



HAL
open science

Getting to know each other: PPIMem, a novel approach for predicting transmembrane protein-protein complexes

Georges Khazen, Aram Gyulkhandanian, Tina Issa, Rachid Maroun

► To cite this version:

Georges Khazen, Aram Gyulkhandanian, Tina Issa, Rachid Maroun. Getting to know each other: PPIMem, a novel approach for predicting transmembrane protein-protein complexes. Computational and Structural Biotechnology Journal, 2021, 19, pp.5184-5197. <10.1016/j.csbj.2021.09.013>. <hal-03354041>

HAL Id: hal-03354041

<https://hal.science/hal-03354041v1>

Submitted on 24 Sep 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire HAL, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization



Getting to know each other: PPIMem, a novel approach for predicting transmembrane protein-protein complexes



Georges Khazen^{a,*}, Aram Gyulkhandanian^b, Tina Issa^a, Rachid C. Maroun^{b,*}

^a Computer Science and Mathematics Department, Lebanese American University, Byblos, Lebanon

^b Inserm U1204/Université d'Evry/Université Paris-Saclay, Structure-Activité des Biomolécules Normales et Pathologiques, 91025 Evry, France

ARTICLE INFO

Article history:

Received 7 June 2021

Received in revised form 23 August 2021

Accepted 12 September 2021

Available online 17 September 2021

Keywords:

Protein-protein interactions

Bioinformatics

Integral-to-the-membrane α -helical proteins

Molecular recognition

Data mining

Transmembrane protein complexes

ABSTRACT

Because of their considerable number and diversity, membrane proteins and their macromolecular complexes represent the functional units of cells. Their quaternary structure may be stabilized by interactions between the α -helices of different proteins in the hydrophobic region of the cell membrane. Membrane proteins equally represent potential pharmacological targets par excellence for various diseases. Unfortunately, their experimental 3D structure and that of their complexes with other intramembrane protein partners are scarce due to technical difficulties. To overcome this key problem, we devised PPIMem, a computational approach for the specific prediction of higher-order structures of α -helical transmembrane proteins. The novel approach involves proper identification of the amino acid residues at the interface of molecular complexes with a 3D structure. The identified residues compose then non-linear interaction motifs that are conveniently expressed as mathematical regular expressions. These are efficiently implemented for motif search in amino acid sequence databases, and for the accurate prediction of intramembrane protein-protein complexes. Our template interface-based approach predicted 21,544 binary complexes between 1,504 eukaryotic plasma membrane proteins across 39 species. We compare our predictions to experimental datasets of protein-protein interactions as a first validation method. The online database that results from the PPIMem algorithm with the annotated predicted interactions are implemented as a web server and can be accessed directly at <https://transint.univ-evry.fr>.

© 2021 The Author(s). Published by Elsevier B.V. on behalf of Research Network of Computational and Structural Biotechnology. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

1. Introduction

Proteins represent the core of cell machinery in organisms, and membrane proteins (MPs) encompass a broad variety of functions, including but not limited to activation, transport, degradation, stabilization, apoptosis, cell signaling and participation in the production of other proteins. More precisely, proteins exert their effects through mutual interactions in networks that represent functional units of cells. It is thus fundamental to understand the interactions between their components. More specifically, knowledge of the 3D structure of the MPs and interfaces involved in macromolecular complex formation remains a fundamental phenomenon governing all processes of life [1]. Therefore, MPs represent ultimate potential pharmacological targets in human health and disease because they include many families involved in protein-protein interaction (PPI) networks, leading to different physiological processes. Current estimates suggest that more than half of all drugs

target MP accessible on the surface of cells, whereas MPs make up less than one third of the human proteome and less than 1% of all solved protein crystal structures [2]. If we define the term TM as a transmembrane α -helical protein, transmembrane intermolecular α -helix – α -helix interactions through the lipid-embedded domains lead to oligomer formation and guide the assembly and function of many cell proteins and receptors [3]. In addition, the assembly of MPs may lead to emergent properties, as a relationship exists between oligomerization and function [4,5]. On the other hand, it is important to understand the functional effects of mutations at the interface because these may modify the assembly of MPs to form multiprotein complexes, leading to disruptions in networks of interactions and phenotypic changes. These mutations may possibly lead to the development of various genetic diseases [6]. Thus, the role of TM domains in protein function is crucial. These TM proteins span the cell membrane in its entirety and represent approximately one-third of the proteomes of organisms [7]. In eukaryotes, they come mostly in the form of α -helical bundles that cross different types of cell membranes and are approximately 3×10^5 in number, excluding polymorphisms or rare mutations [8]. The class of α -helical TM proteins has a higher diversity than

* Corresponding authors.

E-mail addresses: gkhazen@lau.edu.lb (G. Khazen), charbel.maroun@inserm.fr (R.C. Maroun).

their β -barrel counterparts. Thus, the number of possible binary and multiple interactions between them is vastly larger [8]. Estimates of the total number of human PPIs range from 130,000 to 600,000 [9–11], to 3,000,000 [12], which is several orders of magnitude larger than that of the *D. melanogaster* interactome.

High-throughput experimental and theoretical approaches are being used to build protein–protein interaction (PPI) networks. The data covering PPIs are increasing exponentially. Indeed, in the year 2012 more than 235,000 binary interactions were reported [13]. Most protein interaction databases (DBs; non-exhaustive list) [12–32] offer general information about experimentally validated PPIs of all types. The IMEx Consortium and the PSICQUIC service (<http://www.ebi.ac.uk>) group all data dealing with non-redundant protein interactions at one interface [33]. Nevertheless, the DBs are mostly concerned with water-soluble globular proteins and the juxtamembrane interactions of MPs. However, unlike globular proteins, MPs are water-insoluble, their exterior is much more hydrophobic than the interior because of allosteric interactions with the lipid environment, and they lose their native structure when removed from their natural membrane environment. Their thorough investigation has thus lagged due to this technical difficulty [34]. This is reflected in the fact that α -helical transmembrane proteins, as listed in the last version of PDBTM [35] represent only 5511 out of a total of 156,208 protein-only structures in the PDB, that is \sim 3.6%. Consequently, the molecular complexes formed by these proteins is represented by an even lower ratio, of about only 1.4% (619/44,700) protein–protein non-covalent dimer complexes as defined by the PDBePISA v1.52 server (https://www.ebi.ac.uk/msd-srv/prot_int/pistart.html).

Proteome-wide maps of the human interactome network have been generated in the past [25,29,36]. Traditional experimental techniques such as yeast two-hybrid (Y2H) assays [37] are not well suited for identifying MP interactions. Other assays, such as Y2H, are depleted of MPs or unreliable [38]. A new biochemical technique has been developed (MYTH); however, only a limited number of membrane-embedded complexes have been hitherto determined employing it [29,36]. This procedure has been significantly extended (MaMTH) [39]. However, to the best of our knowledge, MaMTH has not been used as a systematic screening assay to map the human MP interactome. Another approach for the identification of integral membrane PPIs in yeast used integral membrane proteins as baits [40]. An alternative method advanced a novel MYTH yeast proteomics technology, allowing the characterization of interaction partners of full-length GPCRs in a drug-dependent manner [41].

On the other hand, a variety of recent methods for the specific prediction of PPIs have flourished based on: i) machine learning and classifiers based on sequence alone [42,43], and deep learning [44]; ii) template structures [45,46]. Again, most of the approaches are parameterized on soluble globular proteins only, as in [47–49]. *Ab-initio* prediction of MP interfaces is rendered difficult, as these include amino acid compositions that are not radically different from the rest of the protein surface in terms of average hydrophobicity; nevertheless, they are better conserved than the rest of the surface [50]. This conservation is embedded in the invariable PPI-Mem nonlinear interface contacting binding motifs that we used.

To circumvent this problem, we developed a 3D structure knowledge-based approach to reliably predict complexes between TM proteins. Indeed, we incorporate the 3D structure as it provides [supplementary information](#) (surface accessibility, residue neighbors, etc.) not readily present in the amino acid sequence alone. The innovative approach is based on the detection of TM non-bonded contact residues between separate chains in experimental 3D structures of MP multimers reported in the PDB [51], and validated by the OPM DB [52]. Querying the PDBsum DB [53] there-

after, we obtained the atomic details of the membrane protein–membrane protein interfaces, that is, the contacts at the intermolecular interface of the complexes, through the PDBsum DB. We then gathered those amino acids at the recognition interface to generate regular expressions or patterns that represented in a linear form the interaction motifs in space. The regular expressions include wildcard regions between the interface contact residues representing the residues not in contact, even though they may be exposed to the membrane lipids. With this information in hand, we proceeded to search in the UniProtKB [54] of protein sequences the obtained motifs in other MPs. To extend our search, we allowed certain degrees of mutation in the membrane-embedded interface contact residues. We reasoned that (i) the interface residues in the experimentally determined structures of complexes between α -helical transmembrane proteins are responsible for the interaction; and (ii) homologs of precisely the template interface motifs are expected to interact analogously [55]. This latter assumption is opposed to a global approach that limits the search to functionally related partners or to overall sequence homolog partners without paying attention to the specific sequence at the interface. In all cases, it is reasonable to assume that the number of interface motifs is limited in nature [56]. Thus, we do not focus on the overall sequence homology, such as in other template-based predictions [45,48]. The linear 1D motifs we obtain are useful for quick local sequence searches, represent 3D epitopes implicitly and have the advantage of encompassing implicitly experimental tertiary and quaternary structures, as opposed to other approaches using only the primary structure.

In this work, we focus on the plasma membrane proteome of integral-to-membrane α -helical transmembrane proteins, ensuring that we probe direct interactions between proteins within the same subcellular localization. In addition, the TMs that compose a complex belong always to the same organism.

2. Materials and methods

2.1. Algorithm

2.1.1. Data mining and filtering

Fig. 1 succinctly summarizes the steps followed for collecting, filtering, and processing the input data from several sources and generating the output data based on regular expressions. We started by obtaining a list of all eukaryote “reviewed” proteins from UniProtKB, a manually annotated and reviewed DB of proteins with experimental evidence of their existence [54]. We mined proteins with cellular component annotations that matched the following GO annotations [57]: “integral component of membrane” (GO:0016021) and “plasma membrane” (GO:0005886); or “integral component of plasma membrane” (GO:0005887). Thus, we considered only proteins characterized as membrane proteins. From the resulting list, we identified the subset of proteins with an experimental 3D structure in the PDB, a resource containing the 3D shapes of biological macromolecules in the form of a complex spanning the TM region. We retained those with at least six buried interacting residues in each monomer (eight buried residues is typical of biological interfaces, [58]). This amount is the minimum for presenting an accessible surface area leading to quaternary structure. To obtain a high-quality homogeneous dataset, we adopted stringent selection conditions. Thus, for the PDB structure to be considered valid, it had to have a high-resolution limit of 3.5 Å or better if it was obtained by X-ray crystallography or cryoEM. We took in all different conformational states, regardless of the presence or absence of ligand(s), pH, symmetry group, apo or holo form, and allosteric effects, unless the differences modified the set of interface residues. We eliminated PDB MPs presenting engineered mutations, insertions, and deletions in the TM segment with respect to

