

# SCIENTIFIC REPORTS

OPEN

## PiPred – a deep-learning method for prediction of $\pi$ -helices in protein sequences

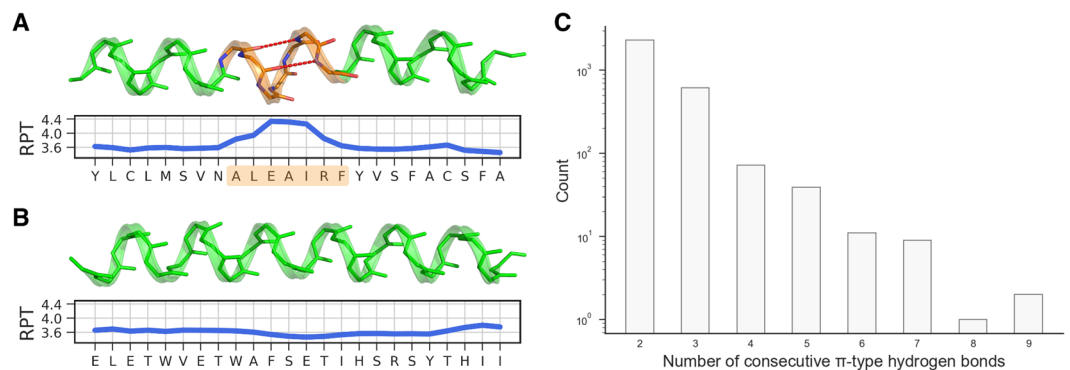
Jan Ludwiczak<sup>1,2</sup>, Aleksander Winski<sup>1</sup>, Antonio Marinho da Silva Neto<sup>1</sup>, Krzysztof Szczepaniak<sup>1</sup>, Vikram Alva<sup>3</sup> & Stanislaw Dunin-Horkawicz<sup>1</sup>

Canonical  $\pi$ -helices are short, relatively unstable secondary structure elements found in proteins. They comprise seven or more residues and are present in 15% of all known protein structures, often in functionally important regions such as ligand- and ion-binding sites. Given their similarity to  $\alpha$ -helices, the prediction of  $\pi$ -helices is a challenging task and none of the currently available secondary structure prediction methods tackle it. Here, we present PiPred, a neural network-based tool for predicting  $\pi$ -helices in protein sequences. By performing a rigorous benchmark we show that PiPred can detect  $\pi$ -helices with a per-residue precision of 48% and sensitivity of 46%. Interestingly, some of the  $\alpha$ -helices mispredicted by PiPred as  $\pi$ -helices exhibit a geometry characteristic of  $\pi$ -helices. Also, despite being trained only with canonical  $\pi$ -helices, PiPred can identify 6-residue-long  $\alpha/\pi$ -bulges. These observations suggest an even higher effective precision of the method and demonstrate that  $\pi$ -helices,  $\alpha/\pi$ -bulges, and other helical deformations may impose similar constraints on sequences. PiPred is freely accessible at: <https://toolkit.tuebingen.mpg.de/#/tools/quick2d>. A standalone version is available for download at: <https://github.com/labstructbioinf/PiPred>, where we also provide the CB6133, CB513, CASP10, and CASP11 datasets, commonly used for training and validation of secondary structure prediction methods, with correctly annotated  $\pi$ -helices.

Helices, dominant protein secondary structure elements, are defined by the recurring pattern of the hydrogen bonds between the amide hydrogen (NH) and the carbonyl oxygen (CO) atoms. In  $\alpha$ -helices, the most abundant type of helices, this interaction occurs between residues in positions  $i$  and  $i + 4$  in the amino acid sequence. Unlike  $\alpha$ -helices,  $\pi$ -helices, a less frequent type of helices, contain hydrogen bonds between residues in positions  $i$  and  $i + 5$  (Fig. 1). Canonical  $\pi$ -helices are characterized by the presence of at least two  $\pi$ -type ( $i \rightarrow i + 5$ ) hydrogen bonds and thus the minimal length of a  $\pi$ -helix is seven residues<sup>1</sup>. Shorter, six-residue-long segments containing a single  $\pi$ -type hydrogen bond are referred to as  $\alpha/\pi$ -bulges and should not be confused with canonical  $\pi$ -helices. The relative instability of  $\pi$ -helices is attributed to unfavorable dihedral angles<sup>2</sup>, weaker van der Waals interactions in the core of the helix<sup>3</sup>, and the large entropic cost of aligning five residues to form  $\pi$ -type hydrogen bonds<sup>3</sup>. As a consequence, most  $\pi$ -helices comprise just seven residues (two consecutive  $\pi$ -type hydrogen bonds); however, examples comprising even 14 residues (nine consecutive  $\pi$ -type hydrogen bonds) have also been identified (Fig. 1).

