

PLOS ONE

Predicting hotspots for disease-causing single nucleotide variants using sequences-based coevolution, network analysis, and machine learning --Manuscript Draft--

Manuscript Number:	PONE-D-23-33350R1
Article Type:	Research Article
Full Title:	Predicting hotspots for disease-causing single nucleotide variants using sequences-based coevolution, network analysis, and machine learning
Short Title:	Predicting nSNVs hotspots with sequences-based machine learning
Corresponding Author:	Wenjun Zheng University at Buffalo, The State University of New York Buffalo, NY UNITED STATES
Keywords:	Centrality, Coevolution, Disease mutations, Machine learning, Protein residue contact network, Single nucleotide variant
Abstract:	<p>To enable personalized medicine, it is important yet highly challenging to accurately predict disease-causing mutations in target proteins at high throughput. Previous computational methods have been developed using evolutionary information in combination with various biochemical and structural features of protein residues to discriminate neutral vs. deleterious mutations. However, the power of these methods is often limited because they either assume known protein structures or treat residues independently without fully considering their global interactions. To address the above limitations, we build upon recent progress in machine learning, network analysis, and protein language models, and develop a sequences-based variant site prediction workflow based on the protein residue contact networks: 1. We employ and integrate various methods of building protein residue networks using state-of-the-art coevolution analysis tools (RaptorX, DeepMetaPSICOV, and SPOT-Contact) powered by deep learning. 2. We use machine learning algorithms (Random Forest, Gradient Boosting, and Extreme Gradient Boosting) to optimally combine 20 network centrality scores to jointly predict key residues as hot spots for disease mutations. 3. Using a dataset of 107 proteins rich in disease mutations, we rigorously evaluate the network scores individually and collectively (via machine learning). This work supports a promising strategy of combining an ensemble of network scores based on different coevolution analysis methods (and optionally predictive scores from other methods) via machine learning to predict candidate sites of disease mutations, which will inform downstream applications of disease diagnosis and targeted drug design.</p>
Order of Authors:	Wenjun Zheng
Opposed Reviewers:	
Response to Reviewers:	see attached file.
Additional Information:	
Question	Response
Financial Disclosure	Yes
Enter a financial disclosure statement that describes the sources of funding for the work included in this submission. Review the submission guidelines for detailed requirements. View published research articles from PLOS ONE for specific examples.	

This statement is required for submission and **will appear in the published article** if the submission is accepted. Please make sure it is accurate.

Funded studies

Enter a statement with the following details:

- Initials of the authors who received each award
- Grant numbers awarded to each author
- The full name of each funder
- URL of each funder website
- Did the sponsors or funders play any role in the study design, data collection and analysis, decision to publish, or preparation of the manuscript?

Did you receive funding for this work?

Please add funding details.
as follow-up to "**Financial Disclosure**

Enter a financial disclosure statement that describes the sources of funding for the work included in this submission. Review the [submission guidelines](#) for detailed requirements. View published research articles from [PLOS ONE](#) for specific examples.

This statement is required for submission and **will appear in the published article** if the submission is accepted. Please make sure it is accurate.

Funded studies

Enter a statement with the following details:

- Initials of the authors who received each award
- Grant numbers awarded to each author
- The full name of each funder
- URL of each funder website
- Did the sponsors or funders play any role in the study design, data collection and analysis, decision to publish, or preparation of the manuscript?

Did you receive funding for this work?"

This study is funded by a grant from NIH.

Please select the country of your main research funder (please select carefully as in some cases this is used in fee calculation).

as follow-up to "**Financial Disclosure**

Enter a financial disclosure statement that describes the sources of funding for the work included in this submission. Review the [submission guidelines](#) for detailed requirements. View published research articles from [PLOS ONE](#) for specific examples.

This statement is required for submission and **will appear in the published article** if the submission is accepted. Please make sure it is accurate.

Funded studies

Enter a statement with the following details:

- Initials of the authors who received each award
- Grant numbers awarded to each author
- The full name of each funder
- URL of each funder website
- Did the sponsors or funders play any role in the study design, data collection and analysis, decision to publish, or preparation of the manuscript?

Did you receive funding for this work?"

Competing Interests

Use the instructions below to enter a competing interest statement for this submission. On behalf of all authors, disclose any [competing interests](#) that could be perceived to bias this work—acknowledging all financial support and any other relevant financial or non-financial competing interests.

This statement is **required** for submission and **will appear in the published article** if the submission is accepted. Please make sure it is accurate and that any funding

UNITED STATES - US

The authors have declared that no competing interests exist.

sources listed in your Funding Information later in the submission form are also declared in your Financial Disclosure statement.

View published research articles from [PLOS ONE](#) for specific examples.

NO authors have competing interests

Enter: *The authors have declared that no competing interests exist.*

Authors with competing interests

Enter competing interest details beginning with this statement:

I have read the journal's policy and the authors of this manuscript have the following competing interests: [insert competing interests here]

* typeset

Ethics Statement

N/A

Enter an ethics statement for this submission. This statement is required if the study involved:

- Human participants
- Human specimens or tissue
- Vertebrate animals or cephalopods
- Vertebrate embryos or tissues
- Field research

Write "N/A" if the submission does not require an ethics statement.

General guidance is provided below. Consult the [submission guidelines](#) for detailed instructions. **Make sure that all information entered here is included in the Methods section of the manuscript.**

Format for specific study types

Human Subject Research (involving human participants and/or tissue)

- Give the name of the institutional review board or ethics committee that approved the study
- Include the approval number and/or a statement indicating approval of this research
- Indicate the form of consent obtained (written/oral) or the reason that consent was not obtained (e.g. the data were analyzed anonymously)

Animal Research (involving vertebrate animals, embryos or tissues)

- Provide the name of the Institutional Animal Care and Use Committee (IACUC) or other relevant ethics board that reviewed the study protocol, and indicate whether they approved this research or granted a formal waiver of ethical approval
- Include an approval number if one was obtained
- If the study involved *non-human primates*, add *additional details* about animal welfare and steps taken to ameliorate suffering
- If anesthesia, euthanasia, or any kind of animal sacrifice is part of the study, include briefly which substances and/or methods were applied

Field Research

Include the following details if this study involves the collection of plant, animal, or other materials from a natural setting:

- Field permit number
- Name of the institution or relevant body that granted permission

Data Availability

Authors are required to make all data underlying the findings described fully available, without restriction, and from the time of publication. PLOS allows rare exceptions to address legal and ethical concerns. See the [PLOS Data Policy](#) and [FAQ](#) for detailed information.

Yes - all data are fully available without restriction

A Data Availability Statement describing where the data can be found is required at submission. Your answers to this question constitute the Data Availability Statement and **will be published in the article**, if accepted.

Important: Stating 'data available on request from the author' is not sufficient. If your data are only available upon request, select 'No' for the first question and explain your exceptional situation in the text box.

Do the authors confirm that all data underlying the findings described in their manuscript are fully available without restriction?

Describe where the data may be found in full sentences. If you are copying our sample text, replace any instances of XXX with the appropriate details.

- If the data are **held or will be held in a public repository**, include URLs, accession numbers or DOIs. If this information will only be available after acceptance, indicate this by ticking the box below. For example: *All XXX files are available from the XXX database (accession number(s) XXX, XXX).*
- If the data are all contained **within the manuscript and/or Supporting Information files**, enter the following: *All relevant data are within the manuscript and its Supporting Information files.*
- If neither of these applies but you are able to provide **details of access elsewhere**, with or without limitations, please do so. For example:

Data cannot be shared publicly because of [XXX]. Data are available from the XXX Institutional Data Access / Ethics Committee (contact via XXX) for researchers who meet the criteria for access to confidential data.

The data underlying the results presented in the study are available from (include the name of the third party

All relevant data and code are available at <https://simtk.org/projects/hotspots>

and contact information or URL).

- This text is appropriate if the data are owned by a third party and authors do not have permission to share the data.

* typeset

Additional data availability information:

PONE-D-23-33350

Predicting hotspots for disease-causing single nucleotide variants using sequences-based coevolution, network analysis, and machine learning

Dear Editor,

Thank you for reviewing my manuscript. Following your and the reviewers' comments, I have revised the manuscript accordingly to address the points raised during the review process.

In this resubmission, I have included the following items as requested:

- A rebuttal letter (labeled 'Response to Reviewers') that responds to each point raised by the academic editor and reviewer(s).
- A marked-up copy of the manuscript that highlights changes made to the original version. It is labeled 'Revised Manuscript with Track Changes'.
- An unmarked version of the revised paper without tracked changes. It is labeled 'Manuscript'.

Figure files were uploaded separately from the manuscript (Fig1.tif etc).

Sincerely,

Wenjun Zheng, PhD

1

2 **Predicting hotspots for disease-causing single nucleotide variants**
3 **using sequences-based coevolution, network analysis, and machine**
4 **learning**

5 **Short title:** Predicting nSNVs hotspots with sequences-based machine learning

6

7

Wenjun Zheng ^{1*}

8 ¹ Department of Physics, State University of New York at Buffalo, NY 14260, United States of
9 America

10

11 * Corresponding author

12 E-mail: wjzheng@buffalo.edu (WZ)

13

14 **Keywords:** Centrality, Coevolution, Disease mutations, Machine learning, Protein residue
15 contact network, Single nucleotide variant

16

17 **Abstract**

18 To enable personalized medicine, it is important yet highly challenging to accurately
19 predict disease-causing mutations in target proteins at high throughput. Previous computational
20 methods have been developed using evolutionary information in combination with various
21 biochemical and structural features of protein residues to discriminate neutral vs. deleterious
22 mutations. However, the power of these methods is often limited because they either assume
23 known protein structures or treat residues independently without fully considering their global
24 interactions. To address the above limitations, we build upon recent progress in machine
25 learning, network analysis, and protein language models, and develop a sequences-based variant
26 site prediction workflow based on the protein residue contact networks: 1. We employ and
27 integrate various methods of building protein residue networks using state-of-the-art coevolution
28 analysis tools (RaptorX, DeepMetaPSICOV, and SPOT-Contact) powered by deep learning. 2.
29 We use machine learning algorithms (Random Forest, Gradient Boosting, and Extreme Gradient
30 Boosting) to optimally combine 20 network centrality scores to jointly predict key residues as
31 hot spots for disease mutations. 3. Using a dataset of 107 proteins rich in disease mutations, we
32 rigorously evaluate the network scores individually and collectively (via machine learning). This
33 work supports a promising strategy of combining an ensemble of network scores based on
34 different coevolution analysis methods (and optionally predictive scores from other methods) via
35 machine learning to predict candidate sites of disease mutations, which will inform downstream
36 applications of disease diagnosis and targeted drug design.

37

38 **Introduction**

39 The holy grail of structural biology is to solve high-resolution biomolecular structures at
40 the genomic scale to inform mechanistic studies of their functions. Thanks to recent revolutions
41 in computational structural biology (accurate protein structure prediction by AlphaFold [1] and
42 RoseTTAFold [2]), it is now feasible to predict native structures for many proteins given their
43 sequences (with some caveats, see [3]), thus practically solving the protein folding problem [4].
44 However, it remains challenging to predict dynamic structural ensembles [5] and mutation-
45 induced effects [6] to meet the demand of mechanistic studies of protein functions and
46 dysfunctions. While the public databases of protein sequences and variations increase rapidly
47 owing to genomic/metagenomic sequencing efforts (the MetaClust database contains about 1.6
48 billion protein sequence fragments [7]), the growth of experimental protein structures [8] and
49 predicted structures remains to catch up (the AlphaFold database contains over 200 million
50 predicted structures [9]). Such sequences-structures gap has motivated the development of new
51 computational tools that make functional sense of protein sequences without directly using
52 structural information (for example, by using deep learning to train large protein language
53 models [10]). Recently, AlphaMissense attained state of the art prediction of missense variant
54 pathogenicity by adapting AlphaFold fine-tuned on human and primate variant population
55 frequency databases [11].

56 A major interest in personalized medicine is in understanding novel genetic variations
57 through genotype-phenotype association studies in relation to diseases. Particularly, a rapidly
58 growing number of non-synonymous single nucleotide variants (nSNVs) have been uncovered in
59 protein coding regions that can adversely impact protein function and cause diseases [12].

60 Various computational methods were developed using evolutionary conservation and phylogeny
61 in combination with biochemical and structural properties of amino acids to discriminate neutral
62 vs. deleterious nSNVs [13-22]. Protein structural dynamics has also proven useful in discovering
63 functionally important residues [23,24] which could constitute hot spots for disease-causing
64 nSNVs [25,26]. However, the requirement of 3D structures has limited the number of nSNVs
65 that can be analyzed by existing structure-based computational tools, although such constraint
66 has been significantly alleviated by recent progress in protein structure prediction [27].

67 As alternatives to structure-based methods, sequences-based coevolution analysis has
68 become increasingly powerful in predicting structural couplings between pairs of contacting
69 residues [28-31] , thanks to the development of direct coupling methods that can overcome the
70 confounding indirect coupling effects [29,32,33] . In principle, coevolving pairs of residues can
71 be identified from a sufficiently large multiple sequence alignment, allowing the prediction of
72 close spatial proximity in the native structures. Boosted by deep learning and other algorithmic
73 developments, this coevolution analysis has led to accurate prediction of residue contacts which
74 make *de novo* protein structure prediction possible [28] . Furthermore, coevolution analysis
75 (enhanced by deep learning) has also been used to study various aspects of protein functional
76 interactions such as allostery [34] . For example, RaptorX uses an ultra-deep neural network
77 combining coevolution information with sequence conservation information to infer 3D contacts
78 with higher accuracy than previous methods [35,36]. DeepMetaPSICOV [37] combines the input
79 feature sets used by earlier methods (MetaPSICOV [38] and DeepCov [39]) as input to a deep,
80 fully convolutional residual neural network. SPOT-Contact predicts protein contact maps by
81 stacking residual convolutional networks with two-dimensional residual bidirectional recurrent
82 LSTM networks, and using both one-dimensional sequence-based and two-dimensional

83 evolutionary coupling based information [40]. These three state-of-the-art coevolution analysis
84 methods are employed in this study to construct protein residue contact maps for network
85 analysis (see below).

86 Another line of protein research is based on the treatment of a protein as a network where
87 amino acid residues are nodes and their bonded/non-bonded interactions form edges [41]. Such
88 models can be readily built upon 3D native structures so that a whole suite of network analysis
89 tools (see <https://networkx.org/>) can be applied. For example, Amitai et al [42] used network
90 analysis of protein structures (using closeness centrality) to identify functional residues. Going
91 beyond network analysis, deep-learning-based study of protein graph neural networks is an
92 active area of research [43].

93 In a recent paper, Butler et al [44] proposed a sequence-based Gaussian network model
94 (Seq-GNM) to calculate the dynamic profile of a protein without a 3D structure. They used
95 coevolution analysis to build a network model which connects residues predicted to be in contact
96 via evolutionary couplings. Their work built on previous studies that shown crystallographic B-
97 factors are useful in predicting the impact of nSNVs on protein function [45,46] : rigid sites with
98 low B-factors are more susceptible to destabilizing nSNVs than flexible sites with high B-
99 factors. Indeed, existing computational tools to diagnose neutral and deleterious nSNVs (such as
100 PolyPhen-2 [47]) use crystallographic B-factors along with other evolutionary and structural
101 features. More specifically, Butler et al used Seq-GNM to compute B-factors for protein
102 residues, and they found that deleterious nSNVs are overabundant at low B-factor sites, while
103 neutral nSNVs are overabundant at high B-factor sites. Mechanistically, low B-factors may
104 indicate that a site is crucial for maintaining structural stability and/or modulating functional
105 motions (as a hinge) and thus susceptible to mutations. In contrast, high B-factors are associated

106 with flexible regions with minimal interactions, which are thus more robust to mutations. Based
107 on these observations, they proposed that the sequences-based predicted B-factors can
108 discriminate between deleterious and neutral nSNVs without structural information.

109 Inspired by the above study and recent progress in machine learning, network analysis,
110 and protein language models, we further develop the sequences-based protein residue network
111 analysis in the following directions: 1. We build protein residue networks using three different
112 coevolution analysis tools (RaptorX, DeepMetaPSICOV, and SPOT-Contact) as enabled by deep
113 learning. 2. We exploit three machine learning algorithms (Random Forest, Gradient Boosting,
114 and Extreme Gradient Boosting) to optimally combine 20 distinct network node centrality scores
115 as calculated from the contact probability matrices to predict hot spot residues for disease
116 mutations. 3. Based on a dataset of 107 proteins with known deleterious/neutral mutations, we
117 evaluate our sequences-based network scores both individually and in combination, and then
118 compare with alternative structures-based network scores and a physics force field based
119 method. By optimally combining three coevolution analysis methods and the resulting 20 network
120 scores by machine learning, we are able to discriminate deleterious and neutral mutation sites
121 accurately (AUC of ROC \sim 0.84), which is on par with structure-based network scores (AUC \sim
122 0.83). Furthermore, by combining our method with a state-of-the-art predictor of the functional
123 effects of sequence variation based on large protein language models (ESM [48]), we have
124 significantly improved the prediction of disease variant sites (AUC \sim 0.89).

125 In the following sections, we first describe the detailed methodology in the order of
126 the proposed workflow, then we report the results of evaluation of our network-based scores both
127 individually and collectively (via machine learning), finally we discuss specific case studies of
128 four proteins to illustrate the usage of our method.

129 **Materials and methods**

130 Here is a summary of the workflow of our sequences-based method:

- 131 a. Collect datasets of protein sequences and variants (see Section 1)
- 132 b. Run co-evolution analysis of a given target protein sequence to build a residue
133 contact map P (see Section 2)
- 134 c. Use NetworkX to calculate node centrality scores based on P (see Section 3)
- 135 d. Use sequence-based GNM to calculate additional node scores (see Section 4)
- 136 e. (optional) Use protein language model (ESM) to predict variant importance (see
137 Section 5)
- 138 f. (optional) Use AlphaFold and FoldX to predict variant importance (see Section 6 and
139 7)
- 140 g. Use machine learning to optimally combine the above scores for classifying
141 deleterious vs neutral variant sites (see Section 8)

142 **1. Datasets of protein sequences and variants**

143 A dataset of 107 protein sequences with ≤ 500 residues and ≥ 20 annotated
144 deleterious/neutral variants were collected from the HumVar database [47] (sources: humvar-
145 2011_12.deleterious.pph.input and humvar-2011_12.neutral.pph.input from
146 ftp://genetics.bwh.harvard.edu/pph2/training/training-2.2.2.tar.gz). Their UniProt ids and
147 sequences can be accessed at <https://simtk.org/projects/hotspots>. This diverse dataset contains 97
148 proteins with their pairwise sequence identity $< 30\%$.

149

150 The HumVar dataset consists of 13,032 human disease-causing mutations from UniProt
151 and 8,946 human nonsynonymous single-nucleotide polymorphisms (nsSNPs) without annotated
152 involvement in disease. This dataset was previously used to train and test PolyPhen-2 [47] for
153 predicting damaging effects of missense mutations, and was used by Butler et al [44] in
154 benchmarking their seq-GNM method for predicting deleterious/neutral nSNVs.

155
156 Since this dataset is highly imbalanced (there are 4040 deleterious mutation sites but only
157 120 neutral mutation sites) [49], we have added 3403 additional neutral sites with very low
158 conservation scores (i.e. grade ≤ 2 as assessed by the ConSurf program [50]). Our objective is to
159 train and test a binary classifier of residues in these proteins as deleterious or neutral. To this
160 end, we split 107 proteins into training and testing sets (with 79 and 28 proteins, respectively),
161 and perform evaluations based on the testing set. The main metric of evaluation is the ROC
162 curves and associated area under the curve (AUC). AUC is a standard metric for evaluating
163 binary classifiers based on the ROC curve of sensitivity and specificity. The ROC curves are also
164 used in other computational papers for variant prediction (see [47]).

165

166 **2. Sequences-based coevolution analysis and protein contact map**

167 **construction**

168 We perform coevolution analysis using three state-of-the-art methods: the RaptorX server
169 (<http://raptorx.uchicago.edu>), the DeepMetaPSICOV server (<http://bioinf.cs.ucl.ac.uk/psipred/>),
170 and the SPOT-Contact server (<https://sparks-lab.org/server/spot-contact/>). A sequence length limit
171 (500) is imposed by the capacity of coevolution analysis servers, and may be circumvented if
172 installing and running coevolution analysis locally.

173 These methods use multiple sequence alignments to compute the probability P_{ij} of residue
174 pair (i, j) forming spatial contact. Based on the matrix of predicted P_{ij} , a protein residue contact
175 map can be built with residues as nodes and pairwise contacts as edges weighted by P_{ij} . By default,
176 we do not apply any threshold cutoff to P_{ij} for defining contacts (unless networks with unweighted
177 edges are required by some node centrality algorithms in NetworkX, where we remove edges with
178 $P_{ij} < 0.1$, and set weight to 1 for the remaining edges).

179

180 **3. Network analysis of protein contact map**

181 By treating a protein contact map as a network of nodes and edges, we calculate various
182 node centrality scores to predict key residues as hotspots for disease mutations.

183 A simple score to measure node centrality is a weighted node degree that accounts for the
184 nearest neighbor interactions (denoted W_1):

$$185 \quad W_{1,i} = \sum_{k \neq i} P_{ik} \quad (1)$$

186 To include indirect couplings beyond the nearest neighbors, we calculate the node degree
187 based on the n'th power of the contact probability matrix (denoted W_n):

$$188 \quad W_{n,i} = \sum_{k \neq i} P_{ik} W_{n-1,k} = \sum_{k \neq i} P^n_{ik} \quad (2)$$

189 As n goes to infinity, W_n converges to the eigenvector of P matrix with the highest
190 eigenvalue λ_{\max} (denoted W_∞):

191 $PW_\infty = \lambda_{\max} W_\infty$ (3)

192 Among various W_n , W_2 can be interpreted as the node degrees of a new network based on
 193 a neighborhood similarity matrix S as follows (denoted W_s):

194 $S_{ij} = \sum_{k \neq i, j} P_{ik} P_{jk}, W_{s,i} = \sum_{k \neq i} S_{ik}$ (4)

195 In this study we use five network scores (W_1, W_2, W_3, W_∞ and W_s) as predictive features
 196 for node importance. Additionally, we exploit 13 network centrality metrics as calculated by the
 197 NetworkX package (see Table 1). To allow meaningful comparison of scores between proteins,
 198 the scores of each protein are sorted and their ranking percentiles are linearly transformed to
 199 values between 0 and 1.

200 **Table 1. Network centrality scores as implemented in the NetworkX package**

201 (see <https://networkx.org/documentation/stable/reference/algorithms/centrality.html>)

Symbol	Centrality name	Definition
C1	degree_centrality	Corresponding to W_1
C2	eigenvector_centrality	Corresponding to W_∞
C3	closeness_centrality	Closeness centrality of a node u is the reciprocal of the average shortest path distance to u over all $n-1$ reachable nodes.
C4	betweenness_centrality	Betweenness centrality of a node u is the sum of the fraction of all-pairs shortest paths that pass through u .
C5	current_flow_closeness_centrality	Current-flow closeness centrality is a variant of closeness centrality based on effective resistance between nodes in a network.
C6	current_flow_betweenness_centrality	Current-flow betweenness centrality is based on an electrical current model for information spreading.
C7	communicability_betweenness_centrality	Communicability betweenness centrality is based on the number of walks connecting every pair of nodes.
C8	load_centrality	Load centrality of a node u is the fraction of all shortest paths that pass through u .
C9	subgraph_centrality	Subgraph centrality of a node u is the sum of weighted closed walks of all lengths starting and ending at u .
C10	harmonic_centrality	Harmonic centrality of a node u is the sum of the reciprocal of the shortest path distances from all other nodes to u .

C11	second_order_centrality	Second order centrality of a node u is the standard deviation of the return times to u of a perpetual random walk on G .
C12	laplacian_centrality	Laplacian Centrality of a node u is measured by the drop in the Laplacian Energy after deleting u from the graph.
C13	katz_centrality_numpy	Katz centrality computes the centrality for a node u based on the centrality of its neighbors. It is a generalization of the eigenvector centrality.

202

203 4. Sequences-based GNM

204 For comparison, we implemented Bulter et al's sequence-based GNM [44]. The original
205 structure-based Gaussian network model (GNM) represents a protein structure as an elastically
206 connected network of residues to obtain the equilibrium fluctuations of residues. In the absence
207 of a structure, the sequence-based GNM (Seq-GNM) treats coevolving residue pairs as
208 contacting pairs.

209 To construct the Kirchhoff matrix (denoted K), each non-bonded residue pair is assigned
210 a value of -1 times its contact probability. The bonded residue pairs $(i, i+1)$ are assigned -1 to
211 enforce local chain connectivity. The diagonal elements of K are assigned so that the sum of each
212 row and column is zero:

$$213 K_{ij} = \begin{cases} -P_{ij} & i \neq j \\ \sum_{k \neq i} P_{ik} & i = j \end{cases} \quad (5)$$

214 The vibrational thermal fluctuations of residues are evaluated by inverting the Kirchhoff
215 matrix (or summing over the modes as weighted by $1/\lambda_m$). The per-residue mean-square
216 fluctuations (MSF), which are proportional to the crystallographic B factors, are given as
217 follows:

218
$$MSF_i \propto K_{ii}^{-1} = \sum_{m>0} \frac{V_{mi}^2}{\lambda_m} \quad (6)$$

219 where the eigen-decomposition of K gives eigenvectors V_m and eigenvalues λ_m that satisfy:

220
$$KV_m = \lambda_m V_m \quad (7)$$

221 Low-MSF residues correspond to rigid cores or hinges of dynamical importance [44].

222 As an alternative way to evaluate node importance using GNM, we perform a
 223 perturbation-based hotspot analysis as follows: For mode m , calculate how much its eigenvalue
 224 changes ($\delta\lambda_{m,i}$) in response to a perturbation at a chosen residue position i [23,24,51] (i.e., by
 225 uniformly weakening the contacts with residue i). Then compute $\delta\lambda_i = \sum_m \delta\lambda_{m,i}$ to assess the
 226 dynamic importance of this residue position [52]. High- $\delta\lambda_i$ residues correspond to sites highly
 227 sensitive to local perturbations that mimic mutations.

228 The above two GNM-based scores are combined with the other network scores for
 229 machine learning.

230

231 **5. ESM based variant prediction**

232 For comparison with our method, we use a deep-learning variant predictor based on a
 233 large protein language model (ESM). We downloaded and installed the ESM package and
 234 pretrained models from <https://github.com/facebookresearch/esm>. Since our dataset consists of
 235 known variants (from HumVar) and added non-conserved sites (with specific mutations
 236 unknown), we simulate the mutational effects on each site by introducing Alanine substitution if

237 the wildtype residue is not an Alanine and Glycine substitution otherwise. Then we process the
238 mutated sequence with 5 pretrained ESM models (esm1v_t33_650M_UR90S_1,
239 esm1v_t33_650M_UR90S_2, esm1v_t33_650M_UR90S_3, esm1v_t33_650M_UR90S_4, and
240 esm1v_t33_650M_UR90S_5), which predict the difference in the probability of observing the
241 wildtype residue and the mutant residue at a given site [48]. We record the predictions of five
242 ESM models as separate features to be optimally integrated via machine learning.

243

244 **6. AlphaFold for structural prediction**

245 We downloaded predicted structures for the 107 proteins from AlphaFold DB
246 (<https://alphafold.ebi.ac.uk/>). A residue contact probability matrix is constructed based on the
247 predicted structures as follows:

$$248 \quad P_{ij} = \frac{1}{1 + e^{d_{ij} - 10}} \quad (8)$$

249 where d_{ij} is the distance between residues i and j , and 10 \AA is used as a soft cutoff distance. We
250 then use this contact probability matrix to perform the same network analysis as in the
251 sequences-based method and for optimization with machine learning.

252

253 **7. Foldx for structural refinement and Alanine scanning analysis**

254 FoldX program [53] was downloaded from <https://foldxsuite.crg.eu/>. We use the
255 RepairPDB command to refine the AlphaFold-predicted models (by fixing bad torsion angles

256 and Van der Waals clashes). Then we use the AlaScan command to mutate each residue to Ala
257 and calculate the resulting changes in Gibbs free energies which are then used as a feature to
258 predict hotspots of disease mutations.

