression database focusing on germline development, meiosis and gametogenesis as well as the mitotic cell cycle	(20)
Gramene	CSHL, USA	Agriculturally important grass genomes	(15)
HapMap	NCBI, USA	Multi-country effort to identify and catalog genetic similarities and differences in human beings	(21)
HGNC	EBI, UK	Repository of human gene nomenclature and associated resources	(10)
IKMC	WTSI, UK	Data on mutant products (mice, ES cells and vectors) generated and made available by members of the International Knockout Mouse Consortium	(22)
InterPro	EBI, UK	Integrated database of predictive protein 'signatures' used for the classification and automatic annotation of proteins and genomes	(11)
IntOGen	UPF, Spain	Integrated multi-dimensional data for the identification of genes and groups of genes involved in cancer development	(23)
KazusaMart	Kazusa, Japan	Cyanobase, rhizobia and plant genome databases	(24)
MGI	Jackson Laboratory, USA	Mouse genome features, locations, alleles and orthologues	(25)
Pancreatic Expression Database	Barts Cancer Institute, UK	Results from published pancreatic cancer papers	(8)
Paramecium DB	CNRS, France	Paramecium genome database	(26)
PepSeeker	University of Manchester, UK	Database of proteome peptide identifications for investigating fragmentation patterns	(27)
Phytozome	JGI/CIG, USA	Comparative genomics of green plants	(28)
Potato Database	CIP, Peru	Potato and sweet potato phenotypic and genomic information	(29)
PRIDE	EBI, UK	Repository for protein and peptide identifications	(30)
Reactome	OICR, Canada; EBI, UK; NYU Medical Center, USA	Curated pathway annotation database	(12)
Rice-Map	Peking University, China	Rice ( <i>japonica</i> and <i>indica</i> ) genome annotation database	(31)
SalmonDB	CMM, Chile	Genomic information for Atlantic salmon, rainbow trout and related species	(32)
SDxMart	UCLA, USA	Saliva diagnostics for high-impact human diseases	(33)
sigReannot	Rennes, France	Aquaculture and farm animal species EST contigs	(34)
UniProt	EBI, UK	Protein sequence and functional information	(35)
VectorBase	University of Notre Dame, USA	Genome information for invertebrate vectors of human pathogens	(36)
VEGA	WTSI, UK	Manual annotation of vertebrate genome sequences	(37)
WormBase	California Institute of Technology, USA; CSHL, USA; EBI, UK; Washington University, USA	<i>Caenorhabditis elegans</i> and related nematode genomic information	(38)
WTSI Mouse Genetics	WTSI, UK	Mouse phenotyping and expression data captured from mutant mouse lines	(39)



### GENE REPORT

Ensembl Gene ID(s)
ENSG00000146648
Go »

GENE INFO

**Ensembl Gene ID:** [ENSG00000146648](#) [EGFR]

<b>Description:</b> epidermal growth factor receptor <small>[Source:HGNC Symbol;Acc:3236]</small>	<b>Chromosome:</b> 7
<b>Gene Start (bp):</b> 55086714	<b>Gene End (bp):</b> 55324313
<b>Band:</b> p11.2	<b>Strand:</b> 1
<b>Gene Biotype:</b> protein_coding	<b>Status:</b> KNOWN
<b>Transcript count:</b> 13	

ALIASES

<b>HGNC symbol:</b> <a href="#">EGFR</a>	<b>EntrezGene ID:</b> <a href="#">1956</a>
<b>VEGA gene ID(s) (OTTG):</b> <a href="#">OTTHUMG00000023661</a>	<b>RefSeq DNA ID:</b> <a href="#">NM_201283</a> , <a href="#">NM_005228</a> , <a href="#">NM_201284</a> , <a href="#">NM_201282</a>
<b>UniProt/SwissProt</b> <a href="#">P00533</a>	
<b>Accession:</b>	

PATHWAY ANNOTATION

**Pathway name (Reactome):** Signaling by EGFR, Grb2 events in EGFR signaling, Gab1 signalosome, Shc events in EGFR signaling, EGFR downregulation, EGFR interacts with phospholipase C-gamma, L1CAM interactions, Axon guidance, Signal transduction by L1