Evolutionarily,  $\pi$ -helices are thought to have originated through the insertion of single residues into  $\alpha$ -helical regions<sup>1</sup>. Despite the fact that such insertions are estimated to have a destabilizing effect, they have been identified in about 15% of all known protein structures, suggesting that they must provide a functional advantage. Indeed,  $\pi$ -helices are frequently found in functionally important regions such as ligand-binding sites<sup>1,4</sup> as well as transmembrane helical domains such as those of G-protein coupled receptors<sup>5,6</sup>. Given the functional importance of  $\pi$ -helices, it is crucial to develop computational tools for their prediction in structures or based on sequence information. The most popular tools for structure-based annotation of secondary elements, such as DSSP<sup>7</sup> or STRIDE<sup>8</sup>, tend to favor  $\alpha$ -helices over  $\pi$ -helices, frequently resulting in the absence of  $\pi$ -helices in assignments (recently DSSP was corrected to account for this problem<sup>9</sup>). However, several dedicated methods for

<sup>1</sup>Laboratory of Structural Bioinformatics, Centre of New Technologies, University of Warsaw, Banacha 2c, 02-097, Warsaw, Poland. <sup>2</sup>Laboratory of Bioinformatics, Nencki Institute of Experimental Biology, Pasteura 3, 02-093, Warsaw, Poland. <sup>3</sup>Department of Protein Evolution, Max-Planck-Institute for Developmental Biology, Max-Planck-Ring 5, 72076, Tübingen, Germany. Jan Ludwiczak and Aleksander Winski contributed equally. Correspondence and requests for materials should be addressed to S.D.-H. (email: [s.dunin-horkawicz@cent.uw.edu.pl](mailto:s.dunin-horkawicz@cent.uw.edu.pl))



**Figure 1.** Exemplary backbone structures of (A) a  $\pi$ -helix (PDB code: 1MXR, chain A, residues 194–217) and (B) an  $\alpha$ -helix (PDB code: 1MXR, A, 103–126). Plots below the structures indicate the number of residues per helical turn (RPT). The  $\pi$ -helical region is shown in orange and  $i \rightarrow i + 5$  hydrogen bonds are indicated with red dashed lines ( $i \rightarrow i + 4$  hydrogen bonds are not shown). (C) The length distribution of  $\sim 3,000$  representative  $\pi$ -helices. Note that the counts are shown in logarithmic scale.

the annotation of  $\pi$ -helices in protein structures have been developed<sup>11,10,11</sup>, providing the possibility of identifying  $\pi$ -helices that are missed by the general-purpose methods.

In comparison to  $\alpha$ -helices,  $\pi$ -helices exhibit different amino acid preferences<sup>12</sup> – aromatic and large aliphatic amino acids are preferred at the termini, whereas polar amino acids, particularly asparagines, tend to be present in the center<sup>11</sup>. Moreover, proline residues are frequently found directly after  $\pi$ -helices and have been suggested to promote the termination of the  $\pi$ -helical structure. Considering the presence of such sequence hallmarks, the prediction of  $\pi$ -helices directly from sequence appears to be possible.

