

# The structural effects of mutations can aid in differential phenotype prediction of beta-myosin heavy chain (Myosin-7) missense variants

Nouf S. Al-Numair<sup>1</sup>, Luis Lopes<sup>2</sup>, Petros Syrris<sup>2</sup>, Lorenzo Monserrat<sup>3</sup>, Perry Elliott<sup>2</sup> and Andrew C.R. Martin<sup>1,\*</sup>

<sup>1</sup>Institute of Structural and Molecular Biology, Division of Biosciences, University College London, Darwin Building, Gower Street, London WC1E 6BT; <sup>2</sup>Institute of Cardiovascular Science, UCL, London; <sup>3</sup>Complejo Hospitalario Universitario de A Coruña, Instituto de Investigación Biomédica Coruña, Spain

Received on XXXXX; revised on XXXXX; accepted on XXXXX

Associate Editor: XXXXXXXX

## ABSTRACT

**Motivation:** High-throughput sequencing platforms are increasingly used to screen patients with genetic disease for pathogenic mutations, but prediction of the effects of mutations remains challenging. Previously we developed SAAPdap (Single Amino Acid Polymorphism Data Analysis Pipeline) and SAAPpred (Single Amino Acid Polymorphism Predictor) that use a combination of rule-based structural measures to predict the effect of missense genetic variants on protein function. Here we investigate whether the same methodology can be used to develop a differential phenotype predictor able to distinguish between the two major clinical phenotypes (hypertrophic cardiomyopathy, HCM, and dilated cardiomyopathy, DCM) associated with mutations in the beta-myosin heavy chain (*MYH7*) gene product (Myosin-7).

**Results:** A random forest predictor trained on rule-based structural analyses together with structural clustering data gave a Matthews' correlation coefficient (MCC) of 0.53 (accuracy, 75%). A *post hoc* removal of machine learning models that performed particularly badly, increased the performance (MCC=0.61, Acc=79%). This suggests that methods used for pathogenicity prediction can be extended for use in differential phenotype prediction.

**Contact:** andrew@bioinf.org.uk –or– andrew.martin@ucl.ac.uk

**Supplementary Information:** Supplementary File 1; Supplementary File 2.

## 1 INTRODUCTION

Mutations in proteins generally result in loss of function, but in some cases can lead to a gain of function. Generally this is not gain of a *novel* function, but an increased activity, often through loss of some type of control mechanism. In general predictors of pathogenicity do not try to distinguish between loss-of-function and gain-of-function mutations, but simply predict whether or not there will be *some* effect on function leading to a pathogenic state.

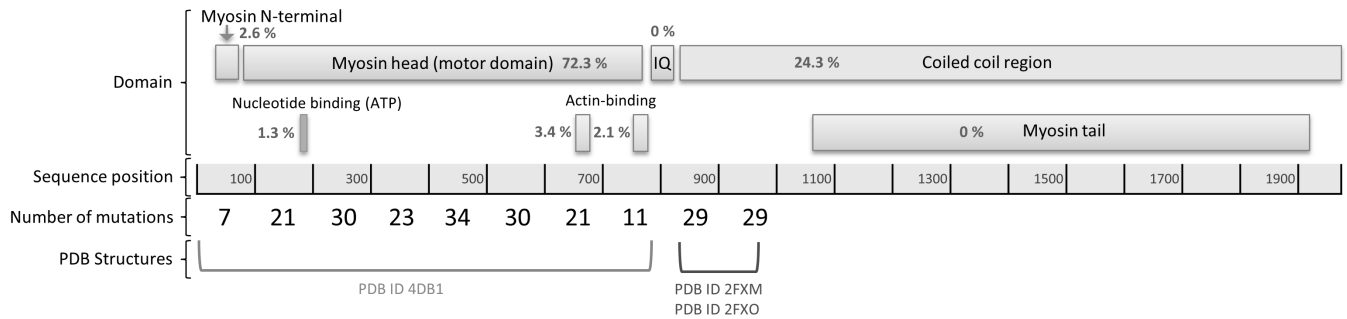
In some cases however, the situation is more complex, with mutations in a single protein leading to a number of distinct

phenotypes. For example, inherited heart muscle diseases, or cardiomyopathies, which are a major cause of sudden cardiac death in the young and an important cause of heart failure at all ages (Hughes and McKenna, 2005) are, as a group, very heterogeneous in genotype and phenotype. Radically different phenotypes can result from mutations in the same gene (Arad *et al.*, 2002).

The widespread application of Single Nucleotide Polymorphism (SNP) chips and high-throughput sequencing has generated an urgent need for informatics tools that can help predict the effects of the many sequence variants that these platforms identify. More than 20 groups have devised methods to predict whether a given mutation will have a deleterious effect (Yue *et al.*, 2006; Uzun *et al.*, 2007; Yip *et al.*, 2004; Dantzer *et al.*, 2005; Karchin *et al.*, 2005; Stitzel *et al.*, 2004; Reumers *et al.*, 2005; Bao *et al.*, 2005; Reva *et al.*, 2011; Schwarz *et al.*, 2010; Bromberg and Rost, 2007; Bromberg *et al.*, 2008; González-Pérez and López-Bigas, 2011; Shihab *et al.*, 2013; Al-Numair and Martin, 2013; Li *et al.*, 2009; Kircher *et al.*, 2014; Calabrese *et al.*, 2009; Worth *et al.*, 2011; Yates *et al.*, 2014), the best known methods being SIFT (Ng and Henikoff, 2003), an evolutionary method which calculates a sophisticated residue conservation score from multiple alignment, and PolyPhen-2 (Adzhubei *et al.*, 2010, 2013), which uses machine learning on a set of eight sequence- and three structure-based features. A more complete list of methods is provided on our web site at <http://www.bioinf.org.uk/saap/methods/>. However, these tools are generally not validated for individual diseases where most available datasets are too small to train machine-learning methods and tend to be heavily unbalanced. An additional problem is that it is often very difficult to obtain reliable validated data on neutral mutations. One of the few cases where a predictor has been produced for an individual class of proteins is the work on voltage-gated potassium channels by Stead *et al.* (2011).

Attempting to distinguish between mutations in a single protein that result in different pathogenic phenotypes is a difficult problem that, unlike pathogenicity prediction, has not been widely addressed. There have been a small number of attempts to distinguish loss-of-function and gain-of-function mutations at a molecular level, but (as stated above) typically

\*to whom correspondence should be addressed



**Fig. 1.** Annotated regions of the Myosin-7 sequence. Regions for which structures are known are indicated, together with the number of known mutations from Table 2 in each 100 amino acids of the sequence. The percentage of the total 235 mutations that map to known structure in each region is indicated: two of the mutations (at positions 82 and 838) do not correspond to any annotated regions. **Myosin N-terminal** Pfam annotation, residues 34–75; **Myosin head (motor domain)** Pfam and InterPro annotation, residues 85–778; **IQ motif** UniProtKB/SwissProt and InterPro annotation, residues 781–810, SMART annotation, residues 780-802; **Coiled coil region** UniProtKB/SwissProt annotation, residues 839–1935, SMART annotation, residues 841-1927; **Nucleotide binding (ATP) region** UniProtKB/SwissProt annotation, residues 178–185; **Actin-binding region** UniProtKB/SwissProt annotation, residues 655–677; **Actin-binding region** UniProtKB/SwissProt annotation, residues 757–771; **Myosin tail** Pfam and InterPro annotation, residues 1068–1926.

