

1

2 **Supplementary Information for**

3 **Comprehensive characterization of amino acid positions in protein structures reveals** 4 **molecular effect of missense variants**

5 **Sumaiya Iqbal, Eduardo Pérez-Palma, Jakob B. Jespersen, Patrick May, David Hoksza, Henrike O. Heyne, Shehab S. Ahmed,**
6 **Zaara T. Rifat, M. Sohel Rahman, Kasper Lage, Aarno Palotie, Jeffrey R. Cottrell, Florence F. Wagner, Mark J. Daly, Arthur J.**
7 **Campbell and Dennis Lal**

8 **Sumaiya Iqbal, Arthur J. Campbell, and Dennis Lal.**

9 **E-mail: sumaiya@broadinstitute.org, arthurc@broadinstitute.org, and lald@ccf.org**

10 **This PDF file includes:**

- 11 Supplementary text
- 12 Figs. S1 to S11
- 13 Tables S1 to S3
- 14 Legends for Dataset S1 to S5
- 15 SI References

16 **Other supplementary materials for this manuscript include the following:**

- 17 Datasets S1 to S5

18 Supporting Information Text

19 **Disease-Associated Genes with Structure (DAGS1330) and Variant Set Preparation.** The Protein Data Bank (PDB) (1) was
20 searched in January 2018 for all available human (in full or chimeric) protein structures, resulting in 43,805 experimentally
21 solved structures. From UniProt (2) release 2018_02, 5,870 unique protein identifiers and 5,850 unique gene names for these
22 43,805 structures were collected. Mapping of coordinates of amino acid residues in the 3D structure to linear protein sequence
23 was derived from the SIFTS database (3).

24 Protein-coding single nucleotide variants (SNVs) in the general population were retrieved from gnomAD, public release
25 2.0.2 (4), which we refer to as population variant/variation in this study. Missense variant annotations were extracted from
26 gnomAD Variant Call Format (VCFs) (<http://gnomad.broadinstitute.org/downloads>) file using vcftools (5). Entries from the
27 canonical transcript (CSQ canonical = “YES” flag), passing gnomAD standard quality controls (Filter = “PASS” flag) were
28 extracted. The canonical transcript is defined as the longest CCDS translation with no stop codons according to Ensembl (6).

29 We collected patient missense variant data from two sources: the ClinVar database (7), February 2018 release and the
30 Human Gene Mutation Database (HGMD®) Professional release 2017.2 (8). ClinVar variants were downloaded directly from
31 the ftp site (<ftp://ftp.ncbi.nlm.nih.gov/pub/clinvar/>) in table format and were subsequently filtered to keep missense variants with
32 “Pathogenic” and/or “Likely Pathogenic” clinical consequence. Raw HGMD files were filtered to collect missense mutations
33 with high confidence (confidence = “HIGH” flag) and disease-causing state (variantType = “DM” flag, indicating Disease
34 Mutation). All annotations refer to the human reference genome version GRCh37.p13/hg19. The combined set of variants
35 from ClinVar and HGMD databases is hereafter referred to as pathogenic variants and the corresponding single amino-acid
36 substitutions as pathogenic variations.

37 Altogether, we collected 1,485,579, 16,570, and 47,036 amino acid-altering missense variants from gnomAD, ClinVar
38 and HGMD databases, respectively. These variants were from 5,850 human genes encoding proteins with 43,805 structures
39 available in the Protein Data Bank (PDB) (1). We filtered out genes for which (i) canonical isoform protein sequences were not
40 translated from canonical transcripts and (ii) no altered amino acid residue was mappable to a protein structure (*SI Appendix*,
41 Table S1). The resulting dataset included 1,330 genes, for which there was at least one pathogenic and population missense
42 variation that we could map on one or more protein structure. Therefore, we included only these 1,330 genes with 164,915
43 population and 32,923 pathogenic variations mappable on 14,270 human protein structures in our final dataset, referred to
44 as Disease-Associated Genes with Structure (DAGS1330) dataset in this study. The list of genes in DAGS1330 is available
45 in Dataset S1. The missense variations in the DAGS1330 set were used to perform the statistical analyses and the outcome
46 (characteristic features of pathogenic and population missense variation position in 3D) was validated on independent test sets.

47 **Validation Dataset Preparation.** A validation set of pathogenic and benign missense variant was collected from the ClinVar
48 (February 2019 release) and HGMD® (professional release 2019.2) databases. All variants present in the DAGS1330 set were
49 removed. The remaining set included 4,712 ClinVar benign and likely benign variants (hereafter collectively referred to as
50 benign) and 17,983 pathogenic variants (ClinVar pathogenic and likely pathogenic variants and HGMD-reported disease
51 mutations). The validation set contained benign and pathogenic variants from 1,286 genes out of the 1,330 genes in the
52 DAGS1330 set. Further, high-throughput mutagenesis readouts, classifying loss-of-function or damaging variations from neutral
53 ones in BRCA1 and PTEN were collected from literature (9, 10)

54 **Feature set mining and annotation.** We annotated the amino acid residues positions in proteins encoded by 1,330 genes with
55 forty (a combination of structural, physicochemical and functional) features from seven main groups. In the following, we
56 introduce the features, and describe their curation and annotation process.

- 57 1. Three-class secondary structures (feature count: 3). Protein secondary structure (SS) is the three-dimensional (3D)
58 form of local segments of proteins that are defined by the hydrogen bonds between the side chain atoms of amino acids
59 in the peptide backbone. We obtained precomputed SS assignment values for the protein structures from the DSSP
60 repository (11). Based on the hydrogen bond energy, eight different types (see below) of SS were assigned to the amino
61 acid residues in the protein structure. These eight different types of SSs can be grouped into three broader categories: (i)
62 β -sheet/strand: β -strand and extended β -sheet; (ii) Helices – 3_{10} -helix, α -helix, π -helix; (iii) Coils – turn, bend, and
63 random loops. For proteins with multiple available structures, we derived a consensus annotation which represents the
64 maximum agreement of SS type assignments out of all structures.
- 65 2. Eight-class secondary structures (feature count: 8). We have also looked into more detailed types of SSs to check for
66 associations with pathogenic or population variants. The eight structure types collected from DSSP are: (i) β -strand –
67 residues in isolated β -bridge; (ii) β -sheet – residues in extended strand that participates in β -ladder; (iii) 3_{10} -helix –
68 hydrogen bond between i^{th} and $(i+3)^{th}$ residues to build each helical turn; (iv) α -helix – hydrogen bond between i^{th}
69 and $(i+4)^{th}$ residues to build each helical turn; (v) π -helix – hydrogen bond between i^{th} and $(i+5)^{th}$ residues to build
70 each helical turn; (vi) turn – hydrogen bonded turn; (vii) bend – residues of high curvature where the angle between
71 $C_i C_{i+2}$ and $C_{i-2} C_i$ is at least 70° ; (viii) loop – random loop where no other rule applies. Similar to the three-class SS
72 annotation, for proteins with multiple structures, we considered the maximum agreement of SS type assignments from all
73 the available structures as the final annotation.
- 74 3. Residue exposure level (feature count: 5). The solvent accessible surface area (ASA) of each amino acid residue in the
75 protein structures was also collected from DSSP repository (11). The ASA is measured in units of square of Angstroms

76 (\AA^2). The Relative solvent accessible area (RSA) of a residue “X” was calculated by normalizing the ASA of that residue
77 by the ASA of the same residue in a reference tripeptide state (Gly-X-Gly), collected from (12). When multiple structures
78 were available for an amino acid, we considered the median of all RSA values for that residues. Based on the value of
79 RSA, we labeled each amino acid with one of the following five types: (i) core (RSA < 5%); (ii) buried (5% ≤ RSA <
80 25%); (iii) medium-buried (25% ≤ RSA < 50%); (iv) medium-exposed (50% ≤ RSA < 75%); (v) exposed (RSA ≥ 75%).

81 4. Physicochemical properties of amino acids (feature count: 8). While every amino acid have an amino and a carboxyl
82 group, their side chains are all different and have distinct physicochemical properties that influence protein conformation
83 and function. Based on the property of the side chain, the 20 natural amino acids are divided into eight groups: (i)
84 aliphatic – alanine (Ala/A), isoleucine (Ile/I), leucine (Leu/L), methionine (Met/M), valine (Val/V); (ii) aromatic
85 – phenylalanine (Phe/F), tryptophan (Trp/W), tyrosine (Tyr/Y); (iii) hydrophobic – aliphatic and aromatic amino
86 acids; (iv) positively-charged – arginine (Arg/R), histidine (His/H), lysine (Lys/K); (v) negatively-charged – aspartic
87 acid (Asp/D), glutamic acid (Glu/E); (vi) neutral – asparagine (Asn/N), glutamine (Gln/Q), serine (Ser/S), threonine
88 (Thr/T); (vii) polar – positively-charged, negatively-charged, and neutral amino acids; (viii) “special” – proline (Pro/P):
89 has a cyclic side chain and cannot make backbone hydrogen bonds; glycine (Gly/G): does not have a side chain and
90 allows flexibility, cysteine (Cys/C): the side chain has a reactive sulfhydryl group. In our analysis, we considered only the
91 amino acid residues that have defined coordinates in protein 3D structures.

92 5. Protein-protein interactions (feature count: 4). The PDBsum database (13) was curated to annotate amino acid residues
93 with residue-residue interaction types between amino acid pairs at distance ≤ 4 Å on the protein-protein interface of
94 complex structures. An amino acid can possibly be involved in different bonds in multiple protein complexes. Here,
95 we kept records of all the interactions that one residue was involved in different conformations. The four types of
96 interactions/bonds that we considered are: (i) disulfide bond – a covalent bond between cysteine side chains, yielding a
97 disulfide bridge, which is an key determinant of protein 3D structure; (ii) salt-bridge interaction – an ionic bond between
98 oppositely charged residues; (iii) hydrogen bond – electrostatic interaction between two atoms bearing partial negative
99 charges, that share a partially positively charged hydrogen; (iv) Van der Waals interaction – a bond formed by the
100 transient and weak Van der Waals attraction between two close atoms.

101 6. Post-translational modifications (feature count: 6). Post-translational modification (PTM) refers to the covalent and
102 enzymatically-mediated modification of proteins to form the mature protein. We collected amino acid positions of
103 six different PTM types from the PhosphoSitePlus database (14): (i) acetylation – introduces an acetyl group; (ii)
104 methylation – addition of a methyl group; (iii) O.GlcNAc – also known as O-linked N-acetylglucosamine, is a form of
105 protein glycosylation; (iv) phosphorylation – attachment of a phosphoryl group; (v) SUMOylation – addition of SUMO
106 (small ubiquitin-like modifiers) molecule; (vi) ubiquitination – attachment of ubiquitin. Upon mapping the PTM types
107 onto the protein structure, we computed the spatial distances between the C_α atoms of each amino acid and different
108 PTM sites in the structure. We annotated each amino acid with the nearest distance to the six PTM types, where
109 available. For our analysis, we considered the amino acids that are spatially close to a PTM site, i.e. if the residue is
110 located within < 10 Å distance from the PTM site in the structure.