GO BIOLOGICAL PROCESS

**GO Biological Process:** protein phosphorylation, transmembrane receptor protein tyrosine kinase signaling pathway, regulation of peptidyl-tyrosine phosphorylation, positive regulation of MAP kinase activity, cell morphogenesis, neuron projection morphogenesis, positive regulation of nitric oxide biosynthetic process, epidermal growth factor receptor signaling pathway, positive regulation of epithelial cell proliferation, activation of MAPKK activity, response to oxidative stress, response to lipid, response to calcium ion, positive regulation of catenin protein nuclear translocation, response to stress, activation of phospholipase C activity, translation, salivary gland morphogenesis, tongue development, positive regulation of synaptic transmission, glutamatergic, positive regulation of protein kinase B signaling cascade, positive regulation of cell migration, signal transduction, positive regulation of cell proliferation, response to osmotic stress, negative regulation of mitotic cell cycle, protein autophosphorylation, ossification, cell surface receptor linked signaling pathway, circadian rhythm, ovulation cycle, negative regulation of apoptosis, astrocyte activation, protein insertion into membrane, cell proliferation, embryonic placenta development, hair follicle development, positive regulation of cyclin-dependent protein kinase activity involved in G1/S, positive regulation of phosphorylation, cerebral cortex cell migration, positive regulation of smooth muscle cell proliferation, response to UV-A, regulation of nitric-oxide synthase activity, cell-cell adhesion, activation of phospholipase A2 activity by calcium-mediated signaling, intracellular protein kinase cascade, morphogenesis of an epithelial fold

GO CELLULAR COMPONENT


**GO Cellular Component:** membrane, plasma membrane, apical plasma membrane, extracellular region, extracellular space, basolateral plasma membrane, endosome, cytoplasm, Shc-EGFR complex, nucleus, integral to membrane, endocytic vesicle, AP-2 adaptor complex, intracellular

GO MOLECULAR FUNCTION

**GO Molecular Function:** protein kinase activity, ATP binding, protein serine/threonine kinase activity, transmembrane receptor protein tyrosine kinase activity, protein tyrosine kinase activity, protein heterodimerization activity, protein binding, receptor signaling protein tyrosine kinase activity, actin filament binding, MAP/ERK kinase kinase activity, epidermal growth factor receptor activity, transmembrane receptor activity, nitric-oxide synthase regulator activity, transferase activity, double-stranded DNA binding, protein phosphatase binding, identical protein binding, nucleotide binding

Figure 4. The Gene Report page for EGFR, displaying data federated from several sources.

SEQUENCES



<input type="radio"/> Unspliced (Transcript)	<input type="radio"/> 3' UTR
<input type="radio"/> Unspliced (Gene)	<input type="radio"/> Exon Sequences
<input type="radio"/> Flank (Transcript)	<input checked="" type="radio"/> cDNA Sequences
<input type="radio"/> Flank (Gene)	<input type="radio"/> Coding Sequences
<input type="radio"/> Flank-coding region (Transcript)	<input type="radio"/> Protein
<input type="radio"/> 5' UTR	

Upstream Flank:

Downstream Flank:

FILTERS

**GENE:**

Limit to genes ...:

Entries with following ID(s):

[upload file](#)

Type:

scrRNA\_pseudogene

snoRNA

snoRNA\_pseudogene

snRNA

snRNA\_pseudogene

tRNA\_pseudogene

TR\_C\_gene

TR\_J\_gene

TR\_V\_gene

TR\_V\_pseudogene

**PROTEIN DOMAINS:**

Limit to genes ...:

Transmembrane domains:  Only  Excluded

Signal domains:  Only  Excluded

HEADER INFORMATION

**GENE INFORMATION**

<input checked="" type="checkbox"/> Ensembl Gene ID	<input type="checkbox"/> Description	<input type="checkbox"/> Associated Gene Name
<input type="checkbox"/> Associated Gene DB	<input type="checkbox"/> Chromosome Name	<input type="checkbox"/> Gene Start (bp)
<input type="checkbox"/> Gene End (bp)	<input type="checkbox"/> Ensembl Protein Family ID(s)	

**TRANSCRIPT INFORMATION**

<input type="checkbox"/> CDS start (within cDNA)	<input type="checkbox"/> CDS end (within cDNA)	<input type="checkbox"/> 5' UTR Start
<input type="checkbox"/> 5' UTR End	<input type="checkbox"/> 3' UTR Start	<input type="checkbox"/> 3' UTR End
<input checked="" type="checkbox"/> Ensembl Transcript ID	<input type="checkbox"/> Ensembl Protein ID	<input type="checkbox"/> Strand
<input type="checkbox"/> Transcript Start (bp)	<input type="checkbox"/> Transcript End (bp)	

Figure 5. The sequence retrieval plug-in page.



## Query examples

One of the great strengths of Central Portal is that it allows cross-database searches that any individual resource would not. Here are some examples of the possibilities afforded by this feature.

Query #1: 'Find insertion-frameshift mutations in the COSMIC database that affect genes involved in Apoptosis'.

