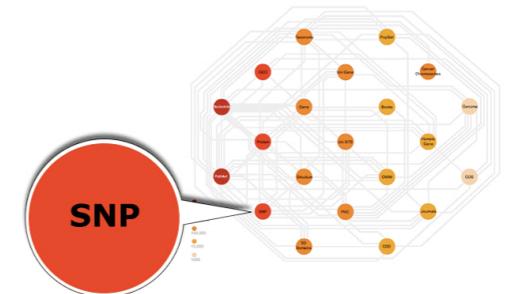
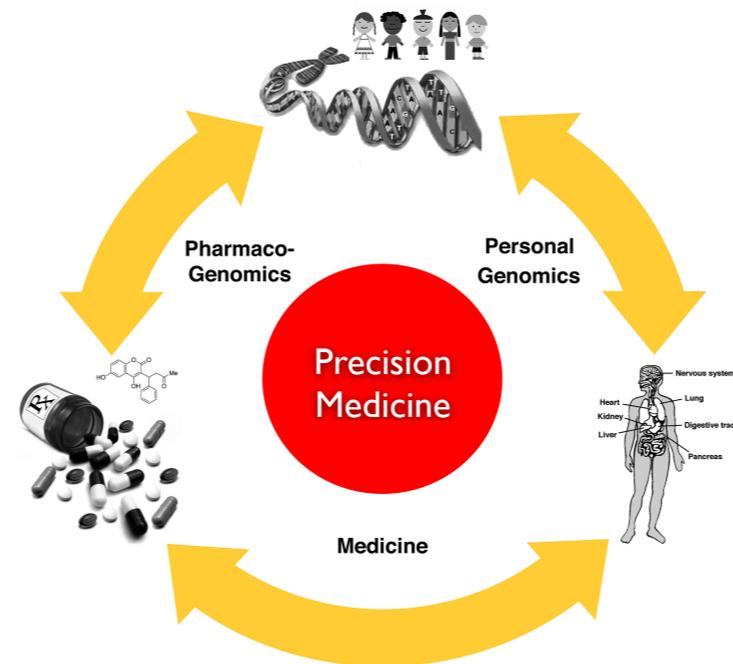
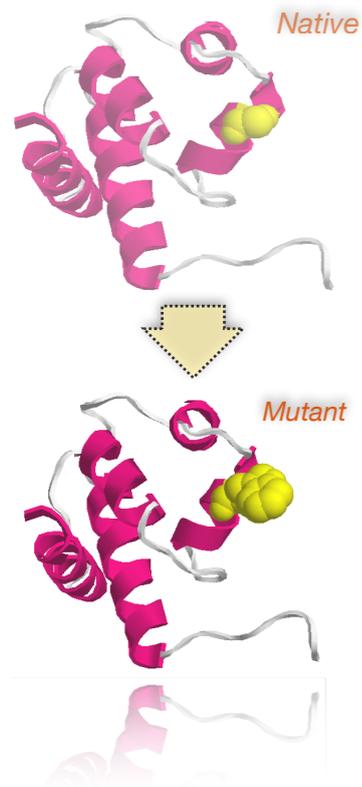


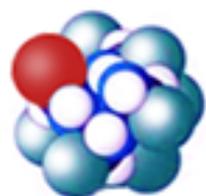
Computational methods for genome interpretation

MCCMB 2021

Monday August 2nd, 2021



Emidio Capriotti
<http://biofold.org/>



**Biomolecules
Folding and
Disease**

Department of Pharmacy
and Biotechnology (FaBiT)
University of Bologna



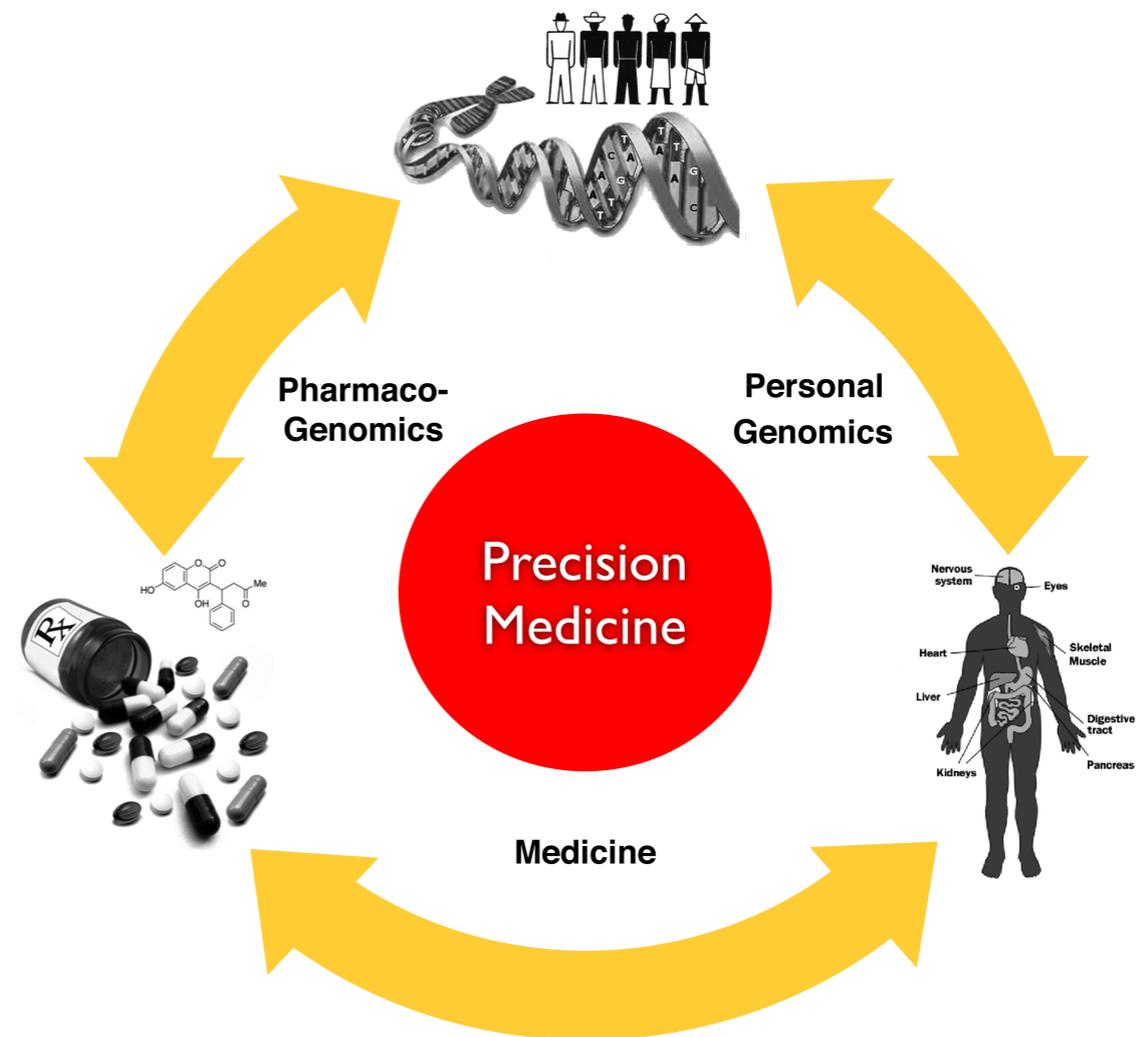
Presentation outline

- **Introduction:** Precision Medicine and Variant interpretation.
- **Protein variants:** sequence and structural features.
- **Meta prediction:** selection of highly-accurate predictions.
- **Impact of noncoding variants:** conservation in noncoding regions.
- **Prediction assessment:** The CAGI experiments.

Precision medicine

In the last decade, the cost of a **whole genome** sequencing experiment dropped below **\$1000**. The increasing amount of sequencing data is **raising important bioinformatics challenges**.

1. **Robust** sequencing data **processing methods**
2. **Interpretation** of the functional effect and the impact of genomic variations
3. Integrating the molecular mechanisms and data to **capture complexity of the system**
4. Make the data **clinically relevant**



Single Nucleotide Variants

Single Nucleotide Variants (SNVs)

is a DNA sequence variation occurring when a single nucleotide A, T, C, or G in the genome differs between members of the species.

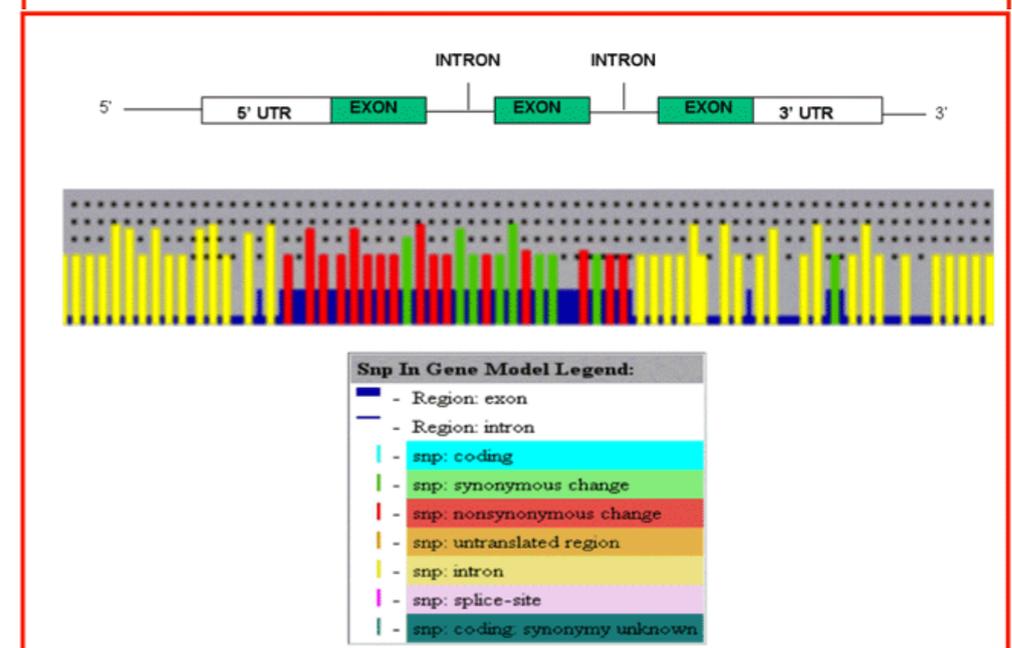
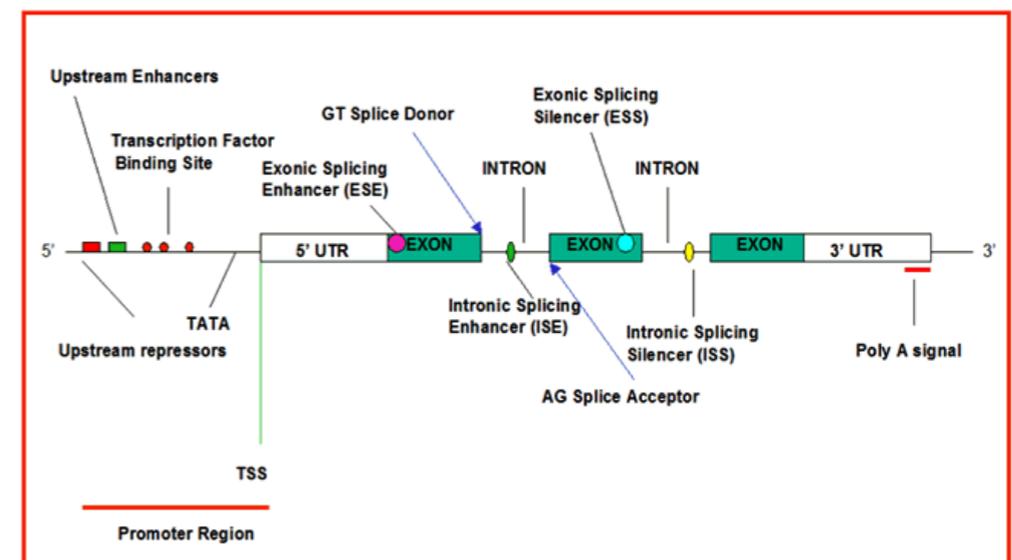
It is used to refer to Polymorphisms when the population frequency is $\geq 1\%$

SNVs occur at any position and can be classified on the base of their locations.

Coding SNVs can be subdivided into two groups:

Synonymous: when single base substitutions do not cause a change in the resultant amino acid

Non-synonymous or Single Amino Acid Variants (SAVs): when single base substitutions cause a change in the resultant amino acid.



1000 Genomes

The 1000 Genomes Project aims to create the **largest public catalogue of human variations and genotype data**. Last version released the genotype of **~2,500 individuals**.

