# MAVISp: A Modular Structure-Based Framework for Protein Variant Effects

Matteo Arnaudi[1,2†], Mattia Utichi[1,2†], Kristine Degn[1,2†], Matteo Tiberti[1†], Ludovica Beltrame[1,2], Karolina Krzesińska[1,2], Pablo Sánchez-Izquierdo Besora[1,2], Eleni Kiachaki[1,2], Simone Scrima[1,2], Laura Bauer[1], Katrine Meldgård[1,2], Anna Melidi[1,2], Lorenzo Favaro[1,2], Anu Oswal[1,2], Guglielmo Tedeschi[6], Terézia Dorčaková[2,], Alberte Heering Estad[1,2], Joachim Breitenstein[1,2], Jordan Safer[4], Paraskevi Saridaki[1], Valentina Sora[1,2], Francesca Maselli[1,3], Philipp Becker[2], Jérémy Vinhas[1], Alberto Pettenella[1], Matteo Lambrughi[1], Claudia Cava[7], Anna Rohlin[8,9], Mef Nilbert[10,11], Sumaiya Iqbal[4], Peter Wad Sackett[2], Burcu Aykac Fas[5], Elena Papaleo[1,2]*

[1]Cancer Structural Biology, Danish Cancer Institute, 2100, Copenhagen, Denmark
[2]Cancer Systems Biology, Section for Bioinformatics, Department of Health and Technology, Technical University of Denmark, 2800, Lyngby, Denmark
[3]Institute of Molecular Bioimaging and Physiology, National Research Council (IBFM-CNR), Via F. Cervi 93, Segrate-Milan, 20054, Milan, Italy
[4]Bioinformatic and Computational Biology group, the Center for Development of Therapeutics, Broad Institute of MIT and Harvard, 415 Main St, Cambridge, MA 02142
[5]Laboratoire de Biochimie Théorique, CNRS (UPR9080), Université Paris Cité, F-75005 Paris, France
[6]Department of Biochemistry and Microbiology, University of Chemistry and Technology Prague, Prague 6 166 28, Czech Republic
[7]Department of Science, Technology and Society, Scuola Universitaria IUSS, Istituto Universitario di Studio Superiori, Piazza della Vittoria 15, 27100 Pavia, Italy
[8]Department of Clinical Genetics and Genomics, Sahlgrenska university hospital, Gothenburg, Sweden
[9] Department of Laboratory Medicine, Institute for Biomedicine, Sahlgrenska Academy, University of Gothenburg, Gothenburg, Sweden
[10]Institute of Clinical Medicine, Dept. Oncology and Pathology, Lund University, Sweden
[11] Department of Clinical Research, Copenhagen University Hospital Amager and Hvidovre, Copenhagen, Denmark

† These authors equally contributed to the work

* Corresponding author: Elena Papaleo, elpap@dtu.dk, elenap@cancer.dk

The role of genomic variants in disease has expanded significantly with the advent of advanced sequencing techniques. The rapid increase in identified genomic variants has led to many variants being classified as Variants of Uncertain Significance or as having conflicting evidence, posing challenges for their interpretation and characterization. Additionally, current methods for predicting pathogenic variants often lack insights into the underlying molecular mechanisms. Here, we introduce MAVISp (Multi-layered Assessment of VarIants by Structure for proteins), a modular structural framework for variant effects, accompanied by a web server (https://services.healthtech.dtu.dk/services/MAVISp-1.0/) to enhance data accessibility, consultation, and re-usability. MAVISp currently provides data over 1000 proteins, encompassing more than eight million variants. A team of biocurators regularly analyzes and updates protein entries using standardized workflows, incorporating free energy calculations or biomolecular simulations. We illustrate the utility of MAVISp through selected case studies. The framework facilitates the analysis of variant effects at the protein level and has the potential to advance the understanding and application of mutational data in disease research.

**Keywords:** variant effects, cancer genomics, protein structures, free energy calculations, protein stability, protein function, long-range structural communication

# 1 Introduction

We are witnessing unprecedented advances in cancer genomics, sequencing[1], structural biology[2], and high-throughput multiplex-based assays[3,4]. While sequencing approaches can identify alterations in the genome, understanding the molecular mechanisms of these variants remains a challenge. Although many variants in human genes associated with disease are currently known, the identification of their effects on human health is lagging behind[5]. Substantial evidence, which is necessary to classify variants according to their effects, is often lacking or contradictory in nature. Consequently, Variants of Uncertain Significance (VUS) or variants found to have conflicting evidence are continuously identified and reported in variant databases[67-11]. VUS remain an outstanding problem which complicate diagnosis and lead to suboptimal diagnosis or choice of therapy [12].

At the same time, the bioinformatics community has developed various approaches for predicting the impact of variants on human health, many of which are benchmarked against or complemented by experimental data and cellular readouts[13-17] In this context, experimental multiplex assays deliver good quality and high-throughput assessment of the effect of variants on different readouts and have effectively been used to aid clinical variant interpretation. [18,19]These computational and experimental approaches allow to classify variants for their potential pathogenic or benign effects, which are then reported in different repositories and compendia[7-10]. In fact, computational methods are currently considered supporting evidence for variant classification, according to recent revisions of the American College of Medical Genetics and Genomics/Association for Molecular Pathology (ACMG/AMP) variant classification guidelines[20]. Variant effect predictors (VEPs), methods designed to predict the effect of a mutation at the genome or protein level, have made considerable progress, as outlined in recent reviews[21-23]. VEPs have classically relied on sequence data and variants with known classifications.

Nonetheless, in recent years, the advent of AlphaFold2[2,24,25] and other similar methodologies has enabled the prediction of accurate three-dimensional (3D) protein structures and complexes, often with a quality comparable to experiments. This, in turn, enabled the inclusion of information about protein structure in machine learning models, which are among the best-performing available VEPs[21]. A well-known example of this is AlphaMissense[26], which is based on a deep learning model similar to AlphaFold2. Additionally, it simultaneously learns to perform structure prediction and trains an unsupervised protein language model, thereby incorporating structural information into the prediction. The latter was then fine-tuned for a variant classification task. Approaches based on protein language models (such as ESM-1b[27] or, more recently, ESM-2[28] and ESM-3[29]), which are unsupervised models of protein sequence, have also shown good performance when used in variant effect prediction tasks[29,30]. ESM-3[29] already incorporates structural information into its training, through specialized tokens, whereas protein sequence models have been used in conjunction with structural information in various ways[31,32]. Even a model such as GEMME, which is an epistatic model entirely based on sequence conservation, has been supplemented with structural information as structure-derived features in ESCOTT[33]. Rhapsody-2 is a VEP that incorporates features derived from protein structure and dynamics within a machine learning framework[34]. Finally, the ability to perform long and accurate biomolecular simulations and robust physical models allows the exploration of conformational changes and protein dynamics across different timescales[35].

In previous pilot projects, we explored structure-based methods to analyze the impact of variants in coding regions of cancer-related genes, focusing on their consequences on the protein product[36-38]. We propose that these methodologies could be widely applied to study disease-associated variants. When formalized and standardized, this approach can complement existing methods for predicting pathogenic variants, such as the aforementioned AlphaMissense[26]. Most available VEPs estimate the likelihood of damaging effects of variants, but do not provide evidence of variant effects in relation to specific altered protein functions at the cellular level. On the contrary, with this contribution, we aim to link the effects of variants to specific underlying molecular mechanisms[38]. A mechanistic understanding of variant effects can help the design of strategies in disease prevention, genetic counseling, clinical care, and treatment. Moreover, from a fundamental research perspective, mechanistic knowledge is also essential for designing and prioritizing experiments to investigate the underlying molecular causes of disease.

52 Considering this, we developed MAVISp (Multi-layered Assessment of VarIants by Structure for proteins) to
53 enable high-throughput variant analysis within standardized workflows. MAVISp integrates results from VEPs
54 and structure-based predictions of variant effects on several protein properties. The data are accessible
55 through a Streamlit-based website for consultation and download (https://services.healthtech.dtu.dk/ser-
56 vices/MAVISp-1.0/). Additionally, we maintain a Gitbook resource with detailed reports for individual proteins
57 (https://elelab.gitbook.io/mavisp/).

58 With this publication, we provide data on *in silico* saturation mutagenesis for all possible variants at each
59 mutation site with structural coverage for 1096 proteins and over eight million variants. New data and updates
60 of existing entries will be continuously released. Currently, we are capable of processing up to 20 new pro-
61 teins weekly, which are deposited in a local version of the database. The public database is updated quarterly.
62 Based on recent statistics (https://elelab.gitbook.io/mavisp/documentation/coverage-and-statistics), we an-
63 ticipate providing 80-100 new proteins with each update, along with additional modules for existing entries.
64 In this manuscript, we provide an overview of the methodology and show examples of data analysis and
65 application.

66
67

# Results

## *Overview of MAVISp and its database*

MAVISp performs a set of independent predictions, each assessing the effect of a specific amino acid substitution on a different aspect of protein function and structural stability, starting from one or more protein structures. These independent predictions are executed by the so-called MAVISp modules (**Fig. 1a**). MAVISp can be applied to individual three-dimensional (3D) protein structures and their complexes (*simple mode*) or to an ensemble of structures generated through various approaches (*ensemble mode*). The framework is modular, allowing all the modules or only a selected subset to be applied, depending on the case study. Each module relies on Snakemake, Dask workflows, or Python scripts, all of which are supported by specific virtual environments. The modules are divided into two main categories: (i) modules to retrieve and select structures for analyses (shown in orange in **Fig. 1a**), (ii) modules to perform analyses related to variant assessment or annotations (shown in blue in **Fig. 1a**). Each module includes a strictly defined protocol for computational analysis that can be carried out either step by step or automatically embedded in more comprehensive pipelines (Methods). They are designed to ensure consistency across all the proteins under investigation and to enhance reproducibility and repeatability. Our prediction modules are also complemented by available experimental data or already available predictions that can be integrated in the MAVISp dataset, such as those for VEPs (shown in green in **Fig. 1a**). All the resources used in the MAVISp framework are reported in **Table S1**, some of which have been developed within this work.

The modules are used in the context of the overall MAVISp workflow (**Fig. 1b**), which is designed to enable multiple biocurators to work concurrently and independently on distinct proteins. Data managers defined a priority list of targets that are analyzed in batches by biocurators, depending on the specific research project requirements. Additional targets of interest for the research community can be requested, as explained in the documentation on GitBook.
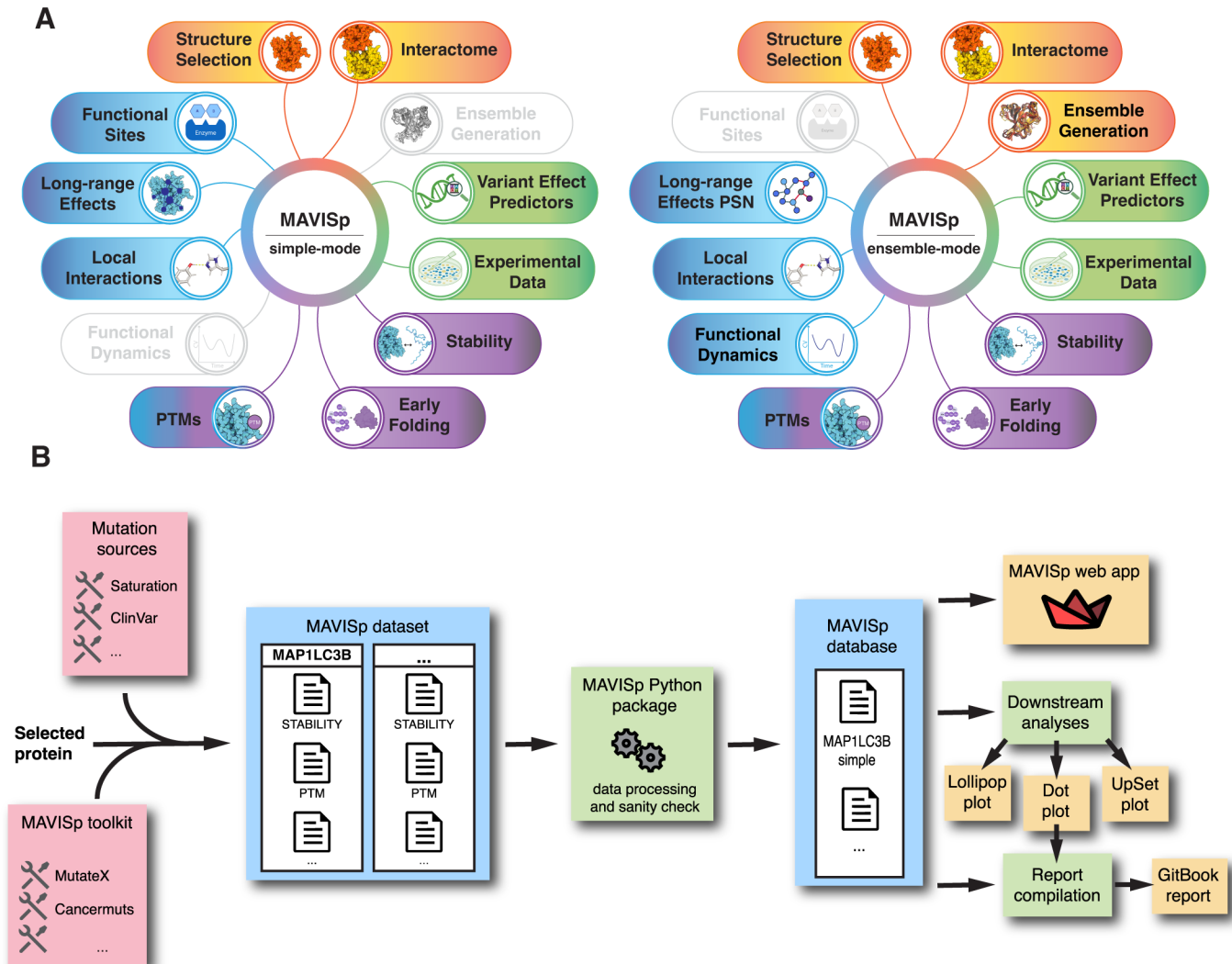
The workflow is designed as a set of consecutive steps that act on a protein of interest at a time. As the first step, once a protein of interest has been selected, a biocurator retrieves structural and functional information about it, along with key identifiers (e.g., gene name, UniProt AC, RefSeq identifier) for the next steps. Additionally, the biocurator proposes a trimming strategy for the protein, e.g., identifying one or more sets of contiguous residues in the protein structure that can effectively serve as input for the prediction steps. This step entails considering only well-structured and high-accuracy regions of our proteins, which is crucial since most MAVISp modules are not designed to handle large intrinsically disordered regions. In selected cases, to avoid potential bias in our structural calculations, the curator may edit the structure by removing long disordered inclusions in structured regions. Furthermore, in the MAVISp *ensemble mode*, where he ENSEMBLE GENERATION module should be carried out, the biocurator identifies the initial structures for the simulations to be performed on the protein target in its free or bound state with other biomolecules and performs the necessary simulations to obtain the final structural ensemble. Once the protein structure or structural ensemble, depending on the mode, is available, the biocurator works with each available module and obtains: i) a list of variants that MAVISp will annotate (see Materials and Methods for details) and ii) the final predictions for each module. To do so, biocurators adhere to strict workflows for data collection based on a set of procedures codified in each module, which is mostly automated via the use of Snakemake pipelines. Once this is completed, the MAVISp data managers will import and aggregate the data using the MAVISp Python package (https://github.com/ELELAB/MAVISp). This step also allows to perform sanity checks, per-module data classifications, and write the results in a human-readable table format, constituting the MAVISp database. The database files are the first product of MAVISp and contain the relevant collected data and metadata for each of the identified variants (https://services.healthtech.dtu.dk/services/MAVISp-1.0/).

The datasets from the MAVISp database can then be further used in two ways. First, biocurators or data managers can perform a set of analyses, referred to as downstream analyses, which are generated downstream of database creation. These analyses result in the generation of publication-ready figures that summarize the predicted effects for each variant and assist results interpretation.

Furthermore, the biocurators use data from the downstream analysis to create a report in GitBook (https://elelab.gitbook.io/mavisp/), using a standard Markdown template and a semi-automated procedure. Biocurators and data managers also act as reviewers for reports created by their peers. A review status is assigned to each GitBook entry to guide users regarding the quality and integrity of the curated data. To achieve this, we defined four review status levels (i.e., stars) for each protein entry (https://elelab.gitbook.io/mavisp/documentation/mavisp-review-status).