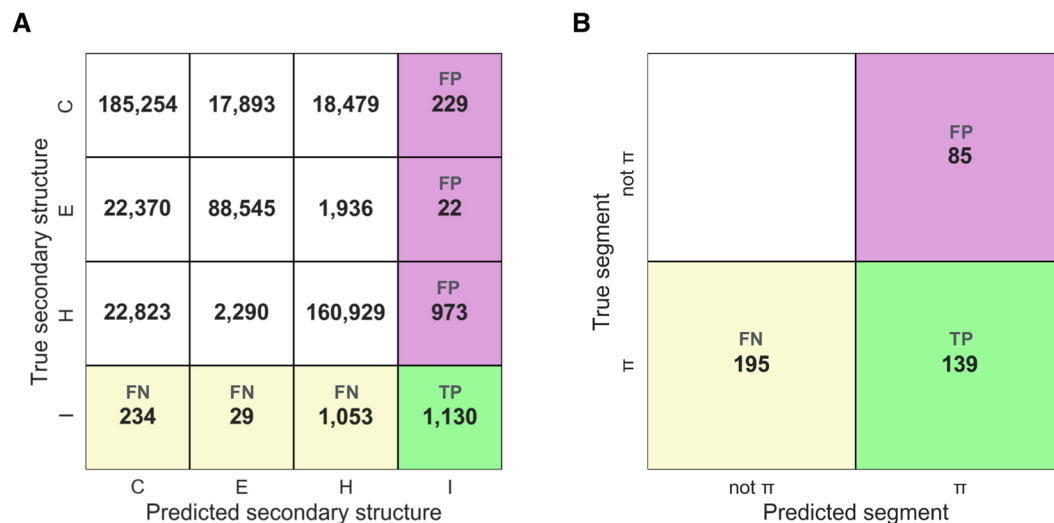
Modern prediction methods, frequently utilizing neural networks and deep learning approaches, achieve accuracies in the range of 75% to 85% for the 3-state secondary structure prediction problem (H:  $\alpha$ -helix, E:  $\beta$ -strand, and C: coil) and up to 70% for the 8-state case that considers five additional secondary structure elements: G:  $3_{10}$ -helix, I:  $\pi$ -helix, T: turn, B:  $\beta$ -bridge, and S: bend)<sup>13,14</sup>. However, none of the state-of-the-art 8-state predictors such as DeepCNP<sup>15</sup>, DCRNN<sup>16</sup>, CNNH\_PSS<sup>17</sup>, C8-Scorpion<sup>18</sup>, GSN<sup>19</sup>, RaptorXss<sup>20</sup>, and SSPro8<sup>21</sup> is capable of predicting  $\pi$ -helices, i. e. the accuracy of these predictors for  $\pi$ -helix class (“P”) is zero. This can be attributed to the properties of the datasets commonly used in the secondary structure prediction problems, like CB6133<sup>19,22</sup> or CB513<sup>19,23</sup>, which contain only a small number of  $\pi$ -helices due to inaccuracies in DSSP<sup>9</sup>. For example the original CB6133 dataset includes only 42 sequences containing a  $\pi$ -helix, which constitutes 0.7% of all sequences in the dataset, about 20 times lower than the fraction of the  $\pi$ -helices in PDB structures (for more details see Supplementary Table 2). The only method that is capable of predicting  $\pi$ -helices is limited to those occurring in transmembrane proteins<sup>24</sup>.

Here, we describe PiPred, a rigorously validated deep learning-based tool, trained on 20,295 diverse sequences containing 3,032 canonical  $\pi$ -helices with seven or more residues. For a given protein sequence, PiPred predicts the per-residue probability of the occurrence of  $\pi$ -helices and the three other basic secondary structure elements ( $\alpha$ -helices,  $\beta$ -strands, and unstructured regions). Moreover, we show that despite being trained with canonical  $\pi$ -helices, PiPred can also predict some six-residue-long  $\pi/\alpha$ -bulges and other helical distortions. PiPred is part of the Quick2D tool offered by the MPI Bioinformatics Toolkit<sup>25</sup> and a standalone version can be obtained from <https://github.com/labstructbioinf/PiPred>.

Additionally, for some commonly used datasets in the development of secondary structure prediction methods (e.g., CB6133, CB513), we offer updated versions with correctly annotated  $\pi$ -helices (Supplementary Table 2; <https://lbs.cent.uw.edu.pl/pipred>). The use of these updated datasets in the development of new secondary structure prediction methods should result in improved detection of  $\pi$ -helices. In fact, as a proof of concept, using the corrected datasets we retrained two state-of-the-art 8-state secondary structure prediction methods, CNNH\_PSS and DCRNN, to be able to detect  $\pi$ -helices, yet with accuracies considerably worse than that of PiPred.