259

260 **8. Machine learning algorithms**

261 We use the following machine learning methods of the scikit-learn package
262 (<https://scikit-learn.org/stable/>) to learn optimal combinations of all features to predict if a given
263 site is deleterious or neutral mutation site:

264 Random Forest Classifier (RF) (`sklearn.ensemble.RandomForestClassifier`): A random
265 forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of
266 the dataset and uses averaging to improve the predictive accuracy and control over-fitting. We
267 tune the following hyper-parameters: `max_depth`, `n_estimators`, `max_features`.

268 Gradient Boosting Classifier (GB) (`sklearn.ensemble.GradientBoostingClassifier`): This
269 algorithm builds an additive model in a forward stage-wise fashion. In each stage a regression
270 tree is fit on the negative gradient of the loss function, e.g. binary log loss. We tune the following
271 hyper-parameters: `n_estimators`, `max_depth`, `max_features`.

272 Extreme Gradient Boosting Classifier (XGB) (`xgboost.XGBClassifier`): This algorithm is
273 an optimized distributed version of gradient boosting designed to be highly efficient, flexible and
274 portable. We tune the following hyper-parameters: `n_estimators`, `max_depth`, `reg_alpha`,
275 `reg_lambda`.

276 These three methods were chosen because they have performed successfully in machine
277 learning contests in Kaggle (see [https://www.packtpub.com/product/the-kaggle-
279 book/9781801817479](https://www.packtpub.com/product/the-kaggle-
278 book/9781801817479)). They are also relatively cheap to train and optimize compared with the
279 deep learning methods.

280 We use Optuna (<https://optuna.org/>) for hyper-parameter tuning of the above algorithms.
281 We have run Optuna multiple times to ensure the resulting best metric is reproducible.

282

283 **Results and discussion**

284

285 This study explores how to systematically utilize the coevolution information from multiple
286 sequence alignments to model and analyze a protein as a residue contact network beyond the
287 scope of GNM. To this end, we first use coevolution analysis to construct a protein residue
288 contact map with edges weighted by the predicted contact probability; then we exploit an array
289 of 20 network-based scores to assess the node importance as predictors for disease mutation
290 sites; finally we evaluate the predictive power of these scores individually and collectively (using
291 machine learning) based on a subset of 107 protein sequences and their variants from the
292 HumVar database. For comparison, we also evaluate alternative methods based on predicted
293 protein structures, a physics-based force field, and protein language models.

294 **1. Evaluation of individual network scores**

295 Based on the protein residue contact maps built from three coevolution analysis tools
296 (DeepMetaPSICOV, RaptorX, and SPOT-Contact), we applied network analysis to calculate 20
297 network scores (see Table 2), measuring node centrality using various different algorithms (see
298 Methods). These scores include simple weighted node degrees for n-hop nearest neighbors (see
299 Methods) and more sophisticated centrality metrics (see Table 1), along with 2 seq-GNM based
300 scores (MSF and $\delta\lambda$, see Methods). We evaluate the performance of each score using the AUC
301 of ROC for the testing set, which provides a balanced evaluation of sensitivity and specificity
302 (see Table 2). More specifically, we sort all testing-set variants by a particular score and predict a

303 variant deleterious/neutral if its score is above/below a cutoff value. This results in an ROC curve
 304 from which we have calculated AUC (see Table 2).

305 Overall, DeepMetaPSICOV (max AUC=0.80) and SPOT-Contact (max AUC=0.81)
 306 perform slightly better than RaptorX (max AUC=0.78). Interestingly, simple weighted node
 307 degrees (W_1 , W_2 , and W_3) perform better than those more complex centrality scores (see Table
 308 2). When computing node degrees, going beyond the nearest neighbors seems to improve the
 309 prediction slightly (see Table 2). Two GNM-based scores perform similarly but slightly worse
 310 than the weighted node degrees (see Table 2). Among those NetworkX-based scores (see Table
 311 1), C5, C11 and C12 outperform the others, while those betweenness-based scores (C4, C6, and
 312 C8) underperform (see Table 2). Therefore, the functional importance of a node/residue pertains
 313 more to its role as a highly-connected hub than as an information bottleneck of the shortest paths.

314 **Table 2. Evaluation of 20 network scores based on protein residue contact maps**
 315 **constructed from 3 coevolution analysis tools (DeepMetaPSICOV, RaptorX, and SPOT-**
 316 **Contact) and AlphaFold-predicted structures**

Score	AUC* of DeepMetaPSICOV	AUC* of RaptorX	AUC* of SPOT-Contact	AUC* of AlphaFold
C1	0.74	0.76	0.73	0.82
C2	0.73	0.74	0.76	0.77
C3	0.76	0.73	0.69	0.73
C4	0.64	0.54	0.60	0.58
C5	0.78	0.76	0.79	0.80
C6	0.63	0.58	0.67	0.60
C7	0.75	0.61	0.72	0.74
C8	0.64	0.54	0.60	0.58
C9	0.77	0.76	0.74	0.78
C10	0.75	0.73	0.68	0.75
C11	0.79	0.76	0.77	0.80
C12	0.77	0.77	0.79	0.83
C13	0.73	0.73	0.76	0.76
$\delta\lambda$	0.79	0.76	0.78	0.83
MSF	0.79	0.76	0.78	0.80

W_1	0.79	0.77	0.78	0.83
W_2	0.80	0.78	0.80	0.83
W_3	0.80	0.78	0.81	0.82
W_∞	0.80	0.74	0.79	0.77
W_s	0.80	0.78	0.80	0.83
FoldX				0.68

317 * The AUC is calculated based on the ROC for all variants of the 28 testing set proteins.

318 Alternatively, we also calculated AUCs based on the ROCs of individual proteins and their
319 summary statistics (see Table S1).

320 For comparison with alternative methods, we evaluated the performance of variant
321 prediction by five pre-trained protein language models (ESM, see Methods), and the resulting
322 AUC varies between 0.79 and 0.81, which are comparable to the network scores (see Table 2).
323 For further comparison with structures-based methods, we also performed network analysis
324 based on protein structures as predicted by AlphaFold (see Methods). Overall, the structures-
325 based scores (max AUC=0.83) perform slightly better than the sequences-based scores. This may
326 be partly due to the structure-based contact maps (see Eq. 8) being more sharply defined than the
327 fuzzier contact-probability-based contact maps. Notably, when structural information is used, our
328 network analysis performs significantly better than a physics-based force field (FoldX) with
329 AUC=0.68. Taken together, these findings support the usefulness of individual sequences-based
330 network centrality scores in predicting important residues on par with alternative more
331 sophisticated methods.

332 To further understand the different accuracies of the above scores, we explore the
333 relationships between them by evaluating the pairwise Pearson correlations (PC) (see Table 3).
334 W_1 , W_2 , W_3 , W_∞ , W_s , MSF and $\delta\lambda$ are highly correlated (with $PC \geq 0.93$ for DeepMetaPSICOV,
335 $PC \geq 0.84$ for SPOT-Contact, $PC \geq 0.86$ for AlphaFold), although their correlations are somewhat
336 weaker for RaptorX. Among the NetworkX-based scores (see Table 1), C5, C11 and C12 are

337 also highly correlated with the above scores. Such strong correlations support the attribution of
 338 higher AUC of these scores (see Table 2) to their capturing the same essential features (i.e. high
 339 node degrees) of those important nodes. In contrast, the betweenness-based scores (C4, C6, and
 340 C8) do not correlate well with the above scores, which is consistent with their lower AUC (see
 341 Table 2).

342 **Table 3. Pearson correlations between network scores (row 1, 2, 3 and 4 correspond to**
 343 **results of DeepMetaPSICOV, SPOT-Contact, RaptorX and AlphaFold, respectively)**

	C1	C2	C3	C4	C5	C6	C7	C8	C9	C10	C11	C12	C13	W ₁	W ₂	W ₃	W _∞	W _s	MSF	δλ
C1	1.00	0.57	0.74	0.63	0.72	0.40	0.82	0.63	0.84	0.79	0.88	0.84	0.54	0.87	0.85	0.83	0.80	0.85	0.87	0.87
C2	0.57	1.00	0.43	0.21	0.75	0.18	0.47	0.21	0.66	0.42	0.67	0.76	0.98	0.65	0.72	0.75	0.80	0.71	0.65	0.65
C3	0.74	0.43	1.00	0.64	0.63	0.36	0.72	0.64	0.73	0.96	0.81	0.57	0.43	0.63	0.63	0.63	0.61	0.64	0.64	0.63
C4	0.63	0.21	0.64	1.00	0.34	0.42	0.64	1.00	0.40	0.65	0.53	0.37	0.21	0.46	0.42	0.40	0.38	0.43	0.47	0.46
C5	0.72	0.75	0.63	0.34	1.00	0.46	0.67	0.34	0.74	0.58	0.91	0.86	0.76	0.80	0.83	0.84	0.84	0.83	0.80	0.80
C6	0.40	0.18	0.36	0.42	0.46	1.00	0.63	0.42	0.21	0.29	0.45	0.33	0.20	0.45	0.41	0.38	0.34	0.42	0.46	0.45
C7	0.82	0.47	0.72	0.64	0.67	0.63	1.00	0.64	0.73	0.71	0.80	0.68	0.46	0.76	0.75	0.74	0.71	0.76	0.77	0.76
C8	0.63	0.21	0.64	1.00	0.34	0.42	0.64	1.00	0.40	0.65	0.53	0.37	0.21	0.46	0.42	0.40	0.37	0.43	0.47	0.46
C9	0.84	0.66	0.73	0.40	0.74	0.21	0.73	0.40	1.00	0.79	0.85	0.79	0.63	0.77	0.81	0.82	0.82	0.81	0.78	0.77
C10	0.79	0.42	0.96	0.65	0.58	0.29	0.71	0.65	0.79	1.00	0.79	0.59	0.41	0.65	0.65	0.64	0.63	0.66	0.66	0.65
C11	0.88	0.67	0.81	0.53	0.91	0.45	0.80	0.53	0.85	0.79	1.00	0.85	0.68	0.84	0.85	0.85	0.84	0.86	0.85	0.84
C12	0.84	0.76	0.57	0.37	0.86	0.33	0.68	0.37	0.79	0.59	0.85	1.00	0.74	0.91	0.94	0.94	0.92	0.92	0.90	0.91
C13	0.54	0.98	0.43	0.21	0.76	0.20	0.46	0.21	0.63	0.41	0.68	0.74	1.00	0.63	0.69	0.73	0.78	0.69	0.63	0.63
W ₁	0.87	0.65	0.63	0.46	0.80	0.45	0.76	0.46	0.77	0.65	0.84	0.91	0.63	1.00	0.98	0.97	0.93	0.98	1.00	1.00
W ₂	0.85	0.72	0.63	0.42	0.83	0.41	0.75	0.42	0.81	0.65	0.85	0.94	0.69	0.98	1.00	1.00	0.97	1.00	0.99	0.98
W ₃	0.83	0.75	0.63	0.40	0.84	0.38	0.74	0.40	0.82	0.64	0.85	0.94	0.73	0.97	1.00	1.00	0.99	0.99	0.97	0.97
W _∞	0.80	0.80	0.61	0.38	0.84	0.34	0.71	0.37	0.82	0.63	0.84	0.92	0.78	0.93	0.97	0.99	1.00	0.97	0.93	0.93
W _s	0.85	0.71	0.64	0.43	0.83	0.42	0.76	0.43	0.81	0.66	0.86	0.92	0.69	0.98	1.00	0.99	0.97	1.00	0.99	0.98
MSF	0.87	0.65	0.64	0.47	0.80	0.46	0.77	0.47	0.78	0.66	0.85	0.90	0.63	1.00	0.99	0.97	0.93	0.99	1.00	1.00
C1	1.00	0.58	0.68	0.55	0.73	0.53	0.85	0.55	0.86	0.73	0.90	0.74	0.57	0.80	0.78	0.77	0.67	0.80	0.82	0.80
C2	0.58	1.00	0.44	0.20	0.88	0.40	0.51	0.20	0.61	0.38	0.74	0.82	0.99	0.70	0.76	0.79	0.88	0.74	0.68	0.70
C3	0.68	0.44	1.00	0.61	0.55	0.38	0.69	0.61	0.73	0.96	0.78	0.42	0.45	0.50	0.52	0.53	0.52	0.55	0.56	0.50
C4	0.55	0.20	0.61	1.00	0.29	0.46	0.63	1.00	0.37	0.62	0.49	0.26	0.20	0.37	0.35	0.34	0.30	0.37	0.40	0.37
C5	0.73	0.88	0.55	0.29	1.00	0.57	0.67	0.29	0.71	0.49	0.88	0.91	0.87	0.84	0.88	0.89	0.87	0.87	0.82	0.84
C6	0.53	0.40	0.38	0.46	0.57	1.00	0.73	0.46	0.33	0.31	0.54	0.55	0.40	0.64	0.61	0.59	0.49	0.62	0.66	0.64
C7	0.85	0.51	0.69	0.63	0.67	0.73	1.00	0.63	0.75	0.68	0.82	0.66	0.51	0.76	0.75	0.73	0.64	0.77	0.80	0.76
C8	0.55	0.20	0.61	1.00	0.29	0.46	0.63	1.00	0.37	0.62	0.49	0.26	0.20	0.37	0.35	0.34	0.30	0.37	0.40	0.37
C9	0.86	0.61	0.73	0.37	0.71	0.33	0.75	0.37	1.00	0.77	0.87	0.66	0.59	0.68	0.71	0.72	0.68	0.73	0.71	0.68
C10	0.73	0.38	0.96	0.62	0.49	0.31	0.68	0.62	0.77	1.00	0.76	0.38	0.38	0.47	0.49	0.49	0.47	0.52	0.54	0.47
C11	0.90	0.74	0.78	0.49	0.88	0.54	0.82	0.49	0.87	0.76	1.00	0.79	0.73	0.80	0.82	0.83	0.78	0.83	0.82	0.80
C12	0.74	0.82	0.42	0.26	0.91	0.55	0.66	0.26	0.66	0.38	0.79	1.00	0.80	0.91	0.93	0.92	0.83	0.90	0.86	0.91
C13	0.57	0.99	0.45	0.20	0.87	0.40	0.51	0.20	0.59	0.38	0.73	0.80	1.00	0.69	0.75	0.78	0.88	0.73	0.67	0.69
W ₁	0.80	0.70	0.50	0.37	0.84	0.64	0.76	0.37	0.68	0.47	0.80	0.91	0.69	1.00	0.98	0.97	0.84	0.98	0.98	1.00
W ₂	0.78	0.76	0.52	0.35	0.88	0.61	0.75	0.35	0.71	0.49	0.82	0.93	0.75	0.98	1.00	1.00	0.90	1.00	0.97	0.98
W ₃	0.77	0.79	0.53	0.34	0.89	0.59	0.73	0.34	0.72	0.49	0.83	0.92	0.78	0.97	1.00	1.00	0.92	0.99	0.96	0.97
W _∞	0.67	0.88	0.52	0.30	0.87	0.49	0.64	0.30	0.68	0.47	0.78	0.83	0.88	0.84	0.90	0.92	1.00	0.90	0.86	0.84

W_s	0.80 0.74 0.55 0.37 0.87 0.62 0.77 0.37 0.73 0.52 0.83 0.90 0.73 0.98 1.00 0.99 0.90 1.00 0.98 0.98
MSF	0.82 0.68 0.56 0.40 0.82 0.66 0.80 0.40 0.71 0.54 0.82 0.86 0.67 0.98 0.97 0.96 0.86 0.98 1.00 0.98
C1	1.00 0.63 0.58 0.22 0.64 0.22 0.29 0.22 0.75 0.65 0.71 0.78 0.55 0.88 0.87 0.85 0.58 0.87 0.68 0.73
C2	0.63 1.00 0.64 0.15 0.82 0.23 0.24 0.15 0.66 0.66 0.78 0.70 0.93 0.60 0.65 0.67 0.82 0.64 0.70 0.55
C3	0.58 0.64 1.00 0.44 0.79 0.35 0.39 0.44 0.64 0.93 0.88 0.51 0.68 0.48 0.50 0.50 0.64 0.50 0.69 0.47
C4	0.22 0.15 0.44 1.00 0.32 0.70 0.67 1.00 0.11 0.38 0.39 0.06 0.22 0.18 0.14 0.14 0.30 0.15 0.37 0.17
C5	0.64 0.82 0.79 0.32 1.00 0.45 0.43 0.32 0.65 0.72 0.96 0.71 0.86 0.63 0.65 0.66 0.79 0.65 0.81 0.62
C6	0.22 0.23 0.35 0.70 0.45 1.00 0.70 0.70 0.12 0.25 0.45 0.17 0.29 0.32 0.28 0.26 0.40 0.28 0.51 0.31
C7	0.29 0.24 0.39 0.67 0.43 0.70 1.00 0.67 0.34 0.34 0.48 0.28 0.29 0.30 0.28 0.27 0.40 0.29 0.49 0.39
C8	0.22 0.15 0.44 1.00 0.32 0.70 0.67 1.00 0.11 0.38 0.39 0.06 0.22 0.18 0.14 0.14 0.30 0.15 0.37 0.17
C9	0.75 0.66 0.64 0.11 0.65 0.12 0.34 0.11 1.00 0.71 0.72 0.71 0.61 0.66 0.69 0.69 0.62 0.69 0.66 0.64
C10	0.65 0.66 0.93 0.38 0.72 0.25 0.34 0.38 0.71 1.00 0.83 0.55 0.67 0.52 0.53 0.53 0.66 0.54 0.72 0.51
C11	0.71 0.78 0.88 0.39 0.96 0.45 0.48 0.39 0.72 0.83 1.00 0.68 0.82 0.64 0.65 0.65 0.77 0.65 0.83 0.63
C12	0.78 0.70 0.51 0.06 0.71 0.17 0.28 0.06 0.71 0.55 0.68 1.00 0.64 0.76 0.77 0.77 0.63 0.77 0.67 0.70
C13	0.55 0.93 0.68 0.22 0.86 0.29 0.29 0.22 0.61 0.67 0.82 0.64 1.00 0.53 0.58 0.61 0.85 0.58 0.72 0.52
W_1	0.88 0.60 0.48 0.18 0.63 0.32 0.30 0.18 0.66 0.52 0.64 0.76 0.53 1.00 0.98 0.97 0.69 0.98 0.81 0.85
W_2	0.87 0.65 0.50 0.14 0.65 0.28 0.28 0.14 0.69 0.53 0.65 0.77 0.58 0.98 1.00 0.99 0.73 1.00 0.80 0.83
W_3	0.85 0.67 0.50 0.14 0.66 0.26 0.27 0.14 0.69 0.53 0.65 0.77 0.61 0.97 0.99 1.00 0.75 0.99 0.80 0.82
W_∞	0.58 0.82 0.64 0.30 0.79 0.40 0.40 0.30 0.62 0.66 0.77 0.63 0.85 0.69 0.73 0.75 1.00 0.73 0.88 0.67
W_s	0.87 0.64 0.50 0.15 0.65 0.28 0.29 0.15 0.69 0.54 0.65 0.77 0.58 0.98 1.00 0.99 0.73 1.00 0.81 0.83
MSF	0.68 0.70 0.69 0.37 0.81 0.51 0.49 0.37 0.66 0.72 0.83 0.67 0.72 0.81 0.80 0.80 0.88 0.81 1.00 0.79
C1	1.00 0.87 0.77 0.38 0.95 0.42 0.82 0.38 0.90 0.81 0.96 0.98 0.85 0.97 0.97 0.96 0.86 0.97 0.95 0.97
C2	0.87 1.00 0.78 0.26 0.91 0.32 0.81 0.26 0.98 0.79 0.90 0.89 0.99 0.87 0.91 0.93 0.99 0.91 0.91 0.87
C3	0.77 0.78 1.00 0.52 0.79 0.37 0.74 0.52 0.82 0.97 0.83 0.72 0.78 0.71 0.74 0.75 0.78 0.74 0.79 0.71
C4	0.38 0.26 0.52 1.00 0.33 0.61 0.55 1.00 0.30 0.51 0.39 0.30 0.27 0.30 0.30 0.29 0.27 0.30 0.33 0.30
C5	0.95 0.91 0.79 0.33 1.00 0.46 0.85 0.33 0.92 0.79 0.99 0.96 0.91 0.95 0.96 0.96 0.91 0.96 1.00 0.95
C6	0.42 0.32 0.37 0.61 0.46 1.00 0.67 0.61 0.33 0.31 0.46 0.40 0.33 0.41 0.39 0.37 0.33 0.39 0.46 0.41
C7	0.82 0.81 0.74 0.55 0.85 0.67 1.00 0.55 0.83 0.72 0.85 0.81 0.81 0.80 0.82 0.83 0.81 0.82 0.84 0.80
C8	0.38 0.26 0.52 1.00 0.33 0.61 0.55 1.00 0.30 0.51 0.39 0.30 0.27 0.30 0.30 0.29 0.27 0.30 0.33 0.30
C9	0.90 0.98 0.82 0.30 0.92 0.33 0.83 0.30 1.00 0.84 0.93 0.91 0.96 0.89 0.93 0.95 0.97 0.93 0.92 0.89
C10	0.81 0.79 0.97 0.51 0.79 0.31 0.72 0.51 0.84 1.00 0.84 0.75 0.79 0.73 0.77 0.78 0.79 0.77 0.80 0.73
C11	0.96 0.90 0.83 0.39 0.99 0.46 0.85 0.39 0.93 0.84 1.00 0.95 0.90 0.94 0.95 0.95 0.90 0.95 0.99 0.94
C12	0.98 0.89 0.72 0.30 0.96 0.40 0.81 0.30 0.91 0.75 0.95 1.00 0.88 1.00 1.00 0.99 0.89 0.99 0.96 1.00
C13	0.85 0.99 0.78 0.27 0.91 0.33 0.81 0.27 0.96 0.79 0.90 0.88 1.00 0.85 0.90 0.92 0.99 0.90 0.91 0.85
W_1	0.97 0.87 0.71 0.30 0.95 0.41 0.80 0.30 0.89 0.73 0.94 1.00 0.85 1.00 0.99 0.97 0.86 0.98 0.95 1.00
W_2	0.97 0.91 0.74 0.30 0.96 0.39 0.82 0.30 0.93 0.77 0.95 1.00 0.90 0.99 1.00 1.00 0.91 1.00 0.96 0.99
W_3	0.96 0.93 0.75 0.29 0.96 0.37 0.83 0.29 0.95 0.78 0.95 0.99 0.92 0.97 1.00 1.00 0.93 1.00 0.96 0.97
W_∞	0.86 0.99 0.78 0.27 0.91 0.33 0.81 0.27 0.97 0.79 0.90 0.89 0.99 0.86 0.91 0.93 1.00 0.91 0.91 0.86
W_s	0.97 0.91 0.74 0.30 0.96 0.39 0.82 0.30 0.93 0.77 0.95 0.99 0.90 0.98 1.00 1.00 0.91 1.00 0.96 0.98
MSF	0.95 0.91 0.79 0.33 1.00 0.46 0.84 0.33 0.92 0.80 0.99 0.96 0.91 0.95 0.96 0.96 0.91 0.96 1.00 0.95

344

345

346 In summary, by evaluating 20 network scores individually, we have found a wide range
347 of performance with AUC varying from 0.54 to 0.81 (see Table 2). The top-performing scores
348 seem to correlate strongly with each other, so they must have captured a common aspect of node
349 centrality that is relevant to functional importance (e.g. high local connectivity instead of high

350 betweenness). Interestingly, the two GNM-based scores, despite measuring distinct dynamic
351 properties (MSF measures thermal fluctuations while $\delta\lambda$ measures sensitivity to local
352 perturbations), are also strongly correlated with each other and those degree-based network
353 scores. Therefore, to speed up the variant prediction workflow we only need to compute those
354 simpler weighted node degrees as features without significantly losing accuracy.

355

356 **2. Combining all network scores to predict variant hotspots by** 357 **machine learning**

358 To optimize the predictive power of the above network-based scores based on three
359 coevolution analysis methods (or AlphaFold), we have employed machine learning algorithms
360 (see Methods) to take them as input features, train a binary classifier which predicts if a residue
361 position is linked to neutral or deleterious variants (using first 79 proteins as training set), and
362 then test its prediction using the remaining 28 proteins as testing set. We use the AUC of ROC as
363 the metric for assessing the prediction quality of the trained classifier.

364 To evaluate the protein residue contact maps constructed by each method, we combine all
365 network scores based on the contact maps predicted by the same method (see Table 2) for
366 machine learning. The resulting AUC of each coevolution analysis method (DeepMetaPSICOV,
367 RaptorX, and SPOT-Contact) is 0.81, 0.80, and 0.82, respectively (see Table 4), which are
368 slightly better than the best AUC of individual scores (0.78~0.81, see Table 2). The lack of
369 substantial improvement may be due to high correlations among the scores (see Table 3) which
370 could reduce the effectiveness of ensemble learning. For comparison, we also trained and tested
371 classifiers using the AlphaFold-predicted contact maps, and alternative classifiers based on

372 protein language models (see Methods). Both alternative methods give comparable yet slightly
 373 better AUC (0.83). Similar to our finding, Butler et al reported AUC of 0.81 after combining the
 374 B-factors of Seq-GNM with evolutionary features [44].

375 **Table 4. Evaluation of classifiers trained by 3 machine learning algorithms (RF, GB and**
 376 **XGB, see Methods) based on the protein residue contact maps constructed from 3**
 377 **coevolution analysis tools (DeepMetaPSICOV, RaptorX, and SPOT-Contact), AlphaFold-**
 378 **predicted structures, and protein language models (ESM).**

Sources of input features	AUC of RF	AUC of GB	AUC of XGB
DeepMetaPSICOV	0.81	0.81	0.81
RaptorX	0.80	0.80	0.80
SPOT-Contact	0.82	0.82	0.82
AlphaFold	0.83	0.83	0.83
ESM	0.83	0.83	0.83
All 3 coevolution methods	0.84	0.84	0.84
All 3 coevolution methods (w/o C1-C13)	0.82	0.82	0.83
All 3 coevolution methods and ESM	0.89	0.89	0.89
AlphaFold and ESM	0.88	0.88	0.88

379

380 To further boost the prediction performance, we have sought to combine the network
 381 scores of all three coevolution analysis methods for machine learning, resulting in better AUC
 382 (0.84) which slightly outperform both AlphaFold and ESM (0.83). To assess the added value of
 383 including 13 NetworkX-based centrality scores (see Table 1), we have performed an ablation
 384 study that excludes them in machine learning, and found slightly lower AUC (0.82~0.83). So it
 385 is possible to speed up the calculation without significantly reducing accuracy. Taken together,
 386 our findings support the power of combining an array of different network scores from different
 387 coevolution analysis tools to optimize the prediction in the framework of ensemble learning.

388 To further explore how well our method complements alternative methods, we have
389 combined all the network scores with the ESM scores in machine learning. Encouragingly, we
390 have obtained markedly improved AUC (0.89), which is comparable to machine learning that
391 combines the AlphaFold-based network scores with the ESM scores (AUC=0.88).

392 For comparison with other studies, Butler et al showed that Seq-GNM combined with
393 evolutionary parameters attained a sensitivity of 0.84 and a specificity of 0.66 [44]. PolyPhen-2
394 achieved a sensitivity of 0.73 and a specificity of 0.8 on the HumVar datasets [47]. While using
395 different training and testing datasets, we have attained competitive results with a sensitivity of
396 0.82 and a specificity of 0.80 (using all the network scores from three coevolution analysis tools
397 and the ESM scores). For more direct comparison, we also evaluated PolyPhen-2 based on the
398 same 28 testing-set proteins and their variants, and obtained an AUC of 0.85, which is close to
399 our method (see Table 4). However, this metric is likely positively biased since PolyPhen-2 has
400 been trained on the HumVar dataset.

401 In summary, via extensive machine learning, we have demonstrated the power of using
402 an ensemble of sequences-based network scores calculated by different co-evolution analysis
403 tools to accurately predict deleterious mutation sites. Although some network scores are highly
404 correlated (see Table 3) and they vary widely in accuracy (see Table 2), these scores seem to be
405 sufficiently diverse to allow effective ensemble learning when combined.