*Structure-based assessment of variants with MAVISp*

**A**



**B**



**Fig. 1. Overview of MAVISp components.** (A) MAVISp includes different modules, each managed by workflow engines or dedicated tools. The modules highlighted in orange handle the selection and collection of protein structures, while the modules in blue and purple are dedicated to structural analyses of variant effects in relation to protein functional- or stability-related properties. Additionally, the framework provided modules with results from VEPs and scores derived by experiments, such as deep mutational scans (green). The procedure begins with a gene name, its UniProt and RefSeq identifiers and the desired structural coverage. For each gene, all the steps can be conducted on a standard server with 32-64 CPUs. The only exceptions are: i) the ENSEMBLE GENERATION module, which includes all-atom MD simulations, and ii) Rosetta-based calculations on binding free energies and folding/unfolding free energy calculations. Depending on the simulation length and system size, these might require access to HPC facilities. On the left, the *simple mode* for the assessment is illustrated, which uses single experimental structures or models from AlphaFold2 or AlphaFold3. On the right, the *ensemble mode* is schematized in which a conformational ensemble for the target protein or its complexes is applied. Hereby, we consider a conformational ensemble a collection of 3D conformations of the protein generated by a sampling method such as molecular dynamics or provided by NMR structures in the PDB (B) Scheme of the current workflow for the MAVISp database and webserver. Biocurators apply specific workflows and protocols within each MAVISp module to generate structure-based predictions of changes linked to variants in each protein target. In doing so, they take advantage of the MAVISp toolkit as well as our mutation sources. The results are gathered into a text-based database. The data are further processed by the MAVISp Python package, which performs consistency checks, aggregate the data and outputs human-readable CSV table files, that make up the MAVISp database. These CSV files are imported by the Streamlit web app, powering the MAVISp webserver (https://services.healthtech.dtu.dk/services/MAVISp-1.0/), where the data are available for interactive visualization and download. In addition, the MAVISp database can be used to generate graphical representations of the data, such us dot plots, lollipop plots, and UpSet plots. Finally, based on the information gathered so far, we provide GitBook reports to facilitate the interpretation of the results: https://elelab.gitbook.io/mavisp/.

Finally, the MAVISp database is presented through a user-friendly Streamlit-based website (https://services.healthtech.dtu.dk/services/MAVISp-1.0/). The web app includes various visualizations to aid the interpretation of MAVISp results that are essentially equivalent to the downstream analyses outlined above: (a) a dot plot displaying classifications for each variant across MAVISp modules, experimental data (if available), and the VEP results, (b) a lollipop plot aggregating relevant mechanistic indicators (i.e., MAVISp-identified effects at the structural level) associated with potentially pathogenic variants, and (c) an interactive representation on the 3D structure, showing the localization of mutation sites identified in (b). These features are designed to support the interpretation of results and facilitate the identification of variants with specific

5

156  mechanisms and multiple effects. The source code for the MAVISp Python package and the web application
157  are available on GitHub (https://github.com/ELELAB/MAVISp), while the complete dataset can be down-
158  loaded from an OSF repository (https://osf.io/ufpzm/ ). The OSF repository also include previous version of
159  the database. Both source code and data are freely available and released under open-source or free li-
160  censes.
161  We invite requests on targets or variants that are not yet available in MAVISp or scheduled for curation. We
162  also welcome contributors as biocurators or developers, pending training and adherence to our guidelines
163  (https://elelab.gitbook.io/mavisp). To facilitate entrance into the MAVISp community of biocurators and de-
164  velopers, we organize training events, research visits and workshops.
165  Notably, a comprehensive update will be conducted annually to incorporate new versions of external tools or
166  resources used by MAVISp, ensuring that resources remain current. Moreover, we continuously expand our
167  toolkit and develop new modules to enable even more comprehensive assessments. The criteria for including
168  new methods and approaches in the framework are detailed in the GitBook documentation (https://elelab.git-
169  book.io/mavisp/documentation/how-to-contribute-as-a-developer).
170
171  
172  *MAVISp modules for structure collection and selection*

173  MAVISp includes various modules to select and model the structures of interest in both *ensemble* and *simple*
174  *mode* (**Fig. 1a**).
175  The STRUCTURE SELECTION module enables biocurators to identify the starting structure for their study,
176  both for models of the free and bound states of the protein of interest. This module includes structure retrieval
177  from the Protein Data Bank (PDB)[39], the AlphaFold Protein Structure Database[25], or through the generation
178  of initial models with AlphaFold3[40], AlphaFold2[2] and, AlphaFold-multimer[24]. In addition, it streamlines the
179  selection of structures in terms of structural quality, experimental resolution, missing residues, amino acidic
180  substitutions with respect to the UniProt reference sequence, as well as the AlphaFold per-residue confi-
181  dence score (pLDDT), integrating tools such as PDBminer[41]. Using AlphaFill[42] further assists in identifying
182  cofactors to be included in the model structure or to identify mutation sites that should be flagged, if located
183  in the proximity of a missing cofactor in the structure model. When necessary, a workflow is available to
184  reconstruct missing residues or design linkers to replace large, disordered loops within structured domains
185  (Methods).
186  According to the protocol established for the generation of the models, we retain 3D structures with reason-
187  able accuracy based on parameters such as pLDDT, Predicted Aligned Error (PAE), and pDOCKQ2[43]. In
188  addition, the module includes protocols based on AlphaFold [24,44] or comparative modeling[45,46] when the com-
189  plex between the protein target and the interactor involves Short Linear Motifs (SLiMs).
190  The INTERACTOME module aids the identification of protein interactors for the target protein and their com-
191  plex structures by querying the Mentha database[47], the PDB, and experimentally validated proteome-wide
192  AlphaFold models [48], as well as the STRING database[49] (Methods). Once a suitable set of interactors has
193  been identified, the information is used to predict protein complex structures, which are then utilized in the
194  subsequent steps (i.e., the LOCAL_INTERACTIONS module, see below).
195  The ENSEMBLE GENERATION module allows the use of structural ensembles from different sources, such
196  as NMR structures deposited in PDB, coarse-grained models for protein flexibility (e.g., CABS-flex[50]) or all-
197  atom Molecular Dynamics (MD) simulations (with GROMACS[51] and PLUMED[52,53]) of the protein structure or
198  its complexes. The choice of the method to be used is based on the required accuracy of the generated
199  ensemble and the available computational resources. Once individual structures or structural ensembles for
200  the protein candidate are selected – either alone or with interactors - the analysis modules can be used.
201
202  
203  *MAVISp modules for structural analysis*

204  MAVISp integrates different analysis modules for both *ensemble* and *simple mode* (**Fig.1a**). The minimal set
205  of data required to import a protein target and its variants into the MAVISp database includes the results from
206  the STABILITY and PTM modules, along with predictions from VEPs. The STABILITY module is devoted to
207  estimating the effects of the variants on the protein structural stability using folding free energy calculations
208  (Methods). This module leverages workflows for high throughput *in silico* mutagenesis scans[54,55] and a newly
209  implemented protocol for RaSP[56] (Methods). All the methods used in this module predict change of free
210  energy of folding upon the insertion of an amino acid substitution, and predictions are performed using FoldX,
211  Rosetta, or RaSP. Once these predictions have been collected, MAVISp applies a consensus approach to
212  classify the effect of the variants (Methods). The defined thresholds for changes in free energy are based on

6

213 evidence that shows that variants with changes in folding free energy below 3 kcal/mol do not exhibit a
214 marked decrease in stability at the cellular level [57,58]. Thus, MAVISp defines the following classes for changes
215 in stability: stabilizing (ΔΔG ≤ - 3 kcal/mol with both methods, FoldX and Rosetta or RaSP), destabilizing
216 (ΔΔG ≥ 3 kcal/mol), neutral ( -2 < ΔΔG < 2 kcal/mol), and uncertain (-3 < ΔΔG ≤ -2 kcal/mol  or 2 ≤ ΔΔG < 3
217 kcal/mol). A variant is also classified as uncertain if the two methods would classify the effect of the variant
218 differently. Since March 2024, we adopted the consensus between RaSP and FoldX as a default for data
219 collection, after performing a benchmark using the MAVISp datasets (**Supplementary Text S1** and
220 https://github.com/ELELAB/MAVISp_RaSP_benchmark). RaSP provides a suitable solution for high-
221 throughput data collection compared to the CPU-intensive scans based on Rosetta. In low-throughput stud-
222 ies, where we focus in detail on a target protein, we can include Rosetta data, which are computationally
223 more demanding.

224 The LOCAL INTERACTION module can be applied if the STRUCTURE SELECTION and INTERACTOME
225 modules identify at least a suitable structure of the complex between the target protein and another biomol-
226 ecule. The LOCAL INTERACTION module is based on estimating of changes in binding free energy for
227 variants at protein sites within 10 Å of the interaction interface, using protocols and consensus strategies that
228 mirror those for STABILITY. In this case, we use a combination of FoldX and Rosetta calculations (Methods).
229 Binding free energy thresholds are set based on the expected error margins of the predictors, approximately
230 ±1 kcal/mol, as outlined by the authors of the methods and in accordance with general good practice in the
231 literature. This approach addresses the scarcity of experimental datasets on amino acid substitutions that
232 impacting protein-protein interactions[59–61], which are often constrained by system heterogeneity, limited mu-
233 tation numbers, or both, thereby complicating reliable benchmarking. We rely on a consensus approach be-
234 tween the results of FoldX and Rosetta on changes in binding free energies upon amino acid substitution.
235 We classify a variant as stabilizing (both methods predict ΔΔG <= -1 kcal/mol), neutral (-1 kcal/mol < ΔΔG <
236 1 kcal/mol) or destabilizing (ΔΔG >=1 kcal/mol). Cases in which the two methods disagree on the classifica-
237 tion, or for which we do not have a prediction for both methods, and the side chain relative solvent accessible
238 area of the residue is >= 25%, are classified as uncertain. This is because, in high-throughput data collection,
239 we cannot exclude the possibility that the site interacts if it is solvent exposed, as often in structural biology,
240 only part of the 3D structures of protein-protein complexes are available or can be modelled. We also included
241 support for LOCAL INTERACTION for protein and DNA interactions, as well as for homodimers. Notably, a
242 strength of our approach is to provide annotations for the effects of protein variants on various biological
243 interfaces for the same target protein.

244 In the *ensemble mode*, the STABILITY and LOCAL INTERACTION modules are used on ensembles of at
245 least 20-25 structures from the simulations or on the three main representative structures upon clustering,
246 depending on the free energy calculation scheme to apply. The results obtained for each structure are then
247 averaged, and classification is performed with the same strategies we use in *simple mode* using these aver-
248 age values. This approach is used to mitigate limitations due to lack of backbone flexibility when these free
249 energy methods are applied to just one single 3D structure[38,54,62,63].

250 The LONG-RANGE module applies coarse-grained models to estimate allosteric free energy changes upon
251 amino acid substitution based on AlloSigMA2[64]. The protocol followed by the LONG-RANGE module has
252 recently been updated and benchmarked using experimental data from deep mutational scans[65].  Details on
253 the parameters and steps for analysis are also provided in the Methods. Variants are annotated as destabi-
254 lizing (positive changes in allosteric free energy), stabilizing (negative changes in allosteric free energy),
255 mixed effects (both conditions occur), or neutral if the variant does not cause any significant change. Addi-
256 tionally, variants that do not cause a significant change in residue side-chain volume are annotated as un-
257 certain. In the *ensemble mode*, we applied graph theory metrics based on changes in the shortest commu-
258 nication paths using atomic contact-based Protein Structure Network[66] . This analysis, combined with the
259 AlloSigMA2 data, allows pinpointing variants with long-range effects to functional sites or protein pockets that
260 could serve as interfaces to recruit interactors or ligands.

261 The FUNCTIONAL SITES module in *simple mode* allow to evaluate the effect of variants at (or in the proximity
262 of) the active site of enzymes or cofactor binding sites of proteins and it is based on analyses of contacts with
263 the second sphere of coordination of the residues belonging to these sites (see Methods).

264 The FUNCTIONAL DYNAMICS module in *ensemble mode* includes enhanced sampling simulations to fur-
265 ther assess the local or long-range effects of a variant. As a first example, we applied this class of methods
266 to validate the long-range effects predicted for p53 variants on the DNA-binding loops[38], and included such
267 results in the MAVISp database.

268 The PTM module currently supports phosphorylation only, annotating the effect of variants at phosphorylata-
269 ble sites. It evaluates how the loss or changes of phosphorylation sites may impact protein regulation, stabil-
270 ity, or interaction with partners. To this goal, the module collects analyses and annotations such as solvent
271 accessibility of the mutation site, inclusion of the site in phosphorylatable linear motif, comparison between
272 predicted changes in folding or binding free energy upon amino acid substitution or upon phosphorylation at
273 the site of interest. In the module, we applied a custom decision logic (**Supplementary Text S2**) to derive
274 the classification for each variant as neutral, damaging, unknown effect, potentially damaging or uncertain.
275 The identification of the phosphorylation sites in the PTM module is based on known experimental phospho-
276 sites and SLiMs, as retrieved by Cancermuts[67]. These data are complemented by a manually curated selec-
277 tion of phospho-modulated SLiMs (https://github.com/ELELAB/MAVISp/blob/main/mavisp/data/phospho-
278 SLiMs_09062023.csv). For solvent-inaccessible phosphorylatable residues, the effects are classified as un-
279 certain in the *simple mode*. In these cases, the *ensemble mode* is required to investigate wheatear a cryptic
280 phosphorylated site may become accessible upon conformational changes[68,69]. Of note, the current version
281 of the PTM module has been designed based on fundamental principles on how phosphorylation can affect
282 the protein structure and should be used to identify variants for further investigation, particularly for experi-
283 mental research. Benchmarking the effectiveness of this module would be difficult at present time, given the
284 relatively small number of amino acid substitutions that can affect phosphorylation currently present in the
285 MAVISp database, especially considering those for which experimental data is available. To this purpose,
286 we are currently in the process of curating and including more proteins relevant to benchmarking the PTM
287 module. These will include experimental data on protein stability and protein-protein interactions upon phos-
288 phorylation [70] [71].
289 MAVISp includes further analyses and annotations, such as predictions on regions involved in early folding
290 events[72], pLDDT score, secondary structure, and side-chain solvent accessibility, which can assist in the
291 interpretation of the results.
292
293

294 *Variant Effect Predictors included in MAVISp*
295

296 MAVISp provides annotations for the variant interpretation reported in ClinVar[9], or calculated with REVEL[73],
297 DeMaSk[74], GEMME[14], EVE (Evolutionary model of variant effect)[75], and AlphaMissense[26]. In MAVISp, each
298 of them is handled by a separate module. The results of these VEPs can be combined with the results from
299 the MAVISp structure-based modules to understand variant effects and to prioritize variants for other studies,
300 as detailed in the examples below.
301

302 *Sources of variants supported by MAVISp*
303

304 By default, we apply in silico saturation mutagenesis, which means that we provide predicted effects for each
305 variant of a target protein at any position that has a structural coverage. Additionally, all variants reported for
306 the target protein in COSMIC, cBioPortal, and ClinVar are annotated within MAVISp. We routinely update
307 and maintain the entries in the MAVISp database to include up-to-date annotations using Cancermuts[67]. All
308 Cancermuts annotations for MAVISp and other protein targets are also available at the Cancermuts web-
309 server, https://services.healthtech.dtu.dk/services/Cancermuts-1.0/. In addition, annotations from lists of var-
310 iants from other studies, such as data on cohort-based or nationwide studies or other disease-related ge-
311 nomic initiatives, can be manually introduced.
312 Currently, MAVISp includes data on eight+ million variants from 1096 proteins (at the date of 20/11/2025).
313 An overview of the currently available data and how to use them to address different research questions is
314 described in detail in the next sections. The first targeted studies in which MAVISp has been applied to
315 understand variants impact in rare genetic diseases[76] or involved in cancer hallmarks[77,78] are also suitable
316 examples
317
318

319 *Interpretation of the results of MAVISp*
320

321 MAVISp provides a comprehensive set of results for many variants; therefore, we have devised a few strat-
322 egies that can be useful to make sense of the MAVISp data for a few common use cases that users might
323 encounter.
324

