

# The 2014 *Nucleic Acids Research* Database Issue and an updated NAR online Molecular Biology Database Collection

Xosé M. Fernández-Suárez<sup>1</sup>, Daniel J. Rigden<sup>2</sup> and Michael Y. Galperin<sup>3,\*</sup>

<sup>1</sup>Life Technologies, Inchinnan Business Park, Paisley PA4 9RF, UK, <sup>2</sup>Institute of Integrative Biology, University of Liverpool, Crown Street, Liverpool L69 7ZB, UK and <sup>3</sup>National Center for Biotechnology Information (NCBI), National Library of Medicine, National Institutes of Health, Bethesda, MD 20894, USA

Received November 17, 2013; Accepted November 18, 2013

## ABSTRACT

The 2014 *Nucleic Acids Research* Database Issue includes descriptions of 58 new molecular biology databases and recent updates to 123 databases previously featured in *NAR* or other journals. For convenience, the issue is now divided into eight sections that reflect major subject categories. Among the highlights of this issue are six databases of the transcription factor binding sites in various organisms and updates on such popular databases as CAZy, Database of Genomic Variants (DGV), dbGaP, DrugBank, KEGG, miRBase, Pfam, Reactome, SEED, TCDB and UniProt. There is a strong block of structural databases, which includes, among others, the new RNA Bricks database, updates on PDBe, PDBsum, ArchDB, Gene3D, ModBase, Nucleic Acid Database and the recently revived iPfam database. An update on the NCBI's MMDB describes VAST+, an improved tool for protein structure comparison. Two articles highlight the development of the Structural Classification of Proteins (SCOP) database: one describes SCOPe, which automates assignment of new structures to the existing SCOP hierarchy; the other one describes the first version of SCOP2, with its more flexible approach to classifying protein structures. This issue also includes a collection of articles on bacterial taxonomy and metagenomics, which includes updates on the List of Prokaryotic Names with Standing in Nomenclature (LPSN), Ribosomal Database Project (RDP), the SILVA/LTP project and several new metagenomics resources. The NAR online Molecular Biology Database Collection, <http://www.oxfordjournals.org/nar/database/c/>, has been expanded to 1552 databases. The entire

Database Issue is freely available online on the *Nucleic Acids Research* website (<http://nar.oxfordjournals.org/>).

## NEW AND UPDATED DATABASES

The 21st annual *Nucleic Acids Research* Database Issue is the largest ever. It includes 185 articles that provide (i) descriptions of the database resources at the NCBI, European Bioinformatics Institute (EBI) and the US Department of Energy Joint Genome Institute (JGI); (ii) 58 new molecular biology databases (Table 1); (iii) updates on 100 databases previously featured in *NAR*; and (iv) updated descriptions of 23 databases that had been previously described in other journals (Table 2). For the past several years, the order of articles in the Database Issue reflected the categorization of the databases in the NAR online Molecular Biology Database Collection (<http://www.oxfordjournals.org/nar/database/c/>). Acting on the advice of many readers, we have now made the categories visible and divided the entire Database Issue into the following eight sections: (i) nucleic acid sequence and structure, transcriptional regulation; (ii) protein sequence and structure, motifs and domains, protein–protein interactions; (iii) metabolic and signalling pathways, enzymes, protein modification; (iv) viruses, bacteria, protozoa and fungi; (v) human genome, model organisms, comparative genomics; (vi) genomic variation, diseases and drugs; (vii) plant databases; and (viii) other molecular biology databases. Although each of these sections unifies several of the categories and/or subcategories of the NAR online Database Collection, we believe that they provide an easy-to-use guide to navigate this huge volume and help placing related databases next to each other.

The first section, in addition to the annual descriptions of GenBank, the European Nucleotide Archive and the DNA Data Bank of Japan, includes update papers on

\*To whom correspondence should be addressed. Tel: +1 301 435 5910; Fax: +1 301 435 7793; Email: [nardatabase@gmail.com](mailto:nardatabase@gmail.com) or [galperin@ncbi.nlm.nih.gov](mailto:galperin@ncbi.nlm.nih.gov)