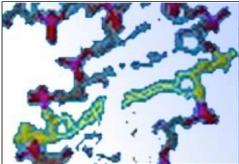
Table 1 | Variants discovered by project, type, population and novelty

a Summary of project data including combined exon populations

Statistic	Low coverage				Trios			Exon (total)	Union across projects
	CEU	YRI	CHB+JPT	Total	CEU	YRI	Total		
Samples	60	59	60	179	3	3	6	697	742
Total raw bases (Gb)	1,402	874	596	2,872	560	615	1,175	845	4,892
Total mapped bases (Gb)	817	596	468	1,881	369	342	711	56	2,648
Mean mapped depth (x)	4.62	3.42	2.65	3.56	43.14	40.05	41.60	55.92	NA
Bases accessed (% of genome)	2.43 Gb (86%)	2.39 Gb (85%)	2.41 Gb (85%)	2.42 Gb (86.0%)	2.26 Gb (79%)	2.21 Gb (78%)	2.24 Gb (79%)	1.4 Mb	NA
No. of SNPs (% novel)	7,943,827 (33%)	10,938,130 (47%)	6,273,441 (28%)	14,894,361 (54%)	3,646,764 (11%)	4,502,439 (23%)	5,907,699 (24%)	12,758 (70%)	15,275,256 (55%)
Mean variant SNP sites per individual	2,918,623	3,335,795	2,810,573	3,019,909	2,741,276	3,261,036	3,001,156	763	NA
No. of indels (% novel)	728,075 (39%)	941,567 (52%)	666,639 (39%)	1,330,158 (57%)	411,611 (25%)	502,462 (37%)	682,148 (38%)	96 (74%)	1,480,877 (57%)
Mean variant indel sites per individual	354,767	383,200	347,400	361,669	322,078	382,869	352,474	3	NA
No. of deletions (% novel)	ND	ND	ND	15,893 (60%)	6,593 (41%)	8,129 (50%)	11,248 (51%)	ND	22,025 (61%)
No. of genotyped deletions (% novel)	ND	ND	ND	10,742 (57%)	ND	ND	6,317 (48%)	ND	13,826 (58%)
No. of duplications (% novel)	259 (90%)	320 (90%)	280 (91%)	407 (89%)	187 (93%)	192 (91%)	256 (92%)	ND	501 (89%)
No. of mobile element insertions (% novel)	3,202 (79%)	3,105 (84%)	1,952 (76%)	4,775 (86%)	1,397 (68%)	1,846 (78%)	2,531 (78%)	ND	5,370 (87%)
No. of novel sequence insertions (% novel)	ND	ND	ND	ND	111 (96%)	66 (86%)	174 (93%)	ND	174 (93%)

Variant databases

dbSNP @ NCBI



dbSNP
dbSNP contains human single nucleotide variations, microsatellites, and small-scale insertions and deletions along with publication, population frequency, molecular consequence, and genomic and RefSeq mapping information for both common variations and clinical mutations.

<http://www.ncbi.nlm.nih.gov/snp>

Single Nucleotide Variants

Homo sapiens 917,705,245

Clinvar @ NCBI



ClinVar
ClinVar aggregates information about genomic variation and its relationship to human health.

<https://www.ncbi.nlm.nih.gov/clinvar/>

Single Nucleotide Variants

Homo sapiens 872,786

Pathogenic 58,167

Benign 119,050

humavar @ UniProt



<https://www.uniprot.org/docs/humavar>

Single Amino acid Variants

Homo sapiens 79,745

Pathogenic 31,398

Benign 39,584

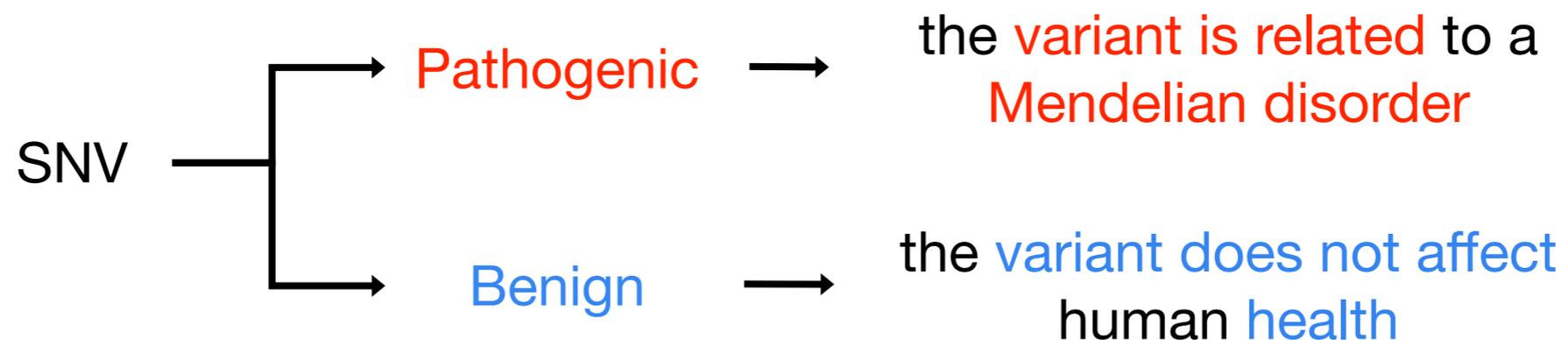
Effects of variants

Impact of **coding variants**

- Physico-chemical properties of the substituted residue
- Evolutionary important residues in specific protein sites
- Sequence–function relationships
- Structure–function relationships

Impact of **noncoding variants**

- Transcription
- Pre-mRNA splicing
- MicroRNA binding
- Altering post-translational modification sites



Protein variants

Sequence, Structure & Function

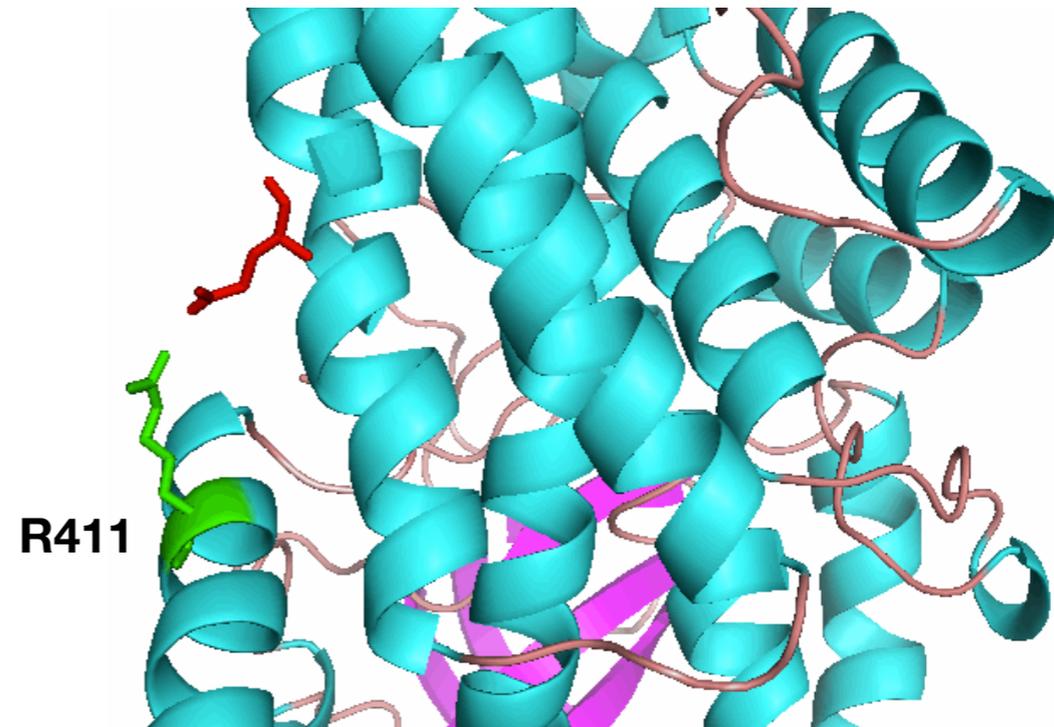
Genomic **variants in sequence motifs can affect protein function.**

Mutation S362A of P53 affect the interaction with hydrolase USP7 and the deubiquitination of the protein.



A **nonsynonymous variant** can affect the **protein structure causing the loss of stability** of the protein.

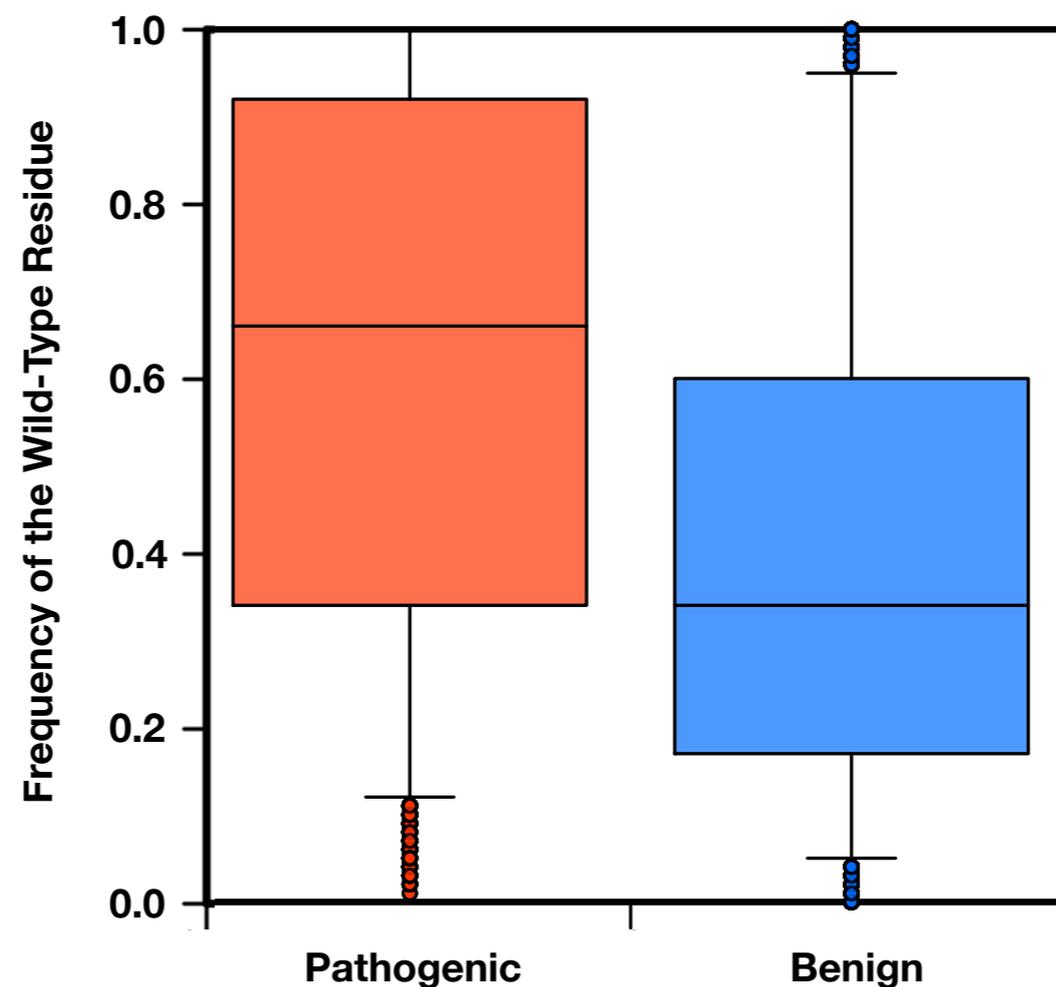
Mutation R411L results in the loss of a salt bridge, destabilizing the structure of the IVD dehydrogenase.



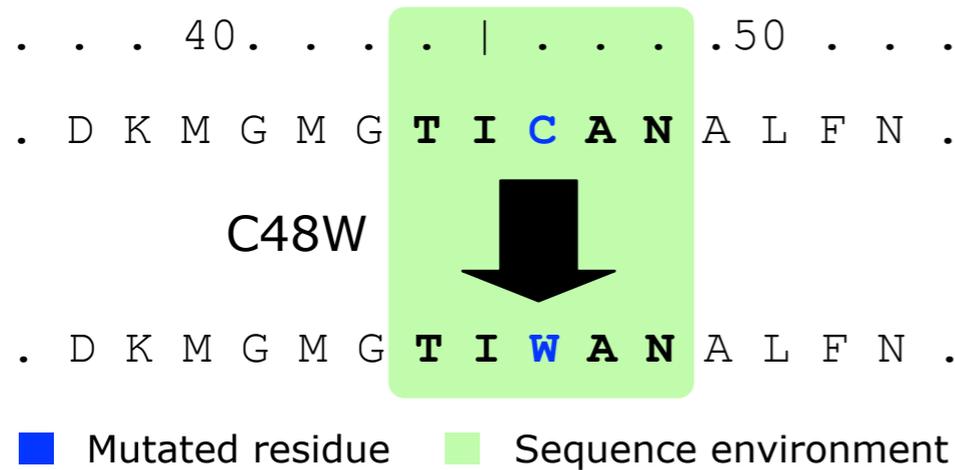
Sequence profile

The protein **sequence profile** is calculated running **BLAST on the UniRef90** dataset and selecting only the hits with e-value $< 10^{-9}$.

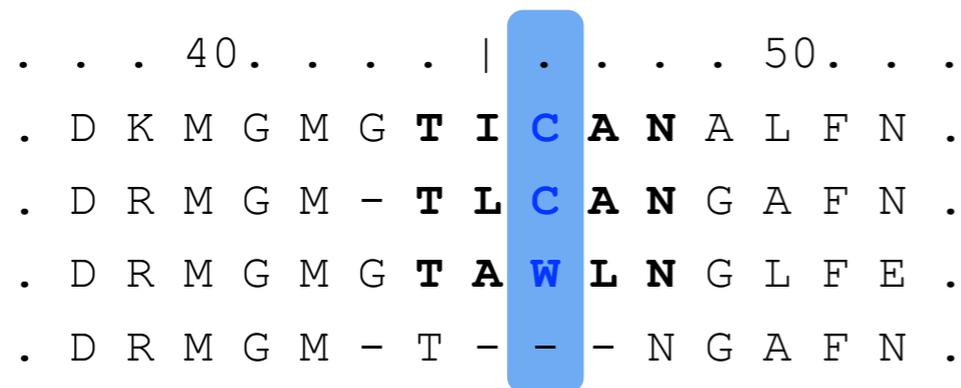
The **distributions of the frequency of the wild-type residues** for Pathogenic and Benign variants are significantly different.