111 7. Functional features (feature count: 6). UniProt database (2) was mined to collect the sequence annotations, describing
112 regions or sites of interest in protein in terms of 25 features. Genomic annotations were download as bed files from
113 the UniProt webpage (ftp://ftp.uniprot.org/pub/databases/uniprot/current_release/knowledgebase). All hg38 coordinates
114 were converted to hg19 coordinates using the UCSC liftover tool (<https://genome.ucsc.edu/cgi-bin/hgLiftOver>). Only
115 correctly lifted-over regions were then reformatted to ANNOVAR (15) region format files for hg19 ([http://annovar.
116 openbioinformatics.org/en/latest/user-guide/input/](http://annovar.openbioinformatics.org/en/latest/user-guide/input/)). Pathogenic and population VCF files were annotated using ANNOVAR
117 with the different UniProt features. Guided by EMBL-EBI online material ([https://www.ebi.ac.uk/training/online/course/
118 UniProt-exploring-protein-sequence-and-functional-information/sequence-features](https://www.ebi.ac.uk/training/online/course/UniProt-exploring-protein-sequence-and-functional-information/sequence-features)) and similarity, these 25 features were
119 coarse grained into six categories to increase the power of the statistical association tests performed in this study. These six
120 feature categories are (i) functional site – active site, metal binding site, binding site, site; (ii) functional/binding region –
121 zinc finger, DNA binding region, nucleotide phosphate binding region, calcium binding region; (iii) sequence motif/region
122 – region, repeat, coiled coil, motif; (iv) modular domain – domain, topological domain, transmembrane, intramembrane;
123 (v) molecular processing-related region – peptide, transit peptide, signal peptide, propeptide; (vi) modified residue –
124 modified residue, lipidation, disulfide bond, cross-link, glycosylation. If an amino acid residue is part of any subtype of a
125 category, that residue is annotated by the corresponding category.

126 **Protein Class Annotation.** To identify the properties of pathogenic and population missense variations on structure at different
127 protein function level, we grouped the 1,330 genes into twenty-four major protein classes (main text, Fig. 1, Step 3) using
128 the PANTHER database (Protein ANalysis THrough Evolutionary Relationships (16), 13.1 release, published in February
129 2018). After automatic annotation, 624 genes were not assigned to any specific protein class. We then collected the protein
130 class information for these 624 genes from (i) the Ensemble family description (version 93) using BioMart (17) and (ii)
131 molecular function and/or biological process annotation available in UniProt (release 2018_02) as defined by the Gene Ontology
132 consortium (18). The list of genes, number of protein 3D structures and missense variations mapped on structures for the
133 protein classes are given in Dataset S2. Note that a protein may have multiple functions and so can be assigned into multiple
134 protein classes.

135 **Statistical Analysis.** The two-sided Fisher's exact test of association was performed for each of the forty 3D features, taking the
136 counts of pathogenic and population missense variations with and without a feature, to quantify the burden of pathogenic or
137 population variations for each feature (main text, Fig. 1, Step 4). An estimate of enrichment or burden (odds ratio, OR), 95%
138 confidence interval (CI) of the OR value, and the p -value (p) showing the significance of the observed burden or association,
139 were obtained from the test output. The described analysis was performed on the full DAGS1330 set and separately on subsets
140 of variations in genes grouped in twenty-four protein functional classes. Thus, we performed a total of 40 (number of features)
141 \times 25 (number of datasets) = 1000 tests. Subsequently, the p -values were corrected to generate a corrected value of p that we
142 called " q ", calculated as $p \times 1000$ according to the Bonferroni correction for multiple testing in statistical analysis. Therefore,
143 a 3D feature is considered to be a characteristic feature of pathogenic variants when the test outputs $OR > 1$ and $q < 0.05$. In
144 contrast, when the test outputs $OR < 1$ and $q < 0.05$, the feature is referred to as a characteristic 3D feature of population
145 variants. This approach for the characterization of variants by comparative enrichment analysis taking both pathogenic and
146 population variants into account in the two-sided Fisher's exact test, reduces the possibility of obtaining a result simply because
147 of the abundance of a certain 3D feature in a protein class. For example, the fraction of residues located in the protein core
148 ("core" residues, see "Residue exposure level" ranking in *Materials and Methods*) that are mutated in both pathogenic and
149 population variants (39% vs. 23%) were used to derive the enrichment of pathogenic variations in the protein core (>2 -fold, Fig.
150 2). Thanks to our comparative method we can be reasonably confident that if we see an enrichment of pathogenic variations in
151 protein core positions, it is not only because our sample contains a high proportion of compact structures where the majority
152 of residues are buried.

153 For further verification of our protein class-specific results by OR enrichment analysis (main text, Fig. 3A), we computed
154 the "relative risk (RR)" (19) of a mutation to be pathogenic given the altered residue has a 3D feature (for all forty features)
155 across the full dataset (DAGS1330 set) and for individual protein classes: Notably, the RR values were strongly correlated
156 (Pearson $r^2 = 94\%$) with the odds ratios (ORs, main text, Fig.3A), indicating that the ORs effectively approximate the RRs
157 for our study (20, 21).

158 **Computation of Pathogenic 3D Feature Index (P3DFi) per Amino Acid.** For each amino acid residue of the proteins encoded by
159 the 1,330 disease-associated genes, we generated the 3D feature annotations (available in MISCAST) and counted the number
160 of pathogenic and population variant-associated 3D features of the amino acid, denoted as $3DF^{PATH}$ and $3DF^{POP}$, respectively.
161 Thereafter, the pathogenic 3D feature index (P3DFi) per amino acid is computed as $3DF^{PATH}$ minus $3DF^{POP}$ (P3DFi > 0 thus
162 indicates a 3D mutational hotspot). Note that we identified the pathogenic and population variant-associated 3D features both
163 for all 1,330 genes analysed together as one pool (main text, Fig. 2) and also for twenty-four different protein classes, grouping
164 genes encoding for similar functions (main text, Fig. 3A). Therefore, P3DFi can be derived using the full DAGS1330-based
165 3D features ($P3DFi_{DAGS1330}$) and also using protein class-specific ($P3DFi_{Protein\ class}$) 3D features. When a protein belongs to
166 more than one protein class, a single $P3DFi_{Protein\ class}$ per amino acid was decided by taking the majority agreement from all
167 $P3DFi_{Protein\ class}$ and $P3DFi_{DAGS1330}$ values. The $P3DFi_{DAGS1330}$ and all $P3DFi_{Protein\ class}$ values per amino acid of the 1,330 proteins
168 are made available through the MISCAST web server (<http://miscast.broadinstitute.org/>)

169 **Development of Ensemble Model.** Annotations of variants in the DAGS1330 set (used as the training set to develop the model)
170 and the validation set with predicted pathogenicity scores from SIFT (22) (SIFT_score), PolyPhen2 (23) (Polyphen2_HVAR_score)
171 and CADD (24) (CADD_phred) were generated using ANNOVAR (15). The categorical prediction (1: deleterious/pathogenic,
172 -1: neutral/benign) was determined using the guidelines described in the Annovar documentation (15). Two ensemble models
173 were developed using the scores from the three existing methods, and separately with $P3DFi_{DAGS1330}$ and $P3DFi_{Protein\ class}$ values.
174 An additional ensemble model was trained only with the scores from existing methods (i.e. without any P3DFi). All the models
175 were developed using the classical random forest method using: Number of estimators or decision tree classifiers = 2,000,
176 quality measure as "gini", and the maximum depth of the trees = 10. To eliminate the imbalanced proportion of pathogenic
177 and benign variants from the training dataset (ratio of pathogenic to benign variants > 5), and ensure robust training of the
178 models, a random oversampling of pathogenic variants was performed. The numerical prediction outputs generated by the
179 models were categorized by applying the default threshold of 0.5 (>0.5 : deleterious/pathogenic, ≤ 0.5 : neutral/benign). Both
180 the training and test datasets, along with the predictions generated by the ensemble models, are available in Dataset S4 and S5.