325 One of the most important outputs from the downstream analyses, of MAVISp is the so-called dot plot, which
326 is available on the GitBook reports or released within the target studies of specific proteins (see below for
327 examples). A dotplot can also be generated within the MAVISp webserver in the "Classification" tab, for up
328 to 50 variants of choice simultaneously. This plot showcases i) the classification of the different VEPs inte-
329 grated in MAVISp, ii) the classification performed by each MAVISp module, iii) the classification of variants
330 in ClinVar, when available, as variant label colors. The code to generate dot plots from MAVISp csv file is
331 also available in GitHub (https://github.com/ELELAB/mavisp_accessory_tools/tree/main/tools). The MAVISp
332 modules classification has a different meaning depending on the considered module: a variant classified as
333 damaging for a VEP usually means it is predicted as functionally damaging or pathogenic (depending on the
334 predictor), while a variant classified as damaging for stability just means that the variant is predicted to com-
335 promise the structural stability of the protein, and one classified as damaging by the long range module is
336 predicted to have significant long-range effects, and so on. Another representation which depends on further
337 processing of a text output created by dot_plot.py (i.e., alphamissense_out.csv) provides a concise repre-
338 sentation of the classes of mechanistic indicators found for each variant in the form of lolliplots. Lolliplots are
339 also available in the GitBook report or in the "Damaging mutation" tab on the website, that shows only those
340 variants that are at the same time: i) classified as pathogenic for AlphaMissense, ii) classified as loss-of-
341 fitness or gain-of-fitness by DeMaSk and iii) damaging for the respective structure-based module of MAVISp.
342 The downstream analysis toolkit also provides the code to prepare upset plots or venn diagrams for the
343 variant source (as reported in Gitbook).
344 Consulting the available dot plot for an entry of interest is therefore the most straightforward place to start to
345 access MAVISp data. To identify a subset of variants of interest, we have defined the following strategy for
346 a data-driven discovery of variants of interest with little other information (i.e. VUS, conflicting evidence or
347 variants not reported in ClinVar). In this case, the dot plot allows to understand first which variants are pre-
348 dicted to be pathogenic, by using the AlphaMissense classification; these are the ones reported as Damaging
349 in the AlphaMissense row. For these, we also consider the output of DeMaSk, that define whether the variant
350 is classified as gain-of-fitness or loss-of-fitness. If a variant fullfil these criteria, we then consider the structure-
351 based MAVISp predictions for mechanistic indicators, that give us one or more explanations of the reason
352 for the effect of the variant. For instance, the variant could be destabilizing the protein structure and will be
353 reported with an altered stability as mechanistic indicator. Another common use case is to use MAVISp to
354 get a mechanistic interpretation of variants already known in ClinVar. In this case, if the variant already has
355 an interpretation of Pathogenic, Likely pathogenic, Benign, or Likely benign, we can just refer to the MAVISp
356 mechanistic interpretation.
357 Importantly, researchers should always refer to specific biological or phenotypical contexts when interpreting
358 predictions from MAVISp, including their knowledge of the biological role the protein investigation has or
359 concerning the nature of the disease of interest. For instance, predictions might lead to different conclusions
360 if the protein under consideration is from a tumor suppressor or from an oncogene.
361 In the next section we illustrate some of the applications of data collected with MAVISp through case stud-
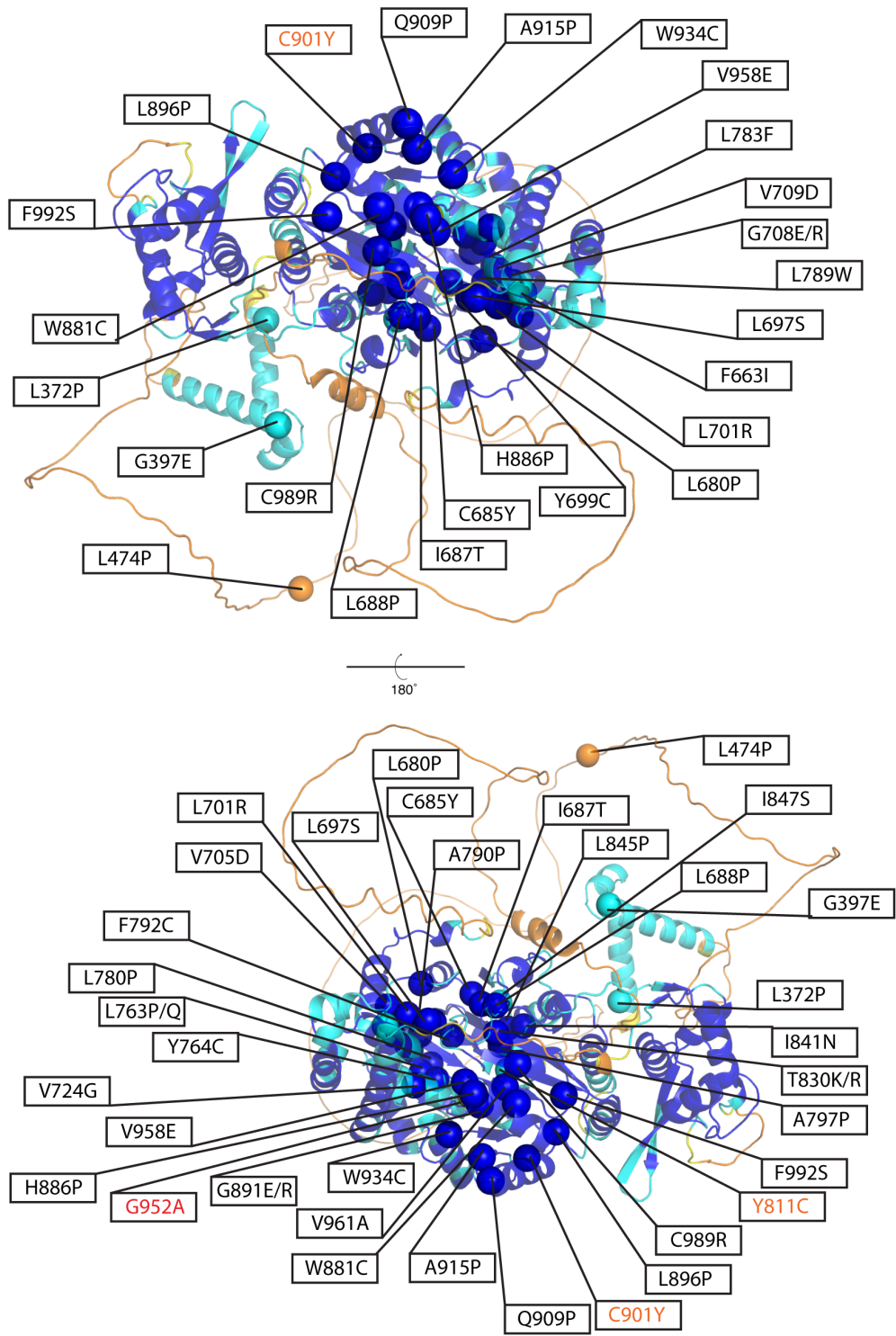362 ies (**Table S2** for mapping of case studies and modules).
363
364 *COSMIC Tumor Suppressor Genes and Oncogenes*
365
366 At first, we prioritized MAVISp data collection of known driver genes in cancer, i.e., tumor suppressors and
367 oncogenes. To this goal, we collected data for the COSMIC Tumor Suppressor Genes (COSMIC v96), while
368 the collection of the COSMIC Oncogene and Dual Role targets is ongoing. Furthermore, we have been in-
369 cluding genes reported as a candidate driver by the Network of Cancer Genes (NGC)[79].
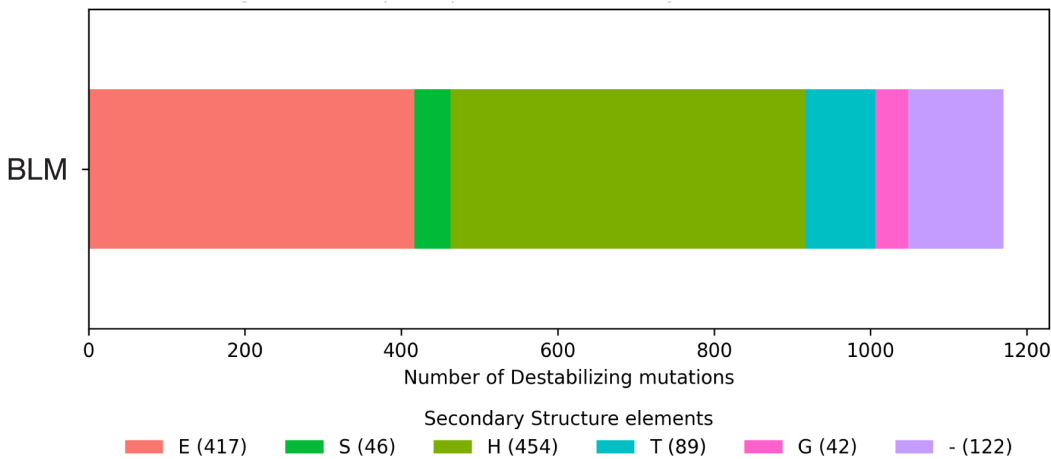370 The MAVISp datasets on cancer driver genes can assist the identification of molecular mechanisms of pre-
371 dicted or known pathogenic variants in these genes, as well as to aid the characterization of Variants of
372 Uncertain Significance (VUS). A recent example is the study we performed on BRCA2[78].In this study, we
373 analyzed BRCA2 variants reported in ClinVar, comparing the predictions from the STABILITY and LOCAL
374 INTERACTIONS modules of MAVISp with results from a multiplex assay which measured the impact of these
375 variants on cell viability. We were able to explain the effect of 84 BRCA2 variants, which were classified as
376 non-functional by the assay, and for which MAVISp predicted effects on protein stability or binding to the
377 binding partner SEM1.
378

9

*Structure-based assessment of variants with MAVISp*

a



b

379

**Fig. 2. Variants with effects on structural stability in the tumor suppressor protein BLM.** (a) The cartoon representation shows the trimmed mode BLM$_{368-1290}$ and the spheres highlight the C$_\alpha$ atom of the 41 positions harboring 45 variants predicted as destabilizing by the MAVISp STABILITY module (RaSP/FoldX consensus) and annotated in ClinVar. Among these, Y764C, G891E, and L896P are also reported in CBioPortal, whereas F663I, L845P and C901Y are also reported in COSMIC. The two views correspond to the same domain rotated by 180°. The backbone and spheres are colored according to the AlphaFold pLDDT scores, i.e., blue - very high (pLDDT > 90), cyan - confident (70 < pLDDT <= 90), yellow - low (50 < pLDDT <= 70), and orange - very low (pLDDT <= 50). The labels indicate the mutation sites and the corresponding variants and are colored by ClinVar classification., uncertain significance (black), conflicting interpretation of pathogenicity (orange), and likely pathogenic (red). (b) The stacked bar plot shows the distribution of destabilizing BLM variants across secondary structure elements as defined by DSSP ( (i.e., H = α-helix, B = residue in isolated β-bridge, E = extended strand, participates in β ladder, G = 3-helix (3$_{10}$ helix), I = 5-helix (π-helix), T = hydrogen bonded turn, S = bend, and "-" = no secondary structure identified). The results refer to the data available in the MAVISp database on 12th September 2025. More information about BLM analyses with MAVISp can be found in the corresponding GitBook report: https://elelab.gitbook.io/mavisp/proteins/blm

In the case of tumor suppressors, the identification of variants that might lead to loss of function is particularly important. Given structure-function relationship in proteins, structural stability represents a key determinant that can be disrupted by amino acid substitutions, potentially resulting in local or more drastic misfolding and loss of function[80] As an example of loss of function due to changes in stability, we report the analysis of the MAVISp entry for the tumor suppressor BLM, a DNA helicase involved in DNA replication, recombination and repair[81]. We identified a total of 1170 predicted destabilizing variants according to the STABILITY module, of which 45 annotated in ClinVar (**Fig. 2a**). Among these, 82% destabilizing variants was found in structured regions of the protein, while the remaining 18% are located in disordered residue stretches (**Fig. 2b**). Of the ClinVar-reported variants, 42 are classified as VUS. Y811C and C901Y are reported with conflicting interpretations and only G952A is reported as likely pathogenic.

These results provide a starting point for variant characterization and prioritization. As suggested in the previous section, our results can be used to guide the selection of a subset of variants that have a predicted pathogenic impact from AlphaMissense, with a loss-of-fitness signature according to DeMaSk and that we predict damaging for stability., These would be suitable candidates for experimental validation. Concerning BLM, MAVISp identifies 41 ClinVar VUS or variants with conflicting evidence that could be prioritized according to these criteria (**Table S3**).

For example, depending on the size of the library to validate, methods such as flow cytometry sorting or cycloheximide chase assays[82,83] or use approaches based on multiplex technologies[84–87] would be useful to validate our predictions

## *Integration of MAVISp data with experimental data*

A useful feature of MAVISp is a dedicated module to curate and import experimentally derived scores on the effects of the variants on different biological readouts (i.e., the EXPERIMENTAL DATA module, **Fig. 1a**). These data can be directly compared with the structural properties we predict with MAVISp, for a variety of purposes. For example, they can serve as additional layer of information respect to the structure-based mechanistic indicators themselves. Additionally, as done in the aforementioned BRCA2 study, they can be used as a source of information for variants with a known detrimental effect that can depend on different mechanisms of action for each variant, which can be investigate using MAVISp. In cases such as this, MAVISp helps identifying the possible mechanism for which variants have an effect, for further in-depth investigation.

Experimental data can also be used to validate the results of certain MAVISp modules, for cases in which the predicted structural properties are related to the experimentally tested biological readouts. Deep mutational scans can also be used to benchmark or tune the thresholds used for classification performed by the MAVISp modules, including structural properties. In this context, the format of MAVISp database files is handy for further data processing, for example using biostatistical models or machine learning. In the case of PTEN, we included data from available deep mutational scans, reporting on the effect of mutations on cellular abundance or phosphatase activity[84,88,89], in its MAVISp entry. Cellular abundance represents a critical property that is often perturbed by missense mutations, and that can be altered by changes in protein structural stability. We therefore compared predictions from the MAVISp STABILITY module—based on a consensus of RaSP and FoldX—with protein abundance scores obtained from VAMP-seq assays [84,89,90]. To compare the classification obtained by the stability module with the experimental data, we considered how the abundance score from the experiment have been classified. Multiple classification strategies have been used for these data: on one side, the ProteinGym benchmark dataset[91] applies a threshold based on the median of the abundance score (i.e., 0.77). Variants that scored lower than this threshold (>=22% reduction of abundance relative to the wild-type) were classified as low abundant, whereas those that scored higher

438 were considered to be similar to the wild-type. The second classification followed the original PTEN study
439 deposited in MaveDB (MaveDB ID urn:mavedb:00000013-a-1), which defines four abundance classes. In
440 this scheme, the 5% lowest-abundance synonymous variants corresponded to a score of 0.71[84], and variants
441 were classified into bin 1 (low-abundant, both score and confidence interval < 0.71), bin 4 (WT-like abundant,
442 both scores > 0.71), bin 2 (likely low-abundant, score < 0.71 but confidence interval > 0.71), and bin 3 (likely
443 WT-like abundant, score > 0.71 but confidence interval < 0.71). For this analysis, we retained only variants
444 in bins 1 and 4, to ensure an unambiguous classification. After applying these filters and excluding uncertain
445 variants defined by the STABILITY module, the MaveDB-based classification contained 1690 variants. To
446 enable a direct comparison between the two classification strategies, the ProteinGym-based dataset, which
447 initially comprised 3211 variants, was filtered to include the same 1690 variants as the filtered MaveDB da-
448 taset. The two classification schemes were found to be largely concordant, differing only for variants with
449 abundance scores between 0.71 and 0.77, which were considered damaging by ProteinGym and neutral by
450 MaveDB.

451 **Fig.3 and Table 1** illustrate the performance of the MAVISp STABILITY classification against the classifica-
452 tion of experimental data on protein abundance for PTEN. In this first comparison, we applied the same
453 threshold suggested for this dataset from the benchmarking dataset ProteinGym[91], which is based on the
454 median value of the DMS scores. The consensus approach provided by the STABILITY module of MAVISp
455 (accuracy 0.814) has an overall better performance in identifying variants that are found to be damaging in
456 the assay than those predicted to cause damaging effects according to GEMME or DeMaSk (**Fig. 3b**). Nev-
457 ertheless, this approach has a lower sensitivity (0.66) compared to GEMME. We thus wondered if the rela-
458 tively low sensitivity we obtained was due to cases with experimental scores too close to the median (**Fig.**
459 **3c**). Additionally, in the original study for PTEN and as deposited in the MaveDB[92,93], a different classification
460 for the variant scoring based on four abundance levels was proposed, as detailed above. We thus performed
461 a comparison of the MAVISp results with the experimental dataset for the PTEN experiment from MaveDB
462 using the abundance level classes as a threshold (**Fig. 3c)**, resulting in increases sensitivity for the methods
463 applied within MAVISp. The results on PTEN from MAVISp fits nicely with recent computational studies of
464 PTEN variants using Rosetta calculations of protein stability and analyses of sequence conservation[16,17].