**Table 1.** Descriptions of new online databases in the 2014 NAR Database issue

Database name	URL	Brief description
1000 Genomes Selection Browser	<a href="http://hsb.upf.edu">http://hsb.upf.edu</a>	Signatures of selection in the human genomes
AgeFactDB	<a href="http://agefactdb.jenage.de">http://agefactdb.jenage.de</a>	<u>A</u> geing <u>F</u> actors, phenotypes and lifespan data
AVPdb	<a href="http://crdd.osdd.net/servers/avpdb">http://crdd.osdd.net/servers/avpdb</a>	A database of experimentally validated <u>A</u> nti <u>V</u> iral <u>P</u> eptides
BacDive	<a href="http://bacdive.dsmz.de">http://bacdive.dsmz.de</a>	<u>B</u> acterial <u>D</u> iversity metadatabase
BacMet	<a href="http://bacmet.biomedicine.gu.se">http://bacmet.biomedicine.gu.se</a>	<u>A</u> ntibacterial biocide and <u>M</u> etal resistance <u>G</u> enes
BloodChIP	<a href="http://149.171.101.136/python/BloodChIP">http://149.171.101.136/python/BloodChIP</a>	Transcription factor binding profiles in human haematopoietic stem/progenitor cells
bNAber	<a href="http://bnabs.org">http://bnabs.org</a>	A database of broadly <u>N</u> eutralizing <u>H</u> IV-1 <u>A</u> ntibodies
CellFinder	<a href="http://www.cellfinder.org">http://www.cellfinder.org</a>	Gene and protein expression, phenotype and images mapped to the cell types
ClinVar	<a href="http://www.ncbi.nlm.nih.gov/clinvar">http://www.ncbi.nlm.nih.gov/clinvar</a>	Genomic <u>V</u> ariation of potential <u>C</u> linical importance
CollecTF	<a href="http://collectf.umbc.edu">http://collectf.umbc.edu</a>	<u>C</u> ollection of verified bacterial <u>T</u> ranscription <u>F</u> actor binding sites
CR Cistrome	<a href="http://compbio.tongji.edu.cn/cr">http://compbio.tongji.edu.cn/cr</a>	<u>C</u> hromatin <u>R</u> egulators and histone modifications in human and mouse
dbPSHP	<a href="http://jjwanglab.org/dbpsdp">http://jjwanglab.org/dbpsdp</a>	A database of recent <u>P</u> ositive <u>S</u> election across <u>H</u> uman <u>P</u> opulations
DRPR	<a href="http://syslab.nchu.edu.tw/DRPR">http://syslab.nchu.edu.tw/DRPR</a>	<u>P</u> henotype-specific <u>R</u> egulatory <u>P</u> rograms derived from TF binding data
DriverDB	<a href="http://ngs.yu.edu.tw/driverdb/">http://ngs.yu.edu.tw/driverdb/</a>	<u>C</u> ancer <u>d</u> river genes/mutations deduced from cancer exome-seq results
EBI metagenomics	<a href="https://www.ebi.ac.uk/metagenomics/">https://www.ebi.ac.uk/metagenomics/</a>	An automated pipeline for the analysis and archiving of metagenomic data
EKPD	<a href="http://ekpd.biocuckoo.us">http://ekpd.biocuckoo.us</a>	<u>E</u> karyotic protein <u>K</u> inase and <u>P</u> hosphatase <u>D</u> atabase
ExoLocator	<a href="http://exolocator.eopsf.org">http://exolocator.eopsf.org</a>	<u>P</u> rotein-coding <u>e</u> xons from complete vertebrate genomes
GoMapMan	<a href="http://www.gomapman.org">http://www.gomapman.org</a>	Unified plant-specific gene ontology
GWIPS-viz	<a href="http://gwips.ucc.ie">http://gwips.ucc.ie</a>	<u>G</u> enome- <u>W</u> ide <u>I</u> nformation on <u>P</u> rotein <u>S</u> ynthesis <i>in vivo</i> using ribosome profiling
Hemolytik	<a href="http://crdd.osdd.net/raghava/hemolytik">http://crdd.osdd.net/raghava/hemolytik</a>	Haemolytic and non-haemolytic peptides
HoPaCI-DB	<a href="http://mips.helmholtz-muenchen.de/HoPaCI/">http://mips.helmholtz-muenchen.de/HoPaCI/</a>	<u>H</u> ost- <u>P</u> athogen <u>I</u> nteractions of <i>Pseudomonas aeruginosa</i> and <i>Coxiella</i> spp.
HRaP	<a href="http://bioinfo.protres.ru/hrap">http://bioinfo.protres.ru/hrap</a>	<u>H</u> omo <u>R</u> epeats and <u>P</u> atterns
InvFEST	<a href="http://invfestdb.uab.cat">http://invfestdb.uab.cat</a>	Polymorphic <u>i</u> nversions in the human genome
IUPHAR/BPS guide to pharmacology	<a href="http://www.guidetopharmacology.org">http://www.guidetopharmacology.org</a>	Properties of established and potential drug targets: GPCRs, ion channels, nuclear hormone receptors, catalytic receptors, transporters and enzymes
LenVarDB	<a href="http://caps.ncbs.res.in/lenvardb">http://caps.ncbs.res.in/lenvardb</a>	<u>L</u> ength <u>V</u> ariation in protein domains
LoQATe	<a href="http://www.weizmann.ac.il/molgen/loqate">http://www.weizmann.ac.il/molgen/loqate</a>	<u>L</u> ocalization and <u>Q</u> uantitation <u>A</u> tlas of the yeast proteome
Lynx	<a href="http://lynx.ci.uchicago.edu">http://lynx.ci.uchicago.edu</a>	Genomic and clinical data on complex heritable disorders
Manteia	<a href="http://manteia.igbmc.fr">http://manteia.igbmc.fr</a>	Embryonic development of the mouse, chicken, zebrafish and human
MCDRiceProt	<a href="http://www.genomeindia.org/biocuration">http://www.genomeindia.org/biocuration</a>	<u>M</u> anually <u>C</u> urated <u>D</u> atabase of <u>R</u> ice <u>P</u> roteins