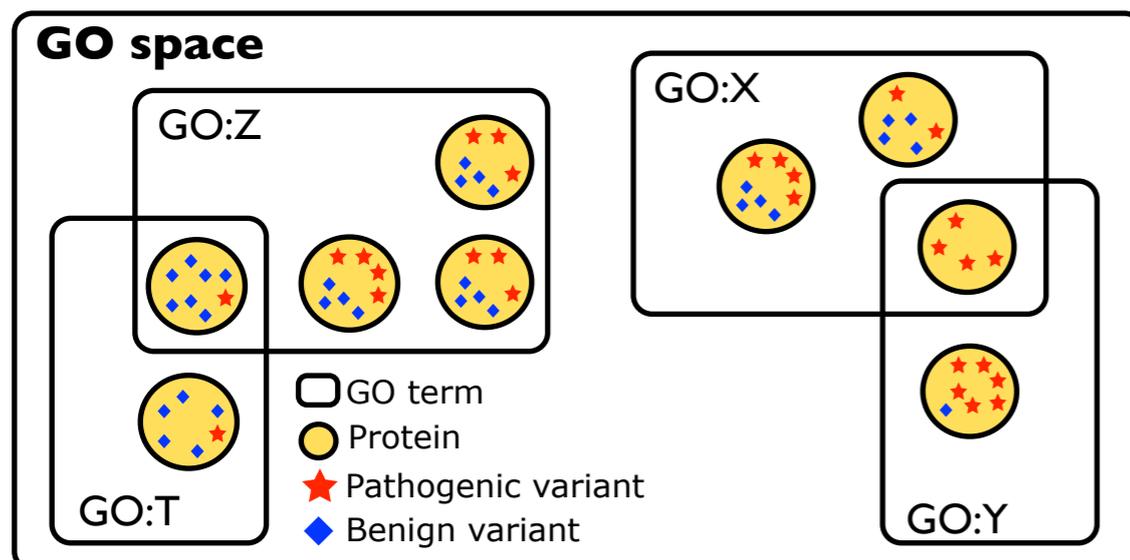
SNPs&GO input features



Sequence information is encoded in 2 vectors each one composed by 20 elements. The **first vector encodes for the mutation** and the **second one for the sequence environment**



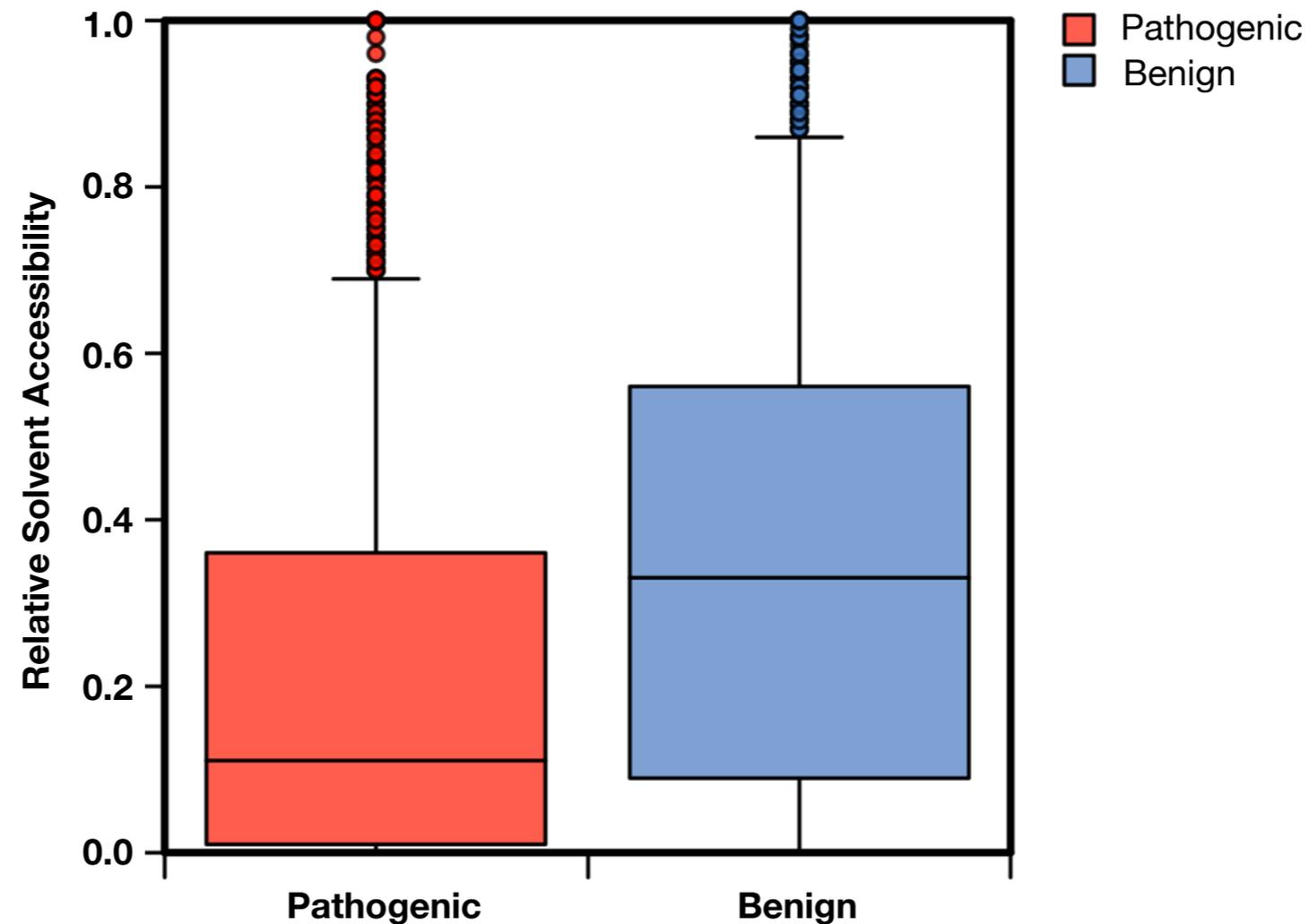
Protein sequence **profile information derived from a multiple sequence alignment**. It is encoded in a **5 elements vector** corresponding to different features general and local features



The **GO information** are encoded in a **2 elements vector** corresponding to the **number unique of GO terms** associated to the protein sequences and the **sum of the logarithm of the total number of Pathogenic and Benign variants for each GO term**.

Structure environment

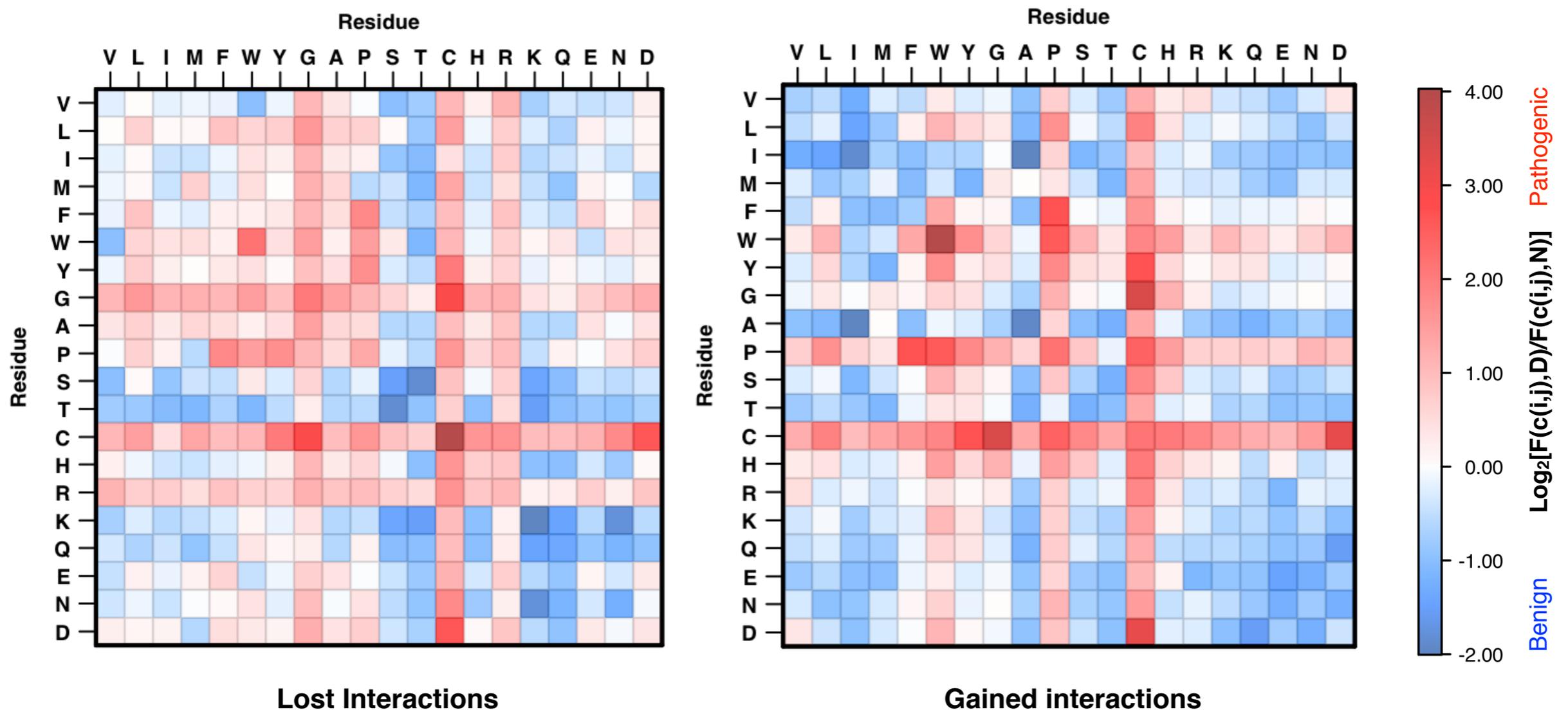
There is a **significant difference** between the **distributions of the Relative Solvent Accessibility for Pathogenic and Benign** variants. The median values of their distributions are ~ 0.1 and 0.35 respectively.



Analysis of the 3D interactions

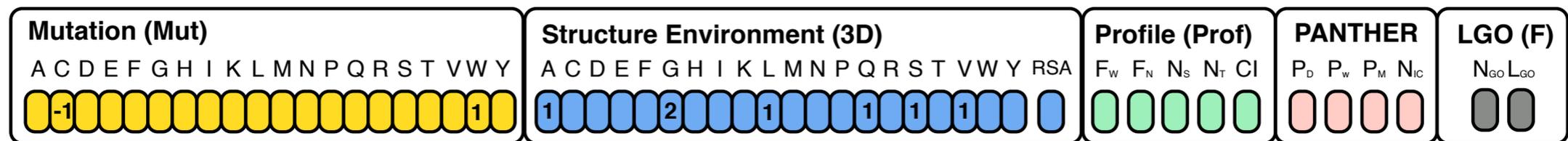
Using the **whole set of SAVs with known structure**, we calculate the **log odd score** of the **ratio** between the **frequencies of the interaction between residue i and j** for **Pathogenic and Benign variants**.

$$LC = \log_2 \left[\frac{n(i, j, Pathogenic) / N(Pathogenic)}{n(i, j, Benign) / N(Benign)} \right]$$



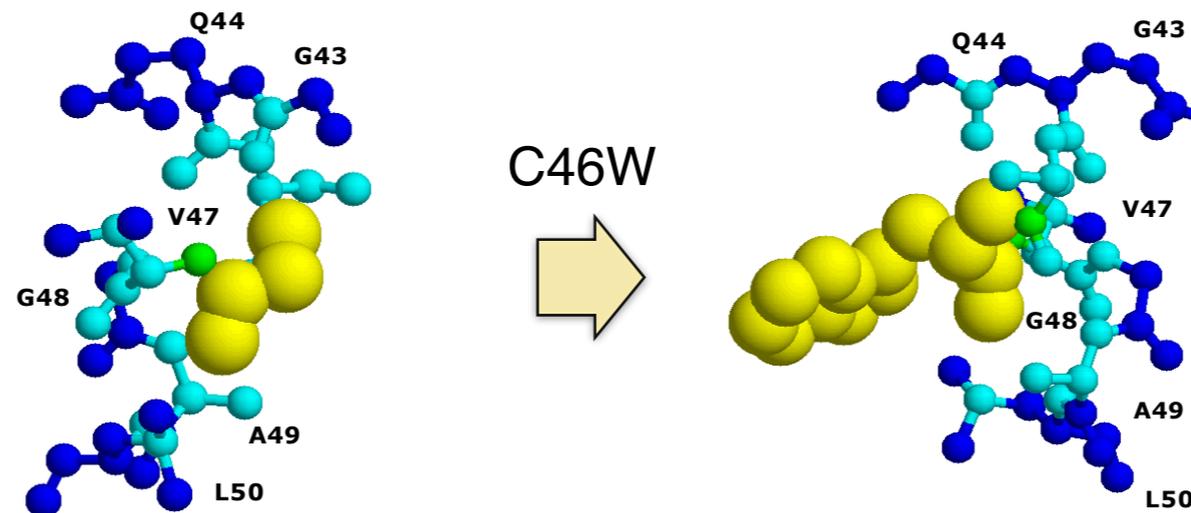
The structure-based method

The method takes as input a 52-element vector encoding for mutation; structure environment, sequence profile and functional score based on GO terms.