gain-of-function mutations result from loss of regulation making the protein constitutively active. For example, mutations that cause the VAB-1 tyrosine kinase to become constitutively active cause severe axon defects (Mohamed and Chin-Sang, 2006). Some of the challenges in the ‘Comparative Assessment of Genome Interpretation’ (CAGI) experiment have required the prediction of the level of enzyme activity (e.g. [genomeinterpretation.org/content/4-NAGLU](http://genomeinterpretation.org/content/4-NAGLU)) and some have been related to familial combined hyperlipidemia or channelopathies ([genomeinterpretation.org/content/FCH](http://genomeinterpretation.org/content/FCH), [genomeinterpretation.org/content/scn5a](http://genomeinterpretation.org/content/scn5a)), but, to our knowledge, there have been no clear cases where predictions have focused on mutations in the same protein resulting in different phenotypes other than through loss of function vs. loss of regulation.

Initially our own focus was on trying to understand the effects that mutations have on protein structure and then to use this information to compare the effects of non-pathogenic mutations and pathogenic deviations (Hurst *et al.*, 2009). Our approach has been to map mutations onto protein structure and to perform a rule-based analysis of the likely structural effects of these mutations in order to ‘explain’ the known functional effect (if any) of the mutation. Since we map mutations to structure, we only consider mutations in proteins for which a structure has been solved. With the recent growth in the amount of mutation data, we have moved from updating a database of analysis of mutations, to providing a server (SAAPdap — Single Amino Acid Polymorphism Data Analysis Pipeline) for analysis of the effects of mutations (<http://www.bioinf.org.uk/saap/dap/>) (Al-Numair and Martin, 2013). We have also developed SAAPpred (Single Amino Acid Polymorphism Predictor) which takes the results of the structural analysis and uses a Random Forest machine-learning method to

predict whether mutations are pathogenic (Al-Numair and Martin, 2013). SAAPpred is restricted to analyzing mutations in proteins for which a native structure is available, but appears to outperform methods such as SIFT (Ng and Henikoff, 2003), PolyPhen-2 (Adzhubei *et al.*, 2010, 2013) and FATHMM (Shihab *et al.*, 2013).

SAAPdap uses a combination of rule-based structural measures to assess whether a mutation is likely to alter the local structural environment. SAAPpred exploits this information to predict whether the function of a protein will be affected and, in turn, lead to disease. The approach has been used to study structural differences between disease-causing mutations and neutral polymorphisms (Hurst *et al.*, 2009; Al-Numair and Martin, 2013), and to analyse mutations in glucose-6-phosphate dehydrogenase (Kwok *et al.*, 2002) and in the tumour suppressor P53 (Martin *et al.*, 2002).

In this paper we investigate whether the approach that we developed for SAAPdap and SAAPpred can be used for differential phenotype prediction specifically for mutations in the beta-myosin heavy chain (Myosin-7, UniProtKB/SwissProt accession P12883, <http://www.uniprot.org/uniprot/P12883>), encoded by the *MYH7* gene (OMIM \*160760), and leading to hypertrophic cardiomyopathy (HCM, OMIM #192600) or dilated cardiomyopathy (DCM, OMIM #613426).

Myosin-7 is part of the force-generating molecular motor of the sarcomere and parts of the structure have been solved. It is divided into three main domains: a globular ‘head’, which includes the ATP-binding site and the actin-binding site; the ‘neck’ which is composed of an  $\alpha$ -helical domain to which the myosin light chains bind and which is further subdivided into a converter region and a lever arm involved in the amplification of mechanical energy; and the ‘tail’ or ‘rod’ region. Together with *MYBPC3* (the gene encoding myosin binding protein C), mutations in *MYH7* are the major cause of HCM as well as being a cause of DCM and left ventricular non-compaction (LVNC) (Haas *et al.*, 2014). In contrast

to *MYBPC3*, where most pathogenic variants cause mRNA and protein truncation, the large majority of *MYH7* variants are missense (Carrier *et al.*, 1997; Richard *et al.*, 2003) which often makes prediction of pathogenicity problematic (Walsh *et al.*, 2010; Kumar *et al.*, 2013).

## 2 MATERIALS AND METHODS

### 2.1 Dataset of variants

A dataset of *MYH7* variants was built from a) disease-causing or likely-pathogenic variants for which phenotypic data are available in the Human Genome Mutation Database (HGMD) (Stenson *et al.*, 2002); b) variants found in a curated dataset extracted from the literature and used for commercial gene testing reports (*Health in Code SL*); and c) variants detected in a cohort of consecutively evaluated unrelated HCM/DCM patients at UCLH. Genetic analysis was approved by the UCLH review board (IRB) and informed written consent was obtained from all subjects (Lopes *et al.*, 2013). Although there are no co-segregation data or functional studies that can ‘prove’ the causality of mutations, selected variants from all three datasets were rare as defined by a minor allele frequency (MAF) < 0.5% in the ESP6500 NIH Heart, Lung and Blood Institute (NHLBI) exome sequencing project database (Pan *et al.*, 2012; Andreasen *et al.*, 2013). This dataset is larger and more comprehensive than the data available from other sources and contains approximately twice the number of Myosin-7 mutations available in Swissvar/Humsavar. The complete dataset has been provided as Supplementary File 1. Proprietary data from HGMD, where the mutations are not available in other datasets, have been indicated solely by their HGMD accession code.

### 2.2 SAAPdap structural analysis and SAAPpred

Our previous software, SAAPdap (Al-Numair and Martin, 2013) performs a set of 14 structural analyses (using software written in Perl and C), plus the calculation of solvent accessibility (Lee and Richards, 1971). SAAPdap provides cutoffs for each of the analyses to suggest whether these are likely to be damaging (Al-Numair and Martin, 2013; Hurst *et al.*, 2009). To predict pathogenicity, a total of 47 features are derived from these analyses (Table 1) and are used as input to SAAPpred, a machine learning method that uses Random Forests to predict whether a mutation is pathogenic (Al-Numair and Martin, 2013). In this paper, the same methodology is used but, rather than using a dataset of pathogenic and phenotypically silent mutations, a dataset of HCM and DCM mutations in Myosin-7 is used.

### 2.3 Features and machine learning for differential phenotype prediction

In addition to the 47 features used in SAAPpred, three other features were derived that represent distances from structural cluster centres. These were identified by clustering the coordinates of HCM and DCM mutations using single linkage clustering and finding the number of clusters that gave the most significant separation of HCM and DCM mutations between the clusters ( $\chi^2$  test). See Section 3.3.

Subsets (listed in the legend to Table 7) of the 50 features were then used to train Random Forest predictors implemented in WEKA

version 3.6.7 (Witten *et al.*, 2011) using the default classification threshold to separate mutations associated with HCM and DCM.