MetaRef	<a href="http://bioref.org">http://bioref.org</a>	<u>R</u> eference clade-specific microbial genes for <u>M</u> etagenomic studies
MitoBreak	<a href="http://mitobreak.portugene.com">http://mitobreak.portugene.com</a>	<u>M</u> itochondrial <u>D</u> N <u>A</u> <u>B</u> reakpoints in human, mouse and rat
MP:PD	<a href="http://proteininformatics.charite.de/mppd">http://proteininformatics.charite.de/mppd</a>	<u>M</u> embrane <u>P</u> roteins: <u>P</u> acking <u>D</u> ensities, packing defects and internal water molecules
MultiTaskDB	<a href="http://wallace.uab.es/multitask">http://wallace.uab.es/multitask</a>	Moonlighting proteins database
mVOC	<a href="http://bioinformatics.charite.de/mvoc">http://bioinformatics.charite.de/mvoc</a>	<u>M</u> icrobial <u>V</u> olatile <u>O</u> rganic <u>C</u> ompounds
NECTAR	<a href="http://cardiodb.org/nectar">http://cardiodb.org/nectar</a>	Disease-related non-synonymous mutations
Network Portal	<a href="http://networks.systemsbio.net">http://networks.systemsbio.net</a>	A database of gene transcription regulatory networks
NeXO	<a href="http://nexontology.org/">http://nexontology.org/</a>	<u>N</u> etwork <u>E</u> xtracted gene <u>O</u> ntology database
NHGRI GWAS Catalog	<a href="http://www.genome.gov/gwastudies">http://www.genome.gov/gwastudies</a> , <a href="http://www.ebi.ac.uk/fgpt/gwas">http://www.ebi.ac.uk/fgpt/gwas</a>	A catalog of published <u>G</u> enome- <u>W</u> ide <u>A</u> ssociation <u>S</u> tudies, maintained at the NHGRI and EBI
OnTheFly	<a href="http://bhapp.c2b2.columbia.edu/OnTheFly">http://bhapp.c2b2.columbia.edu/OnTheFly</a>	<u>D</u> N <u>A</u> -binding specificities of transcription factors in <i>Drosophila</i>
pE-DB	<a href="http://pedb.vib.be">http://pedb.vib.be</a>	<u>P</u> rotein <u>E</u> nsemble <u>D</u> ata <u>B</u> ase: ensembles of intrinsically disordered and unfolded proteins
P-MITE	<a href="http://pmite.hzau.edu.cn/django/mite">http://pmite.hzau.edu.cn/django/mite</a>	<u>P</u> lant <u>M</u> iniature <u>I</u> nverted-repeat <u>T</u> ransposable <u>E</u> lements (MITEs)
POGO-DB	<a href="http://pogo.ece.drexel.edu">http://pogo.ece.drexel.edu</a>	<u>P</u> airwise comparisons <u>O</u> f <u>G</u> enomes and universal <u>O</u> rthologous genes
PortEco	<a href="http://porteco.org">http://porteco.org</a>	<i>Escherichia coli</i> K-12 knowledgebase <u>P</u> ortal
RADAR	<a href="http://rnaedit.com">http://rnaedit.com</a>	A <u>R</u> igorously <u>A</u> nnnotated <u>D</u> atabase of <u>A</u> -to- <u>I</u> <u>R</u> NA editing
RepeatsDB	<a href="http://repeatsdb.bio.unipd.it">http://repeatsdb.bio.unipd.it</a>	Repeats in protein structures
RhizoBase	<a href="http://genome.microbedb.jp/rhizobase">http://genome.microbedb.jp/rhizobase</a>	Manually curated annotations for <u>r</u> hizobial genomes
RiceWiki	<a href="http://ricewiki.big.ac.cn">http://ricewiki.big.ac.cn</a>	<u>W</u> iki-based open-content platform for community curation of <u>r</u> ice genes
RNA Bricks	<a href="http://iimcb.genesilico.pl/rnabricks">http://iimcb.genesilico.pl/rnabricks</a>	<u>R</u> NA structural modules and their interactions
rSNPBase	<a href="http://rsnp.psych.ac.cn">http://rsnp.psych.ac.cn</a>	<u>A</u> nnnotated <u>S</u> N <u>P</u> s within regulatory DNA elements
SABdab	<a href="http://opig.stats.ox.ac.uk/webapps/sabdab">http://opig.stats.ox.ac.uk/webapps/sabdab</a>	<u>S</u> tructural <u>A</u> ntibody <u>D</u> atabase
SMMRNA	<a href="http://www.smmrna.org">http://www.smmrna.org</a>	<u>S</u> mall <u>M</u> olecule inhibitors of <u>R</u> NA
SporeWeb	<a href="http://sporeweb.molgenrug.nl">http://sporeweb.molgenrug.nl</a>	<u>R</u> egulatory pathways during the <u>s</u> porulation cycle of <i>Bacillus subtilis</i>
SuperPain	<a href="http://bioinformatics.charite.de/superpain">http://bioinformatics.charite.de/superpain</a>	Compounds that stimulate or relieve pain
TFBSShape	<a href="http://rohslab.cmb.usc.edu/TFBSShape">http://rohslab.cmb.usc.edu/TFBSShape</a>	<u>D</u> N <u>A</u> <u>s</u> hape features of <u>T</u> ranscription <u>F</u> actor <u>B</u> inding <u>S</u> ites
TISdb	<a href="http://tisdb.human.cornell.edu">http://tisdb.human.cornell.edu</a>	<u>A</u> lternative <u>T</u> ranslation <u>I</u> nitiation <u>S</u> ites
Transformer	<a href="http://bioinformatics.charite.de/transformer">http://bioinformatics.charite.de/transformer</a>	<u>B</u> iotransformation of drugs and food ingredients by human enzymes
uORFdb	<a href="http://cbdm.mdc-berlin.de/tools/uorfdb">http://cbdm.mdc-berlin.de/tools/uorfdb</a>	<u>U</u> pstream <u>O</u> R <u>F</u> s and their effect of translation of downstream CDSs
WormQTL <sup>HD</sup>	<a href="http://www.wormqtl-hd.org">http://www.wormqtl-hd.org</a>	<u>L</u> inks from <u>h</u> uman <u>d</u> isease to natural variation data in <i>Caenorhabditis elegans</i>