RBF Kernel

Output

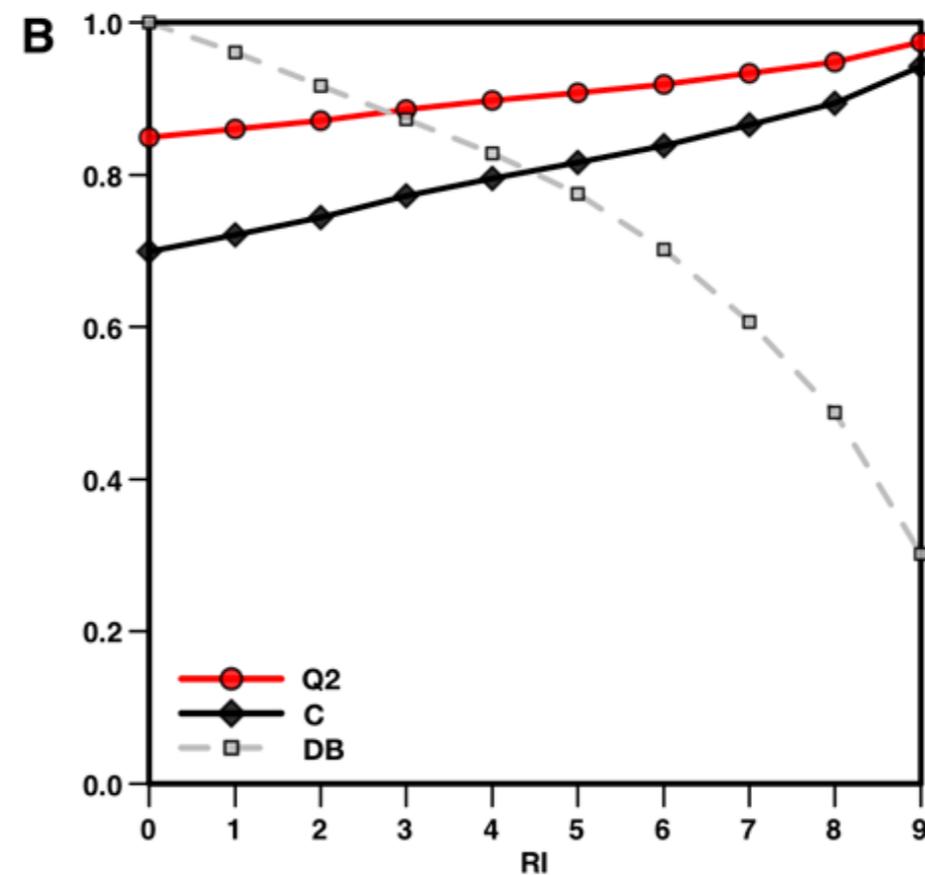
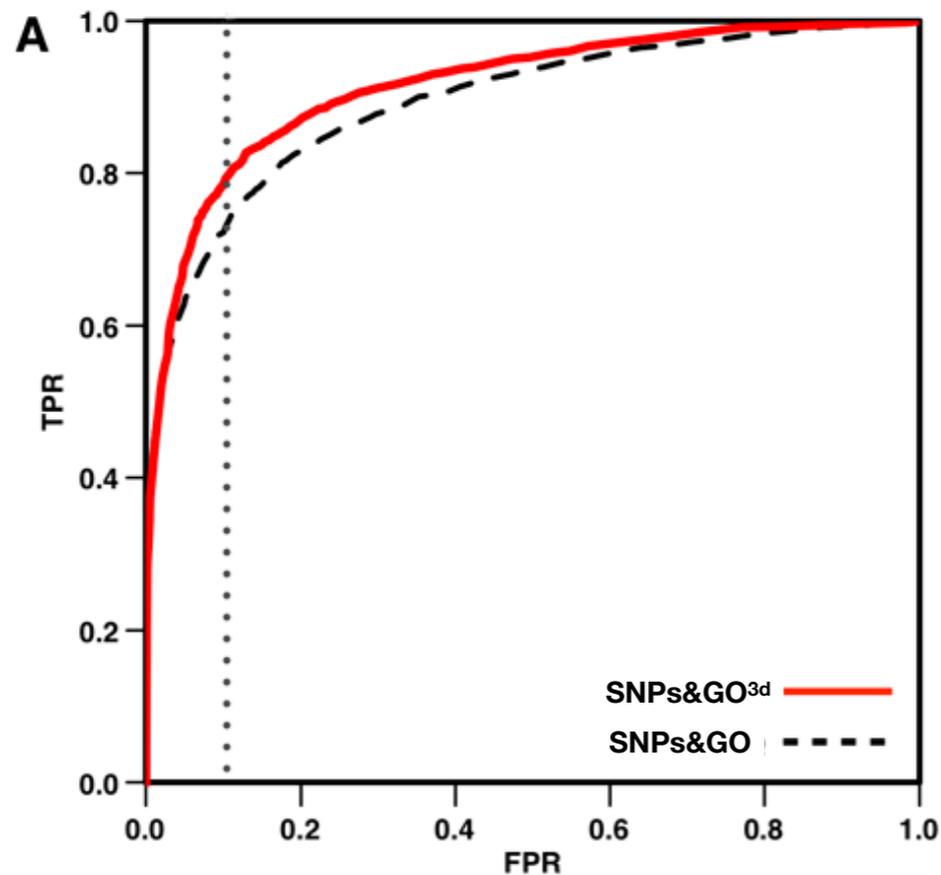


■ Mutated Aminoacid ■ 0 < R < 2Å ■ 2 < R < 4Å ■ 4 < R < 6Å

Sequence vs structure

The structure-based method results in better accuracy with respect to the sequence-based one. Structure based prediction are 3% more accurate and correlation coefficient increases of 0.06. If 10% of FP are accepted the TPR increases of 7%.

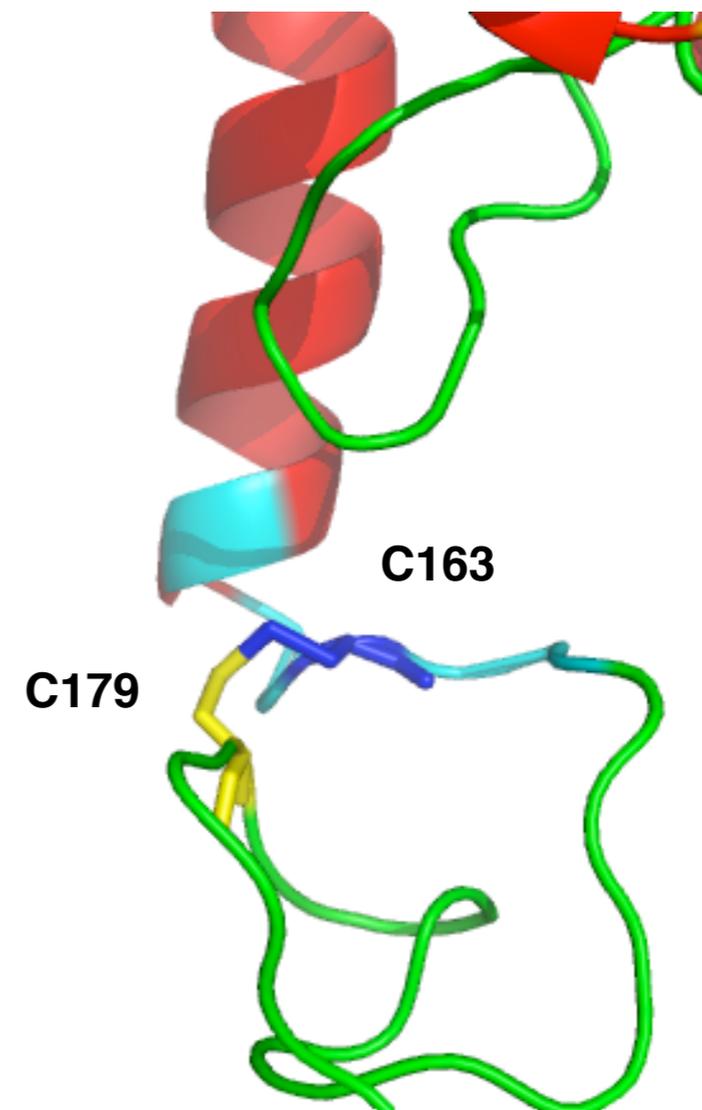
	Q2	P[D]	S[D]	P[N]	S[N]	C	AUC
SNPs&GO	0.82	0.81	0.83	0.82	0.81	0.64	0.89
SNPs&GO ^{3d}	0.85	0.84	0.87	0.86	0.83	0.70	0.92



Prediction example

Damaging missing Cys-Cys interaction in the Glycosylasparaginase. The mutation p.Cys163Ser results in the loss of the disulfide bridge between Cys163 and Cys179. This amino acid variant is responsible for Aspartylglucosaminuria.

1APY: Chain A, Res: 2.0 Å



Meta prediction approach

Protein variant predictors

Many predictor of the effect of Single Amino acid Variants (SAVs) are available. They mainly use **information from multiple sequence alignment** to predict the effect of a given mutation. In this study we consider

- **PhD-SNP**: Support Vector Machine-based method using sequence and profile information (Capriotti et al. 2006).
- **PANTHER**: Hidden Markov Model-based method using a HMM library of protein families (Thomas and Kejariwal 2004).
- **SNAP**: Neural network based method to predict the functional effect of single point mutations (Bromberg et al. 2008).
- **SIFT**: Probabilistic method based on the analysis of multiple sequence alignments (Ng and Henikoff 2003).

Predictors accuracy

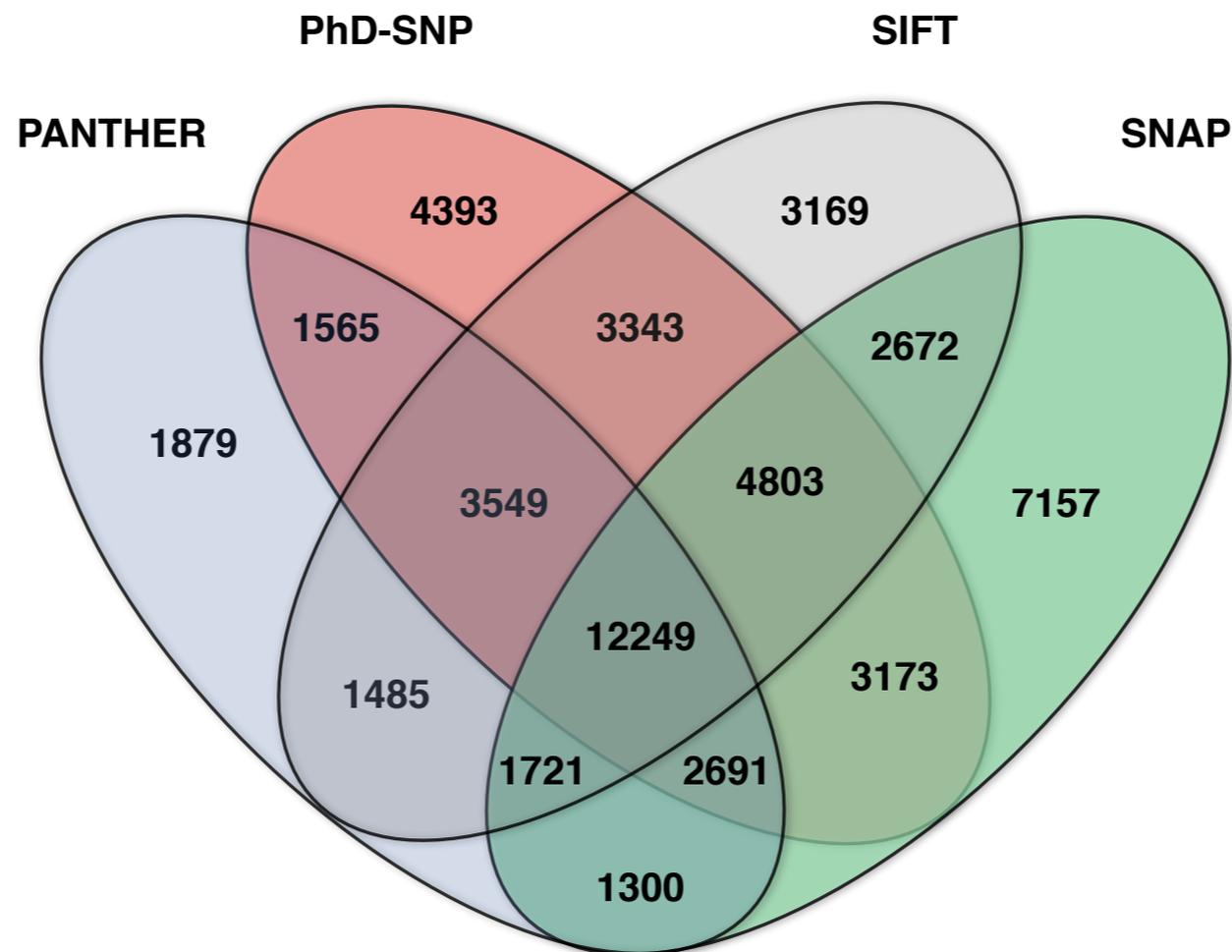
The accuracy of each predictor has been tested on a set of 35,986 mutations equally distributed between Pathogenic and Benign variants. **PhD-SNP results in better accuracy but is the only one optimized** using a cross-validation procedure. **SNAP** shows lowest accuracy **but it has been developed for a different task.**