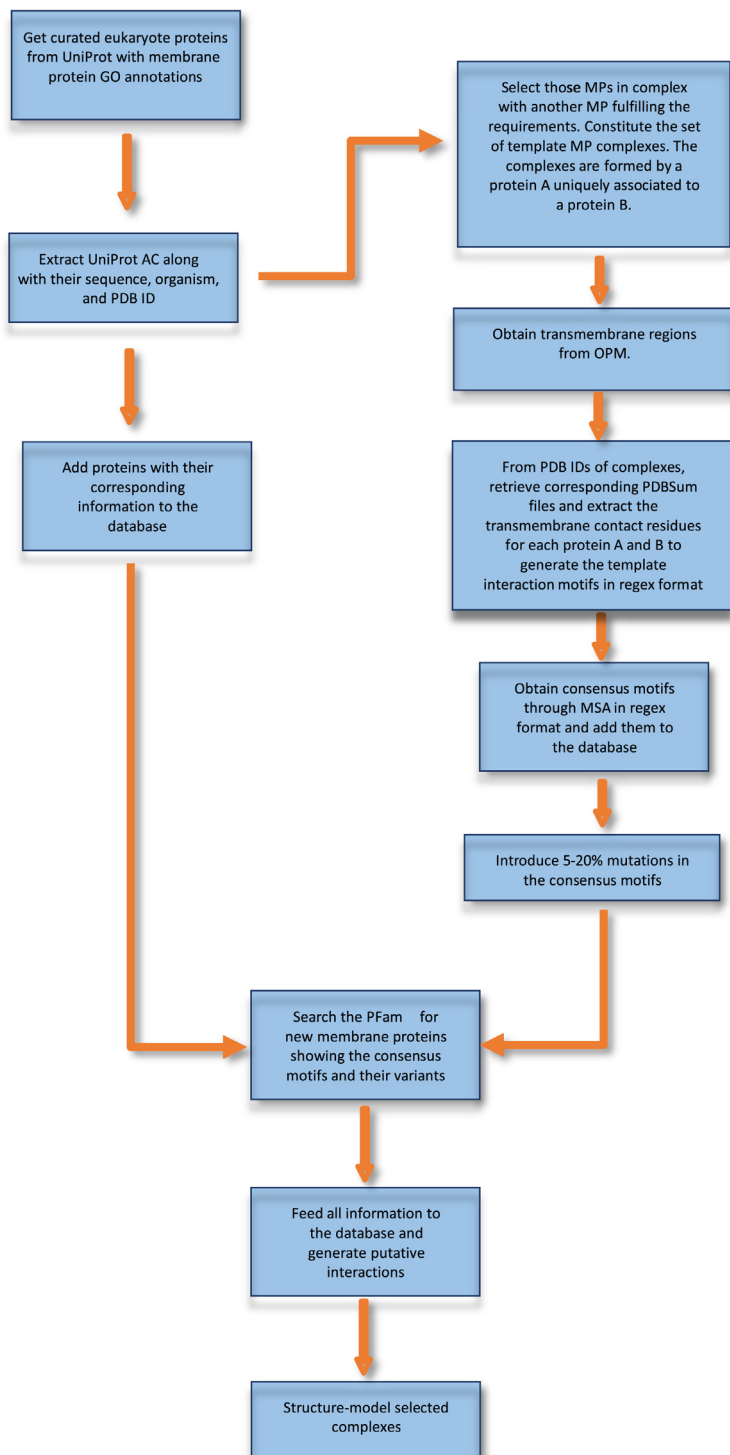


Fig. 1. Flowchart illustrating the PPIMem algorithm: from information retrieval to detection of recognition motifs to generation of putative interactions, to 3D modeling of complexes.

the wild-type or natural variant sequence in UniProtKB, just like chimeras in which the xenophobic part is TM. We also ignored redundant MPs, those whose 3D structures showed no significant TM segments, and pair interactions that are redundant due to symmetry. Manual curation excluded non-parallel configurations of the protomers, head-to-head or head-to-tail orientations (resulting most probably from crystal interfaces), out-of-membrane interactions only, and TM segments that do not interact, as dictated by the cell membrane. Intramolecular interactions within each

protomer of an oligomer, as well as with the lipids and detergents environment are included implicitly. The PDBsum DB contains the PDB molecules that make up a complex and the interactions between them. This information delivers what interface residues of protomer A interact with what interface residues of protomer or chain B, and what protomer A interacts precisely with what protomer B and not another.

Finally, to ensure that the oligomer structures we considered are quaternary structures with biological sense, we used EPPIC, a

protein–protein interface classifier [58–60], and PRODIGY, a classifier of biological interfaces in protein complexes [61,62] to distinguish between crystallographic and biological assemblies.

2.1.2. Motif extraction

To carefully choose the PDB structures to work on, we referred to the OPM DB [63], which provides the orientation of known spatial arrangements of unique structures of representative MPs coming from the PDB with respect to the hydrocarbon core of the lipid bilayer. We chose all the PDB structures that mapped to the UniProtKB extracted MPs and pulled out all available PDBsum files of these structures. We double-checked the chosen PDB structures with the MPStruc DB of MPs of known 3D structure (<https://blanco.biomol.uci.edu/mpstruc/>). PDBsum contains atomic non-bonded contacts between amino acid residues at the interface of molecules in a multimer complex. We used the information in PDBsum to extract the intermolecular contact residues. We filled in with the non-contact residues in the sequence, that is, those between the contact residues, as wildcards. Thereafter, we defined the nonlinear binding motifs by obtaining the corresponding linear sequences (motif instance). From the PDBsum file listing the contacts, we formulated two motifs, one corresponding to partner protein A and the other one to its partner protein B. Because we are only interested in the recognition site at the TM interface region, we ensured that each interacting residue belonged to the TM part of the sequence. We represented our motifs using the Regex format (<https://www.regular-expressions.info/>) and denoted the TM contact residues by their one-letter symbol, the residues in between (representing the wildcard regions) by a dot, and the curly braces for the consecutive occurrences of the wildcard. As an example, we take UniProtKB membrane protein P04626 (receptor tyrosine-protein kinase ErbB2). The code of the spatial structure of the ErbB2 dimeric transmembrane domain in the PDB is 2N2A. Its PDBsum entry contains non-bonded contacts across the surface for each of the chains A and B. The residues in between constitute the non-contact residues wildcard and are represented by curly braces. PDBsum lists chain A residues I12, V16, L20, V23, L24, V27, and L31 as establishing contacts with chain B. The result of the contact motif for chain A is $^{12}IX_3VX_3LX_2VLX_2VX_3L^{31}$ (Fig. 2), equivalent in the regular expression notation to $I.\{3\}V.\{3\}L.\{2\}VL.\{2\}V.\{3\}L$. As this is a homodimer and the motif is the only one, the contact motif is the same for chain B. The wild cards $\{\}$ represent extramembrane loops or inner protein subdomains of different lengths that do not establish membrane interface contacts. We do not focus on these variable segments although of course, those loops may be involved in interactions in the cytosol or in the extracellular milieu; consequently, we do not look in our data mining at whether those segments are conserved in the amino acid sequence. We focus only on the intramembrane interface contact residues. If the resulting motif is an intramembrane binding motif, the types of variable residues or the length of the interdomains should not be essential.

Finally, to verify the consistency of our efficient algorithm and the validity of the obtained binding motifs (that they will not match with many sequences just by random), we sent several BLAST (<https://blast.ncbi.nlm.nih.gov>) searches against the PDB of several motifs or their parts thereof (in case they were separated by long stretches of non-contact amino acid residues between the contact residues). We always recovered the corresponding template PDB files, as well as related PDB codes, among the maximum BLAST scores and with query covers of 100 (not shown).

2.1.3. Searching for identified motifs in other protein sequences

To eliminate redundancy in our data, we grouped motifs derived from mutating the surface residues in the native motifs to build a consensus motif for each cluster derived from multiple

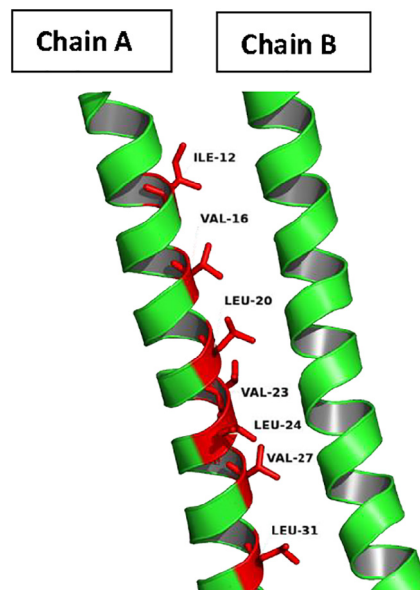


Fig. 2. Contact nonlinear interface residues of chain A of the ErbB2 dimeric transmembrane domain (UniProtKB P04626) of sequence $^{12}\text{Ile-Ser-Ala-Val-Val-Gly-Ile-Leu-Leu-Val-Val-Val-Leu-Gly-Val-Val-Phe-Gly-Ile-Leu}^{31}$. Contact residues are in bold and compose the motif $I.\{3\}V.\{3\}L.\{2\}VL.\{2\}V.\{3\}L$.

sequence alignments (MSA). This included all recovered sequences. This allowed us to retain just conserved interface residues, preventing therefore potentially deleterious effects on protein assembly. The mutation rates of the contact residues ranged from 0% (exact match) to 20%, with increments of 5%. In this way, we wished to amplify our search, allowing us to find other protein sequences with homologous motifs. The mutation procedure was applied independently to each protomer, on each motif. Most of the time, the two motifs from a homodimer are identical. However, in some cases, one of the monomers of a homodimer may reveal another interacting surface. In this case, the corresponding motif is different.

In this manner, we queried our consensus motifs against the Pfam dataset (<https://pfam.xfam.org/>). We defined a COST parameter as the number of amino acid mutations allowed anywhere in the motif, depending on the number of contact residues it contains and the mutation rate for the run. We assigned a score of 0 to both insertions and deletions to ensure that no contact residue is lost. For instance, when generating new sites from a valid motif with eight contact residues, $\text{COST} = 2$ for a mutation rate of 20% ($8 \times 0.20 = 1.6$ rounded up to 2). The values of COST_A vary from 0 to 4 and those of COST_B from 0 to 4, 7–8, reflecting the fact that there are TMs whose sequences contain more than 30 contact residues in their interaction motifs and accommodating up to 20% mutations with respect to the consensus motifs. An example of this is UniProtKB AC Q920B6 (Potassium channel subfamily K member 2) from *R. norvegicus*. In which case, the contact residues are shared by motifs in different TM α -helices of the same protein.

To intentionally keep track of which A motif interacts with which motif B, we kept the motifs in separate pools. In other words, the rationale for keeping the two sets A and B of motifs distinct is that a given motif A matches a given motif B and not just any other one. Subsequently, we paired the predicted motifs from novel interactions based on the PDBsum validated interactions. In this way, we were certain that pattern A from the new protein A binds to its complementary pattern B of new protein B. Because motifs can be found fortuitously anywhere in the sequence, we considered only those motifs belonging to the TM region. We also checked which TMs showed the motifs in their PDB structures, if available.

2.1.4. Implementation

2.1.4.1. The PPIMem database. Thereafter, we naturally built a heterogeneous relational DB named PPIMem, which contains all the found interactions. The DB is the result of the effective implementation of the automated template-based recognition site search pipeline. We used a MySQL DB to maintain all the information collected in a structured manner. To access the DB and search for our data, we built a web interface using PHP and HTML5, which allows the user to query the needed information. Users can query the DB for obtaining motifs by entering a UniProtKB AC, a gene name, a type of organism, a mutation percentage, or a motif of interest or part thereof using the Regex format. A link to the UniProtKB site for each UniProtKB AC will be available shortly. The user can choose more than one filter option when querying and will exclusively obtain interactions thought to occur in TM regions between plasma membrane proteins of the same species. The user can also adjust the values of the following parameters: Number of contact residues, Mutation rate, COST and Valid, as an interval or a fixed value. Homodimers of the TMs dealt with do not appear in the database, as all PPIMem TM proteins are expected to form them. We will update our DB following each release of the UniProtKB and OPM datasets and then regenerate all statistics.