406

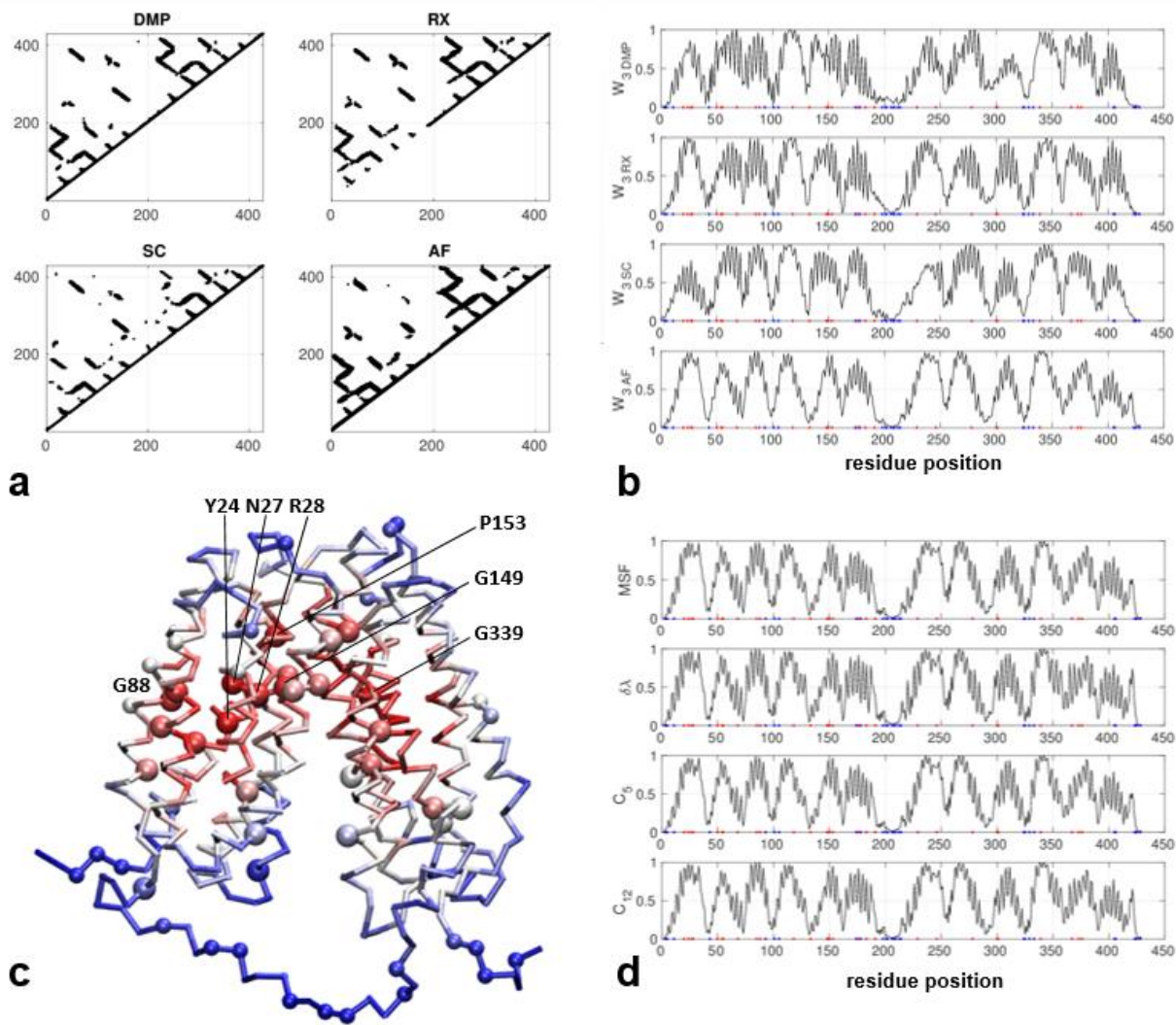
407 **3. Case studies:**

408 To illustrate the biomedical significance of our predictions of variant sites with network scores,
409 we discuss in details the following four proteins from our dataset.

410 **Glucose-6-phosphate exchanger** (Uniprot id: O43826): As an inorganic phosphate and
411 glucose-6-phosphate antiporter, it transports cytoplasmic glucose-6-phosphate into the lumen of
412 the endoplasmic reticulum and translocates inorganic phosphate in the opposite direction. Being
413 involved in glucose production through glycogenolysis and gluconeogenesis, it plays a central
414 role in homeostatic regulation of blood glucose levels. It is linked to diseases like congenital
415 disorder of glycosylation and glycogen storage disease (see
416 <https://www.uniprot.org/uniprotkb/O43826/entry#function>).

417 The AlphaFold-predicted structure forms a dimer of transmembrane helical domains with
418 most deleterious mutation sites concentrated inside the central core while those non-conserved
419 residues (i.e. neutral mutation sites) are mostly located on the periphery (see Fig 1c). The contact
420 maps predicted by three coevolution analysis tools all agree well with the contact map based on
421 the AlphaFold structure (see Fig 1a) (except that RaptorX omitted many local contacts in
422 residues 1-200). As a result, the network centrality scores (W_3) also agree well between these
423 methods (see Fig 1b), although the coevolution-based network scores are generally noisier (with
424 more spikes) than the structure-based scores (see Fig 1b). Different network scores calculated
425 from the same contact map are also highly similar (see Fig 1d) despite being based on different
426 algorithms. For example, scores of $\delta\lambda$ and MSF agree very well (see Fig 1d). Encouragingly,
427 those residues identified with high network scores are primarily within the central core (inside
428 each domain or in the inter-domain hinge region), thus overlapping with most deleterious
429 mutations (see Fig 1c). Among those top-10% hotspot residues (see Fig 1c), mutations Y24H,
430 N27K, R28H, G88D, G149E, P153L, and G339C were implicated in causing glycogen storage
431 disease [54] . Two of these mutations (R28H and G149E) were found to exhibit undetectable

432 microsomal glucose-6-phosphate transport activity in transient expression studies[55], thus
433 confirming their functional importance.



434
435 **Figure 1. Results for Glucose-6-phosphate exchanger (Uniprot id: O43826):** (a) Four contact
436 maps constructed from coevolution analysis by DeepMetaPSICOV (DMP), RaptorX (RX),
437 SPOT-Contact (SC), and the predicted structure by AlphaFold (AF) (only those contacts with
438 probability >0.1 are shown). (b) W_3 scores for all residue positions based on the contact maps in
439 (a), where red and blue dots mark residues with deleterious and neutral mutations, respectively.
440 (c) Predicted structure by AlphaFold as colored by W_3 scores (red/blue for high/low values),

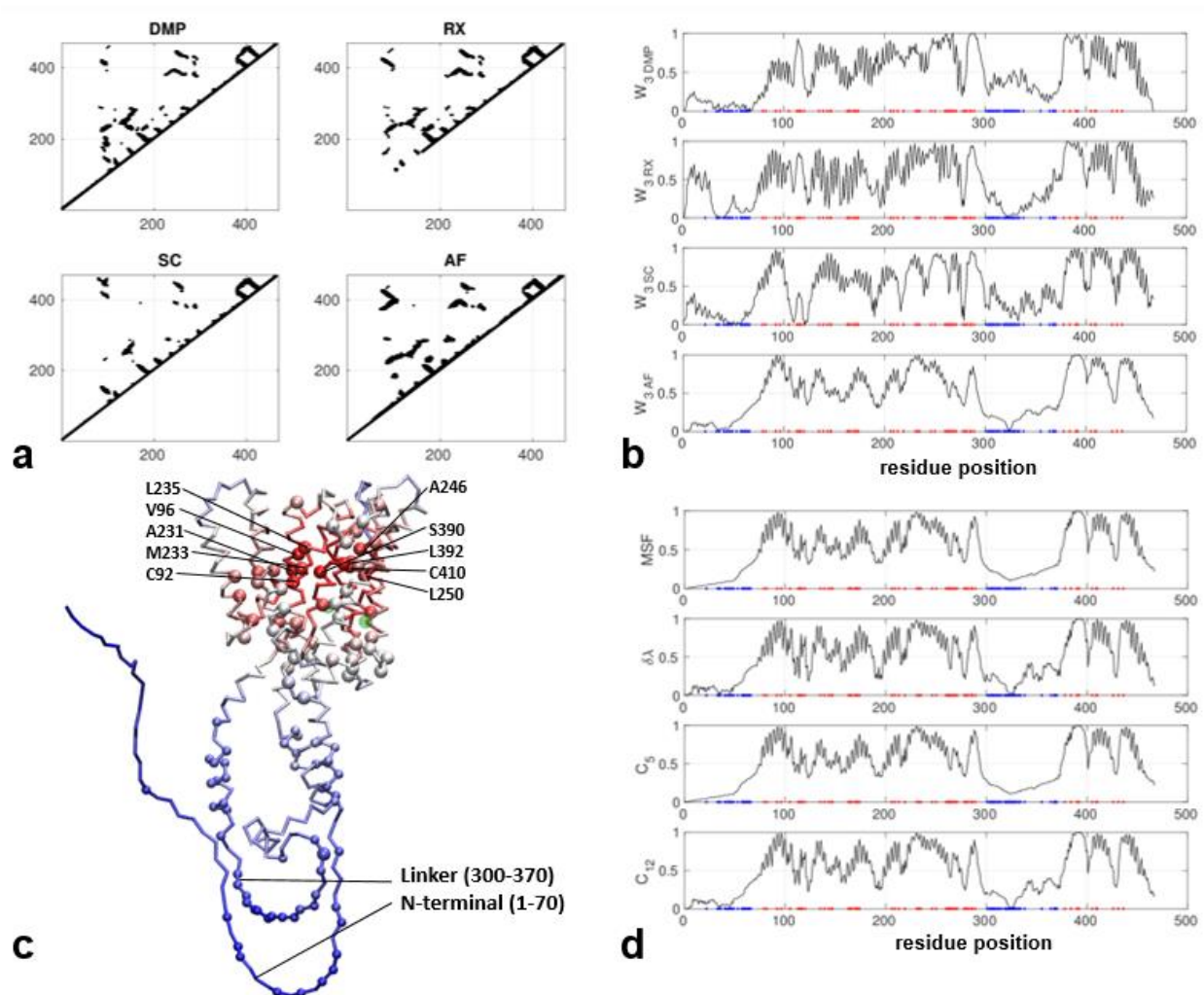
441 where residues with deleterious and neutral mutations are shown as large and small balls,
442 respectively (G20, Y24, N27, R28, G50, S54, S55, G68, L85, G88, W118, Q133, A148, G149,
443 G150, P153, C176, C183, P191, L229, W246, I278, R300, H301, G339, A367, A373, G376, see
444 <https://www.uniprot.org/uniprotkb/O43826/variant-viewer>). (d) Four other network scores (MSF,
445 $\delta\lambda$, C5 and C12) for all residue positions based on the contact maps in (a).

446

447 **Presenilin-1** (Uniprot id: P49768): As the catalytic subunit of the gamma-secretase complex, it
448 catalyzes the intramembrane cleavage of integral membrane proteins such as Notch receptors. It
449 is involved in various diseases including a familial early-onset form of Alzheimer disease and a
450 form of frontotemporal dementia (see [https://www.uniprot.org/uniprotkb/](https://www.uniprot.org/uniprotkb/P49768/entry#function)
451 [P49768/entry#function](https://www.uniprot.org/uniprotkb/P49768/entry#function)).

452 The AlphaFold-predicted structure consists of two closely packed helical domains with
453 most deleterious mutations clustered inside the core domain while the non-conserved residues
454 are mostly located on the N-terminal loop (residues 1-70) and the inter-domain linker (residues
455 300-370) (see Fig 2c). The active site [56] (D257 and D385) is also located in the core domain
456 (colored green in Fig 2c). The contact maps predicted by three coevolution analysis methods all
457 resemble the contact map based on the predicted structure (see Fig 2a) (except that RaptorX
458 omitted local contacts in residues 1-100). As a result, the network scores agree well between
459 them in the helical domains (see Fig 2c), but with more differences in the flexible regions
460 (residues 1-70 and 300-370). Reassuringly, those residues identified by high network scores are
461 primarily clustered within the central core overlapping with most deleterious mutations, while
462 the flexible N-terminal and linker feature low scores consistent with low sequence conservation
463 (see Fig 2c). Among those top 10% hotspot residues (see Fig 2c), mutations at C92, V96, A231,

464 M233, L235, A246, L250, S390, L392, and C410 were found to cause loss of function and
 465 altered amyloid-beta production [57] : C92S led to loss of protease function and increased
 466 Abeta42 levels. V96F caused loss of protease activity. A231T/V and M233T led to decreased
 467 protease activity, altered amyloid-beta production and increased amyloid-beta 42/amyloid-beta
 468 40 ratio. L235P/R and S390I abolished protease activity. A246E and L250S abolished protease
 469 activity and increased amyloid-beta 42/amyloid-beta 40 ratio. L392V resulted in reduced
 470 proteolysis, altered amyloid-beta production and increased amyloid-beta 42/amyloid-beta 40
 471 ratio. C410I reduced proteolysis. Since most of these residues are not near the active site, their
 472 effects on protease activity are likely allosteric.



473

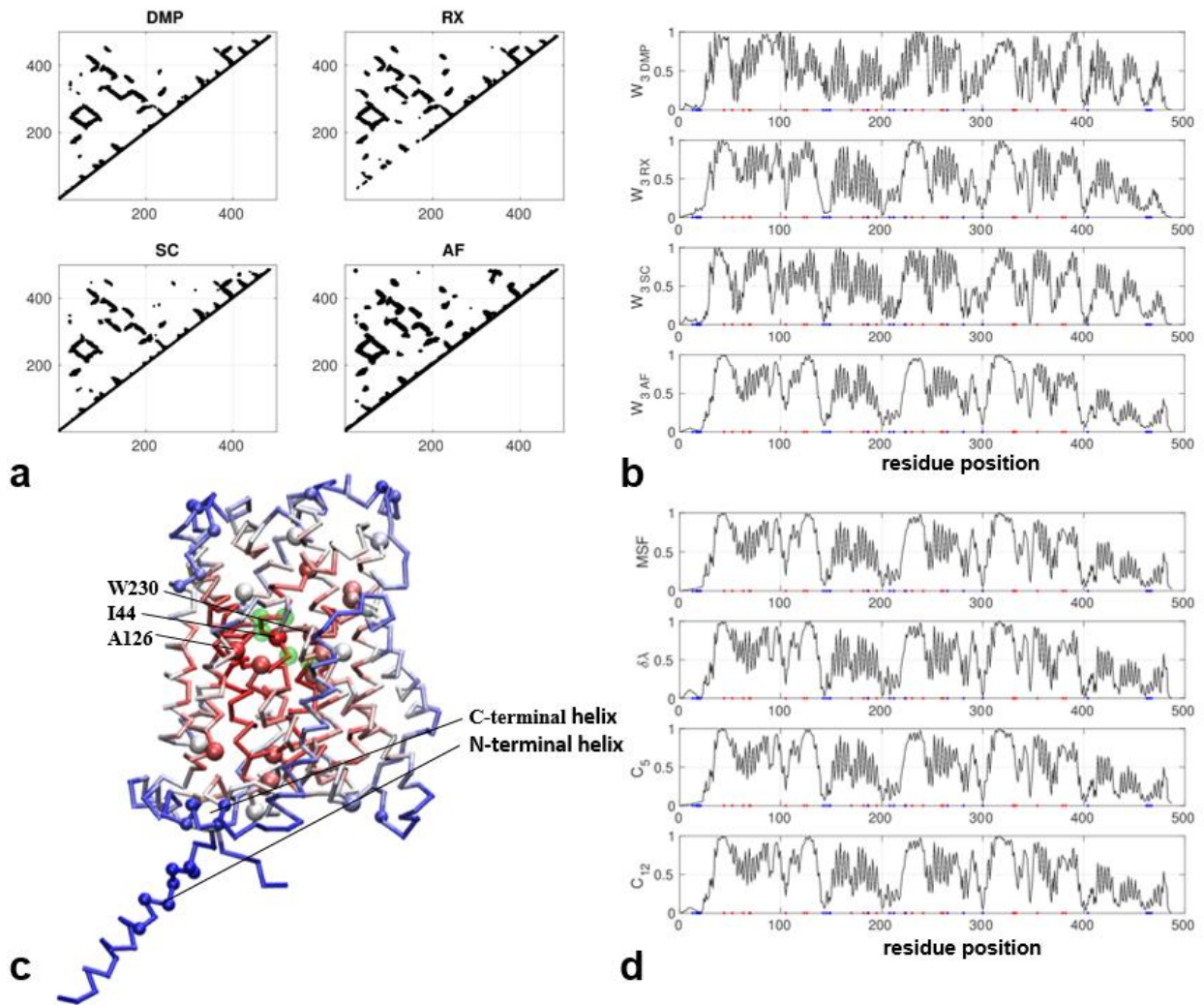
474 **Figure 2. Results for Presenilin-1 (Uniprot id: P49768):** (a) Four contact maps constructed
475 from coevolution analysis by DeepMetaPSICOV (DMP), RaptorX (RX), SPOT-Contact (SC),
476 and the predicted structure by AlphaFold (AF) (only those contacts with probability >0.1 are
477 shown). (b) W_3 scores for all residue positions based on the contact maps in (a), where red and
478 blue dots mark residues with deleterious and neutral mutations, respectively. (c) Predicted
479 structure by AlphaFold as colored by W_3 scores (red/blue for high/low values), where residues
480 with deleterious and neutral mutations are shown as large and small balls, respectively (A79,
481 V82, C92, V96, F105, L113, Y115, T116, P117, E120, N135, M139, I143, M146, T147, H163,
482 W165, L166, S169, L171, L173, L174, G206, G209, I213, L219, A231, M233, L235, A246,
483 L250, A260, L262, C263, P264, G266, P267, R269, L271, R278, E280, L282, A285, L286,
484 S289, D333, G378, G384, S390, L392, N405, A409, C410, A426, A431, P436, see
485 <https://www.uniprot.org/uniprotkb/P49768/variant-viewer>), and active-site residues are colored
486 in green. (d) Four other network scores (MSF, $\delta\lambda$, C5 and C12) for all residue positions based on
487 the contact maps in (a).

488

489 **b(0,+)-type amino acid transporter 1** (Uniprot id: P82251): It forms a functional
490 transporter complex that mediates the electrogenic exchange between cationic amino acids and
491 neutral amino acids. Its dysfunction is linked to Cystinuria, an autosomal disorder characterized
492 by impaired epithelial cell transport of cystine and dibasic amino acids in the proximal renal
493 tubule and gastrointestinal tract (see <https://www.uniprot.org/uniprotkb/P82251/entry#function>).

494 The AlphaFold-predicted structure consists of a helical domain with deleterious
495 mutations concentrating inside the core domain while those non-conserved residues are mostly
496 located on the domain periphery (N-terminal and C-terminal helices) (see Fig 3c). The active site

497 consists of residues 43-47 and 233 and is also located in the core domain (colored green in Fig
498 3c). The contact maps predicted by three coevolution analysis tools are all similar to the contact
499 map based on the AlphaFold structure (see Fig 3a) (except that RaptorX omitted some local
500 contacts in residues 1-200). As a result, the network scores agree well between these methods
501 (see Fig 3b). Reassuringly, those residues identified with high network scores are primarily
502 within the central core and overlap with most deleterious mutations, while the peripheral regions
503 feature low scores consistent with low sequence conservation. Among those top-10% hotspot
504 residues (see Fig 3c), mutations I44T, A126T, and W230R were implicated in Cystinuria. *In*
505 *vitro* measurements showed W230R has almost no transport activity, and it was proposed that
506 W230 serves as a gate between two substrate-binding pockets and undergoes conformational
507 changes to enable amino acid transport [58] . Although the A126T mutation is mildly
508 dysfunctional [59], it is notable among a cluster of conserved residues with small sidechains in
509 the contact regions of transmembrane helices, hinting for its possible role in helix-helix
510 association and relative motions.



511

512 **Figure 3. Results for amino acid transporter 1 (Uniprot id: P82251):** (a) Four contact maps
 513 constructed from coevolution analysis by DeepMetaPSICOV (DMP), RaptorX (RX), SPOT-
 514 Contact (SC), and the predicted structure by AlphaFold (AF) (only those contacts with
 515 probability >0.1 are shown). (b) W_3 scores for all residue positions based on the contact maps in
 516 (a), where red and blue dots mark residues with deleterious and neutral mutations, respectively.
 517 (c) Predicted structure by AlphaFold as colored by W_3 scores (red/blue for high/low values),
 518 where residues with deleterious and neutral mutations are shown as large and small balls,
 519 respectively (V142,L223,I44,P52,G63,W69,A70,G105,T123,A126,V170,A182,I187,G195,
 520 A224,W230,I241,G259,P261,V330,A331,R333,A354,S379,A382, see

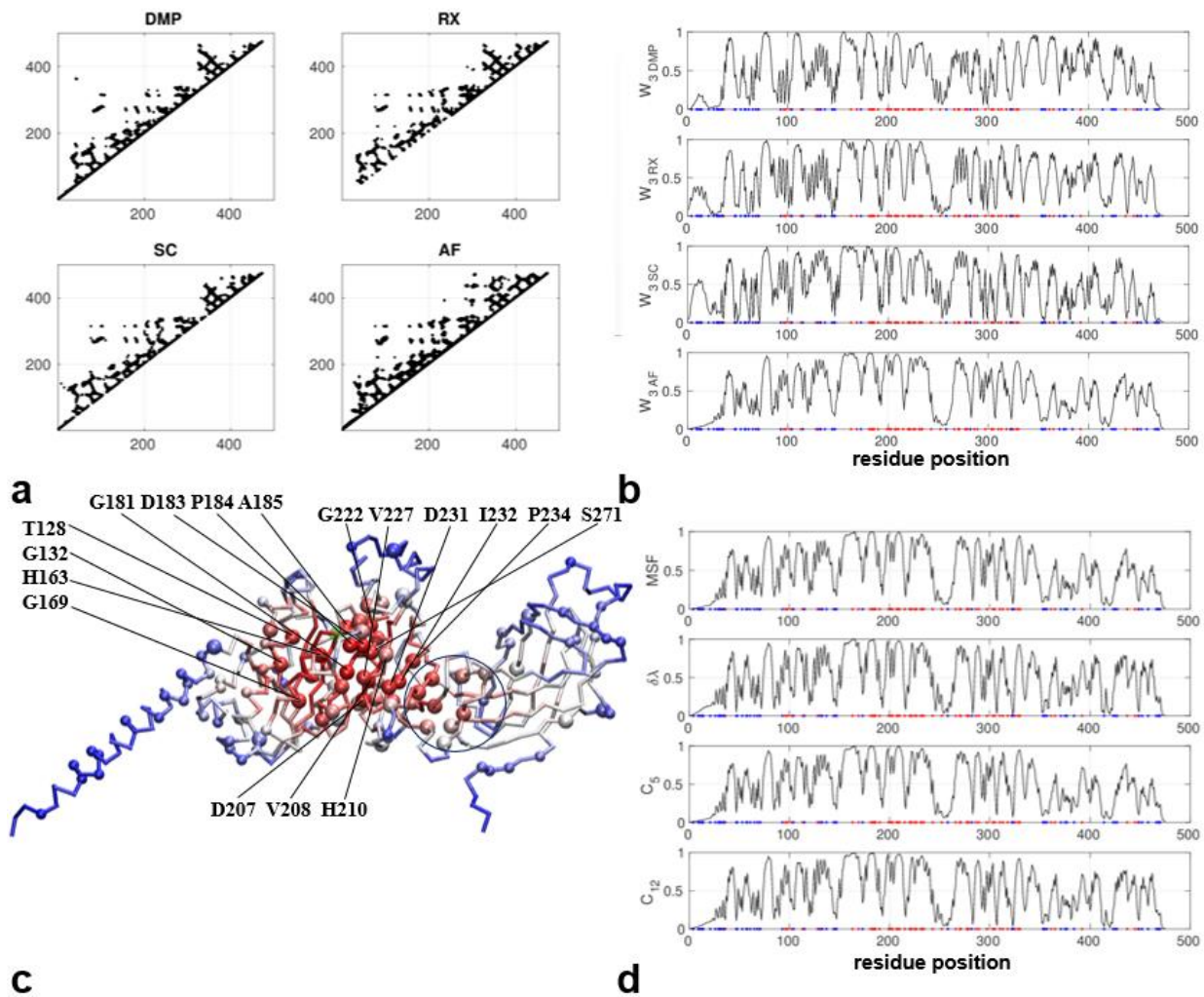
521 <https://www.uniprot.org/uniprotkb/P82251/variant-viewer>), and active-site residues are colored
522 in green. (d) Four other network scores (MSF, $\delta\lambda$, C5 and C12) for all residue positions based on
523 the contact maps in (a).

524

525 **Lipoprotein lipase** (Uniprot: P06858): As a key enzyme in triglyceride metabolism, it
526 catalyzes the hydrolysis of triglycerides from circulating chylomicrons and very low density
527 lipoproteins, thus playing an important role in lipid clearance from the blood stream, lipid
528 utilization and storage (see <https://www.uniprot.org/uniprotkb/P06858/entry#function>).

529 The AlphaFold-predicted structure consists of an N-terminal helix, a central α/β domains,
530 and a C-terminal β domain. Most deleterious mutations are concentrated inside the central
531 domain while the non-conserved residues are mostly located on the periphery (including N-
532 terminal helix and C-terminal domain) (see Fig 4c). The active site is comprised of a catalytic
533 triad of S159, D183, and H268 [60] in the central domain (colored green in Fig 4c). The contact
534 maps predicted by three coevolution analysis methods are similar to the contact map based on
535 the AlphaFold structure (see Fig 4a). As a result, the network scores agree well between these
536 methods (see Fig 4b) with minor differences in peripheral regions (such as the N-terminal helix).
537 As predicted, those residues identified with high network scores are primarily within the central
538 domain overlapping with most deleterious mutations, while the peripheral N-terminal helix and
539 C-terminal domain feature low scores consistent with low sequence conservation (see Fig 4c).
540 Notably, some of them are found at the interface between the central domain and the C-terminal
541 domain (circled in Fig 4c), possibly mediating inter-domain motions. Among those top-10%
542 hotspot residues (see Fig 4c), T128, G132, H163, G169, G181, D183, P184, A185, D207, V208,
543 H210, G222, V227, D231, I232, P234 and S271 are known to harbor pathogenic mutations in

544 Hyperlipoproteinemia 1, an autosomal recessive metabolic disorder characterized by defective
545 breakdown of dietary fats. Both H163 and G169 lie in helix 4 that constitutes part of the highly
546 conserved beta-epsilon serine-alpha folding motif which is near S159 of the active site.
547 Supporting their functional relevance, mutations H163R and G169E were found to abolish the
548 enzymatic activity [61] . Near D183 (one of the catalytic triad), mutations G181S and P184R
549 were found to abolish the catalytic activity [62] . Further from D183, conserved substations
550 D207E and H210Q abolished the enzyme activity [63], and mutations D231E, I232S and P234L
551 led to loss of the catalytic function [64] . These mutations may disrupt allosteric interactions with
552 the central catalytic domain. Another conservative mutation S271T (near D183) also led to loss
553 of enzyme activity [65]. Taken together, these residues may function by directly or indirectly
554 coupling to the active site.



555

556 **Figure 4. Results for Lipoprotein lipase (Uniprot id: P06858):** (a) Four contact maps
 557 constructed from coevolution analysis by DeepMetaPSICOV (DMP), RaptorX (RX), SPOT-
 558 Contact (SC), and the predicted structure by AlphaFold (AF) (only those contacts with
 559 probability >0.1 are shown). (b) W_3 scores for all residue positions based on the contact maps in
 560 (a), where red and blue dots mark residues with deleterious and neutral mutations, respectively.

561 (c) Predicted structure by AlphaFold as colored by W_3 scores (red/blue for high/low values),
 562 where residues with deleterious and neutral mutations are shown as large and small balls,
 563 respectively

564 (H71,A427,D36,N70,V96,A98,R102,W113,T128,G132,H163,G169,G181,D183,P184,A185,

565 G186,E190,S199,D201,A203,D207,V208,H210,G215, S220,I221,G222, K225,V227,D231,
566 I232,P234,C243,I252,C266,R270,S271,D277,S278,L279,S286,Y289,F297,L303,C305,
567 C310,L313,N318,S325,M328,L330,A361,S365,L392,E437,E437,C445,E448, see
568 <https://www.uniprot.org/uniprotkb/P06858/variant-viewer>), and active-site residues are colored
569 in green. (d) Four other network scores (MSF, $\delta\lambda$, C5 and C12) for all residue positions based on
570 the contact maps in (a).

