

VRG: a database of vascular dysfunctions related genes

June 1, 2006

Sara Zanivan^{1,3a}, Davide Corà^{2,3a}, Michele Caselle² and Federico Bussolino¹

¹ Department of Oncological Sciences and Institute for Cancer Research and Treatment, University of Torino School of Medicine, Strada Provinciale 142 Km 3.95, 10060 Candiolo (TO), Italy.

² Department of Theoretical Physics, University of Torino and INFN, Istituto Nazionale di Fisica Nucleare, Via P. Giuria 1 - I 10125 Turin, Italy.

^aCorresponding authors, email: sara.zanivan@ircc.it, cora@to.infn.it.

³ Equal contribution

Abstract

Heart and vascular defects occur in a large number of hereditary and sporadic human diseases as result of a complex interplay of genetic factors. Since genome sequencing of many organisms disclosed similarities among genomes, animal models are crucial for the discovery of genes involved in those pathological processes. Therefore we propose *VRG* database, in which human data have been manually managed and integrated with mouse information in order to create a catalogue of genes involved in vascular diseases.

Keywords: genome, database, human vascular diseases, animal models.

1 Introduction

Heart and blood vessels dysfunctions are predominant causes of disability and death in humans [1]. In 2003, the World Health Organization [2] estimated that 16.7 million people die each year of cardiovascular (CV) diseases, covering 29% of all deaths around the globe, thus highlighting the value of devoting resources to study the pathogenesis of these settings.

Phenotypic alterations result from a complex network of genetic modifications [3] and a huge effort has to be done to better understand their mechanisms at the molecular scale. Even though experimental approaches could be relevant [4, 5, 6, 7], computer science provides an important support as well. It is essential, in fact, a comprehensive collection and integration of public available data generated by the research community [8, 9, 10] to allow an easy access and recovery of information. To date, different databases collect information of known genes involved in human vascular diseases. For instance, the Online Mendelian Inheritance in Man (OMIM) [10] catalogues phenotypes and genotypes classified by phenotypic features and mutated gene while the Cardiovascular Comparative Genomic Database (CVCGD) [11, 12] collects well known and comparatively annotated cardiovascular genes. However, it is not straightforward to recover a complete list of genes related to vascular dysfunctions. OMIM, for example, allows only a free-text search of clinical synopsis thus rendering difficult the selection of genes involved in vascular dysfunctions. Therefore, in order to generate an accurate collection of human genes, it would be useful to review all phenotypic OMIM entries, select those directly or indirectly related to vascular dysfunctions, and collect the corresponding genes into a catalogue as complete as possible.

The increasing availability of genomic sequences from many organisms [13, 14, 15] provided the opportunity of orthologous-sequences comparison [16, 17].

This analysis revealed that sequences performing important functions are frequently conserved among evolutionarily distant species [18, 19]. Moreover, animal models yield the identification of novel genes that often cause defects when mutated in humans and therefore they are used as a tool to investigate both mono and multifactorial human diseases [6, 8, 20, 21, 22]. Because they can be subject of large scale mutation screenings [23, 24], animals are a simple model to apply genomic strategies. For instance, knock-out (KO) mice proved useful in elucidating gene function and provided many insights into human biology and diseases [19, 23, 24]. Therefore, genome-wide collections of KO mice can significantly contribute to biomedical discovery [23]. A prominent example is represented by the Mouse Genome Database (MGD) [9, 25], a publicly available resource of KO mice which collects genomic, genetic, functional, and phenotypic data about mouse genes in order to identify candidate genes associated with complex phenotypes [26]. Since MGD has generated an ontology in order to catalogue the altered phenotypes of mutant mice, it results simple to select genes related to particular vascular dysfunctions.

Gene regulation is one of the major mechanisms underpinning the correct cardiovascular system morphogenesis and function [1, 24] and it has been demonstrated that similar molecular mechanisms are involved in its regulation both in physiological (i.e. during embryogenesis) and pathological conditions [27]. Therefore, an in depth analysis of the embryonic vascular development could yield useful information to understand the pathogenesis of vascular diseases in adults [28]. Also in this case animal models, i.e. zebrafish [29] or xenopus [30], will be crucial to improve our knowledge.