2.1.5. Molecular docking

he predicted interface residues can be intentionally used as constraints to reconstruct the structure of dimers through docking. We consider successful docking as a partial validation procedure. To illustrate our approach, we generated 3D binary complexes derived from selected predicted pairs of TMs with the goal of looking at their architecture. To do so, we searched for a protein–protein docking program that would allow us to perform a steered docking simulation using the epitopes extracted from the molecular interface of the complex. We processed and analyzed a considerable number of experimental 3D TM structure files using several docking programs. As we were aware of the docking interfaces which we used for the restraint-driven docking, we did not perform in general *ab initio* calculations. Neither were we concerned whether the docking program was trained on sets composed primarily of cytosolic proteins. Even though several tested docking programs were sufficiently precise, we decided to use GRAMM-X [64,65] for creating the novel protein–protein 3D complexes. GRAMM-X has an option in which the user can submit those residues that might form the interface between the “receptor” and the “ligand.” The program also has options for listing the residues required to be in contact. To verify the performance of GRAMM-X for TMs, we benchmarked it against several TM complexes in the PDB. GRAMM-X was indeed able to reproduce many of the experimental TM complexes. For molecular docking simulation and identification of the PPIMem-predicted complexes, we chose examples in which the 3D PDB structures of proteins were already available or represented highly homologous templates. After interface-driven docking, we manually curated the output by filtering out non-parallel, perpendicular, or oblique protomer pairs, regardless of the calculated energy. We also considered the topology of the TM in the membrane. We visualized the obtained 3D structures of the complexes using PyMol 2.4 (www.pymol.org).

3. Results

3.1. Automated multi-stage pipeline to predict and explore thousands of novel transmembrane protein–protein direct interactions (TMPPI)

UniProtKB provided us with 13,352 MPs that include the GO annotations mentioned in the S&M section. Overall, these proteins mapped to 954 distinct oligomer MP PDB structures. As we focused

on structures that satisfy the requirements signified in the S&M section, we were able to validate only ~50 PDB files of TM-TM complexes. After carefully checking which corresponding PDBsum files to consider, we ended up with 53 non-redundant template interactions, associated with 48 unique reviewed UniProtKB entries across species and forming our template set. Fifty of these interactions correspond to structural homomers and three to structural heteromers (Table S1, yellow highlighting). The set includes X-ray, solution NMR, and electron microscopy structures. The TMs include families like receptor tyrosine kinases, TLRs, ion channels, Cys-Loop and immune receptors, gap junctions, transporters, and GPCRs. In these, besides *H. sapiens*, a taxonomically diverse set of organisms across the eukaryotic kingdom of the tree of life is adequately represented: *A. thaliana*, *B. taurus*, *G. gallus*, *M. musculus*, *O. sativa*, *R. norvegicus*, *S. cerevisiae*, and several others. On another hand, there are 21 structures of complexes between bitopic proteins, 32 between polytopic proteins, and one mixed bitopic-polytopic protein complex in our template set. The oligomerization order is adequately represented, with 23 homo 2-mers, five homo 3-mers, 12 homo 4-mers, one hetero 4-mer, four homo 5-mers, three hetero 6-mers, one hetero 10-mer, and one homo 12-mer. The occurrences of the number of TM helices for each of the protomers of the reference complexes shows that bitopic single-span helix complexes are the most preponderant; these are followed by polytopic TMs composed of 6, 4, 7, 9, 2, 8, and 10 and 12 helices. The following EC numbers are represented in this set: EC 2 10 (transferases), EC 3 (one hydrolase), and EC 7 (one translocase). When verifying the protein–protein interfaces in the complexes for the type of assembly they form (crystallographic or biological), we found that all the X-ray or electron microscopy complexes were classified as biological (Table S2). We did not submit the NMR-determined complexes to the test. Therefore, this is the number of experimentally solved structures of TM protein complexes that we used as the reference template set.

The PPIMem pipeline results indicate there are 1504 unique TMs involved in the predicted complexes, going from UniProtKB accession codes A0PJK1 to Q9Y6W8. Of these, 417 are human (Table S3). The combined number of motifs found after removing redundancies due to different chains of the same structures interacting more than once was 98 (Table S4).

In many complexes, the motifs in each subunit are identical, i.e., motif A is coupled to an identical motif B. But in other instances, the two paired motifs are different (Table S5). In this situation, the complexes are typically comprised by different monomeric subunits, although if both subunits are the same protein, but contain more than one interacting surface, the complex is a homodimer in which the contact surfaces are not the same. Of course, even when motifs A and B are the same, the complex may not be necessarily a homodimer.

We observed that some amino acid residues were more favored than others in the TM recognition sites. For instance, the hydrophobic side chains Leu, Ile, Val, and Phe were the most abundant, followed by Ala and Gly. This residue distribution has been found in the past [66]. Leu was found more than 300 times, making about one fifth of all contact residues, like reported before [67] (Fig. S1). As expected, the physicochemical properties of TMPPI-binding sites are different from the exposed sites of soluble proteins. The amino acid residue abundances for helix interactions we found in the motifs match those in the literature [68,69]. Fig. S2 shows, as a “heat map,” couples of contact residues at the interface for our template set of complexes, as they come out from PDBsum. We can see that the largest value corresponds to the contact between two Leu residues, followed by contact between Phe residues, and then the Leu-Phe pair. The least observed interactions include His-His for a homotypic pair and Trp-Cys for a heterotypic pair. This outcome suggests that residues tend to contact

other residues sharing the same or similar physicochemical properties, and agrees with the statistics obtained for inter-residue interactions in the MP Bundles DB for α -helical MPs [68], as well as with the amino acid abundances in the middle region of the membrane [70]. The statistical trends in the contacts imply specificity at the interface of the TM helices, as well as correlated mutations of coupled contact residues and paired motifs.

The number of TM atomic non-bonded contacts per complex covers a broad range: 11–87 for single-span TM proteins, and 8–186 for multiple-span TM complexes (Fig. 3). As expected, the latter show in average a larger number of non-bonded contacts.

Finally, we looked at the Pfam-A set of protein families [71] to which the 63 different proteins in the template complexes of Table S1 belong. We found that the most populated families were PF07714 (Pfam family PK_Tyr_Ser-Thr) and PF00520 (Pfam family Ion_trans), with occurrences of 10/110 and 6/110, respectively. Forty-two proteins out of the 63 had only one Pfam occurrence each (not shown).

3.2. α -Helix packing motifs

Walters and DeGrado [72] have examined helix-packing motifs in membrane proteins involved in the folding of a helical membrane protein, i.e. interactions between helices within the same protein. As we are dealing with interhelical interactions between different proteins embedded in the membrane, helix-pairing motifs are not directly comparable. However, we can look at the pairing motifs of TM proteins composed of a single TM α -helical domain, such as glycophorin A, presenting a G. $\{3\}$ G packing submotif [73–75], in which Gly or other small residue space four residues that mediate parallel interhelix protein interactions. Indeed, we find this small motif and its variants as part of the extracted binding motifs from our set of template TM complexes (Table S1). The structures include parallel single α -helix homodimers (Table S1, PDB 2L2T and 2LOH), multiple α -helix homooligomers with a different motif in each component (Table S1; PDB 5AEX, chain A; PDB 5CTG, chain A). The G. $\{3\}$ G submotif appears as well in the non-redundant PPIMem motifs (Table S4, in bold orange).

An antiparallel coiled-coil submotif of α -helices with Ala or another small residue in every seventh position has been termed Alacoil [76]. It can be found in human aquaporin 5 (PDB 3D9S) as the PPIMem binding motif A. $\{2\}$ T. $\{2\}$ QA, and as A. $\{2\}$ VF.LA in the TRPV1 ion channel in complex with double-knot toxin and

resiniferatoxin (PDB 5IRX) of our template protein complex set (Table S1; Table S4 in bold green).

Another specific pattern is the Aromatic. $\{2\}$ Aromatic motif [77], like in F. $\{2\}$ F of connexin-26 (PDB 2ZW3), of the mouse TRPC4 ion channel (PDB 5Z96), and of TRPM4 (PDB 6BCO, 6BQV). In our reference set, this submotif appears as part of the TrkA transmembrane domain (PDB 2 N90). In the SWEET transporter (PDB 5CTG), it takes the form F. $\{2\}$ Y; in polycystic kidney disease protein 2 (PKD2) (PDB 5T4D) and polycystic kidney disease-like channel PKD2L1 (PDB 5Z1W), it becomes F. $\{2\}$ W (Table S1; Table S4 bold blue).

Another submotif that appears very frequently is the L. $\{6\}$ L heptad, analogous to the Leu zipper that mediates protein complex formation in water-soluble proteins [78]. As this pattern is fundamental, it is presented in Table S6 and Table S4 (bold gold).

PPIMem reveals precisely original packing motifs through its non-redundant consensus motifs, like the frequent motif A. $\{2\}$ (A/L/V)(L/F) (Table S4). The following motifs with Ile are also well populated: for position $i-1$, (F/L/W/W)I; for positions $i-2$ and $i-1$, (F/L)CI; for positions $i+3$ and following, I. $\{2\}$ (A/L/V), I. $\{3\}$ (L/V), I. $\{5\}$ I, and I. $\{6\}$ A. PPIMem also uncovers Ile composite motifs, like I. $\{2\}$ (A/F)I. $\{2\}$ (T/L), and I. $\{3\}$ I. $\{3\}$ L. The most frequent motifs are those beginning with Leu: L. $\{2\}$ (A/F/G/I/L/N/T/V)(A/G/F/I/L/V), and L. $\{3\}$ (F/V). Leu composite motifs include L. $\{2\}$ (F/I)L. $\{2\}$ G, L. $\{2\}$ IF. $\{2\}$ LL. $\{2\}$ F, L. $\{3\}$ L. $\{2\}$ FF, and L. $\{3\}$ L. $\{2\}$ INP. $\{2\}$ L. $\{3\}$ V. The occurrence of these patterns is naturally reflected in the interhelical contacts between proteins of the Ile-Ile, Leu-Leu and Ile-Leu residue pairs (Fig. S2). The composition of these residues shows, as if it was necessary, that hydrophobic residues are enriched at the interface between TMs.

A statistical analysis of amino acid patterns in TM helices [75,79] results in 30 over-represented ($p \lll$, odds ratio greater than 1), and 30 under-represented ($p \lll$, odds ratio < 1) pairs. Thus, G. $\{3\}$ G (GG4 in Senes et al. notation) is the most significant pair among over-represented pairs. Indeed, this and other significant over-represented pairs (21 out of 30), such as I. $\{3\}$ I, G. $\{3\}$ A, IG, I.G, V.G, I. $\{3\}$ V, IP, V. $\{3\}$ V, V. $\{3\}$ I, AV, and G. $\{2\}$ L, have their equivalents in non-redundant PPIMem motifs (Table S4). On the other hand, several under-represented pairs (10 out of 30) are in fact absent from PPIMem motifs in Table S4, such as I.I, F. $\{3\}$ I, and I. $\{3\}$ G. Moreover, several most significant over-represented triplets of Senes et al. appear in our motifs, lending support to our template-based approach. Lastly, it is worth mentioning that the presence of known sequence motifs alone does not guarantee interactions [80].