	Q2	P[D]	S[D]	P[N]	S[N]	C	PM
PhD-SNP	0.76	0.78	0.74	0.75	0.78	0.53	100
PANTHER	0.74	0.79	0.73	0.69	0.74	0.48	74
SNAP	0.64	0.59	0.90	0.79	0.38	0.33	100
SIFT	0.70	0.74	0.64	0.68	0.76	0.41	92

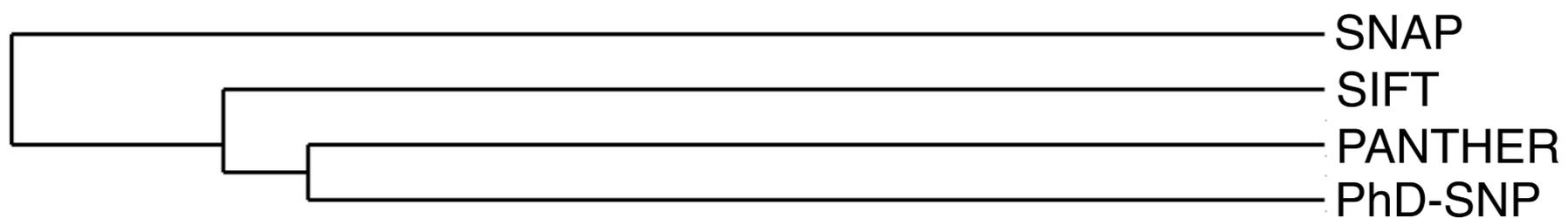
DB: Benign (N) 17883 and Pathogenic (D) 17883

Predictors tree

Using the prediction similarity we can build the predictors tree



UPGMA tree based on correlations



Prediction analysis

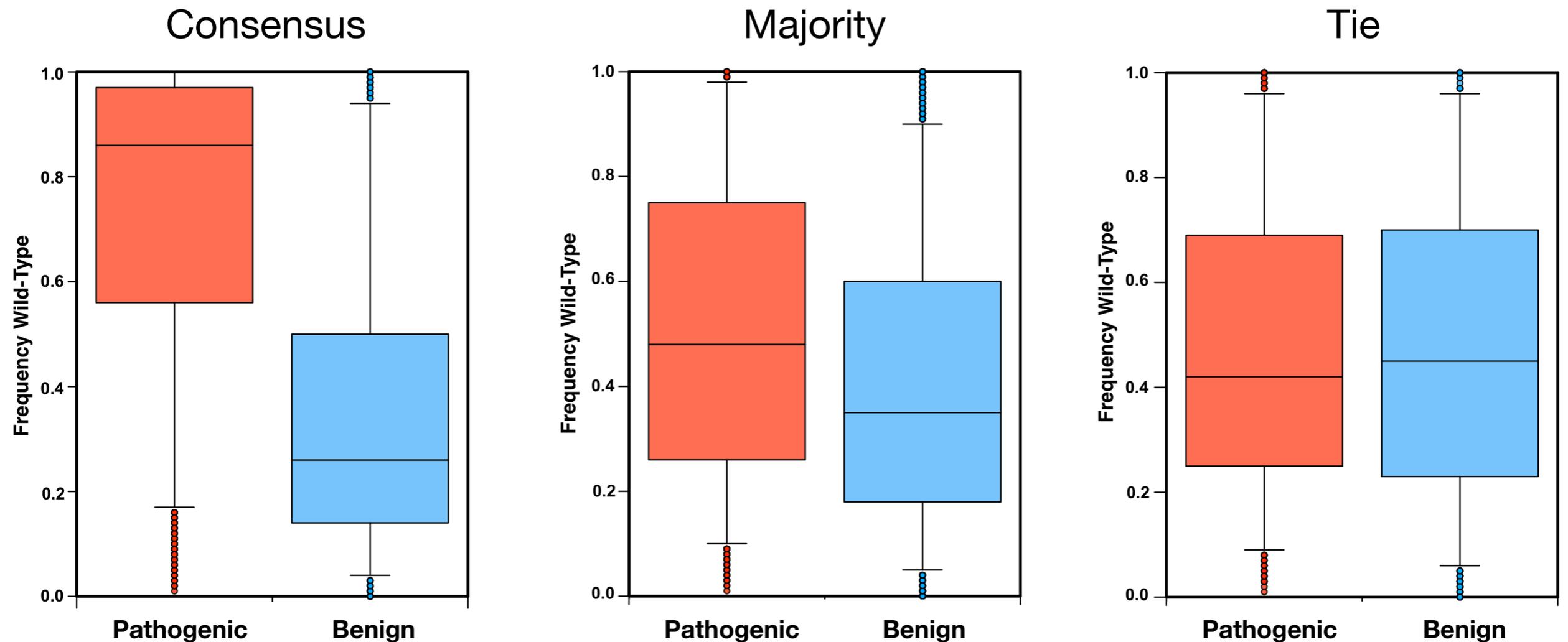
The accuracy of the predictions has been evaluated considering three different subset

- **Consensus:** all the predictions returned by the methods are in agreement.
- **Tie:** equal number of methods predicting Pathogenic and Benign
- **Majority:** One of the two possible classes is predominant

	Q2	P[D]	S[D]	P[N]	S[N]	C	AUC	%DB
PhD-SNP	0.76	0.78	0.74	0.75	0.78	0.53	0.84	100
Consensus	0.87	0.87	0.92	0.87	0.79	0.73	0.89	46
Majority	0.70	0.67	0.56	0.72	0.80	0.37	0.82	40
Tie	0.61	0.51	0.43	0.66	0.73	0.16	0.67	14

Subset conservation

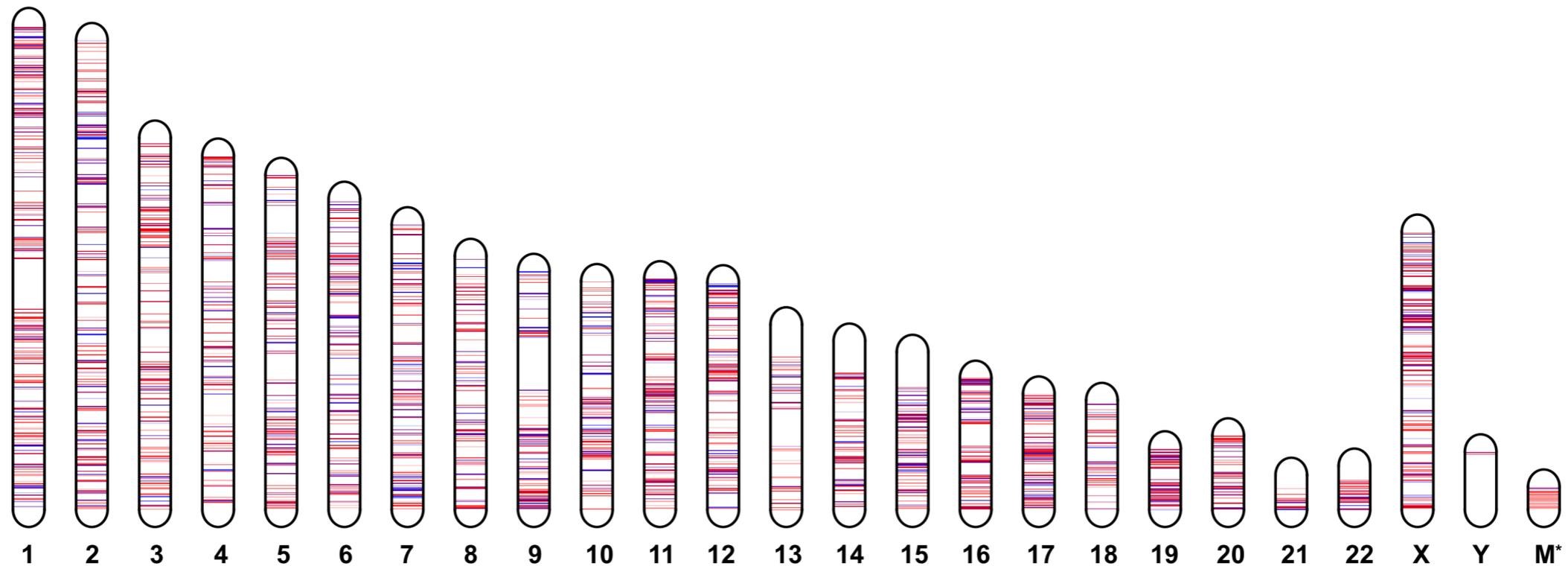
The **distributions of the wild-type frequencies** for Pathogenic and Benign variants **on the *Consensus* subset have very little overlap.**



From coding to noncoding

Whole-genome predictions

Most of the genetic variants occur in noncoding region that represents >98% of the whole genome.

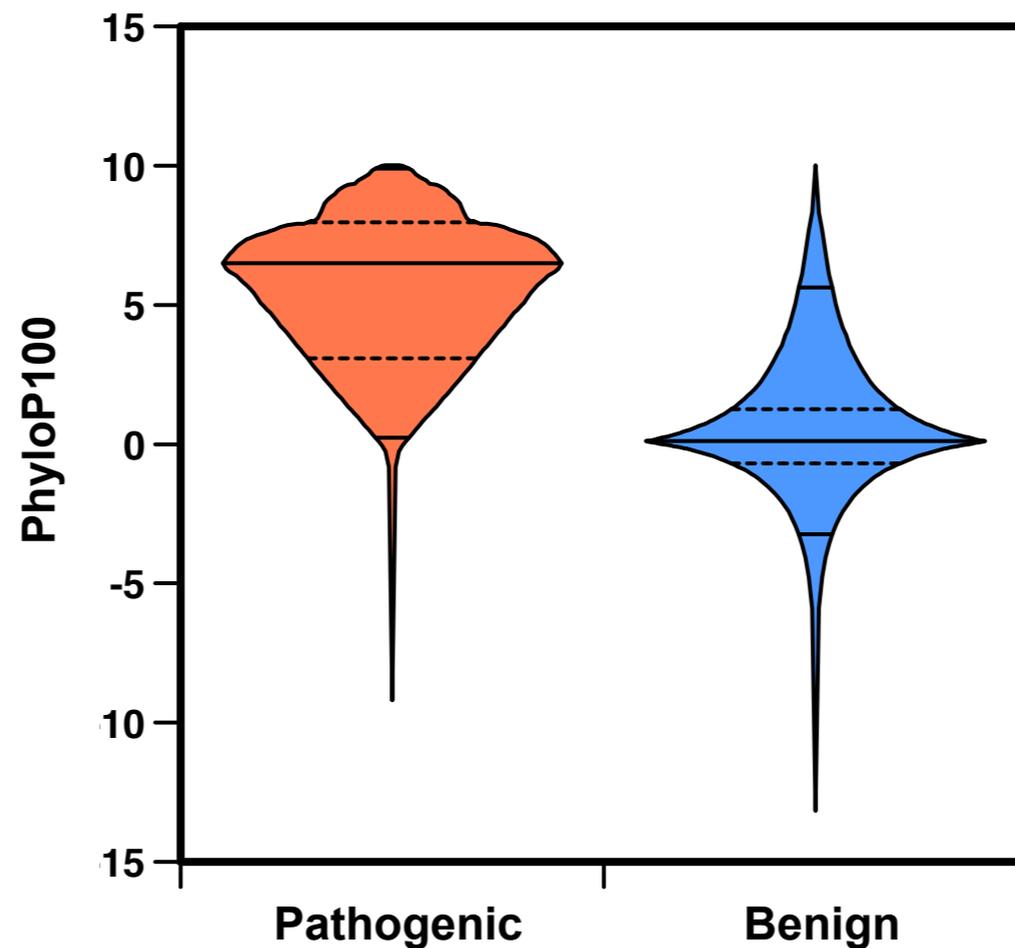


Predict the effect of SNVs in noncoding region is a challenging task because conservation is more difficult to estimate.

The sequence alignment is more complex task for noncoding regions.