## Results and Discussion

**Functional importance of  $\pi$ -helices in protein structures.** To assess the functional role of  $\pi$ -helices, we surveyed 2,555 representative  $\pi$ -helices present in protein structures co-crystallized with ligands and found that 24% of them interact with at least one ligand, most frequently with protoporphyrin IX and its derivatives (e.g. heme, chlorophyll), nucleoside derivatives (e.g. NAD, NADP, FAD), and various ions (e.g. phosphate, sulfate, zinc, and iron). Moreover, examination of 237 representative transmembrane (TM) structures obtained from PDBTM database<sup>26</sup> revealed that 45% of them contain at least one  $\pi$ -helix. Curiously,  $\pi$ -helices found in TM domains frequently (42%) also interact with ligands such as retinal and chlorophyll. To systematically investigate the association between the presence of  $\pi$ -helices and biological functions, we performed Gene Ontology (GO) enrichment analysis, with a focus on identifying GO terms overrepresented in proteins containing  $\pi$ -helices. We found that structures containing one or more  $\pi$ -helices are enriched with GO terms such as “oxidation-reduction process” (p-value =  $3e-61$ ), “heme binding” (p-value =  $7e-18$ ), “nucleotide binding” (p-value =  $1e-12$ ), and “metal ion binding” (p-value =  $7e-8$ ).



**Figure 2.** The performance of PiPred in detecting canonical  $\pi$ -helices comprising seven or more residues and containing at least two  $\pi$ -type  $i \rightarrow i + 5$  hydrogen bonds. In the confusion matrices, true positives (TP), false negatives (FN), and false positives (FP) are shown in green, yellow, and purple, respectively. **(A)** Performance of the per-residue predictions. Letters indicate secondary structure elements: I:  $\pi$ -helix, H:  $\alpha$ -helix, E:  $\beta$ -strand, and C: coil. Numbers indicate residue counts. **(B)** Performance of the per-segment prediction. Numbers indicate segment counts.

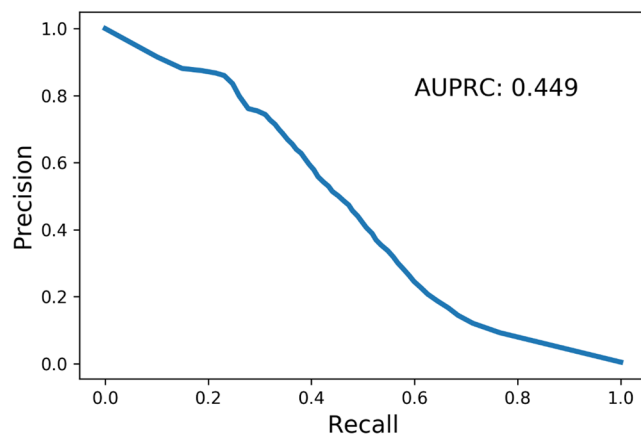
The above results illustrate that  $\pi$ -helices play an important role in binding ligands and in the functioning of helical transmembrane domains. These results are in agreement with previous observations that helical deformations, including  $\pi$ -helices, are frequently present in the proximity of NAD-based cofactors and heme groups<sup>4</sup>, and that  $\pi$ -helices participate in the coordination of ions<sup>10</sup>.

**PiPred predictor.** We scanned the structures in the Protein Data Bank for the presence of  $\pi$ -helices and subsequently constructed two sets: a training set comprising 20,295 structures used to train PiPred, a deep learning-based method for predicting  $\pi$ -helices in protein sequences; and an independent test set comprising 2,215 structures (Test set “7”) used to validate its performance. Importantly, to assure fairness of the validation procedure, the test set comprised only sequences that share no more than 30% identity with the sequences of the training set (for detailed statistics on the training and test sets see Supplementary Table 1).

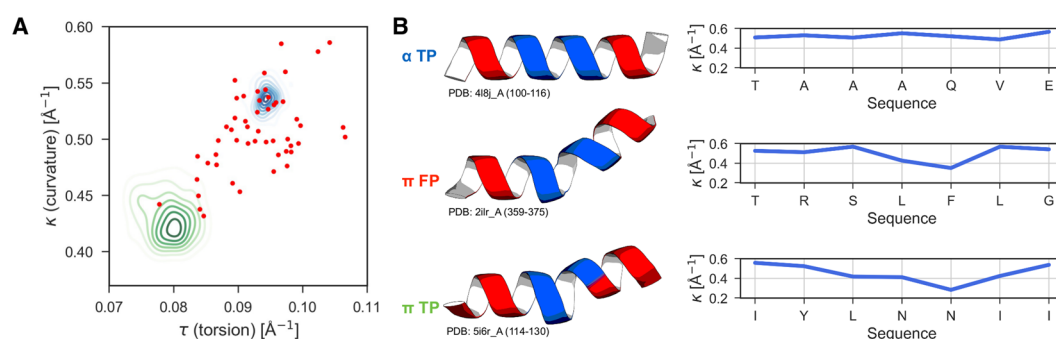
First, we determined the performance of PiPred in predicting  $\pi$ -helices at single-residue resolution. Every residue of each test set sequence was assigned a label corresponding to the predicted secondary structure type (I –  $\pi$ -helix, H –  $\alpha$ -helix, E –  $\beta$ -strand, C – unstructured region). The predicted labels were subsequently compared to the true labels defined based on their three-dimensional structure (for detailed procedure see Methods). True positives were defined as correctly predicted  $\pi$ -helical residues, whereas false negatives and false positives were defined as  $\pi$ -helical residues predicted as non- $\pi$ -helical and non- $\pi$ -helical predicted as  $\pi$ -helical, respectively. Using these values, we obtained a precision of 48% and a sensitivity of 46% (Fig. 2A). Moreover, the area under the precision-recall curve (AUPRC) analysis clearly indicated that PiPred is far-better than random assignment and that it is effective at various probability cut-offs (Fig. 3).