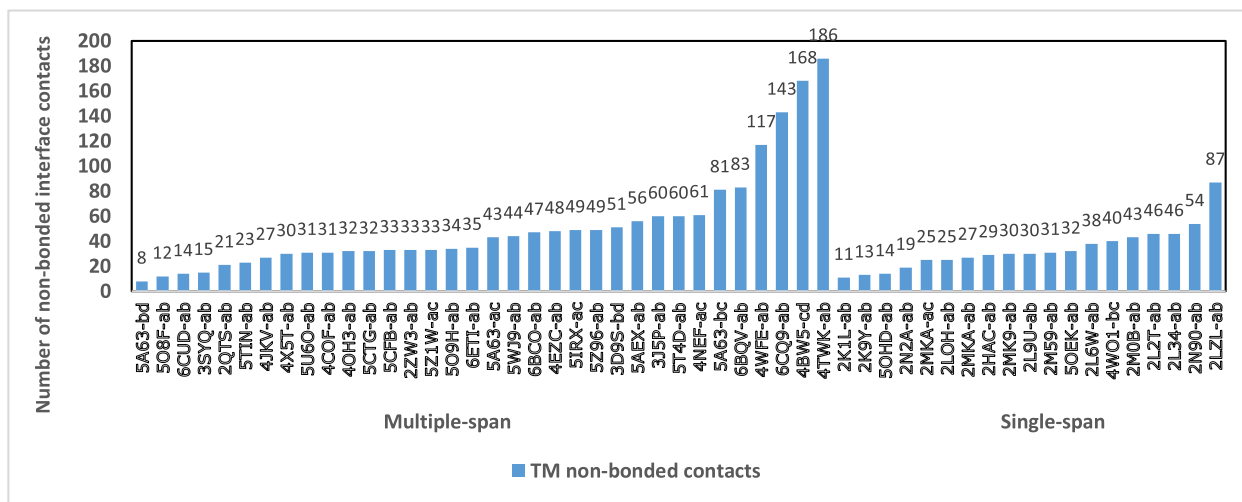


Fig. 3. Number of PDBsum non-bonded interface contacts for each of the experimental template complex PDB structures of Table S1.

We then wished to look at the number of UniProtKB ACs resulting for each number n of contact residues for the biological species recorded in the PDB and UniProtKB. As seen in Fig. 4, the count occurrence of contact residues in the motifs is largest for $n = 7$ for both protein A and protein B of the interacting pair. The quasi-periodicity of the interface recognition residues reflects a spatial arrangement corresponding to a heptad α -helical pattern. This is valid also for motifs separated by values of the linker residues ≥ 4 . For example, the long motif DWN.{2}IA.{2}V.{2}LI.{24}F.{3}F.{2}A.{2}YLV.{2}FM (Table S4) contains two motifs of 14 and 17 residues, separated by 24 linker residues. A motif may attain a total of 20 residues, like in E.{2}TL.VGF.IL.{2}SL.{2}TY (Table S4), possessing enough length to span a bilayer. Our motifs therefore may be rather long, implying residue correlations beyond the pair level. Moreover, a supplementary information is conveyed by the PPIMem motifs: the residues that compose them may be not only on the same face of the helix but are interfacial residues. The corresponding statistics for *H. sapiens* show similar trends: small motifs with six to seven buried contact residues are the most abundant (peak at $n = 7$, combined number of redundant proteins A and B, 9655). The range of contact residues in the motifs was 6–32. As the number of contact residues increased in the motifs, the number of hits decreased drastically, with only one to four predictions for both proteins in the 18–33 contact residue range (Fig. 4). These findings suggest a limited set of binding motifs in nature.

The protein–protein docking benchmark 5.0 PPDB [81] assembles non-redundant high-resolution complex structures, for which the X-ray or NMR unbound structures of the constituent proteins are also available (<https://zlab.umassmed.edu/benchmark/>). However, none of the complexes in PPDB v5.0 correspond to MPs. On our side, from our 53 non-redundant template structures of TM-TM complexes, we were able to extract a subset of them along with the unbound structures of their components to define a TM-TM complex benchmark made up of ten sets of structures (Table S7).

When comparing our benchmark to the “Dimeric complexes” table of the Membrane Protein Complex Docking Benchmark, MemCplxDB [82], we only recover the PDB 5A63 complex. The reason is that MemCplxDB shows many interactions between MP and non-MP, cytosolic proteins (antibodies, peripheral proteins, etc.), which we do not deal with. MemCplxDB includes interactions as well between oligomers within a multimer complex, and prokaryotic MPs of β -barrel structure. Our benchmark represents thus a standard set of true positives for integral-to-membrane proteins interacting through their α -helical TM segments. As the 3D

structures of more TM complexes appear, the benchmark will grow and could serve for a machine learning approach of prediction of membrane PPIs.

3.3. Predicted interactions

The number of PPIMem-predicted membrane heteromer interactions for the 39 species dealt with is 21,544 among 1504 TMs. The homodimers are hence 1504 in number and represent 6.5% of all complexes. Of the total heteromer interactions, 9797 among 417 TMs correspond to *H. sapiens*, the homodimers representing 4.1%.

PPIMem predicts interactions thus for 417 human genes, including many disease genes. The Mendelian Inheritance in Man (MIM) number contains valuable information on the known mendelian disorders caused by variants affecting the gene represented in the entry and focuses on phenotype-genotype relationships (<https://www.omim.org/>). Table S3 shows that 101 out of 417 PPIMem TMs are involved in human disease according to MIM. This set provides 83 TMs in PPIMem extracted from complexes with nil mutations (mutation rate = 0%) in their binding motifs. We looked at their missense variants in the index of human variants curated from literature reports in UniProtKB (<https://www.uniprot.org/docs/humsavar>), and focused on the following categories as defined by the ACMG/AMP terminology [83]: likely pathogenic, pathogenic or of uncertain significance. Table S3 shows that several mutants involved PPIMem interface contact residues belonging to TM α -helices, suggesting a destabilization of the complex as the molecular basis of the disease. The genes of the three membrane proteins are the α -1 subunit of the glycine receptor (GLRA1; P23415), the β -3 subunit of the GABA receptor (GABRB3; P28472) and the β -2 subunit of the GABA receptor (GABRB2; P47870), proteins of considerable biological interest. PPIMem predicts complexes in which one or two of the subunits are associated to known mendelian disorders (Table S8).

At different mutation rates, we defined and identified many potential recognition sites and novel binary complexes. As mentioned in the S&M section, we built a consensus amino acid contact motif for all the matched sequences of a given binding site in the 0–20% mutation range considering the contact residues only. This led us to detect conserved amino acid residues among the contact residues of the nonlinear binding motifs. The most prevalent consensus motif A found was AV.{2}GL.{2}GA.{2}L, illustrated by the instance sequence ⁴²³AVFSGLICVAMYL⁴³⁶ of the

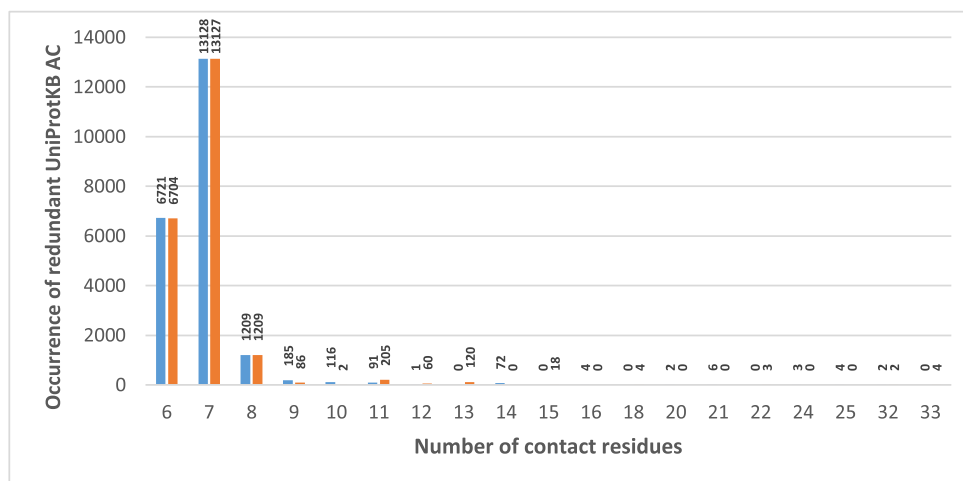


Fig. 4. Occurrence of (redundant) UniProtKB ACs as a function of interface contact residues for Proteins A (blue) and corresponding Proteins B (orange). Mutation rate 0–20%. Organisms extracted from the PDB and UniProtKB. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

sodium-dependent proline transporter (UniProtKB Q99884) at a 15% mutation allowance (Fig. 5a). For consensus motif B, the same motif was the most frequent (Fig. 5b). The least frequent motifs are characterized by hefty linker segments between the contact residues and thus by substantial sequence lengths. An example is motif B: LL.AS.{90}LV.EA.FAI.NI.{2}S.{2}L.{3}F.{19}I.{100}I, found in the short transient receptor potential channel 4 (UniProtKB Q9UBN4). In general, the interface contact residues are part of conserved sequences of TM regions.

We cannot know at this stage whether or not the amino acids composing our binding motifs represent hot-spots, i.e., those that contribute to the binding free energy [84].

Mutation rates of 0% for motif A and motif B result in those proteins whose contact residue sequences conform exactly to the consensus motifs. The PPIMem outcome provides 79 entries of heterooligomer complexes across species, of which 25 are *H. sapiens*. Retaining motif A only, PPIMem predicts 346 interactions for *H. sapiens* at a mutation rate of 0%. For the mutations rates 5, 10, 15, and 20%, the number of interactions is 16, 4439, 4272, and 724, respectively.

As mentioned, we derived a consensus sequence of the binding motif from multiple alignments of the contact residues of all sequences found for a given motif. Thus, for example, the consensus sequence resulting from up to a 20% mutation rate of the VV.{2}A.{2}A.{2}VL.{2}I.{3}I motif of length 19 is shown in Fig. S3, leading to VVX₂AX₂AX₂VL as the fingerprint of the nonlinear binding motif.

A striking example of how PPIMem focuses exclusively on the sequence homology at the level of the binding motif, as opposed to other template-based predictions based on overall sequence homology, is illustrated by the pair of TM proteins T-cell surface glycoprotein CD3 ζ chain (P20963) and the high affinity immunoglobulin epsilon receptor subunit γ (P30273). These proteins are predicted to form a complex by our algorithm. The corresponding sequences show an overall sequence identity of 16%, implying these proteins are not related by sequence, even though possessing the same the consensus binding motif C.{2}LD.{2}L.{2}YG.{2}LT.LF, whereas the sequence identity at the level of the consensus motif is 90%. Thus, our method sampling the binding interface is more robust and specific as it recovers unrelated proteins.

Some of the proteins, like the ligand-gated chloride channel human glycine receptor subunit α -3 (O75311), show several distinctly different interface motifs, suggesting a promiscuous binding behavior: A.{3}V.{3}I.{3}L.{6}S.{2}R.{19}L.{3}F.{2}L, Y.{2}I.{7}L.{2}I L.{16}GL.{2}T.{2}LT.{2}T.{2}SG.R, Q.{6}LI.IL.{5}WI.{6}A.{2}AL.{2}T, and I.{3}L.{2}T.{6}R.{12}D.{2}MA.{6}F.{2}LL. It may happen thus that PPIMem presents a pair of TMs forming the same complex more than once. This is because one or the two proteins of the putative complex might present more than one binding site. The O75311 – P18505 (γ -aminobutyric acid receptor subunit β -1) couple illustrates this situation, in which each of the four predicted complexes presents different binding sites. This finding agrees with the multiplicity of protein binding modes and their multiple functionalities.

3.4. PPIMem and other datasets

The PPI template-based prediction algorithm and server PRISM2.0 [85] also uses the 3D structural similarity of interfaces with respect to templates to propose other protein complexes. PRISM is not specific for MPs and requires the 3D structure of the targets to propose an interaction, whereas PPIMem is specific for TMs and does not require the 3D structure of the subunits composing the putative complex. Thus, when having an interface template corresponding to an MP, PRISM may propose not only MP protein complexes but also globular protein complexes. Therefore,

many of our TM template interfaces are absent in the PRISM dataset.