Entry point	Filters
Gene retrieval > cancer genes	COSMIC: Mutation type-AA: Insertion-frameshift KEGG: KEGG Pathway: apoptosis

By integrating data from the COSMIC and KEGG databases, Central Portal allows users to identify COSMIC mutations specific to their pathways of interest. The Pathway title links back to the KEGG web site and mutation ID links back to the COSMIC web site, providing the ability to obtain more detailed information on the pathway or on the mutation, respectively.

Query #2: 'Retrieve the cDNA sequences of protein-coding human genes that have HGNC IDs' (Figure 5).

Entry point	Data sets	Filters/attributes
Sequence retrieval > Ensembl	<i>Homo sapiens</i> gene (GRCh37.p2)	Sequences: cDNA sequences  Filters: Limit to genes: with HGNC ID(s) Type: protein_coding Header information: Ensembl Gene ID Ensembl Transcript ID

By combining the sequence retrieval tool with search capabilities, BioMart reduces what is often a two-step process—retrieving a list of genes, and then retrieving the sequences of those genes—into a single query.

## Future directions

BioMart Central Portal is constantly evolving thanks to the efforts of the community that supports it and contributes data. To make joining Central Portal easier, we are creating BioMart Central Registry. With this resource, database administrators will be able to create an account, add their data sources and suggest categorization for them. Once registered, participants will also be able to make changes to their databases and notify Central Portal of updates.

In addition to including new data sets, Central Portal will evolve, as new tools are developed and added. Such tools will perform deeper analysis, such as detecting enrichment of certain properties (e.g. GO terms) within a given set of genes or calculating consequences given a list of SNP terms. BioMart plug-ins developed by other community members may also be incorporated, further strengthening the project as a whole.

## Acknowledgements

BioMart Central Portal is a collaborative, community effort and as such it is the product of the efforts of dozens, if not hundreds, of people. Creating a biological database is a multi-step process: experimenters must collect the data, database managers must create data models and administer databases and bioinformaticians must create methods for analysing the data. Additionally, over the years many programmers have contributed to the BioMart project codebase. We would like to acknowledge all the hard work of the many contributors to the projects that BioMart comprises.

## Funding

The development of the BioMart software and the creation and hosting of BioMart Central Portal was supported by the Ontario Institute for Cancer Research and the Ontario Ministry for Research and Innovation. The individual data sources that Central Portal comprises are funded separately and independently.

*Conflict of interest.* None declared.

## References

- Haider,S., Ballester,B., Smedley,D. et al. (2009) BioMart Central Portal—unified access to biological data. *Nucleic Acids Res.*, **37**(Web Server issue), W23–W27.
- Smedley,D., Haider,S., Ballester,B. et al. (2009) BioMart—biological queries made easy. *BMC Genom.*, **10**, 22.
- Zhang,J., Haider,S., Guberman,J.M. et al. (2011) BioMart: a data federation framework for large collaborative projects. *Database*, (This issue).
- Zhang,J., Baran,J., Cros,A. et al. (2011) International Cancer Genome Consortium Data Portal: a One Stop-Shop for Cancer Genomics Data. *Database*, (This issue).
- Kanehisa,M. and Goto,S. (2000) KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.*, **28**, 27–30.
- Kanehisa,M., Goto,S., Furumichi,M. et al. (2010) KEGG for representation and analysis of molecular networks involving diseases and drugs. *Nucleic Acids Res.*, **38**(Database issue), D355–D360.
- Kanehisa,M., Goto,S., Hattori,M. et al. (2006) From genomics to chemical genomics: new developments in KEGG. *Nucleic Acids Res.*, **34**(Database issue), D354–D357.