In order to reach a more comprehensive knowledge on vascular dysfunctions it could be useful to generate a complete list of human and animal genes involved in vascular phenotypic alterations. Due to the huge amount of information collected in different databases, computerized systems could simplify data recovery,

management and analysis. Therefore, we generated VRG [31], a publicly available database aimed at integrating mouse and human information, which were manually curated to create a wide catalogue of genes involved in vascular diseases.

2 DATABASE

2.1 Human section

We started with the information contained in the ftp section of the OMIM database [3] (March 2005 version). We selected two different sources of data:

- the morbidmap file
- the complete OMIM report flat-file

In the first one all the possible diseases presented in the OMIM catalogue are annotated with the corresponding related gene ids. Overall, the correspondence between a certain genetic disease and a certain pool of genes is a one-to-many relationship, meaning that it is possible to have one or more gene for a specific disease as well one or more diseases associated to the same gene id.

The second one is essentially constituted by the transposition in an ASCII computerized flat file version of the original catalogue provided by Victor McKusick's book, Mendelian Inheritance in Man. For our purposes it is necessary to specify that:

- each OMIM entry is given a unique six-digit number whose first digit indicates the mode of inheritance of the gene involved;
- each OMIM entry is characterized by a special field, called Clinical Synopsis (CS) which reports a description of the observed phenotype for the corresponding disease.

The OMIM database is actually not meant to be organized into a relational database, so an automated managing of the information contained is not straightforward. We choose the following algorithm:

For each entry of the morbid-map file, in which there is a clear correspondence between a certain gene and a certain disease, we selected the correspond-

ing CS from the OMIM file. All the possible CS collected in this catalogue were then manually-curated and divided into four main categories, where possible:

1. VASCULAR: alterations related to the vascular system
2. NEURO: alterations related to the nervous system
3. NEURO-VASCULAR: alterations related both to the nervous and the vascular system
4. INTERESTING: metabolic and/or mitochondrial alterations, so important to be related with neuro-vascular disturbs.

We then searched the Ensembl database (version 25) [32] to make a direct connection, where possible, between the gene ids internally used by OMIM and human gene ids provided by Ensembl. If available, we selected the corresponding mouse and zebrafish orthologous for each of the disease genes.

We finally grouped and stored into a MySQL relational database such information, namely:

- disease-name: name of the disease according to OMIM
- disease-id: disease identifier according to OMIM
- gene-name: gene name according to OMIM
- ENSG-id: human gene name according to Ensembl
- gene-mol-descr: description of the molecular activity of the gene
- gene-id: gene identifier according to OMIM
- location: location of the gene involved in the disease
- cs: the Clinical Synopses for the disease
- mouse-ortholog: mouse ortholog according to Ensembl
- zebrafish-ortholog: zebrafish ortholog according to Ensembl

2.2 Mouse section

The second source of data for our work was the Mouse Genome Informatics - MGI database [25], 3.44 version. The MGI database is a collection of a large amount of data related to all the aspects of the mouse biology and genomics. In particular we concentrated our attention on the manually curated list of mouse KO experiments recorded in the ftp section of the database. In this section, each mouse gene is annotated together with the outcome of the corresponding KO experiment, if available, and the complete list of phenotypes is then organized in a fixed and controlled vocabulary provided by the curators with unique identifiers. The vocabulary of phenotypes is internally organized in an ontology-based way. Unlike from the human case, in which a similar collection of information does not exist, the list of genes / KO phenotype can be hence handled by computational means in a rigorous manner. The MGI database also provides a connection with external databases, in particular with the Ensembl database and also includes annotations of human/mouse orthology.

For our purposes, we extracted from the MGI a list of genes related to three particular different phenotype annotations:

- MP:0005385 cardiovascular
- MP:0003631 nervous system
- MP:0005386 behavior/neurological phenotype

We identified the corresponding human and zebrafish orthologous for each of the genes selected as described above.

We finally grouped and stored into a MySQL relational database such information, namely:

- MGI-id: internal gene identifier of the MGI
- gene-id: gene identifier according to MGI

- ENSMUSG-id: mouse gene name according to Ensembl
- phenotype-id: knockout phenotype according to MGI
- human-ortholog: human ortholog according to Ensembl
- zebrafish-ortholog: zebrafish ortholog according to Ensembl

for each of three phenotypes previously selected.