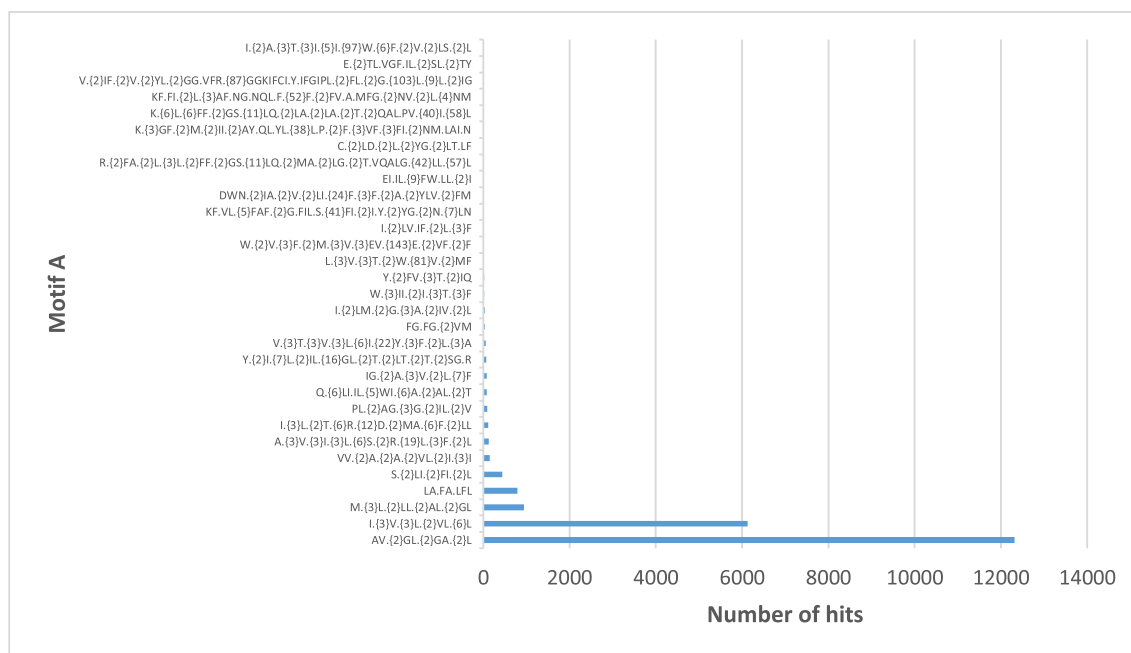
Although most datasets based on experimental approaches cover the entire human proteome, again, MPs are under-represented (for an overview of the major high-throughput experimental methods used to detect PPIs, see [86,87]). Thus, the experimental work of Rolland and colleagues on the human interactome network [25] found only 41 interactions between MPs. Twenty-eight of these proteins were found to interact in the IntAct DB. Nevertheless, none of the interactions we extracted from the structural PDBsum DB were found among the 41 interactions above. It seems that some of these interactions bear between the juxtamembrane regions of the MPs reported by Rolland et al. for MPs. We did not find either any of our predictions in their results. In the HI-III-20/HuRI updated dataset [12], none of the high-quality binary PPIs are MPs (log₂ odds ratio < 0, Extended Data Figure 7a of Luck et al.)!

To verify our predictions of TM-TM complexes, we confronted our PPIMem predictions at 0–20% mutation rate for human proteins to several datasets like FpClass of predicted human PPIs, BioPlex 3.0 of experimentally determined PPIs (uses the TAP-MS technology), MENTHA, IID, HMRI, and BIPS. In Table S9 we compare our PPIMem predictions to FpClass and to BioPlex with the aim of representing an independent qualitative validation of our method. For example, PPIMem predicts the interaction between isoform 2 of the human γ -aminobutyric acid receptor subunit β -3 (GABRB3, P28472-2) and subunit α -6 (GABRA6, Q16445). In the BioPlex 3.0 dataset, the interaction is detected with a probability ≥ 0.99 . The ComplexPortal DB (<https://www.ebi.ac.uk/complexportal/home>) reports the same interaction as Complex AC: CPX-2164 and CPX-2951. Given that BioPlex reports ~ 2000 MP interactions (all membranes included) with score ≥ 0.78 , the present non-specific correspondence between the two datasets is of $\sim 2\%$, considering the 35 interactions described in Table S9. Searching the PPIMem human proteins in FpClass results in 8857 predicted PPIs with a score ≥ 0.5 , out of 26,456 interactions for 278 PPIMem proteins. In the end, we obtained 74 interactions in common between the two datasets, representing a correspondence of 0.8%. As a reminder, the BioPlex and PPIMem sets are independent, i.e., there is no intended intersection between them, and the BioPlex dataset is for the HEK293T cell, whereas our dataset is for no specific cell at all. In addition, the FpClass and PPIMem sets are orthogonal and 139 PPIMem human proteins were not found in FpClass.

Despite the difficulties of comparing different datasets, we show how some of our predicted data for *H. sapiens* is found in other datasets:

- The MENTHA experimentally determined direct protein interactions DB presents also the P28472-Q16445, as well as the P28472-P47870 (GABRB2) interaction.
- The IID DB validates experimentally the Vascular endothelial growth factor receptor 2 (P35968) – receptor 3 (P35916) interaction.
- The HMRI DB, which seems to list only heteromers and for which not all the interactions are between MPs, shows a correspondence for the heterotypic pair TYROBP- KLRC2 (p value = 0,034341).
- BIPS [30] predicts putative interactions and is based on sequence homology between proteins found in PPI DBs and templates. We find several correlations between BIPS and PPIMem. For instance, we propose an interaction between T-cell surface glycoprotein CD3 ζ chain (P20963) and high immunity immunoglobulin ϵ receptor subunit γ (P30273). BIPS predicts a similar pair between T-cell surface glycoprotein CD3 ζ chain (P20963) and low-affinity immunoglobulin γ Fc region receptor III-A (P08637).

a)



b)

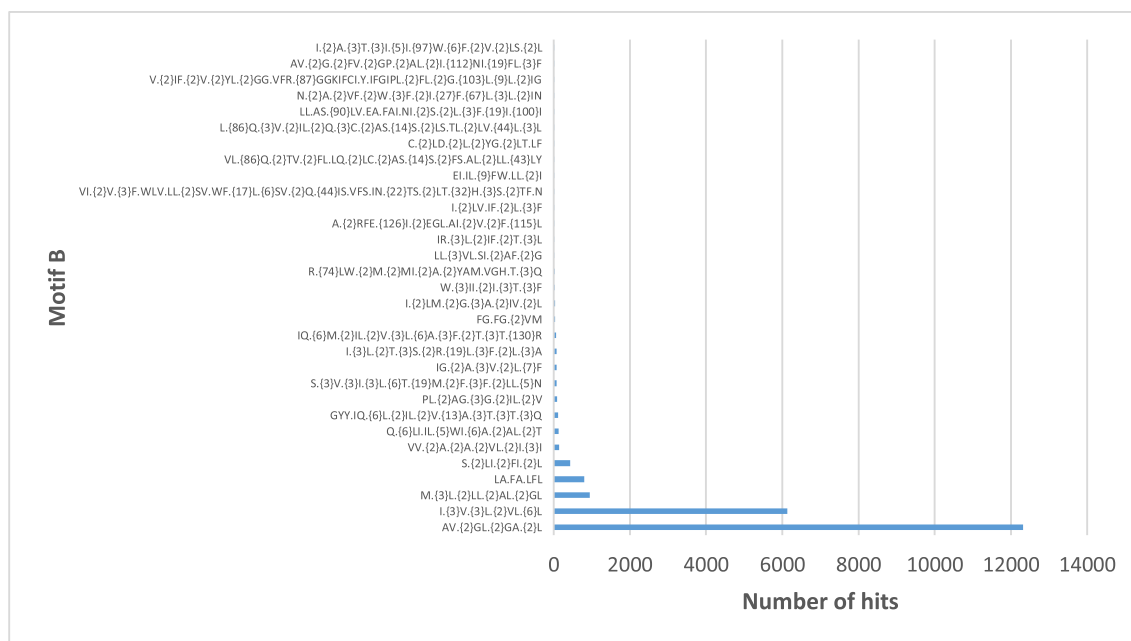


Fig. 5. Number of hits (i.e., number of times the motif appears) per motif for all species. a) motif A, b) motif B. Mutation rate 0–20%.

Finally, as the template set (Table S1) is extremely small to serve as a training set, we could not evaluate quantitatively interface predictions using standard criteria (ROC plot, precision, recall, PCC, etc.). In addition, we do not count with a “negative” set, i.e., a set reporting α -transmembrane proteins shown not to interact.

3.5. Interaction networks

Networks link overlapping pairs of proteins, from which it is possible to propose multimer complexes if the binding sites

are independent and non-overlapping. The architecture of a network reflects biological processes and molecular functions. Thus, from the predictions, *de novo* connections can be found, linking the network to a disease pathway, and proposing innovative possible cellular roles for some of the complexes. We illustrate below a subnetwork of PPIMem-predicted TMPPIs for *H. sapiens*:

fx1

Q9NY15 (*STAB1*) Stabilin 1; O14494 (*PLPP1*) Phospholipid phosphatase 1; Q15758 (*SLC1A5*) Solute carrier family 1 member 5 or