In addition to the per-residue performance statistics, in which all residues are treated separately, we used an alternative approach that is based on secondary structure segments. All true  $\pi$ -helical segments present in the test set and predicted by PiPred were pooled together. An overlap by at least one residue between true and predicted  $\pi$ -helices was considered to be a correct prediction (true positive). True  $\pi$ -helices and predicted  $\pi$ -helices that did not overlap with other segments were considered to be false negatives and false positives, respectively. The segment-based approach yielded a precision of 62% and a sensitivity of 42% (Fig. 2B). Most of the incorrect predictions (false negatives and false positives) resulted from the prediction of  $\pi$ -helices as  $\alpha$ -helices and *vice versa*: Among 195  $\pi$ -helical segments missed by PiPred (false negatives), 110 (56%) were mispredicted as  $\alpha$ -helices. Similarly, among 85 non- $\pi$ -helical regions predicted by PiPred as  $\pi$ -helices (false positives), 62 (73%) are  $\alpha$ -helices according to structure-based DSSP assignment.

**Common features of  $\pi$ -helices and other helical deformations.** To further investigate cases in which PiPred confuses  $\alpha$ - and  $\pi$ -helices with each other, we compared the geometry (for details see “Differential geometry analyses” in Methods) and hydrogen bonding patterns of  $\alpha$ -helices correctly predicted as  $\alpha$ -helices (“ $\alpha$  TP”; Fig. 4),  $\alpha$ -helices mispredicted as  $\pi$ -helices (“ $\pi$  FP”), and  $\pi$ -helices correctly predicted as  $\pi$ -helices (“ $\pi$  TP”). The geometry of helices was represented as the average per-residue curvature and torsion values: The curvature indicates how much a backbone curve deviates from a straight line at a given point, whereas torsion expresses how much the curve deviates from a plane at a given point. We found that  $\alpha$ -helices mispredicted as  $\pi$ -helices (red points in Fig. 4A) frequently show curvature and torsion values lying between those typical for  $\alpha$ - and  $\pi$ -helices. In such cases, the curvature and torsion values are lower than those of canonical  $\alpha$ -helices, which can



**Figure 3.** Precision-Recall plot for the detection of  $\pi$ -helices with PiPred.



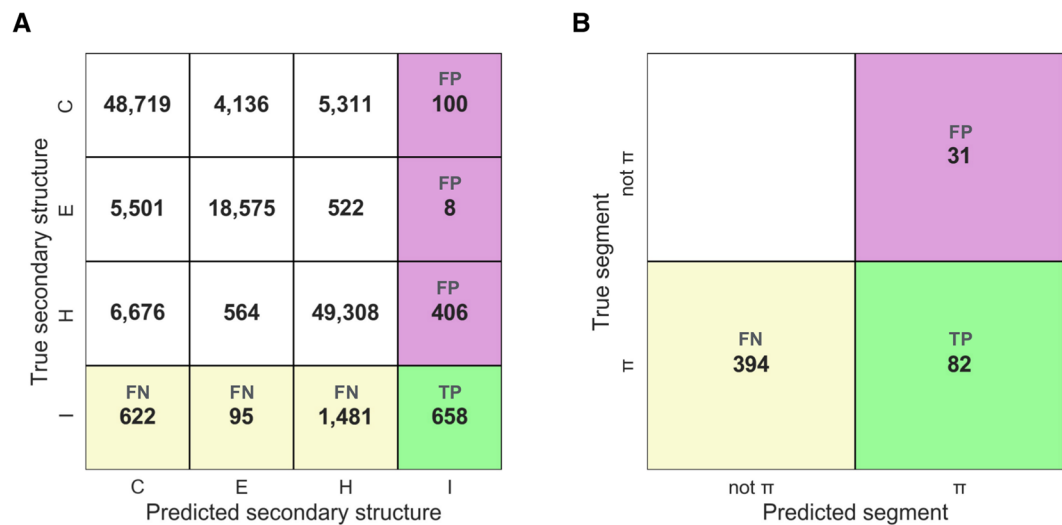
**Figure 4.** The geometry of helices in the test set.  $\alpha$  TP (true positives),  $\pi$  FP (false positives), and  $\pi$  TP indicate  $\alpha$ -helices correctly predicted by PiPred as  $\alpha$ -helices,  $\alpha$ -helices mispredicted as  $\pi$ -helices, and  $\pi$ -helices correctly predicted as  $\pi$ -helices, respectively. **(A)** Curvature and torsion values for  $\alpha$  TP and  $\pi$  TP are shown as kernel density estimate plots, and for  $\pi$  FP as points. **(B)** Exemplary structures and corresponding curvature plots. Residue ranges shown in the plots are colored blue in the structures.