466 **Table 1.** Performances of MAVISp modules and VEP predictors against experimental measurements of protein abundance and
467 phosphatase activity for PTEN.

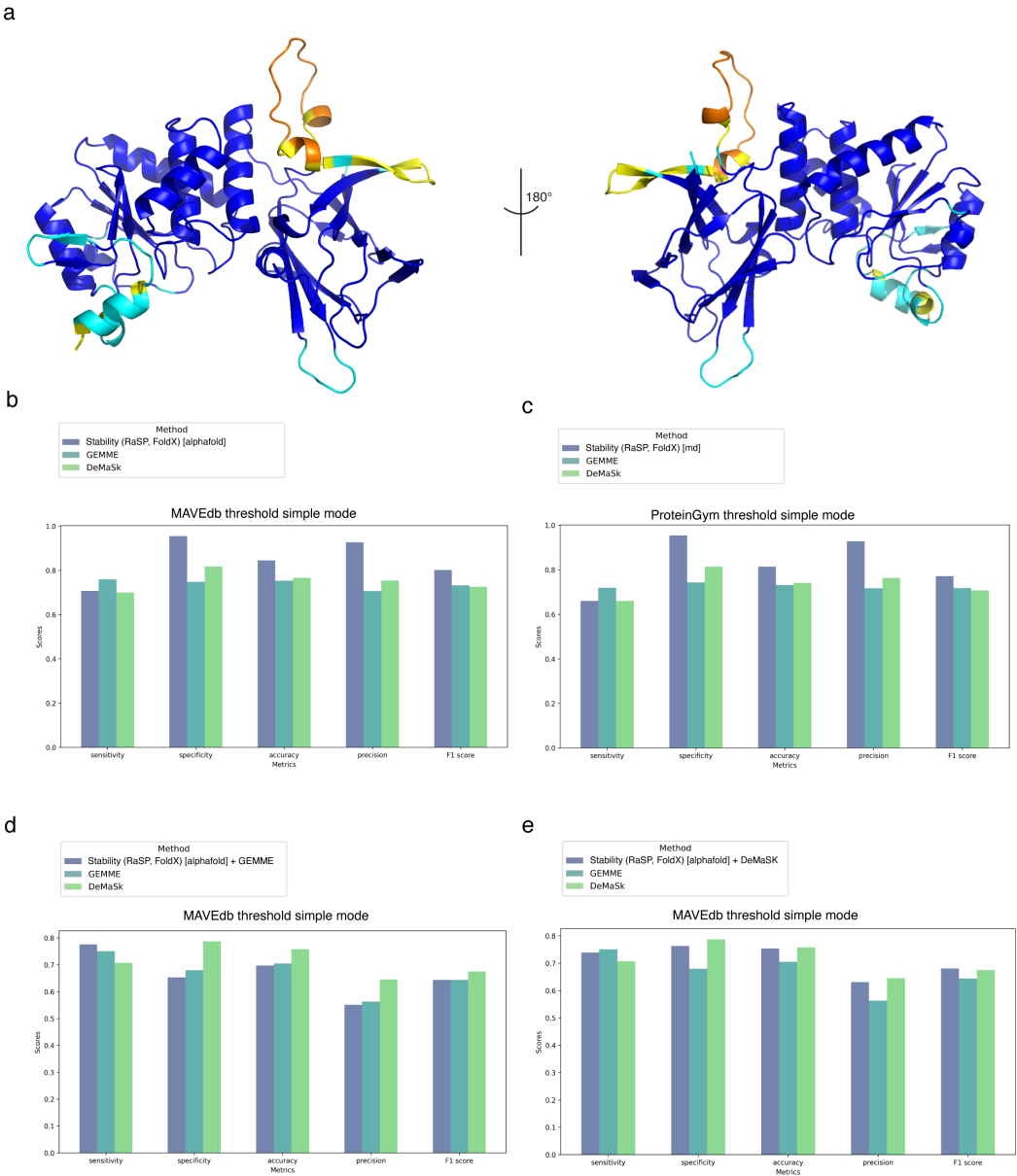| Assay | column comparison | threshold_mode | sensitivity | specificity | accuracy | precision | F1 score |
|---|---|---|---|---|---|---|---|
| protein abundance assay | RaSP/FoldX consensus | ProteinGym | 0,666 | 0,954 | 0,814 | 0,928 | 0,771 |
| | GEMME | | 0,719 | 0,743 | 0,731 | 0,717 | 0,718 |
| | DeMaSk | | 0,66 | 0,814 | 0,741 | 0,763 | 0,708 |
| | RaSP/FoldX consensus | MaveDB | 0,707 | 0,955 | 0,845 | 0,927 | 0,802 |
| | GEMME | | 0,759 | 0,748 | 0,753 | 0,706 | 0,732 |
| | DeMaSk | | 0,7 | 0,818 | 0,766 | 0,754 | 0,726 |
| phosphatase assay | RaSP/FoldX + GEMME | MaveDB | 0,776 | 0,653 | 0,697 | 0,551 | 0,644 |
| | RaSP/FoldX + DeMaSk | | 0,739 | 0,763 | 0,754 | 0,631 | 0,681 |
| | GEMME | | 0,751 | 0,68 | 0,705 | 0,563 | 0,644 |
| | DeMaSk | | 0,707 | 0,787 | 0,758 | 0,645 | 0,675 |

472 We next assessed whether MAVISp could also inform predictions of variant effects on PTEN phosphatase
473 activity. Experimental data from a cellular phosphatase assay (MaveDB ID urn:mavedb:00000054-a-1) were
474 classified as reduced (< 0.89), wildtype-like (0.89–1), or hyperactive (> 1). Here, we investigated whether
475 integrating MAVISp STABILITY data with VEP results could enhance predictive power, since reduced stability
476 is not the only possible mechanism for loss of phosphatase activity in mutated variants. We combined the
477 STABILITY results with GEMME or DeMaSk, applying a priority logic in which damaging calls from
478 GEMME/DeMaSk were given priority. This strategy produced performance comparable to GEMME or
479 DeMaSk alone (**Fig.3d-e, Table 1**). Notably, combining changes in folding free energies with GEMME

480 increased sensitivity (0.78) but reduced specificity (0.653), yielding an F1 score comparable to that of
481 GEMME but lower than the one of DeMaSk (0.64; **Fig.3d-e, Table 1**).

482
483 Overall, with the examples in this section, we illustrate examples on how to use MAVISp data to compare
484 predictions and experiments, as well as how to integrate MAVISp modules on structural properties with VEP
485 results.



**Fig. 3. Comparison of GEMME, DeMaSk, and MAVISp STABILITY module predictions with experimentally-derived scores for protein abundance and phosphatase activity of PTEN.** (a) The trimmed AlphaFold structure (residues 1-351) of PTEN used for MAVISp stability module calculations is shown as a cartoon, colored according to pLDDT scores. (b-e) Histograms with performances of MAVISp STABILITY module, DeMaSk, and GEMME in predicting the effect of variants using VAMP-seq scores with ProteinGym (b) and MaveDB (c) thresholds. (d-e) illustrates the performances of the same tools or their combination against an experimental functional readout that assess the phosphatase activity at the cellular level.

## *Proteins involved in cancer hallmarks*

498 To expand the contents of the MAVISp database, we have also been focusing on protein targets related to
499 cancer hallmarks[94], and in particular on proteins involved in cancer hallmarks related to protein clearance at
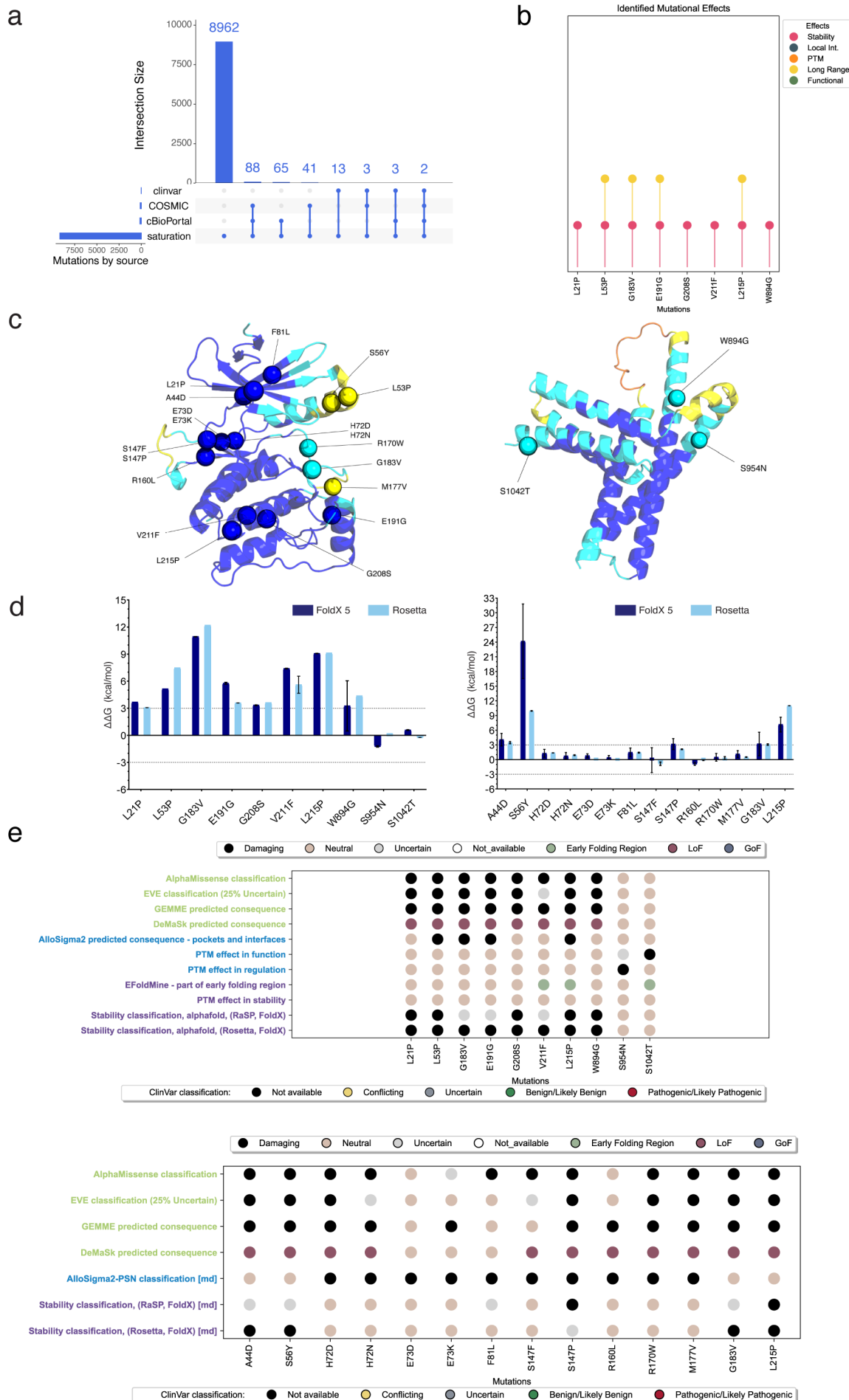
500   the cellular level, i.e. the ability to escape cell death through apoptosis and autophagy, as well as kinases or
501   transcription factors involved in the regulation of cellular proliferation. Mitochondrial apoptosis is tightly regu-
502   lated by a network of protein-protein interactions between pro-survival and pro-apoptotic proteins. We inves-
503   tigated the mutational landscape in cancer of this group of proteins in a previous study[95], which includes
504   structural analyses with different *simple mode* MAVISp modules for both the pro-survival proteins BCL2,
505   BCL2L1, BCL2L2, BCL2L10, MCL1 and BCL2A1, as well as the pro-apoptotic members of the family BOK,
506   BAX and BAK1. In these analyses, the C-terminal transmembrane helix has been removed since the current
507   version of our approach does not support transmembrane proteins or domains, illustrating an example on
508   how the STRUCTURE SELECTION module works.

509   Autophagy is a clearance mechanism with a dual role in cancer. The autophagy pathway relies on approxi-
510   mately 40 proteins, constituting the core autophagy machinery[96]. As an example of the application of MAVISp
511   to this group of proteins, we applied the *simple mode* to the markers of autophagosome formation MAP1LC3B
512   and the central kinase ULK1, building on the knowledge provided by previous work[36,37].

513   In the case of ULK1, we expanded our analysis to cover a larger part of the structure of the protein, meaning
514   that both the N-terminal (residues 7-279) and C-terminal domains (837-1046) have been used for the MAVISp
515   assessment. ULK1 also serves as an example of how to customize the trimming of an AlphaFold model to
516   exclude disordered regions or linkers with residues featuring low pLDDT scores, in *simple mode*. In fact, their
517   inclusion could lead to predictions of questionable quality. Disordered regions cannot be properly represented
518   by a single conformation, and the *ensemble mode* would be necessary to derive more reliable conclusions.

519   ULK1 featured 215 variants reported in COSMIC, cBioPortal and/or ClinVar, as shown in its dot plot (**Fig.**
520   **4A**), which was generated using the downstream analysis tools of MAVISp. Using the *simple mode*, 59 vari-
521   ants had predicted long-range mixed effects. Furthermore, eight had a damaging effect on stability, one had
522   a damaging PTM effect on regulation (S954N), one had a possible damaging PTM effect in function
523   (S1042T), and four variants (L53P, G183V, E191G, and L215P) are characterized by both effects on stability
524   and long-range communication (**Fig. 4B-D**). Most of the variants that were predicted to have long-range and
525   structure-destabilizing effects are in the N-terminal kinase domain of the protein, suggesting that mutations
526   in this domain could result in the inactivation of ULK1 by compromising its 3D architecture. We then per-
527   formed a one-microsecond MD simulation of the ULK1 N-terminal kinase domain (residues 3-279, PDB ID:
528   5CI7) to generate a structural ensemble for the MAVISp *ensemble mode*. In this case, we used an approach
529   based on graph analysis from a contact-based PSN (Methods), as provided by the LONG-RANGE module,
530   which verified if long-range communication occurs between mutation and response sites predicted by the
531   coarse grain model used in the *simple mode*. The *ensemble mode* also validates the prediction on the effect
532   of variants on stability that were done in *simple mode*, as it compensates for the none or limited mobility of
533   the protein main chain that characterize the used in the STABILITY module. Overall, the application of the
534   *ensemble* mode allowed to validate five variants with predicted long-range damaging effects (H72D, H72N,
535   E73D, E73K, and R160L) and two variants with a damaging effect on stability (G183V and L215P). The
536   predicted destabilizing (**Fig. 4E**) variant L215P has been also identified in samples from The Cancer Genome
537   Atlas (TCGA) [36].