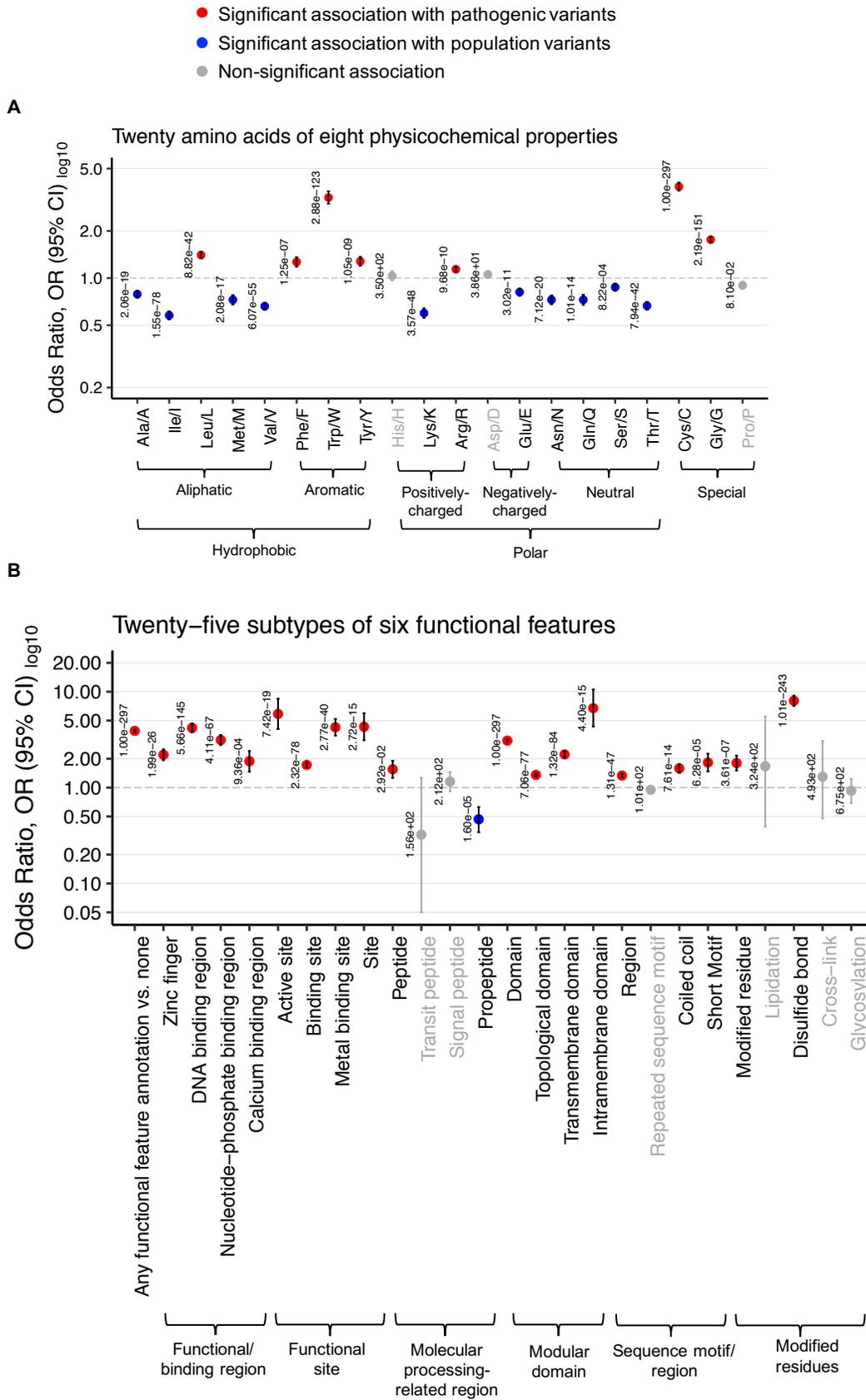


Fig. S1

Fig. S1. (Previous page.) Association of pathogenic and population missense variants from 1,330 disease-associated genes (DAGS1330 set) with (A) twenty natural amino acids (grouped into eight physicochemical properties and analyzed in the main text, **Fig. 2**), (B) 3D sites with twenty-five functional features collected from UniProt (coarse grained into six features and analyzed in the main text, **Fig. 2**). The plots show the results of two-sided Fisher's exact tests of association between 32,923 pathogenic and 164,915 population amino acid variations with the features. Circles show the odds ratio (OR) and are labelled with the corrected p or " q " values, showing the significance of the association (a value of $1.0e-297$ is to be read as $<1.0e-297$, indicating maximum significance), and the horizontal bars show the 95% confidence interval (CI). The $OR > 1$ and $OR < 1$, along with $q < 0.05$, indicate that the corresponding feature (y-axis) has an enrichment of pathogenic (red circle) and population (blue circle) variants, respectively. The $OR = 1$ (dashed vertical line) indicates that there is no association between a variant type (pathogenic or population) and a feature. To facilitate the visualization, minimum and maximum values of OR along x-axis are set to 0.05 and 20.0, respectively. For non-significant association ($q \geq 0.05$), the circle, CI bar and feature names are grey. Out of the three amino acids that are categorized to have "special" physicochemical properties in this study, Cys and Gly had significant enrichment of pathogenic variations (in A). The Cys enrichment is linked to the strong association between disulfide bond and pathogenic variants (main text, **Fig. 2**). Glycine is the smallest amino acid and is predominantly found in flexible parts of the protein structure (more than 67% of the Gly reference amino acid residues in our dataset form pliable coils). Association between the substitution of Gly amino acid and pathogenic variants (2-fold) can be due to the steric clashes introduced by the substitution, causing destabilization of the protein structure. While all the aromatic amino acids were found enriched for pathogenic mutations, all the aliphatic amino acids, only except Leu, were found tolerant to substitution. Out of the positively-charged amino acids, Arginine (Arg/R, $OR = 1.1$, $q < 9.7e-10$) alone had a modest but significant burden of pathogenic mutations (in A). The supplemental analysis of the twenty-five individual features collected from UniProt (in B) allowed us to quantify the burden of pathogenic variants for non-membrane topological domains ($OR = 1.4$, $q = 7.1e-77$), transmembrane ($OR = 2.2$, $q = 1.3e-84$), and intramembrane domains ($OR = 6.7$, $q = 4.4e-15$) of membrane-spanning proteins. In addition to recapitulating features like short biologically-relevant motifs, DNA binding sites, and zinc fingers, that were known to be frequently impaired in pathogenic variants (25, 26), our study revealed novel associations as well. Some notable examples of features that were not previously linked with the germline, primarily non-cancer (91% in our dataset, **Fig. S10**) pathogenic variants are active sites ($OR = 5.9$, $q = 7.4e-19$), metal-binding sites ($OR = 4.3$, $q = 2.8e-40$), coiled-coils ($OR = 1.6$, $q = 7.6e-14$) and short, biologically-active peptides ($OR = 1.6$, $q = 2.9e-02$), highlighting the importance of these 3D sites for the function of a protein and their putative intolerance to perturbation by missense variation.

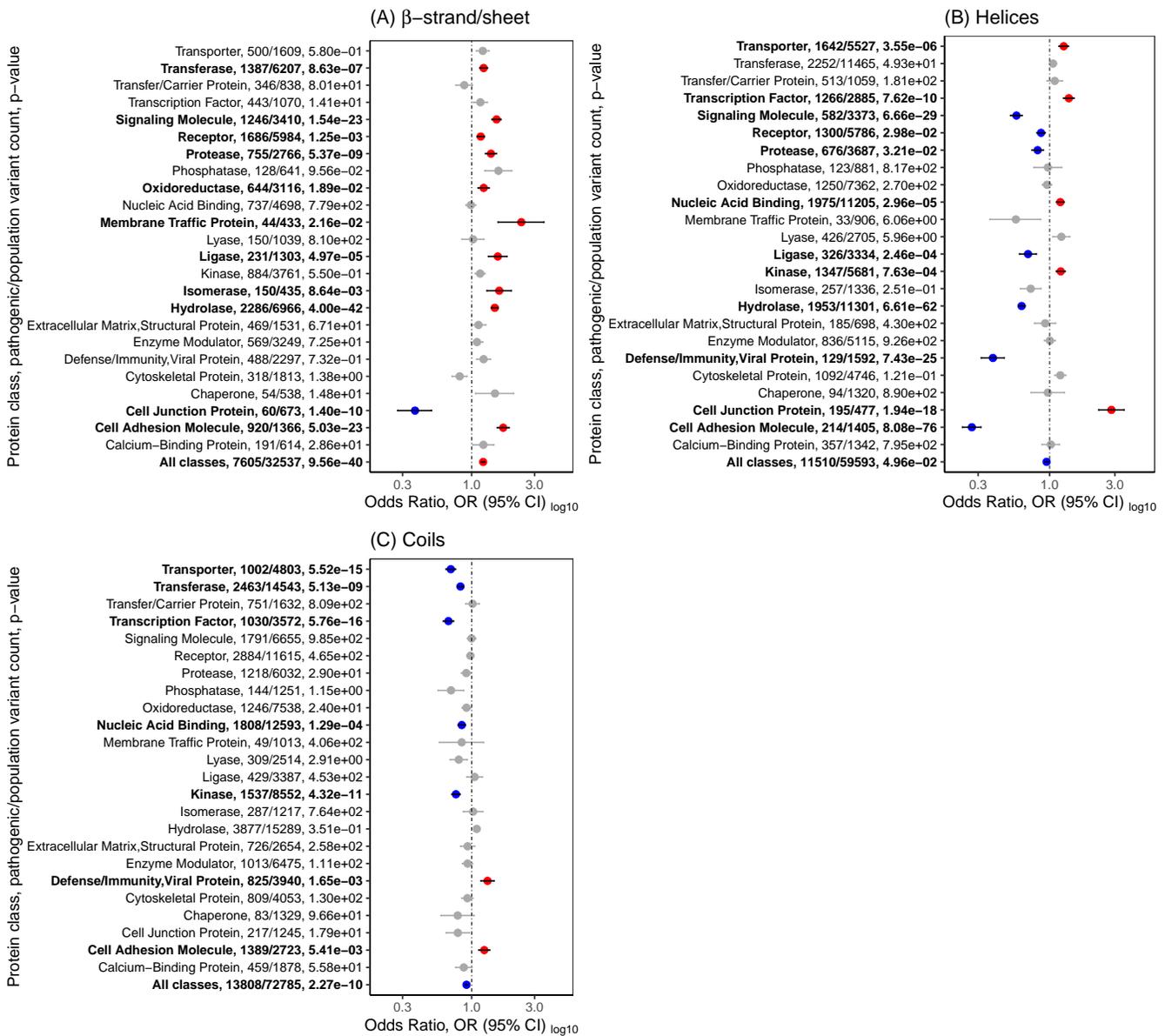


Fig. S2. Enrichment of pathogenic and population missense variants for “All classes” and twenty-four protein functional-classes in the three-class secondary structure types. The plot shows the results of two-sided Fisher’s exact tests of association. Circles show the odds ratio (OR) and the horizontal bars show the 95% confidence interval (CI). The y-axis is labelled with the protein class name, counts of pathogenic and population variants with the corresponding feature, and corrected p or “ q ” values, showing the significance of the association (a value of 1.0×10^{-297} should be read as $<1.0 \times 10^{-297}$, indicating maximum significance). The $OR > 1$ and $OR < 1$, along with $q < 0.05$, indicate that the corresponding feature (plot label) has an enrichment of pathogenic (red circle) and population (blue circle) variants, for the respective protein class (y-axis). The $OR = 1$ (dashed vertical line) indicates that there is no association between a variant type (pathogenic or population) and a feature. To facilitate the visualization, minimum and maximum values of OR along x-axis are set to 0.05 and 20.0, respectively. For non-significant association ($q \geq 0.05$), the circle and CI bar are grey.