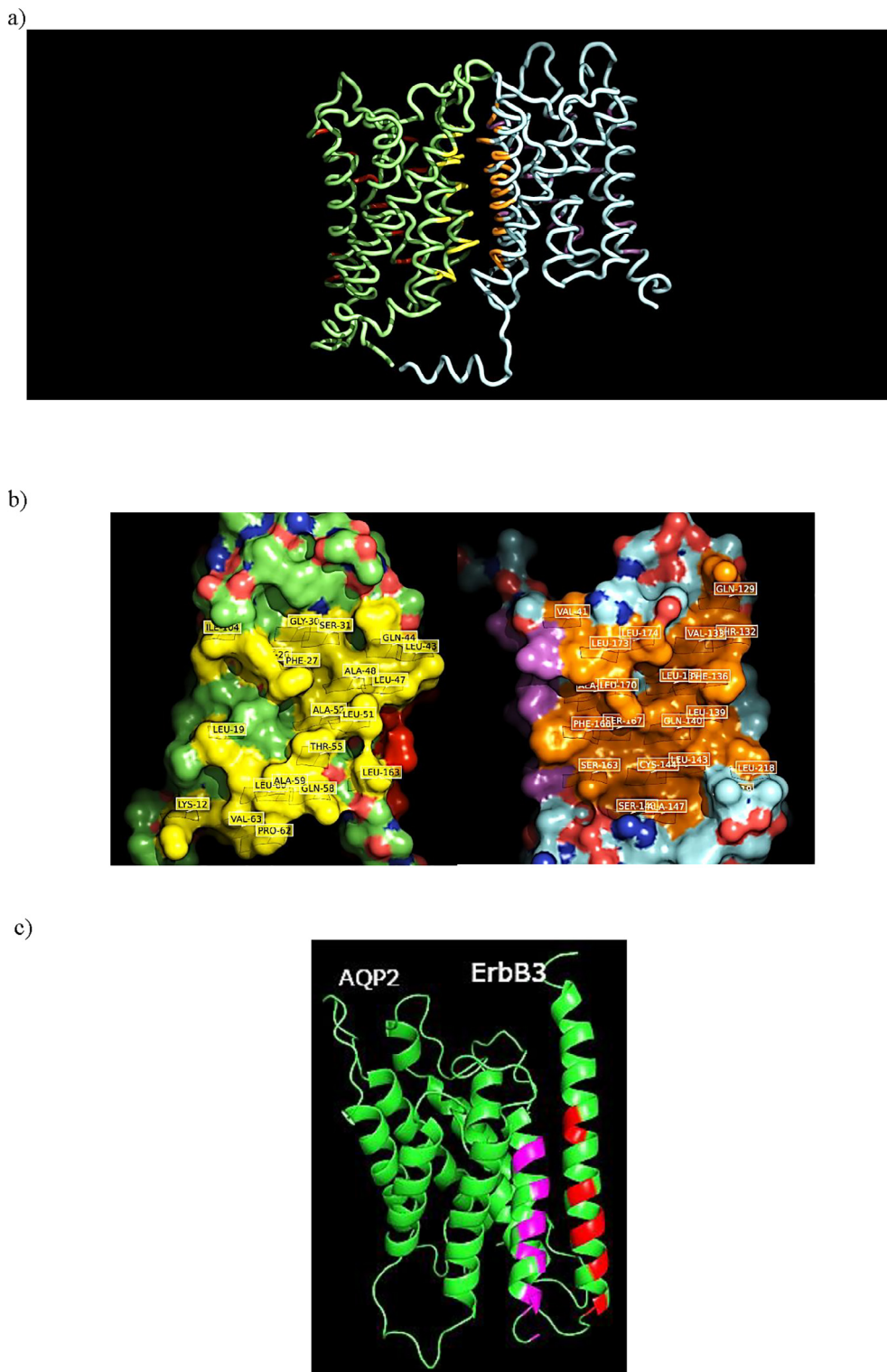


Fig. 6. Low-resolution cartoon structural model of predicted PPIMem TM-TM complexes obtained by molecular docking. a) *H. sapiens'* protomers of aquaporin 5 (UniProtKB P55064, PDB 3D9S, green) and aquaporin 2 (P41181, 4NEF, cyan) in complex, with PPIMem interface residues in yellow and orange, respectively. The binding motif for the former is K.(6)L.(6)FF.(2)GS.(11)LQ.(2)LA.(2)LA.(2)T.(2)QAL.PV.(40)I.(58)L, whereas that for the latter is R.(2)FA.(2)L.(3)L.(2)FF.(2)GS.(11)LQ.(2)MA.(2)LG.(2)T.VQALG.(42)LL.(57)L. Regions in red and magenta in the rear side of the molecules correspond to a second binding motif; b) the same complex in a solvent-accessible surface representation in which each chain has been rotated 90° towards the viewer, revealing the contact interface residues (labeled) –yellow for P55064 and orange for P41181; c) 3D model of the docking complex between *R. norvegicus'* protomers of aquaporin 2 (P34080) and the receptor tyrosine-protein kinase ErbB-3 (Q62799) in which the AV.(2)GL.(2)GA.(2)L binding motif, present on both protein surfaces, was used to direct the docking. AQP2 interface residues are in purple and ErbB-3 residues in red. (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)

Neutral amino acid transporter B(0); Q9NPY3 (CD93) Complement component C1q receptor.

In support of the proposed subnetwork, we found that all four TMs were present in two tissues of *H. sapiens* – adipose tissue (major) and breast (minor) as reported in the Human Protein Atlas (<https://www.proteinatlas.org/>). Furthermore, a CD93 - STAB1 interaction in humans has been reported in the String DB of PPI networks, and CD93 and PLPP1 are co-expressed in *G. gallus* [88].

3.6. Negative interactions

The importance of recording negative results of PPI assays in interatomic DBs, i.e., those indicating the tested proteins do not interact, has been raised [89]. However, their identification is less straightforward. This feature should eventually lead us to define a set of true negative interactions with the goal of training a predictor of TMPPIs as the sampling of negatives is crucial for optimal performance [90,91]. Looking for negative interactions in the IntAct DB for our PPIMem proteins, we find two negative interactions of two PPIMem proteins, but with two non-MPs. Analogously, the Negatome datasets [92,93] and Stelzl [94] compile sets of protein pairs that are unlikely to engage in direct physical interactions. We observed that spanning the Negatome dataset with our predicted positive interactions for *H. sapiens* gives no results. In other words, the Negatome DB does not report any of our complexes as a negative interaction. Conversely, several MPs in the Negatome set are absent from PPIMem. Even though this information is not conclusive, it goes along with our results; obviously, we will continue to probe our dataset as more negative interaction data become available.

3.7. Molecular docking simulations

To support the predictions from a structural point of view, we selected several TM pairs for molecular docking simulations for which PPIMem predicted an interaction based on the interface epitopes. To begin with, we took human AQP5 (P55064) and AQP2 (P41181). In this case, the crystal structures exist for both TMs (PDBs 3D9S and 4NEF, respectively). For the docking simulation, we took one of the PPIMem predicted patterns for AQP5 as the binding site (K.{6}L.{6}FF.{2}GS.{11}LQ.{2}LA.{2}LA.{2}T.{2}QALPV.{40}I.{58}L; Table S1), and one of the predicted patterns for AQP2 (R.{2}FA.{2}L.{3}L.{2}FF.{2}GS.{11}LQ.{2}MA.{2}LG.{2}T.VQALG.{42}LL.{57}L; Table S1). It is interesting to note that both motifs show contact residues separated by regions with many wildcard residues (40 and 58 for AQP5; 42 and 57 for AQP2). This is a direct consequence of the spatial distribution of the contact residues that may belong to different helices but still be involved in the same intermolecular interaction. Moreover, not all residues of a given exposed transmembrane helix are necessarily contact residues, but only a few, such as one Ile and one Leu at the end of the AQP5 motif and belonging to varied helices. The other example is the three Leu at the end of the AQP2 motif, two of which belong to one helix and the third one to another helix. That is why the primary structure of the PPIMem binding motif codes also for the nonlinear tertiary structure. Fig. 6a shows a heterodimer among the top 10 GRAMM-X docking predictions that involves exactly the membrane-exposed interface regions of each TM. In Fig. 6b each chain has been rotated 90° towards the viewer to reveal the contact interface residues.

On another hand, we selected the predicted rat AQP2 (P34080)-ErbB-3 (Q62779) pair, both subunits interacting through the AV.{2}GL.{2}GA.{2}L pattern on both protein surfaces. Since the experimental structures are unavailable for either TM, we homology-modeled them from human AQP2 (PDB 4NEF) and human ErbB-3 (PDB 2L9U), respectively. The sequences of both aquaporins (rat and human), are more than 30% identical in the TM region, just like

both ErbB-3's; the resulting individual 3D models are thus highly reliable. Fig. 6c shows that the modeled rat AQP2-ErbB-3 heterodimer respects the query interface, suggesting that the complex is viable. In both complexes, the contact residues are indeed at the interface of the complex and may lead to its formation.

Lastly, in the PPIMem user interface, when Valid = 2, the 3D structures of each of the isolated protomers making up a putative complex are available. It is therefore possible to perform the directed docking implemented above to obtain 3D structures of any one of the corresponding 67 complexes (63 for *H. sapiens*; Table S10). The PDB IDs can be found in the UniProtKB through the future UniProtKB link in our database.

4. Conclusions and discussion

In this work, we developed PPIMem, a wide-scope, interface residue template-based protocol destined to predict at large scale TM complexes resulting from direct physical interactions among their α -helical TM segments. PPIMem is a model-driven biological discovery tool to be queried for the discovery of verified and potential interactions, and for obtaining varied types of necessary information about them. It contains TM oligomerization recognition sites based on the key assumption that homologous structural interacting motifs always interact in similar orientations and with equivalent interfaces. In this work, we exclusively report interactions taking place in the eukaryotic plasma membrane interactome in which the binding sites specifically involve TM regions. PPIMem predicts therefore an TM protein interactome with thousands of *de novo* interactions, including multiple recognition sites, i.e., TMs with more than one interface, important for multimer formation. The obtained nonlinear sequence motifs identify homo and heterodimer interface amino acid residues that represent the first step to generating higher-order macromolecular edifices. Albeit our benchmark set is small (Table S7), it is of excellent quality and does not mix distinct types of organisms or membranes as other sets do. We have not assessed the effects of intramembrane mutations in the structure or function of the TMs, as the goal of our logical approach implementing different degrees of mutations was that of finding other TMs with homologous interfaces and thus forecasting new complexes that could lead to protein function annotation.

The uniqueness of the PPIMem approach resides on our focusing on the local, membrane-exposed interface residues, largely responsible for the formation of a complex between transmembrane proteins. The resulting TM interactome represents “first draft” predictions and contains 21,544 unique entries for all species dealt with, of which 9798 for *H. sapiens*. The considerable number of protein partners we uncover suggests that even distantly related or unrelated TM proteins make use of regions of their surface with similar sequences and arrangements to bind to other proteins. Thus, the TMs Q9JLF1 and Q5J316 predicted to form PPIMem complexes present the same interaction motif A.{3}V.{3}L.{6}S.{2}R.{19}L.{3}F.{2}L but the full sequences do not present significant identities to their template TM P23415 or to each other. The local binding motif sequence identities are two times greater (Table S11). The predicted interaction partners can lead to generating low-resolution 3D structures for several complexes, especially for those whose 3D structure of the individual subunits are available. In general, if complex formation is feasible as the interacting surfaces of the individual proteins manage to face each other in the docked complex provides a partial validation of the method.

There are many caveats of any analysis comparing PPIs from multiple sources [95], as there are large discrepancies and dramatic differences in the content between experimental PPI data collected by the same or different techniques, reflecting inherent limitations to each detection method, such as errors and ambiguities leading to false positives (FP). Consequently, our attempt to

compare our predictions to experimental data is more haphazard. Indeed, the intersections between various interaction maps that employ altogether diverse approaches are very small [96,97]. For example, even though the interactions in our template set exist as PDB structures, surprisingly most of them do not show in PPI databases. Moreover, the presence of orthologs and splice isoforms makes the research more cumbersome, not to mention that reported interactions are often functional associations and not necessarily direct, physical interactions between proteins leading to a complex [49]. To this, we have to add differences in curation policies [98]. Thus, the comparison between our predicted set of TM complexes and other datasets is necessarily qualitative. Disagreements occur also because TMs are not well probed in experiments, so that the screening of the huge potential interaction space is not complete. Because of all these factors, we do not have a large, validated dataset to compare to and assess the FPs. Nevertheless, to reduce the FPs we considered only motifs with six or more contact residues. Searching for a given protein in the database may result in too many entries in the PPIMem webpage. But if we consider that these entries include orthologs among different species and paralogs within the same species, the number of “primary” interactions, i.e., interactions between families of proteins, is much less. PPIMem thus allows to examine potential complexes and generate hypotheses for further investigation. But, for the time being, it is not possible to perform proper cross-validation, as our initial starting set (the template complexes) is small and includes only true interactions, no negative interactions. In few words, we cannot quantify the performance of PPIMem through a ROC curve, for example, and attribute scores for the ranking of the predicted complexes. Despite these limitations, we found several of our forecast interactions in different high-throughput experimental PPI DBs, validating in part our approach.

Complementary to the sequence-based co-evolution PPI prediction methods [8,99–101], our approach encodes 3D into 1D and thus adds the spatial dimension to a given TM interactome. This may lead to pioneering biological hypotheses concerning PPIs at the membrane level, to genotype-phenotype relationships, to investigation of the effect of pathological mutations on the interaction between TMs, and to propose molecular mechanisms of action. Recovering PPIMem predictions for human TM complexes in several experimental PPI datasets obtained by different methods (BioPlex 3.0, MENTHA, IID, HMRI) highlights the pertinence of our approach.

The predicted TM-TM interactions could be verified experimentally with specific techniques like video microscopy, FRET [102], or by our exclusive Microtubule Bench approach [103]. By applying machine learning methods, the PPIMem method can be improved by insuring that the TMs belong to the same developmental stage, tissue, cell type, site of expression, reaction and metabolic pathways; that they display functional similarity, and do not show a gene distance of more than 20 [104].

Incidentally, the PPIMem algorithm can be applied to soluble proteins and to other cell membranes (mitochondria, nucleoplasm, endoplasmic reticulum, Golgi apparatus) and across the tree of life, provided 3D structures of corresponding protein complexes are available. Our developed methodology can equally be extended by properly introducing side-chain and main-chain h-bond and electrostatic information at the interface, important for MPPIs [105]. Other applications of our approach include homologous protein networks in other organisms.