571

572 **Conclusion**

573 To conclude, we have combined machine learning, network analysis, and protein
574 language models to develop a sequences-based variant site prediction method based on the
575 protein residue contact networks which incorporate sequential, structural, dynamic, and
576 interaction information simultaneously:

- 577 1. We build protein residue networks by exploiting three different state-of-the-art coevolution
578 analysis tools (RaptorX, DeepMetaPSICOV, and SPOT-Contact) that complement each other.
- 579 2. We use three powerful machine learning algorithms (Random Forest, Gradient Boosting, and
580 Extreme Gradient Boosting) to optimally combine 20 network centrality scores to accurately
581 predict key residues as hot spots for disease mutations.
- 582 3. We train and validate our method using a dataset of 107 proteins rich in disease mutations,
583 demonstrating its high accuracy in distinguishing between deleterious and neutral sites (with
584 AUC of ROC ~ 0.84). Further improvement can be achieved after combining our method with
585 the ESM-based method.

586 This study has established a useful strategy of combining an ensemble of network scores
587 based on different coevolution analysis methods via machine learning to predict key variants
588 sites of relevance to disease mutations. The code and dataset are made available to public to
589 enable future developments and applications (see <https://simtk.org/projects/hotspots>).

590 For future work, it will be interesting to go beyond contact map predictions by integrating
591 other scores derived from the co-evolution analysis (for example, see refs [66-68]) in our
592 workflow, which may further boost the accuracy of variant site prediction.

593

594

References

596

597

- 598 1. Jumper J, Evans R, Pritzel A, Green T, Figurnov M, et al. (2021) Highly accurate protein structure
599 prediction with AlphaFold. *Nature* 596: 583-589.
- 600 2. Baek M, DiMaio F, Anishchenko I, Dauparas J, Ovchinnikov S, et al. (2021) Accurate prediction of
601 protein structures and interactions using a three-track neural network. *Science* 373: 871-876.
- 602 3. Terwilliger TC, Liebschner D, Croll TI, Williams CJ, McCoy AJ, et al. (2023) AlphaFold predictions are
603 valuable hypotheses and accelerate but do not replace experimental structure determination.
604 *Nat Methods*.
- 605 4. Al-Janabi A (2022) Has DeepMind's AlphaFold solved the protein folding problem? *Biotechniques* 72:
606 73-76.
- 607 5. Medina E, D RL, Sanabria H (2021) Unraveling protein's structural dynamics: from configurational
608 dynamics to ensemble switching guides functional mesoscale assemblies. *Curr Opin Struct Biol*
609 66: 129-138.
- 610 6. Pak MA, Markhieva KA, Novikova MS, Petrov DS, Vorobyev IS, et al. (2023) Using AlphaFold to predict
611 the impact of single mutations on protein stability and function. *PLoS One* 18: e0282689.
- 612 7. Steinegger M, Soding J (2018) Clustering huge protein sequence sets in linear time. *Nat Commun* 9:
613 2542.
- 614 8. Burley SK, Berman HM, Bhikadiya C, Bi C, Chen L, et al. (2019) RCSB Protein Data Bank: biological
615 macromolecular structures enabling research and education in fundamental biology,
616 biomedicine, biotechnology and energy. *Nucleic Acids Res* 47: D464-D474.
- 617 9. Varadi M, Anyango S, Deshpande M, Nair S, Natassia C, et al. (2022) AlphaFold Protein Structure
618 Database: massively expanding the structural coverage of protein-sequence space with high-
619 accuracy models. *Nucleic Acids Res* 50: D439-D444.
- 620 10. Bepler T, Berger B (2021) Learning the protein language: Evolution, structure, and function. *Cell Syst*
621 12: 654-669 e653.
- 622 11. Cheng J, Novati G, Pan J, Bycroft C, Zengulyte A, et al. (2023) Accurate proteome-wide missense
623 variant effect prediction with AlphaMissense. *Science* 381: eadg7492.
- 624 12. Hassan MS, Shaalan AA, Dessouky MI, Abdelnaiem AE, ElHefnawi M (2019) A review study:
625 Computational techniques for expecting the impact of non-synonymous single nucleotide
626 variants in human diseases. *Gene* 680: 20-33.
- 627 13. Niroula A, Urolagin S, Vihinen M (2015) PON-P2: prediction method for fast and reliable
628 identification of harmful variants. *PLoS One* 10: e0117380.
- 629 14. Pejaver V, Urresti J, Lugo-Martinez J, Pagel KA, Lin GN, et al. (2020) Inferring the molecular and
630 phenotypic impact of amino acid variants with MutPred2. *Nat Commun* 11: 5918.
- 631 15. Qi H, Zhang H, Zhao Y, Chen C, Long JJ, et al. (2021) MVP predicts the pathogenicity of missense
632 variants by deep learning. *Nat Commun* 12: 510.
- 633 16. Yue P, Melamud E, Moulton J (2006) SNPs3D: candidate gene and SNP selection for association studies.
634 *BMC Bioinformatics* 7: 166.
- 635 17. Tang H, Thomas PD (2016) Tools for Predicting the Functional Impact of Nonsynonymous Genetic
636 Variation. *Genetics* 203: 635-647.
- 637 18. Katsonis P, Koire A, Wilson SJ, Hsu TK, Lua RC, et al. (2014) Single nucleotide variations: biological
638 impact and theoretical interpretation. *Protein Sci* 23: 1650-1666.

- 639 19. Singh A, Olowoyeye A, Baenziger PH, Dantzer J, Kann MG, et al. (2008) MutDB: update on
640 development of tools for the biochemical analysis of genetic variation. *Nucleic Acids Res* 36:
641 D815-819.
- 642 20. Bromberg Y, Rost B (2007) SNAP: predict effect of non-synonymous polymorphisms on function.
643 *Nucleic Acids Res* 35: 3823-3835.
- 644 21. Ng PC, Henikoff S (2003) SIFT: Predicting amino acid changes that affect protein function. *Nucleic*
645 *Acids Res* 31: 3812-3814.
- 646 22. Pers TH, Timshel P, Hirschhorn JN (2015) SNPsnap: a Web-based tool for identification and
647 annotation of matched SNPs. *Bioinformatics* 31: 418-420.
- 648 23. Zheng W, Tekpinar M (2009) Large-scale evaluation of dynamically important residues in proteins
649 predicted by the perturbation analysis of a coarse-grained elastic model. *BMC Struct Biol* 9: 45.
- 650 24. Zheng W, Brooks BR, Doniach S, Thirumalai D (2005) Network of dynamically important residues in
651 the open/closed transition in polymerases is strongly conserved. *Structure* 13: 565-577.
- 652 25. Ponzoni L, Bahar I (2018) Structural dynamics is a determinant of the functional significance of
653 missense variants. *Proc Natl Acad Sci U S A* 115: 4164-4169.
- 654 26. Butler BM, Gerek ZN, Kumar S, Ozkan SB (2015) Conformational dynamics of nonsynonymous
655 variants at protein interfaces reveals disease association. *Proteins* 83: 428-435.
- 656 27. Jumper J, Evans R, Pritzel A, Green T, Figurnov M, et al. (2021) Applying and improving AlphaFold at
657 CASP14. *Proteins* 89: 1711-1721.
- 658 28. Marks DS, Hopf TA, Sander C (2012) Protein structure prediction from sequence variation. *Nat*
659 *Biotechnol* 30: 1072-1080.
- 660 29. Morcos F, Pagnani A, Lunt B, Bertolino A, Marks DS, et al. (2011) Direct-coupling analysis of residue
661 coevolution captures native contacts across many protein families. *Proc Natl Acad Sci U S A* 108:
662 E1293-1301.
- 663 30. Hopf TA, Scharfe CP, Rodrigues JP, Green AG, Kohlbacher O, et al. (2014) Sequence co-evolution
664 gives 3D contacts and structures of protein complexes. *Elife* 3.
- 665 31. Marks DS, Colwell LJ, Sheridan R, Hopf TA, Pagnani A, et al. (2011) Protein 3D structure computed
666 from evolutionary sequence variation. *PLoS One* 6: e28766.
- 667 32. Burger L, van Nimwegen E (2010) Disentangling direct from indirect co-evolution of residues in
668 protein alignments. *PLoS Comput Biol* 6: e1000633.
- 669 33. Jones DT, Buchan DW, Cozzetto D, Pontil M (2012) PSICOV: precise structural contact prediction
670 using sparse inverse covariance estimation on large multiple sequence alignments.
671 *Bioinformatics* 28: 184-190.
- 672 34. Halabi N, Rivoire O, Leibler S, Ranganathan R (2009) Protein sectors: evolutionary units of three-
673 dimensional structure. *Cell* 138: 774-786.
- 674 35. Wang S, Sun S, Li Z, Zhang R, Xu J (2017) Accurate De Novo Prediction of Protein Contact Map by
675 Ultra-Deep Learning Model. *PLoS Comput Biol* 13: e1005324.
- 676 36. Ma J, Wang S, Wang Z, Xu J (2015) Protein contact prediction by integrating joint evolutionary
677 coupling analysis and supervised learning. *Bioinformatics* 31: 3506-3513.
- 678 37. Kandathil SM, Greener JG, Jones DT (2019) Prediction of interresidue contacts with
679 DeepMetaPSICOV in CASP13. *Proteins* 87: 1092-1099.
- 680 38. Jones DT, Singh T, Kosciolk T, Tetchner S (2015) MetaPSICOV: combining coevolution methods for
681 accurate prediction of contacts and long range hydrogen bonding in proteins. *Bioinformatics* 31:
682 999-1006.
- 683 39. Jones DT, Kandathil SM (2018) High precision in protein contact prediction using fully convolutional
684 neural networks and minimal sequence features. *Bioinformatics* 34: 3308-3315.

- 685 40. Hanson J, Paliwal K, Litfin T, Yang Y, Zhou Y (2018) Accurate prediction of protein contact maps by
686 coupling residual two-dimensional bidirectional long short-term memory with convolutional
687 neural networks. *Bioinformatics* 34: 4039-4045.
- 688 41. Yan W, Yu C, Chen J, Zhou J, Shen B (2020) ANCA: A Web Server for Amino Acid Networks
689 Construction and Analysis. *Front Mol Biosci* 7: 582702.
- 690 42. Amitai G, Shemesh A, Sitbon E, Shklar M, Netanel D, et al. (2004) Network analysis of protein
691 structures identifies functional residues. *J Mol Biol* 344: 1135-1146.
- 692 43. Velickovic P (2023) Everything is connected: Graph neural networks. *Curr Opin Struct Biol* 79:
693 102538.
- 694 44. Butler BM, Kazan IC, Kumar A, Ozkan SB (2018) Coevolving residues inform protein dynamics profiles
695 and disease susceptibility of nSNVs. *PLoS Comput Biol* 14: e1006626.
- 696 45. Chasman D, Adams RM (2001) Predicting the functional consequences of non-synonymous single
697 nucleotide polymorphisms: structure-based assessment of amino acid variation. *J Mol Biol* 307:
698 683-706.
- 699 46. Gerek ZN, Ozkan SB (2011) Change in allosteric network affects binding affinities of PDZ domains:
700 analysis through perturbation response scanning. *PLoS Comput Biol* 7: e1002154.
- 701 47. Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, et al. (2010) A method and server
702 for predicting damaging missense mutations. *Nat Methods* 7: 248-249.
- 703 48. Meier J, Rao R, Verkuil R, Liu J, Sercu T, et al. (2021) Language models enable zero-shot prediction of
704 the effects of mutations on protein function. *bioRxiv*: 2021.2007.2009.450648.
- 705 49. Vihinen M (2020) Problems in variation interpretation guidelines and in their implementation in
706 computational tools. *Mol Genet Genomic Med* 8: e1206.
- 707 50. Ashkenazy H, Erez E, Martz E, Pupko T, Ben-Tal N (2010) ConSurf 2010: calculating evolutionary
708 conservation in sequence and structure of proteins and nucleic acids. *Nucleic Acids Res* 38:
709 W529-533.
- 710 51. Zheng W, Brooks BR, Thirumalai D (2006) Low-frequency normal modes that describe allosteric
711 transitions in biological nanomachines are robust to sequence variations. *Proc Natl Acad Sci U S*
712 *A* 103: 7664-7669.
- 713 52. Zheng W (2016) Probing the structural dynamics of the SNARE recycling machine based on coarse-
714 grained modeling. *Proteins*.
- 715 53. Delgado J, Radosky LG, Cianferoni D, Serrano L (2019) FoldX 5.0: working with RNA, small molecules
716 and a new graphical interface. *Bioinformatics* 35: 4168-4169.
- 717 54. Veiga-da-Cunha M, Gerin I, Chen YT, de Barsey T, de Lonlay P, et al. (1998) A gene on chromosome
718 11q23 coding for a putative glucose- 6-phosphate translocase is mutated in glycogen-storage
719 disease types Ib and Ic. *Am J Hum Genet* 63: 976-983.
- 720 55. Hiraiwa H, Pan CJ, Lin B, Moses SW, Chou JY (1999) Inactivation of the glucose 6-phosphate
721 transporter causes glycogen storage disease type 1b. *J Biol Chem* 274: 5532-5536.
- 722 56. Wolfe MS, Xia W, Ostaszewski BL, Diehl TS, Kimberly WT, et al. (1999) Two transmembrane
723 aspartates in presenilin-1 required for presenilin endoproteolysis and gamma-secretase activity.
724 *Nature* 398: 513-517.
- 725 57. Sun L, Zhou R, Yang G, Shi Y (2017) Analysis of 138 pathogenic mutations in presenilin-1 on the in
726 vitro production of Abeta42 and Abeta40 peptides by gamma-secretase. *Proc Natl Acad Sci U S A*
727 114: E476-E485.
- 728 58. Yan R, Li Y, Shi Y, Zhou J, Lei J, et al. (2020) Cryo-EM structure of the human heteromeric amino acid
729 transporter b(0,+)-AT-rBAT. *Sci Adv* 6: eaay6379.
- 730 59. Font MA, Feliubadalo L, Estivill X, Nunes V, Golomb E, et al. (2001) Functional analysis of mutations
731 in SLC7A9, and genotype-phenotype correlation in non-Type I cystinuria. *Hum Mol Genet* 10:
732 305-316.

- 733 60. Emmerich J, Beg OU, Peterson J, Previato L, Brunzell JD, et al. (1992) Human lipoprotein lipase.
734 Analysis of the catalytic triad by site-directed mutagenesis of Ser-132, Asp-156, and His-241. *J*
735 *Biol Chem* 267: 4161-4165.
- 736 61. Reina M, Brunzell JD, Deeb SS (1992) Molecular basis of familial chylomicronemia: mutations in the
737 lipoprotein lipase and apolipoprotein C-II genes. *J Lipid Res* 33: 1823-1832.
- 738 62. Bruin T, Tuzgol S, van Diermen DE, Hoogerbrugge-van der Linden N, Brunzell JD, et al. (1993)
739 Recurrent pancreatitis and chylomicronemia in an extended Dutch kindred is caused by a
740 Gly154-->Ser substitution in lipoprotein lipase. *J Lipid Res* 34: 2109-2119.
- 741 63. Haubenwallner S, Horl G, Shachter NS, Presta E, Fried SK, et al. (1993) A novel missense mutation in
742 the gene for lipoprotein lipase resulting in a highly conservative amino acid substitution
743 (Asp180-->Glu) causes familial chylomicronemia (type I hyperlipoproteinemia). *Genomics* 18:
744 392-396.
- 745 64. Gotoda T, Yamada N, Kawamura M, Kozaki K, Mori N, et al. (1991) Heterogeneous mutations in the
746 human lipoprotein lipase gene in patients with familial lipoprotein lipase deficiency. *J Clin Invest*
747 88: 1856-1864.
- 748 65. Hata A, Emi M, Luc G, Basdevant A, Gambert P, et al. (1990) Compound heterozygote for lipoprotein
749 lipase deficiency: Ser----Thr244 and transition in 3' splice site of intron 2 (AG----AA) in the
750 lipoprotein lipase gene. *Am J Hum Genet* 47: 721-726.
- 751 66. Bisardi M, Rodriguez-Rivas J, Zamponi F, Weigt M (2022) Modeling Sequence-Space Exploration and
752 Emergence of Epistatic Signals in Protein Evolution. *Mol Biol Evol* 39.
- 753 67. Cocco S, Feinauer C, Figliuzzi M, Monasson R, Weigt M (2018) Inverse statistical physics of protein
754 sequences: a key issues review. *Rep Prog Phys* 81: 032601.
- 755 68. Rodriguez-Rivas J, Croce G, Muscat M, Weigt M (2022) Epistatic models predict mutable sites in
756 SARS-CoV-2 proteins and epitopes. *Proc Natl Acad Sci U S A* 119.
- 757

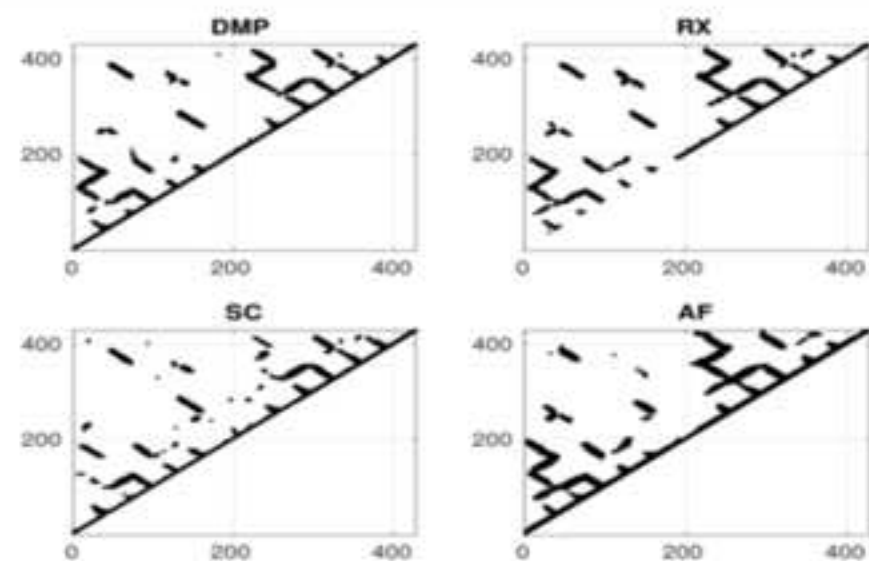
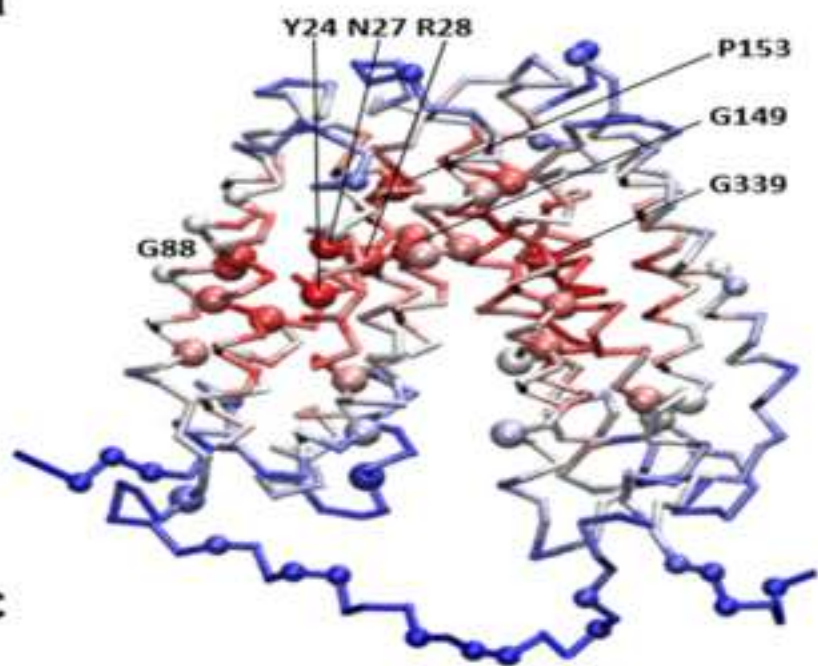
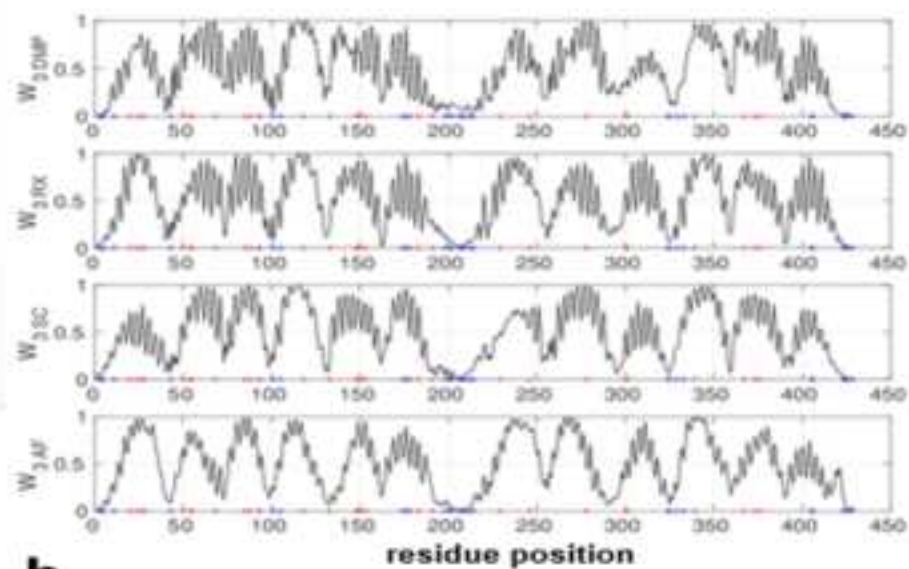
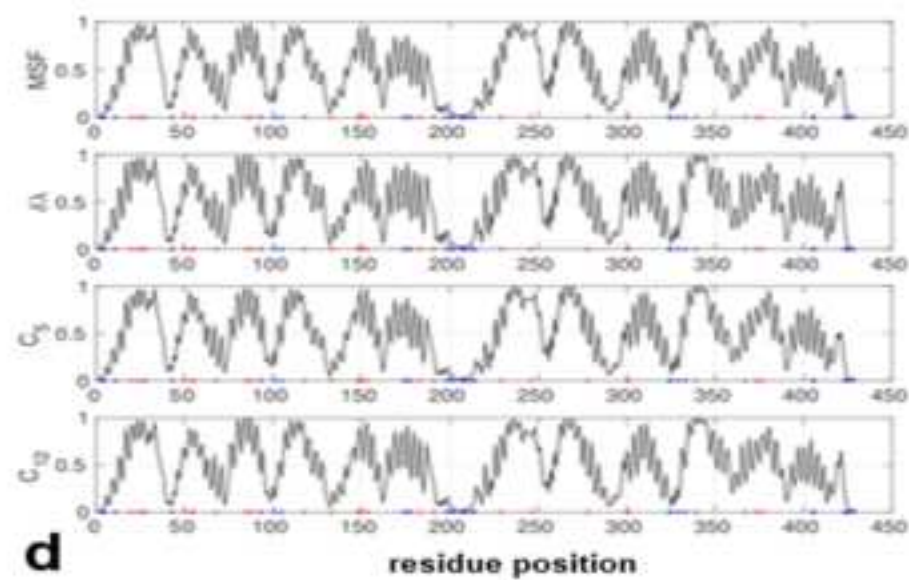
758

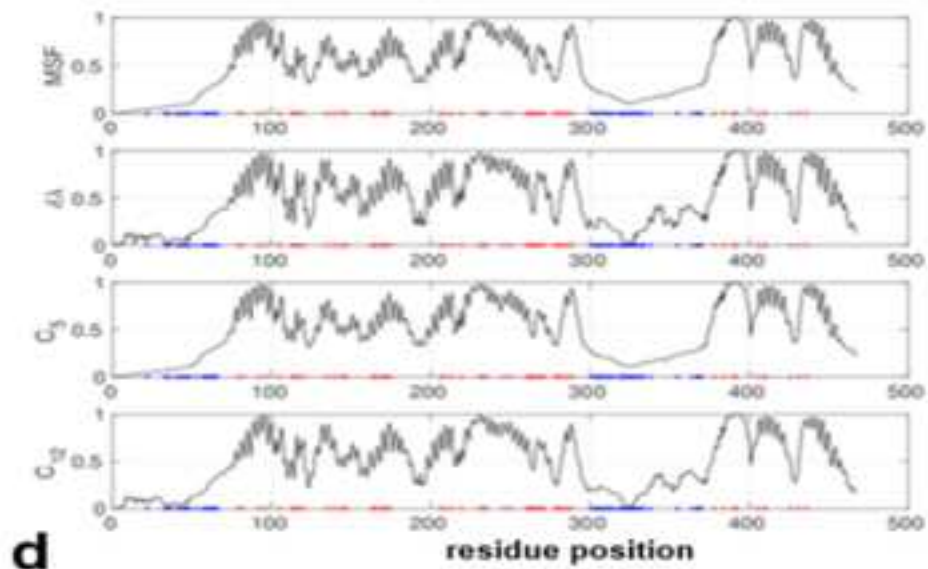
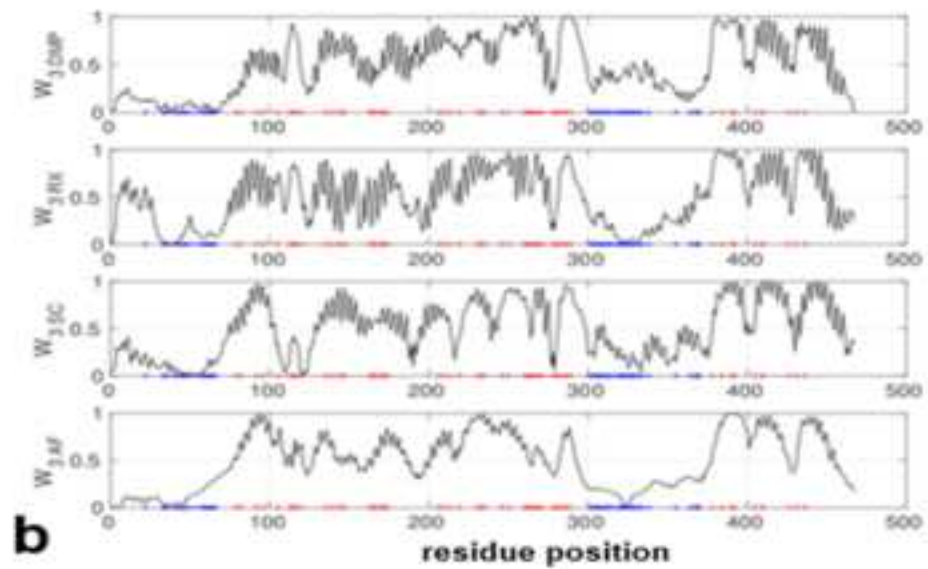
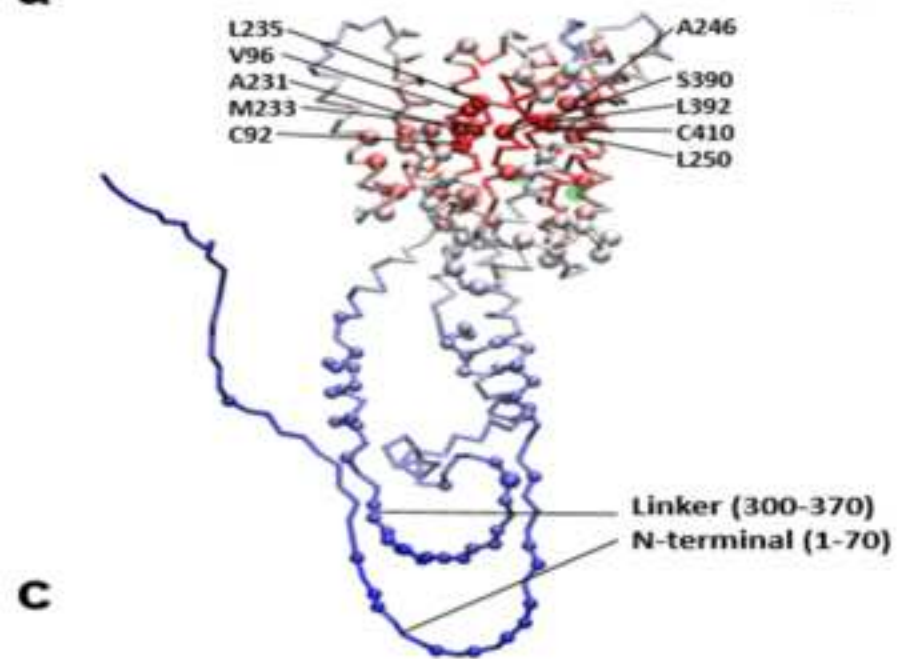
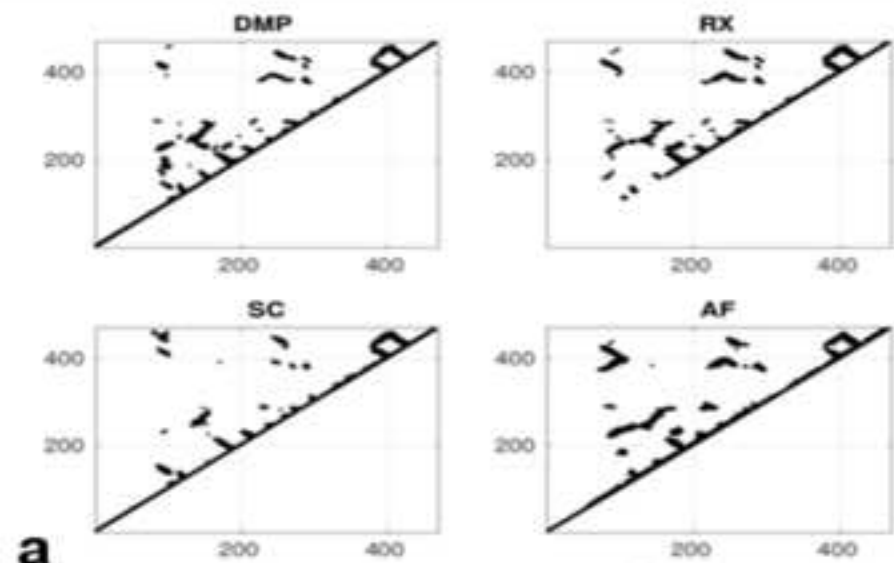
759