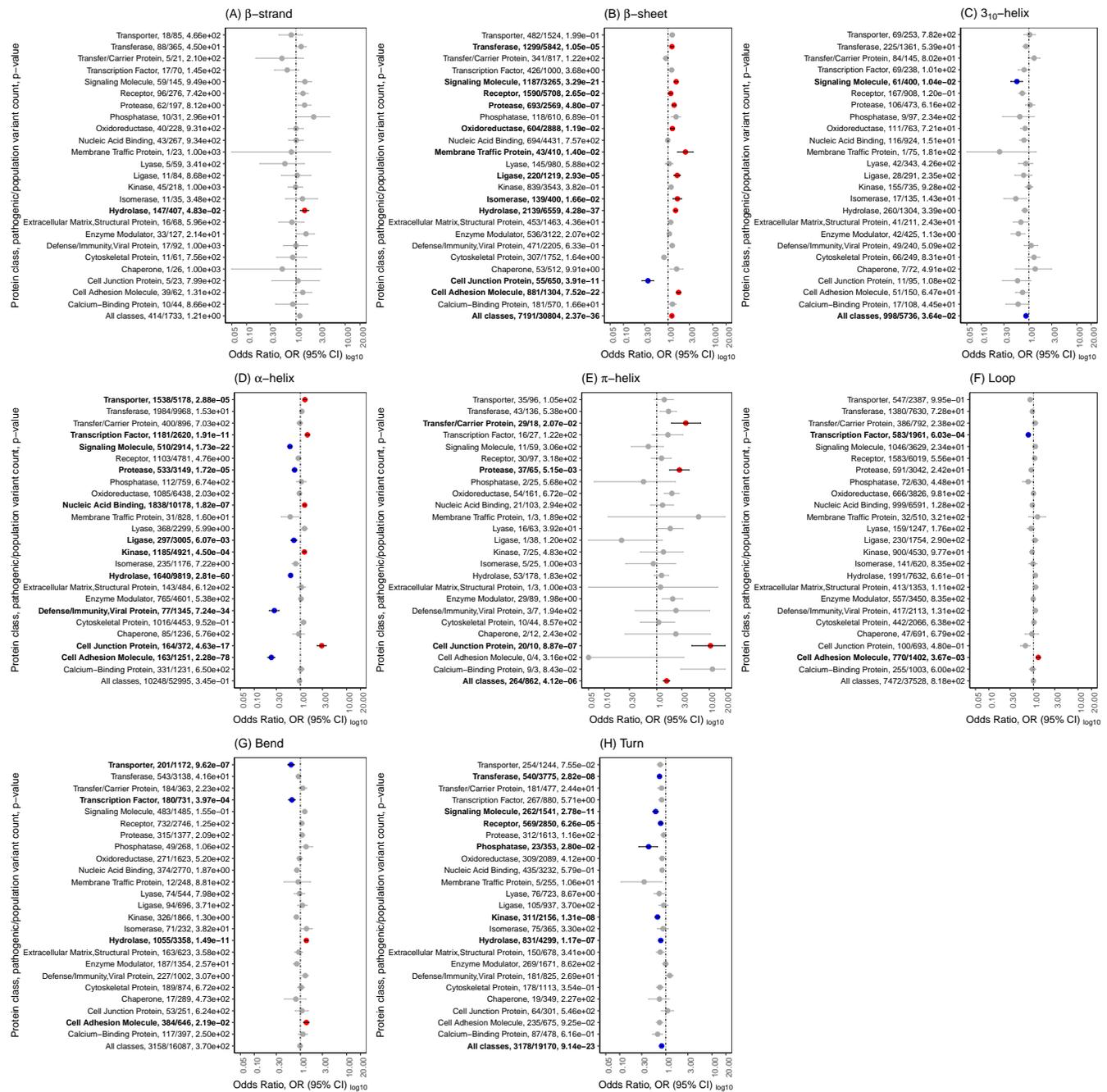


Fig. S3. Enrichment of pathogenic and population missense variants for “All classes” and twenty-four protein functional-classes in the eight-class secondary structure types. The plot shows the results of two-sided Fisher’s exact tests of association. Circles show the odds ratio (OR) and the horizontal bars show the 95% confidence interval (CI). The y-axis is labelled with the protein class name, counts of pathogenic and population variants with the corresponding feature, and corrected p or “ q ” values, showing the significance of the association (a value of $1.0e-297$ is should be read as $<1.0e-297$, indicating maximum significance). The $OR > 1$ and $OR < 1$, along with $q < 0.05$, indicate that the corresponding feature (plot label) has an enrichment of pathogenic (red circle) and population (blue circle) variants for the respective protein class (y-axis). The $OR = 1$ (dashed vertical line) indicates that there is no association between a variant type (pathogenic or population) and a feature. To facilitate the visualization, minimum and maximum values of OR along x-axis are set to 0.05 and 20.0, respectively. For non-significant association ($q \geq 0.05$), the circle and CI bar are grey.

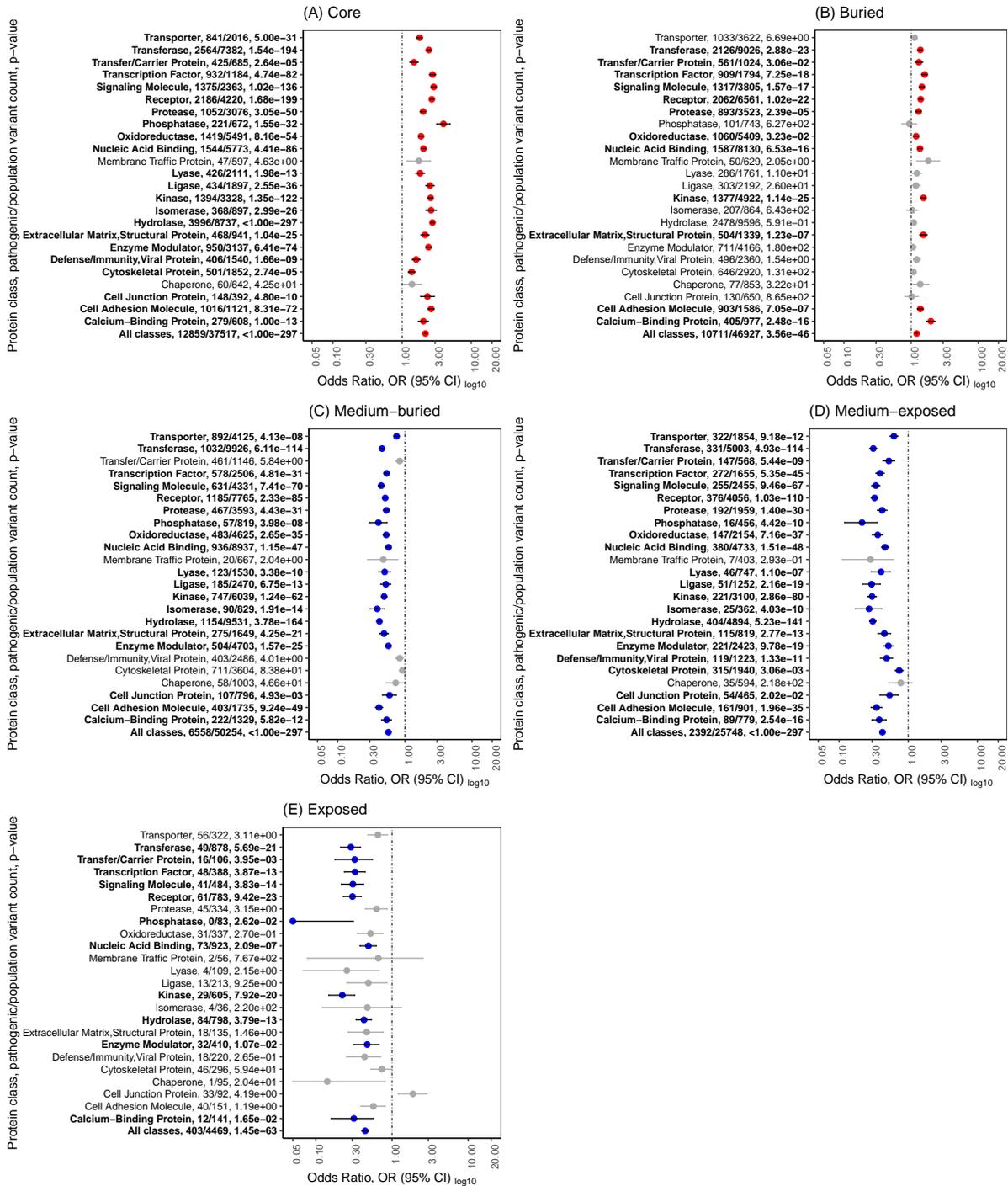


Fig. S4. Enrichment of pathogenic and population missense variants for “All classes” and twenty-four protein functional-classes in amino acid residues with different exposure levels. The plot shows the results of two-sided Fisher’s exact tests of association. Circles show the odds ratio (OR) and the horizontal bars show the 95% confidence interval (CI). The y-axis ticks are labelled with the protein class name, counts of pathogenic and population variants with the corresponding feature, and corrected p or “ q ” values, showing the significance of the association (a value of $1.0e-297$ is should be read as $<1.0e-297$, indicating maximum significance). The $OR > 1$ and $OR < 1$, along with $q < 0.05$, indicate that the corresponding feature (plot label) has an enrichment of pathogenic (red circle) and population (blue circle) variants, for the respective protein class (y-axis). The $OR = 1$ (dashed vertical line) indicates that there is no association between a variant type (pathogenic or population) and a feature. To facilitate the visualization, minimum and maximum values of OR along x-axis are set to 0.05 and 20.0, respectively. For non-significant association ($q \geq 0.05$), the circle and CI bar are grey.