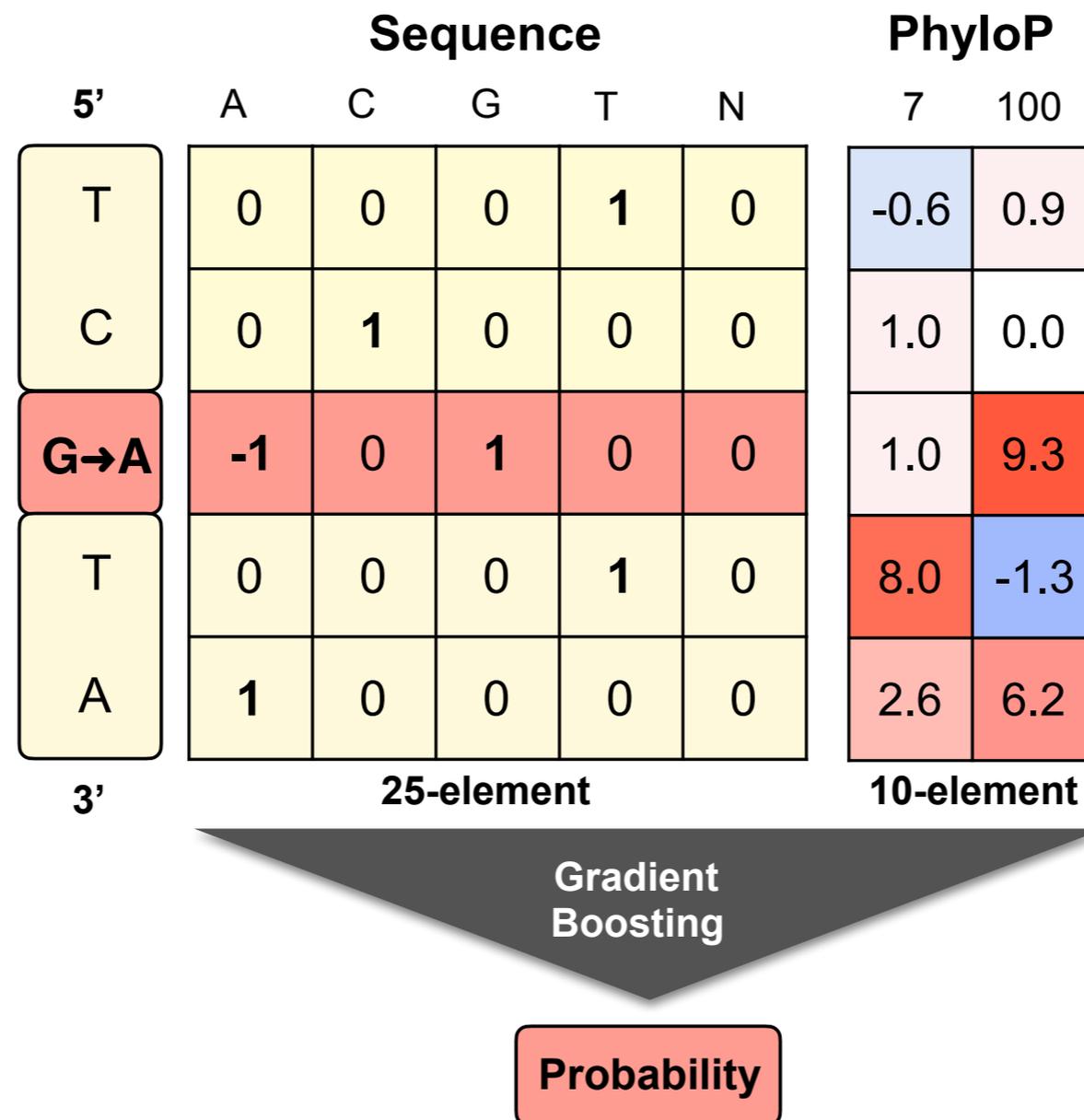
PhyloP100 score

Conservation analysis based on the pre-calculated score available at the UCSC revealed a **significant difference between the distribution of the PhyloP100 scores in Pathogenic and Benign SNVs.**



PhD-SNPg

PhD-SNPg is a simple method that takes in input **35 sequence-based features** from a window of 5 nucleotides around the mutated position.

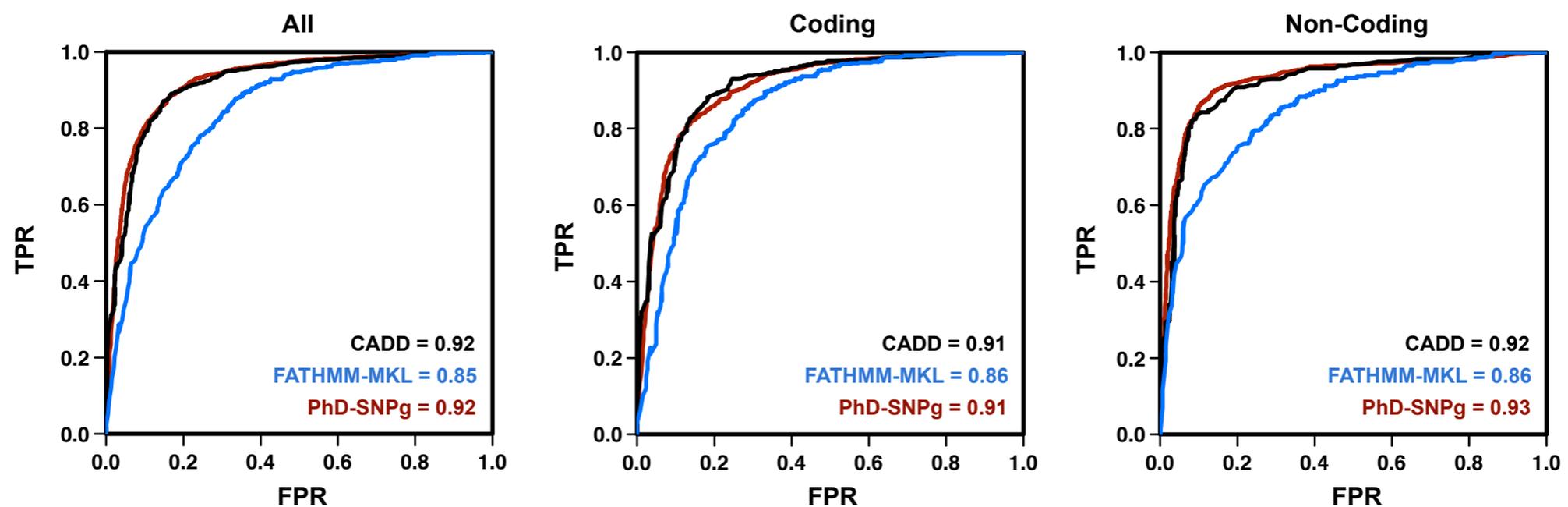


<http://snps.biofold.org/phd-snp-g/>

Benchmarking

PhD-SNP^g has been tested in cross-validation on a set of 35,802 SNVs and on a blind set of 1,408 variants recently annotated.

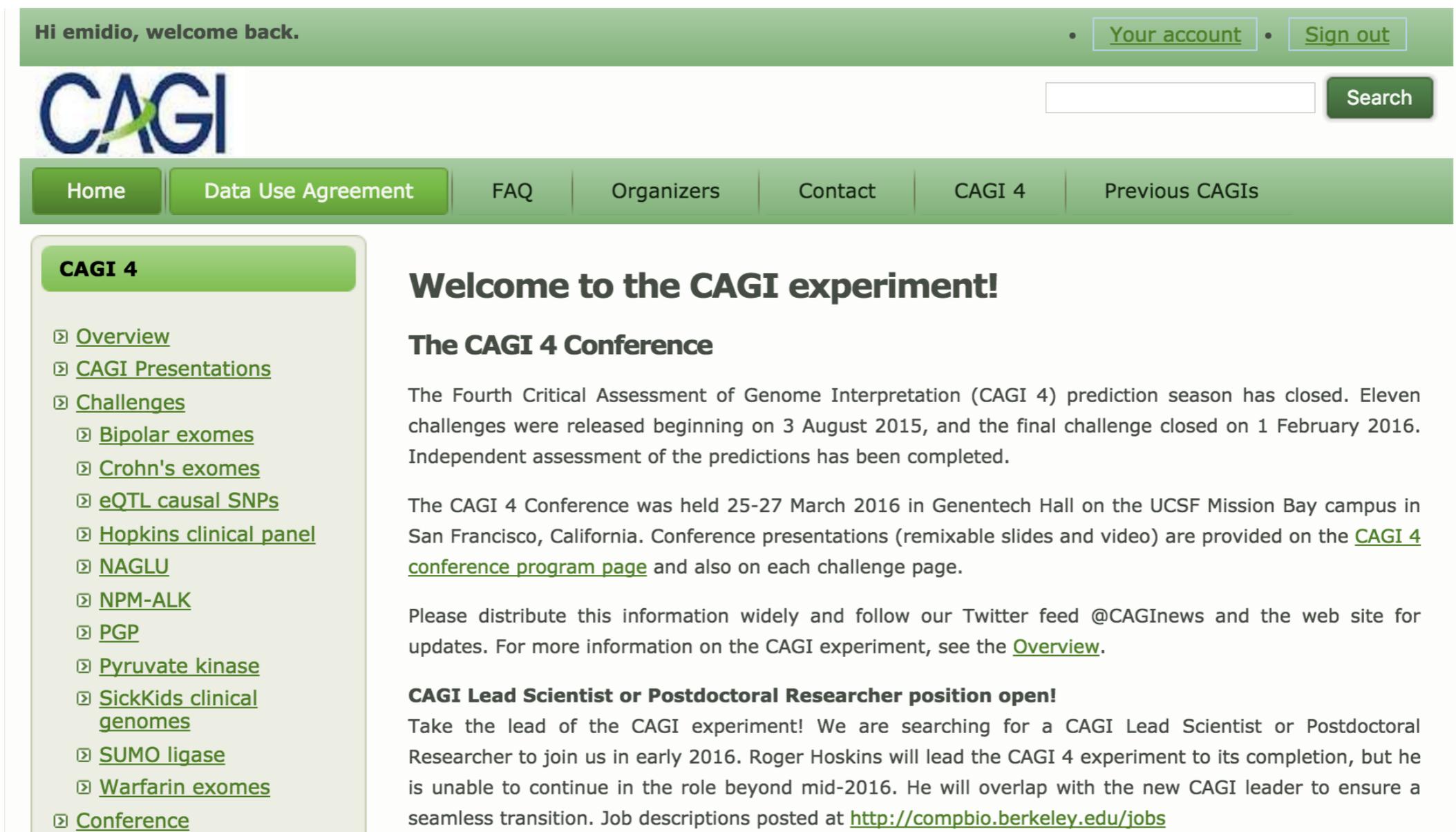
	Q2	TNR	NPV	TPR	PPV	MCC	F1	AUC
PhD-SNP^g	0.861	0.774	0.884	0.925	0.847	0.715	0.884	0.924
Coding	0.849	0.671	0.845	0.938	0.850	0.651	0.892	0.908
Non-Coding	0.876	0.855	0.911	0.901	0.839	0.753	0.869	0.930



Blind testing

CAGI experiments

The Critical Assessment of Genome Interpretation is a community experiment to objectively assess computational methods for predicting the phenotypic impacts of genomic variation.



The screenshot shows the CAGI website homepage. At the top, a green banner displays a user greeting: "Hi emidio, welcome back." followed by links for "Your account" and "Sign out". Below this is the CAGI logo and a search bar with a "Search" button. A navigation menu contains links for "Home", "Data Use Agreement", "FAQ", "Organizers", "Contact", "CAGI 4", and "Previous CAGIs". The main content area features a "CAGI 4" section header, a list of links for various challenges (e.g., Overview, CAGI Presentations, Challenges, Bipolar exomes, Crohn's exomes, eQTL causal SNPs, Hopkins clinical panel, NAGLU, NPM-ALK, PGP, Pyruvate kinase, SickKids clinical genomes, SUMO ligase, Warfarin exomes, Conference), and three main text blocks: "Welcome to the CAGI experiment!", "The CAGI 4 Conference", and "CAGI Lead Scientist or Postdoctoral Researcher position open!".

Hi emidio, welcome back. • [Your account](#) • [Sign out](#)

CAGI [Search](#)

[Home](#) [Data Use Agreement](#) [FAQ](#) [Organizers](#) [Contact](#) [CAGI 4](#) [Previous CAGIs](#)

CAGI 4

- [Overview](#)
- [CAGI Presentations](#)
- [Challenges](#)
 - [Bipolar exomes](#)
 - [Crohn's exomes](#)
 - [eQTL causal SNPs](#)
 - [Hopkins clinical panel](#)
 - [NAGLU](#)
 - [NPM-ALK](#)
 - [PGP](#)
 - [Pyruvate kinase](#)
 - [SickKids clinical genomes](#)
 - [SUMO ligase](#)
 - [Warfarin exomes](#)
- [Conference](#)

Welcome to the CAGI experiment!

The CAGI 4 Conference

The Fourth Critical Assessment of Genome Interpretation (CAGI 4) prediction season has closed. Eleven challenges were released beginning on 3 August 2015, and the final challenge closed on 1 February 2016. Independent assessment of the predictions has been completed.

The CAGI 4 Conference was held 25-27 March 2016 in Genentech Hall on the UCSF Mission Bay campus in San Francisco, California. Conference presentations (remixable slides and video) are provided on the [CAGI 4 conference program page](#) and also on each challenge page.

Please distribute this information widely and follow our Twitter feed @CAGInews and the web site for updates. For more information on the CAGI experiment, see the [Overview](#).

CAGI Lead Scientist or Postdoctoral Researcher position open!