## 3 RESULTS

### 3.1 *MYH7* mutation data analysis and prediction of pathogenicity

*MYH7* mutations associated with various cardiomyopathy phenotypes are shown in Table 2. Note that it is not possible to know whether variants are truly pathogenic; rather we treat mutations *associated* with an HCM or DCM cardiomyopathy phenotype in the above-mentioned databases, or in the literature, as actual positives. A total of 403 mutations were identified in the *MYH7* gene. More than two-thirds of them have previously been published in the literature as being associated with disease and the others are novel variants.

Since we map mutations to protein structure and therefore require a structure to be solved of the protein of interest, we are not able to analyse all mutations. Of the 396 unique mutations (i.e. distinct mutations, different from one another at the protein level) in *MYH7*, 166 (41.9%) did not map to structure and therefore could not be analysed (see Table 2). This situation should improve as further structures become available. 385 of the 396 unique mutations had a recorded phenotype and of these 230 mapped to at least one Protein DataBank (PDB) chain. Table 3 lists three PDB structures which were identified for human Myosin-7. Two other PDB files (IDs 1ik2 and 3dtp) were eliminated since one was a 3D homology model and the other was a human-chicken fusion protein. Most mutations were associated with HCM ( $n = 298$ ), whereas all other phenotypes were associated with fewer than 50 mutations each, including DCM with the next highest number of mutations ( $n = 46$ ). The majority of mutations in both HCM and DCM were unique (292 and 46 respectively, see Table 2). Since mutations related to these phenotypes were the most abundant, further analyses were restricted to HCM and DCM, grouping the remaining phenotypes as ‘other’.

The distribution of the variants amongst the structural and functionally-annotated domains of the beta-myosin heavy chain protein was analysed. Figure 1 shows the regions for which structures are known and the distribution of observed mutations together with the domains of the Myosin-7 sequence as annotated by UniProtKB/SwissProt (UniProt Consortium, 2014) ([http://www.uniprot.org/uniprot/P12883#section\\_features](http://www.uniprot.org/uniprot/P12883#section_features)), Pfam (Finn *et al.*, 2014) (<http://pfam.xfam.org/protein/P12883>), SMART (Letunic *et al.*, 2012) ([http://smart.embl.de/smart/show\\_motifs.pl?ID=P12883](http://smart.embl.de/smart/show_motifs.pl?ID=P12883)), and InterPro (Hunter *et al.*, 2012) (<http://www.ebi.ac.uk/interpro/protein/P12883>). All of the 235 unique variants that mapped to structure were located in the myosin globular ‘head’ domain or the ‘neck’ region with no mutations seen in the ‘tail’ or ‘IQ motif’ regions. 99.1% of mutations were in annotated domains or regions, while just two mutations (0.9%, at positions 82 and 838) were in un-annotated parts of the sequence.

The individual structural effects for the 230 unique mutations which mapped to structure and for which a phenotype was also recorded (see Table 2) were analyzed using SAAPdap. 175 variants (76.1%) were classified as likely to be damaging by one or more

Analysis	Features	Type
Binding	Is the residue involved in binding (defined by presence of specific contacts with another protein chain or ligand)?	Boolean
Interface	Is the residue in an interface (defined by change in solvent accessibility between complexed and uncomplexed forms)?	Boolean
SProtFT	Is the residue annotated with a functionally relevant SwissProt feature? Which of 12 SwissProt features appear? (ACT_SITE, BINDING, CA_BIND, DNA_BIND, NP_BIND, METAL, MOD_RES, CARBOHYD, MOTIF, LIPID, DISULFID, CROSSLNK)?	Boolean 12 x Boolean
RelAccess	Relative solvent accessibility of the residue	Percentage
ImPACT	ImPACT conservation score for the residue if it is found to be significantly conserved	Real
HBond	If the native residue was involved in a hydrogen bond, the difference in hydrogen bonding pseudo-energy	Real
SurfacePhobic	The difference in hydrophobicity if the residue is on the surface and the hydrophobicity has increased	Real
CorePhobic	The difference in hydrophobicity if the residue is buried and the hydrophobicity has decreased	Real
BuriedCharge	The difference in charge if the residue is buried	Integer
SSGeom	Was the native residue involved in a disulphide bond?	Boolean
Void	The difference in size of the largest void The sizes of the 10 largest voids in the native protein The sizes of the 10 largest voids in the mutant protein	Real 10 x Real 10 x Real
Clash	The sum of the van der Waals and torsional energy for the minimum perturbation protocol modelled sidechain replacement	Real
Glycine	If the native residue was a glycine, the Ramachandran pseudo-energy difference of the mutation	Real
Proline	If the mutant residue was a proline, the Ramachandran pseudo-energy difference of the mutation	Real
CisPro	Was the native residue a cis-proline?	Boolean

**Table 1.** The 47 features used in SAAPred machine learning derived from the 14 structural analyses in SAAPdap.

Disease (Phenotype)	Total* mutations	Unique† mutations	Mutations mapped to PDB
HCM	298	292	188
DCM	46	46	21
RCM	1	1	1
LVNC	17	17	1
LVNC/ASD	1	1	1
DCM/Endocardial Fibroelastosis	1	1	1
DCM/LVNC	3	3	2
HCM/LVNC	1	1	1
HCM/DCM/LVNC	2	2	2
HCM/DCM	3	3	3
HCM/RCM/DCM	2	2	2
Laing distal myopathy	4	4	2
Ebstein	5	5	1
Cardiomyopathy and distal myopathy	3	3	2
Myosin storage myopathy	3	3	1
Hyaline body myopathy	1	1	1
No recorded phenotype	11	11	5
Total	403	396	235

**Table 2.** Numbers of *MYH7* mutations for each phenotype. Abbreviations: PDB, Protein DataBank; DCM, Dilated Cardiomyopathy; HCM, Hypertrophic Cardiomyopathy; RCM, Restrictive Cardiomyopathy; LVNC, Left Ventricular Non-compaction; ASD, Atrial Septal Defect. The mutations for which there was no recorded phenotype were excluded from structural analysis, meaning that only 230 mutations which mapped to PDB structures could be analysed. For the novel differential phenotype predictor, only the 209 unique HCM and DCM mutations that mapped to PDB structures were used. \*Total mutations represents the total count of amino acid mutations. Sometimes the same mutation may be observed multiple times because the DNA level mutation is different or because of redundancy between different data sources. †Unique mutations represents the number of non-redundant mutations at the protein level.

PDB ID	Resolution	Description	Residues
2fxm	2.7Å	Structure of the human beta-myosin S2 fragment	A: 838–961 B: 850–961
2fxo	2.5Å	Structure of the human beta-myosin S2 fragment	A: 838–963 B: 838–961 C: 838–962 D: 838–963
4db1	2.6Å	Cardiac human myosin S1DC, beta isoform complexed with Mn-AMPPNP	A: 2-777 B: 2-775

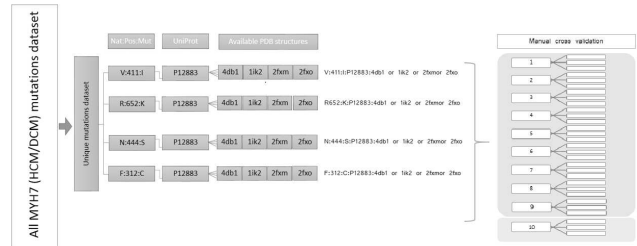
**Table 3.** PDB structures for UniProtKB/SwissProt accession code P12883. PDB files may be accessed at <http://www.pdb.org/> or viewed using PDBSum (<http://www.ebi.ac.uk/thornton-srv/databases/pdbsum/>). Note that PDB file 2fxo contains a mutation Glu924Lys.