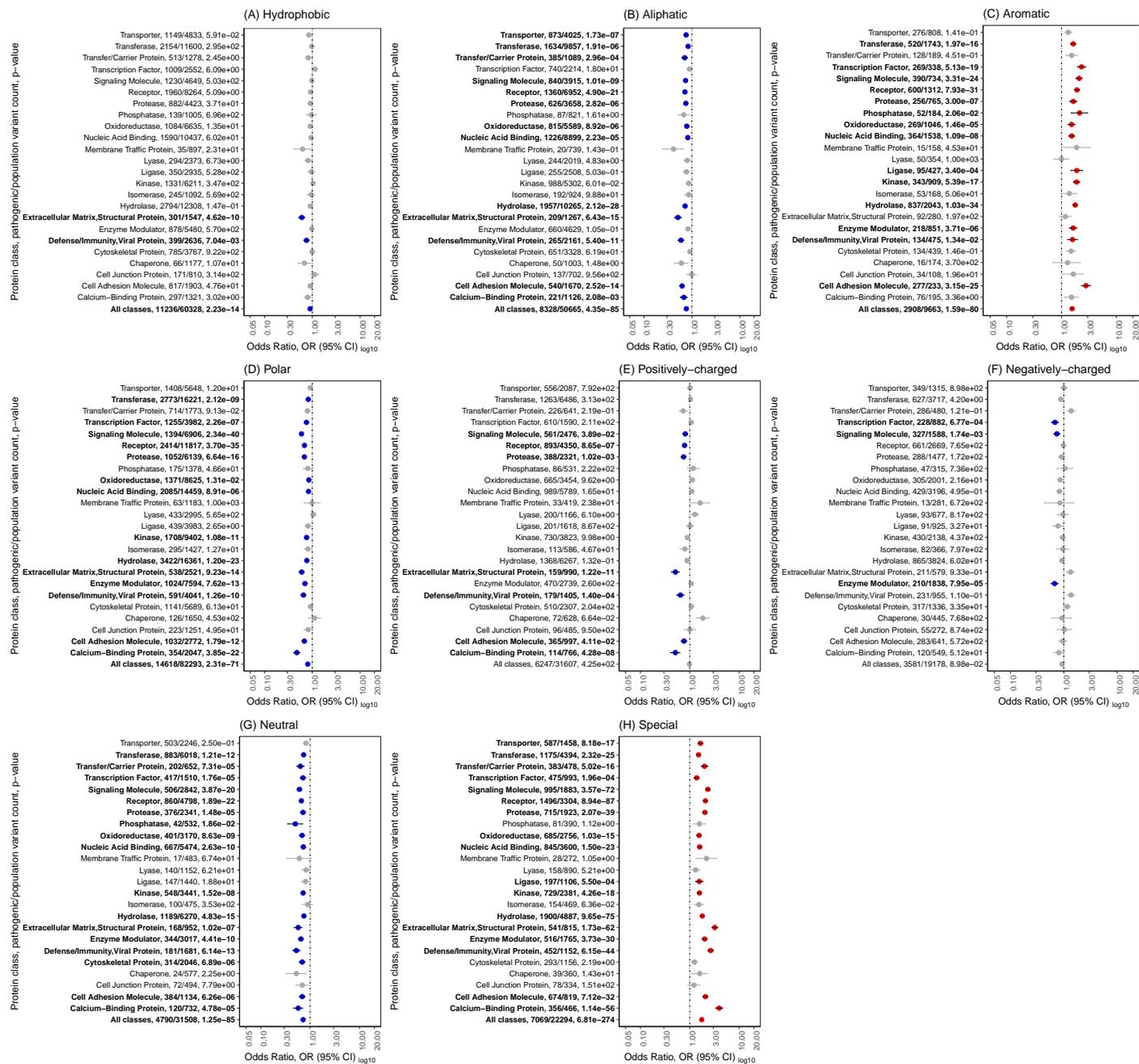


Fig. S5. Enrichment of pathogenic and population missense variants for “All classes” and twenty-four protein functional-classes in amino acids of eight different physicochemical properties. The plot shows the results of two-sided Fisher’s exact tests of association. Circles show the odds ratio (OR) and the horizontal bars show the 95% confidence interval (CI). The y-axis ticks are labelled with the protein class name, counts of pathogenic and population variants with the corresponding feature, and corrected p or “ q ” values, showing the significance of the association (a value of $1.0e-297$ is should be read as $<1.0e-297$, indicating maximum significance). The $OR > 1$ and $OR < 1$, along with $q < 0.05$, indicate that the corresponding feature (plot label) has an enrichment of pathogenic (red circle) and population (blue circle) variants, for the respective protein class (y-axis). The $OR = 1$ (dashed vertical line) indicates that there is no association between a variant type (pathogenic or population) and a feature. To facilitate the visualization, minimum and maximum values of OR along x-axis are set to 0.05 and 20.0, respectively. For non-significant association ($q \geq 0.05$), the circle and CI bar are grey.

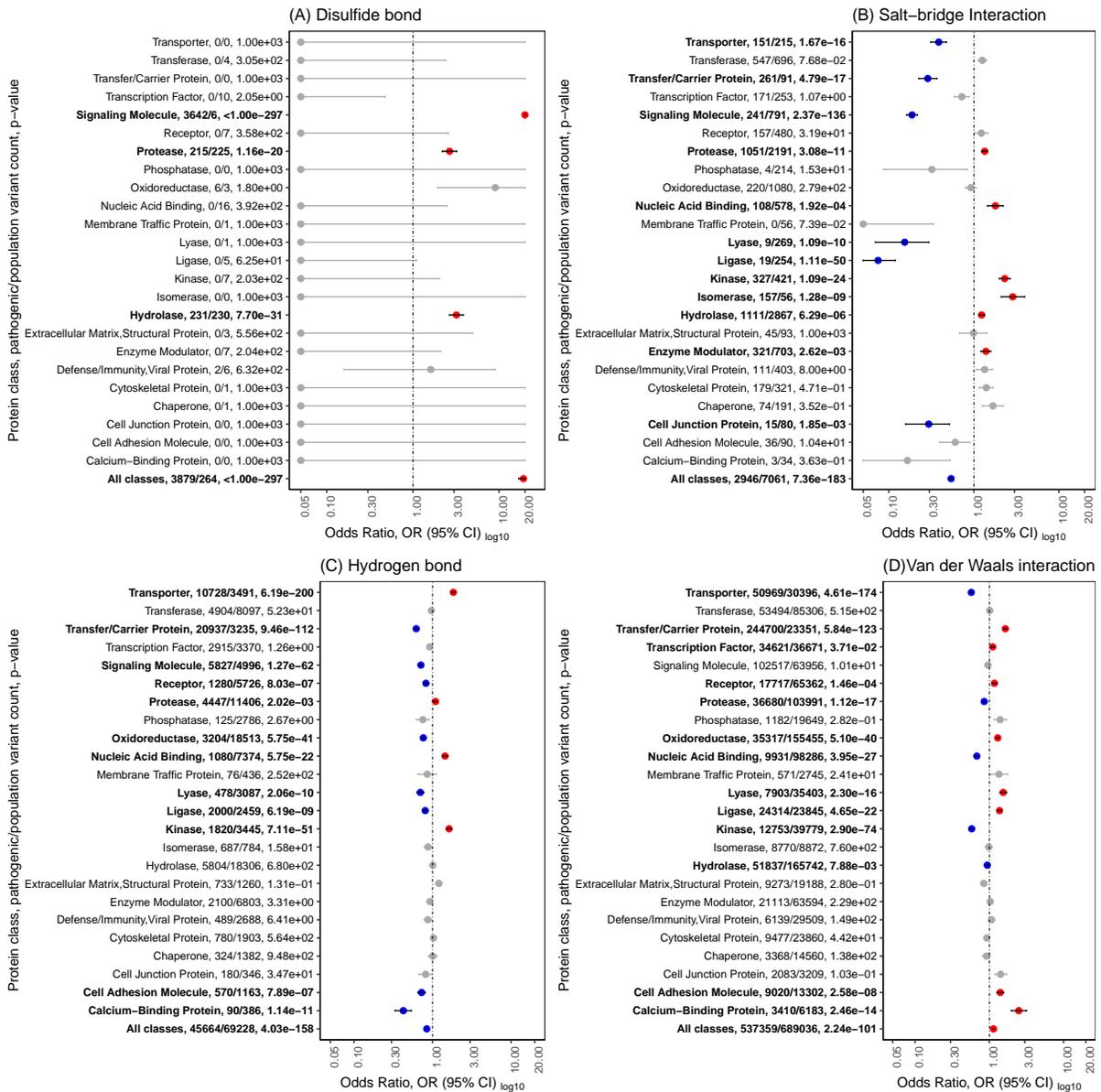


Fig. S6. Enrichment of pathogenic and population missense variants for “All classes” and twenty-four protein functional-classes in protein-protein interaction types. The plot shows the results of two-sided Fisher’s exact tests of association. Circles show the odds ratio (OR) and the horizontal bars show the 95% confidence interval (CI). The y-axis ticks are labelled with the protein class name, counts of pathogenic and population variants with the corresponding feature, and corrected *p* or “*q*” values, showing the significance of the association (a value of 1.0e-297 is should be read as <1.0e-297, indicating maximum significance). The OR > 1 and OR < 1, along with *q* < 0.05, indicate that the corresponding feature (plot label) has an enrichment of pathogenic (red circle) and population (blue circle) variants, for the respective protein class (y-axis). The OR = 1 (dashed vertical line) indicates that there is no association between a variant type (pathogenic or population) and a feature. To facilitate the visualization, minimum and maximum values of OR along x-axis are set to 0.05 and 20.0, respectively. For non-significant association (*q* ≥ 0.05), the circle and CI bar are grey.

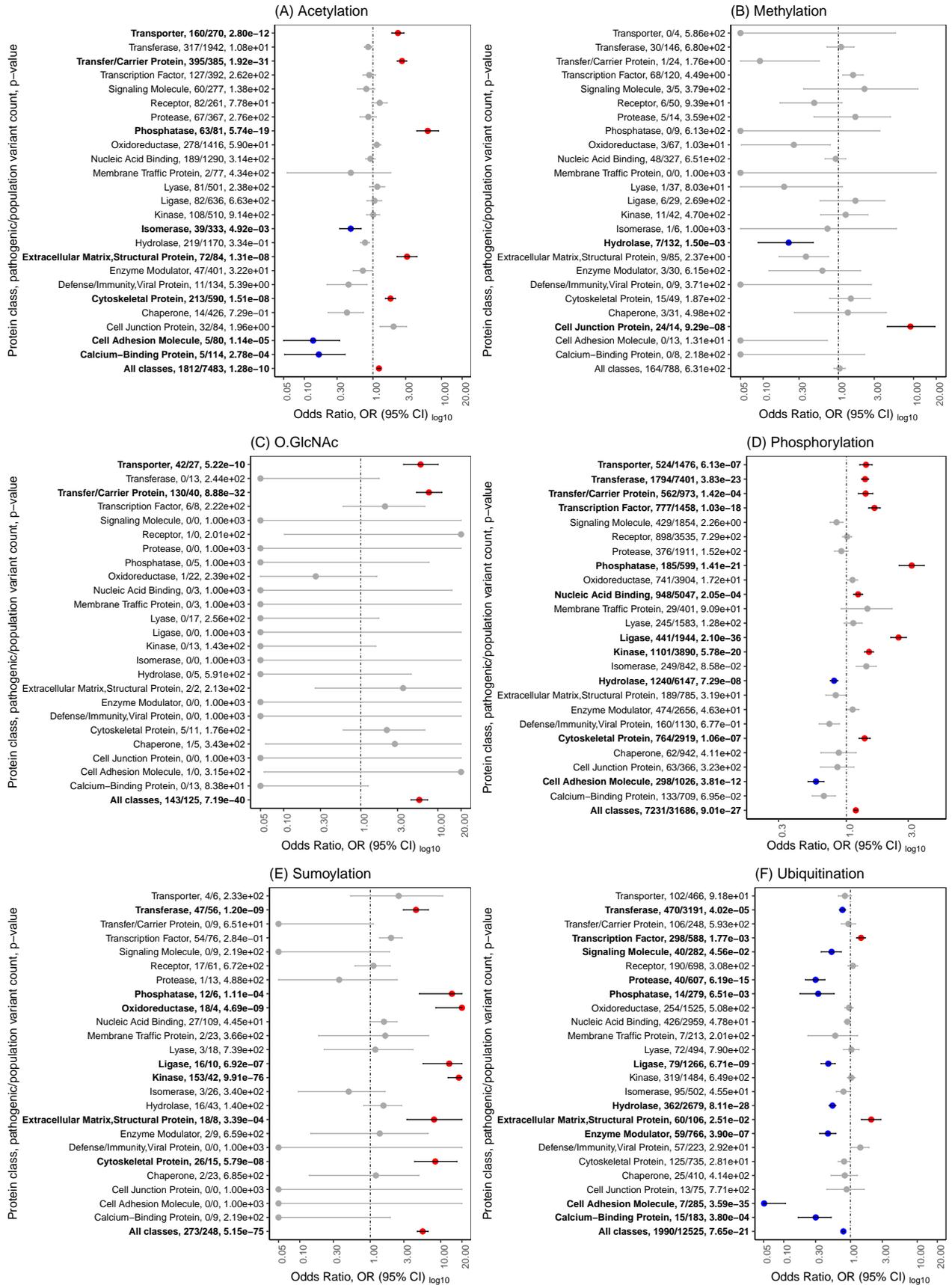


Fig. S7

Fig. S7. (Previous page.) Enrichment of pathogenic and population missense variants for “All classes” and twenty-four protein functional-classes in post-translational modification types. The plot shows the results of two-sided Fisher’s exact tests of association. Circles show the odds ratio (OR) and the horizontal bars show the 95% confidence interval (CI). The y-axis ticks are labelled with the protein class name, counts of pathogenic and population variants with the corresponding feature, and corrected p or “ q ” values, showing the significance of the association (a value of $1.0e-297$ is should be read as $<1.0e-297$, indicating maximum significance). The $OR > 1$ and $OR < 1$, along with $q < 0.05$, indicate that the corresponding feature (plot label) has an enrichment of pathogenic (red circle) and population (blue circle) variants, for the respective protein class (y-axis). The $OR = 1$ (dashed vertical line) indicates that there is no association between a variant type (pathogenic or population) and a feature. To facilitate the visualization, minimum and maximum values of OR along x-axis are set to 0.05 and 20.0, respectively. For non-significant association ($q \geq 0.05$), the circle and CI bar are grey.