Take the lead of the CAGI experiment! We are searching for a CAGI Lead Scientist or Postdoctoral Researcher to join us in early 2016. Roger Hoskins will lead the CAGI 4 experiment to its completion, but he is unable to continue in the role beyond mid-2016. He will overlap with the new CAGI leader to ensure a seamless transition. Job descriptions posted at <http://compbio.berkeley.edu/jobs>

<https://genomeinterpretation.org/>

The NAGLU challenge

NAGLU is a lysosomal glycohydrolyase which deficiency causes a rare disorder referred as Sanfilippo B disease

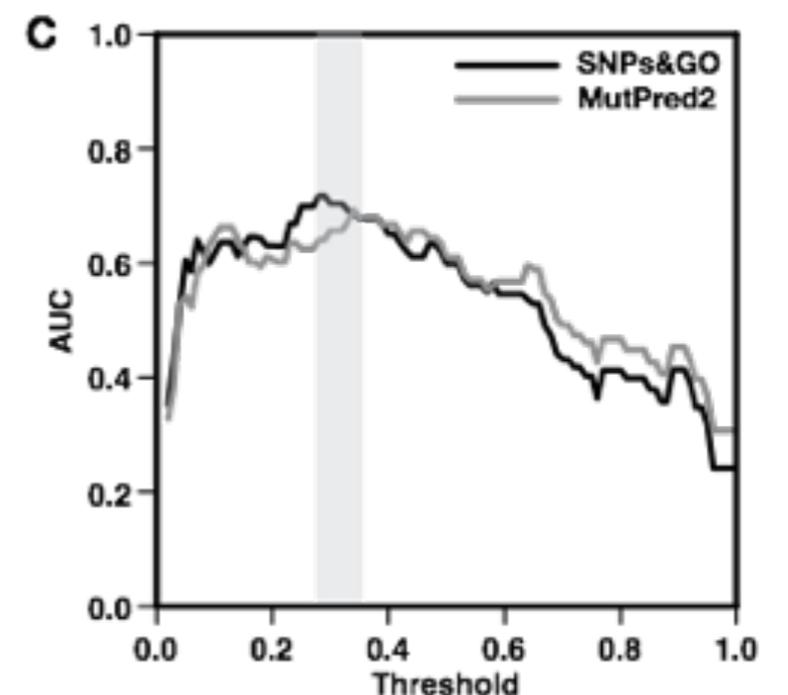
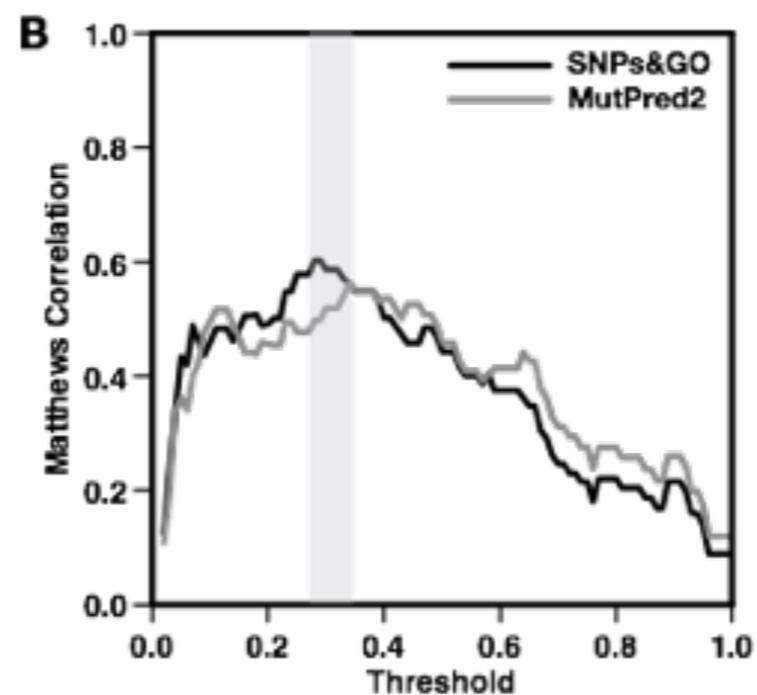
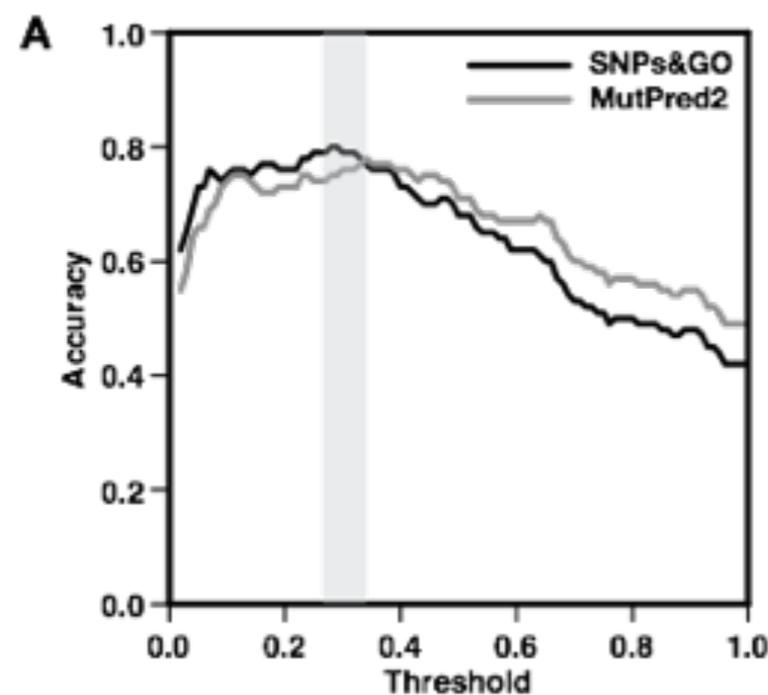
Challenge: Predict the effect of the 165 variants on NAGLU enzymatic activity.

The submitted prediction should be a **numeric value ranging from 0 (no activity) to 1 (wild-type level of activity)**.

A posteriori evaluation

An evaluation of the performance shows that **SNPs&GO** reaches similar accuracy than the best method (MutPred2)

Method	Q2	AUC	MC	RMSE	rPearson	rSpearman	rKendallTau
MutPred2	0.780	0.850	0.565	0.30	0.595	0.619	0.443
SNPs&GO	0.800	0.854	0.603	0.33	0.575	0.616	0.445
SNPs&GO ⁰⁹	0.750	0.749	0.499	0.46	0.477	0.495	0.409



Conclusions

- Evolutionary information is an important feature for the prediction of deleterious variants. The **pathogenic variants tend to occur in conserved protein sites.**
- Structural information encoded through the **relative solvent accessibility and the structure environment improves** the predictions of pathogenic variants.
- The implementation of **meta-prediction based approach allows to select highly-accurate predictions.**
- **Nucleotide conservation** is an important feature to **predict the impact of SNVs also in noncoding regions.**

Acknowledgments

Structural Genomics @CNAG

Marc A. Marti-Renom
Francois Serra

Spain

Computational Biology and Bioinformatics Research Group (UIB)

Jairo Rocha

Spain

Division of Informatics at UAB

Malay Basu

Division Clinical Immunology
& Rheumatology

Harry Schroeder

Mohamed Khass

USA

Helix Group (Stanford University)

Russ B. Altman

Jennifer Lahti

Tianyun Liu

Grace Tang

USA

Bologna Biocomputing Group

Rita Casadio

Pier Luigi Martelli

University of Torino

Piero Fariselli

University of Camerino

Mario Compiani

Italy

Mathematical Modeling of Biological Systems (University of Düsseldorf)

Markus Kollmann

Linlin Zhao

Germany

Other Collaborations

Yana Bromberg, Rutgers University, NJ (USA)

Hannah Carter, UCSD, CA (USA)

Francisco Melo, Universidad Catolica, (Chile)

Cedric Notredame, CRG Barcelona (Spain)

Gustavo Parisi, Univ. de Quilmes (Argentina)

Frederic Rousseau, KU Leuven (Belgium)

Joost Schymkowitz, KU Leuven (Belgium)

FUNDING

Italian MIUR: PRIN 2017

NIH: 1R21 AI134027- 01A1

Italian MIUR: FFABR 2017

UNIBO: International Cooperation

Startup funding Dept. of Pathology UAB

NIH:3R00HL111322-04S1 Co-Investigator

EMBO Short Term Fellowship

Marie Curie International Outgoing Grant

Marco Polo Research Project

BIOSAPIENS Network of Excellence

SPINNER Consortium

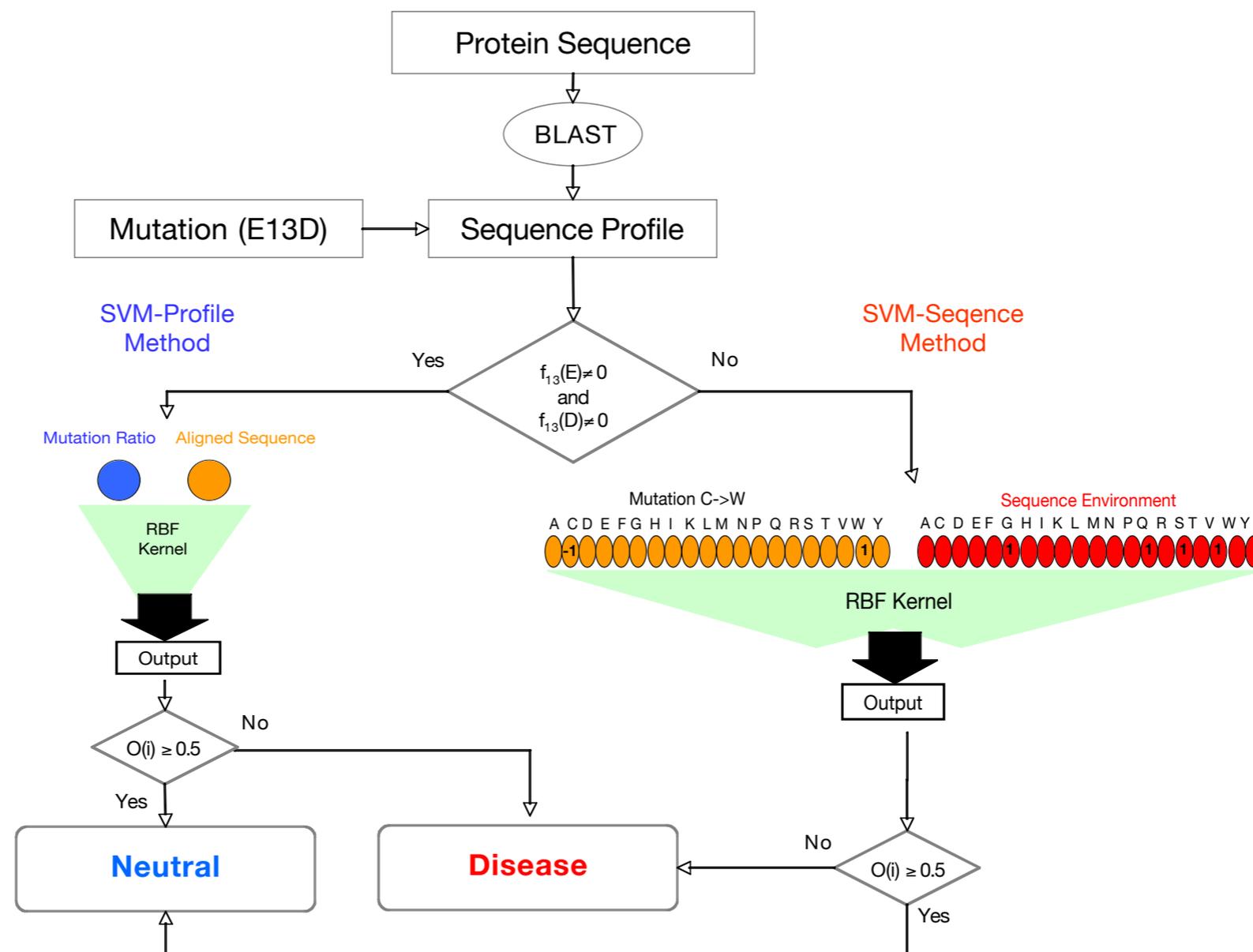
Biomolecules, Folding and Disease



<http://biofold.org/>

Hybrid method structure

Hybrid Method is based on a decision tree with **SVM-Sequence** coupled to **SVM-Profile**. Tested on more than 21,000 variants our method reaches 74% of accuracy and 0.46 correlation coefficient.



Classification results

SVM-Sequence is more accurate in the prediction of **disease related mutations** and **SVM-Profile** is more accurate in the prediction of **neutral polymorphism**.
Both methods have the **same Q2 level**.

	Q2	P[D]	Q[D]	P[N]	Q[N]	C
SVM-Sequence	0.70	0.71	0.84	0.65	0.46	0.34
SVM-Profile	0.70	0.74	0.49	0.68	0.86	0.39
HybridMeth	0.74	0.80	0.76	0.65	0.70	0.46

D = Disease related N = Neutral

The Hybrid Method have higher accuracy than the previous two methods **increasing the accuracy** up to 74% **and the correlation coefficient** up to 0.46.

<http://snps.biofold.org/phd-snp>

Gene Ontology

The **Gene Ontology project** is a major bioinformatics initiative with the aim of standardizing the **representation of gene and gene product attributes across species** and databases. The project provides a controlled vocabulary of terms for describing gene product characteristics and gene product annotation data.



<http://www.geneontology.org/>

The ontology is represented by a **direct acyclic graph covers three domains;**

- **cellular component**, the parts of a cell or its extracellular environment;
- **molecular function**, the elemental activities of a gene product at the molecular level, such as binding or catalysis
- **biological process**, operations or sets of molecular events with a defined beginning and end, pertinent to the functioning of integrated living units: cells, tissues, organs and organisms.

The P16 challenge

CDKN2A is the most common, high penetrance, susceptibility gene identified to date in **familial malignant melanoma**. **p16^{INK4A}** is one of the two **oncosuppressor** which promotes cell cycle arrest by inhibiting cyclin dependent kinase (CDK4/6).