SAAPdap Structural Analysis	Number of mutations
No PDB structure available	166
No individual significant structural effect	55
At least one significant structural effect	175
• HBond	42
• BuriedCharge	31
• SProtFT	2
• Interface	48
• Clash	14
• Proline	2
• ImPACT	138
• Binding	20
• Void	0
• SurfacePhobic	15
• Glycine	8
• CisPro	1
• CorePhilic	26
• SSGeom	0

**Table 4.** SAAPdap Structural Analysis for the 230 unique Myosin-7 mutations with a recorded phenotype which mapped to structure (see Table 2).

individual SAAPdap analyses while, for 55 variants, no significant individual structural effect was detected (see Table 4). The most frequent features affected were: mutation of a highly conserved residue (ImPACT) occurring in 138 variants; mutation of an interface amino acid occurring in 48 variants; and disruption of hydrogen-bonds occurring in 42 variants. Other significant mutation effects occurred less frequently, with no observed mutations causing voids.

The output from SAAPdap for the 230 unique mutations that mapped to structure was then fed into SAAPpred (Al-Numair and Martin, 2013) and 92.7% were predicted as pathogenic (i.e.  $S_n=0.927$ ). This compares with 69.51% predicted to be pathogenic using SIFT and 90% predicted to be pathogenic using PolyPhen-2. Other metrics such as specificity (Sp), accuracy (Acc), the F1-score and the Matthews’ Correlation Coefficient (MCC) could not be calculated since no set of validated non-pathogenic single amino acid mutations is available — even in the ESP 5K and 1000 Genomes data there are very few missense variants in MYH7 with a frequency  $> 5\%$  that could comfortably be classified as benign.



**Fig. 2.** MYH7 (HCM/DCM) dataset selection for machine learning. A unique mutation level filtering is used, where the same mutation (UniProtKB/SwissProt:Native:Number:Mutant) does not occur in training and testing sets. This was achieved using a ‘manual’ (non-WEKA) cross-validation that splits the dataset into  $N$  sets, each one in turn was chosen as the testing set and the remaining  $N - 1$  were used for training.

### 3.2 A machine-learning approach for MYH7 differential phenotype prediction

All mutations associated with multiple phenotypes, or causing phenotypes other than HCM or DCM were discarded leaving the 188 unique HCM and 21 unique DCM mutations which map to structure.

Using the results of the SAAPdap structural analysis described above, of the 47 ‘features’ used to describe the mutations, 14 were found to be redundant (i.e. they had the same value for all examples in the dataset: the 13 UniProtKB/SwissProt features and the disulphide (SSGeom) analysis), thus reducing the number of features to 33. Although a single structure was used with SAAPpred, because of the limited size of the available dataset for differential phenotype prediction, it was desirable to exploit multiple structures to enrich the dataset. PDB files 4db1 and 2fxm contain two copies of the protein while 2fxo contains four copies (Table 3). These data were then used to train Random Forest models in WEKA. The use of multiple structures for each mutation meant that cross-validation could not be performed within WEKA since it is possible that WEKA could select the same mutation (in a different structure) to be in both training and testing sets.

To address the cross-validation problem and to deal with the severe imbalance of the dataset (there being many more HCM mutations than DCM), Perl code was written to limit the size of each class by selecting examples at random and to divide the 188 HCM and 21 DCM unique mutations with available PDB structures into sets of approximately the same size. For example, if the data were split into 21 sets, each of these 21 sets in turn (each containing one DCM mutation) was chosen as a test set and the remaining 20 sets



Number of folds / models	$T$	$m_{try}$	Acc	MCC
10	1000	10	0.6229	0.2463
10	1000	15	0.6750	0.3590
<b>10</b>	<b>1000</b>	<b>20</b>	<b>0.7000</b>	<b>0.4103</b>
10	1000	25	0.6916	0.3851
10	50	20	0.6833	0.3681
10	100	20	0.6916	0.3872
10	500	20	0.6937	0.4023
<b>10</b>	<b>1000</b>	<b>20</b>	<b>0.7000</b>	<b>0.4103</b>
10	2000	20	0.6812	0.3686
10	5000	20	0.7000	0.4005

**Table 5.** Exploring the number of features and number of trees in HCM vs. DCM prediction.  $T$  is the number of trees;  $m_{try}$  is the number of randomly chosen attributes in every split. Initially  $m_{try}$  was explored using  $T = 1000$  and an optimum value of 20 was identified (shown in bold).  $T$  was then explored retaining the optimum value of 1000. Performance measures: accuracy (Acc) and Matthews’ correlation coefficient (MCC). All scores are averaged over 10-folds of ‘manual’ (non-WEKA) cross-validation.

Number of clusters	Significance	Percentage of Expecteds < 5
2	$p < 0.4384$	0
<b>3</b>	$p < \mathbf{0.0003755}$	16.7%
4	$p < 0.001256$	37.5%
5	$p < 0.002577$	50%
6	$p < 0.005057$	50%
7	$p < 0.01013$	50%
8	$p < 0.01778$	56.25%
9	$p < 0.03044$	55.56%
10	$p < 0.03116$	60%

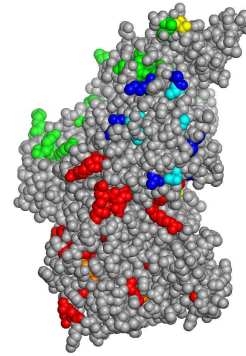
**Table 6.** Significance calculated from  $\chi^2$  tests on the ability of 3D clustering to separate HCM from DCM mutations. The highest significance result is shown in bold. For the  $p$ -value to be reliable, there must be no more than 20% of expected counts less than five. Consequently the  $p$ -values for  $\geq 4$  clusters will be over-estimated.

(each containing the remaining 20 DCMs) were used for training. In each case, the data sets were enlarged with all the available PDB chain structures and balanced datasets were generated by retaining all the DCM mutations and randomly drawing the same number of mutations from the HCM dataset (see Figure 2). The random draws from the HCM dataset were taken 10 times over to provide a representative sample of the HCM class and the results from the trained predictors were averaged.

The parameter space described by the number of features used in each tree decision point ( $m_{try}$ ) and the number of trees ( $T$ ) was explored and, as shown in Table 5, the best results were obtained using 1000 trees with 20 features (accuracy of 70% and MCC=0.41).

### 3.3 Structural clustering of mutations

Anecdotal evidence suggested that HCM- and DCM-associated mutations tend to be distributed differently across the Myosin-7



**Fig. 3.** Clustering Myosin-7 mutations in the N-terminal region using PDB file 4db1. For the three clusters, HCM mutations are shown in 1: red, 2: green and 3: blue, while DCM mutations are shown in 1: orange, 2: yellow and 3: cyan. DCM mutations are over-represented in cluster 3 (cyan); when they appear in clusters 1 and 2, (orange and yellow) they are mostly buried.

structure. This observation was exploited in an attempt to improve the results.