Once equipped with those two relational databases, respectively built from the OMIM and MGI data, we developed a set of tools devoted to the automatic extraction of data (queries) from these databases themselves and to the automatic generation of a set of html web pages building up the VRG Disease Database, available at: <http://www.to.infn.it/ftbio/VRG-database/main.html>.

In Fig. 1 and 2 we reported two snapshots of the VRG database.

3 Perspectives

Here we propose a publicly available web based database collecting vascular related genes. Earlier studies have shown the advantage of comparative approaches to better understand human biological processes and their altered counterparts [19, 23, 24]. By using OMIM [3] and MGD [9, 25], we integrated human and mouse information to create the Vascular Related Genes (VRG) database, which is structured to easily access and recover information about genes involved in vascular dysfunctions from multiple species. This allows determining similar pathogenetic mechanisms or selected specificities characterizing the role of a gene in different species or combining them thus having a more complete view of the molecular mechanisms regulating disease onset and maintaining. Moreover, it allows the exploitation of lower organisms to retrieve information that could give useful insights on human genes function and eventually to discover new diagnostic tools or therapeutic targets to treat vascular dysfunctions. However, to reach these goals additional work is required to include data coming from other animal models. Indeed, we are going to integrate zebrafish data for its qualities as a model for studies of vertebrate genetics including cardiovascular diseases [29].

Based on the new emerging parallels in the development of vascular and nervous systems [33, 34, 35], we decided to integrate vascular and neural phenotypes within the VRG database. To this aim, human and mouse genes determining phenotypic alterations of the nervous system were included in the VRG database. In particular, genes responsible of both nervous and vascular dysfunctions are collected into a separate list and could be considered candidates as new molecules regulating nerves and vessels behavior.

Hence, we propose the VRG database as a comprehensive and easily accessible catalogue of genes related to vascular and nervous dysfunctions.

4 Acknowledgments

We thank Guido Serini and Marco Arese for critical reading the manuscript. This study was supported by Associazione Italiana per la Ricerca sul Cancro, Istituto Superiore di Sanit . (AIDS National Projects), Ministero dell Universit  e della Ricerca (MIUR) (60% and PRIN 2004 projects), Ministero della Salute (Ricerca Finalizzata 2002,2003 and 2004), Regione Piemonte, European Community (LSHM-CT-2003-503251) (<http://www.engv.org>), FIRB (Fondo per gli Investimenti della Ricerca di Base) from the Italian Ministry of the University and Scientific Research, number RBNE03B8KK006.

References

- [1] Kumar V, Abbas AK
N. F: Robbins and Cotran pathologic basis of disease, 7th.
edn: Elsevier Saunders; 2004.
- [2] WHO [<http://www.who.int>]
- [3] Boyadjiev SA, Jabs EW
Online Mendelian Inheritance in Man (OMIM) as a knowledgebase for human developmental disorders.
Clin Genet 2000, 57(4):253-266.
- [4] Cullen P, Rauterberg J, Lorkowski S
The pathogenesis of atherosclerosis.
Handb Exp Pharmacol 2005(170):3-70.
- [5] Alvarez-Garcia I, Miska EA
MicroRNA functions in animal development and human disease.
Development 2005, 132(21):4653-4662.
- [6] D'Orleans-Juste P, Honore JC, Carrier E, Labonte J
Cardiovascular diseases: new insights from knockout mice.
Curr Opin Pharmacol 2003, 3(2):181-185.
- [7] Lambrechts D, Carmeliet P
Genetics in zebrafish, mice, and humans to dissect congenital heart disease: insights in the role of VEGF.
Curr Top Dev Biol 2004, 62:189-224.
- [8] Sprague J, Doerry E, Douglas S, Westerfield M
The Zebrafish Information Network (ZFIN): a resource for genetic, genomic and developmental research.
Nucleic Acids Res 2001, 29(1):87-90.