8. Cutts,R.J., Gadaleta,E., Lemoine,N.R. *et al.* (2011) Using BioMart as a framework to create a cancer-specific database. *Database*, Published online 11 June 2011, doi:10.1093/database/bar024.
9. Kinsella,R., Kahari,A., Haider,S. *et al.* (2011) Ensembl BioMarts: a hub for data retrieval across the taxonomic space. *Database*, Published online 23 July 2011, doi:10.1093/database/bar030.
10. Povey,S., Lovering,R., Bruford,E. *et al.* (2001) The HUGO Gene Nomenclature Committee (HGNC). *Hum. Genet.*, **109**, 678–680.
11. Jones,P., Binns,D., McMenamin,C. *et al.* (2011) The InterPro BioMart: Powerful, federated query and web-service access to the InterPro resource. *Database*, Published online 23 July 2011, doi:10.1093/database/bar033.
12. Haw,R., Croft,D., Yung,C.K. *et al.* (2011) The Reactome BioMart. *Database*, (This issue).
13. Stevenson,P., Richardson,L., Venkataraman,S. *et al.* (2011) The BioMart interface to the eMouseAtlas gene expression EMAGE. *Database*, (This issue).
14. Shepherd,R., Forbes,S.A., Beare,D. *et al.* (2011) Data mining using the Catalogue of Somatic Mutations in Cancer BioMart (COSMICMart). *Database*, Published online 23 May 2011, doi:10.1093/database/bar018.
15. Spooner,W., Youens-Clark,K. and Ware,D. (2011) GrameneMart: The BioMart Data Portal for the Gramene Project. *Database*, This issue.
16. Arnaiz,O., Malinowska,A., Klotz,C. *et al.* (2009) Cildb: a knowledge-base for centrosomes and cilia. *Database*, **2009**, bap022 Published online 7 December 2009, doi: 10.1093/database/bap022.
17. Wilkinson,P., Sengerova,J., Matteoni,R. *et al.* (2010) EMMA–mouse mutant resources for the international scientific community. *Nucleic Acids Res.*, **38**(Database issue), D570–D576.
18. Diez-Roux,G., Banfi,S., Sultan,M. *et al.* (2011) A high-resolution anatomical atlas of the transcriptome in the mouse embryo. *PLoS Biol.*, **9**, e1000582.
19. Blake,A. (2011) The EuroPhenome BioMart: A mouse phenotyping resource. *Database*, (This issue).
20. Lardenois,A., Gattiker,A., Collin,O. *et al.* (2010) GermOnline 4.0 is a genomics gateway for germline development, meiosis and the mitotic cell cycle. *Database*, **2010**, baq030 Published online 10 December 2010, doi:10.1093/database/baq030.
21. International HapMap Consortium. (2003) The International HapMap Project. *Nature*, **426**, 789–796.
22. Oakley,J., Iyer,V., Skarnes,W.C. *et al.* (2011) BioMart as an integration solution for the International Knockout Mouse Consortium. *Database*, (This issue).
23. Perez-Llamas,C., Gundem,G. and Lopez-Bigas,N. (2011) Integrative Cancer Genomics (IntOGen) in BioMart. *Database*, (This issue).
24. KazusaMart. <http://mart.kazusa.or.jp/biomart/martview/> (5 August 2011, date last accessed).
25. Shaw,D.R. (2009) Searching the Mouse Genome Informatics (MGI) resources for information on mouse biology from genotype to phenotype. *Curr. Protoc. Bioinformatics*, **Chapter 1**, Unit 1.7.
26. Arnaiz,O. and Sperling,L. (2011) ParameciumDB in 2011: new tools and new data for functional and comparative genomics of the model ciliate Paramecium tetraurelia. *Nucleic Acids Res.*, **39**(Database issue), D632–D636.
27. McLaughlin,T., Siepen,J.A., Selley,J. *et al.* (2006) PepSeeker: a database of proteome peptide identifications for investigating fragmentation patterns. *Nucleic Acids Res.*, **34**(Database issue), D649–D654.
28. Phytozome. <http://www.phytozome.net/biomart/martview> (5 August 2011, date last accessed).
29. International Potato Center. <http://germplasmb.cip.cgiar.org/biomart/martview> (5 August 2011, date last accessed).
30. Vizcaino,J.A., Cote,R., Reisinger,F. *et al.* (2010) The proteomics identifications database: 2010 update. *Nucleic Acids Res.*, **38**(Database issue), D736–D742.
31. Rice Map. <http://ricemart.cbi.edu.cn/biomart/martview> (5 August 2011, date last accessed).
32. SalmonDB. <http://genomicasalmones.dim.uchile.cl:9002/biomart/martview/> (5 August 2011, date last accessed).
33. Ai,J., Hu,S., Kasprzyk,A. *et al.* (2011) SDxMart: The BioMart Data Portal for Saliva Diagnostics. *Database*, This Issue.
34. Moreews,F., Klopp,C., Rauffet,G. *et al.* (2011) SigReannot-mart: A query environment for expression microarray probe re-annotations. *Database*, (This issue).
35. UniProt Consortium. (2010) The Universal Protein Resource (UniProt) in 2010. *Nucleic Acids Res.*, **38** (Database issue), D142–D148.
36. Lawson,D., Arensburger,P., Atkinson,P. *et al.* (2007) VectorBase: a home for invertebrate vectors of human pathogens. *Nucleic Acids Res.*, **35**(Database issue), D503–D505.
37. Wilming,L.G., Gilbert,J.G., Howe,K. *et al.* (2008) The vertebrate genome annotation (Vega) database. *Nucleic Acids Res.*, **36**(Database issue), D753–D760.
38. Harris,T.W., Antoshechkin,I., Bieri,T. *et al.* (2010) WormBase: a comprehensive resource for nematode research. *Nucleic Acids Res.*, **38**(Database issue), D463–D467.
39. Mouse Resources Portal. <http://www.sanger.ac.uk/htgt/biomart/martview> (5 August 2011, date last accessed).