PDB file 2fxm, which represents the C-terminal region, contains only two DCM mutations compared with 35 HCM, indicating that DCM mutations are very rare in this domain. For the N-terminal domain (PDB file 4db1), the  $C\alpha$  positions of the mutated residues were clustered using single linkage hierarchical clustering. For each of 2 ... 10 clusters, a  $\chi^2$  test was performed to see how well the clustering separated HCM from DCM mutations, as shown in Table 6. Apart from two clusters, these are all clearly significant at the  $p < 0.05$  level. However, as the number of clusters gets larger, one needs to take care with the significance levels, because no more than 20% of expected values should be  $< 5$  and none  $< 1$  (significance will be over-estimated if either of these is true). For  $\geq 4$  clusters, the first of these fails and for  $\geq 6$  clusters the second also fails. However, between three and six clusters the significance is so good, that (while it will be over-estimated for 4–6 clusters) it is clearly still better than  $p < 0.05$  with 3 clusters giving the most significant result and passing both of the validity criteria even if a Bonferroni correction is made for multiple testing. Consequently we clearly have clusters of residues in the N-terminal region that are over/under populated with DCM and HCM mutations compared with what is expected.

Figure 3 illustrates the three clusters in the N-terminal domain contained in PDB file 4db1. Cluster members are listed in Supplementary File 2. In particular, DCM is highly over-represented in the third (blue/cyan) cluster. DCM mutations in clusters 1 and 2 (orange and yellow) are hardly visible and therefore mostly buried. On the other hand, the DCM mutations in cluster 3 (cyan) are largely on the surface.

As a control, to ensure that the significance of the clustering was not a random effect, we also permuted the labels randomly for the three clusters 1000 times over and calculated the average random  $p$ -value ( $p = 0.5133$ ,  $\sigma_{n-1} = 0.2859$ ) from a  $\chi^2$  test. This is clearly not significant and compares with the true labels which gave a  $p < 0.0003755$ . This  $p$ -value is 1.794 standard deviations away from the mean on the distribution of random  $p$ -values which is significant at the  $p < 0.05$  level.

To use this information in machine learning, the centroid of each cluster was calculated and the feature vector for each mutation was expanded by the addition of the distances from the C-alpha of the mutated residue to each of the three centroids. Mutations that were in the C-terminal domain (and mapped to PDB files 2fxm and 2fxo rather than 4db1) were given distances of 100.0Å, 100.0Å, 100.0Å from the three clusters.

### 3.4 Optimizing the machine learning

As described above, initial training to explore the number of trees and features was performed using 10 machine-learning models (equivalent to cross-validation folds, each with a random selection of the HCM data) with the prediction results averaged across the 10. Using a larger number of models allows more of the HCM data to be exploited in each model while maintaining balanced datasets. Using 21 machine-learning models, only one unique DCM mutation can be held back from training for test purposes.

After determining the optimum number of features and trees, the most informative features were explored together with different numbers of machine-learning models (5, 11 and 21 models). Odd numbers were used to allow a jury vote in predictions. Addition of the ‘clustering’ feature described above was also explored. The different feature sets are described in detail together with summary results in Table 7. In brief, the feature subsets were as follows: ‘All’ the full standard set of 33 informative features (47 from SAAPdap, but with the 14 redundant features, which were identical for all mutations, removed); ‘Top 5 voids’ uses only the top five largest voids (before and after mutation) instead of the standard 10; ‘Delta voids’ uses differences between void sizes in native and mutant structures rather than absolute values; ‘Set1’ was a selection of the five features found to be most discriminatory using  $\chi^2$  tests on each of the features; ‘Set2’ and ‘Set3’ were sets of features randomly generated within WEKA, ‘Set2’ being based on the ‘All’ dataset and ‘Set3’ being based on the ‘Delta voids’ set.

Initially, the number of machine-learning models was tested using the full feature set (‘All’), plus those feature sets that reduced the amount of void data (‘Top 5 voids’ and ‘Delta voids’), with and without the clustering features. Having established that 11 models was the most effective, the reduced feature sets were explored using a smaller value of  $m_{try}$  owing to the much reduced number of features.

As shown in Table 7, the best performance was obtained using 11 machine-learning models with ‘Set2’ plus the clustering features. Cross-validation with 11 models used 19 of the 21 DCMs in each training set with 2 held back for testing. This gave an accuracy of 75% and MCC=0.531. By removing two machine-learning models that performed particularly badly and did not predict any DCM mutations (whether correct or incorrect), this increased to an accuracy of 79% and MCC=0.61. It appears that these particularly bad machine-learning models have failed to learn the characteristics of DCM mutations. To apply the method to novel mutations, we would remove these two bad machine-learning models and use the remaining nine to make predictions.

### 3.5 Control Experiments

To ensure that the performance of the predictor does not come only from the structural clustering, we also tested the performance of using the structural clusters alone. Using the 2–10 structural clusters

described above, each cluster was assigned as a DCM or HCM cluster based on that phenotype having a higher observed/expected ratio in that cluster. An additional cluster was created to represent the mutations that map to the C-terminal domain (PDB code 2fxm) which has a very small number of DCM mutations. Each mutation was then predicted as DCM or HCM based on its cluster membership. Performance was then calculated for each level of clustering with best performance being achieved with three clusters (plus the C-terminal domain cluster): MCC=0.33, ACC=0.89,  $Sn_{HCM}$ =0.95,  $Sn_{DCM}$ =0.33. For a real prediction problem, cluster membership would need to be assigned based on the distance to the closest cluster centre (average linkage) or closest cluster member (single linkage). Clearly this performance is considerably worse than our full predictor as judged by MCC (full predictor MCC=0.53, or MCC=0.61 with the worst machine-learning models removed).

This is also a good example to illustrate the well-known problem in machine learning that accuracy is a poor indicator of performance with unbalanced datasets (the cluster-only prediction gives ACC=0.89 while the full predictor gives ACC=0.75, or ACC=0.79 with the worst models machine-learning removed). However, simply predicting everything as HCM would give ACC=0.90 and, by definition,  $Sn_{HCM}$ =1.00 and  $Sn_{DCM}$ =0.00, while the MCC would be a much better indicator of overall performance giving a value of MCC=0.12 (adding 1 to TP,FP,TN,FN — treating HCM as positive and DCM as negative — since TN=FN=0 results in a divide-by-zero error).

As a control on the overall prediction, the testing was repeated using two of the test sets, but the labels were randomly shuffled five times over. As expected, the prediction performance was essentially random with an MCC=−0.123 for the first test set and MCC=−0.115 for the second test set.

## 4 DISCUSSION

It is logical to assume that the functional consequences of mutations in the same gene depend on the specific domain or region where the variant is localized (Woo *et al.*, 2003), but the hypothesis that the structural impact of a missense variant influences differential pathogenic phenotype or outcome has not previously been tested.