- [9] Eppig JT, Bult CJ, Kadin JA, Richardson JE, Blake JA, Anagnostopoulos A, Baldarelli RM, Baya M, Beal JS, Bello SM et al
The Mouse Genome Database (MGD): from genes to mice—a community resource for mouse biology.
Nucleic Acids Res 2005, 33(Database issue):D471-475.
- [10] OMIM [<http://www.ncbi.nlm.nih.gov/omim/>]
- [11] CVCGD [<http://pga.lbl.gov/cvcgd.html>]
- [12] O’Kane DJ, Weinshilboum RM, Moyer TP
Pharmacogenomics and reducing the frequency of adverse drug events.
Pharmacogenomics 2003, 4(1):1-4.
- [13] International Human Genome Sequencing Consortium.
Finishing the euchromatic sequence of the human genome.
Nature 2004, 431(7011):931-945.
- [14] Venter JC, Adams MD, Myers EW, Li PW, Mural RJ, Sutton GG, Smith HO, Yandell M, Evans CA, Holt RA et al
The sequence of the human genome.
Science 2001, 291(5507):1304-1351.
- [15] Waterston RH, Lindblad-Toh K, Birney E, Rogers J, Abril JF, Agarwal P, Agarwala R, Ainscough R, Alexandersson M, An P et al
Initial sequencing and comparative analysis of the mouse genome.
Nature 2002, 420(6915):520-562.
- [16] Miller W, Makova KD, Nekrutenko A, Hardison RC
Comparative genomics.
Annu Rev Genomics Hum Genet 2004, 5:15-
- [17] Varki A, Altheide TK
Comparing the human and chimpanzee genomes: searching for needles in

- a haystack.*
Genome Res 2005, 15(12):1746-1758.
- [18] Curry BB
Animal models used in identifying gender-related differences.
Int Jour toxicology 2001, 20:153-160.
- [19] Paigen K, Eppig JT
A mouse phenome project.
Mamm Genome 2000, 11(9):715-717.
- [20] Pennacchio LA, Rubin EM
Comparative genomic tools and databases: providing insights into the human genome.
J Clin Invest 2003, 111(8):1099-1106.
- [21] Towbin JA, Casey B, Belmont J
The molecular basis of vascular disorders. Am J Hum Genet 1999, 64(3):678-684.
- [22] Lopez-Bigas N, Blencowe BJ, Ouzounis CA
Highly consistent patterns for inherited human diseases at the molecular level.
Bioinformatics 2006, 22(3):269-277.
- [23] Austin CP, Battey JF, Bradley A, Bucan M, Capecchi M, Collins FS, Dove WF, Duyk G, Dymecki S, Eppig JT et al
The knockout mouse project.
Nat Genet 2004, 36(9):921-924.
- [24] Chien KR
Genomic circuits and the integrative biology of cardiac diseases.
Nature 2000, 407(6801):227-232.
- [25] MGI [<http://www.informatics.jax.org>]

- [26] Blake JA, Eppig JT, Bult CJ, Kadin JA, Richardson JE
The Mouse Genome Database (MGD): updates and enhancements.
Nucleic Acids Res 2006, 34(Database issue):D562-567.
- [27] Folkman J
Angiogenesis in cancer, vascular, rheumatoid and other disease.
Nat Med 1995, 1(1):27-31.
- [28] Ware JA, Simons M
Angiogenesis and cardiovascular disease.
New York: Oxford University Press; 1999.
- [29] Weinstein BM, Fishman MC
Cardiovascular morphogenesis in zebrafish.
Cardiovasc Res 1996, 31 Spec No:E17-24.
- [30] Ny A, Autiero M, Carmeliet P
Zebrafish and Xenopus tadpoles: small animal models to study angiogenesis and lymphangiogenesis.
Exp Cell Res 2006, 312(5):684-693.
- [31] VRG [http://www.to.infn.it/ftbio/VRG_database/main.html]
- [32] Birney E, Andrews D, Caccamo M, Chen Y, Clarke L, Coates G, Cox T, Cunningham F, Curwen V, Cutts T et al
Ensembl 2006.
Nucleic Acids Res 2006, 34(Database issue):D556-561.
- [33] Serini G, Bussolino F
Common cues in vascular and axon guidance.
Physiology (Bethesda) 2004, 19:348-354.
- [34] Weinstein BM
Vessels and nerves: marching to the same tune.
Cell 2005, 120(3):299-302.