be related to the increase of the helical radius and the number of residues per turn, features typical for  $\pi$ -helices and  $\pi/\alpha$ -bulges (Figs 1A and 4B). Despite this resemblance, these structures do not meet the criteria for canonical  $\pi$ -helices, as they lack two (or more) consecutive  $i \rightarrow i + 5$  hydrogen bonds with energy lower than alternative  $i \rightarrow i + 4$  bonds (see “Construction of the training and test sets” in the Methods).

Intrigued by this observation, we decided to perform a systematic analysis of cases in which non- $\pi$ -helical structures are predicted by PiPred as  $\pi$ -helices. To this end, we constructed an additional test set, Test set “6”, comprising 449 structures with a total of 476  $\pi/\alpha$ -bulges, i.e. six-residue-long secondary structure motifs characterized by the presence of a single  $i \rightarrow i + 5$  hydrogen bond (Supplementary Table 1). Analogously to the benchmark with the standard test set, we assessed the performance of PiPred in per-residue and per-segment detection of  $\pi/\alpha$ -bulges (Fig. 5). The per-residue predictions yielded precision of 56% and sensitivity of 23%, while the per-segment predictions resulted in precision and sensitivity of 73% and 17%, respectively.

The poor sensitivity of PiPred in detecting non- $\pi$ -helical deformations is to be expected as the method was trained using canonical  $\pi$ -helices comprising seven or more residues. However, the fact that PiPred was exclusively trained using canonical  $\pi$ -helices and yet can detect other helical deformations (Figs 4 and 5) suggests that canonical  $\pi$ -helices,  $\pi/\alpha$ -bulges as well as helical deformations that do not involve  $\pi$ -type  $i \rightarrow i + 5$  hydrogen bonding share similar sequence features to some extent. This can be interpreted in the context of the observation that the transition between  $\alpha$  and  $\pi$  conformations do occur in protein structures and can be achieved by shifting hydrogen bonding patterns. Such transitions were for example described in transmembrane helices<sup>5,27</sup> and monooxygenases<sup>1</sup>. In the latter case, they were termed “peristaltic-like shifts” and were associated with binding of a ligand. It is thus possible that some of the false positive predictions indicate that a given region has a propensity for the formation of a canonical  $\pi$ -helix, even though it does not assume such a conformation in an experimental structure.

**Application of the method.** To explore the possible applications of PiPred, we used it to scan 7,700 Pfam families with no structural data and comprising at least 30 sequences in the seed alignment. We found that 1200 of them may harbor uncharacterized  $\pi$ -helices, 82 of which are associated with the GO term “integral component of membrane” (enrichment p-value =  $5e-12$ ). In addition, a manual inspection of the Pfam descriptions revealed that a further 73 families were membrane-bound proteins, 54 various transporters, and nine G-protein-coupled



**Figure 5.** The performance of PiPred in detecting  $\pi/\alpha$ -bulges comprising six residues and containing a single  $\pi$ -type  $i \rightarrow i + 5$  hydrogen bond. For details see caption of Fig. 2.