In practice, a novel mutation would be tested for predicted pathogenicity before an HCM/DCM prediction was performed. We confirmed that the SAAPpred approach performs well in identifying pathogenic mutations in *MYH7* and went on to test a machine-learning method that discriminated between pathogenic variants associated with an HCM or DCM phenotype (accuracy of 75% and MCC=0.531). This was achieved by averaging 11 machine-learning models using feature Set2 (Binding, RelAccess, SurfacePhobic, CorePhilic, Voids, MutantLargestVoid1, NativeLargestVoid1, Clash, Proline, CisPro and Clustering) and using 1000 trees with 5 features. These results are surprisingly good considering the limited size of the dataset used in training. Indeed the results are as good as the overall performance of some methods used for general pathogenicity prediction — for example, our assessment (Al-Numair and Martin, 2013) of MutationAssessor showed an overall accuracy of 69.8% and MCC=0.453, while SIFT showed an overall accuracy of 76.3% and MCC=0.528. Clearly these results are comparable with what we are able to achieve for HCM/DCM differential phenotype prediction which is a more difficult problem

Number of folds / models	Features used	$T$	$m_{try}$	$Sn_{HCM}$	$Sn_{DCM}$	F1	Acc	MCC
5	All	1000	20	0.572	0.611	0.576	0.576	0.152
5	All + Clustering	1000	20	0.755	0.481	0.679	0.648	0.311
5	Top 5 voids + Clustering	1000	20	0.735	0.611	0.688	0.681	0.368
5	10 delta void + Clustering	1000	20	0.785	0.407	0.676	0.608	0.205
11	All	1000	20	0.705	0.648	0.673	0.682	0.429
11	All + Clustering	1000	20	0.739	0.463	0.662	0.608	0.220
11	Top 5 voids + Clustering	1000	20	0.830	0.481	0.741	0.699	0.427
11	10 delta voids + Clustering	1000	20	0.830	0.519	0.730	0.676	0.521
21	All	1000	20	0.619	0.648	0.585	0.631	0.357
21	All + Clustering	1000	20	0.746	0.463	0.684	0.623	0.293
21	Top 5 voids + Clustering	1000	20	0.690	0.463	0.610	0.627	0.374
21	10 delta voids + Clustering	1000	20	0.619	0.426	0.584	0.560	0.133
11	Set1 + Clustering	1000	5	0.659	0.593	0.603	0.625	0.314
<b>11</b>	<b>Set2 + Clustering</b>	<b>1000</b>	<b>5</b>	<b>0.795</b>	<b>0.574</b>	<b>0.737</b>	<b>0.750</b>	<b>0.531</b>
11	Set3 + Clustering	1000	5	0.852	0.519	0.746	0.699	0.520

**Table 7.** Summary results of machine learning performance using different features of HCM/DCM dataset and using different numbers of folds of cross-validation. The best performing predictor is shown in bold. ( $T$ : the number of trees;  $m_{try}$ : the number of randomly chosen attributes in every split;  $Sn_{HCM}$ : Sensitivity for HCM mutations;  $Sn_{DCM}$ : Sensitivity for DCM mutations; F1: The F1-score; Acc: Accuracy; MCC: Matthews’ Correlation Coefficient)

- ‘All’: BuriedCharge, Binding, CorePhilic, CisPro, Clash, Glycine, HBond, ImPACT, Interface, MutantLargestVoid1...MutantLargestVoid10, NativeLargestVoid1...NativeLargestVoid10, Proline, RelAccess, SurfacePhobic, Void.
- ‘Top 5 voids’: BuriedCharge, Binding, CorePhilic, CisPro, Clash, Glycine, HBond, ImPACT, Interface, MutantLargestVoid1...MutantLargestVoid5, NativeLargestVoid1...NativeLargestVoid5, Proline, RelAccess, SurfacePhobic, Void.
- ‘Delta Voids’: BuriedCharge, Binding, CorePhilic, CisPro, Clash, Glycine, HBond, ImPACT, Interface, DeltaLargestVoid1...DeltaLargestVoid10, Proline, RelAccess, SurfacePhobic, Void.
- ‘Set1’: Uses the most informative features based on  $\chi^2$  tests: Binding, RelAccess, ImPACT and Glycine.
- ‘Set2’: A WEKA randomly selected dataset: Binding, RelAccess, SurfacePhobic, CorePhilic, TotalVoidVolume, MutantLargestVoid, NativeLargestVoid, Clash, Proline, CisPro.
- ‘Set3’: A WEKA randomly selected dataset based on the ‘Delta Voids’ set: Binding, Interface, RelAccess, ImPACT, Hbond, BuriedCharge, DeltaVoidTotal, DeltaVoidLargest1...DeltaVoidLargest5, Clash, Glycine.

owing to the small imbalanced dataset. By removing two machine-learning models that performed particularly badly, the performance was increased to an accuracy of 79% and MCC=0.61.

Because the SAAPdap structural analysis relies on having a crystal structure of the protein in question, our predictions are limited to mutations in regions of the protein for which a structure has been solved. Consequently, we are only able to look at 188 of 292 unique mutations leading to HCM and 21 of 46 mutations leading to DCM. If structures become available for more of the protein, then this situation will improve. However, for mutations that are present in disordered regions of structure, different methods of prediction will be required. It is also possible that the performance of the method may be further improved by taking into account missing parts of the structure. However, since all the structural parameters included in the prediction are the results of local interactions, this is unlikely to have a significant effect.

Our analysis of the structural distribution of HCM- and DCM-associated mutations showed that there was a highly statistically significant difference in the locations of these mutations. Referring to Figure 3, DCM is highly over-represented in the blue/cyan cluster and largely on the surface, while DCM mutations present in the

remaining clusters are mostly buried. The functional consequences of this distribution warrant further *in vitro* studies.

#### 4.1 Conclusions and future directions

Missense single nucleotide variants in *MYH7* lead to a dominant negative effect in which the mutated protein is not degraded but rather integrates into the sarcomere, leading to the disease phenotype. The various effects of individual variants on fibre contractile velocity, force and calcium sensitivity have been proposed as an explanation for the existence of dramatically different phenotypes arising from genetic variation in the same molecule. A paradigm has been proposed whereby mutations that increase motor activity and power output lead to HCM, while those that diminish motor function and decrease power output lead to DCM (Spudich, 2014).

This work confirms the hypothesis that structural data can be used with machine learning to create a differential phenotype predictor, in this case able to distinguish between HCM and DCM mutations in *MYH7*. The performance exceeds that of the well-known SIFT program in the problem of predicting pathogenic vs. neutral mutations. Differential phenotype prediction has all the



challenges of pathogenicity prediction with the added complications of having a small unbalanced dataset. This work provides the basis for differential phenotype prediction and with further work could be used to guide clinical genetic testing strategies and further clinical investigations.

## 5 ACKNOWLEDGMENTS

NSAN was funded by the Saudi Arabian Ministry of Higher Education (MOHE) and the King Faisal Specialist Hospital & Research Centre (KFSH&RC). LRL was supported by a grant from the Gulbenkian Doctoral Programme for Advanced Medical Education, sponsored by Fundação Calouste Gulbenkian, Fundação Champalimaud, Ministério da Saúde and Fundação para a Ciência e Tecnologia, Portugal. This work was undertaken at UCLH/UCL who received a proportion of funding from the UK Department of Health's National Institute for Health Research Biomedical Research Centres funding scheme. LM received funding from grant FIS 2011: P111/02604, Instituto de Salud Carlos III, Madrid, Spain. LM is a shareholder in *Health in Code SL*. The remaining authors have no interests/relationships to declare.