**Table 2.** Updated descriptions of online databases not previously featured in the NAR Database issue

Database name	URL	Brief description
COLOMBOS	<a href="http://colombos.net">http://colombos.net</a>	<u>C</u> ollections of <u>M</u> icroarrays for <u>B</u> acterial <u>O</u> rganisms
Consensus CDS	<a href="http://www.ncbi.nlm.nih.gov/projects/CCDS">http://www.ncbi.nlm.nih.gov/projects/CCDS</a>	A collaborative effort to identify a core set of human protein coding regions
CottonGen	<a href="http://www.cottongen.org">http://www.cottongen.org</a>	<u>C</u> otton <u>G</u> enomics, genetics and breeding
Database of Genomic Variants	<a href="http://dgv.tcag.ca/dgv/app/home">http://dgv.tcag.ca/dgv/app/home</a>	Curated catalog of human genomic structural variation
dbGaP	<a href="http://www.ncbi.nlm.nih.gov/gap">http://www.ncbi.nlm.nih.gov/gap</a>	Database of Genotypes and Phenotypes
DECIPHER	<a href="http://decipher.sanger.ac.uk">http://decipher.sanger.ac.uk</a>	Database of Chromosomal Imbalance and Phenotype in Humans Using Ensembl Resources
GEISHA	<a href="http://geisha.arizona.edu">http://geisha.arizona.edu</a>	<u>G</u> allus <u>E</u> xpression <u>I</u> n <u>S</u> itu <u>H</u> ybridization <u>A</u> nalysis
GeneProf	<a href="http://www.geneprof.org">http://www.geneprof.org</a>	Human and mouse gene expression data from RNA-seq and ChIP-seq
GGBN	<a href="http://www.ggbn.org/dataportal">http://www.ggbn.org/dataportal</a>	The <u>G</u> lobal <u>G</u> enome <u>B</u> iodiversity <u>N</u> etwork portal
HMDD	<a href="http://cmbi.bjmu.edu.cn/hmdd">http://cmbi.bjmu.edu.cn/hmdd</a>	<u>H</u> uman <u>M</u> icroRNA and <u>D</u> isease <u>A</u> sociations <u>D</u> atabase
Human Phenotype Ontology	<a href="http://www.human-phenotype-ontology.org">http://www.human-phenotype-ontology.org</a>	Standardized vocabulary of phenotypic abnormalities in human disease
IMPC	<a href="http://www.mousephenotype.org">http://www.mousephenotype.org</a>	<u>I</u> nternational <u>M</u> ouse <u>P</u> henotyping <u>C</u> onsortium portal
iPfam	<a href="http://www.ipfam.org">http://www.ipfam.org</a>	Protein family interactions mapped to Pfam domains
KBDOCK	<a href="http://kbdock.loria.fr">http://kbdock.loria.fr</a>	<u>K</u> nowledge- <u>B</u> ased <u>D</u> ocking: protein domain interfaces
Locus Reference Genomic	<a href="http://www.lrg-sequence.org">http://www.lrg-sequence.org</a>	Each LRG is a stable genomic DNA sequence for a region of the human genome
LPSN	<a href="http://www.bacterio.net">http://www.bacterio.net</a>	<u>L</u> ist of <u>P</u> rokaroyotic names with <u>S</u> tanding in <u>N</u> omenclature
NDB	<a href="http://ndbserver.rutgers.edu">http://ndbserver.rutgers.edu</a>	<u>N</u> ucleic <u>A</u> cid <u>D</u> ata <u>B</u> ase, nucleic acids-containing structures
Plasma Proteome Database	<a href="http://www.plasmaproteomedatabase.org">http://www.plasmaproteomedatabase.org</a>	Quantitative information on proteins in human plasma and serum
Progenetix	<a href="http://www.progenetix.org">http://www.progenetix.org</a>	Copy number abnormalities in human cancer
SEED	<a href="http://pubseed.theseed.org">http://pubseed.theseed.org</a> or <a href="http://www.theseed.org">http://www.theseed.org</a>	Genome annotations based on curated functional systems
SFLD	<a href="http://sfld.rbvi.ucsf.edu">http://sfld.rbvi.ucsf.edu</a>	<u>S</u> tructure- <u>F</u> unction <u>L</u> inkage <u>D</u> atabase: sequence conservation in enzyme superfamilies
SoyKB	<a href="http://soykb.org">http://soykb.org</a>	<u>S</u> oybean <u>K</u> nowledge <u>B</u> ase
Virus variation	<a href="http://www.ncbi.nlm.nih.gov/genomes/VirusVariation">http://www.ncbi.nlm.nih.gov/genomes/VirusVariation</a>	Variation data on influenza, dengue and West Nile viruses
YeastNet	<a href="http://www.inetbio.org/yeastnet">http://www.inetbio.org/yeastnet</a>	Functional gene networks for <i>Saccharomyces cerevisiae</i>