	PiPred	DCRNN		CNNH_PSS	
	" $\pi$ " F1 score <sup>(a)</sup>	Q8 accuracy <sup>(b)</sup>	" $\pi$ " F1 score	Q8 accuracy	" $\pi$ " F1 score
Test set "7"	0.471	0.697	0.314	0.702	0.223
CB513 <sup>(c)</sup>	0.538	0.699	0.348	0.701	0.343
CASP10 <sup>(c)</sup>	0.557	0.720	0.328	0.725	0.286
CASP11 <sup>(c)</sup>	0.559	0.699	0.476	0.705	0.409

**Table 1.** Performance of retrained DCRNN and CNNH\_PSS methods on the PiPred "7" test set, and the updated CB513, CASP10 and CASP11 datasets. <sup>(a)</sup>F1 – harmonic average of the precision and sensitivity; <sup>(b)</sup>Q8 – index used for the evaluation of secondary structure prediction methods; <sup>(c)</sup>updated versions of the datasets were used. For details refer to Supplementary Table 2.

receptors. Of the remaining families, 342 were annotated as domains of unknown function. A flat-file containing these predictions is available at: <https://lbs.cent.uw.edu.pl/pipred>.

**Correction of the CB6133, CB5926, CB513, CASP10, and CASP11 datasets.** The most routinely used datasets for training and benchmarking secondary structure prediction methods (i.e., *CB6133*, *CB5926*, *CB513*, *CASP10*, and *CASP11 datasets*) contain only a very limited number of  $\pi$ -helices due to the limitations of DSSP ("Original datasets" in Supplementary Table 2). To address this issue, we corrected the secondary structure assignments in these datasets utilizing an established method for  $\pi$ -helix annotation<sup>1</sup>. Strikingly, these corrections resulted in a nearly 10-fold increase in the number of  $\pi$ -helices ("Updated datasets" in Supplementary Table 2) and did not introduce changes to the distribution of the other secondary structure labels (Supplementary Fig. 1). This prompted us to verify, as a proof of concept, whether two of the current state-of-the-art methods DCRNN and CNNH\_PSS, which are not capable of detecting  $\pi$ -helices, will be able to predict them after retraining with the updated datasets. Indeed, after retraining both methods became capable of detecting  $\pi$ -helices, while maintaining the same overall prediction accuracy as their original versions (Table 1). These updated datasets (available at: <https://lbs.cent.uw.edu.pl/pipred>) should therefore be useful for the development of general-purpose 8-state secondary structure prediction methods that are also capable of detecting  $\pi$ -helices.

To compare the new, retrained versions of DCRNN and CNNH\_PSS with PiPred, we tested them on the test set used to assess the performance of PiPred (Test set "7"; Supplementary Table 1) as well as on the updated CB513, CASP10, and CASP11 datasets (Supplementary Table 2). PiPred significantly outperformed both methods in all four tests (Table 1): PiPred "7" test set (p-value =  $1e-14$  and p-value =  $1e-30$  for DCRNN and CNNH\_PSS, respectively), updated CB513 set (p-value = 0.003 and p-value = 0.0004), updated CASP10 set (p-value = 0.02 and p-value = 0.006), and updated CASP11 set (p-value = 0.06 and p-value = 0.003).

## Conclusions

Our results illustrate that PiPred can be used for the annotation of potential functional sites in proteins since  $\pi$ -helices frequently contribute to protein-ligand interactions. Moreover, the prediction of  $\pi$ -helices and related helical distortions will be helpful for modeling the tertiary structure of transmembrane domains<sup>28</sup>. Finally, we envision that PiPred will also be useful for protein design tasks focused on the creation of ligand-binding pockets and new proteins with ligand-binding potential.