5. Data availability

Fig. S4 shows a screen capture of PPIMem’s first page. The corresponding opensource code and instructions for running the pre-

diction algorithm have been deposited at <https://github.com/PPIMem>. The database with the annotated predicted interactions is implemented as a web application that supports sorting and filtering. The output data can be downloaded as a csv file and the predictions can be accessed at <https://transint.univ-evry.fr>.

Several pertinent notes are to be found in the Supplemental Material.

Funding

This work was supported by:

- Hubert Curien CEDRE program, Grant 13 Santé/L2.
- Fonds pour le Rayonnement de la Recherche. Université d’Evry-Val-d’Essonne, Actions 2 and 3 – Incoming and outgoing mobilities.
- Fonds pour le Rayonnement de la Recherche, Université d’Evry-Val-d’Essonne, Action 1 – Financing for supporting the emergence of innovative projects in the framework of the evolution of scientific policy.
- French Embassy in Armenia. Fellowship of French Government.

CRediT authorship contribution statement

Georges Khazen: Methodology, Software, Validation, Formal analysis, Resources, Data curation, Visualization, Funding acquisition. **Aram Gyulkhandanian:** Investigation, Visualization. **Tina Issa:** Software, Investigation. **Rachid C. Maroun:** Conceptualization, Methodology, Validation, Formal analysis, Resources, Writing - original draft, Writing - review & editing, Visualization, Supervision, Project administration, Funding acquisition.

Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgments

We thank Prof. D. Pastré for critical reading of the manuscript. We thank A. Couturier and S. Castillon of the Direction des Systèmes d’Information of the Université d’Evry for implementing and installing the PPIMem database in a public server.

Appendix A. Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.csbj.2021.09.013>.

References

- [1] Kastritis PL, Bonvin AMJJ. On the binding affinity of macromolecular interactions: daring to ask why proteins interact. *J R Soc Interface* 2013;10: 20120835.
- [2] Yin H, Flynn AD. Drugging membrane protein interactions. *Annu Rev Biomed Eng* 2016;18:51–76.
- [3] Bocharov EV, Mineev KS, Pavlov KV, Kuznetsov AS, Efremov RG, Arseniev AS. Helix-helix interactions in membrane domains of bitopic proteins: Specificity and role of lipid environment. *Biochim Biophys Acta (BBA) – Biomembranes* 2017;1859:561–76.
- [4] Yamamoto K, Caporini MA, Im S-C, Waskell L, Ramamoorthy A. Transmembrane interactions of full-length mammalian bitopic cytochrome-P450-cytochrome-b5 complex in lipid bilayers revealed by sensitivity-enhanced dynamic nuclear polarization solid-state NMR spectroscopy. *Sci Rep* 2017;7:4116.
- [5] Guidolin D, Marcoli M, Tortorella C, Maura G, Agnati LF. G protein-coupled receptor-receptor interactions give integrative dynamics to intercellular communication. *Rev Neurosci* 2018.

- [6] Jubb HC, Pandurangan AP, Turner MA, Ochoa-Montaño B, Blundell TL, Ascher DB. Mutations at protein-protein interfaces: small changes over big surfaces have large impacts on human health. *Prog Biophys Mol Biol* 2017. <https://doi.org/10.1016/j.pbiomolbio.2016.10.002>.
- [7] Stevens TJ, Arkin IT. Do more complex organisms have a greater proportion of membrane proteins in their genomes? *Proteins Struct Funct Genet* 2000. [https://doi.org/10.1002/\(SICI\)1097-0134\(20000601\)39:4<417::AID-PROT140>3.0.CO;2-Y](https://doi.org/10.1002/(SICI)1097-0134(20000601)39:4<417::AID-PROT140>3.0.CO;2-Y).
- [8] Lage K. Protein-protein interactions and genetic diseases: The interactome. *CBiochim Biophys Acta (BBA) – Mol Basis Dis* 2004;1842:1971–80.
- [9] Bork P, Jensen LJ, von Mering C, Ramani AK, Lee I, Marcotte EM. Protein interaction networks from yeast to human. *Curr Opin Struct Biol* 2004;14:292–9.
- [10] Stumpf MPH, Thorne T, de Silva E, Stewart R, An HJ, Lappe M, et al. Estimating the size of the human interactome. *PNAS* 2008;105:6959–64.
- [11] Venkatesan K, Rual J-F, Vazquez A, Stelzl U, Lemmens I, Hirozane-Kishikawa T, et al. An empirical framework for binary interactome mapping. *Nat Methods* 2009;6:83–90.
- [12] Luck K, Kim D-K, Lambourne L, Spirohn K, Begg BE, Bian W, et al. A reference map of the human binary protein interactome. *Nature* 2020;580:402–8.
- [13] Licata L, Briganti L, Peluso D, Perfetto L, Iannuccelli M, Galeota E, et al. MINT, the molecular interaction database: 2012 Update. *Nucleic Acids Res* 2012. <https://doi.org/10.1093/nar/gkr930>.
- [14] Chatri-aryamontri A, Oughtred R, Boucher L, Rust J, Chang C, Kolas NK, et al. The BioGRID interaction database: 2017 update. *Nucleic Acids Res* 2017;45:D369–79.
- [15] Alfarano C, Andrade CE, Anthony K, Bahroos N, Bajec M, Bantoft K, et al. The Biomolecular Interaction Network Database and related tools 2005 update. *Nucleic Acids Res* 2005.
- [16] Szklarczyk D, Morris JH, Cook H, Kuhn M, Wyder S, Simonovic M, et al. The STRING database in 2017: quality-controlled protein-protein association networks, made broadly accessible. *Nucleic Acids Res* 2017;45:D362–8.
- [17] Orchard S, Ammari M, Aranda B, Breuza L, Briganti L, Broackes-Carter F, et al. The MINTAct project - IntAct as a common curation platform for 11 molecular interaction databases. *Nucleic Acids Res* 2014. <https://doi.org/10.1093/nar/gkt1115>.
- [18] Calderone A, Castagnoli L, Cesareni G. mentha: a resource for browsing integrated protein-interaction networks. *Nat Methods* 2013;10:690–1.
- [19] Sokolina K, Kittanakom S, Snider J, Kotlyar M, Maurice P, Gandía J, et al. Systematic protein-protein interaction mapping for clinically relevant human GPCRs. *Mol Syst Biol* 2017;13:918.
- [20] Peri S, Navarro JD, Amanchy R, Kristiansen TZ, Jonnalagadda CK, Surendranath V, et al. Development of human protein reference database as an initial platform for approaching systems biology in humans. *Genome Res* 2003;13:2363–71.
- [21] Keshava Prasad TS, Goel R, Kandasamy K, Keerthikumar S, Kumar S, Mathivanan S, et al. Human protein reference database–2009 update. *Nucleic Acids Res* 2009;37:D767–72.
- [22] Das J, Yu H. HINT: high-quality protein interactomes and their applications in understanding human disease. *BMC Syst Biol* 2012;6:1–12.
- [23] Giurgiu M, Reinhard J, Brauner B, Dunger-Kaltenbach I, Fobo G, Frishman G, et al. CORUM: the comprehensive resource of mammalian protein complexes–2019. *Nucleic Acids Res* 2019;47:D559–63.
- [24] Alonso-López G, Campos-Laborie FJ, Gutiérrez MA, Lambourne L, Calderwood MA, De Las Rivas J. APID database: redefining protein-protein interaction experimental evidences and binary interactomes. *Database (Oxford)* 2019;2019.
- [25] Hwang H, Petrey D, Honig B. A hybrid method for protein-protein interface prediction. *Protein Sci* 2016;25:159–65.
- [26] Huttlin EL, Ting L, Bruckner RJ, Gebreab F, Gygi MP, Szpyt J, et al. The BioPlex network: a systematic exploration of the human interactome. *Cell* 2015;162:425–40.
- [27] Kotlyar M, Pastrello C, Malik Z, Jurisica I. IID 2018 update: context-specific physical protein-protein interactions in human, model organisms and domesticated species. *Nucleic Acids Res* 2019;47:D581–9.
- [28] Alanis-Lobato G, Andrade-Navarro MA, Schaefer MH. HIPPIE v2.0: enhancing meaningfulness and reliability of protein-protein interaction networks. *Nucleic Acids Res* 2017;45:D408–14.
- [29] Miller JP, Lo RS, Ben-Hur A, Desmarais C, Stagljar I, Noble WS, et al. Large-scale identification of yeast integral membrane protein interactions. *Proc Natl Acad Sci* 2005;102:12123–8.
- [30] Garcia-García J, Schleker S, Klein-Seetharaman J, Oliva B. BIPS: BIANA Interolog Prediction Server. A tool for protein-protein interaction inference. *Nucleic Acids Res* 2012;40:W147–51.
- [31] Sarkar D, Jana T, Saha S. LMPID: a manually curated database of linear motifs mediating protein-protein interactions. *Database (Oxford)* 2015.
- [32] Ruepp A, Waegele B, Lechner M, Brauner B, Dunger-Kaltenbach I, Fobo G, et al. CORUM: the comprehensive resource of mammalian protein complexes–2009. *Nucleic Acids Res* 2010;38:D497–501.
- [33] Orchard S, Kerrien S, Abbani S, Aranda B, Bhate J, Bidwell S, et al. Protein interaction data curation: the International Molecular Exchange (IMEx) consortium. *Nat Methods* 2012;9:345–50.
- [34] Carpenter EP, Beis K, Cameron AD, Iwata S. Overcoming the challenges of membrane protein crystallography. *Curr Opin Struct Biol* 2008;18:581–6.
- [35] Kozma D, Simon I, Tusnády GE. PDBTM: Protein Data Bank of transmembrane proteins after 8 years. *Nucleic Acids Res* 2013;41:D524–9.
- [36] Mosca R, Pons T, Céol A, Valencia A, Aloy P. Towards a detailed atlas of protein-protein interactions. *Curr Opin Struct Biol* 2013;23:929–40.
- [37] Iyer K, Burkle L, Auerbach D, Thaminy S, Dinkel M, Engels K, et al. Utilizing the split-ubiquitin membrane yeast two-hybrid system to identify protein-protein interactions of integral membrane proteins. *Sci Signaling* 2005. <https://doi.org/10.1126/stke.2752005pl3>.
- [38] Rual J-F, Venkatesan K, Hao T, Hirozane-Kishikawa T, Dricot A, Li N, et al. Towards a proteome-scale map of the human protein-protein interaction network. *Nature* 2005;437:1173–8.
- [39] Petschnigg J, Groisman B, Kotlyar M, Taipale M, Zheng Y, Kurat CF, et al. The mammalian-membrane two-hybrid assay (MaMTH) for probing membrane-protein interactions in human cells. *Nat Methods* 2014;11:585–92.
- [40] Hubsman M, Yudkovsky G, Aronheim A. A novel approach for the identification of protein-protein interaction with integral membrane proteins. *Nucleic Acids Res* 2001;29.
- [41] Kittanakom S, Barrios-Rodiles M, Petschnigg J, Arnoldo A, Wong V, Kotlyar M, et al. CHIP-MYTH: a novel interactive proteomics method for the assessment of agonist-dependent interactions of the human β_2 -adrenergic receptor. *Biochem Biophys Res Commun* 2014;445:746–56.
- [42] You Z-H, Lei Y-K, Zhu L, Xia J, Wang B. Prediction of protein-protein interactions from amino acid sequences with ensemble extreme learning machines and principal component analysis. *BMC Bioinf* 2013;14(Suppl 8):S10.
- [43] You Z-H, Zhu L, Zheng C-H, Yu H-J, Deng S-P, Ji Z. Prediction of protein-protein interactions from amino acid sequences using a novel multi-scale continuous and discontinuous feature set. *BMC Bioinf* 2014;15(Suppl 15):S9.
- [44] Hashemifar S, Neyshabur B, Khan AA, Xu J. Predicting protein-protein interactions through sequence-based deep learning. *Bioinformatics* 2018;34:i802–10.
- [45] Szilagyi A, Zhang Y. Template-based structure modeling of protein-protein interactions. *Curr Opin Struct Biol* 2014;24:10–23.
- [46] Maheshwari S, Brylinski M. Template-based identification of protein-protein interfaces using eFindSitePPI. *Methods* 2016;93:64–71.
- [47] Keskin O, Tuncbag N, Gursoy A. Predicting protein-protein interactions from the molecular to the proteome level. *Chem Rev* 2016;116.
- [48] Zhang QC, Petrey D, Deng L, Qiang L, Shi Y, Thu CA, et al. Structure-based prediction of protein-protein interactions on a genome-wide scale. *Nature* 2012;490:556–60.
- [49] Zhang QC, Petrey D, Garzón JI, Deng L, Honig B. PrePPI: a structure-informed database of protein-protein interactions. *Nucleic Acids Res* 2013;41:D828–33.
- [50] Li B, Mendenhall J, Meiler J. Interfaces between alpha-helical integral membrane proteins: characterization, prediction, and docking. *Comput Struct Biotechnol J* 2019;17:699–711.
- [51] Burley SK, Berman HM, Kleywegt GJ, Markley JL, Nakamura H, Velankar S. Protein data bank (PDB): the single global macromolecular structure archive. *Methods Mol Biol* 2017;1607:627–41.
- [52] Lomize MA, Lomize AL, Pogozheva ID, Mosberg HI. OPM: orientations of proteins in membranes database. *Bioinformatics* 2006;22:623–5.
- [53] de Beer TAP, Berka K, Thornton JM, Laskowski RA. PDBsum additions. *Nucleic Acids Res* 2014;42:D292–6.
- [54] UniProt Consortium. UniProt: a hub for protein information. *Nucleic Acids Res* 2015;43:D204–12.
- [55] Aloy P, Ceulemans H, Stark A, Russell RB. The relationship between sequence and interaction divergence in proteins. *J Mol Biol* 2003;332:989–98.
- [56] Keskin O, Nussinov R. Favorable scaffolds: proteins with different sequence, structure and function may associate in similar ways. *Protein Eng Des Sel* 2005;18:11–24.
- [57] Gene Ontology Consortium. Gene Ontology Consortium: going forward. *Nucleic Acids Res* 2015;43:D1049–56.
- [58] Duarte JM, Srebnik A, Schärer MA, Capitani G. Protein interface classification by evolutionary analysis. *BMC Bioinf* 2012;13:334.
- [59] Capitani G, Duarte JM, Baskaran K, Bliven S, Somody JC. Understanding the fabric of protein crystals: computational classification of biological interfaces and crystal contacts. *Bioinformatics* 2016;32:481–9.
- [60] Duarte JM, Biyani N, Baskaran K, Capitani G. An analysis of oligomerization interfaces in transmembrane proteins. *BMC Struct Biol* 2013;13:21.
- [61] Elez K, Bonvin AMJJ, Vangone A. Distinguishing crystallographic from biological interfaces in protein complexes: role of intermolecular contacts and energetics for classification. *BMC Bioinf* 2018;19:438.
- [62] Jiménez-García B, Elez K, Koukos PI, Bonvin AM, Vangone A. PRODIGY-crystal: a web-tool for classification of biological interfaces in protein complexes. *Bioinformatics* 2019;35:4821–3.
- [63] Lomize AL, Pogozheva ID. Solvation models and computational prediction of orientations of peptides and proteins in membranes. *Methods Mol Biol* 2013;1063:125–42.
- [64] Tovchigrechko A, Vakser IA. GRAMM-X public web server for protein-protein docking. *Nucleic Acids Res* 2006;34:W310–4.
- [65] Tovchigrechko A, Vakser IA. Development and testing of an automated approach to protein docking. *Proteins* 2005;60:296–301.
- [66] Ulmschneider MB, Sansom MSP. Amino acid distributions in integral membrane protein structures. *Biochim Biophys Acta (BBA) – Biomembranes* 2001;1512:1–14.
- [67] Jha AN, Vishveshwara S, Banavar JR. Amino acid interaction preferences in proteins. *Protein Sci* 2010;19:603–16.