Challenge: Evaluate how different variants of p16 protein impact its ability to block cell proliferation.

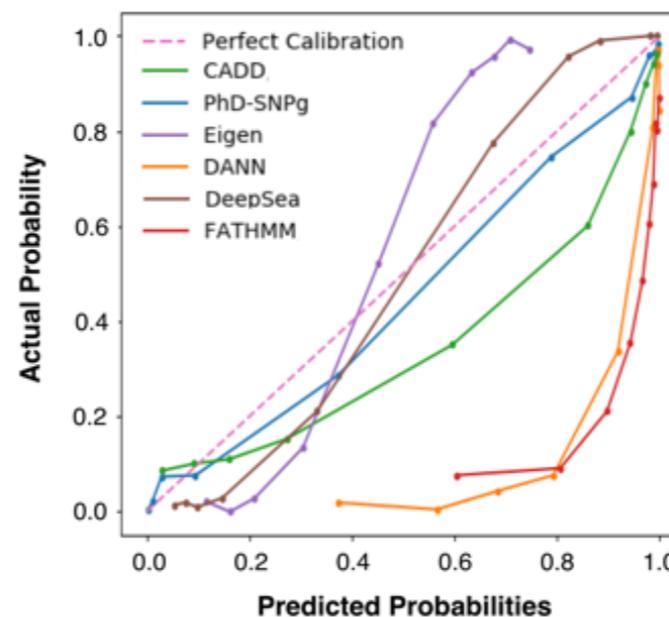
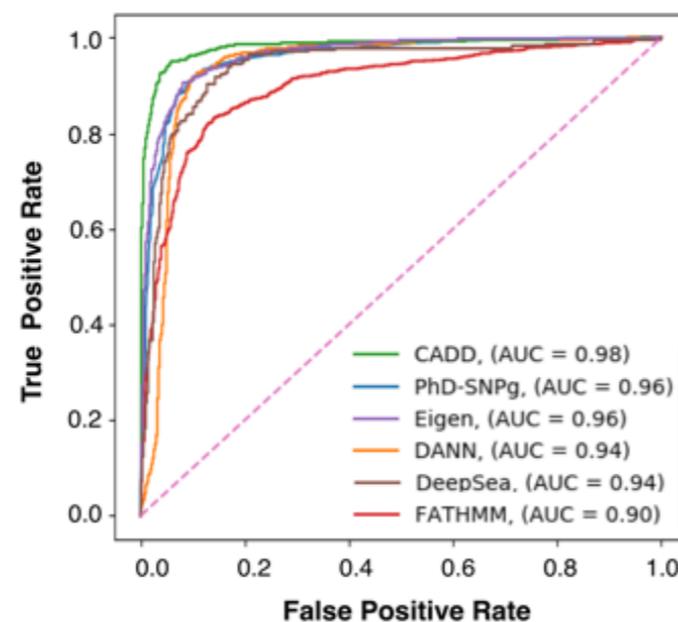
Provide a number between **50%** that represent the normal **proliferation rate of control cells** and **100%** the maximum proliferation rate in case cells.

Prediction calibration

The calibration refers to the correspondence between the probabilistic output of the method and the observed fraction of positive cases.

On ~2,000 newly annotated variants **PhD-SNPg** and **CADD** among the most accurate and calibrated methods with AUC > 0.96 and Brier Score < 0.07. Nevertheless CADD output needs to be transformed to be calibrated.

	BS_{All}	BS_{Coding}	BS_{Noncoding}
PhD-SNPg	0.07	0.10	0.03
CADD*	0.05	0.05	0.04



SNPs&GO prediction

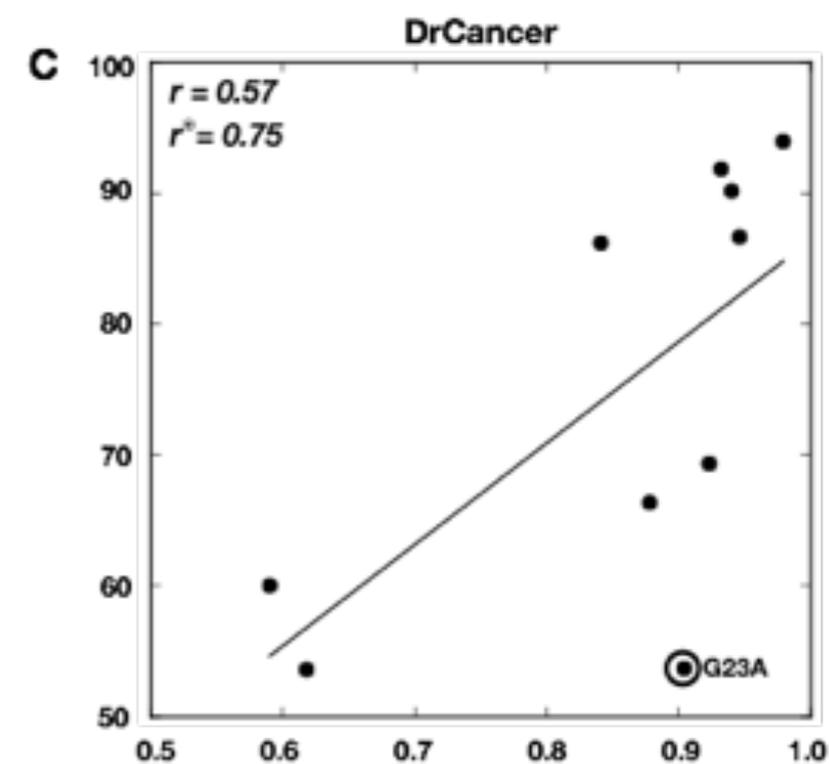
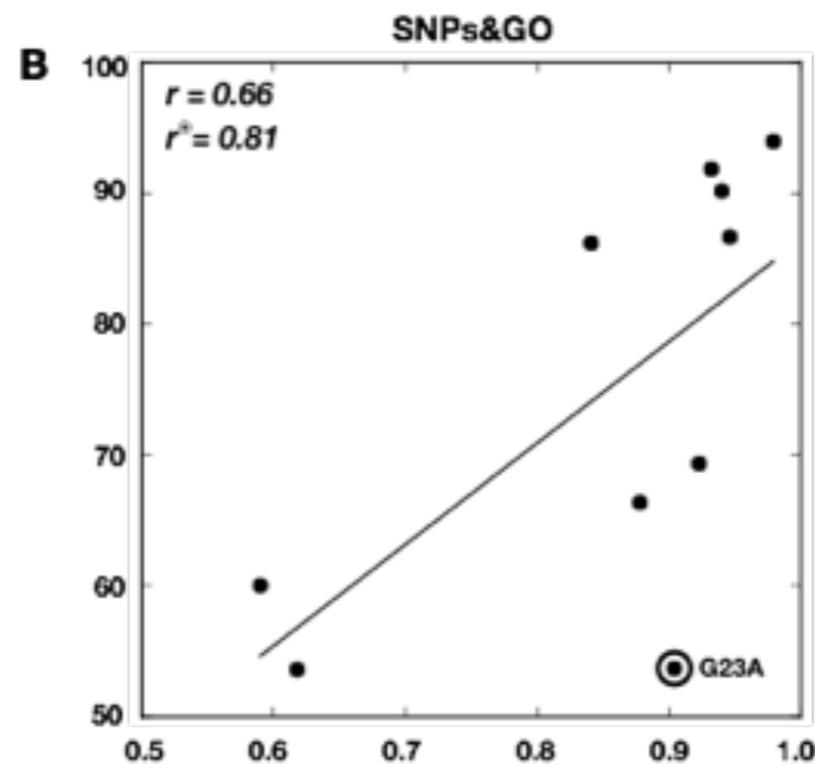
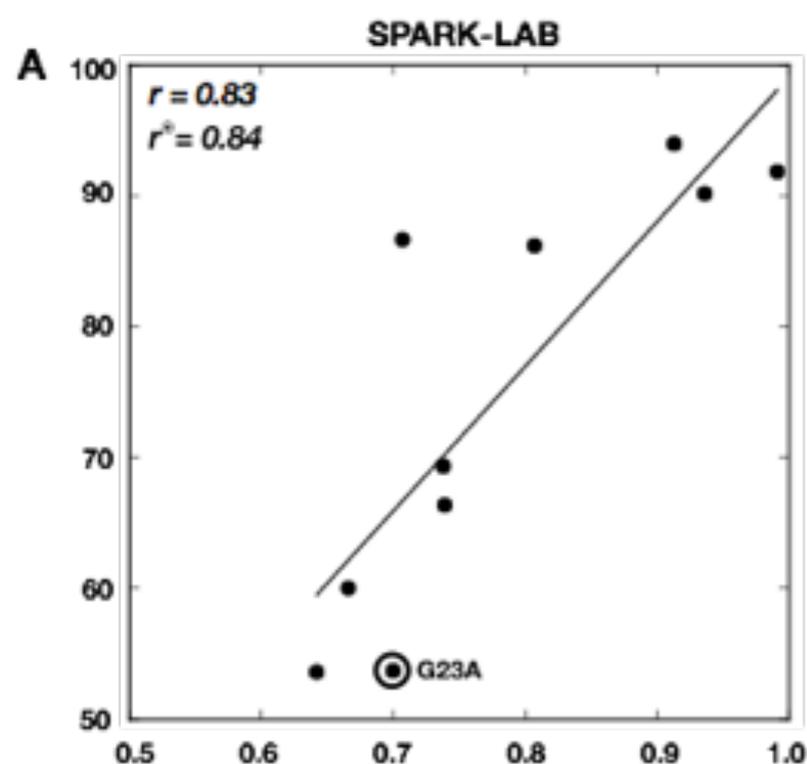
Proliferation rates predicted using the **output of SNPs&GO** without any optimization.

Variant	Prediction	Real	Δ	%WT	%MUT
G23R	0.932	0.918	0.014	84	0
G23S	0.923	0.693	0.230	84	1
G23V	0.940	0.901	0.039	84	0
G23A	0.904	0.537	0.367	84	2
G23C	0.946	0.866	0.080	84	0
G35E	0.590	0.600	0.010	12	14
G35W	0.841	0.862	0.021	12	0
G35R	0.618	0.537	0.081	12	4
L65P	0.878	0.664	0.214	15	1
L94P	0.979	0.939	0.040	56	0

P16 predictions

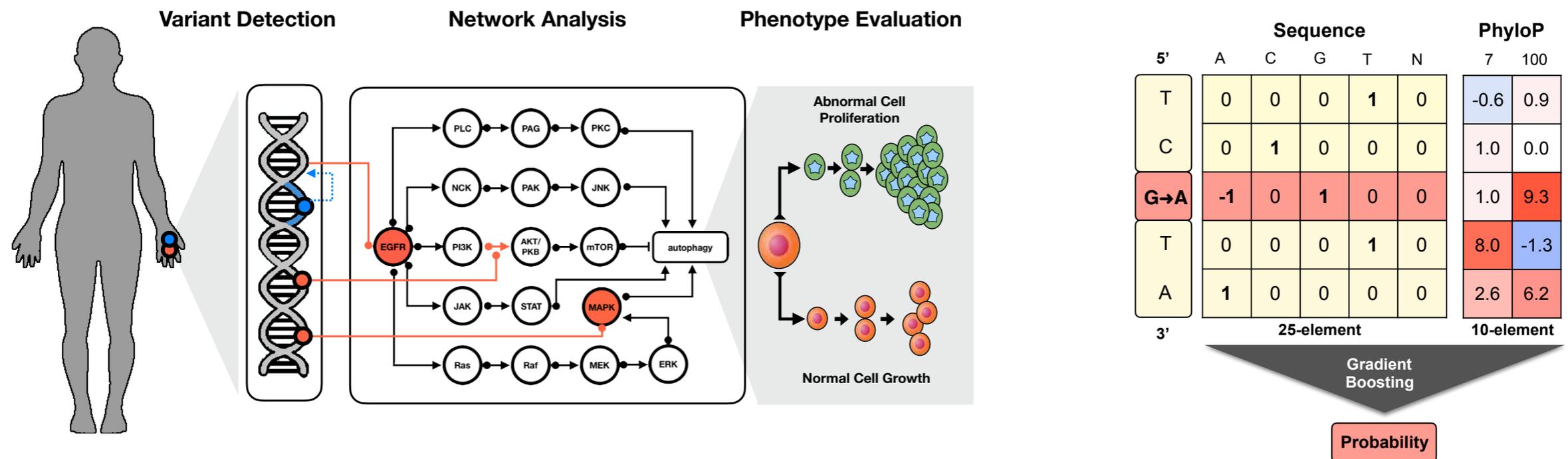
SNPs&GO resulted among the best methods for predicting the impact of P16INK4A variants on cell proliferation.

Method	Q2	AUC	MC	RMSE	rPearson	rSpearman	rKendallTau
SPARK-LAB	0.900	0.920	0.816	0.30	0.595	0.619	0.443
SNPs&GO	0.700	0.880	0.500	0.33	0.575	0.616	0.445
DrCancer	0.600	0.840	0.333	0.46	0.477	0.495	0.409



PhD-SNPg

Variations in regulatory regions can **perturb gene networks** changing the topology or the edge weight of the biological network



PhD-SNPg implements a gradient-boosting algorithms that can run relying only on web resources

<http://snps.biofold.org/phd-snp-g>