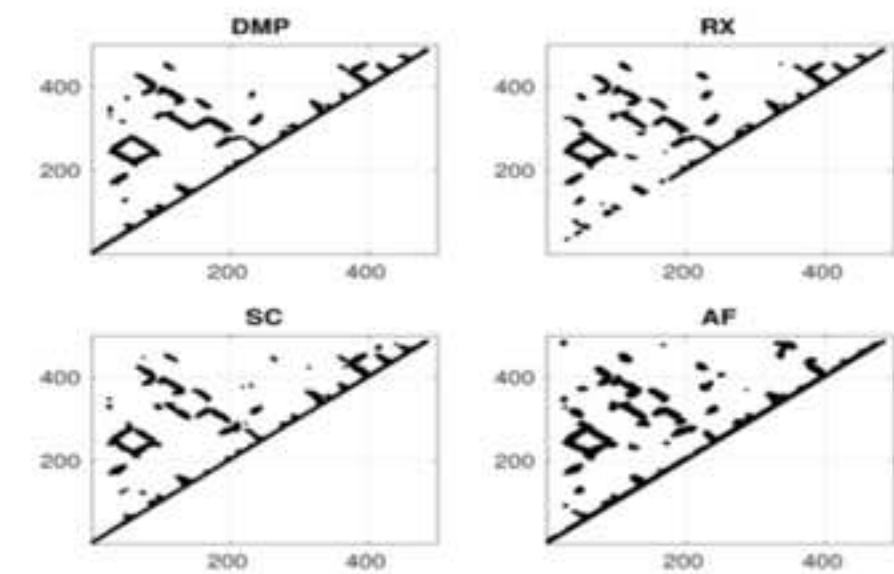
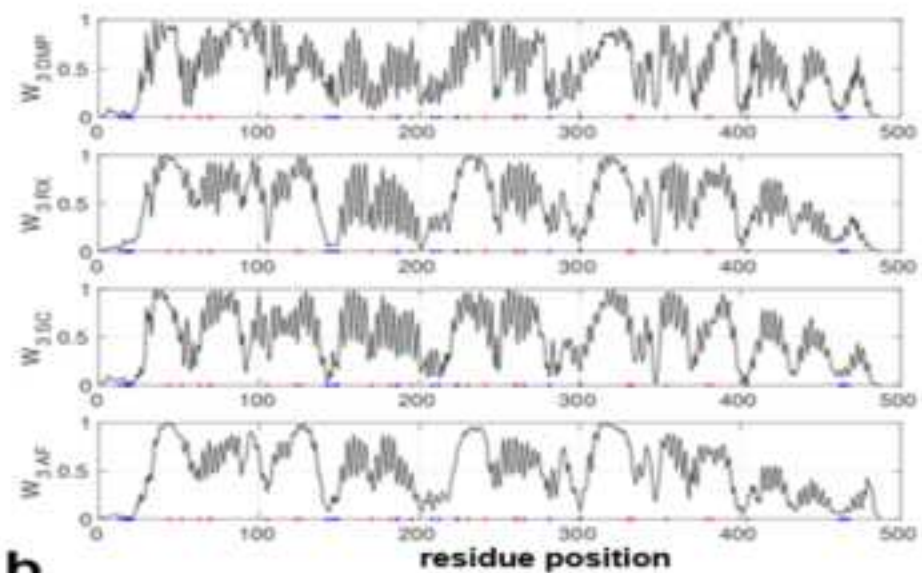
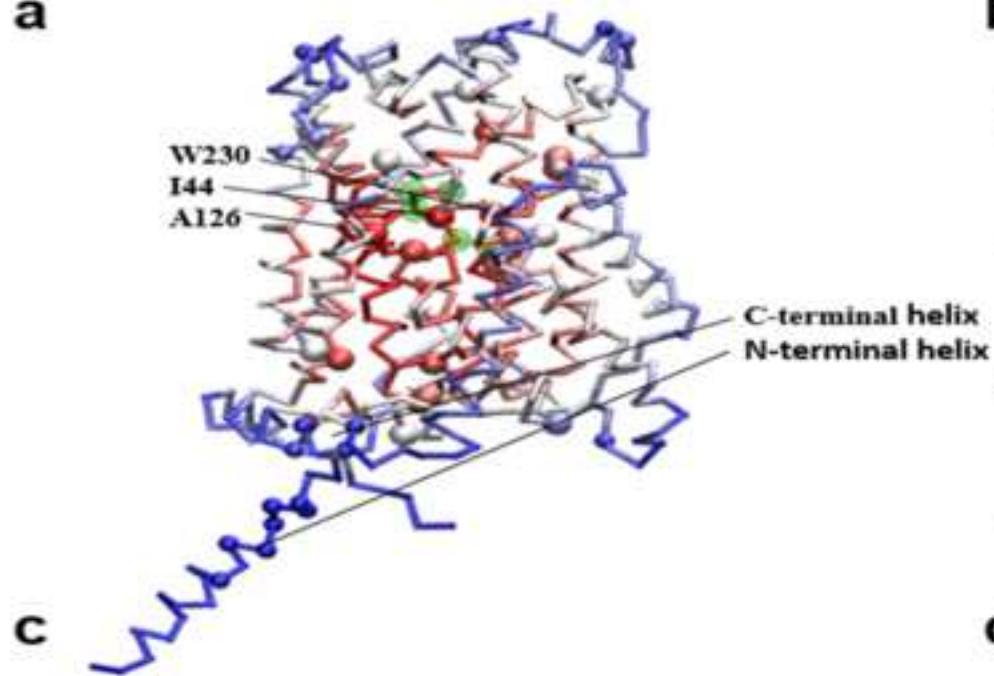
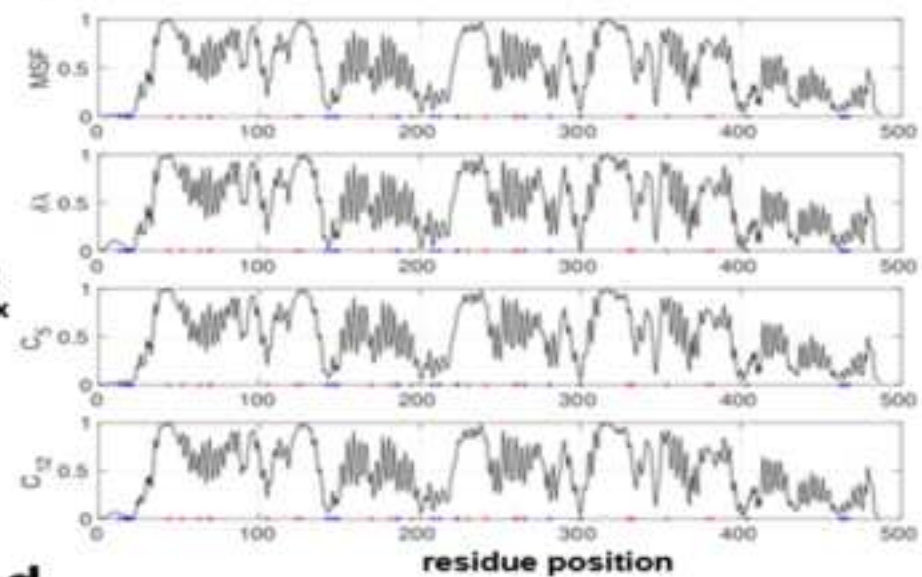
760 **Supporting Information**761 **S1 Table. Evaluation of 20 network scores based on protein residue contact maps**762 **constructed from 3 coevolution analysis tools (DeepMetaPSICOV, RaptorX, and SPOT-**763 **Contact)**

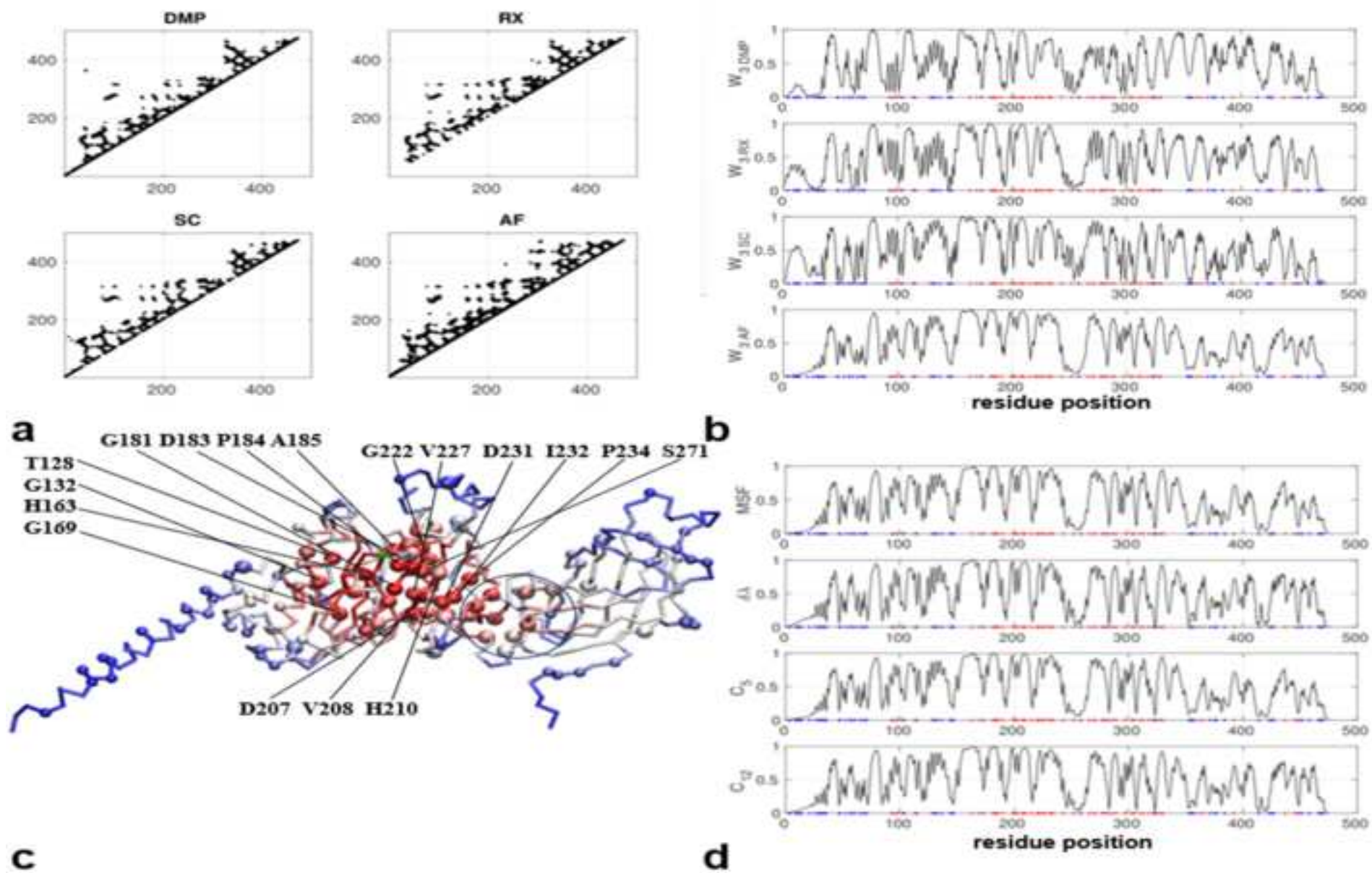
Score	AUC* of DeepMetaPSICOV	AUC* of RaptorX	AUC* of SPOT-Contact
C1	0.75±0.13	0.78±0.14	0.75±0.15
C2	0.75±0.17	0.75±0.19	0.76±0.19
C3	0.78±0.10	0.75±0.15	0.70±0.15
C4	0.65±0.12	0.52±0.15	0.61±0.07
C5	0.78±0.17	0.78±0.17	0.78±0.18
C6	0.63±0.16	0.56±0.19	0.65±0.17
C7	0.77±0.11	0.61±0.20	0.72±0.16
C8	0.65±0.12	0.52±0.15	0.61±0.07
C9	0.79±0.13	0.78±0.16	0.75±0.16
C10	0.76±0.12	0.75±0.14	0.69±0.15
C11	0.80±0.13	0.78±0.17	0.78±0.17
C12	0.78±0.16	0.80±0.14	0.80±0.17
C13	0.75±0.17	0.74±0.18	0.76±0.19
$\delta\lambda$	0.80±0.13	0.78±0.14	0.78±0.16
MSF	0.81±0.14	0.79±0.15	0.80±0.17
W_1	0.81±0.14	0.79±0.15	0.80±0.17
W_2	0.81±0.15	0.77±0.15	0.79±0.17
W_3	0.81±0.14	0.78±0.15	0.80±0.16
W_∞	0.80±0.13	0.77±0.14	0.78±0.16
W_s	0.80±0.13	0.78±0.13	0.78±0.16

764 * mean ± standard-deviation

**a****c****b****d**



**a****b****c****d**



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16

Predicting hotspots for disease-causing single nucleotide variants using sequences-based coevolution, network analysis, and machine learning

Short title: Predicting nSNVs hotspots with sequences-based machine learning

Wenjun Zheng ^{1*}

¹ Department of Physics, State University of New York at Buffalo, NY 14260, United States of America

* Corresponding author

E-mail: wjzheng@buffalo.edu (WZ)

Keywords: Centrality, Coevolution, Disease mutations, Machine learning, Protein residue contact network, ~~Protein language model~~, Single nucleotide variant

Formatted: Font: Not Bold

17 Abstract

18 To enable personalized ~~genetics and~~ medicine, it is important yet highly challenging to
19 accurately predict disease-causing mutations in target proteins at high throughput. Previous
20 computational methods have been developed using evolutionary ~~phylogeny information~~ in
21 combination with ~~various~~ biochemical and structural ~~properties/features~~ of ~~amino acids~~ protein
22 ~~residues~~ to discriminate neutral vs. deleterious mutations. However, the power of these methods
23 is often limited because they ~~either assume known protein structures or do not fully~~
24 ~~incorporate~~ treat residues independently, ~~structural, dynamic, and interaction information critical~~
25 ~~for protein functions. To address these limitations without fully considering their global~~
26 ~~interactions.~~ To address the above limitations, we build upon recent progress in machine
27 learning, network analysis, and protein language models, and develop a sequences-based variant
28 site prediction workflow based on the protein residue contact networks: 1. We employ and
29 integrate various methods of building protein residue networks using state-of-the-art coevolution
30 analysis tools (~~e.g.,~~ RaptorX, DeepMetaPSICOV, and SPOT-Contact) powered by deep learning.
31 2. We use machine learning algorithms (~~e.g.,~~ Random Forest, Gradient Boosting, and Extreme
32 Gradient Boosting) to optimally combine ~~13-20~~ network centrality scores (~~calculated by~~
33 ~~NetworkX~~) with ~~7 other network scores calculated from the contact probability matrices~~ to
34 jointly predict key residues as hot spots for disease mutations. 3. Using a dataset of 107 proteins
35 rich in disease mutations, we rigorously evaluate the network scores individually and collectively
36 (~~via machine learning~~) in comparison with alternative structures-based network scores (~~using~~
37 ~~predicted structures by AlphaFold~~). ~~By optimally combining three coevolution analysis methods~~
38 ~~and the resulting network scores by machine learning, we are able to discriminate deleterious and~~
39 ~~neutral mutation sites accurately (AUC of ROC—0.84). Furthermore, by combining our method~~

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

40 ~~with a state-of-the-art predictor of the functional effects of sequence variations based on large~~
41 ~~protein language models, we have significantly improved the prediction of disease variant sites~~
42 ~~(AUC \rightarrow 0.89).~~ This work supports a promising strategy of combining an ensemble of network
43 scores based on different coevolution analysis methods (and optionally predictive scores from
44 other methods) via machine learning to predict candidate sites of disease mutations, which will
45 inform downstream applications of disease diagnosis and targeted drug design.

46

47 Introduction

48 The holy grail of structural biology is to solve high-resolution biomolecular structures at
49 the genomic scale to inform mechanistic studies of their functions. Thanks to recent revolutions
50 in computational structural biology (~~e.g.~~ accurate protein structure prediction by AlphaFold [1]
51 and ~~other deep learning based methods~~ [RoseTTAFold](#) [2]), it is now feasible to predict ~~static~~
52 [native](#) structures for ~~many many~~ proteins ~~of interest~~ given their sequences (with some caveats,
53 see [3]), thus practically solving the protein folding problem [4]. However, it remains
54 challenging to predict dynamic structural ensembles [5] and mutation-induced ~~structural~~
55 ~~changes~~ [effects](#) [6] to meet the demand of mechanistic studies of protein functions and
56 dysfunctions. While the public databases of protein sequences and variations increase rapidly
57 owing to genomic/metagenomic sequencing efforts (~~e.g.~~ the MetaClust database contains about
58 1.6 billion protein sequence fragments [7]), the growth of experimental protein structures [8] and
59 predicted structures remains to catch up (~~e.g.~~ the AlphaFold database contains over 200 million
60 predicted structures [9]). Such sequences-structures gap has motivated the development of new
61 computational tools that make functional sense of protein sequences without directly using
62 structural information (for example, by using deep learning to train large protein language
63 models [10]). Recently, AlphaMissense attained state of the art prediction of missense variant
64 pathogenicity by adapting AlphaFold fine-tuned on human and primate variant population
65 frequency databases [11].

66 A major interest in personalized ~~genetics and~~ medicine is in understanding novel genetic
67 variations through genotype-phenotype association studies in relation to diseases. Particularly, a
68 rapidly growing number of non-synonymous single nucleotide variants (nSNVs) have been

69 uncovered in protein coding regions that can adversely impact protein function and cause
70 diseases [12]. Various computational methods were developed using evolutionary conservation
71 and phylogeny in combination with biochemical and structural properties of amino acids to
72 discriminate neutral vs. deleterious nSNVs [13-22]. Protein structural dynamics has also proven
73 useful in discovering functionally important residues [23,24] which could constitute hot spots
74 for disease-causing nSNVs [25,26]. However, the requirement of 3D structures has limited the
75 number of nSNVs that can be analyzed by existing structure-based computational tools, although
76 such constraint has been significantly alleviated by recent progress in protein structure prediction
77 [27].

78 As alternatives to structure-based methods, sequences-based coevolution analysis has
79 become increasingly powerful in predicting structural couplings between pairs of contacting
80 residues [28-31], thanks to the development of direct coupling methods that can overcome the
81 confounding indirect coupling effects [29,32,33]. In principle, coevolving pairs of residues can
82 be identified from a sufficiently large multiple sequence alignment, allowing the prediction of
83 close spatial proximity in the native structures. Boosted by deep learning and other algorithmic
84 developments, this coevolution analysis has led to accurate prediction of residue contacts which
85 make *de novo* protein structure prediction possible [28]. Furthermore, coevolution analysis
86 (enhanced by deep learning) has also been used to study various aspects of protein functional
87 interactions such as allostery [34]. For example, RaptorX uses an ultra-deep neural network
88 combining coevolution information with sequence conservation information to infer 3D contacts
89 with higher accuracy than previous methods [35,36]. DeepMetaPSICOV [37] combines the input
90 feature sets used by earlier methods (~~e.g.~~ MetaPSICOV [38] and DeepCov [39]) as input to a
91 deep, fully convolutional residual neural network. SPOT-Contact predicts protein contact maps

92 by stacking residual convolutional networks with two-dimensional residual bidirectional
93 recurrent LSTM networks, and using both one-dimensional sequence-based and two-dimensional
94 evolutionary coupling based information [40]. These three state-of-the-art coevolution analysis
95 methods are employed in this study to construct protein residue contact maps for network
96 analysis (see below).

97 Another line of protein research is based on the treatment of a protein as a network where
98 amino acid residues are nodes and their bonded/non-bonded interactions form edges [41]. Such
99 models can be readily built upon 3D native structures so that a whole suite of network analysis
100 tools (see <https://networkx.org/>) can be applied. For example, Amitai et al [42] used network
101 analysis of protein structures (~~using e.g.~~ closeness centrality) to identify functional residues.
102 Going beyond network analysis, deep-learning-based study of protein graph neural networks is
103 an active area of research [43].

104 In a recent paper, Butler et al [44] proposed a sequence-based Gaussian network model
105 (Seq-GNM) to calculate the dynamic profile of a protein without a 3D structure. They used
106 coevolution analysis to build a network model which connects residues predicted to be in contact
107 via evolutionary couplings. Their work built on previous studies that shown crystallographic B-
108 factors are useful in predicting the impact of nSNVs on protein function [45,46] : rigid sites with
109 low B-factors are more susceptible to destabilizing nSNVs than flexible sites with high B-
110 factors. Indeed, existing computational tools to diagnose neutral and deleterious nSNVs
111 (~~e.g. such as~~ PolyPhen-2 [47]) use crystallographic B-factors along with other evolutionary and
112 structural features. More specifically, Butler et al used Seq-GNM to compute B-factors for
113 protein residues, and they found that deleterious nSNVs are overabundant at low B-factor sites,
114 while neutral nSNVs are overabundant at high B-factor sites. Mechanistically, low B-factors may

115 indicate that a site is crucial for maintaining structural stability and/or modulating functional
116 motions (~~e.g.,~~ as a hinge) and thus susceptible to mutations. In contrast, high B-factors are
117 associated with flexible regions (~~e.g., loops~~) with minimal interactions, which are thus more
118 robust to mutations. Based on these observations, they proposed that the sequences-based
119 predicted B-factors can discriminate between deleterious and neutral nSNVs without structural
120 information.

121 Inspired by the above study and recent progress in machine learning, network analysis,
122 and protein language models, we further develop the sequences-based protein residue network
123 analysis in the following directions: 1. We ~~exploit and integrate various methods of building~~
124 protein residue networks using ~~three state-of-the-art different~~ coevolution analysis tools (~~e.g.,~~
125 RaptorX, DeepMetaPSICOV, and SPOT-Contact) ~~as powered-enabled~~ by deep learning. 2. We
126 ~~use-exploit three~~ machine learning algorithms (~~e.g.,~~ Random Forest, Gradient Boosting, and
127 Extreme Gradient Boosting) to optimally combine ~~13-20 distinct~~ network ~~node~~ centrality scores
128 (~~calculated by the NetworkX package~~) with ~~as 7 other network scores~~ calculated from the contact
129 probability matrices to ~~jointly~~ predict ~~key residues as~~ hot spot ~~residuess~~ for disease mutations. 3.
130 ~~Using-Based on~~ a dataset of 107 proteins ~~rich in with known disease deleterious/neutral~~
131 mutations, we ~~rigorously~~ evaluate ~~the our sequences-based sequences-based~~ network scores ~~both~~
132 individually and in combination, ~~and then (via machine learning) in comparison~~ with alternative
133 structures-based network scores and a physics force field based method (~~using predicted~~
134 ~~structures by AlphaFold~~). By optimally combing three coevolution analysis methods and the
135 resulting 20 network scores by machine learning, we are able to discriminate deleterious and
136 neutral mutation sites accurately (AUC of ROC ~ 0.84), which is on par with structure-based
137 network scores (AUC ~ 0.83). Furthermore, by combining our method with a state-of-the-art

138 predictor of the functional effects of sequence variation based on large protein language models
139 (ESM [48]), we have significantly improved the prediction of disease variant sites (AUC ~ 0.89).

140 ~~In the following sections, we first describe the detailed methodology in the order of~~
141 ~~the proposed workflow, then we report the results of evaluation of our network-based scores both~~
142 ~~individually and collectively (via machine learning), finally we discuss specific case studies of~~
143 ~~four proteins to illustrate the usage of our method. This work supports the strategy of combining~~
144 ~~an ensemble of network scores based on different coevolution analysis methods via machine~~
145 ~~learning to predict candidate sites for disease mutations which will inform many downstream~~
146 ~~biomedical applications.~~

147 **Materials and methods**

148 Here is a summary of the workflow of our sequences-based method:

- 149 a. Collect datasets of protein sequences and variants (see Section 1)
- 150 b. Run co-evolution analysis of a given target protein sequence to build a residue
151 contact map P (see Section 2)
- 152 c. Use NetworkX to calculate node centrality scores based on P (see Section 3)
- 153 d. Use sequence-based GNM to calculate additional node scores (see Section 4)
- 154 e. (optional) Use protein language model (ESM) to predict variant importance (see
155 Section 5)
- 156 f. (optional) Use AlphaFold and FoldX to predict variant importance (see Section 6 and
157 7)
- 158 g. Use machine learning to optimally combine the above scores for classifying
159 deleterious vs neutral variant sites (see Section 8)

Formatted: List Paragraph, Numbered + Level: 1 +
Numbering Style: a, b, c, ... + Start at: 1 + Alignment:
Left + Aligned at: 0.5" + Indent at: 0.75"

Formatted: Font: (Default) Times New Roman, 12 pt

Formatted: Font: 12 pt, Not Bold

1. Training and testing Datasets of protein sequences and variants

A dataset of 107 protein sequences with ≤ 500 residues and ≥ 20 annotated deleterious/neutral variants were collected from the HumVar database [47] (sources: humvar-2011_12.deleterious.pph.input and humvar-2011_12.neutral.pph.input from <ftp://genetics.bwh.harvard.edu/pph2/training/training-2.2.2.tar.gz>). Their UniProt ids and sequences are as follows:

<https://simtk.org/projects/hotspots>. This diverse dataset contains 97 proteins with their pairwise sequence identity < 30%.

<https://simtk.org/projects/hotspots>. This diverse dataset contains 97 proteins with their pairwise sequence identity < 30%.

<https://simtk.org/projects/hotspots>. This diverse dataset contains 97 proteins with their pairwise sequence identity < 30%.

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

The HumVar dataset consists of 13,032 human disease-causing mutations from UniProt and 8,946 human nonsynonymous single-nucleotide polymorphisms (nsSNPs) without annotated involvement in disease. This dataset was previously used to train and test PolyPhen-2 [47] for predicting damaging effects of missense mutations, and was used by Butler et al [44] in

183 benchmarking their seq-GNM method for predicting deleterious/neutral nSNVs. ~~However,~~
184 ~~because they used different subsets of HumVar to train and test their methods, it is not possible to~~
185 ~~directly compare the performance between our method and theirs.~~

186
187 Since this dataset is highly imbalanced (there are 4040 deleterious mutation sites but only
188 120 neutral mutation sites) [49], we have added 3403 additional neutral sites with very low
189 conservation scores (i.e. grade ≤ 2 as assessed by the ConSurf program [50]). Our objective is to
190 train and test a binary classifier of residues in these proteins as deleterious or neutral. To this
191 end, we split 107 proteins into training and testing sets (with 79 and 28 proteins, respectively),
192 and perform evaluations based on the testing set. The main metric of evaluation is the ROC
193 curves and associated area under the curve (AUC). AUC is a standard metric for evaluating
194 binary classifiers based on the ROC curve of sensitivity and specificity. The ROC curves are also
195 used in other computational papers for variant prediction (see [47]).