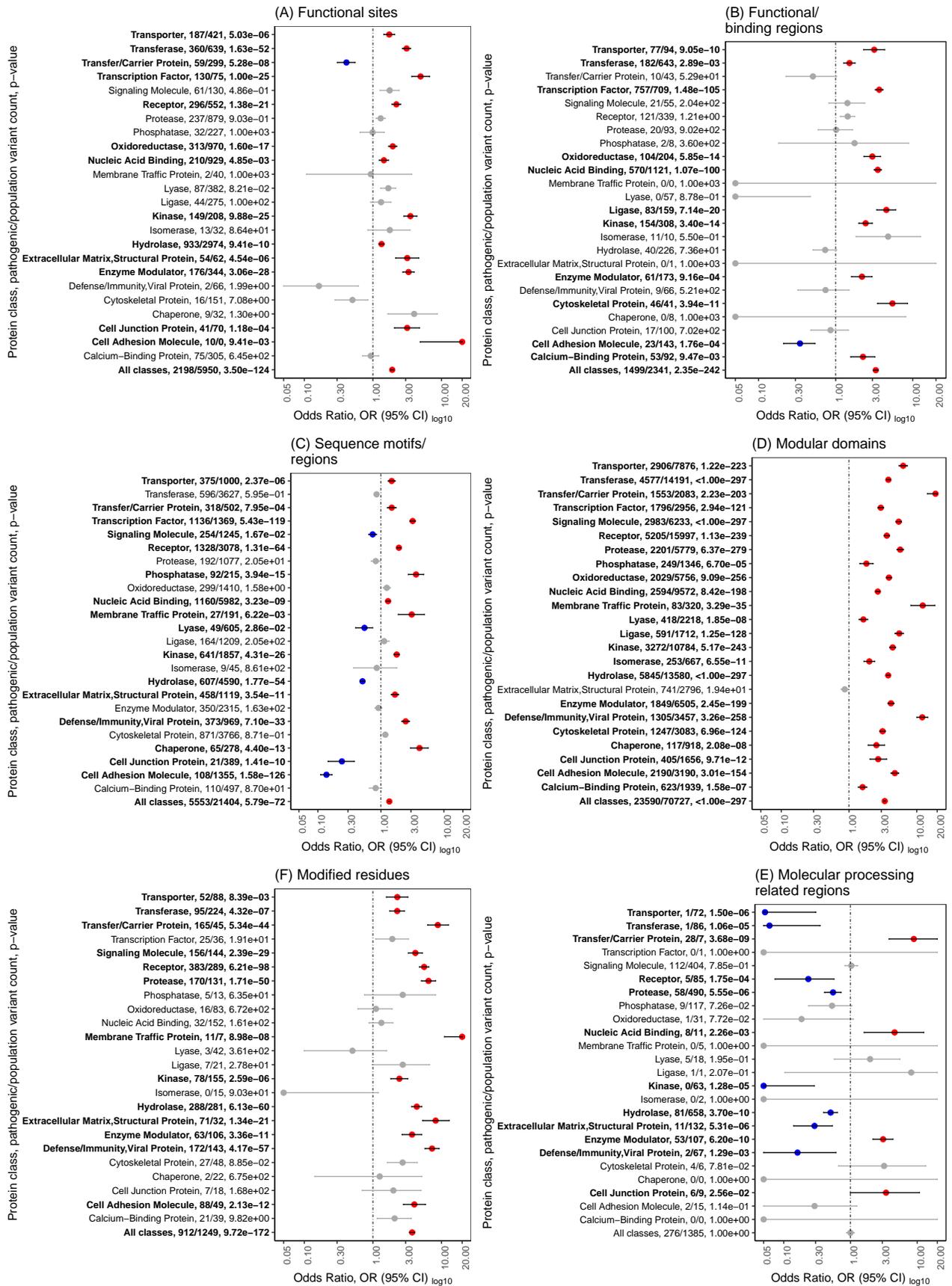


Fig. S8

Fig. S8. (Previous page.) Enrichment of pathogenic and population missense variants for “All classes” and twenty-four protein functional-classes in six UniProt-based functional features. The plot shows the results of two-sided Fisher’s exact tests of association. Circles show the odds ratio (OR) and the horizontal bars show the 95% confidence interval (CI). The y-axis ticks are labelled with the protein class name, counts of pathogenic and population variants with the corresponding feature, and corrected p or “ q ” values, showing the significance of the association (a value of $1.0e-297$ is should be read as $<1.0e-297$, indicating maximum significance). The $OR > 1$ and $OR < 1$, along with $q < 0.05$, indicate that the corresponding feature (plot label) has an enrichment of pathogenic (red circle) and population (blue circle) variants, for the respective protein class (y-axis). The $OR = 1$ (dashed vertical line) indicates that there is no association between a variant type (pathogenic or population) and a feature. To facilitate the visualization, minimum and maximum values of OR along x-axis are set to 0.05 and 20.0, respectively. For non-significant association ($q \geq 0.05$), the circle and CI bar are grey.

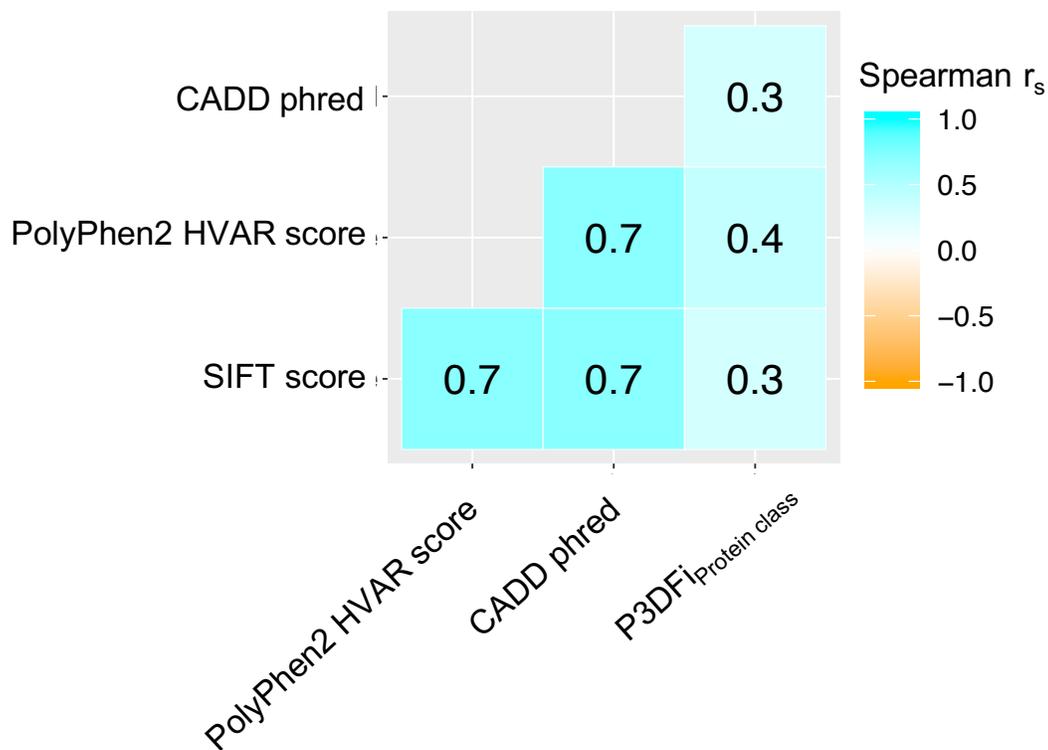


Fig. S9. Pairwise correlation (Spearman’s non-parametric test) between the protein class specific pathogenic 3D Feature index (P3DFi_{Protein class}) derived in this study and the three existing variant pathogenicity prediction scores: SIFT (22), PolyPhen2 (23) and CADD (24). Note that, unlike the other scores P3DFi was not generated by a learning model to predict pathogenicity, instead it quantifies the difference between the pathogenic and population variant associated 3D features for each amino acid. P3DFi > 0 thus indicates the 3D mutational hotspots. The correlation plot shows that P3DFi has the maximum correlation of only 40% with the PolyPhen2 score. We thus speculate that P3DFi, derived using the 3D structure related information only, can serve as an orthogonal determinant of pathogenicity when combined with these existing methods (see **Table 1** and **Fig. 5** in the main text for relevant results).