## REFERENCES

- Adzhubei, I. A., Schmidt, S., Peshkin, L., Ramensky, V. E., Gerasimova, A., Bork, P., Kondrashov, A. S., and Sunyaev, S. R. (2010). A method and server for predicting damaging missense mutations. *Nature Methods*, **7**(4), 248–249.
- Adzhubei, I. A., Jordan, D. M., and Sunyaev, S. R. (2013). Predicting functional effect of human missense mutations using PolyPhen-2. *Curr Protoc Hum Genet*, **76**, 7.20.
- Al-Numair, N. S. and Martin, A. C. R. (2013). The SAAP pipeline and database: tools to analyze the impact and predict the pathogenicity of mutations. *BMC Genomics*, **14**(3), 1–11.
- Andreasen, C., Nielsen, J. B., Refsgaard, L., Holst, A. G., Christensen, A. H., Andreasen, L., Sajadieh, A., Haunsø, S., Svendsen, J. H., and Olesen, M. S. (2013). New population-based exome data are questioning the pathogenicity of previously cardiomyopathy-associated genetic variants. *Eur J Hum Genet*, **21**, 918–928.
- Arad, M., Seidman, J., and Seidman, C. E. (2002). Phenotypic diversity in hypertrophic cardiomyopathy. *Hum Molec Genet*, **11**(20), 2499–2506.
- Bao, L., Zhou, M., and Cui, Y. (2005). nsSNPAnalyzer: identifying disease-associated nonsynonymous single nucleotide polymorphisms. *Nucleic Acids Res*, **33**, W480–W482.
- Bromberg, Y. and Rost, B. (2007). SNAP: predict effect of non-synonymous polymorphisms on function. *Nucleic Acids Res*, **35**, 3823–3835.
- Bromberg, Y., Yachdav, G., and Rost, B. (2008). SNAP predicts effect of mutations on protein function. *Bioinformatics*, **24**, 2397–2398.
- Calabrese, R., Capriotti, E., Fariselli, P., Martelli, P. L., and Casadio, R. (2009). Functional annotations improve the predictive score of human disease-related mutations in proteins. *Human Mutation*, **30**, 1237–1244.
- Carrier, L., Bonne, G., Bahrend, E., Yu, B., Richard, P., Niel, F., Hainque, B., Cruaud, C., Gary, F., Labeit, S., Bouhour, J.-B., Dubourg, O., Desnos, M., Hagege, A. A., Trent, R. J., Komajda, M., Fiszman, M., and Schwartz, K. (1997). Organization and sequence of human cardiac myosin binding protein C gene (MYBPC3) and identification of mutations predicted to produce truncated proteins in familial hypertrophic cardiomyopathy. *Circulation Res*, **80**(3), 427–434.
- Dantzer, J., Moad, C., Heiland, R., and Mooney, S. (2005). MutDB services: Interactive structural analysis of mutation data. *Nucleic Acids Res*, **33**, W311–W314.
- Finn, R. D., Bateman, A., Clements, J., Coggill, P., Eberhardt, R. Y., Eddy, S. R., Heger, A., Hetherington, K., Holm, L., Mistry, J., Sonnhammer, E. L. L., Tate, J., and Punta, M. (2014). Pfam: the protein families database. *Nucleic Acids Res*, **42**, D222–D230.
- González-Pérez, A. and López-Bigas, N. (2011). Improving the assessment of the outcome of nonsynonymous SNVs with a consensus deleteriousness score. *Condel. Am J Hum Genet*, **88**, 440–449.
- Haas, J., Frese, K. S., Peil, B., Kloos, W., Keller, A., Nietsch, R., Feng, Z., Müller, S., Kayvanpour, E., Vogel, B., Sedaghat-Hamedani, F., Lim, W.-K., Zhao, X., Fradkin, D., Köhler, D., Fischer, S., Franke, J., Marquart, S., Barb, I., Li, D. T., Amr, A., Ehlermann, P., Mereles, D., Weis, T., Hassel, S., Kremer, A., King, V., Wirsz, E., Isnard, R., Komajda, M., Serio, A., Grasso, M., Syrris, P., Wicks, E., Plagnol, V., Lopes, L., Gadgaard, T., Eiskjær, H., Jørgensen, M., Garcia-Gustiniani, D., Ortiz-Genga, M., Crespo-Leiro, M. G., Deprez, R. H. L. D., Christiaans, I., van Rijsingen, I. A., Wilde, A. A., Waldenstrom, A., Bolognesi, M., Bellazzi, R., Möner, S., Bermejo, J. L., Monserrat, L., Villard, E., Mogensen, J., Pinto, Y. M., Charron, P., Elliott, P., Arbustini, E., Katus, H. A., and Meder, B. (2014). Atlas of the clinical genetics of human dilated cardiomyopathy. *Eur Heart J*, **36**, 1123–1135.
- Hughes, S. E. and McKenna, W. J. (2005). New insights into the pathology of inherited cardiomyopathy. *Heart*, **91**, 257–264.
- Hunter, S., Jones, P., Mitchell, A., Apweiler, R., Attwood, T. K., Bateman, A., Bernard, T., Binns, D., Bork, P., Burge, S., de Castro, E., Coggill, P., Corbett, M., Das, U., Daugherty, L., Duquenne, L., Finn, R. D., Fraser, M., Gough, J., Haft, D., Hulo, N., Kahn, D., Kelly, E., Letunic, I., Lonsdale, D., Lopez, R., Madera, M., Maslen, J., McAnulla, C., McDowall, J., McMenamin, C., Mi, H., Mutowo-Muullenet, P., Mulder, N., Natale, D., Orengo, C., Pesce, S., Punta, M., Quinn, A. F., Rivoire, C., Sangrador-Vegas, A., Selengut, J. D., Sigrist, C. J. A., Scheremetjew, M., Tate, J., Thimmajananthan, M., Thomas, P. D., Wu, C. H., Yeats, C., and Yong, S.-Y. (2012). InterPro in 2011: new developments in the family and domain prediction database. *Nucleic Acids Res*, **40**, D306–D312.
- Hurst, J. M., McMillan, L. E. M., Porter, C. T., Allen, J., Fakorede, A., and Martin, A. C. R. (2009). The SAAPdb web resource: a large-scale structural analysis of mutant proteins. *Hum Mutat*, **30**, 616–624.
- Karchin, R., Diekhans, M., Kelly, L., Thomas, D. J., Pieper, U., Eswar, N., Haussler, D., and Sali, A. (2005). LS-SNP: large-scale annotation of coding non-synonymous SNPs based on multiple information sources. *Bioinformatics*, **21**, 2814–2820.
- Kircher, M., Witten, D. M., Jain, P., O’Roak, B. J., Cooper, G. M., and Shendure, J. (2014). A general framework for estimating the relative pathogenicity of human genetic variants. *Nature Genetics*, **46**, 310–315.
- Kumar, A., Rajendran, V., Sethumadhavan, R., and Purohit, R. (2013). Roadmap to determine the point mutations involved in cardiomyopathy disorder: a Bayesian approach. *Gene*, **519**, 34–40.
- Kwok, C. J., Martin, A. C. R., Au, S. W. N., and Lam, V. M. S. (2002). G6PDdb, an integrated database of glucose-6-phosphate dehydrogenase (G6PD) mutations. *Hum Mutat*, **19**, 217–224.
- Lee, B. K. and Richards, F. M. (1971). The interpretation of protein structures: Estimation of static accessibility. *J Mol Biol*, **55**, 379–400.
- Letunic, I., Doerks, T., and Bork, P. (2012). SMART 7: Recent updates to the protein domain annotation resource. *Nucleic Acids Res*, **40**, D302–D305.
- Li, B., Krishnan, V. G., Mort, M. E., Xin, F., Kamati, K. K., Cooper, D. N., Mooney, S. D., and Radivojac, P. (2009). Automated inference of molecular mechanisms of disease from amino acid substitutions. *Bioinformatics*, **25**, 2744–2750.
- Lopes, L. R., Zekavati, A., Syrris, P., Hubank, M., Giambartolomei, C., Dalageorgou, C., Jenkins, S., McKenna, W., UK10k Consortium, Plagnol, V., and Elliott, P. M. (2013). Genetic complexity in hypertrophic cardiomyopathy revealed by high-throughput sequencing. *J Med Genet*, **50**, 228–239.
- Martin, A. C. R., Facchiano, A. M., Cuff, A. L., Hernandez-Boussard, T., Olivier, M., Hainaut, P., and Thornton, J. M. (2002). Integrating mutation data and structural analysis of the TP53 tumor-suppressor protein. *Hum Mutat*, **19**, 149–164.
- Mohamed, A. M. and Chin-Sang, I. D. (2006). Characterization of loss-of-function and gain-of-function Eph receptor tyrosine kinase signaling in *C. elegans* axon targeting and cell migration. *Developmental Biology*, **290**, 164–176.
- Ng, P. C. and Henikoff, S. (2003). SIFT: Predicting amino acid changes that affect protein function. *Nucleic Acids Res*, **31**, 3812–3814.
- Pan, S., Caleshu, C. A., Dunn, K. E., Foti, M. J., Moran, M. K., Soyinka, O., and Ashley, E. A. (2012). Cardiac structural and sarcomere genes associated with cardiomyopathy exhibit marked intolerance of genetic variation. *Circ Cardiovasc Genet*, **5**, 602–610.
- Reumers, J., Schymkowitz, J., Ferkinghoff-Borg, J., Stricher, F., Serrano, L., and Rousseau, F. (2005). SNPeff: a database mapping molecular phenotypic effects of human non-synonymous coding SNPs. *Nucleic Acids Res*, **33**, D527–D532.
- Reva, B., Antipin, Y., and Sander, C. (2011). Predicting the functional impact of protein mutations: Application to cancer genomics. *Nucleic Acids Res*, **39**, e118–e118.
- Richard, P., Charron, P., Carrier, L., Ledeuil, C., Cheav, T., Pichereau, C., Benaiche, A., Isnard, R., Dubourg, O., Burban, M., Gueffet, J.-P., Millaire, A., Desnos, M., Schwartz, K., Hainque, B., Komajda, M., and for the EUROGENE Heart Failure Project (2003). Hypertrophic cardiomyopathy: Distribution of disease genes, spectrum of mutations, and implications for a molecular diagnosis strategy. *Circulation*, **107**(17), 2227–2232.