five microRNA databases: miRBase, miRNEST, mirTarBase, PolymiRTS and starBASE, and on the NONCODE database of various types of non-coding RNA. There are also several databases of transcription factor (TF) binding sites (TFBSs), including updates on JASPAR and YEASTRACT and new databases of TFBS in *Escherichia coli*, *Drosophila* and human haematopoietic stem cells (1–5). An interesting work, chosen by the reviewers and NAR editors as a ‘Breakthrough paper’, describes TFBSshape (6), a database of DNA structural features (minor groove width, roll, propeller twist and helix twist) of TFBSs for various TFs, which have been collected from the JASPAR (1) and UniProbe (7) databases. The TFBSshape website, <http://rohslab.cmb.usc.edu/TFBSshape/>, allows the users to submit their own aligned TFBS sequences, which could be used, for example, to compare the DNA binding specificities of closely related TFs (6).

The protein database section includes annual updates on UniProt and KEGG, as well as updates of such popular databases as Pfam, eggNOG, ELM, FireDB, SEED, SIMAP and Transporter protein Classification DataBase (TCDB). Two new databases, HRaP and RepeatsDB, collect information on protein repeats, the former at the sequence level (runs of the same amino acid residue) and the latter at the structural level (8,9).

As in previous years, this Database Issue features an impressive selection of structural databases. Two of

them deal with nucleic acid structure: an update on the well-known Nucleic Acids Database (NDB) and RNA Bricks, a new database of RNA 3D motifs and their contacts (10,11). The block of protein structure databases includes, among others, updates on Protein Data Bank in Europe (PDBe), PDBsum, ArchDB, Gene3D, ModBase, SCOP and the recently revived iPfam database. Diverse improvements at PDBe include comprehensive visualization and analysis of the rapidly growing collection of electron microscopy-derived structures, whereas PDBsum now offers facilities to browse domain architectures and new connections to ligand and SNP data (12). For the past 18 years, the NCBI’s Molecular Modeling Database (MMDB) displayed the lists of structural neighbours of a given protein, calculated using the Vector Alignment Search Tool (VAST) (13). The MMDB update paper describes VAST+, a recent extension of that tool, which allows calculation of structural similarity at the level of ‘biological assemblies’ (hetero- or homo-oligomeric protein complexes). Accordingly, for macromolecular complexes, MMDB now displays precalculated lists of similar protein complexes ranked by the extent of similarity (14). Several databases, including iPfam, 3did and UniHI, reflect current interest in the structural basis of protein interaction networks and take on the challenge of presenting complex protein–protein and protein–ligand interaction data in clear and useful ways (15–17). The aptly named Negatome database (18) provides a useful

benchmark, documenting protein pairs that definitely do not interact, and could be used as negative control for the constantly growing protein 'interactome'. A pair of papers published back-to-back highlight two different directions in the development of the Structural Classification of Proteins (SCOP) database: one of them describes SCOPe, an extension of SCOP that focuses on regularly and automatically assigning new structures to the existing SCOP hierarchy, whereas the other one describes the birth of SCOP2, with its more flexible graph-based approach to classifying protein structures and documenting the subtleties of their relationships (19,20).

The section on enzymes and metabolism includes updates on three metabolic pathway databases, MetaCyc, Reactome and the Small Molecule Pathway Database (21–23). This issue also features updates of two excellent databases of the active sites in various enzyme superfamilies, the Catalytic Site Atlas and the Structure-Function Linkage Database (SFLD) (24,25). There are also updates of the Carbohydrate-Active enzymes database (CAZy) and the protease database MEROPS, as well as new databases: EKPD, a database of eukaryotic protein kinases and protein phosphatases, and MultiTaskDB, a database of 'moonlighting' enzymes (26–29).

The increased interest in microbial genomics, fuelled in part by the Human Microbiome Project, led to several important developments in database construction. Many databases now emphasize improved genome annotation for a variety of microbes, including human pathogens (IMG, PATRIC, SEED), and for selected free-living microorganisms (CyanoBase, PortEco, Rhizobase, SubtiWiki). A number of databases, such as JGI'S IMG/M (30) and the newly created EBI metagenomics resource (31), strive to capture microbial diversity in natural environments. The rapid growth of completely or partially sequenced microbial genomes makes particularly important their proper classification, which increasingly relies on such resources as the Ribosomal Database Project (RDP), the SILVA/LTP project, BacDive at the DSMZ-German Collection of Microorganisms and Cell Cultures and the List of Prokaryotic Names with Standing in Nomenclature (LPSN) (32–35). The new MetaRef database collects from reference microbial genomes clade-specific genes that could be useful for taxonomic assignments of metagenomic reads (36).