- [35] Autiero M, De Smet F, Claes F, Carmeliet P
Role of neural guidance signals in blood vessel navigation.
Cardiovasc Res 2005, 65(3):629-638.

VASCULAR diseases

disease_NAME	OMIM_disease_ID	gene_NAME	Ensembl_Id	MIM_gene_ID	LOCATION	gene_MOL_DESCR	MOUSE_ORTHO	ZEBRAFISH_ORTHO
1 Adrenal hyperplasia, congenital, due to 11-beta-hydroxylase deficiency (3)	-	CYP11B1, P450C11, FHI	ENS000000160882	202010	8q21	cytochrome P450, family 11, subfamily B, polypeptide 1	-	-
2 Afriomerenia, 202400 (3)	202400	FGA	ENS000000171560	134820	4q28	fibronogen, A alpha polypeptide	ENS00000000028001	ENS00000000020741
3 Afriomerenia, congenital, 202400 (3)	202400	FGB	ENS000000171564	134830	4q28	fibronogen, B beta polypeptide	ENS00000000033831	ENS0000000008969
4 Alkaptonuria, 203300 (3)	203300	HGD, AKU	ENS000000113924	607474	3q21-q23	homogentisate 1,2-dioxygenase (homogentisate oxidase)	ENS00000000022821	ENS00000000017934
5 Anemia, Diamond-Blackfan, 105630 (2)	105630	DBA2	-	606129	8p23.3-p22	-	-	-
6 Anemia, Diamond-Blackfan, 105630 (3)	105630	RPS19, DBA	ENS000000103372	603474	19q13.2	ribosomal protein S19	ENS0000000004952	ENS00000000030602
7 Anemia, hemolytic, due to PK deficiency (3)	-	PKLR, PK1	ENS000000145627	266200	1q21	pyruvate kinase, liver and RBC	ENS00000000041237	-
8 Anemia, hemolytic, due to UMPH1 deficiency; 266120 (3)	266120	NTSC3, UMPH1, PSN1	ENS000000122643	606224	7p15-p14	5'-nucleotidase, cytosolic III	ENS00000000029780	-
9 Ankylosing spondylitis (2)	-	AS, ANS	-	106300	6p21.3	ankylosing spondylitis	-	-
10 Antithrombin III deficiency (3)	-	AT3	-	107300	1q23-q25	-	-	-
11 Arrhythmogenic right ventricular dysplasia 2, 600996 (3)	600996	RYR2, VTSIP	ENS000000198626	180902	1q42.1-q43	ryanodine receptor 2 (cardiac)	ENS00000000021313	ENS00000000011422
12 Arrhythmogenic right ventricular dysplasia-1 (2)	-	ARVD1	-	107970	14q23-q24	arrhythmogenic right ventricular dysplasia 1	-	-
13 Arrhythmogenic right ventricular dysplasia-2 (2)	-	ARVD2	-	600996	1q42-q43	arrhythmogenic right ventricular dysplasia 2	-	-
14 Arrhythmogenic right ventricular dysplasia-3 (2)	-	ARVD3	-	602036	14q12-q22	arrhythmogenic right ventricular dysplasia 3	-	-
15 Arrhythmogenic right ventricular dysplasia-4 (2)	-	ARVD4	-	602037	2q32.1-q32.3	arrhythmogenic right ventricular dysplasia 4	-	-
16 Atransferrinemia, 209300 (3)	209300	TF	ENS000000091513	190000	3q21	transferrin	ENS00000000032354	ENS00000000016771
17 Atrial fibrillation, familial (2)	-	ATFB2	-	608383	10q22-q24	-	-	-
18 Atrial fibrillation, familial, 607554 (3)	607554	KCNQ1, KCNA9, LQT1, KYLQTT1, ATFB1	ENS000000039118	607542	11p15.5	potassium voltage-gated channel, KCQT-like subfamily, member 1	ENS00000000009545	ENS00000000024926
19 Atrioventricular canal defect, 600309 (2)	600309	AVSD1, AVCD	-	606215	1p31-p21	atrioventricular septal defect 1	-	-
20 Barth syndrome, 302060 (3)	302060	TAFZ, EPEZ, BTHS, CMD3A	ENS000000102125	300394	Xq28	tafazzin (cardiomyopathy, dilated 3A (X-linked), endocardial fibroelastosis 2, Barth syndrome)	ENS00000000009992	ENS00000000041421
21 Barthel syndrome, type 3, 607364 (3)	607364	CLCNKB	ENS000000184908	602023	1p36	chloride channel Kb	-	-
22 Beckwith-Wiedemann syndrome, 130630 (3)	130630	CDKN1C, KIP2, PWS	ENS000000129737	600836	11p15.5	cyclin-dependent kinase inhibitor 1C (p57, Kip2)	ENS00000000037664	ENS00000000010878