538
539

*Structure-based assessment of variants with MAVISp*



**Fig. 4. MAVISp *ensemble mode* to identify damaging variants in the autophagy kinase ULK1.** a) We examined the central autophagy kinase ULK1 using MAVISp, generating a saturation of all possible variants within the N-terminal (residues 7-279) and C-terminal domains (residues 837-
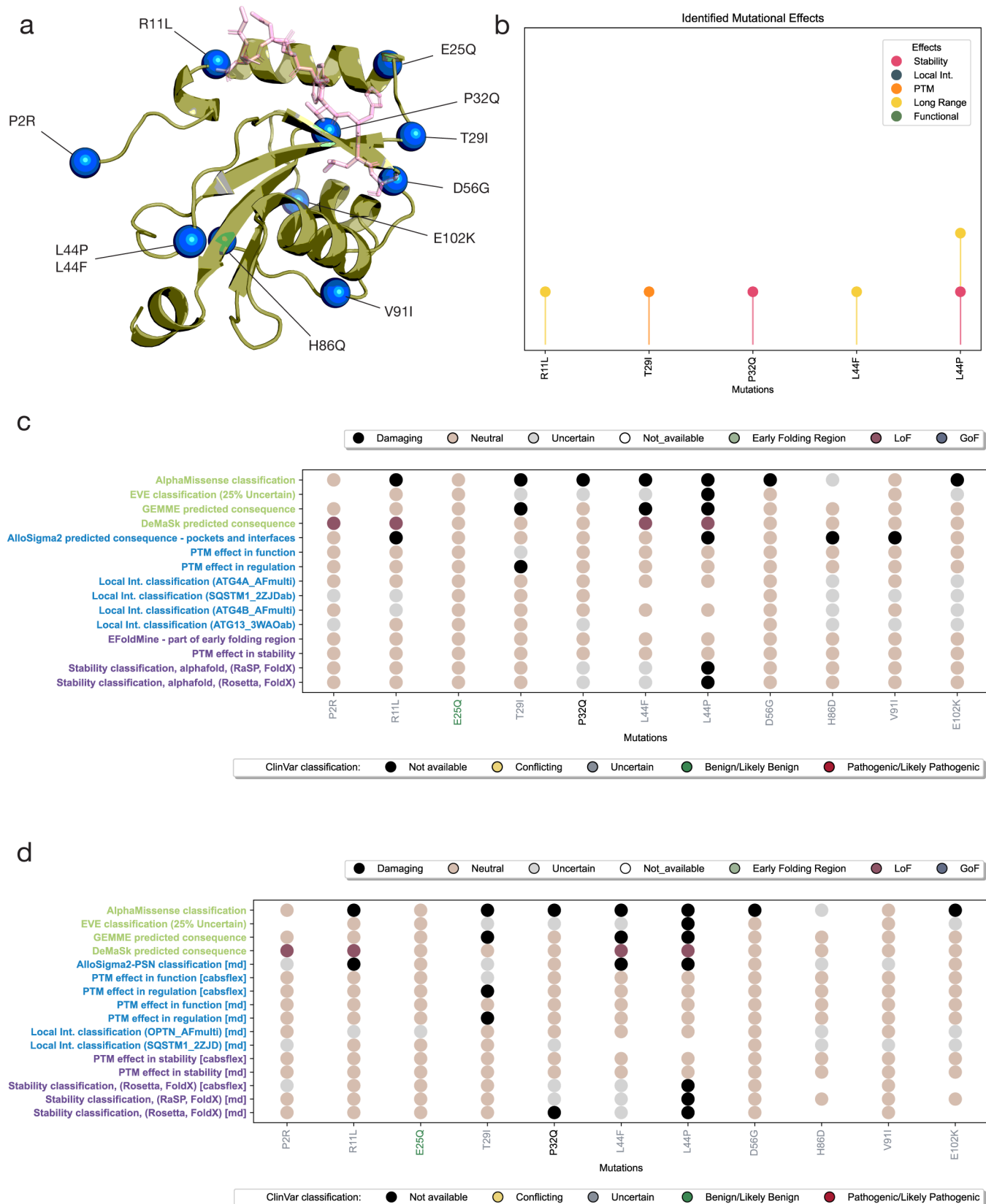
15

544    1046), leading to a total of 8,962 variants. Of these, 215 variants have been identified in COSMIC, cBioPortal, and/or ClinVar databases. b) Among
545    the ones reported in the previous databases, eight variants were reported as pathogenic by AlphaMissense (L21P, L53P, G183V, E191G, G208S,
546    V211F, L215P, and W894G) and among these, four variants are predicted to have a damaging effect on both protein stability and long-range com-
547    munication (L53P, G183V, E191G, and L215P). c) Using MAVISp *simple* and *ensemble* modes, we identified 22 variants with destabilizing effects in
548    terms of folding free energy, long-range effects, or PTM effects in regulation or in function. The mutation sites are highlighted with spheres on the
549    AlphaFold models of the ULK1 N-terminal (left) and C-terminal (right) domains. d) We showed the predicted changes in folding free energy upon
550    amino acid substitution for each of the 22 variants as calculated by the STABILITY module of MAVISp with MutateX and RosettaDDGPrediction with
551    the *simple mode* (left) or with the *ensemble mode* (right). Interestingly, most of the variants that alter structural stability are located in the catalytic
552    domain of the enzyme. This suggests potential mechanisms for ULK1 inactivation. e) Summary of the predicted effects on the 22 variants of ULK1
553    that have been found damaging with at least one MAVISp module with the *simple mode* (upper) or with the *ensemble mode* (lower) using the dot plot
554    representation provided by the MAVISp toolkit for downstream analyses. Of note, the lower legend refers to the color of variants on the X-axis which
555    are related to the ClinVar effect category.

558    The MAVISp entry of the autophagy marker MAP1LC3B provides an example on how the data for the LOCAL
559    INTERACTION module can be obtained in a case of a protein that interacts with a functional motif embedded
560    in intrinsically disordered proteins, i.e., a short linear motif (SLiM). MAP1LC3B in fact is able to bind to pro-
561    teins harboring a so called LC3-interacting region (LIR)[97]. In MAVISp, we report the results for the effect on
562    binding affinity of variants in MAP1LC3B or in its binding partners using three examples of this mode of
563    interaction modeling the binding of MAP1LC3B with the LIR regions of its binding partner SQSTM1 (**Fig.5a**),
564    ATG13, and Optineurin. In this case, we first applied the protocols for (phospho)-SLiM identification devel-
565    oped within the MAVISp framework (Methods) and PDBminer to identify possible starting structures. In the
566    case of optineurin, we further model the flanking regions[77]. We identified ten variants annotated in ClinVar:
567    nine reported as VUS (E102K, H86D, T29I, V91I, P2R, L44P, L44F, D56G, and R11L) and one as benign,
568    i.e., E25Q (**Fig.5a**). MAVISp managed to predict a putative mechanistic explanation for the effect of four
569    variants (**Fig.5b-d**): T29I is predicted to disrupt regulation by phosphorylation, L44P has an effect on both
570    structural stability and long-range effects to distal sites, L44F and R11L have long-range effects (**Fig 5b**).
571    Additionally, a variant found in cancer studies, P32Q, is predicted to have a detrimental effect on structural
572    stability, confirming previous experimental results which showed propensity for aggregation[37]. Of note, this
573    variant is identified with an uncertain prediction for the effect on stability in MAVISp simple mode, whereas
574    two different approaches for generating a conformational ensembles accounting for protein dynamics pre-
575    dicted a destabilizing effect (**Fig 5d**). Additionally, all the variants with a mechanistic indicator from MAVISp
576    are also predicted as pathogenic by AlphaMissense (**Fig.5c-d**) and are good candidates to further experi-
577    mental studies for their effects on the autophagy flux or other functional readouts. V91I is likely to be benign
578    variants since all the predictors identified neutral effects (**Fig. 5c-d**).

*Structure-based assessment of variants with MAVISp*



**Fig. 5**. **Analysis of MAP1LC3B VUS Variants from ClinVar.** a) A structural model (PDB ID: 2ZJD) of the MAP1LC3B (green) interaction with the LIR motif of SQSTM1(pink) highlights ten ClinVar-reported variants (E102K, H86D, T29I, V91I, P2R, L44P, L44F, D56G, R11L and E25Q) along with the cancer-related variant P32Q. These variants are depicted as blue spheres on the structure. (b) Among these variants, five (R11L, T29I, P32Q, L44F and L44P) are predicted as damaging by AlphaMissense. Interestingly, L44P shows a predicted damaging effect on both long-range communication and stability.

(c-d) Summary of the predicted effects on the 11 variants of MAP1LC3B as reported by MAVISp dot plot with the simple mode (c) or with the ensemble mode (d) using the dot plot representation provided by the MAVISp toolkit for downstream analyses. Of note, the lower legend refers to the color of variants on the X-axis which are related to the ClinVar effect category.
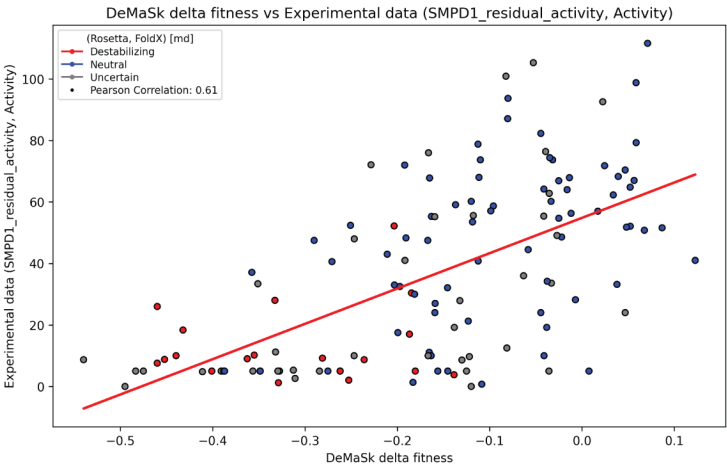
17

## *Application of MAVISp to transmembrane proteins and to variants associated to other diseases*

The STABILITY and LOCAL INTERACTION modules do not support predictions for variants in transmembrane regions. A survey on methods to predict folding free energy changes induced by amin on transmembrane proteins suggested that existing protocols, based on FoldX or Rosetta, are suitable for soluble proteins[98]. Therefore, the protocols implemented in the MAVISp modules for transmembrane proteins only retain those variants that are not in contact with the membrane. An example of a MAVISp entry for this class of proteins is PILRA, which has a low pLDDT score in the transmembrane region, and has been therefore excluded from the model, focusing on the analyses on the variants in the 32-153 region. In addition, we included other transmembrane proteins in the database such as ATG9A and EGFR.
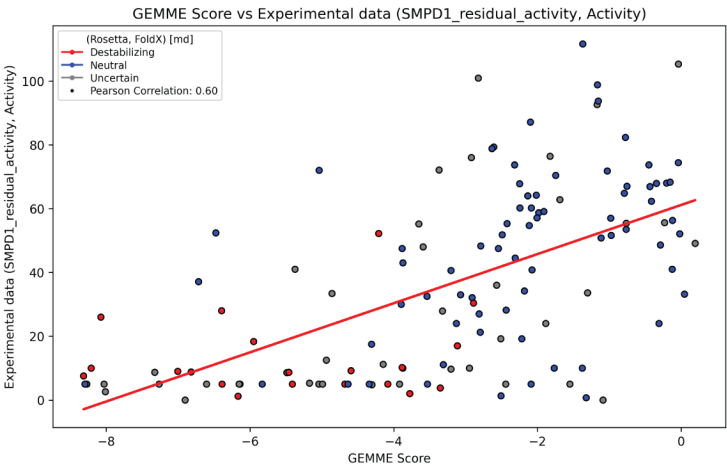
PILRA is a protein target connected to neurodegenerative diseases[99], along with KIF5A, CFAP410, and CYP2R1, illustrating the broad applicability of MAVISp to proteins involved in different diseases. Proteins associated with other diseases, such as TTR, SOD1, and SMPD1, have also been included in the MAVISp database. SMPD1 has been recently investigated in a targeted study using the *ensemble mode* of MAVISp together with other methodologies, validating our results by means of experimental data measuring the residual catalytic activity of enzyme variants[76]. As previously stated, MAVISp integrates curated experimental data for specific target proteins, which can be analyzed together with the results from the computational modules. To this goal, the dot plot representation provided by the downstream analyses toolkit of MAVISp and by the MAVISp database achieves a complete overview of both the experimental and the computational results (**Fig. 6a-c)** for SMPD1. Additionally, when a set of experimental data is available, it is possible to evaluate the correlation between predictions and experimental data (**Fig. 6d-e**). For SMPD1, we have obtained data on the residual catalytic activity of the enzyme for 135 variants[76], available in the literature. Thanks to the MAVISp protocol, we predicted the effect of amino acid substitutions on changes in folding free energies as well as data for predicted functional effects from VEPs, which can be compared with the experimental data. The score values produced by the VEPs were mildly correlated with the residual activity measurements (Pearson correlation coefficient ~0.6). Of note, most of the variants that have a predicted destabilizing effects on the stability are found at values of experimental residual activity lower than 20%, confirming what observed in our previous study[76] and suggesting that changes in stability for SMPD1 can be help identifying damaging variants of this enzyme. Nonetheless, in this case, the experimental readout cannot be explained by stability changes alone. Thus, variants found with low residual activity and functionally damaging (GEMME and DeMaSk scores lower than -3 and -0.25, respectively) and that are neutral for stability according to MAVISp are good candidates for further investigation. For example, biomolecular simulations or computational chemistry methods could be used to investigate the effects of these variants on the catalytic mechanism of the enzyme and its lipid transport. Finally, variants, such as Y500H, which have a low residual activity, high loss-of-fitness scores and are uncertain for the STABILITY module, can be analyzed for their propensity to fall in early folding regions (see entry in the MAVISp database) and could be investigated in the *ensemble mode* using enhanced sampling simulations to accurately estimate their folding free energy profiles.

a



b



**Fig. 6. MAVISp, GEMME, and DeMaSk predictions on the impact of SMPD1 variant subset.** A subset of SMPD1 variants, for which experimental data on enzyme activity have been selected, is shown with predictions from MAVISp, GEMME, and DeMaSk **a-b)** Scatter plots comparing DeMaSk (left) and GEMME (right) scores against experimental assay scores for enzymatic activity. The red line represents the regression, while the dotted line marks the threshold below which enzyme activity is considered inactive. Dots are colored based on the MAVISp STABILITY module classification (Rosetta/FoldX consensus): Destabilizing, Neutral, or Uncertain.

# Conclusions and Future Perspective

MAVISp provides a multi-layered assessment of the effects of variants found in cancer studies or other diseases using structural methods. MAVISp results are especially useful for variant interpretation and prioritization. These results can be useful as a complementary resource to available pathogenic scores or high-throughput experiments. MAVISp can help to pinpoint the effects linked to a pathogenic variant for further studies.

A significant advantage of MAVISp is its comprehensive coverage, expanding beyond clinically identified variants, by including novel variants yet to be characterized in other databases. This makes MAVISp a valuable resource for researchers and clinicians, facilitating the exploration of novel variants and their underlying pathogenic mechanisms. MAVISp can help on one side to associated mechanistic indicators to variants that are known or predicted pathogenic, as well as to aid in the characterization of the effects of VUS or variants with conflict evidence at the molecular level. Finally, we envision that MAVISp could become, in time, a community-driven effort and serve as a repository of data for the effects of disease-related variants more broadly. The results reported in MAVISp will provide an atlas of functional annotations for disease-related variants.

We have previously framed MAVISp in the context of others computational frameworks that collect data from different sources or integrate different structure-based methods to characterize variants[100]. This idea has led

659   to the production of different attempts. Missense3D[101] predicts the impact of variants on an array of structural
660   features; ADDRESS[102] includes predictions on stability and intermolecular contacts for variant found in Uni-
661   Prot humsavar; MUTATIONEXPLORER[103] uses Rosetta and RaSP to predict the effect of amino acid sub-
662   stitutions on stability or binding, on user-provided structures; VUStruct[104] selects relevant protein structures
663   for the protein of interest and performs a wide array of predictions, including on the effect of variants on
664   stability, binding surface, PTMs. The Genomics 2 Proteins[105] portal includes data from several sources, in-
665   cluding some overlapping with MAVISp such as Phosphosite or MaveDB, as well as features calculated on
666   the protein structure. ProtVar[106] also aggregates variant from different sources and includes both variant
667   effect predictors, prediction of change on stability upon amino acid substitution, as well as prediction of com-
668   plex structures. MAVISp, is, to our knowledge, the first resource to integrate data on binding free-energies,
669   data derived from molecular dynamics simulations, as well as experimental data from different sources, and
670   the first to integrate predictions on long-range effects as a database. While MAVISp has a lower coverage
671   than others, it includes carefully curated manual steps, such as during protein structure preparation and
672   simulation.
673   As the database grows, it will provide high quality data on different structural properties that can also be used
674   for benchmarking purposes or as features in machine learning models. To this goal, the stringent data col-
675   lection that we designed and present here is pivotal to build meaningful and accurate predictive models.
676   We would like to highlight previous studies that have demonstrated the usefulness of MAVISp and its proto-
677   cols. For example, we have showcased the versatility of MAVISp in characterizing the effects induced by a
678   redox post-translational modification of Cysteine (*S*-nitrosylation) using structural methods[107] . We focused
679   on variants found in cancer samples for their capability to alter the propensity of cysteine to be *S*-nitrosylated,
680   or a population-shift mechanism induced by the PTM. The collection of data using MAVISp modules has
681   been pivotal to aggregate variants for each target of interest in the study on *S*-nitrosylation. The pipelines
682   developed in the study of *S*-nitrosylation will be integrated within the MAVISp PTM module, extending it
683   beyond support for phosphorylation, which is currently supported by MAVISp.
684   Alterations in transcription factors are often linked to aberrant gene expression, including processes such as
685   proliferation, cell death, and other cancer hallmarks[108]. Different mechanisms are at the base of alterations
686   in the activity of transcription factors in cancer, including point mutations. A previous study on TP53 served
687   as a platform to develop different modules currently available in MAVISp [38]. We thus aim to expand the
688   MAVISp database to include more transcription factors. To this goal, one of the datasets under data collection
689   covers the protein targets from the TRRUST2 database[109], which includes experimentally characterized tran-
690   scription factors and their targets, of which 150 have been already processed and included in the MAVISp
691   database.
692   Furthermore, MAVISp provides pre-calculated values of changes in folding or binding free energies and other
693   metrics that can also be reanalyzed in the context of other research projects. With the examples on PTEN
694   and SMPD1 provided here, we introduced the curation of experimental data in MAVISp, as a source of ex-
695   perimental validation. The implementation of additional modules for MAVISp (e.g., degron[110] and aggregation
696   propensity[111]) would likely improve coverage of the diverse mechanisms regulating protein abundance. Of
697   note, MAVISp supports either data from multiplex assays of variant effects or experimental data from litera-
698   ture mining of the biocurators. The purpose of collecting experimental data is to validate our findings, update
699   protocols, and continuously improve the included methodologies. The reliability of our predictions depends
700   on their alignment with experimental results, which can be used as reference data to benchmark and improve
701   our predictions over time. The database currently includes 16 protein entries with experimental data.
702   At this stage, MAVISp can provide annotations for variants of transmembrane proteins exclusively in regions
703   that are not in contact with the membrane. Recently published approaches[112] could enable the application of
704   the STABILITY module to transmembrane regions as well. In addition, we will include support to intrinsically
705   disordered regions in the ensemble mode, designing new modules to reflect the most important properties of
706   these regions.
707   We foresee that MAVISp will provide a large amount of data on structure-based properties related to the
708   changes that  can exert at the protein level, which could be exploited for design of experimental biological
709   readouts, also towards machine-learning applications for variant assessment and classification or to under-
710   stand the importance of specific variants in connection with clinical variables, such as drug resistance, risk
711   of relapse and more.
712
713   # Methods
714