Formatted: Font: (Default) Times New Roman, 12 pt,
Font color: Auto

Formatted: Font: (Default) Times New Roman, 12 pt,
Font color: Auto

Formatted: Font: (Default) Times New Roman, 12 pt,
Font color: Auto

197 **S2. Sequences-based coevolution analysis and protein contact map** 198 **construction**

199 We perform coevolution analysis using three state-of-the-art methods: the RaptorX server
200 (<http://raptorx.uchicago.edu>), the DeepMetaPSICOV server (<http://bioinf.cs.ucl.ac.uk/psipred/>),
201 and the SPOT-Contact server (<https://sparks-lab.org/server/spot-contact/>). A sequence length limit
202 (500) is imposed by the capacity of coevolution analysis servers, and may be circumvented if
203 installing and running coevolution analysis locally.

204 These methods use multiple sequence alignments to compute the probability P_{ij} of residue
205 pair (i, j) forming spatial contact. Based on the matrix of predicted P_{ij} , a protein residue contact
206 map can be built with residues as nodes and pairwise contacts as edges weighted by P_{ij} . By default,
207 we do not apply any threshold cutoff to P_{ij} for defining contacts (unless networks with unweighted
208 edges are required by some node centrality algorithms in NetworkX, where we remove edges with
209 $P_{ij} < 0.1$, and set weight to 1 for the remaining edges).

210

211 **3. Network analysis of protein contact map**

212 By treating a protein contact map as a network of nodes and edges, we calculate various
213 node centrality scores to predict key residues as hotspots for disease mutations.

214 A simple score to measure node centrality is a weighted node degree that accounts for the
215 nearest neighbor interactions (denoted W_1):

$$216 \quad W_{1,i} = \sum_{k \neq i} P_{ik} \quad (1)$$

217 To include indirect couplings beyond the nearest neighbors, we calculate the node degree
218 based on the n'th power of the contact probability matrix (denoted W_n):

$$219 \quad W_{n,i} = \sum_{k \neq i} P_{ik} W_{n-1,k} = \sum_{k \neq i} P^n_{ik} \quad (2)$$

220 As n goes to infinity, W_n converges to the eigenvector of P matrix with the highest
221 eigenvalue λ_{\max} (denoted W_∞):

222 $PW_\infty = \lambda_{\max} W_\infty$ (3)

223 Among various W_n , W_2 can be interpreted as the node degrees of a new network based on
 224 a neighborhood similarity matrix S as follows (denoted W_s):

225 $S_{ij} = \sum_{k \neq i, j} P_{ik} P_{jk}$, $W_{s,i} = \sum_{k \neq i} S_{ik}$ (4)

226 In this study we use five network scores (W_1 , W_2 , W_3 , W_∞ and W_s) as predictive features
 227 for node importance. Additionally, we exploit 13 network centrality metrics as calculated by the
 228 NetworkX package (see Table 1). To allow meaningful comparison of scores between proteins,
 229 the scores of each protein are sorted and their ranking percentiles are linearly transformed to
 230 values between 0 and 1.

231 **Table 1. Network centrality scores as implemented in the NetworkX package**

232 (see <https://networkx.org/documentation/stable/reference/algorithms/centrality.html>)

Symbol	Centrality name	Definition
C1	degree_centrality	Corresponding to W_1
C2	eigenvector_centrality	Corresponding to W_∞
C3	closeness_centrality	Closeness centrality of a node u is the reciprocal of the average shortest path distance to u over all $n-1$ reachable nodes.
C4	betweenness_centrality	Betweenness centrality of a node u is the sum of the fraction of all-pairs shortest paths that pass through u .
C5	current_flow_closeness_centrality	Current-flow closeness centrality is a variant of closeness centrality based on effective resistance between nodes in a network.
C6	current_flow_betweenness_centrality	Current-flow betweenness centrality is based on an electrical current model for information spreading.
C7	communicability_betweenness_centrality	Communicability betweenness centrality is based on the number of walks connecting every pair of nodes.
C8	load_centrality	Load centrality of a node u is the fraction of all shortest paths that pass through u .
C9	subgraph_centrality	Subgraph centrality of a node u is the sum of weighted closed walks of all lengths starting and ending at u .
C10	harmonic_centrality	Harmonic centrality of a node u is the sum of the reciprocal of the shortest path distances from all other nodes to u .

C11	second_order_centrality	Second order centrality of a node u is the standard deviation of the return times to u of a perpetual random walk on G .
C12	laplacian_centrality	Laplacian Centrality of a node u is measured by the drop in the Laplacian Energy after deleting u from the graph.
C13	katz_centrality_numpy	Katz centrality computes the centrality for a node u based on the centrality of its neighbors. It is a generalization of the eigenvector centrality.

233

234

235

236 **4. Sequences-based GNM**

237 For comparison, we implemented Bulter et al's sequence-based GNM [44]. The original
 238 structure-based Gaussian network model (GNM) represents a protein structure as an elastically
 239 connected network of residues to obtain the equilibrium fluctuations of residues. In the absence
 240 of a structure, the sequence-based GNM (Seq-GNM) treats coevolving residue pairs as
 241 contacting pairs.

242 To construct the Kirchhoff matrix (denoted K), each non-bonded residue pair is assigned
 243 a value of -1 times its contact probability. The bonded residue pairs $(i, i+1)$ are assigned -1 to
 244 enforce local chain connectivity. The diagonal elements of K are assigned so that the sum of each
 245 row and column is zero:

$$246 \quad K_{ij} = \begin{cases} -P_{ij} & i \neq j \\ \sum_{k \neq i} P_{ik} & i = j \end{cases} \quad (5)$$

247 The vibrational thermal fluctuations of residues are evaluated by inverting the Kirchhoff
 248 matrix (or summing over the modes as weighted by $1/\lambda_m$). The per-residue mean-square

249 fluctuations (MSF), which are proportional to the crystallographic B factors, are given as
250 follows:

$$251 \quad MSF_i \propto K_{ii}^{-1} = \sum_{m>0} \frac{V_{mi}^2}{\lambda_m} \quad (6)$$

252 where the eigen-decomposition of K gives eigenvectors V_m and eigenvalues λ_m that satisfy:

$$253 \quad KV_m = \lambda_m V_m \quad (7)$$

254 Low-MSF residues correspond to rigid cores or hinges of dynamical importance [44].

255 As an alternative way to evaluate node importance using GNM, we perform a
256 perturbation-based hotspot analysis as follows: For mode m , calculate how much its eigenvalue
257 changes ($\delta\lambda_{m,i}$) in response to a perturbation at a chosen residue position i [23,24,51] (i.e., by
258 uniformly weakening the contacts with residue i). Then compute $\delta\lambda_i = \sum_m \delta\lambda_{m,i}$ to assess the
259 dynamic importance of this residue position [52]. High- $\delta\lambda_i$ residues correspond to sites highly
260 sensitive to local perturbations that mimic mutations.

261 The above two GNM-based scores are combined with the other network scores for
262 machine learning.

263

264 **5. ESM based variant prediction**

265 For comparison with our method, we use a deep-learning variant predictor based on a
266 large protein language model (ESM). We downloaded and installed the ESM package and

267 pretrained models from <https://github.com/facebookresearch/esm>. Since our dataset consists of
268 known variants (from HumVar) and added non-conserved sites (with specific mutations
269 unknown), we simulate the mutational effects on each site by introducing Alanine substitution if
270 the wildtype residue is not an Alanine and Glycine substitution otherwise. Then we process the
271 mutated sequence with 5 pretrained ESM models (esm1v_t33_650M_UR90S_1,
272 esm1v_t33_650M_UR90S_2, esm1v_t33_650M_UR90S_3, esm1v_t33_650M_UR90S_4, and
273 esm1v_t33_650M_UR90S_5), which predict the difference in the probability of observing the
274 wildtype residue and the mutant residue at a given site [48]. We record the predictions of five
275 ESM models as separate features to be optimally integrated via machine learning.

276

277 **6. AlphaFold for structural prediction**

278 We downloaded predicted structures for the 107 proteins from AlphaFold DB
279 (<https://alphafold.ebi.ac.uk/>). A residue contact probability matrix is constructed based on the
280 predicted structures as follows:

$$281 \quad P_{ij} = \frac{1}{1 + e^{d_{ij} - 10}} \quad (8)$$

282 where d_{ij} is the distance between residues i and j , and 10 \AA is used as a soft cutoff distance. We
283 then use this contact probability matrix to perform the same network analysis as in the
284 sequences-based method and for optimization with machine learning.

285

286 **7. Foldx for structural refinement and Alanine scanning analysis**

287 FoldX program [53] was downloaded from <https://foldxsuite.crg.eu/>. We use the
288 RepairPDB command to refine the AlphaFold-predicted models (by fixing bad torsion angles
289 and Van der Waals clashes). Then we use the AlaScan command to mutate each residue to Ala
290 and calculate the resulting changes in Gibbs free energies which are then used as a feature to
291 predict hotspots of disease mutations.

292

293 **8. Machine learning algorithms**

294 We use the following machine learning methods of the scikit-learn package
295 (<https://scikit-learn.org/stable/>) to learn optimal combinations of all features to predict if a given
296 site is deleterious or neutral mutation site:

297 Random Forest Classifier (RF) (`sklearn.ensemble.RandomForestClassifier`): A random
298 forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of
299 the dataset and uses averaging to improve the predictive accuracy and control over-fitting. We
300 tune the following hyper-parameters: max_depth, n_estimators, max_features.

301 Gradient Boosting Classifier (GB) (`sklearn.ensemble.GradientBoostingClassifier`): This
302 algorithm builds an additive model in a forward stage-wise fashion. In each stage a regression
303 tree is fit on the negative gradient of the loss function, e.g. binary log loss. We tune the following
304 hyper-parameters: n_estimators, max_depth, max_features.

305

306 - Extreme Gradient Boosting Classifier (XGB) (`xgboost.XGBClassifier`): This algorithm is
307 an optimized distributed version of gradient boosting designed to be highly efficient, flexible and

308 portable. We tune the following hyper-parameters: n_estimators, max_depth, reg_alpha,
309 reg_lambda.

310 These three methods were chosen because they have performed successfully in machine
311 learning contests in Kaggle (see <https://www.packtpub.com/product/the-kaggle->
312 book/9781801817479). They are also relatively cheap to train and optimize compared with the
313 deep learning methods.

314 We use Optuna (<https://optuna.org/>) for hyper-parameter tuning of the above algorithms.

315 We have run Optuna multiple times to ensure the resulting best metric is reproducible.

Formatted: Font color: Auto

316

317 **Results and discussion**

318

319 This study explores how to systematically utilize the coevolution information from
320 multiple sequence alignments to model and analyze a protein as a residue contact network
321 beyond the scope of GNM. To this end, we first use coevolution analysis to construct a protein
322 residue contact map with edges weighted by the predicted contact probability; then we exploit an
323 array of 20 network-based scores to assess the node importance as predictors for disease
324 mutation sites; finally we evaluate the predictive power of these scores individually and
325 collectively (using machine learning) based on a subset of 107 protein sequences and their
326 variants from the HumVar database. For comparison, we also evaluate alternative methods based
327 on predicted protein structures, a physics-based force field, and protein language models.

328

329 **1.**

330 **Evaluation of individual network scores**

331 Based on the protein residue contact maps built from three coevolution analysis tools
332 (DeepMetaPSICOV, RaptorX, and SPOT-Contact), we applied network analysis to calculate 20
333 network scores (see Table 2), measuring node centrality using various different algorithms (see
334 Methods). These scores include simple weighted node degrees for n-hop nearest neighbors (see
335 Methods) and more sophisticated centrality metrics (see Table 1), along with 2 seq-GNM based
336 scores (MSF and $\delta\lambda$, see Methods). We evaluate the performance of each score using the AUC
337 of ROC for the testing set, which provides a balanced evaluation of sensitivity and specificity as

Formatted: Indent: First line: 0"

338 ~~functions of the cutoff score~~ (see Table 2). More specifically, we sort all testing-set variants by a
 339 particular score and predict a variant deleterious/neutral if its score is above/below a cutoff value.
 340 This results in an ROC curve from which we have calculated AUC (see Table 2).

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

Formatted: Font color: Auto

341 Overall, DeepMetaPSICOV (max AUC=0.80) and SPOT-Contact (max AUC=0.81)
 342 perform slightly better than RaptorX (max AUC=0.78). Interestingly, simple weighted node
 343 degrees (W_1 , W_2 , and W_3) perform better than those more complex centrality scores (see Table
 344 2). When computing node degrees, going beyond the nearest neighbors seems to improve the
 345 prediction slightly (see Table 2). Two GNM-based scores perform similarly but slightly worse
 346 than the weighted node degrees (see Table 2). Among those NetworkX-based scores (see Table
 347 1), C5, C11 and C12 outperform the others, while those betweenness-based scores (C4, C6, and
 348 C8) underperform (see Table 2). Therefore, the functional importance of a node/residue pertains
 349 more to its role as a highly-connected hub than as an information bottleneck of the shortest paths.

350 **Table 2. Evaluation of 20 network scores based on protein residue contact maps**
 351 **constructed from 3 coevolution analysis tools (DeepMetaPSICOV, RaptorX, and SPOT-**
 352 **Contact) and AlphaFold-predicted structures**

Score	AUC* _{of} DeepMetaPSICOV	AUC* _{of} RaptorX	AUC* _{of} SPOT-Contact	AUC* _{of} AlphaFold
C1	0.74	0.76	0.73	0.82
C2	0.73	0.74	0.76	0.77
C3	0.76	0.73	0.69	0.73
C4	0.64	0.54	0.60	0.58
C5	0.78	0.76	0.79	0.80
C6	0.63	0.58	0.67	0.60
C7	0.75	0.61	0.72	0.74
C8	0.64	0.54	0.60	0.58
C9	0.77	0.76	0.74	0.78
C10	0.75	0.73	0.68	0.75
C11	0.79	0.76	0.77	0.80
C12	0.77	0.77	0.79	0.83
C13	0.73	0.73	0.76	0.76

$\delta\lambda$	0.79	0.76	0.78	0.83
MSF	0.79	0.76	0.78	0.80
W_1	0.79	0.77	0.78	0.83
W_2	0.80	0.78	0.80	0.83
W_3	0.80	0.78	0.81	0.82
W_∞	0.80	0.74	0.79	0.77
W_s	0.80	0.78	0.80	0.83
FoldX				0.68

* The AUC is calculated based on the ROC for all variants of the 28 testing set proteins.

Alternatively, we also calculated AUCs based on the ROCs of individual proteins and their summary statistics (see Table S1).

For comparison with alternative methods, we evaluated the performance of variant prediction by five pre-trained protein language models (ESM, see Methods), and the resulting AUC varies between 0.79 and 0.81, which are comparable to the network scores (see Table 2). For further comparison with structures-based methods, we also performed network analysis based on protein structures as predicted by AlphaFold (see Methods). Overall, the structures-based scores (max AUC=0.83) perform slightly better than the sequences-based scores. This may be partly due to the structure-based contact maps (see Eq. 8) being more sharply defined than the fuzzier contact-probability-based contact maps. Notably, when structural information is used, our network analysis performs significantly better than a physics-based force field (FoldX) with AUC=0.68. Taken together, these findings support the usefulness of individual sequences-based network centrality scores in predicting important residues on par with alternative more sophisticated methods.

To further understand the different accuracies of the above scores, we explore the relationships between them by evaluating the pairwise Pearson correlations (PC) (see Table 3). W_1 , W_2 , W_3 , W_∞ , W_s , MSF and $\delta\lambda$ are highly correlated (with $PC \geq 0.93$ for DeepMetaPSICOV,

C13	0.57 0.99 0.45 0.20 0.87 0.40 0.51 0.20 0.59 0.38 0.73 0.80 1.00 0.69 0.75 0.78 0.88 0.73 0.67 0.69
W ₁	0.80 0.70 0.50 0.37 0.84 0.64 0.76 0.37 0.68 0.47 0.80 0.91 0.69 1.00 0.98 0.97 0.84 0.98 0.98 1.00
W ₂	0.78 0.76 0.52 0.35 0.88 0.61 0.75 0.35 0.71 0.49 0.82 0.93 0.75 0.98 1.00 1.00 0.90 1.00 0.97 0.98
W ₃	0.77 0.79 0.53 0.34 0.89 0.59 0.73 0.34 0.72 0.49 0.83 0.92 0.78 0.97 1.00 1.00 0.92 0.99 0.96 0.97
W _∞	0.67 0.88 0.52 0.30 0.87 0.49 0.64 0.30 0.68 0.47 0.78 0.83 0.88 0.84 0.90 0.92 1.00 0.90 0.86 0.84
W _s	0.80 0.74 0.55 0.37 0.87 0.62 0.77 0.37 0.73 0.52 0.83 0.90 0.73 0.98 1.00 0.99 0.90 1.00 0.98 0.98
MSF	0.82 0.68 0.56 0.40 0.82 0.66 0.80 0.40 0.71 0.54 0.82 0.86 0.67 0.98 0.97 0.96 0.86 0.98 1.00 0.98
C1	1.00 0.63 0.58 0.22 0.64 0.22 0.29 0.22 0.75 0.65 0.71 0.78 0.55 0.88 0.87 0.85 0.58 0.87 0.68 0.73
C2	0.63 1.00 0.64 0.15 0.82 0.23 0.24 0.15 0.66 0.66 0.78 0.70 0.93 0.60 0.65 0.67 0.82 0.64 0.70 0.55
C3	0.58 0.64 1.00 0.44 0.79 0.35 0.39 0.44 0.64 0.93 0.88 0.51 0.68 0.48 0.50 0.50 0.64 0.50 0.69 0.47
C4	0.22 0.15 0.44 1.00 0.32 0.70 0.67 1.00 0.11 0.38 0.39 0.06 0.22 0.18 0.14 0.14 0.30 0.15 0.37 0.17
C5	0.64 0.82 0.79 0.32 1.00 0.45 0.43 0.32 0.65 0.72 0.96 0.71 0.86 0.63 0.65 0.66 0.79 0.65 0.81 0.62
C6	0.22 0.23 0.35 0.70 0.45 1.00 0.70 0.70 0.12 0.25 0.45 0.17 0.29 0.32 0.28 0.26 0.40 0.28 0.51 0.31
C7	0.29 0.24 0.39 0.67 0.43 0.70 1.00 0.67 0.34 0.34 0.48 0.28 0.29 0.30 0.28 0.27 0.40 0.29 0.49 0.39
C8	0.22 0.15 0.44 1.00 0.32 0.70 0.67 1.00 0.11 0.38 0.39 0.06 0.22 0.18 0.14 0.14 0.30 0.15 0.37 0.17
C9	0.75 0.66 0.64 0.11 0.65 0.12 0.34 0.11 1.00 0.71 0.72 0.71 0.61 0.66 0.69 0.69 0.62 0.69 0.66 0.64
C10	0.65 0.66 0.93 0.38 0.72 0.25 0.34 0.38 0.71 1.00 0.83 0.55 0.67 0.52 0.53 0.53 0.66 0.54 0.72 0.51
C11	0.71 0.78 0.88 0.39 0.96 0.45 0.48 0.39 0.72 0.83 1.00 0.68 0.82 0.64 0.65 0.65 0.77 0.65 0.83 0.63
C12	0.78 0.70 0.51 0.06 0.71 0.17 0.28 0.06 0.71 0.55 0.68 1.00 0.64 0.76 0.77 0.77 0.63 0.77 0.67 0.70
C13	0.55 0.93 0.68 0.22 0.86 0.29 0.29 0.22 0.61 0.67 0.82 0.64 1.00 0.53 0.58 0.61 0.85 0.58 0.72 0.52
W ₁	0.88 0.60 0.48 0.18 0.63 0.32 0.30 0.18 0.66 0.52 0.64 0.76 0.53 1.00 0.98 0.97 0.69 0.98 0.81 0.85
W ₂	0.87 0.65 0.50 0.14 0.65 0.28 0.28 0.14 0.69 0.53 0.65 0.77 0.58 0.98 1.00 0.99 0.73 1.00 0.80 0.83
W ₃	0.85 0.67 0.50 0.14 0.66 0.26 0.27 0.14 0.69 0.53 0.65 0.77 0.61 0.97 0.99 1.00 0.75 0.99 0.80 0.82
W _∞	0.58 0.82 0.64 0.30 0.79 0.40 0.40 0.30 0.62 0.66 0.77 0.63 0.85 0.69 0.73 0.75 1.00 0.73 0.88 0.67
W _s	0.87 0.64 0.50 0.15 0.65 0.28 0.29 0.15 0.69 0.54 0.65 0.77 0.58 0.98 1.00 0.99 0.73 1.00 0.81 0.83
MSF	0.68 0.70 0.69 0.37 0.81 0.51 0.49 0.37 0.66 0.72 0.83 0.67 0.72 0.81 0.80 0.80 0.88 0.81 1.00 0.79
C1	1.00 0.87 0.77 0.38 0.95 0.42 0.82 0.38 0.90 0.81 0.96 0.98 0.85 0.97 0.97 0.96 0.86 0.97 0.95 0.97
C2	0.87 1.00 0.78 0.26 0.91 0.32 0.81 0.26 0.98 0.79 0.90 0.89 0.99 0.87 0.91 0.93 0.99 0.91 0.91 0.87
C3	0.77 0.78 1.00 0.52 0.79 0.37 0.74 0.52 0.82 0.97 0.83 0.72 0.78 0.71 0.74 0.75 0.78 0.74 0.79 0.71
C4	0.38 0.26 0.52 1.00 0.33 0.61 0.55 1.00 0.30 0.51 0.39 0.30 0.27 0.30 0.30 0.29 0.27 0.30 0.33 0.30
C5	0.95 0.91 0.79 0.33 1.00 0.46 0.85 0.33 0.92 0.79 0.99 0.96 0.91 0.95 0.96 0.96 0.91 0.96 1.00 0.95
C6	0.42 0.32 0.37 0.61 0.46 1.00 0.67 0.61 0.33 0.31 0.46 0.40 0.33 0.41 0.39 0.37 0.33 0.39 0.46 0.41
C7	0.82 0.81 0.74 0.55 0.85 0.67 1.00 0.55 0.83 0.72 0.85 0.81 0.81 0.80 0.82 0.83 0.81 0.82 0.84 0.80
C8	0.38 0.26 0.52 1.00 0.33 0.61 0.55 1.00 0.30 0.51 0.39 0.30 0.27 0.30 0.30 0.29 0.27 0.30 0.33 0.30
C9	0.90 0.98 0.82 0.30 0.92 0.33 0.83 0.30 1.00 0.84 0.93 0.91 0.96 0.89 0.93 0.95 0.97 0.93 0.92 0.89
C10	0.81 0.79 0.97 0.51 0.79 0.31 0.72 0.51 0.84 1.00 0.84 0.75 0.79 0.73 0.77 0.78 0.79 0.77 0.80 0.73
C11	0.96 0.90 0.83 0.39 0.99 0.46 0.85 0.39 0.93 0.84 1.00 0.95 0.90 0.94 0.95 0.95 0.90 0.95 0.99 0.94
C12	0.98 0.89 0.72 0.30 0.96 0.40 0.81 0.30 0.91 0.75 0.95 1.00 0.88 1.00 1.00 0.99 0.89 0.99 0.96 1.00
C13	0.85 0.99 0.78 0.27 0.91 0.33 0.81 0.27 0.96 0.79 0.90 0.88 1.00 0.85 0.90 0.92 0.99 0.90 0.91 0.85
W ₁	0.97 0.87 0.71 0.30 0.95 0.41 0.80 0.30 0.89 0.73 0.94 1.00 0.85 1.00 0.99 0.97 0.86 0.98 0.95 1.00
W ₂	0.97 0.91 0.74 0.30 0.96 0.39 0.82 0.30 0.93 0.77 0.95 1.00 0.90 0.99 1.00 1.00 0.91 1.00 0.96 0.99
W ₃	0.96 0.93 0.75 0.29 0.96 0.37 0.83 0.29 0.95 0.78 0.95 0.99 0.92 0.97 1.00 1.00 0.93 1.00 0.96 0.97
W _∞	0.86 0.99 0.78 0.27 0.91 0.33 0.81 0.27 0.97 0.79 0.90 0.89 0.99 0.86 0.91 0.93 1.00 0.91 0.91 0.86
W _s	0.97 0.91 0.74 0.30 0.96 0.39 0.82 0.30 0.93 0.77 0.95 0.99 0.90 0.98 1.00 1.00 0.91 1.00 0.96 0.98
MSF	0.95 0.91 0.79 0.33 1.00 0.46 0.84 0.33 0.92 0.80 0.99 0.96 0.91 0.95 0.96 0.96 0.91 0.96 1.00 0.95

381

382

383 In summary, by evaluating 20 network scores individually, we have found a wide range

384 of performance with AUC varying from 0.54 to 0.81 (see Table 2). The top-performing scores

385 seem to correlate strongly with each other, so they must have captured a common aspect of node
386 centrality that is relevant to functional importance (e.g. high local connectivity instead of high
387 betweenness). Interestingly, the two GNM-based scores, despite measuring distinct dynamic
388 properties (MSF measures thermal fluctuations while $\delta\lambda$ measures sensitivity to local
389 perturbations), are also strongly correlated with each other and those degree-based network
390 scores. Therefore, to speed up the variant prediction workflow we only need to compute those
391 simpler weighted node degrees as features without significantly losing accuracy.

392

393 **2. Combining all network scores to predict variant hotspots by** 394 **machine learning**

395 To optimize the predictive power of the above network-based scores based on three
396 coevolution analysis methods (or AlphaFold), we have employed machine learning algorithms
397 (see Methods) to take them as input features, train a binary classifier which predicts if a residue
398 position is linked to neutral or deleterious variants (using first 79 proteins as training set), and
399 then test its prediction using the remaining 28 proteins as testing set. We use the AUC of ROC as
400 the metric for assessing the prediction quality of the trained classifier.

401 To evaluate the protein residue contact maps constructed by each method, we combine all
402 network scores based on the contact maps predicted by the same method (see Table 2) for
403 machine learning. The resulting AUC of each coevolution analysis method (DeepMetaPSICOV,
404 RaptorX, and SPOT-Contact) is 0.81, 0.80, and 0.82, respectively (see Table 4), which are
405 slightly better than the best AUC of individual scores (0.78~0.81, see Table 2). The lack of
406 substantial improvement may be due to high correlations among the scores (see Table 3) which

407 could reduce the effectiveness of ensemble learning. For comparison, we also trained and tested
 408 classifiers using the AlphaFold-predicted contact maps, and alternative classifiers based on
 409 protein language models (see Methods). Both alternative methods give comparable yet slightly
 410 better AUC (0.83). Similar to our finding, Butler et al reported AUC of 0.81 after combining the
 411 B-factors of Seq-GNM with evolutionary features [44].

412 **Table 4. Evaluation of classifiers trained by 3 machine learning algorithms (RF, GB and**
 413 **XGB, see Methods) based on the protein residue contact maps constructed from 3**
 414 **coevolution analysis tools (DeepMetaPSICOV, RaptorX, and SPOT-Contact), AlphaFold-**
 415 **predicted structures, and protein language models (ESM).**

Sources of input features	AUC of RF	AUC of GB	AUC of XGB
DeepMetaPSICOV	0.81	0.81	0.81
RaptorX	0.80	0.80	0.80
SPOT-Contact	0.82	0.82	0.82
AlphaFold	0.83	0.83	0.83
ESM	0.83	0.83	0.83
All 3 coevolution methods	0.84	0.84	0.84
All 3 coevolution methods (w/o C1-C13)	0.82	0.82	0.83
All 3 coevolution methods and ESM	0.89	0.89	0.89
AlphaFold and ESM	0.88	0.88	0.88

416
 417

418 To further boost the prediction performance, we have sought to combine the network
 419 scores of all three coevolution analysis methods for machine learning, resulting in better AUC
 420 (0.84) which slightly outperform both AlphaFold and ESM (0.83). To assess the added value of
 421 including 13 NetworkX-based centrality scores (see Table 1), we have performed an ablation
 422 study that excludes them in machine learning, and found slightly lower AUC (0.82~0.83). So it

423 is possible to speed up the calculation without significantly reducing accuracy. Taken together,
424 our findings support the power of combining an array of different network scores from different
425 coevolution analysis tools to optimize the prediction in the framework of ensemble learning.

426 To further explore how well our method complements alternative methods, we have
427 combined all the network scores with the ESM scores in machine learning. Encouragingly, we
428 have obtained markedly improved AUC (0.89), which is comparable to machine learning that
429 combines the AlphaFold-based network scores with the ESM scores (AUC=0.88).

