

## Supporting Information:

### Insight into neutral and disease-associated human genetic variants through interpretable predictors

Bastiaan A. van den Berg<sup>1,6,7</sup>, Marcel J.T. Reinders<sup>1,6,7</sup>, Dick de Ridder<sup>1,5,6,7</sup>, Tjaart A.P. de Beer<sup>2,3,4,\*</sup>

**1** Delft Bioinformatics Lab, Department of Intelligent Systems, Faculty Electrical Engineering, Mathematics and Computer Science, Delft University of Technology, Mekelweg 4, 2628CD, Delft, The Netherlands

**2** European Bioinformatics Institute (EMBL-EBI) European Molecular Biology Laboratory, Wellcome Trust Genome Campus, Hinxton, Cambridge, CB10 1SD, United Kingdom

**3** Biozentrum, University of Basel, Basel 4056, Switzerland

**4** SIB Swiss Institute of Bioinformatics, Basel 4056, Switzerland

**5** Bioinformatics Group, Wageningen University, Droevendaalsesteeg 1, 6708PB, Wageningen, The Netherlands

**6** Netherlands Bioinformatics Centre, Nijmegen, The Netherlands

**7** Kluyver Centre for Genomics of Industrial Fermentation, Delft, The Netherlands

\* E-mail: Corresponding tjaart.debeer@unibas.ch

## Used web resources

Fasta files with human protein sequences were obtained from the Ensemble FTP server<sup>1</sup>, and from the UniProt website<sup>2</sup>. Human variation data was obtained from the SwissProt human single amino acid variations (humsavar) data base<sup>3</sup>. Only the variants annotated as disease-associated were used for our data set. Predictions for the variations in our data set were obtained using the SIFT website<sup>4</sup> and using the PolyPhen2 website<sup>5</sup>.

## Generating MSA's using HHBlits

A local installation of the software tool HHBlits<sup>6,7</sup> was used to generate MSA's for the set of human proteins in File S4. A script was used to run the following command for each protein in this set:

```
hhblits -i A0A183.fsa -o $OUTPUT/A0A183.hhr \
  -oa3m $OUTPUT/A0A183.a3m -n 1 \
  -d $DATABASES/hhsuite/uniprot20_2013_03/uniprot20_2013_03 \
  -cpu 8 -maxfilt 500000 -maxmem 50GB
```

The raw output files (.a3m) were parsed to generate files with the aligned sequences only, these files were used for calculating the MSA-based features. Zip files with the raw output (4.4 GB) and the generated alignment files (2.9 GB) are available on request.

<sup>1</sup>[http://ftp.ensembl.org/pub/release-71/fasta/homo\\_sapiens/pep/Homo\\_sapiens.GRCh37.71.pep.all.fa.gz](http://ftp.ensembl.org/pub/release-71/fasta/homo_sapiens/pep/Homo_sapiens.GRCh37.71.pep.all.fa.gz)

<sup>2</sup>[http://www.uniprot.org/uniprot/?query=organism:9606+AND+keyword:"Complete+proteome+\(KW-0181\)"+reviewed:yes&force=yes&format=fasta](http://www.uniprot.org/uniprot/?query=organism:9606+AND+keyword:"Complete+proteome+(KW-0181)"+reviewed:yes&force=yes&format=fasta)

<sup>3</sup><http://www.uniprot.org/docs/humsavar>

<sup>4</sup>[http://sift.jcvi.org/www/SIFT\\_enst\\_submit.html](http://sift.jcvi.org/www/SIFT_enst_submit.html)

<sup>5</sup><http://genetics.bwh.harvard.edu/pph2/bgi.shtml>

<sup>6</sup><http://toolkit.tuebingen.mpg.de/hhblits>

<sup>7</sup><http://toolkit.genzentrum.lmu.de/pub/HH-suite/>

## Predicting PFAM domains

A local installation of the software tool PfamScan<sup>8</sup> was used to predict PFAM domains for the set of human proteins in File S4, using the following command:

```
./pfam_scan.pl -fasta protein_sequences.fsa \  
-dir $HOME/tools/pfam_scan/hmm \  
-outfile pfam_raw.txt -as -cpu 1
```

The raw output of this software tool is given in File S5. A parsed version of this file was used to calculate the PFAM-based features.

---

<sup>8</sup><http://ftp.ebi.ac.uk/pub/databases/Pfam/Tools/PfamScan.tar.gz>