- Schwarz, J. M., Rödelberger, C., Schuelke, M., and Seelow, D. (2010). MutationTaster evaluates disease-causing potential of sequence alterations. *Nature Methods*, **7**, 575–576.
- Shihab, H. A., Gough, J., Cooper, D. N., Stenson, P. D., Barker, G. L. A., Edwards, K. J., Day, I. N. M., and Gaunt, T. R. (2013). Predicting the functional, molecular, and phenotypic consequences of amino acid substitutions using hidden Markov models. *Hum Mutat*, **34**, 57–65.
- Spudich, J. A. (2014). Hypertrophic and dilated cardiomyopathy: four decades of basic research on muscle lead to potential therapeutic approaches to these devastating genetic diseases. *Biophys J*, **106**, 1236–1249.
- Stead, L. F., Wood, I. C., and Westhead, D. R. (2011). Kvsnp: accurately predicting the effect of genetic variants in voltage-gated potassium channels. *Bioinformatics*, **27**(16), 2181–2186.
- Stenson, P. D., Ball, E. V., Mort, M., Phillips, A. D., Shaw, K., and Cooper, D. N. (2002). *The Human Gene Mutation Database (HGMD) and Its Exploitation in the Fields of Personalized Genomics and Molecular Evolution*, chapter 1, Unit 1.13. John Wiley & Sons, Inc.
- Stitzel, N. O., Binkowski, T. A., Tseng, Y. Y., Kasif, S., and Liang, J. (2004). topoSNP: a topographic database of non-synonymous single nucleotide polymorphisms with and without known disease association. *Nucleic Acids Res*, **32**, D520–D522.
- UniProt Consortium (2014). Activities at the Universal Protein Resource (UniProt). *Nucleic Acids Res*, **42**, 7486–7486.
- Uzun, A., Leslin, C. M., Abyzov, A., and Ilyin, V. (2007). Structure SNP (StSNP): a web server for mapping and modeling nsSNPs on protein structures with linkage to metabolic pathways. *Nucleic Acids Res*, **35**, W384–W392.
- Walsh, R., Rutland, C., Thomas, R., and Loughna, S. (2010). Cardiomyopathy: a systematic review of disease-causing mutations in myosin heavy chain 7 and their phenotypic manifestations. *Cardiology*, **115**, 49–60.
- Witten, I. H., Frank, E., and Hall, M. A. (2011). *Data Mining: Practical Machine Learning Tools and Techniques*. Morgan Kaufmann, Burlington, MA 01803, USA, third edition.
- Woo, A., Rakowski, H., Liew, J. C., Zhao, M.-S., Liew, C.-C., Parker, T. G., Zeller, M., Wigle, E. D., and Sole, M. J. (2003). Mutations of the beta myosin heavy chain gene in hypertrophic cardiomyopathy: Critical functional sites determine prognosis. *Heart*, **89**, 1179–1185.
- Worth, C., Preissner, R., and Blundell, T. (2011). SDM — a server for predicting effects of mutations on protein stability and malfunction. *Nucleic Acids Res*, **39**, W215–W222.
- Yates, C. M., Filippis, I., Kelley, L. A., and Sternberg, M. J. (2014). SuSPect: enhanced prediction of single amino acid variant (SAV) phenotype using network features. *J Mol Biol*, **426**, 2692–2701.
- Yip, Y. L., Scheib, H., Diemand, A. V., Gattiker, A., Famiglietti, L. M., Gasteiger, E., and Bairoch, A. (2004). The Swiss-Prot variant page and the ModSNP database: a resource for sequence and structure information on human protein variants. *Hum Mutat*, **23**, 464–470.
- Yue, P., Melamud, E., and Moul, J. (2006). SNPs3D: candidate gene and SNP selection for association studies. *BMC Bioinformatics*, **7**, 166–166.