## *Initial structures for MAVISp and STRUCTURE_SELECTION module*

As a default, in the high-throughput data collection, we use models from the AlphaFold2 database[25] for most of the target proteins and trim them to remove regions with pLDDT scores < 70 at the N- or C-termini or very long disordered linkers between folded domains. For proteins coordinating cofactors, in the low-throughput targeted studies, we re-modeled the relevant cofactors upon analyses with AlphaFill[42] and where needed through MODELLER[113]. A summary of the initial structures used for each protein included in the database is reported in OSF (https://osf.io/y3p2x/). In selected cases, we have replaced long disordered loops with short residue stretches using a custom pipeline based on MODELLER (https://github.com/ELELAB/MAVISp_loop_replacer). This was done to avoid potential bias in our structural calculations, due to the arbitrary conformation of such loops and their spurious contacts with the rest of the structure. In addition, for proteins with transmembrane regions, we used the PPM (Positioning of Proteins in Membrane) server 3.0 from OPM (Orientations of Proteins in Membrane)[114,115]. For target proteins larger than 2700 residues, whose structures are not provided by the AlphaFold2 database, we model them using AlphaFold3.

The advantage of using AlphaFold-predicted structures in the default high-throughput data collection of MAVISp lies in their ability to achieve quality comparable to experimental data, as demonstrated in previous work[2], and at the same time circumventing limitations typically associated with experimental approaches, such as artifacts, missing atoms, and incomplete or absent residues.

## *INTERACTOME module*

In the INTERACTOME module, implemented in the freely available PPI2PDB toolkit (https://github.com/ELELAB/PPI2PDB), we identify known interactors of the target protein by extracting data from the Mentha database[47] and match them to available PDB structures, using the *mentha2pdb* script. Mentha2pdb also examines experimentally validated dimeric complexes generated with AlphaFold2 from the HuRI and HuMAP databases by Burke et al.[48]. Mentha2PDB provides annotations of the interactors and generates input files for AlphaFold-Multimer.

Complementarily, we retrieve interactors from the STRING database[49] and process them analogously using our STRING2PDB tool, which maps STRING interactions to available PDB structures. The tool restricts retrieval from the physical subnetwork of STRING with evidence of interaction supported by either curated database annotation or experimental data.

As a final step, we aggregate all interaction data for the target protein into a single table, ranking interactors primarily by Mentha and secondarily by STRING score to prioritize experimentally supported pairs. We then add complexes retrieved directly from the PDB via pdbminer-complexes (https://github.com/ELELAB/MAVISp_automatization/tree/main/mavisp_templates/) to capture interactions not yet reflected in PPI databases.

We also use other methods to identify four different classes of short linear motifs (BRCT, LIR, BH3 and UIM) in our target proteins. Depending on the type, we use a combination of simple regular expression matching, a method designed by us for structure-based identification of short linear motifs SLiMfast (available at https://github.com/ELELAB/SLiMfast) together with another method for predicting changes in secondary structure propensity that may be induced by phosphorylation in the core of putative LIR motifs, phosphor-iLIR (https://github.com/ELELAB/phospho-iLIR), or DeepLoc 2.0[116] for predicting the subcellular localization of the protein, especially useful for BRCT motifs.

## *Free energy calculations for STABILITY, LOCAL INTERACTION and LONG-RANGE modules*

We applied the BuildModel module of FoldX5 suite[117] averaging over five independent runs for the calculations of changes in free energy of folding upon amino acid substitution with MutateX and the FoldX5 method. We used the cartddg2020 protocol for folding free energy calculations with Rosetta suite and the ref2015 energy function. In this protocol, only one structure is generated at the relax step and then optimized in Cartesian space. Five rounds of Cartesian space optimization provide five pairs of wild-type and mutant structures for each variant. The change in folding free energy is then calculated on the pair characterized by the lower value of free energy for the mutant variant, as described in the original protocol[118].

We used MutateX to calculate changes in binding free energy for the LOCAL INTERACTION module using the BuildModel and AnalyzeComplex functions of FoldX5 suite and averaging over five runs. With Rosetta, we used the flexddg protocol as implemented in RosettaDDGPrediction and the talaris2014 energy function. We used 35,000 backrub trials and a threshold for the absolute score for minimization convergence of 1

772 Rosetta Energy Unit (REU). The protocol then generates an ensemble of 35 structures for each mutant var-
773 iant and calculates the average changes in binding free energy. We used Rosetta 2022.11 version for both
774 stability and binding calculations. In the applications with RosettaDDGPrediction the Rosetta Energy Units
775 (REUs) were converted to kcal/mol with available conversion factors[118]. We also applied RaSP using the
776 same protocol provided in the original publication[56] and adjusting the code in a workflow according to MA-
777 VISp-compatible formats (https://github.com/ELELAB/RaSP_workflow). We have included data on 131 com-
778 plexes at the date of 16/10/2025 (https://osf.io/y3p2x/ ).
779 For the calculations of allosteric free energy, we used the structure-based statistical mechanical model of
780 allostery (SBSMMA)[119,120] implemented in AlloSigMA2[64]. The model describes the mutated variants as 'UP'
781 or 'DOWN' mutations depending on difference in steric hindrance upon the substitution. We followed a re-
782 cently updated and benchmarked protocol[65]. In brief, we classified as uncertain those variants for which the
783 absolute changes in the volume of the side chain upon the amino acid substitution was lower than 5 $Å^3$, as
784 recently applied to p53[38]. As a default, we considered as having an effect only variants that were exposed to
785 the solvent (≥25% relative solvent accessibility of the side chain), with associated changes in absolute value
786 of allosteric free energy larger than 2 kcal/mol and considered as remote response sites those that were at
787 a distance higher than 5.5 Å from the mutation site, considering all heavy atoms, and which belongs to pock-
788 ets as identified by Fpocket[121] ( see workflow at https://github.com/ELELAB/MAVISp_allosigma2_workflow/)
789

*Efoldmine*
790
791
792 The EFOLDMINE module, integrated within the simple mode of MAVISp, predicts residues with early folding
793 propensity using the EfoldMine tool[72]. Trained on residue-level hydrogen/deuterium exchange nuclear mag-
794 netic resonance (HDX NMR) folding data from the Start2Fold database[122], this tool uses secondary structure
795 propensity and backbone/side-chain dynamics in a support-vector machine algorithm to predict early folding
796 regions based on the target's sequence.
797 In MAVISp, we incorporated EfoldMine to determine whether point
798 mutations in variants fall within the predicted early folding regions, using a threshold of 0.169 to define resi-
799 dues involved in early folding events as suggested by the developers of the method[72] and considering only
800 regions with a minimum length of three early folding residues to exclude isolated peaks.[71].
801

*FUNCTIONAL SITE module*
802
803
804 The FUNCTIONAL SITES module aids the identification of variants that might impact cofactor binding sites
805 or active site residues, as well as the residues within the second coordination sphere with respect to active
806 site residues of enzymes or their corresponding binding sites. It is based on a contact analysis performed
807 with the Arpeggio software[123]. Before the analysis, the model structure is subjected to energy minimization
808 with Conjugate Gradients[124] in 50 steps, using the MMFF94 force field[125], a van der Waals cutoff of 0.1, an
809 interacting cutoff of 5.0 Å, and a physiological pH of 7.4. Subsequently, the output is further preprocessed to
810 exclude clashes and proximal contacts (https://github.com/ELELAB/mavisp_accessory_tools).
811

*Molecular dynamics simulations for MAVISp ensemble mode*
812
813
814 We used either previously published[37,38,126–130] or newly collected one microsecond all-atom molecular dynam-
815 ics simulations performed using the CHARMM22* or CHARMM36m force fields[131]. All the simulations have
816 been carried out in the canonical ensemble after a final equilibration steps and using explicit solvent and
817 periodic boundary conditions. The templates files used for the simulations are provided in OSF
818 (https://osf.io/y3p2x/).
819 Ensembles generated using simulations are then subject to quality control, either using Mol_Analysis[132] or
820 MetaD_Analysis (https://github.com/ELELAB/MetaD-Analysis) tools.
821 As a first example of how we intend to use metadynamics data for the FUNCTIONAL_DYNAMICS module
822 we used the simulations from TP53 where the effects of amino acid substitutions on an interface for protein-
823 protein interaction (residues 207-213) was investigated. We used a collective variable based on distances
824 between two residues (D208-R156) that were effective in capturing open (active) and closed (inactive) con-
825 formations of the loop. See repositories associated with the enhanced sampling simulations of TP53[38]. All
826 the newly generated trajectories will be deposited as different entries in OSF, and the link is reported in the
827 metadata on the MAVISp webserver. At the date of 01/11/2024, we have included 45 protein targets in the
828 *ensemble mode* using as source of ensemble mostly unbiased MD simulations of 500 ns or one-μs, as

829 detailed in the corresponding metadata on the MAVISp webserver. In some cases, we included ensembles
830 generated by a coarse-grain model of flexibility or using the conformation provided by NMR structures from
831 the PDB (see INPUT STRUCTURES tables in https://osf.io/y3p2x/).

832
833
834 *Protein Structure Networks and path analysis for MAVISp ensemble mode*

835 In the ensemble mode we apply a module building upon the simple mode LONG_RANGE module. It uses
836 AlloSigma2-PSN (https://github.com/ELELAB/MAVISp_allosigma2_workflow/) where we constructed an
837 atomic-contact PSN on the full trajectories using PyInteraph2[66]. Pairs of residues were retained only if their
838 sequence distance exceeded Proxcut threshold of 1 and their edge calculations remained within less than
839 4.5Å, based on the thresholds described in PyInteraph2[66]. We retained edges with an occurrence greater
840 than Pcrit threshold of 50% across the ensemble frames, weighted on the interaction strength Imin of 3.
841 Subsequently, we used the path_analysis function of PyInteraph2 to identify the shortest paths of communi-
842 cation between each pair of AlloSigMA2[64] predicted mutation and respective response sites, using a mini-
843 mum distance threshold of 5.5 Å and retained paths that were four residues or longer.

844
845
846 *CABS-flex ensembles for MAVISp ensemble mode*

847 We used the coarse-grained CABS-flex 2.0 method and software[50] as a part of a Snakemake[133] pipeline,
848 available at https://github.com/ELELAB/MAVISp_CABSflex_pipeline. The pipeline includes the possibility to
849 tune the calculations by different restraints, secondary structure definition, ligand binding and more. It also
850 contains a quality control step to evaluate the secondary structure content of the generated structures with
851 respect to the starting one, using DSSP[134] and the SOV-refine score[135].

852
853 *Variant Effect Prediction*

854
855 We used DeMaSk[74], GEMME[14], EVE[75], REVEL[73] and AlphaMissense[26] as predictors for the effect of any
856 possible amino acid substitution to natural amino acids, on the full protein sequence of the main UniProt[136]
857 isoform of each protein. We used available default parameters for each method unless noted otherwise. We
858 used the standalone version of DeMaSk as available on its public GitHub (commit ID 10fa198), with BLAST+
859 2.13.0. We followed the protocol available on GitHub: we first generated the aligned homologs sequence file
860 by using the demask.homologs module and then calculated fitness impact predictions. Finally, we classified
861 as loss-of-fitness those variants having a DeMaSk delta fitness score in absolute value lower or equal to -
862 0.25, gain-of-fitness if the score is higher than 0.25, and neutral otherwise (**Supplementary Text S3**). We
863 used the available online webserver to obtain variant effect predictions with GEMME, upon setting the num-
864 ber of JET iterations to 5, to obtain more precise results. [126] We have classified variants having a GEMME
865 score <= -3 as damaging, and neutral otherwise. Thresholds were selected according to our benchmarking
866 (**Supplementary Text S3).** To obtain EVE scores, we have used the scripts, protocol and parameters avail-
867 able on the EVE GitHub (commit iD 740b0a7) as part of a custom-built Snakemake[133]-based pipeline, avail-
868 able at https://github.com/ELELAB/MAVISp_EVE_pipeline Using EVE first requires building a protein-spe-
869 cific Bayesian variational autoencoder model, which learns evolutionary constraints between residues from
870 a multiple sequence alignment. In the current MAVISp release, we generated such alignments using EVcou-
871 plings[137], using the Uniref100[138] sequence database released on 01/03/2023, by keeping sequences with at
872 least 50% of coverage with the target protein sequence, alignment positions with a minimum of 70% residue
873 occupancy, and using a bit score threshold for inclusion of 0.5 bits with no further hyperparameter exploration.
874 We then used our pipeline to perform model training, calculation of the evolutionary index, and used a global-
875 local mixture of Gaussian Mixture Models to obtain a pathogenicity score and classification. We have used
876 pre-computed REVEL scores for variants as available in dbSNFP[139,140], accessed through myvari-
877 ants.info[141,142], as implemented in Cancermuts. We have classified as damaging variants that have a REVEL
878 score larger or equal to 0.5[143]. We included AlphaMissense pathogenicity prediction scores and classification
879 as available by the dataset of prediction for all possible amino acid substitutions in UniProt canonical
880 isoforms, release version 2[144].

881
882
883 *Annotations from experimental data for EXPERIMENTAL_DATA module*

884 We developed Python scripts to identify the overlap in coverage between the Mave database (MaveDB)[92]
885 and MAVISp, and to retrieve the score sets associated with the shared entries from the MaveDB[92] database

886 through their API (https://api.mavedb.org/docs). Where available, we also extracted information on methods
887 and classification thresholds. For entries where this information was incomplete, the corresponding publica-
888 tions were manually reviewed to extract thresholds for variant classification.
889 The ProteinGym[91] repository was locally downloaded from GitHub, and a custom Python script was used to
890 process the datasets based on the reference files provided in the repository. The datasets used for the anal-
891 ysis contained the experimental scores and the classification provided by the authors either based on the
892 median of the score distributions or via manual annotation. The scores and their classifications were then
893 integrated into the final database file generated by MAVISp. The aggregated scores, along with their classi-
894 fications, were compiled into the final database file produced by MAVISp through a module dedicated to the
895 experimental data.
896
897 *Identification of RefSeq identifiers*
898
899 To ensure the correct RefSeq annotations in MAVISp, we implemented a Python tool, *compare_seq.py*
900 (https://github.com/ELELAB/mavisp_accessory_tools/), to verify the sequence identity between the ca-
901 nonical UniProt sequence used in our analyses and the corresponding RefSeq protein identifier to be used
902 for the ClinVar search. The Uniprot sequences were retrieved using the UniProt REST API, while the RefSeq
903 protein sequences were fetched from the NCBI Entrez Protein database. We implemented a global pairwise
904 alignment using the Biophyton pairwise2 module with the globalxx scheme to assess sequence identity. Each
905 comparison was classified as either an exact match, a mismatch (identity <100%), or unresolved due to
906 missing or unresolvable sequences. To improve performances, the analyses were parallelized using multi-
907 threading via Python concurrent.futures. The results were logged into structured CSV reports for consultation.
908 This allows data managers to identify exisiting entries in MAVISp with  RefSeq identifiers inconsistent with
909 provided UniProt accession code and assign them to biocurators for entry review.