One of the major highlights of this issue is the block of articles on the improved annotation of human genome and detailed analysis of genome variation and its potential clinical significance. These articles include, among others, updates on the Consensus CDS project, a collaborative effort to identify a core set of human protein-coding regions, and on the dbGaP, a database of genotyping results and related clinically relevant phenotypes (37,38). dbGaP contains openly available study data but requires pre-authorization for access to personal health information, such as individual-level data including phenotypic data tables and genotypes (see <http://www.ncbi.nlm.nih.gov/gap> for details). This issue also includes descriptions of several related databases: Locus Reference Genomic, a set of reference sequences for reporting of clinically relevant sequence variants; the NCBI's ClinVar, a

database documenting these clinically relevant sequence variants; the NHGRI GWAS Catalog, a curated resource of SNP-trait associations; Sanger Institute's DECIPHER, a database of pathogenic single nucleotide variants, indels and copy-number variants; and the Database of Genomic Variants (DGV) at the Toronto's Centre for Applied Genomics (39–43). There are also several more specialized databases (canSAR, DriverDB, FINDBase, HbVar, Lynx, NECTAR, Progenetix) that cover genetic defects leading to various human diseases, including cancer. In addition, three separate databases, Selectome, dbPSHP and 1000 Genomes Selection Browser, report the sites of likely positive selection in human genomes.

Model organism databases featured in this issue include regular updates from *Saccharomyces* Genome Database (SGD), WormBase, FlyBase, Mouse Genome Database (MGD), Mouse Gene Expression Database, Mouse Phenome Database and Vertebrate Genome Annotation (VEGA) database, as well as a description of the International Mouse Phenotyping Consortium (IMPC) web portal.

## NAR ONLINE MOLECULAR BIOLOGY DATABASE COLLECTION

The NAR online Molecular Biology Database Collection (freely available at <http://www.oxfordjournals.org/nar/database/a/>) has been updated by including the databases introduced in the 2014 NAR Database Issue. This list has been expanded by including such databases as CREDO, DoSA, DBATE, RedoxDB and TMBB-DB, whose descriptions had been published in our sister journals *Bioinformatics* and *Database* (Oxford) and are freely available online, as well as selected databases published elsewhere. The entire collection has been carefully curated by checking all non-responsive database websites; coordinators of such databases have been asked to confirm their commitment to maintaining their resources. Based on the received responses (or a lack thereof), URLs of 193 databases have been corrected and 24 obsolete databases have been removed from the list. As a result of these changes, the online collection now includes 1552 databases that are sorted into 14 categories and 41 subcategories.

Suggestions for inclusion in the collection of additional databases are welcome. They should include extended databases summaries in plain text, generally formatted according to the <http://www.oxfordjournals.org/nar/database/summary/1> template, including references to the published database descriptions freely available online, and should be addressed to XMFS at [xose.m.fernandez@gmail.com](mailto:xose.m.fernandez@gmail.com).

## ACKNOWLEDGEMENTS

The authors thank Sir Richard Roberts and Drs Alex Bateman and David Landsman for helpful comments and Dr Martine Bernardes-Silva and the Oxford University Press team led by Jennifer Boyd and Oliver Barham for their help in compiling this issue.

## FUNDING

NIH Intramural Research Program at the National Library of Medicine (to M.Y.G.). Funding for open access charge: Waived by Oxford University Press.

*Conflict of interest statement.* The authors' opinions do not necessarily reflect the views of their respective institutions. X.M.F.S. is an employee of Life Technologies.