430 For comparison with other studies, Butler et al showed that Seq-GNM combined with
431 evolutionary parameters attained a sensitivity of 0.84 and a specificity of 0.66 [44]. PolyPhen-2
432 achieved a sensitivity of 0.73 and a specificity of 0.8 on the HumVar datasets [47]. While using
433 different training and testing datasets, we have attained competitive results with a sensitivity of
434 0.82 and a specificity of 0.80 (using all the network scores from three coevolution analysis tools
435 and the ESM scores). For more direct comparison, we also evaluated PolyPhen-2 based on the
436 same 28 testing-set proteins and their variants, and obtained an AUC of 0.85, which is close to
437 our method (see Table 4). However, this metric is likely positively biased, since PolyPhen-2 has
438 been trained on the HumVar dataset.

Formatted: Font: 12 pt, Font color: Auto

Formatted: Font: 12 pt, Font color: Auto

Formatted: Font: 12 pt, Font color: Auto

Formatted: Font color: Auto

439 In summary, via extensive machine learning, we have demonstrated the power of using
440 an ensemble of sequences-based network scores calculated by different co-evolution analysis
441 tools to accurately predict deleterious mutation sites. Although some network scores are highly
442 correlated (see Table 3) and they vary widely in accuracy (see Table 2), these scores seem to be
443 sufficiently diverse to allow effective ensemble learning when combined.

444

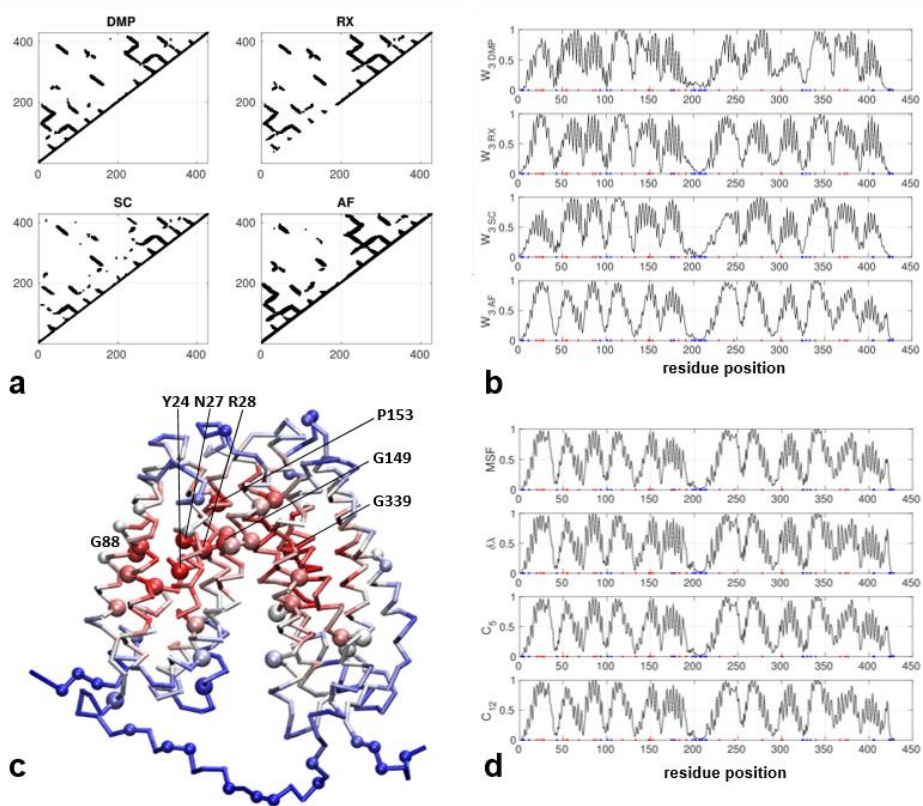
445 **3. Case studies:**

446 To illustrate the biomedical significance of our predictions of variant sites with network scores,
447 we discuss in details the following four proteins from our dataset.

448 **Glucose-6-phosphate exchanger** (Uniprot id: O43826): As an inorganic phosphate and
449 glucose-6-phosphate antiporter, it transports cytoplasmic glucose-6-phosphate into the lumen of
450 the endoplasmic reticulum and translocates inorganic phosphate in the opposite direction. Being
451 involved in glucose production through glycogenolysis and gluconeogenesis, it plays a central
452 role in homeostatic regulation of blood glucose levels. It is linked to diseases like congenital
453 disorder of glycosylation and glycogen storage disease (see
454 <https://www.uniprot.org/uniprotkb/O43826/entry#function>).

455 The AlphaFold-predicted structure forms a dimer of transmembrane helical domains with
456 most deleterious mutation sites concentrated inside the central core while those non-conserved
457 residues (i.e. neutral mutation sites) are mostly located on the periphery (see Fig 1c). The contact
458 maps predicted by three coevolution analysis tools all agree well with the contact map based on
459 the AlphaFold structure (see Fig 1a) (except that RaptorX omitted many local contacts in
460 residues 1-200). As a result, the network centrality scores (e.g. W_3) also agree well between
461 these methods (see Fig 1b), although the coevolution-based network scores are generally noisier
462 (with more spikes) than the structure-based scores (see Fig 1b). Different network scores
463 calculated from the same contact map are also highly similar (see Fig 1d) despite being based on
464 different algorithms. For example, scores of $\delta\lambda$ and MSF agree very well (see Fig 1d).
465 Encouragingly, those residues identified with high network scores are primarily within the
466 central core (inside each domain or in the inter-domain hinge region), thus overlapping with

467 most deleterious mutations (see Fig 1c). Among those top-10% hotspot residues (see Fig 1c),
 468 mutations Y24H, N27K, R28H, G88D, G149E, P153L, and G339C were implicated in causing
 469 glycogen storage disease [54] . Two of these mutations (R28H and G149E) were found to exhibit
 470 undetectable microsomal glucose-6-phosphate transport activity in transient expression
 471 studies[55], thus confirming their functional importance.



472
 473 **Figure 1. Results for Glucose-6-phosphate exchanger (Uniprot id: O43826):** (a) Four contact
 474 maps constructed from coevolution analysis by DeepMetaPSICOV (DMP), RaptorX (RX),
 475 SPOT-Contact (SC), and the predicted structure by AlphaFold (AF) (only those contacts with

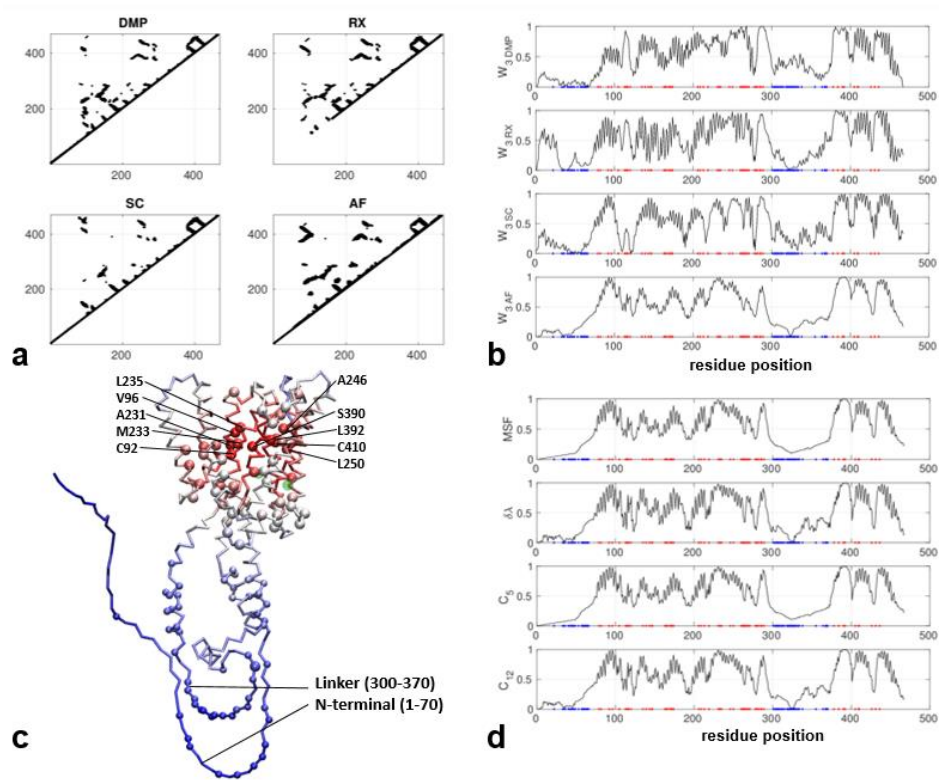
476 probability >0.1 are shown). (b) W_3 scores for all residue positions based on the contact maps in
477 (a), where red and blue dots mark residues with deleterious and neutral mutations, respectively.
478 (c) Predicted structure by AlphaFold as colored by W_3 scores (red/blue for high/low values),
479 where residues with deleterious and neutral mutations are shown as large and small balls,
480 respectively (G20, Y24, N27, R28, G50, S54, S55, G68, L85, G88, W118, Q133, A148, G149,
481 G150, P153, C176, C183, P191, L229, W246, I278, R300, H301, G339, A367, A373, G376, see
482 <https://www.uniprot.org/uniprotkb/O43826/variant-viewer>). (d) Four other network scores (MSF,
483 $\delta\lambda$, C5 and C12) for all residue positions based on the contact maps in (a).

484

485 **Presenilin-1** (Uniprot id: P49768): As the catalytic subunit of the gamma-secretase complex, it
486 catalyzes the intramembrane cleavage of integral membrane proteins such as Notch receptors. It
487 is involved in various diseases including a familial early-onset form of Alzheimer disease and a
488 form of frontotemporal dementia (see [https://www.uniprot.org/uniprotkb/](https://www.uniprot.org/uniprotkb/P49768/entry#function)
489 [P49768/entry#function](https://www.uniprot.org/uniprotkb/P49768/entry#function)).

490 The AlphaFold-predicted structure consists of two closely packed helical domains with
491 most deleterious mutations clustered inside the core domain while the non-conserved residues
492 are mostly located on the N-terminal loop (residues 1-70) and the inter-domain linker (residues
493 300-370) (see Fig 2c). The active site [56] (D257 and D385) is also located in the core domain
494 (colored green in Fig 2c). The contact maps predicted by three coevolution analysis methods all
495 resemble the contact map based on the predicted structure (see Fig 2a) (except that RaptorX
496 omitted local contacts in residues 1-100). As a result, the network scores agree well between
497 them in the helical domains (see Fig 2c), but with more differences in the flexible regions
498 (residues 1-70 and 300-370). Reassuringly, those residues identified by high network scores are

499 primarily clustered within the central core overlapping with most deleterious mutations, while
500 the flexible N-terminal and linker feature low scores consistent with low sequence conservation
501 (see Fig 2c). Among those top 10% hotspot residues (see Fig 2c), mutations at C92, V96, A231,
502 M233, L235, A246, L250, S390, L392, and C410 were found to cause loss of function and
503 altered amyloid-beta production [57] : C92S led to loss of protease function and increased
504 Abeta42 levels. V96F caused loss of protease activity. A231T/V and M233T led to decreased
505 protease activity, altered amyloid-beta production and increased amyloid-beta 42/amyloid-beta
506 40 ratio. L235P/R and S390I abolished protease activity. A246E and L250S abolished protease
507 activity and increased amyloid-beta 42/amyloid-beta 40 ratio. L392V resulted in reduced
508 proteolysis, altered amyloid-beta production and increased amyloid-beta 42/amyloid-beta 40
509 ratio. C410I reduced proteolysis. Since most of these residues are not near the active site, their
510 effects on protease activity are likely allosteric.



511

512 **Figure 2. Results for Presenilin-1 (Uniprot id: P49768):** (a) Four contact maps constructed
 513 from coevolution analysis by DeepMetaPSICOV (DMP), RaptorX (RX), SPOT-Contact (SC),
 514 and the predicted structure by AlphaFold (AF) (only those contacts with probability >0.1 are
 515 shown). (b) W_3 scores for all residue positions based on the contact maps in (a), where red and
 516 blue dots mark residues with deleterious and neutral mutations, respectively. (c) Predicted
 517 structure by AlphaFold as colored by W_3 scores (red/blue for high/low values), where residues
 518 with deleterious and neutral mutations are shown as large and small balls, respectively (A79,
 519 V82, C92, V96, F105, L113, Y115, T116, P117, E120, N135, M139, I143, M146, T147, H163,
 520 W165, L166, S169, L171, L173, L174, G206, G209, I213, L219, A231, M233, L235, A246,

521 L250, A260, L262, C263, P264, G266, P267, R269, L271, R278, E280, L282, A285, L286,
522 S289, D333, G378, G384, S390, L392, N405, A409, C410, A426, A431, P436, see
523 <https://www.uniprot.org/uniprotkb/P49768/variant-viewer>), and active-site residues are colored
524 in green. (d) Four other network scores (MSF, $\delta\lambda$, C5 and C12) for all residue positions based on
525 the contact maps in (a).

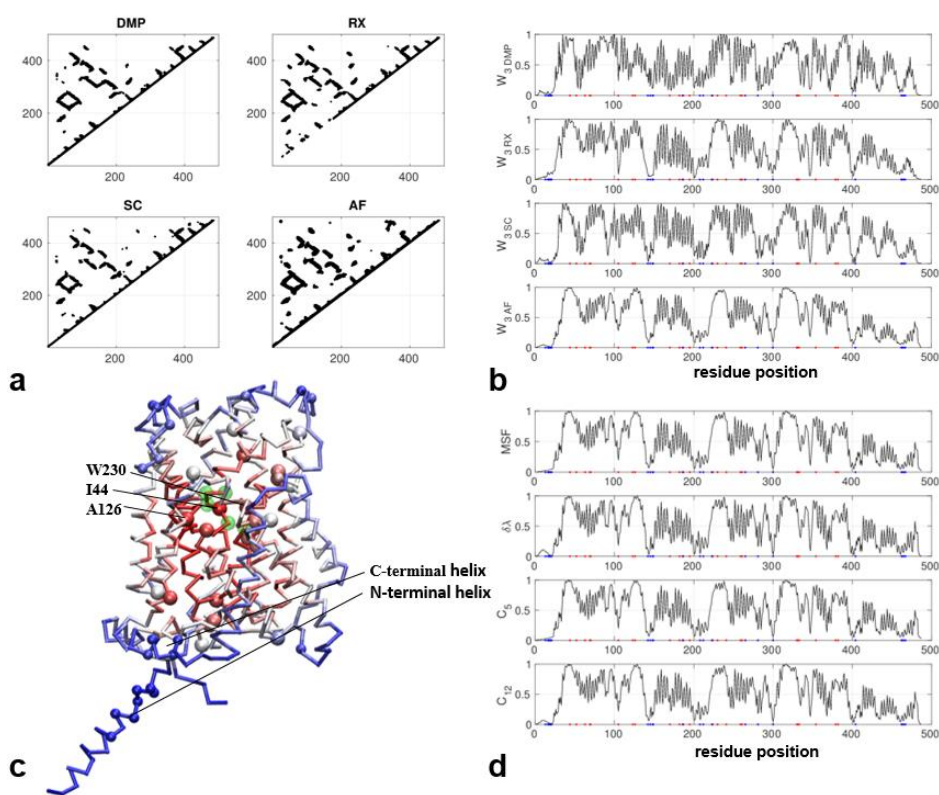
526

527 **b(0,+)-type amino acid transporter 1** (Uniprot id: P82251): It forms a functional
528 transporter complex that mediates the electrogenic exchange between cationic amino acids and
529 neutral amino acids. Its dysfunction is linked to Cystinuria, an autosomal disorder characterized
530 by impaired epithelial cell transport of cystine and dibasic amino acids in the proximal renal
531 tubule and gastrointestinal tract (see <https://www.uniprot.org/uniprotkb/P82251/entry#function>).

532 The AlphaFold-predicted structure consists of a helical domain with deleterious
533 mutations concentrating inside the core domain while those non-conserved residues are mostly
534 located on the domain periphery (e.g. N-terminal and C-terminal helices) (see Fig 3c). The active
535 site consists of residues 43-47 and 233 and is also located in the core domain (colored green in
536 Fig 3c). The contact maps predicted by three coevolution analysis tools are all similar to the
537 contact map based on the AlphaFold structure (see Fig 3a) (except that RaptorX omitted some
538 local contacts in residues 1-200). As a result, the network scores agree well between these
539 methods (see Fig 3b). Reassuringly, those residues identified with high network scores are
540 primarily within the central core and overlap with most deleterious mutations, while the
541 peripheral regions feature low scores consistent with low sequence conservation. Among those
542 top-10% hotspot residues (see Fig 3c), mutations I44T, A126T, and W230R were implicated in
543 Cystinuria. *In vitro* measurements showed W230R has almost no transport activity, and it was

Formatted: Space After: 8 pt

544 proposed that W230 serves as a gate between two substrate-binding pockets and undergoes
 545 conformational changes to enable amino acid transport [58] . Although the A126T mutation is
 546 mildly dysfunctional [59], it is notable among a cluster of conserved residues with small
 547 sidechains in the contact regions of transmembrane helices, hinting for its possible role in helix-
 548 helix association and relative motions.



549 **Figure 3. Results for amino acid transporter 1 (Uniprot id: P82251):** (a) Four contact maps
 550 constructed from coevolution analysis by DeepMetaPSICOV (DMP), RaptorX (RX), SPOT-
 551 Contact (SC), and the predicted structure by AlphaFold (AF) (only those contacts with
 552 probability >0.1 are shown). (b) W_3 scores for all residue positions based on the contact maps in
 553

554 (a), where red and blue dots mark residues with deleterious and neutral mutations, respectively.
555 (c) Predicted structure by AlphaFold as colored by W_3 scores (red/blue for high/low values),
556 where residues with deleterious and neutral mutations are shown as large and small balls,
557 respectively (V142,L223,I44,P52,G63,W69,A70,G105,T123,A126,V170,A182,I187,G195,
558 A224,W230,I241,G259,P261,V330,A331,R333,A354,S379,A382, see
559 <https://www.uniprot.org/uniprotkb/P82251/variant-viewer>), and active-site residues are colored
560 in green. (d) Four other network scores (MSF, $\delta\lambda$, C5 and C12) for all residue positions based on
561 the contact maps in (a).

562

563

564

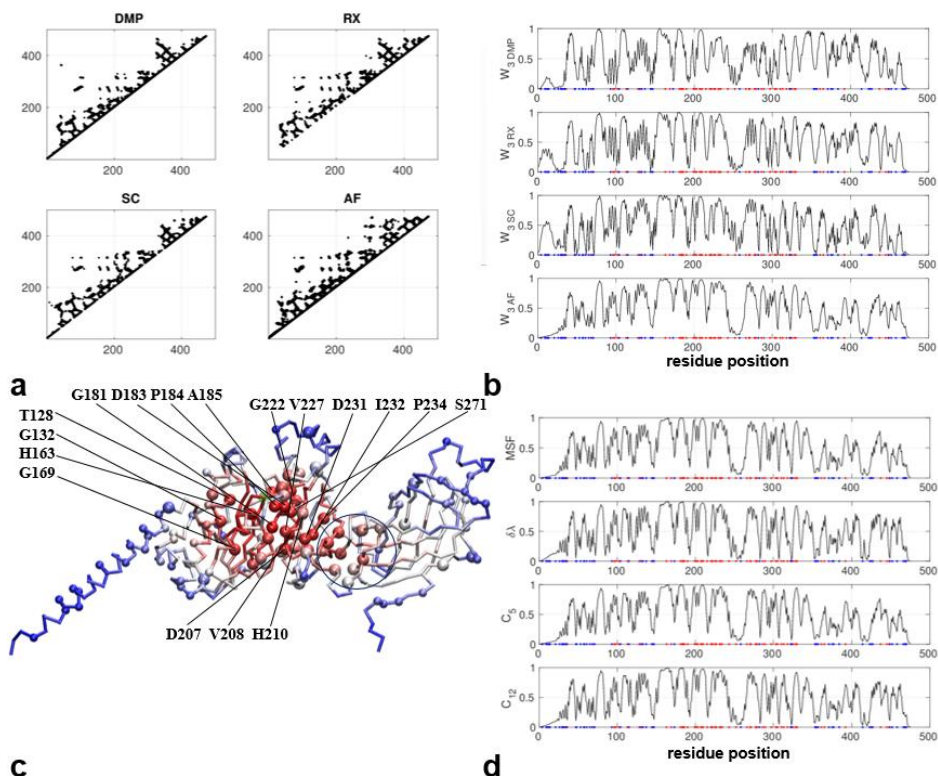
565 **Lipoprotein lipase** (Uniprot: P06858): As a key enzyme in triglyceride metabolism, it
566 catalyzes the hydrolysis of triglycerides from circulating chylomicrons and very low density
567 lipoproteins, thus playing an important role in lipid clearance from the blood stream, lipid
568 utilization and storage (see <https://www.uniprot.org/uniprotkb/P06858/entry#function>).

569 The AlphaFold-predicted structure consists of an N-terminal helix, a central α/β domains,
570 and a C-terminal β domain. Most deleterious mutations are concentrated inside the central
571 domain while the non-conserved residues are mostly located on the periphery (including e.g. N-
572 terminal helix and C-terminal domain) (see Fig 4c). The active site is comprised of a catalytic
573 triad of S159, D183, and H268 [60] in the central domain (colored green in Fig 4c). The contact
574 maps predicted by three coevolution analysis methods are similar to the contact map based on
575 the AlphaFold structure (see Fig 4a). As a result, the network scores agree well between these
576 methods (see Fig 4b) with minor differences in peripheral regions (e.g. insuch as the N-terminal

Formatted

Formatted: Indent: First line: 0"

577 helix). As predicted, those residues identified with high network scores are primarily within the
578 central domain overlapping with most deleterious mutations, while the peripheral N-terminal
579 helix and C-terminal domain feature low scores consistent with low sequence conservation (see
580 Fig 4c). Notably, some of them are found at the interface between the central domain and the C-
581 terminal domain (circled in Fig 4c), possibly mediating inter-domain motions. Among those top-
582 10% hotspot residues (see Fig 4c), T128, G132, H163, G169, G181, D183, P184, A185, D207,
583 V208, H210, G222, V227, D231, I232, P234 and S271 are known to harbor pathogenic
584 mutations in Hyperlipoproteinemia 1, an autosomal recessive metabolic disorder characterized
585 by defective breakdown of dietary fats. Both H163 and G169 lie in helix 4 that constitutes part of
586 the highly conserved beta-epsilon serine-alpha folding motif which is near S159 of the active
587 site. Supporting their functional relevance, mutations H163R and G169E were found to abolish
588 the enzymatic activity [61] . Near D183 (one of the catalytic triad), mutations G181S and P184R
589 were found to abolish the catalytic activity [62] . Further from D183, conserved substitutions
590 D207E and H210Q abolished the enzyme activity [63], and mutations D231E, I232S and P234L
591 led to loss of the catalytic function [64] . These mutations may disrupt allosteric interactions with
592 the central catalytic domain. Another conservative mutation S271T (near D183) also led to loss
593 of enzyme activity [65]. Taken together, these residues may function by directly or indirectly
594 coupling to the active site.



595 **c**

596 **Figure 4. Results for Lipoprotein lipase (Uniprot id: P06858):** (a) Four contact maps

597 constructed from coevolution analysis by DeepMetaPSICOV (DMP), RaptorX (RX), SPOT-

598 Contact (SC), and the predicted structure by AlphaFold (AF) (only those contacts with

599 probability >0.1 are shown). (b) W₃ scores for all residue positions based on the contact maps in

600 (a), where red and blue dots mark residues with deleterious and neutral mutations, respectively.

601 (c) Predicted structure by AlphaFold as colored by W₃ scores (red/blue for high/low values),

602 where residues with deleterious and neutral mutations are shown as large and small balls,

603 respectively

604 (H71,A427,D36,N70,V96,A98,R102,W113,T128,G132,H163,G169,G181,D183,P184,A185,

605 G186,E190,S199,D201,A203,D207,V208,H210,G215, S220,I221,G222, K225,V227,D231,
606 I232,P234,C243,I252,C266,R270,S271,D277,S278,L279,S286,Y289,F297,L303,C305,
607 C310,L313,N318,S325,M328,L330,A361,S365,L392,E437,E437,C445,E448, see
608 <https://www.uniprot.org/uniprotkb/P06858/variant-viewer>), and active-site residues are colored
609 in green. (d) Four other network scores (MSF, $\delta\lambda$, C5 and C12) for all residue positions based on
610 the contact maps in (a).

611

612 **Conclusion**

613 To conclude, we have combined machine learning, network analysis, and protein
614 language models to develop a sequences-based variant site prediction method based on the
615 protein residue contact networks which incorporate sequential, structural, dynamic, and
616 interaction information simultaneously:

Formatted: Font color: Auto

617 1. We build protein residue networks by exploiting three different state-of-the-art coevolution
618 analysis tools (RaptorX, DeepMetaPSICOV, and SPOT-Contact) that complement each other.

Formatted: Font color: Auto

619 2. We use three powerful machine learning algorithms (Random Forest, Gradient Boosting, and
620 Extreme Gradient Boosting) to optimally combine 20 network centrality scores to accurately
621 predict key residues as hot spots for disease mutations.

622 3. We train and validate our method using a dataset of 107 proteins rich in disease mutations,
623 demonstrating its high accuracy in distinguishing between deleterious and neutral sites (with
624 AUC of ROC ~ 0.84). Further improvement can be achieved after combining our method with
625 the ESM-based method.

626 This study has established a useful strategy of combining an ensemble of network scores
627 based on different coevolution analysis methods via machine learning to predict key variants
628 sites of relevance to disease mutations. The code and dataset are made available to public to
629 enable future developments and applications (see <https://simtk.org/projects/hotspots>).

630 For future work, it will be interesting to go beyond contact map predictions by integrating
631 other scores derived from the co-evolution analysis (for example, see refs [66-68]) in our
632 workflow, which may further boost the accuracy of variant site prediction.

633

634

635 References

- 636 reference
- 637
- 638 1. Jumper J, Evans R, Pritzel A, Green T, Figurnov M, et al. (2021) Highly accurate protein structure
639 prediction with AlphaFold. *Nature* 596: 583-589.
 - 640 2. Baek M, DiMaio F, Anishchenko I, Dauparas J, Ovchinnikov S, et al. (2021) Accurate prediction of
641 protein structures and interactions using a three-track neural network. *Science* 373: 871-876.
 - 642 3. Terwilliger TC, Liebschner D, Croll TI, Williams CJ, McCoy AJ, et al. (2023) AlphaFold predictions are
643 valuable hypotheses and accelerate but do not replace experimental structure determination.
644 *Nat Methods*.
 - 645 4. Al-Janabi A (2022) Has DeepMind's AlphaFold solved the protein folding problem? *Biotechniques* 72:
646 73-76.
 - 647 5. Medina E, D RL, Sanabria H (2021) Unraveling protein's structural dynamics: from configurational
648 dynamics to ensemble switching guides functional mesoscale assemblies. *Curr Opin Struct Biol*
649 66: 129-138.
 - 650 6. Pak MA, Markhieva KA, Novikova MS, Petrov DS, Vorobyev IS, et al. (2023) Using AlphaFold to predict
651 the impact of single mutations on protein stability and function. *PLoS One* 18: e0282689.
 - 652 7. Steinegger M, Soding J (2018) Clustering huge protein sequence sets in linear time. *Nat Commun* 9:
653 2542.
 - 654 8. Burley SK, Berman HM, Bhikadiya C, Bi C, Chen L, et al. (2019) RCSB Protein Data Bank: biological
655 macromolecular structures enabling research and education in fundamental biology,
656 biomedicine, biotechnology and energy. *Nucleic Acids Res* 47: D464-D474.
 - 657 9. Varadi M, Anyango S, Deshpande M, Nair S, Natassia C, et al. (2022) AlphaFold Protein Structure
658 Database: massively expanding the structural coverage of protein-sequence space with high-
659 accuracy models. *Nucleic Acids Res* 50: D439-D444.
 - 660 10. Bepler T, Berger B (2021) Learning the protein language: Evolution, structure, and function. *Cell Syst*
661 12: 654-669 e653.
 - 662 11. Cheng J, Novati G, Pan J, Bycroft C, Zemgulyte A, et al. (2023) Accurate proteome-wide missense
663 variant effect prediction with AlphaMissense. *Science* 381: eadg7492.
 - 664 12. Hassan MS, Shaalan AA, Dessouky MI, Abdelnaiem AE, ElHefnawi M (2019) A review study:
665 Computational techniques for expecting the impact of non-synonymous single nucleotide
666 variants in human diseases. *Gene* 680: 20-33.
 - 667 13. Niroula A, Urolagin S, Vihinen M (2015) PON-P2: prediction method for fast and reliable
668 identification of harmful variants. *PLoS One* 10: e0117380.
 - 669 14. Pejaver V, Urresti J, Lugo-Martinez J, Pagel KA, Lin GN, et al. (2020) Inferring the molecular and
670 phenotypic impact of amino acid variants with MutPred2. *Nat Commun* 11: 5918.
 - 671 15. Qi H, Zhang H, Zhao Y, Chen C, Long JJ, et al. (2021) MVP predicts the pathogenicity of missense
672 variants by deep learning. *Nat Commun* 12: 510.
 - 673 16. Yue P, Melamud E, Moulton J (2006) SNPs3D: candidate gene and SNP selection for association studies.
674 *BMC Bioinformatics* 7: 166.
 - 675 17. Tang H, Thomas PD (2016) Tools for Predicting the Functional Impact of Nonsynonymous Genetic
676 Variation. *Genetics* 203: 635-647.
 - 677 18. Katsonis P, Koire A, Wilson SJ, Hsu TK, Lua RC, et al. (2014) Single nucleotide variations: biological
678 impact and theoretical interpretation. *Protein Sci* 23: 1650-1666.