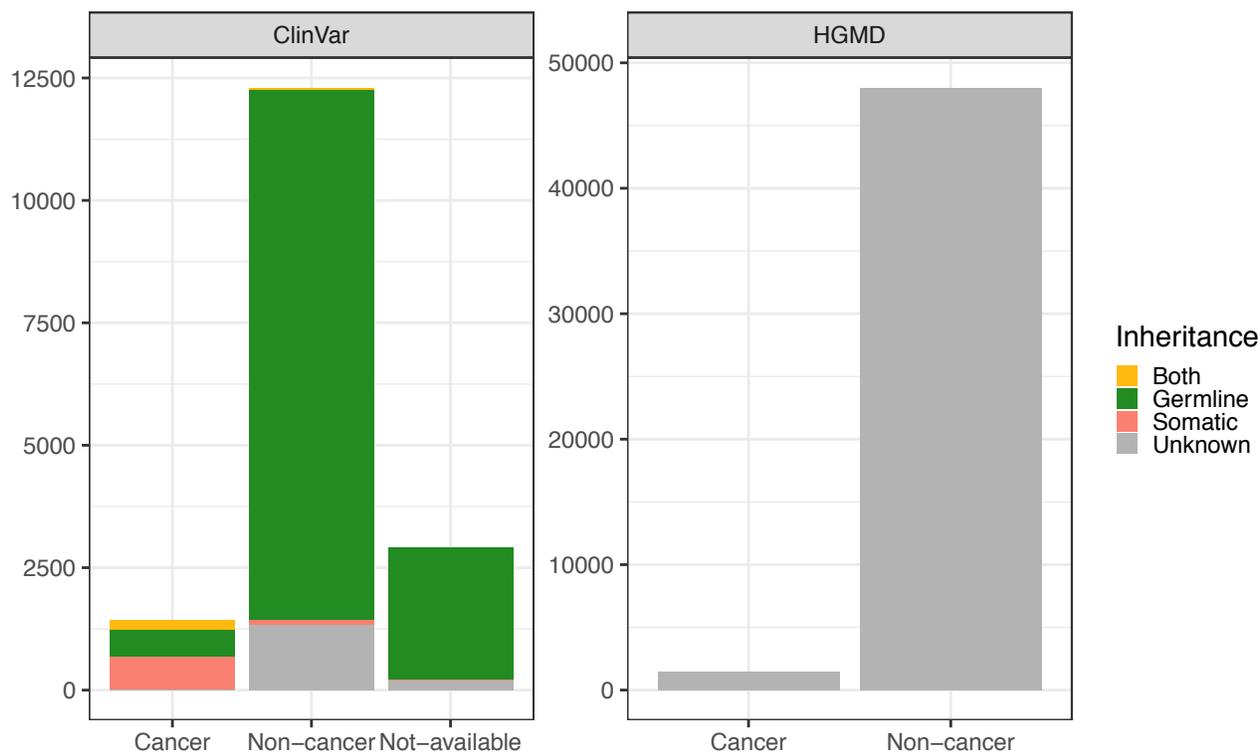
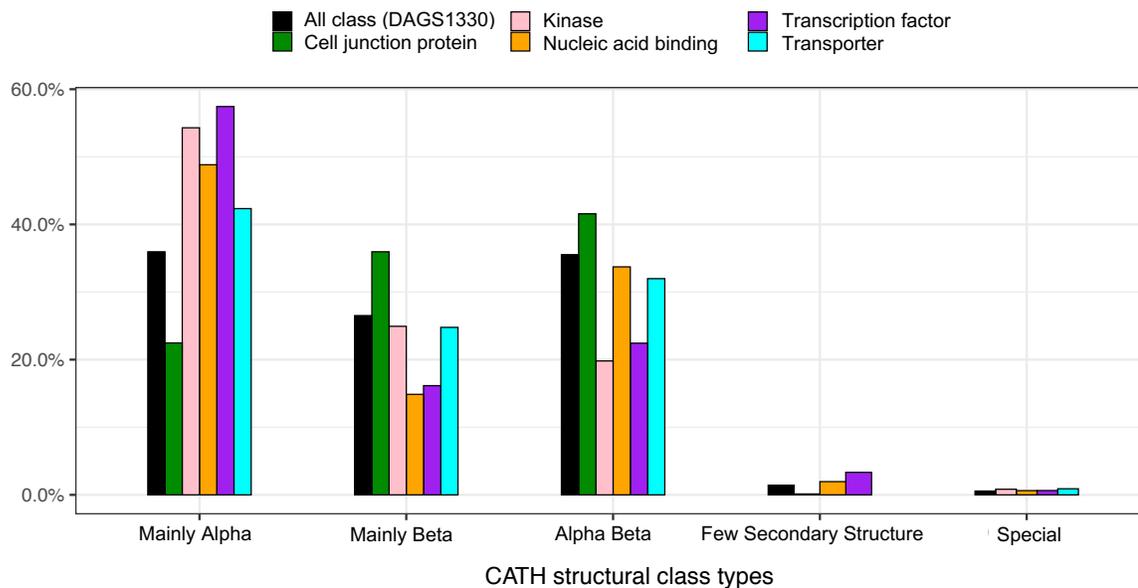


Fig. S10. Distribution of pathogenic variations in 1,330 genes (DAGS1330 set, analyzed in this work) associated with their respective phenotype: cancer, non-cancer and not-provided, and according to the variants' allele origin or inheritance, that is: somatic, germline, both and unknown (as available in the ClinVar database (7)). Variants were considered to be related to cancer if the phenotype description included at least one of the following keywords: "cancer", "carcinoma", "leukemia", "tumor", "neoplasm", "sarcoma", "melanoma", "lymphoma" or "glioma". "Not provided" and/or "Unknown" phenotypes were kept separate as "Not-provided". Out of the 16,638 pathogenic and likely pathogenic variations from ClinVar database (7), 74% (n = 12,301) are associated with non-cancer phenotype, 9% (n = 1,431) are associated with cancer, and for 17% (n = 2,906) the phenotype data are not available. At the same time, out of the 49,501 disease mutations from HGMD database (8), 97% (n = 48,013) are associated with non-cancer phenotype. In total, out of all pathogenic variants, 91% have non-cancer mutations and only about 4% have cancer mutations. Importantly, only about 1% of all variations in these 1,330 genes obtained from ClinVar and HGMD databases are somatic (assuming all HGMD variants are germline according to the literature (8)).

A



B

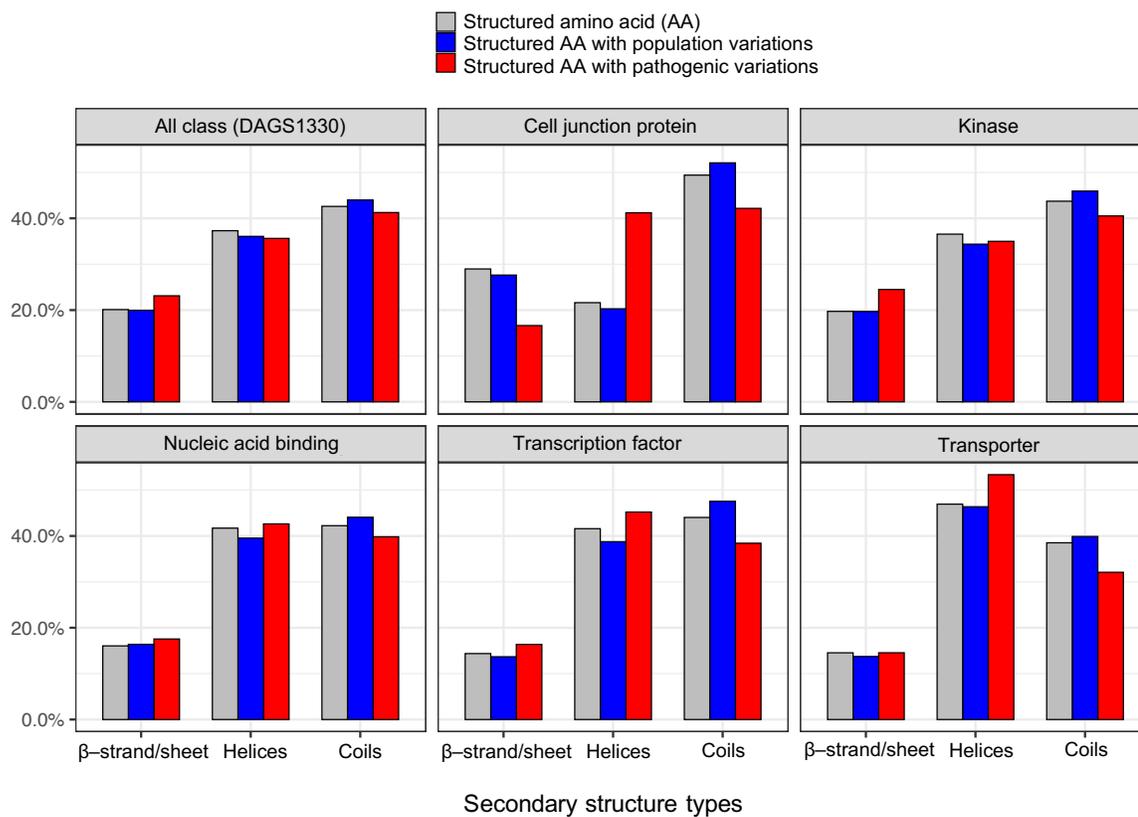


Fig. S11

Fig. S11. (Previous page.) Pathogenic variations were found enriched in α -helix for five protein classes: cell junction protein, kinase, nucleic acid binding, transcription factor, and transporter, in contrast to the “All class (DAGS1330)” set, because of a higher observed rate of pathogenic variations on α -helix than the population variations. (A) Distribution of different folds adopted by proteins from all classes and from the five aforementioned classes (collected from the CATH Protein Structural Classification database (27)). For 1,170 out of 1,330 proteins, CATH annotations were available (version: April 2020) for 5,649 unique chains in a structure from the PDB (1). (B) Proportion of all residues, residues mutated in pathogenic variants and population variants of different secondary structure types (β -strand/sheet, helices, and coils), for all classes and for the five above mentioned protein classes. We observed that a higher proportion of protein chains from all classes adopt a “mainly alpha” fold (~36%) than “mainly beta” fold (~26%) (in A). This is also the case for four out of these five protein classes (kinase, nucleic acid binding, transcription factor, and transporter) where we find an enrichment of pathogenic mutations in helices, meaning that the percentage of pathogenic mutations is higher in this structural motif than that of population mutations. Notably, though, even in cell junction proteins, which adopt a mainly beta fold (~37% beta vs. ~22% alpha), we find an enrichment of pathogenic (42%) over population (20%) mutations in helices (OR = 2.8, *SI Appendix*, Fig. S2B), especially α -helix (OR = 2.9, *SI Appendix*, Fig. S3D) and α -helix (OR = 10.5, *SI Appendix*, Fig. S3E). This exemplifies that our results are robust with respect to the structural composition of proteins and, more in general, to the abundance of a given feature (see *Materials and Methods*).

Table S1. Number of 3D features associated to pathogenic and population missense variants (referred to as the characteristic features identified in this study) for “All classes” by joint analysis of 1,330 genes in the DAGS1330 dataset and individually for twenty-four protein functional classes (class-specific features).

Protein class	Characteristic 3D features of pathogenic variants	Class-specific characteristic features of pathogenic variants outside the “All class” based features (1 st row, 2 nd column)*	Characteristic 3D features of population variants	Class-specific characteristic features of population variants outside the “All class” based features (1 st row, 4 th column)*	Fraction of unique 3D features specific for protein classes outside the “All class” based features (%) (1 st row, 2 nd and 4 th columns)+
All classes (DAGS1330)	18	n/a	14	n/a	n/a
Calcium-Binding Protein	6	0	10	2	12.5
Cell Adhesion Molecule	13	3	14	6	33.3
Cell Junction Protein	7	3	6	3	46.2
Chaperone	2	0	0	0	0
Cytoskeletal Protein	6	0	2	0	0
Defense/Immunity protein	7	1	8	2	20.0
Enzyme Modulator	9	2	7	1	18.8
Extracellular Matrix Protein	9	1	8	2	17.6
Hydrolase	12	3	15	6	33.3
Isomerase	5	1	3	1	25.0
Kinase	15	4	9	2	25.0
Ligase	10	0	7	0	0
Lyase	3	0	5	1	12.5
Membrane Traffic Protein	5	0	0	0	0
Nucleic Acid Binding	13	4	8	1	23.8
Oxidoreductase	11	0	6	0	0
Phosphatase	7	0	6	0	0
Protease	12	2	11	4	26.1
Receptor	11	0	10	1	4.8
Signaling Molecule	9	0	16	3	12.0
Transcription Factor	13	3	9	3	27.3
Transfer/Carrier Protein	12	1	7	1	10.5
Transferase	12	0	10	1	4.5
Transporter	13	3	7	3	30.0

* Number of unique class-specific features of pathogenic and population variants outside the **18** and **14** features found by analyzing all genes, respectively.