## REFERENCES

- Mathelier, A., Zhao, X., Zhang, A.W., Parcy, F., Worsley-Hunt, R., Arenillas, D.J., Buchman, S., Chen, C.Y., Chou, A., Ienasescu, H. *et al.* (2014) JASPAR 2014: an extensively expanded and updated open-access database of transcription factor binding profiles. *Nucleic Acids Res.*, **42**, D142–D147.
- Teixeira, M.C., Monteiro, P.T., Guerreiro, J.F., Goncalves, J.P., Mira, N.P., Dos Santos, S.C., Cabrito, T.R., Palma, M., Costa, C., Francisco, A.P. *et al.* (2014) The YEASTRACT database: an upgraded information system for the analysis of gene and genomic transcription regulation in *Saccharomyces cerevisiae*. *Nucleic Acids Res.*, **42**, D161–D166.
- Kilic, S., White, E.R., Sagitova, D.M., Cornish, J.P. and Erill, J. (2014) CollecTF: a database of experimentally validated transcription factor-binding sites in Bacteria. *Nucleic Acids Res.*, **42**, D156–D160.
- Shazman, S., Lee, H., Socol, Y., Mann, R. and Honig, B. (2014) OnTheFly: a database of *Drosophila melanogaster* transcription factors and their binding sites. *Nucleic Acids Res.*, **42**, D167–D171.
- Chacon, D., Beck, D., Perera, D., Wong, J.W. and Pimanda, J.E. (2014) BloodChIP: a database of comparative genome-wide transcription factor binding profiles in human blood cells. *Nucleic Acids Res.*, **42**, D172–D177.
- Yang, L., Zhou, T., Dror, I., Mathelier, A., Wasserman, W.W., Gordan, R. and Rohs, R. (2014) TFBSshape: a motif database for DNA shape features of transcription factor binding sites. *Nucleic Acids Res.*, **42**, D148–D155.
- Robasky, K. and Bulyk, M.L. (2011) UniPROBE, update 2011: expanded content and search tools in the online database of protein-binding microarray data on protein-DNA interactions. *Nucleic Acids Res.*, **39**, D124–D128.
- Lobanov, M.Y., Sokolovskiy, I.V. and Galzitskaya, O.V. (2014) HRaP: database of occurrence of HomoRepeats and Patterns in proteomes. *Nucleic Acids Res.*, **42**, D273–D278.
- Di Domenico, T., Potenza, E., Walsh, I., Parra, R.G., Giollo, M., Minervini, G., Piovesan, D., Ihsan, A., Ferrari, C., Kajava, A.V. *et al.* (2014) RepeatsDB: a database of tandem repeat protein structures. *Nucleic Acids Res.*, **42**, D352–D357.
- Coimbatore Narayanan, B., Westbrook, J., Ghosh, S., Petrov, A.I., Sweeney, B., Zirbel, C.L., Leontis, N.B. and Berman, H.M. (2014) The Nucleic Acid Database: new features and capabilities. *Nucleic Acids Res.*, **42**, D114–D122.
- Chojnowski, G., Walen, T. and Bujnicki, J.M. (2014) RNA Bricks—a Database of RNA 3D motifs and their interactions. *Nucleic Acids Res.*, **42**, D123–D131.
- de Beer, T.A., Berka, K., Thornton, J.M. and Laskowski, R.A. (2014) PDBsum additions. *Nucleic Acids Res.*, **42**, D292–D296.
- Gibrat, J.F., Madej, T. and Bryant, S.H. (1996) Surprising similarities in structure comparison. *Curr. Opin. Struct. Biol.*, **6**, 377–385.
- Madej, T., Lanczycki, C.J., Zhang, D., Thiessen, P.A., Geer, R.C., Marchler-Bauer, A. and Bryant, S.H. (2014) MMDB and VAST+: tracking structural similarities between macromolecular complexes. *Nucleic Acids Res.*, **42**, D297–D303.
- Finn, R.D., Miller, B.L., Clements, J. and Bateman, A. (2014) iPfam: a database of protein family and domain interactions found in the Protein Data Bank. *Nucleic Acids Res.*, **42**, D364–D373.
- Mosca, R., Ceol, A., Stein, A., Olivella, R. and Aloy, P. (2014) 3did: a catalog of domain-based interactions of known three-dimensional structure. *Nucleic Acids Res.*, **42**, D374–D379.
- Kalathur, R.K., Pinto, J.P., Hernandez-Prieto, M.A., Machado, R.S., Almeida, D., Chaurasia, G. and Futschik, M.E. (2014) UniHI 7: an enhanced database for retrieval and interactive analysis of human molecular interaction networks. *Nucleic Acids Res.*, **42**, D408–D414.
- Blohm, P., Frishman, G., Smialowski, P., Goebels, F., Wachinger, B., Ruepp, A. and Frishman, D. (2014) Negatome 2.0: a database of non-interacting proteins derived by literature mining, manual annotation and protein structure analysis. *Nucleic Acids Res.*, **42**, D396–D400.
- Andreeva, A., Howorth, D., Chothia, C., Kulesha, E. and Murzin, A.G. (2014) SCOP2 prototype: a new approach to protein structure mining. *Nucleic Acids Res.*, **42**, D310–D314.
- Fox, N.K., Brenner, S.E. and Chandonia, J.-M. (2014) SCOPe: Structural Classification of Proteins extended, integrating SCOP and ASTRAL data and classification of new structures. *Nucleic Acids Res.*, **42**, D304–D309.
- Caspi, R., Altman, T., Billington, R., Dreher, K., Foerster, H., Fulcher, C.A., Holland, T.A., Keseler, I.M., Kothari, A., Kubo, A. *et al.* (2014) The MetaCyc database of metabolic pathways and enzymes and the BioCyc collection of Pathway/Genome Databases. *Nucleic Acids Res.*, **42**, D459–D471.
- Croft, D., Fabregat, M., Haw, R., Milacic, M., Weiser, J., Wu, G., Caudy, M., Garapati, P., Gillespie, M., Kaudar, M.R. *et al.* (2014) The reactome pathway knowledgebase. *Nucleic Acids Res.*, **42**, D472–D477.
- Jewison, T., Su, Y., Disfany, F.M., Liang, Y., Knox, C., Maciejewski, A., Poelzer, J., Huynh, J., Zhou, Y., Arndt, D. *et al.* (2014) SMPDB 2.0: big improvements to the small molecule pathway database. *Nucleic Acids Res.*, **42**, D478–D484.
- Akiva, E., Brown, S., Almonacid, D.E., Barber, A.E., Custer, A.F., Hicks, M.A., Huang, C.C., Lauck, F., Mashiyama, S.T., Meng, E.C. *et al.* (2014) The Structure-Function Linkage Database. *Nucleic Acids Res.*, **42**, D521–D530.
- Furnham, N., Holliday, G.L., de Beer, T.A.P., Jacobsen, J.O.B., Pearson, W.R. and Thornton, J.M. (2014) The catalytic site atlas 2.0: cataloguing catalytic sites and residues identified in enzymes. *Nucleic Acids Res.*, **42**, D485–D489.
- Lombard, V., Ramulu, H.G., Drula, E., Coutinho, P.M. and Henrissat, B. (2009) The Carbohydrate-Active Enzymes Database (CAZy) in 2013. *Nucleic Acids Res.*, **42**, D490–D495.
- Rawlings, N.D., Waller, M., Barrett, A.J. and Bateman, A. (2014) MEROPS: the database of proteolytic enzymes, their substrates and inhibitors. *Nucleic Acids Res.*, **42**, D503–D509.
- Wang, Y., Liu, Z., Cheng, H., Gao, T., Pan, Z., Yang, Q., Guo, A. and Xue, Y. (2014) EKPD: a hierarchical database of eukaryotic protein kinases and protein phosphatases. *Nucleic Acids Res.*, **42**, D496–D502.
- Hernández, S., Ferragut, G., Amela, I., Perez-Pons, J.A., Piñol, J., Mozo-Villarias, A., Cedano, J. and Querol, E. (2014) MultitaskProtDB: a database of multitasking proteins. *Nucleic Acids Res.*, **42**, D517–D520.
- Markowitz, V.M., Chen, I.M., Chu, K., Szeto, E., Palaniappan, K., Pillay, M., Ratner, A., Huang, J., Pagani, I., Tringe, S. *et al.* (2014) IMG/M 4 version of the integrated metagenome comparative analysis system. *Nucleic Acids Res.*, **42**, D568–D573.
- Hunter, S., Corbett, M., Denise, H., Fraser, M., Gonzalez-Beltran, A., Hunter, C., Jones, P., Leinonen, R., McAnulla, C., Maguire, E. *et al.* (2014) EBI metagenomics—a new resource for the analysis and archiving of metagenomic data. *Nucleic Acids Res.*, **42**, D600–D606.
- Cole, J.R., Wang, Q., Fish, J.A., Chai, B., McGarrell, D.M., Sun, Y., Brown, C.T., Porras-Alfaro, A., Kuske, C.R. and Tiedje, J.M. (2014) The Ribosomal Database Project: data and tools for high throughput rRNA analysis. *Nucleic Acids Res.*, **42**, D633–D642.
- Yilmaz, P., Parfrey, L.W., Yarza, P., Gerken, J., Prüsse, E., Quast, C., Schweer, T., Peplies, J., Ludwig, W. and Glöckner, F.O. (2014) The SILVA and 'The All Species Living Tree (LTP)' taxonomic frameworks. *Nucleic Acids Res.*, **42**, D643–D648.
- Söhngen, C., Bunk, B., Podstawka, A., Gleim, D. and Overmann, J. (2014) BacDive—The Bacterial Diversity metadatabase. *Nucleic Acids Res.*, **42**, D592–D599.
- Parte, A. (2014) LPSN—list of prokaryotic names with standing in nomenclature. *Nucleic Acids Res.*, **42**, D613–D616.