910 Additionally,    we    provide    the    biocurators    with    a    Python-based    script    (*uniprot2refseq,*    l
911 https://github.com/ELELAB/mavisp_accessory_tools/) that identifies RefSeq IDs for the UniProt canoni-
912 cal protein isoform. For each UniProt AC, we queried the UniProt REST API to obtain RefSeq protein cross-
913 references (NP_* IDs) from the canonical entry in JSON format. Only protein-level RefSeq entries were con-
914 sidered. The canonical UniProt protein sequence was downloaded in FASTA format, and each RefSeq se-
915 quence was retrieved from the NCBI Protein database using Biopython and the Entrez API. Pairwise global
916 alignments were performed using the Biopython pairwise2 module and we estimate the percentage sequence
917 identity as the number of identical residues over the length of the longer sequence. Results were saved in
918 tabular format, including UniProt AC, RefSeq ID, and sequence identity. This approach aids the biocurators
919 to identify the RefSeq IDs for the canonical isoform of the protein undebefore starting with the data collection.
920 The script is expected to be used by the biocurators before each run with the MAVISp automatization work-
921 flow described below.

922
923 *Workflows for automatization and data collection within MAVISp*
924
925 We provide and maintain two Snakemake workflows for the data collection of the default modules of MAVISp.
926 The first is a Snakemake pipeline to automate MutateX runs as much as possible. It is designed to automat-
927 ically download the chosen structure(s) from the AlphaFold structural database, or a custom structure input
928 file, when necessary, trim them as requested, and generate desired MutateX folding free energy scans with
929 a predictable directory structure. It only requires as input a csv file with metadata on the desired scan and a
930 configuration file with details on the run to be performed. It is available at https://github.com/ELELAB/mu-
931 tatex_pipelines/tree/main/custom_collect_scan.
932 Once such a scan is available, it is possible to use a second Snakemake pipeline, called MAVISp_automati-
933 zation, which performs most of the steps that are necessary to annotate a protein for a MAVISp simple mode
934 entry. Similarly to the previous pipeline, it only requires metadata on the target protein to be analyzed, as
935 well as a MutateX mutational scan. It generates a dataset that can then be imported into the MAVISp data-
936 base, except for predictions performed using Rosetta-based methods, since these are much more computa-
937 tionally expensive and need to be performed separately using the RosettaDDGPrediction pipeline[55]. Using a
938 Snakemake pipeline allows to improve efficiency and scalability, allowing to use multi-core system to process
939 several proteins or perform different analyses in parallel. It is available at https://github.com/ELELAB/MA-
940 VISp_automatization.

*Structure-based assessment of variants with MAVISp*

## Data Availability

The data can either be consulted through our web server (https://services.healthtech.dtu.dk/services/MA-VISp-1.0/) or as individual CSV files in the OSF repository https://osf.io/ufpzm/. Other raw data and utilities can be found at the MAVISp extended data OSF repository (https://osf.io/y3p2x/) Reports for several proteins are available at https://elelab.gitbook.io/mavisp/.

## References

1. Xuan, J., Yu, Y., Qing, T., Guo, L. & Shi, L. Next-generation sequencing in the clinic: Promises and challenges. *Cancer Lett* **340**, 284–295 (2013).

2. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021).

3. Weile, J. & Roth, F. P. Multiplexed assays of variant effects contribute to a growing genotype–phenotype atlas. *Human Genetics 2018 137:9* **137**, 665–678 (2018).

4. Fowler, D. M. *et al.* An Atlas of Variant Effects to understand the genome at nucleotide resolution. *Genome Biol* **24**, 147 (2023).

5. Fowler, D. M. & Rehm, H. L. Will variants of uncertain significance still exist in 2030? *Am J Hum Genet* **111**, 5–10 (2024).

6. Burke, W., Parens, E., Chung, W. K., Berger, S. M. & Appelbaum, P. S. The Challenge of Genetic Variants of Uncertain Clinical Significance. *Ann Intern Med* **175**, 994–1000 (2022).

7. Gao, J. *et al.* Integrative analysis of complex cancer genomics and clinical profiles using the cBioPortal. *Sci Signal* **6**, pl1 (2013).

8. Cerami, E. *et al.* The cBio Cancer Genomics Portal: An open platform for exploring multidimensional cancer genomics data. *Cancer Discov* https://doi.org/10.1158/2159-8290.CD-12-0095 (2012) doi:10.1158/2159-8290.CD-12-0095.

9. Landrum, M. J. *et al.* ClinVar: improvements to accessing data. *Nucleic Acids Res* **48**, D835–D844 (2020).

10. Tate, J. G. *et al.* COSMIC: the Catalogue Of Somatic Mutations In Cancer. *Nucleic Acids Res* **47**, D941–D947 (2019).

11. Stenson, P. D. *et al.* The Human Gene Mutation Database (HGMD®): optimizing its use in a clinical diagnostic or research setting. *Hum Genet* **139**, 1197–1207 (2020).

12. Burke, W., Parens, E., Chung, W. K., Berger, S. M. & Appelbaum, P. S. The Challenge of Genetic Variants of Uncertain Clinical Significance: A Narrative Review. *Ann Intern Med* **175**, 994–1000 (2022).

13. Riesselman, A. J., Ingraham, J. B. & Marks, D. S. Deep generative models of genetic variation capture the effects of mutations. *Nat Methods* **15**, 816–822 (2018).

14. Laine, E., Karami, Y. & Carbone, A. GEMME: A Simple and Fast Global Epistatic Model Predicting Mutational Effects. *Mol Biol Evol* **36**, 2604–2619 (2019).

15. Høie, M. H., Cagiada, M., Beck Frederiksen, A. H., Stein, A. & Lindorff-Larsen, K. Predicting and interpreting large-scale mutagenesis data using analyses of protein stability and conservation. *Cell Rep* **38**, 110207 (2022).

16. Cagiada, M. *et al.* Understanding the Origins of Loss of Protein Function by Analyzing the Effects of Thousands of Variants on Activity and Abundance. *Mol Biol Evol* **38**, 3235–3246 (2021).

17. Jepsen, M. M., Fowler, D. M., Hartmann-Petersen, R., Stein, A. & Lindorff-Larsen, K. Classifying disease-associated variants using measures of protein activity and stability. *Protein Homeostasis Diseases* 91–107 (2020) doi:10.1016/B978-0-12-819132-3.00005-1.

18. Gelman, H. *et al.* Recommendations for the collection and use of multiplexed functional data for clinical variant interpretation. *Genome Med* **11**, 1–11 (2019).

19. McEwen, A. E., Tejura, M., Fayer, S., Starita, L. M. & Fowler, D. M. Multiplexed assays of variant effect for clinical variant interpretation. *Nat Rev Genet* https://doi.org/10.1038/S41576-025-00870-X (2025) doi:10.1038/S41576-025-00870-X.

20. Pejaver, V. *et al.* Calibration of computational tools for missense variant pathogenicity classification and ClinGen recommendations for PP3/BP4 criteria. *Am J Hum Genet* **109**, 2163–2177 (2022).

995 21. Gerasimavicius, L., Teichmann, S. A. & Marsh, J. A. Leveraging protein structural information to im-
996 prove variant effect prediction. *Curr Opin Struct Biol* **92**, 103023 (2025).

997 22. Livesey, B. J. & Marsh, J. A. Variant effect predictor correlation with functional assays is reflective
998 of clinical classification performance. *Genome Biol* **26**, 1–27 (2025).

999 23. Rastogi, R. *et al.* Critical assessment of missense variant effect predictors on disease-relevant variant
000 data. *Hum Genet* **144**, 281–293 (2025).

001 24. Evans, R. *et al.* Protein complex prediction with AlphaFold-Multimer. *biorxiv*
002 https://doi.org/10.1101/2021.10.04.463034 (2022) doi:10.1101/2021.10.04.463034.

003 25. Varadi, M. *et al.* AlphaFold Protein Structure Database: massively expanding the structural coverage
004 of protein-sequence space with high-accuracy models. *Nucleic Acids Res* **50**, D439–D444 (2022).

005 26. Cheng, J. *et al.* Accurate proteome-wide missense variant effect prediction with AlphaMissense. *Sci-
006 ence* **381**, eadg7492 (2023).

007 27. Rives, A. *et al.* Biological structure and function emerge from scaling unsupervised learning to 250
008 million protein sequences. *Proceedings of the National Academy of Sciences* **118**, e2016239118
009 (2021).

010 28. Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a language model.
011 *Science* **379**, 1123–1130 (2023).

012 29. Hayes, T. *et al.* Simulating 500 million years of evolution with a language model. *Science* **387**, 850–
013 858 (2025).

014 30. Brandes, N., Goldman, G., Wang, C. H., Ye, C. J. & Ntranos, V. Genome-wide prediction of disease
015 variant effects with a deep protein language model. *Nature Genetics 2023 55:9* **55**, 1512–1522
016 (2023).

017 31. Sun, Y. & Shen, Y. Structure-informed protein language models are robust predictors for variant ef-
018 fects. *Hum Genet* 1–17 (2024) doi:10.1007/S00439-024-02695-W.

019 32. Blaabjerg, L. M., Jonsson, N., Boomsma, W., Stein, A. & Lindorff-Larsen, K. SSEmb: A joint em-
020 bedding of protein sequence and structure enables robust variant effect predictions. *Nature Communi-
021 cations 2024 15:1* **15**, 1–9 (2024).

022 33. Tekpinar, M., David, L., Henry, T. & Carbone, A. PRESCOTT: a population aware, epistatic, and
023 structural model accurately predicts missense effects. *Genome Biol* **26**, 1–42 (2025).

024 34. Banerjee, A., Bogetti, A. T. & Bahar, I. Accurate identification and mechanistic evaluation of patho-
025 genic missense variants with Rhapsody-2. *Proceedings of the National Academy of Sciences* **122**,
026 e2418100122 (2025).

027 35. Hollingsworth, S. A. & Dror, R. O. Molecular Dynamics Simulation for All. *Neuron* **99**, 1129–1143
028 (2018).

029 36. Kumar, M. & Papaleo, E. A pan-cancer assessment of alterations of the kinase domain of ULK1, an
030 upstream regulator of autophagy. *Sci Rep* **10**, 14874 (2020).

031 37. Fas, B. A. *et al.* The conformational and mutational landscape of the ubiquitin-like marker for autoph-
032 agosome formation in cancer. *Autophagy* 1–24 (2020).

033 38. Degn, K. *et al.* Cancer-related Mutations with Local or Long-range Effects on an Allosteric Loop of
034 p53. *J Mol Biol* **434**, 167663 (2022).

035 39. Berman, H. M. *et al.* The Protein Data Bank. *Nucleic Acids Res* **28**, 235–242 (2000).

036 40. Abramson, J. *et al.* Accurate structure prediction of biomolecular interactions with AlphaFold 3. *Na-
037 ture* **630**, 493–500 (2024).

038 41. Degn, K., Beltrame, L., Tiberti, M. & Papaleo, E. PDBminer to Find and Annotate Protein Structures
039 for Computational Analysis. *J Chem Inf Model* **63**, 7274–7281 (2023).

040 42. Hekkelman, M. L., de Vries, I., Joosten, R. P. & Perrakis, A. AlphaFill: enriching AlphaFold models
041 with ligands and cofactors. *Nat Methods* **20**, 205–213 (2022).

042 43. Zhu, W., Shenoy, A., Kundrotas, P. & Elofsson, A. Evaluation of AlphaFold-Multimer prediction on
043 multi-chain protein complexes. *Bioinformatics* **39**, (2023).

044 44. Tsaban, T. *et al.* Harnessing protein folding neural networks for peptide-protein docking. *Nat Com-
045 mun* **13**, (2022).

046 45. Di Rita, A. *et al.* HUWE1 E3 ligase promotes PINK1/PARKIN-independent mitophagy by regulating
047 AMBRA1 activation via IKKα. *Nat Commun* **9**, 3755 (2018).

048   46.   Holdgaard, S. G. *et al.* Selective autophagy maintains centrosome integrity and accurate mitosis by
049         turnover of centriolar satellites. *Nat Commun* **10**, 1–19 (2019).
050   47.   Calderone, A., Castagnoli, L. & Cesareni, G. Mentha: A resource for browsing integrated protein-in-
051         teraction networks. *Nat Methods* **10**, 690–691 (2013).
052   48.   Burke, D. F. *et al.* Towards a structurally resolved human protein interaction network. *Nat Struct Mol*
053         *Biol* **30**, (2023).
054   49.   Szklarczyk, D. *et al.* The STRING database in 2025: protein networks with directionality of regula-
055         tion. *Nucleic Acids Res* **53**, D730–D737 (2025).
056   50.   Kuriata, A. *et al.* CABS-flex 2.0: A web server for fast simulations of flexibility of protein structures.
057         *Nucleic Acids Res* **46**, W338–W343 (2018).
058   51.   Abraham, M. J. *et al.* GROMACS: High performance molecular simulations through multi-level par-
059         allelism from laptops to supercomputers. *SoftwareX* **2**, 19–25 (2015).
060   52.   Tribello, G. A., Bonomi, M., Branduardi, D., Camilloni, C. & Bussi, G. PLUMED 2: New feathers
061         for an old bird. *Comput Phys Commun* **185**, 604–613 (2014).
062   53.   Bonomi, M. *et al.* Promoting transparency and reproducibility in enhanced molecular simulations.
063         *Nat Methods* **16**, 670–673 (2019).
064   54.   Tiberti, M. *et al.* MutateX: an automated pipeline for in silico saturation mutagenesis of protein struc-
065         tures and structural ensembles. *Brief Bioinform* **23**, bbac074 (2022).
066   55.   Sora, V. *et al.* RosettaDDGPrediction for high-throughput mutational scans: from stability to binding.
067         *Protein Science* **32**, e4527 (2023).
068   56.   Blaabjerg, L. M. *et al.* Rapid protein stability prediction using deep learning representations. *Elife* **12**,
069         e82593 (2023).
070   57.   Nielsen, S. V. *et al.* Predicting the impact of Lynch syndrome-causing missense mutations from struc-
071         tural calculations. *PLoS Genet* **13**, e1006739 (2017).
072   58.   Abildgaard, A. B. *et al.* Computational and cellular studies reveal structural destabilization and degra-
073         dation of MLH1 variants in Lynch syndrome. *Elife* **8**, e49138 (2019).
074   59.   Jankauskaite, J., Jiménez-García, B., Dapkunas, J., Fernández-Recio, J. & Moal, I. H. SKEMPI 2.0:
075         An updated benchmark of changes in protein-protein binding energy, kinetics and thermodynamics
076         upon mutation. *Bioinformatics* **35**, 462–469 (2019).
077   60.   Sampson, J. M. *et al.* Robust Prediction of Relative Binding Energies for Protein–Protein Complex
078         Mutations Using Free Energy Perturbation Calculations. *J Mol Biol* **436**, (2024).
079   61.   Sargsyan, K. & Lim, C. Using protein language models for protein interaction hot spot prediction
080         with limited data. *BMC Bioinformatics* **25**, (2024).
081   62.   Sapozhnikov, Y., Patel, J. S., Ytreberg, F. M. & Miller, C. R. Statistical modeling to quantify the un-
082         certainty of FoldX-predicted protein folding and binding stability. *BMC Bioinformatics* **24**, 426
083         (2023).
084   63.   Peccati, F., Alunno-Rufini, S. & Jiménez-Osés, G. Accurate Prediction of Enzyme Thermostabiliza-
085         tion with Rosetta Using AlphaFold Ensembles. *J Chem Inf Model* **63**, 898–909 (2023).
086   64.   Tan, Z. W., Guarnera, E., Tee, W. V. & Berezovsky, I. N. AlloSigMA 2: Paving the way to designing
087         allosteric effectors and to exploring allosteric effects of mutations. *Nucleic Acids Res* **48**, W116–
088         W124 (2020).
089   65.   Krzesińska, K. *et al.* Deciphering Long-Range Effects of Mutations: An Integrated Approach Using
090         Elastic Network Models and Protein Structure Networks. *J Mol Biol* **437**, 169359 (2025).
091   66.   Sora, V. *et al.* PyInteraph2 and PyInKnife2 to Analyze Networks in Protein Structural Ensembles. *J*
092         *Chem Inf Model* **63**, 4237–4245 (2023).
093   67.   Tiberti, M. *et al.* The Cancermuts software package for the prioritization of missense cancer variants:
094         a case study of AMBRA1 in melanoma. *Cell Death Dis* **13**, 872 (2022).
095   68.   Orioli, S., Henning Hansen, C. G. & Lindorff-Larsen, K. Transient exposure of a buried phosphoryla-
096         tion site in an autoinhibited protein. *Biophys J* **121**, 91–101 (2022).
097   69.   Henriques, J. & Lindorff-Larsen, K. Protein Dynamics Enables Phosphorylation of Buried Residues
098         in Cdk2/Cyclin-A-Bound p27. *Biophys J* **119**, 2010–2018 (2020).
099   70.   Potel, C. M. *et al.* Impact of phosphorylation on thermal stability of proteins. *Nat Methods* **18**, 757–
100         759 (2021).