- [68] Mayol E, Campillo M, Cordero A, Olivella M. Inter-residue interactions in alpha-helical transmembrane proteins. *Bioinformatics* 2019;35:2578–84.
- [69] Eilers M, Patel AB, Liu W, Smith SO. Comparison of helix interactions in membrane and soluble alpha-bundle proteins. *Biophys J* 2002;82:2720–36.
- [70] Newport TD, Sansom MSP, Stansfeld PJ. The MemProtMD database: a resource for membrane-embedded protein structures and their lipid interactions. *Nucleic Acids Res* 2019;47:D390–7.
- [71] Mistry J, Chuguransky S, Williams L, Qureshi M, Salazar GA, Sonnhammer ELL, et al. Pfam: The protein families database in 2021. *Nucleic Acids Res* 2021;49:D412–9.
- [72] Walters RFS, DeGrado WF. Helix-packing motifs in membrane proteins. *Proc Natl Acad Sci U S A* 2006;103:13658–63.
- [73] Lemmon MA, Flanagan JM, Treutlein HR, Zhang J, Engelman DM. Sequence specificity in the dimerization of transmembrane alpha-helices. *Biochemistry* 1992;31:12719–25.
- [74] Russ WP, Engelman DM. The GxxxG motif: a framework for transmembrane helix-helix association. *J Mol Biol* 2000;296:911–9.
- [75] Senes A, Gerstein M, Engelman DM. Statistical analysis of amino acid patterns in transmembrane helices: the GxxxG motif occurs frequently and in association with beta-branched residues at neighboring positions. *J Mol Biol* 2000;296:921–36.
- [76] Gernert KM, Surles MC, Labean TH, Richardson JS, Richardson DC. The Alacoil: A very tight, antiparallel coiled-coil of helices. *Protein Sci* 1995;4:2252–60.
- [77] Sal-Man N, Gerber D, Bloch I, Shai Y. Specificity in transmembrane helix-helix interactions mediated by aromatic residues. *J Biol Chem* 2007;282:19753–61.
- [78] Gurezka R, Laage R, Brosig B, Langosch D. A heptad motif of leucine residues found in membrane proteins can drive self-assembly of artificial transmembrane segments. *J Biol Chem* 1999;274:9265–70.
- [79] Liu Y, Engelman DM, Gerstein M. Genomic analysis of membrane protein families: abundance and conserved motifs. *Genome Biol* 2002;3:1818.
- [80] Li E, Wimley WC, Hristova K. Transmembrane helix dimerization: beyond the search for sequence motifs. *Biochim Biophys Acta* 2012;1818:183–93.
- [81] Vreven T, Moal IH, Vangone A, Pierce BG, Kastriitis PL, Torchala M. Updates to the integrated protein-protein interaction benchmarks: docking benchmark version 5 and affinity benchmark version 2. *J Mol Biol* 2015;427.
- [82] Koukos PI, Faro I, van Noort CW, Bonvin AMJJ. A membrane protein complex docking benchmark. *J Mol Biol* 2018;430:5246–56.
- [83] Richards S, Aziz N, Bale S, Bick D, Das S, Gastier-Foster J, et al. Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genet Med* 2015;17:405–24.
- [84] Thanos CD, DeLano WL, Wells JA. Hot-spot mimicry of a cytokine receptor by a small molecule. *PNAS* 2006;103:15422–7.
- [85] Baspinar A, Cukuroglu E, Nussinov R, Keskin O, Gursoy A. PRISM: a web server and repository for prediction of protein-protein interactions and modeling their 3D complexes. *Nucleic Acids Res* 2014;42:W285–9.
- [86] Wodak SJ, Vlasblom J, Turinsky AL, Pu S. (2013) Protein-protein interaction networks: The puzzling riches.
- [87] Rao VS, Srinivas K, Sujini GN, Kumar GNS. Protein-protein interaction detection: methods and analysis. *Int J Proteomics* 2014;2014:147648.
- [88] Szklarczyk D, Gable AL, Lyon D, Junge A, Wyder S, Huerta-Cepas J, et al. STRING v11: protein-protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic Acids Res* 2019;47:D607–13.
- [89] Alvarez-Ponce D. Recording negative results of protein-protein interaction assays: an easy way to deal with the biases and errors of interactomic data sets. *Briefings Bioinf* 2016;18:bbw075.
- [90] Ben-Hur A, Noble WS. Choosing negative examples for the prediction of protein-protein interactions. *BMC Bioinf* 2006;7(Suppl 1):S2.
- [91] Trabuco LG, Betts MJ, Russell RB. Negative protein-protein interaction datasets derived from large-scale two-hybrid experiments. *Methods* 2012;58:343–8.
- [92] Smialowski P, Pagel P, Wong P, Brauner B, Dunger I, Fobo G, et al. The Negatome database: a reference set of non-interacting protein pairs. *Nucleic Acids Res* 2010;38:D540–4.
- [93] Blohm P, Frishman G, Smialowski P, Goebels F, Wachinger B, Ruepp A, et al. Negatome 2.0: a database of non-interacting proteins derived by literature mining, manual annotation and protein structure analysis. *Nucleic Acids Res* 2014;42:D396–400.
- [94] Stelzl U, Worm U, Lalowski M, Haenig C, Brembeck FH, Goehler H, et al. A human protein-protein interaction network: a resource for annotating the proteome. *Cell* 2005;122:957–68.
- [95] Mathivanan S, Periaswamy B, Gandhi T, Kandasamy K, Suresh S, Mohmood R, et al. An evaluation of human protein-protein interaction data in the public domain. *BMC Bioinf* 2006;7:S19.
- [96] Pitre S, Alamgir M, Green JR, Dumontier M, Dehne F, Golshani A. Computational methods for predicting protein-protein interactions. In *Advances in biochemical engineering/biotechnology* 2008;110:247–67.
- [97] Aloy P, Russell RB. The third dimension for protein interactions and complexes. *Trends Biochem Sci* 2002;27:633–8.
- [98] Turinsky AL, Razick S, Turner B, IM Donaldson, Wodak SJ. Literature curation of protein interactions: measuring agreement across major public databases. *Database (Oxford)* 2010;baq026.
- [99] Liu CH, Li KC, Yuan S. Human protein-protein interaction prediction by a novel sequence-based co-evolution method: co-evolutionary divergence. *Bioinformatics* 2013. <https://doi.org/10.1093/bioinformatics/bts620>.
- [100] Hamp T, Rost B. More challenges for machine learning protein interactions. *Bioinformatics* 2015;2:1–5.
- [101] Sun T, Zhou B, Lai L, Pei J. Sequence-based prediction of protein protein interaction using a deep-learning algorithm. *BMC Bioinf* 2017. <https://doi.org/10.1186/s12859-017-1700-2>.
- [102] Chen T, He B, Tao J, He Y, Deng H, Wang X, et al. Application of Förster Resonance Energy Transfer (FRET) technique to elucidate intracellular and In Vivo biofate of nanomedicines. *Adv Drug Deliv Rev* 2019;143:177–205.
- [103] Boca M, Kretov DA, Desforges B, Mephon-Gaspard A, Curmi PA, Pastré D. Probing protein interactions in living mammalian cells on a microtubule bench. *Sci Rep* 2015;5:17304.
- [104] Hopf TA, Schärfe CPI, Rodrigues JPGLM, Green AG, Kohlbacher O, Sander C, et al. Sequence co-evolution gives 3D contacts and structures of protein complexes. *Elife* 2014;3.
- [105] Bowie JU. Membrane protein folding: how important are hydrogen bonds? *Curr Opin Struct Biol* 2011;21:42–9.