36. Huang,K., Brady,A., Mahurkar,A., White,O., Gevers,D., Huttenhower,C. and Segata,N. (2014) MetaRef: a pan-genomic database for comparative and community microbial genomics. *Nucleic Acids Res.*, **42**, D617–D624.
37. Farrell,C.M., O’Leary,N.A., Harte,R.A., Loveland,J.E., Wilming,L.G., Wallin,C., Diekhans,M., Barrell,D., Searle,S.M., Aken,B. *et al.* (2014) Current status and new features of the Consensus Coding Sequence database. *Nucleic Acids Res.*, **42**, D865–D872.
38. Tryka,K.A., Hao,L., Sturcke,A., Jin,Y., Wang,Z.Y., Ziyabari,L., Lee,M., Popova,N., Sharopova,N., Kimura,M. *et al.* (2014) NCBI’s Database of Genotypes and Phenotypes: dbGaP. *Nucleic Acids Res.*, **42**, D975–D979.
39. MacArthur,J.A.L., Morales,J., Tully,R.E., Astashyn,A., Gil,L., Bruford,E.A., Larsson,P., Flicek,P., Dalgleish,R., Maglott,D.R. *et al.* (2014) Locus Reference Genomic: reference sequences for the reporting of clinically relevant sequence variants. *Nucleic Acids Res.*, **42**, D873–D878.
40. Welter,D., MacArthur,J., Morales,J., Burdett,T., Hall,P., Junkins,H., Klemm,A., Flicek,P., Manolio,T., Hindorf,L. *et al.* (2014) The NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Res.*, **42**, D1001–D1006.
41. Landrum,M., Lee,J.M., Riley,G., Jang,W., Rubinstein,W.S., Church,D.M. and Maglott,D. (2014) ClinVar: public archive of relationships among sequence variation and human phenotype. *Nucleic Acids Res.*, **42**, D980–D985.
42. Bragin,E., Chatzimichali,E.A., Wright,C.F., Hurles,M.E., Firth,H.V., Bevan,A.P. and Swaminathan,G.J. (2014) DECIPHER: database for the interpretation of phenotype-linked plausibly pathogenic sequence and copy-number variation. *Nucleic Acids Res.*, **42**, D993–D1000.
43. Macdonald,J.R., Ziman,R., Yuen,R.K., Feuk,L. and Scherer,S.W. (2014) The Database of Genomic Variants: a curated collection of structural variation in the human genome. *Nucleic Acids Res.*, **42**, D986–D992.