Figure 1: Snapshot-1 of a page of the database.

mouse CARDIOVASCULAR genes

<u>MGI_id</u>	<u>gene_id</u>	<u>gene_name</u>	<u>ENSMUSG_id</u>	<u>phenotype_id</u>	<u>phenotype_description</u>	<u>Human ortholog</u>	<u>Zebrafish ortholog</u>
1	MGI:101762	Elk3	ENSMUSG00000003398	MP:0005385	cardiovascular system phenotype	ENSG00000111145	ENSDBARG00000019688
2	MGI:101802	F2r	ENSMUSG00000004876	MP:0005385	cardiovascular system phenotype	ENSG00000181104	ENSDBARG00000003523
3	MGI:101876	Tead1	ENSMUSG00000005320	MP:0005385	cardiovascular system phenotype	ENSG00000187079	ENSG00000187079
4	MGI:101900	Mmp14	ENSMUSG00000000957	MP:0005385	cardiovascular system phenotype	ENSG00000157227	ENSDBARG00000002235
5	MGI:101924	Sic12a2	ENSMUSG00000002497	MP:0005385	cardiovascular system phenotype	ENSG00000064651	ENSDBARG00000010640
6	MGI:102539	Tbx6	ENSMUSG00000003699	MP:0005385	cardiovascular system phenotype	ENSG00000149922	ENSDBARG00000011785
7	MGI:102541	Tbx5	ENSMUSG000000018283	MP:0005385	cardiovascular system phenotype	ENSG00000089225	ENSDBARG00000024894
8	MGI:102548	Tsc2	ENSMUSG00000002496	MP:0005385	cardiovascular system phenotype	ENSG00000103197	ENSG00000103197
9	MGI:102556	Tbx4	ENSMUSG00000000094	MP:0005385	cardiovascular system phenotype	ENSG00000121075	ENSDBARG00000011939
10	MGI:102643	Myh11	ENSMUSG000000018820	MP:0005385	cardiovascular system phenotype	ENSG00000133392	ENSDBARG00000009782
11	MGI:102700	Itga7	ENSMUSG00000002548	MP:0005385	cardiovascular system phenotype	ENSG00000135424	ENSG00000135424
12	MGI:102709	Cav1	ENSMUSG000000007655	MP:0005385	cardiovascular system phenotype	ENSG00000103974	ENSG00000103974
13	MGI:102720	Ednrb	ENSMUSG00000002122	MP:0005385	cardiovascular system phenotype	ENSG00000136160	ENSDBARG00000003773
14	MGI:102768	Mfge8	ENSMUSG000000039605	MP:0005385	cardiovascular system phenotype	ENSG00000072958	ENSDBARG00000020838
15	MGI:102776	Aplm1	ENSMUSG000000003033	MP:0005385	cardiovascular system phenotype	ENSG00000135363	ENSDBARG000000039673
16	MGI:102811	Lmo2	ENSMUSG000000032698	MP:0005385	cardiovascular system phenotype	ENSG00000134571	ENSDBARG00000011615
17	MGI:102844	Myh9c3	ENSMUSG000000002100	MP:0005385	cardiovascular system phenotype	ENSG00000003847	ENSG00000003847
18	MGI:102889	Cspg2	ENSMUSG000000021614	MP:0005385	cardiovascular system phenotype		

Figure 2: Snapshot-2 of a page of the database.