- 679 19. Singh A, Olowoyeye A, Baenziger PH, Dantzer J, Kann MG, et al. (2008) MutDB: update on
680 development of tools for the biochemical analysis of genetic variation. *Nucleic Acids Res* 36:
681 D815-819.
- 682 20. Bromberg Y, Rost B (2007) SNAP: predict effect of non-synonymous polymorphisms on function.
683 *Nucleic Acids Res* 35: 3823-3835.
- 684 21. Ng PC, Henikoff S (2003) SIFT: Predicting amino acid changes that affect protein function. *Nucleic*
685 *Acids Res* 31: 3812-3814.
- 686 22. Pers TH, Timshel P, Hirschhorn JN (2015) SNPsnap: a Web-based tool for identification and
687 annotation of matched SNPs. *Bioinformatics* 31: 418-420.
- 688 23. Zheng W, Tekpinar M (2009) Large-scale evaluation of dynamically important residues in proteins
689 predicted by the perturbation analysis of a coarse-grained elastic model. *BMC Struct Biol* 9: 45.
- 690 24. Zheng W, Brooks BR, Doniach S, Thirumalai D (2005) Network of dynamically important residues in
691 the open/closed transition in polymerases is strongly conserved. *Structure* 13: 565-577.
- 692 25. Ponzoni L, Bahar I (2018) Structural dynamics is a determinant of the functional significance of
693 missense variants. *Proc Natl Acad Sci U S A* 115: 4164-4169.
- 694 26. Butler BM, Gereke ZN, Kumar S, Ozkan SB (2015) Conformational dynamics of nonsynonymous
695 variants at protein interfaces reveals disease association. *Proteins* 83: 428-435.
- 696 27. Jumper J, Evans R, Pritzel A, Green T, Figurnov M, et al. (2021) Applying and improving AlphaFold at
697 CASP14. *Proteins* 89: 1711-1721.
- 698 28. Marks DS, Hopf TA, Sander C (2012) Protein structure prediction from sequence variation. *Nat*
699 *Biotechnol* 30: 1072-1080.
- 700 29. Morcos F, Pagnani A, Lunt B, Bertolino A, Marks DS, et al. (2011) Direct-coupling analysis of residue
701 coevolution captures native contacts across many protein families. *Proc Natl Acad Sci U S A* 108:
702 E1293-1301.
- 703 30. Hopf TA, Scharfe CP, Rodrigues JP, Green AG, Kohlbacher O, et al. (2014) Sequence co-evolution
704 gives 3D contacts and structures of protein complexes. *Elife* 3.
- 705 31. Marks DS, Colwell LJ, Sheridan R, Hopf TA, Pagnani A, et al. (2011) Protein 3D structure computed
706 from evolutionary sequence variation. *PLoS One* 6: e28766.
- 707 32. Burger L, van Nimwegen E (2010) Disentangling direct from indirect co-evolution of residues in
708 protein alignments. *PLoS Comput Biol* 6: e1000633.
- 709 33. Jones DT, Buchan DW, Cozzetto D, Pontil M (2012) PSICOV: precise structural contact prediction
710 using sparse inverse covariance estimation on large multiple sequence alignments.
711 *Bioinformatics* 28: 184-190.
- 712 34. Halabi N, Rivoire O, Leibler S, Ranganathan R (2009) Protein sectors: evolutionary units of three-
713 dimensional structure. *Cell* 138: 774-786.
- 714 35. Wang S, Sun S, Li Z, Zhang R, Xu J (2017) Accurate De Novo Prediction of Protein Contact Map by
715 Ultra-Deep Learning Model. *PLoS Comput Biol* 13: e1005324.
- 716 36. Ma J, Wang S, Wang Z, Xu J (2015) Protein contact prediction by integrating joint evolutionary
717 coupling analysis and supervised learning. *Bioinformatics* 31: 3506-3513.
- 718 37. Kandathil SM, Greener JG, Jones DT (2019) Prediction of interresidue contacts with
719 DeepMetaPSICOV in CASP13. *Proteins* 87: 1092-1099.
- 720 38. Jones DT, Singh T, Kosciolok T, Tetchner S (2015) MetaPSICOV: combining coevolution methods for
721 accurate prediction of contacts and long range hydrogen bonding in proteins. *Bioinformatics* 31:
722 999-1006.
- 723 39. Jones DT, Kandathil SM (2018) High precision in protein contact prediction using fully convolutional
724 neural networks and minimal sequence features. *Bioinformatics* 34: 3308-3315.

- 725 40. Hanson J, Paliwal K, Litfin T, Yang Y, Zhou Y (2018) Accurate prediction of protein contact maps by
726 coupling residual two-dimensional bidirectional long short-term memory with convolutional
727 neural networks. *Bioinformatics* 34: 4039-4045.
- 728 41. Yan W, Yu C, Chen J, Zhou J, Shen B (2020) ANCA: A Web Server for Amino Acid Networks
729 Construction and Analysis. *Front Mol Biosci* 7: 582702.
- 730 42. Amitai G, Shemesh A, Sitbon E, Shklar M, Netanel D, et al. (2004) Network analysis of protein
731 structures identifies functional residues. *J Mol Biol* 344: 1135-1146.
- 732 43. Velickovic P (2023) Everything is connected: Graph neural networks. *Curr Opin Struct Biol* 79:
733 102538.
- 734 44. Butler BM, Kazan IC, Kumar A, Ozkan SB (2018) Coevolving residues inform protein dynamics profiles
735 and disease susceptibility of nSNVs. *PLoS Comput Biol* 14: e1006626.
- 736 45. Chasman D, Adams RM (2001) Predicting the functional consequences of non-synonymous single
737 nucleotide polymorphisms: structure-based assessment of amino acid variation. *J Mol Biol* 307:
738 683-706.
- 739 46. Gerek ZN, Ozkan SB (2011) Change in allosteric network affects binding affinities of PDZ domains:
740 analysis through perturbation response scanning. *PLoS Comput Biol* 7: e1002154.
- 741 47. Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, et al. (2010) A method and server
742 for predicting damaging missense mutations. *Nat Methods* 7: 248-249.
- 743 48. Meier J, Rao R, Verkuil R, Liu J, Sercu T, et al. (2021) Language models enable zero-shot prediction of
744 the effects of mutations on protein function. *bioRxiv*: 2021.2007.2009.450648.
- 745 49. Vihinen M (2020) Problems in variation interpretation guidelines and in their implementation in
746 computational tools. *Mol Genet Genomic Med* 8: e1206.
- 747 50. Ashkenazy H, Erez E, Martz E, Pupko T, Ben-Tal N (2010) ConSurf 2010: calculating evolutionary
748 conservation in sequence and structure of proteins and nucleic acids. *Nucleic Acids Res* 38:
749 W529-533.
- 750 51. Zheng W, Brooks BR, Thirumalai D (2006) Low-frequency normal modes that describe allosteric
751 transitions in biological nanomachines are robust to sequence variations. *Proc Natl Acad Sci U S*
752 *A* 103: 7664-7669.
- 753 52. Zheng W (2016) Probing the structural dynamics of the SNARE recycling machine based on coarse-
754 grained modeling. *Proteins*.
- 755 53. Delgado J, Radusky LG, Cianferoni D, Serrano L (2019) FoldX 5.0: working with RNA, small molecules
756 and a new graphical interface. *Bioinformatics* 35: 4168-4169.
- 757 54. Veiga-da-Cunha M, Gerin I, Chen YT, de Barsey T, de Lonlay P, et al. (1998) A gene on chromosome
758 11q23 coding for a putative glucose- 6-phosphate translocase is mutated in glycogen-storage
759 disease types Ib and Ic. *Am J Hum Genet* 63: 976-983.
- 760 55. Hiraiwa H, Pan CJ, Lin B, Moses SW, Chou JY (1999) Inactivation of the glucose 6-phosphate
761 transporter causes glycogen storage disease type 1b. *J Biol Chem* 274: 5532-5536.
- 762 56. Wolfe MS, Xia W, Ostaszewski BL, Diehl TS, Kimberly WT, et al. (1999) Two transmembrane
763 aspartates in presenilin-1 required for presenilin endoproteolysis and gamma-secretase activity.
764 *Nature* 398: 513-517.
- 765 57. Sun L, Zhou R, Yang G, Shi Y (2017) Analysis of 138 pathogenic mutations in presenilin-1 on the in
766 vitro production of Abeta42 and Abeta40 peptides by gamma-secretase. *Proc Natl Acad Sci U S A*
767 114: E476-E485.
- 768 58. Yan R, Li Y, Shi Y, Zhou J, Lei J, et al. (2020) Cryo-EM structure of the human heteromeric amino acid
769 transporter b(0,+)-AT-rBAT. *Sci Adv* 6: eaay6379.
- 770 59. Font MA, Feliubadalo L, Estivill X, Nunes V, Golomb E, et al. (2001) Functional analysis of mutations
771 in SLC7A9, and genotype-phenotype correlation in non-Type I cystinuria. *Hum Mol Genet* 10:
772 305-316.

- 773 60. Emmerich J, Beg OU, Peterson J, Previato L, Brunzell JD, et al. (1992) Human lipoprotein lipase.
774 Analysis of the catalytic triad by site-directed mutagenesis of Ser-132, Asp-156, and His-241. *J*
775 *Biol Chem* 267: 4161-4165.
- 776 61. Reina M, Brunzell JD, Deeb SS (1992) Molecular basis of familial chylomicronemia: mutations in the
777 lipoprotein lipase and apolipoprotein C-II genes. *J Lipid Res* 33: 1823-1832.
- 778 62. Bruin T, Tuzgol S, van Diermen DE, Hoogerbrugge-van der Linden N, Brunzell JD, et al. (1993)
779 Recurrent pancreatitis and chylomicronemia in an extended Dutch kindred is caused by a
780 Gly154-->Ser substitution in lipoprotein lipase. *J Lipid Res* 34: 2109-2119.
- 781 63. Haubenwallner S, Horl G, Shachter NS, Presta E, Fried SK, et al. (1993) A novel missense mutation in
782 the gene for lipoprotein lipase resulting in a highly conservative amino acid substitution
783 (Asp180-->Glu) causes familial chylomicronemia (type I hyperlipoproteinemia). *Genomics* 18:
784 392-396.
- 785 64. Gotoda T, Yamada N, Kawamura M, Kozaki K, Mori N, et al. (1991) Heterogeneous mutations in the
786 human lipoprotein lipase gene in patients with familial lipoprotein lipase deficiency. *J Clin Invest*
787 88: 1856-1864.
- 788 65. Hata A, Emi M, Luc G, Basdevant A, Gambert P, et al. (1990) Compound heterozygote for lipoprotein
789 lipase deficiency: Ser----Thr244 and transition in 3' splice site of intron 2 (AG----AA) in the
790 lipoprotein lipase gene. *Am J Hum Genet* 47: 721-726.
- 791 66. Bisardi M, Rodriguez-Rivas J, Zamponi F, Weigt M (2022) Modeling Sequence-Space Exploration and
792 Emergence of Epistatic Signals in Protein Evolution. *Mol Biol Evol* 39.
- 793 67. Cocco S, Feinauer C, Figliuzzi M, Monasson R, Weigt M (2018) Inverse statistical physics of protein
794 sequences: a key issues review. *Rep Prog Phys* 81: 032601.
- 795 68. Rodriguez-Rivas J, Croce G, Muscat M, Weigt M (2022) Epistatic models predict mutable sites in
796 SARS-CoV-2 proteins and epitopes. *Proc Natl Acad Sci U S A* 119.

797

798

799

800 Supporting Information

Formatted: Font: 18 pt

801 S1 Table. Evaluation of 20 network scores based on protein residue contact maps

802 constructed from 3 coevolution analysis tools (DeepMetaPSICOV, RaptorX, and SPOT-

803 Contact)

Score	AUC* of DeepMetaPSICOV	AUC* of RaptorX	AUC* of SPOT-Contact
<u>C1</u>	<u>0.75±0.13</u>	<u>0.78±0.14</u>	<u>0.75±0.15</u>
<u>C2</u>	<u>0.75±0.17</u>	<u>0.75±0.19</u>	<u>0.76±0.19</u>
<u>C3</u>	<u>0.78±0.10</u>	<u>0.75±0.15</u>	<u>0.70±0.15</u>
<u>C4</u>	<u>0.65±0.12</u>	<u>0.52±0.15</u>	<u>0.61±0.07</u>
<u>C5</u>	<u>0.78±0.17</u>	<u>0.78±0.17</u>	<u>0.78±0.18</u>
<u>C6</u>	<u>0.63±0.16</u>	<u>0.56±0.19</u>	<u>0.65±0.17</u>
<u>C7</u>	<u>0.77±0.11</u>	<u>0.61±0.20</u>	<u>0.72±0.16</u>
<u>C8</u>	<u>0.65±0.12</u>	<u>0.52±0.15</u>	<u>0.61±0.07</u>
<u>C9</u>	<u>0.79±0.13</u>	<u>0.78±0.16</u>	<u>0.75±0.16</u>
<u>C10</u>	<u>0.76±0.12</u>	<u>0.75±0.14</u>	<u>0.69±0.15</u>
<u>C11</u>	<u>0.80±0.13</u>	<u>0.78±0.17</u>	<u>0.78±0.17</u>
<u>C12</u>	<u>0.78±0.16</u>	<u>0.80±0.14</u>	<u>0.80±0.17</u>
<u>C13</u>	<u>0.75±0.17</u>	<u>0.74±0.18</u>	<u>0.76±0.19</u>
<u>δλ</u>	<u>0.80±0.13</u>	<u>0.78±0.14</u>	<u>0.78±0.16</u>
<u>MSF</u>	<u>0.81±0.14</u>	<u>0.79±0.15</u>	<u>0.80±0.17</u>
<u>W₁</u>	<u>0.81±0.14</u>	<u>0.79±0.15</u>	<u>0.80±0.17</u>
<u>W₂</u>	<u>0.81±0.15</u>	<u>0.77±0.15</u>	<u>0.79±0.17</u>
<u>W₃</u>	<u>0.81±0.14</u>	<u>0.78±0.15</u>	<u>0.80±0.16</u>
<u>W_∞</u>	<u>0.80±0.13</u>	<u>0.77±0.14</u>	<u>0.78±0.16</u>
<u>W_s</u>	<u>0.80±0.13</u>	<u>0.78±0.13</u>	<u>0.78±0.16</u>

Formatted Table

804 * mean ± standard-deviation

Formatted: Subscript

Response to Editor

I have addressed the following additional journal requirements:

1. I have reformatted the manuscript to comply with PLOS ONE's style requirements (journals.plos.org/plosone/s/submission-guidelines), including those for file naming.
2. No human participants are involved in this study.
3. To comply with PLOS ONE guidelines on code sharing I will make the code available at <https://simtk.org/projects/hotspots> .
4. Since I use WORD to write the manuscript, the PLOS LaTeX template is not applicable.
5. I have corrected the grant information so it matches between the 'Funding Information' and 'Financial Disclosure' sections. I also provide the correct grant numbers (3R01NS108750) for the awards in the 'Funding Information' section.
6. In the financial disclosure "This study is funded by a grant from NIH...", I would like to add this amended Role of Funder statement: "The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript."
7. I have made all relevant data available without restriction at <https://simtk.org/projects/hotspots> .

I have revised the paper to address the following additional Editor Comments: In particular, the abstract and introduction should be rewritten based on the provided comments. In addition, the proposed methodology should be better described and justified. Finally results should be better presented as suggested by reviewers.

Response to Reviewers

Response to Reviewer #1

We thank reviewer 1 for the constructive comments and suggestions!

Abstract

•The abstract in this paper did not adequately capture the details expected in an abstract. Background to the research was not introduced. Problem being addressed was not efficiently stated. Previous methods that have tried addressing the problems were not highlighted. For instance, the author stated that “To meet this challenge, we build upon recent progress in machine learning, network analysis, and protein language models, and develop a ...” without actually highlighting the previous work done. Majority of the content in the abstract are on the authors’ finding. The author is expected to provide a basic or simple background to the research. Followed by a brief description of what has been previously done before stating what the problems are from what has been done. The author can be more specific on the keywords which shouldn’t be more than 6.

We have revised Abstract to include more background information (including problem to address, previous methods and their limitations). We have reduced the number of keywords to no more than six.

Introduction

•The authors should minimize the use of ‘e.g’ both in the introduction section and the abstract.
We have removed most ‘e.g’ in the paper.

•Citation in the body of the introduction is quite few. More citations and references should be provided.

We have added more references in Introduction (with total 48 references cited).

•Kindly provide justification for this claim “...it is now feasible to predict static structures for many proteins of interest given their sequences”.

We have added new references to substantiate the above claim with some caveats (line 42).

•It appears that the last paragraph of the introduction is exactly same as some major part of the abstract word for word. Authors are encouraged to avoid self-plagiarism.

We have rewritten this paragraph and Abstract to avoid duplications between them.

•Similar to the abstract, information in the introductory part of this paper is insufficient.

We have tried our best to give a detailed introduction to the background and literature relevant to this study (with total 48 references cited). We will appreciate it if the reviewer could kindly give more specific comments if any information is still missing.

•Also, authors are advised to provide a brief description of the different sections of the manuscript at the end of the introduction.

We have added a brief outline of the Methods and Results sections at the end of Introduction (line 125).

Materials and Methods

- The subsequent section after the introduction should be captioned “Materials and Methods” as opposed to the caption used by the author.

We have renamed the Method section to “Materials and Methods”.

- Each section and subsection should be numbered accordingly. Currently, sections and subsections are not identifiable.

We have numbered all subsections of “Materials and Methods” in the order of the proposed workflow.

- Furthermore, for easier flow and understanding of the methodology, a framework or an algorithm of the methodology could be added. This would provide readers with a conceptualized view of the methodology.

We have added a summary of the workflow of our method at the beginning of “Materials and Methods” (line 131).

- The uniprot Id of the 107 protein sequences collected can be provided as a supplementary file.

Detailed information about the dataset of 107 protein sequences and variants is available at the following site: <https://simtk.org/projects/hotspots>

- The author claimed to have used Random Forest, Gradient Boosting Classifier and Extreme Gradient Boosting Classifier. Kindly provide a justification for the choice of this machine learning methods

These three methods were chosen because they have performed successfully in machine learning contests in Kaggle. They are also relatively cheap to train and optimize compared with the deep learning methods (see line 276).

- In addition, what parameters were tuned by the author? Did the author used Optuna for all the ML methods. How is the result prior to the use of hyper-parameter tuning technique?

We added details of hyper-parameters at lines 266, 270, 274. Yes, we used Optuna for each of the three ML methods. While the resulting improvement is modest (relative to the default parameters), it is a common practice in ML to perform task-specific hyper-parameters tuning. We have run Optuna multiple times to ensure the resulting best metric is reproducible.

- Why was AUC of ROC used as the metric for assessing prediction quality?

The AUC is a standard metric for evaluating binary classifiers based on the ROC curve of sensitivity and specificity. The ROC curves are also used in other computational papers for variant predictions (see line 162).

Results and Discussion

- Most importantly, results obtained by the author should be presented while being discussed instead of being added to the supplementary or being placed at the end of the paper. This makes it difficult to understand the result being discussed. Some of the figures should also be presented as they are being discussed.

We have moved figures and tables to where they were discussed.

Conclusion

- The conclusion section is completely missing in this paper. Author is encourage to provide a conclusion section alongside a summarized and itemized key findings from the research.

We have added Conclusion to summarize our key findings.

Response to Reviewer #2

We thank reviewer 2 for the constructive comments and suggestions!

Major points

- Sequences-based coevolution methods return a prediction of the protein structure from multiple sequence alignments (MSAs). It is not clear to me how these MSAs (or a unique MSA?) are designed. The authors sometimes refer to them (or to it), and sometimes they put the focus on the 107 sequences used in the training and testing of the binary classifier but it is not clear to me how they are chosen, and if they belong to the initial MSA or not. I have appreciated the final focus on four sequences, even though the performances, and a comparison to other techniques, of the new pipeline (for instance the AUC) should be included in the discussion.

To clarify, the coevolution-based methods that we used (RaptorX, DeepMetaPSICOV, and SPOT-Contact) take a protein sequence as input and then predict a residue-residue contact probability matrix (not the protein structure) based on MSA. Our method does not utilize the MSA directly. The 107 protein sequences were chosen from the HumVar database (see line 143), and they do not belong to a specific MSA. Instead, each of these sequences was used to build its own MSA by the above coevolution-based methods.

We have compared our sequences-based method with alternative methods based on structures (AlphaFold) or physical force fields (FoldX) or protein language model (ESM) using a testing set of 28 proteins (see Tables 2 and 4). We also added additional comparison with PolyPhen-2 (see line 397).

- Linked to the first question, if the analysis is performed on a unique dataset, i.e. a unique MSA of a protein family, I encourage the authors to repeat these experiments on different protein domains, as a stress test to their new pipeline.

To clarify, our evaluation is not performed on a unique MSA of a particular protein family. Instead, it is based on a diverse dataset of 107 proteins from various different protein families, which contain 97 dissimilar proteins (with their pairwise sequence identity < 30%) (line 147).

- In the third step of the workflow, the authors split the 107 protein sequences into a train and a test set. How is this subdivision decided? One should in principle check whether the sequences in the two sets are sufficiently “distant” (for instance using clustering analysis) or repeat the procedure for different assignments into the two subsets. Additionally, how does the choice of the size of these two sets affect the results?

The train/test split is purely random: 79 for training and 28 for testing. We have checked the sequence similarity between the two sets, only 1 training-set protein has sequence identity >30% with proteins of the testing set. In fact, 97/107 proteins have pairwise sequence identity <30%. We repeated training and testing on this reduced dataset of 97 dissimilar proteins, and the resulting AUCs are only slightly higher (by <0.03), suggesting that the few similar sequences did not markedly change the training/testing results.

The choice of train/test split ratio (~75/25) is based on common practice of ML (see <https://onlinelibrary.wiley.com/doi/full/10.1002/sam.11583>). We also tried other sizes of testing set (18 and 38) which only slightly increase the resulting AUCs by 0.01~0.04. So the training/testing results are not sensitive to this choice.

- Coevolution methods, together with a proxy for pairwise contact prediction, allow for an estimate of the degree of deleteriousness of point and pairwise mutations. In recent developments (see Refs. 37-42 in [1] and also more recent works in [2-3]), the authors show that the sequence “energy” of the coevolution models can be interpreted as (negative) protein fitness which, indeed, correlates well with deep mutational scanning-based measures. Since this coevolution information is exploited in the first step of the presented pipeline, I would like to ask the following questions:

- Can the authors display how well these methods alone perform compared to the author's workflow?

Although it would be useful to use ΔE_DCA defined in ref[3] to predict variant effects, such calculation is not supported by the co-evolution-based contact prediction methods used here (RaptorX, DeepMetaPSICOV, and SPOT-Contact), and is therefore beyond the scope of this study. This study is limited to the use of residue contact maps as predicted by co-evolution analysis, rather than fully exploiting all scores derived from the co-evolution analysis of MSA.

- Can this information be integrated (together with the pure network scores) in the second and third steps?

Yes, we expect potentially fruitful integration of our network-based methods with other co-evolution-based scores like ΔE_DCA to improve the prediction of variant effects. We have mentioned this and cited ref [1-3] in the discussion of future work (see line 590).

- It is not clear to me why the twenty measures used in the second step (if I have understood correctly, they are associated with network centrality properties) are the correct (or the most informative) metrics to cope with sequence hotspots. I would expect that some important hotspots may be related to the interaction with other proteins, and, therefore, they may be “far” from the core of the folded protein network. Also, other information like the electrical-chemical properties of the amino acids is neglected. Can the authors comment on them?

We agree that not all hotspot residues can be predicted by our method which focused on intra-protein residue contacts, and some hotspot residues may not possess high centrality scores (as mentioned above by the reviewer). However, we have shown that our method is competitive with alternative methods (ESM and PolyPhen-2), supporting the value of centrality scores as informative predictors for disease mutations. Further, our method provides new predictive features that complement other scores (such as ΔE_DCA), and together they may enable more accurate predictions of variant effects (see line 590).

- Pag. 15. The evaluation of each single score (among the twenty?) against the different structures is not well described. The comparison is made according to which measure? The final binary classifier? If this is the case, does it mean that the authors compared the importance of each score independently and then all together?

For the evaluation of each single score, we sort all testing-set variants by that score and predict a variant deleterious/neutral if its score is above/below a cutoff value. This resulted in an ROC curve from which we calculated AUC (line 302).

After assessing the individual scores independently, we then used ML to combine them to train a binary classifier which was then assessed with the AUC of ROC (see Table 4).

- Overall, the presentation of the results is a little confusing to me. Comparisons seem to be made “internally” by changing the metrics in the second step, of the protein structure prediction in the first step. On page 19, the authors mention a comparison between their algorithm and Seq-GNM,

PolyPhen-2 but they run on different datasets, so the final performances may not be compared. I believe that a more “fair” comparison with other state-of-the-art techniques should be presented. We agree that fair comparisons with alternative methods like PolyPhen-2 are desirable. However, since our training/testing datasets were taken from HumVar, and PolyPhen-2 has been trained on this dataset, it is not possible to compare their performance objectively. With this caveat in mind, we have evaluated PolyPhen-2 based on the same 28 testing-set proteins and their variants, and found the AUC of PolyPhen-2 to be 0.85 (see line 397), which is close to our method (after combining 3 co-evolution analyses, see Table 4). Additionally, we have compared our method with another state-of-the-art method based on ESM, which gave AUC~0.83. Therefore, our method is competitive with these alternative methods. To be clear, our main goal is to complement rather than compete with existing methods. Indeed, we have shown that the combination of our method with ESM has yielded better performance than ESM alone (see Table 4).

- The ROC and AUC metrics are presented as final comparison metrics. Is the value presented in the manuscript an average value among the test sequences? If this is the case, what about the standard error associated with it?

To clarify, the AUCs in Table 2 are calculated from the ROC for all variants of the 28 testing set proteins. We also calculated AUCs for each protein alone, and their means and S.D. are shown in Table S1. The cumulative AUCs in Table 2 are comparable to the mean AUCs of individual proteins in Table S1.

Minor points

- Would it be possible to swap the Methods section and the Results section? Also, some parts are frequently repeated in both Methods and Results.

I am afraid not because the PLOS ONE format requires the Methods section precedes the Results section.

We have reduced methodological details in Results to avoid repetitions.

- Within the Methods section, in my opinion, the paper would gain readability if the sections followed the main pipeline of the method (now the third step is described before the first one).

We have numbered the subsections in the order of the workflow/pipeline (starting with the datasets). The workflow is summarized at line 131.

References

[1] Inverse statistical physics of protein sequences: a key issues review S Cocco, C Feinauer, M Figliuzzi, R Monasson, M Weigt Reports on Progress in Physics 81 (3), 032601

[2] Modeling sequence-space exploration and emergence of epistatic signals in protein evolution M Bisardi, J Rodriguez-Rivas, F Zamponi, M Weigt Molecular biology and evolution 39 (1), msab321

[3] Epistatic models predict mutable sites in SARS-CoV-2 proteins and epitopes J Rodriguez-Rivas, G Croce, M Muscat, M Weigt Proceedings of the National Academy of Sciences 119 (4), e2113118119



Click here to download Data Review URL
<http://simtk.org/projects/hotspots>