+ Fraction of unique (out of the total) protein class-specific pathogenic and population variant associated features.

Table S2. Number of genes, missense variants/amino acid substitutions, protein structure in different steps of Disease-Associated Genes with Structure, namely DAGS1330, dataset preparation

	Population variants	Pathogenic variants	
	gnomAD (missense)	ClinVar (likely/-pathogenic missense)	HGMD (disease mutation)
Genes with experimentally solved structure available for the encoded proteins	5,724	1,466	1,673
Missense variants of the genes with structure available for the encoded proteins	1,485,579	16,570	47,036
Amino acid substitutions mapped on the structure	496,869	8,137	30,730
Genes with missense variants mappable on the structure	4,897	1,330*	
Protein structures onto which the variants were mapped	29,870	14,270*	
Missense variants of DAGS1330 genes that were mapped on the structure	164,915	32,923*	

* combined (unique) counts from HGMD and ClinVar databases.

Number of genes, variants, and protein structures included in the **DAGS1330** dataset to perform the statistical analysis is bold faced.

Table S3. Comparison of the most predictive P3DFi (P3DFi > 2 and P3DFi < -2) to three state-of-the-art predictors (SIFT (22), PolyPhen2 (23) and CADD (24)) in stratifying pathogenic variants from benign variants.

Method	Recall/ Sensitivity/ True Positive Rate	Selectivity/ Specificity/ True Negative Rate	Balanced Accuracy	Matthews Correlation Coefficient (MCC)	F1 Score	Precision	Fall-out/ False Positive Rate	Miss Rate/ False Negative Rate
P3DFi _{Protein class}	83.68%	85.65%	84.66%	67.18%	87.53%	91.74%	14.35%	16.32%
P3DFi* _{DAGS1330}	79.58%	<u>83.73%</u>	81.66%	60.90%	84.61%	90.31%	<u>16.27%</u>	20.42%
SIFT (22)	<u>90.54%</u>	73.84%	82.19%	65.79%	88.65%	86.84%	26.16%	<u>9.46%</u>
PolyPhen2 (23)	88.28%	78.95%	<u>83.62%</u>	<u>67.03%</u>	<u>88.58%</u>	88.88%	21.05%	11.72%
CADD (24)	95.56%	61.08%	78.32%	63.08%	88.49%	82.40%	38.92%	4.44%

The best and second best score values in each metric are boldfaced and underlined, respectively.

The performances are evaluated on the variants (total count = 1, 824) with the highest and lowest P3DFi values (see Fig. 4, P3DFi > 2 and P3DFi < -2).

* See computation of P3DFi in *Materials and Methods*.

181 **SI Dataset S1 (Dataset_S1.xlsx)**

182 Information of the genes in DAGS1330 (disease-associated genes with structure) dataset. Gene name (count = 1,330),
183 UniProt protein accession number, protein length, CCDS identifier, Ensemble gene identifier, and HGNC identifier.

184 **SI Dataset S2 (Dataset_S2.xlsx)**

185 Protein functional class information. Protein class name, number of genes and protein structures in each class, number of
186 pathogenic and population variations mapped on structures for each class, and the list of genes annotated in each protein class.

187 **SI Dataset S3 (Dataset_S3.xlsx)**

188 Protein 3D features associated to pathogenic and population missense variants (identified in this study) for “All classes”
189 (joint analysis of 1,330 genes in the DAGS1330 dataset), and individually for twenty-four protein functional classes. The
190 columns represent seven main feature categories and each cell outlines the 3D feature that has a significant burden of pathogenic
191 and/or population variations (two-sided Fisher’s exact test, corrected p -value or “ q ” < 0.05). The ‘x’ indicates that no feature
192 of the corresponding category was significantly associated with pathogenic or population variants.

193 **SI Dataset S4 (Dataset_S4.xlsx)**

194 Training dataset to develop the “Random forest” ensemble models. The first column reports the true class labels (1:
195 deleterious or pathogenic, -1: neutral or benign). The next three columns show the scores generated by the variant pathogenicity
196 predictors: SIFT, PolyPhen2, and CADD. The all protein (DAGS1330 dataset, **Dataset S1**) based and protein functional
197 class specific pathogenic 3D feature index (P3DFi_{DAGS1330} and P3DFi_{Protein class}) derived in this study are reported in the last two
198 columns. For 209,110 variants (31,913 pathogenic and 177,197 benign) in the DAGS1330 dataset, all the existing predictors’
199 scores were available, and therefore, were used for training. Three Random forest classifiers (number of decision trees =
200 2000, evaluation function = “gini”, depth of the trees = 10) were developed: two models were trained separately with the
201 P3DFi_{DAGS1330} and P3DFi_{Protein class} values in addition to the scores from SIFT, PolyPhen2 and CADD, and the third model was
202 trained without any P3DFi values (see **Table 1** and **Fig. 5** in the main text for relevant results).

203 **SI Dataset S5 (Dataset_S5.xlsx)**

204 Test dataset to evaluate the “Random forest” ensemble models and compare their outputs to SIFT, PolyPhen2, and CADD
205 scores. The first column reports the true class labels (1: deleterious or pathogenic, -1: neutral or benign). The next three
206 columns show the scores generated by the variant pathogenicity predictors: SIFT, PolyPhen2, and CADD. Pathogenic 3D
207 feature indices (P3DFi_{DAGS1330} and P3DFi_{Protein class}) derived in this study are reported in the fifth and sixth columns, respectively.
208 For 22,362 (17,707 pathogenic and 4,655 benign) out of 22,695 variants in the validation set, all the existing predictors’
209 scores were available, and therefore, were used for evaluation (see **Table 1** and **Fig. 5** in the main text for relevant results).
210 Predicted probability outputs (>0.5 : deleterious or pathogenic, ≤ 0.5 : neutral or benign) from three ensemble models developed
211 in this study are given in the last three columns.

212 **References**

- 213 1. HM Berman, PE Bourne, J Westbrook, C Zardecki, The protein data bank in *Protein Structure*. (CRC Press), pp. 394–410
214 (2003).
- 215 2. R Apweiler, et al., Uniprot: the universal protein knowledgebase. *Nucleic acids research* **32**, D115–D119 (2004).
- 216 3. S Velankar, et al., Sifts: structure integration with function, taxonomy and sequences resource. *Nucleic acids research* **41**,
217 D483–D489 (2012).
- 218 4. KJ Karczewski, et al., The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**,
219 434–443 (2020).
- 220 5. P Danecek, et al., The variant call format and vcftools. *Bioinformatics* **27**, 2156–2158 (2011).
- 221 6. DR Zerbino, et al., Ensembl 2018. *Nucleic acids research* **46**, D754–D761 (2018).
- 222 7. MJ Landrum, et al., Clinvar: improving access to variant interpretations and supporting evidence. *Nucleic acids research*
223 **46**, D1062–D1067 (2018).
- 224 8. PD Stenson, et al., The human gene mutation database: building a comprehensive mutation repository for clinical and
225 molecular genetics, diagnostic testing and personalized genomic medicine. *Hum. genetics* **133**, 1–9 (2014).
- 226 9. TL Mighell, S Evans-Dutson, BJ O’Roak, A saturation mutagenesis approach to understanding pten lipid phosphatase
227 activity and genotype-phenotype relationships. *The Am. J. Hum. Genet.* **102**, 943–955 (2018).
- 228 10. GM Findlay, et al., Accurate classification of brca1 variants with saturation genome editing. *Nature* **562**, 217–222 (2018).
- 229 11. W Kabsch, C Sander, Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical
230 features. *Biopolym. Orig. Res. on Biomol.* **22**, 2577–2637 (1983).
- 231 12. MZ Tien, AG Meyer, DK Sydykova, SJ Spielman, CO Wilke, Maximum allowed solvent accessibilities of residues in
232 proteins. *PLoS one* **8** (2013).
- 233 13. RA Laskowski, J Jabłońska, L Pravda, RS Vařeková, JM Thornton, Pdbsum: Structural summaries of pdb entries. *Protein*
234 *science* **27**, 129–134 (2018).

- 235 14. PV Hornbeck, et al., Phosphositeplus, 2014: mutations, ptms and recalibrations. *Nucleic acids research* **43**, D512–D520
236 (2015).
- 237 15. K Wang, M Li, H Hakonarson, Annovar: functional annotation of genetic variants from high-throughput sequencing data.
238 *Nucleic acids research* **38**, e164–e164 (2010).
- 239 16. H Mi, et al., Panther version 11: expanded annotation data from gene ontology and reactome pathways, and data analysis
240 tool enhancements. *Nucleic acids research* **45**, D183–D189 (2017).
- 241 17. DR Zerbino, et al., Ensembl 2018. *Nucleic acids research* **46**, D754–D761 (2018).
- 242 18. GO Consortium, The gene ontology resource: 20 years and still going strong. *Nucleic acids research* **47**, D330–D338
243 (2019).
- 244 19. R Beaglehole, R Bonita, T Kjellström, , et al., *Basic epidemiology*. (World Health Organization Geneva), (1993).
- 245 20. A Viera, Odds ratios and risk ratios: what's the difference and why does it matter? *South. medical journal* **101**, 730–734
246 (2008).
- 247 21. SD Simon, Understanding the odds ratio and the relative risk. *J. andrology* **22**, 533–536 (2001).
- 248 22. PC Ng, S Henikoff, Sift: Predicting amino acid changes that affect protein function. *Nucleic acids research* **31**, 3812–3814
249 (2003).
- 250 23. IA Adzhubei, et al., A method and server for predicting damaging missense mutations. *Nat. methods* **7**, 248–249 (2010).
- 251 24. P Rentzsch, D Witten, GM Cooper, J Shendure, M Kircher, Cadd: predicting the deleteriousness of variants throughout
252 the human genome. *Nucleic acids research* **47**, D886–D894 (2019).
- 253 25. O Wagih, et al., A resource of variant effect predictions of single nucleotide variants in model organisms. *Mol. systems*
254 *biology* **14** (2018).
- 255 26. M Hicks, I Bartha, J di Iulio, JC Venter, A Telenti, Functional characterization of 3d protein structures informed by
256 human genetic diversity. *Proc. Natl. Acad. Sci.* **116**, 8960–8965 (2019).
- 257 27. M Knudsen, C Wiuf, The cath database. *Hum. Genomics* **4**, 207–212 (2010).