71. Huang, H. *et al.* iPTMnet: an integrated resource for protein post-translational modification network discovery. *Nucleic Acids Res* **46**, D542–D550 (2018).

72. Raimondi, D., Orlando, G., Pancsa, R., Khan, T. & Vranken, W. F. Exploring the Sequence-based Prediction of Folding Initiation Sites in Proteins. *Sci Rep* **7**, 8826 (2017).

73. Ioannidis, N. M. *et al.* REVEL: An Ensemble Method for Predicting the Pathogenicity of Rare Missense Variants. *Am J Hum Genet* **99**, 877–885 (2016).

74. Munro, D. & Singh, M. DeMaSk: A deep mutational scanning substitution matrix and its use for variant impact prediction. *Bioinformatics* **36**, 5322–5329 (2020).

75. Frazer, J. *et al.* Disease variant prediction with deep generative models of evolutionary data. *Nature* **599**, 91–95 (2021).

76. Scrima, S., Lambrughi, M., Tiberti, M., Fadda, E. & Papaleo, E. ASM variants in the spotlight: A structure-based atlas for unraveling pathogenic mechanisms in lysosomal acid sphingomyelinase. *Biochimica et Biophysica Acta (BBA) - Molecular Basis of Disease* **1870**, 167260 (2024).

77. Antonescu, O. N. *et al.* Decoding Phospho-Regulation and Flanking Regions in Autophagy-Associated Short Linear Motifs: A Case Study of Optineurin-LC3B Interaction. *BioRxiv* 1–43 (2023) doi:10.1101/2023.09.30.560296.

78. Sahu, S. *et al.* AVENGERS: Analysis of Variant Effects using Next Generation sequencing to Enhance BRCA2 Stratification. *BioRxiv* 1–31 (2023) doi:10.1101/2023.12.14.571713.

79. Repana, D. *et al.* The Network of Cancer Genes (NCG): a comprehensive catalogue of known and candidate cancer genes from cancer sequencing screens. *Genome Biol* **20**, 1 (2019).

80. Cagiada, M. *et al.* Understanding the Origins of Loss of Protein Function by Analyzing the Effects of Thousands of Variants on Activity and Abundance. *Mol Biol Evol* **38**, 3235–3246 (2021).

81. Manthei, K. A. & Keck, J. L. The BLM dissolvasome in DNA replication and repair. *Cellular and Molecular Life Sciences 2013 70:21* **70**, 4067–4084 (2013).

82. Miao, Y. *et al.* Cycloheximide (CHX) Chase Assay to Examine Protein Half-life. *Bio Protoc* **13**, (2023).

83. McKinnon, K. M. Flow Cytometry: An Overview. *Curr Protoc Immunol* **120**, 5.1.1-5.1.11 (2018).

84. Matreyek, K. A. *et al.* Multiplex assessment of protein variant abundance by massively parallel sequencing. *Nat Genet* **50**, 874–882 (2018).

85. Levy, E. D., Kowarzyk, J. & Michnick, S. W. High-Resolution Mapping of Protein Concentration Reveals Principles of Proteome Architecture and Adaptation. *Cell Rep* **7**, 1333–1340 (2014).

86. Tsuboyama, K. *et al.* Mega-scale experimental analysis of protein folding stability in biology and design. *Nature* **620**, 434–444 (2023).

87. Yen, H.-C. S., Xu, Q., Chou, D. M., Zhao, Z. & Elledge, S. J. Global Protein Stability Profiling in Mammalian Cells. *Science* **322**, 918–923 (2008).

88. Post, K. L. *et al.* Multi-model functionalization of disease-associated PTEN missense mutations identifies multiple molecular mechanisms underlying protein dysfunction. *Nat Commun* **11**, (2020).

89. Mighell, T. L., Thacker, S., Fombonne, E., Eng, C. & O'Roak, B. J. An Integrated Deep-Mutational-Scanning Approach Provides Clinical Insights on PTEN Genotype-Phenotype Relationships. *Am J Hum Genet* **106**, 818–829 (2020).

90. Post, K. L. *et al.* Multi-model functionalization of disease-associated PTEN missense mutations identifies multiple molecular mechanisms underlying protein dysfunction. *Nat Commun* **11**, (2020).

91. Notin, P. *et al.* ProteinGym: Large-Scale Benchmarks for Protein Fitness Prediction and Design. *Adv Neural Inf Process Syst* **36**, 64331–64379 (2023).

92. Esposito, D. *et al.* MaveDB: An open-source platform to distribute and interpret data from multiplexed assays of variant effect. *Genome Biol* **20**, (2019).

93. Rubin, A. F. *et al.* MaveDB 2024: a curated community database with over seven million variant effects from multiplexed functional assays. *Genome Biol* **26**, 13 (2025).

94. Hanahan, D. Hallmarks of Cancer: New Dimensions. *Cancer Discov* **12**, 31–46 (2022).

95. Kønig, S. M., Rissler, V., Terkelsen, T., Lambrughi, M. & Papaleo, E. Alterations of the interactome of Bcl-2 proteins in breast cancer at the transcriptional, mutational and structural level. *PLoS Comput Biol* **15**, e1007485 (2019).

96. Suzuki, H., Osawa, T., Fujioka, Y. & Noda, N. N. Structural biology of the core autophagy machinery. *Curr Opin Struct Biol* **43**, 10–17 (2017).

155 97.  Rogov, V. V. *et al.* Atg8 family proteins, LIR/AIM motifs and other interaction modes. *Autophagy*
156      *Reports* **2**, (2023).

157 98.  Geng, C., Xue, L. C., Roel-Touris, J. & Bonvin, A. M. J. J. Finding the ΔΔG spot: Are predictors of
158      binding affinity changes upon mutations in protein–protein interactions ready for it? *WIREs Compu-*
159      *tational Molecular Science* **9**, (2019).

160 99.  Li, Y. *et al.* Genomics of Alzheimer's disease implicates the innate and adaptive immune systems.
161      *Cellular and Molecular Life Sciences* **78**, 7397–7426 (2021).

162 100. Arnaudi, M., Utichi, M., Tiberti, M. & Papaleo, E. Predicting the structure-altering mechanisms of
163      disease variants. *Curr Opin Struct Biol* **91**, 102994 (2025).

164 101. Khanna, T., Hanna, G., Sternberg, M. J. E. & David, A. Missense3D-DB web catalogue: an atom-
165      based analysis and repository of 4M human protein-coding genetic variants. *Hum Genet* **140**, 805–
166      812 (2021).

167 102. Woodard, J., Zhang, C. & Zhang, Y. ADDRESS: A Database of Disease-associated Human Variants
168      Incorporating Protein Structure and Folding Stabilities. *J Mol Biol* **433**, 166840 (2021).

169 103. Philipp, M. *et al.* MutationExplorer: a webserver for mutation of proteins and 3D visualization of en-
170      ergetic impacts. *Nucleic Acids Res* **52**, W132–W139 (2024).

171 104. Moth, C. W. *et al.* VUStruct: a compute pipeline for high throughput and personalized structural biol-
172      ogy. *bioRxiv* https://doi.org/10.1101/2024.08.06.606224 (2024) doi:10.1101/2024.08.06.606224.

173 105. Kwon, S. *et al.* Genomics 2 Proteins portal: a resource and discovery tool for linking genetic screen-
174      ing outputs to protein sequences and structures. *Nat Methods* **21**, 1947–1957 (2024).

175 106. Stephenson, J. D. *et al.* ProtVar: mapping and contextualizing human missense variation. *Nucleic Ac-*
176      *ids Res* **52**, W140–W147 (2024).

177 107. Papaleo, E. *et al.* TRAP1 S-nitrosylation as a model of population-shift mechanism to study the ef-
178      fects of nitric oxide on redox-sensitive oncoproteins. *Cell Death Dis* **14**, (2023).

179 108. Hanahan, D. Hallmarks of Cancer: New Dimensions. *Cancer Discov* **12**, 31–46 (2022).

180 109. Han, H. *et al.* TRRUST v2: an expanded reference database of human and mouse transcriptional reg-
181      ulatory interactions. *Nucleic Acids Res* **46**, D380–D386 (2018).

182 110. Larsen, F. B. *et al.* Comprehensive degron mapping in human transcription factors. *bioRxiv*
183      2025.05.16.654404 (2025) doi:10.1101/2025.05.16.654404.

184 111. Zambrano, R. *et al.* AGGRESCAN3D (A3D): server for prediction of aggregation properties of pro-
185      tein structures. *Nucleic Acids Res* **43**, W306–W313 (2015).

186 112. Katarina Sooe Tiemann, J., Zschach, H., Lindorff-Larsen, K. & Stein, A. Interpreting the molecular
187      mechanisms of disease variants in human membrane proteins.
188      https://doi.org/10.1101/2022.07.12.499731 doi:10.1101/2022.07.12.499731.

189 113. Webb, B. & Sali, A. Comparative Protein Structure Modeling Using MODELLER. *Curr Protoc Bio-*
190      *informatics* **54**, 5.6.1-5.6.37 (2016).

191 114. Lomize, A. L., Todd, S. C. & Pogozheva, I. D. Spatial arrangement of proteins in planar and curved
192      membranes by PPM 3.0. *Protein Sci* **31**, 209–220 (2022).

193 115. Lomize, M. A., Pogozheva, I. D., Joo, H., Mosberg, H. I. & Lomize, A. L. OPM database and PPM
194      web server: resources for positioning of proteins in membranes. *Nucleic Acids Res* **40**, D370–D376
195      (2012).

196 116. Thumuluri, V., Almagro Armenteros, J. J., Johansen, A. R., Nielsen, H. & Winther, O. DeepLoc 2.0:
197      multi-label subcellular localization prediction using protein language models. *Nucleic Acids Res* **50**,
198      W228–W234 (2022).

199 117. Delgado, J., Radusky, L. G., Cianferoni, D. & Serrano, L. FoldX 5.0: Working with RNA, small mol-
200      ecules and a new graphical interface. *Bioinformatics* 1–2 (2019) doi:10.1093/bioinformatics/btz184.

201 118. Frenz, B. *et al.* Prediction of Protein Mutational Free Energy: Benchmark and Sampling Improve-
202      ments Increase Classification Accuracy. *Front Bioeng Biotechnol* **8**, (2020).

203 119. Tee, W.-V., Guarnera, E. & Berezovsky, I. N. Reversing allosteric communication: From detecting
204      allosteric sites to inducing and tuning targeted allosteric response. *PLoS Comput Biol* **14**, e1006228
205      (2018).

206 120. Guarnera, E. & Berezovsky, I. N. Structure-Based Statistical Mechanical Model Accounts for the
207      Causality and Energetics of Allosteric Communication. *PLoS Comput Biol* **12**, e1004678 (2016).

208  121. Le Guilloux, V., Schmidtke, P. & Tuffery, P. Fpocket: An open source platform for ligand pocket de-
209       tection. *BMC Bioinformatics* **10**, 1–11 (2009).
210  122. Pancsa, R., Varadi, M., Tompa, P. & Vranken, W. F. Start2Fold: a database of hydrogen/deuterium
211       exchange data on protein folding and stability. *Nucleic Acids Res* **44**, D429–D434 (2016).
212  123. Jubb, H. C. *et al.* Arpeggio: A Web Server for Calculating and Visualising Interatomic Interactions in
213       Protein Structures. *J Mol Biol* **429**, 365–371 (2017).
214  124. Richard, J. & August, S. An Introduction to the Conjugate Gradient Method Without the Agonizing
215       Pain. *Science* (1994).
216  125. Halgren, T. A. Merck molecular force field. I. Basis, form, scope, parameterization, and performance
217       of MMFF94. *J Comput Chem* **17**, (1996).
218  126. Invernizzi, G., Tiberti, M., Lambrughi, M., Lindorff-Larsen, K. & Papaleo, E. Communication routes
219       in ARID domains between distal residues in helix 5 and the DNA-binding loops. *PLoS Comput Biol*
220       **10**, e1003744 (2014).
221  127. Nygaard, M. *et al.* The mutational landscape of the oncogenic MZF1 SCAN domain in cancer. *Front*
222       *Mol Biosci* **3**, (2016).
223  128. Salamanca Viloria, J., Allega, M. F., Lambrughi, M. & Papaleo, E. An optimal distance cutoff for
224       contact-based protein structure networks using side chain center of masses. *Sci Rep* **7**, 2838 (2016).
225  129. Papaleo, E., Sutto, L., Gervasio, F. L. & Lindorff-Larsen, K. Conformational changes and free ener-
226       gies in a proline isomerase. *J Chem Theory Comput* **10**, 4169–4174 (2014).
227  130. Srinivasan, S. *et al.* The conformational plasticity of structurally unrelated lipid transport proteins
228       correlates with their mode of action. *PLoS Biol* **22**, e3002737 (2024).
229  131. Piana, S., Lindorff-Larsen, K. & Shaw, D. E. How robust are protein folding simulations with respect
230       to force field parameterization? *Biophys J* **100**, (2011).
231  132. Lambrughi, M. *et al. Analyzing Biomolecular Ensembles*. *Methods in Molecular Biology* vol. 2022
232       (2019).
233  133. Mölder, F. *et al.* Sustainable data analysis with Snakemake. *F1000Res* **10**, 33 (2021).
234  134. Kabsch, W. & Sander, C. Dictionary of protein secondary structure: Pattern recognition of hydrogen-
235       bonded and geometrical features. *Biopolymers* https://doi.org/10.1002/bip.360221211 (1983)
236       doi:10.1002/bip.360221211.
237  135. Liu, T. & Wang, Z. SOV-refine: A further refined definition of segment overlap score and its signifi-
238       cance for protein structure similarity. *Source Code Biol Med* **13**, 1–10 (2018).
239  136. Consortium, T. U. UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic Acids Res* **51**,
240       D523–D531 (2023).
241  137. Hopf, T. A. *et al.* The EVcouplings Python framework for coevolutionary sequence analysis. *Bioin-*
242       *formatics* **35**, 1582–1584 (2019).
243  138. Suzek, B. E., Wang, Y., Huang, H., McGarvey, P. B. & Wu, C. H. UniRef clusters: a comprehensive
244       and scalable alternative for improving sequence similarity searches. *Bioinformatics* **31**, 926–932
245       (2015).
246  139. Liu, X., Jian, X. & Boerwinkle, E. dbNSFP: A lightweight database of human nonsynonymous SNPs
247       and their functional predictions. *Hum Mutat* **32**, 894–899 (2011).
248  140. Liu, X., Li, C., Mou, C., Dong, Y. & Tu, Y. dbNSFP v4: a comprehensive database of transcript-spe-
249       cific functional predictions and annotations for human nonsynonymous and splice-site SNVs. *Ge-*
250       *nome Med* **12**, 1–8 (2020).
251  141. Xin, J. *et al.* High-performance web services for querying gene and variant annotation. *Genome Biol*
252       https://doi.org/10.1186/s13059-016-0953-9 (2016) doi:10.1186/s13059-016-0953-9.
253  142. Lelong, S. *et al.* BioThings SDK: a toolkit for building high-performance data APIs in biomedical
254       research. *Bioinformatics* **38**, 2077 (2022).
255  143. Hopkins, J. J., Wakeling, M. N., Johnson, M. B., Flanagan, S. E. & Laver, T. W. REVEL Is Better at
256       Predicting Pathogenicity of Loss-of-Function than Gain-of-Function Variants. *Hum Mutat* **2023**, 1–6
257       (2023).
258  144. Cheng, J. *et al.* Predictions for AlphaMissense. https://doi.org/10.5281/ZENODO.8360242 (2023)
259       doi:10.5281/ZENODO.8360242.

260

# Author Contributions (CRediT Classification)

# Acknowledgements