## Methods

**Construction of the training and test sets.** The Protein Data Bank structures, grouped into clusters comprising entries that share at least 50% sequence identity, were obtained from <https://www.rcsb.org/pages/download/ftp> (bc-50.out file). Structures longer than 700 or shorter than 30 residues, with resolution greater than 2.5 Å, and solved by methods other than X-ray crystallography were discarded. This yielded a set comprising 205,527 structures grouped within 24,276 clusters. Secondary structure state was assigned to each residue of each structure using DSSP<sup>7</sup>. Considering the weak performance of DSSP in detecting  $\pi$ -helices, we used a custom  $\pi$ -helix assignment method implemented based on<sup>1</sup>. First, all “T” labels ( $\pi$ -helices) defined by DSSP were substituted with “C” (coils). Next, residue ranges were re-marked as  $\pi$ -helical if (i) at least two  $\pi$ -type hydrogen bonds were present according to DSSP (in the cases where DSSP indicated two alternative hydrogen bonds, e.g.  $i \rightarrow i + 4$  and  $i \rightarrow i + 5$ , the one with lower energy was considered), (ii) at least one of the  $\pi$ -type hydrogen bonds had energy equal or lower than  $-2.0$  kcal/mol, and (iii) torsion angles for all residues in the range fell into the broadly defined helical region ( $-180^\circ < \varphi < 0^\circ$ ,  $-120^\circ < \Psi < 45^\circ$ ). Finally the 8-state assignment was reduced to a 4-state assignment (“H” – helix, containing “G” and “H” labels, “S” – strand, containing “E” and “B” labels, “C” – coil, containing “S”, “T”, and “C” labels, and “T” –  $\pi$ -helix). An analogous approach was used to generate an independent assignment of  $\pi/\alpha$ -bulges, i.e. secondary structure elements comprising just a single  $i \rightarrow i + 5$  hydrogen bond. Such six-residue-long regions were defined based on the presence of one  $\pi$ -type hydrogen bond of energy equal or lower than  $-2.0$  kcal/mol and torsion angles fulfilling the criteria listed above.

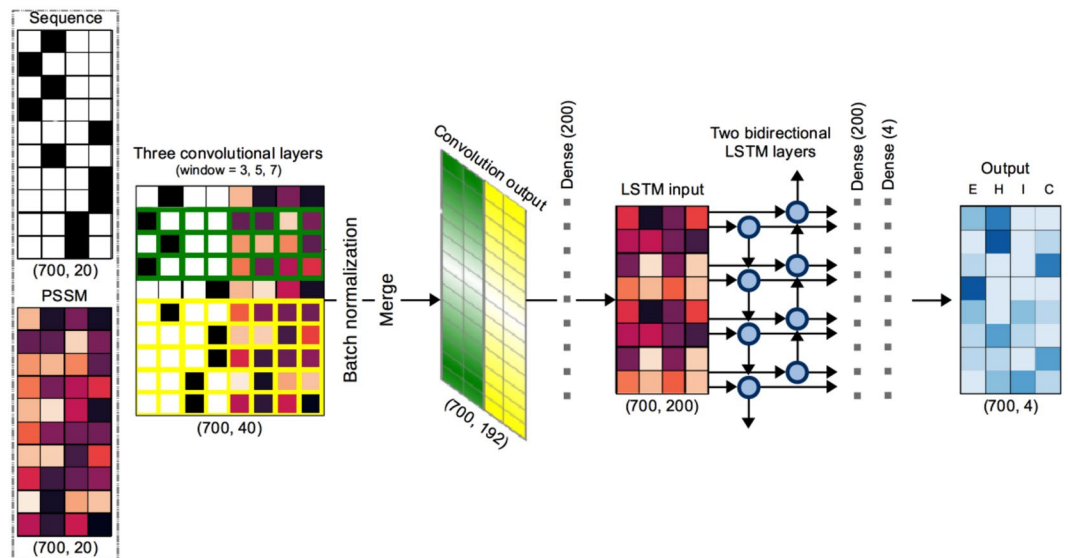
To build a dataset in which the pairwise sequence identity does not exceed 50%, we selected a representative structure from each of the aforementioned 24,276 clusters. To this end, we used the following procedure: First, from each cluster that does not contain any structure with a  $\pi$ -helix and/or a  $\pi/\alpha$ -bulge, a single representative structure with the lowest resolution was selected. From the remaining clusters, representative structures with the lowest resolution and longest  $\pi$ -helical segment were selected; however, those containing also  $\pi/\alpha$ -bulges were discarded. Such a procedure resulted in an initial dataset comprising 22,510 structures, of which 2,985 contained at least one canonical  $\pi$ -helix and did not contain any  $\pi/\alpha$ -bulge (positive examples set), and 19,525 did not contain  $\pi$ -helices as well as  $\pi/\alpha$ -bulges (negative examples set). From the remaining 1,750 structures, those containing at least one  $\pi/\alpha$ -bulge but no  $\pi$ -helices were selected, filtered to 30% sequence identity, and used to generate a separate “ $\pi/\alpha$ -bulge” set comprising 449 structures.

Sequences corresponding to the 22,959 structures (2,985 positive examples set, 19,525 negative examples set, and 449 “ $\pi/\alpha$ -bulge” set) were used as queries in PSI-BLAST<sup>29</sup> (E-value  $< 0.001$ , three iterations) searches of the NCBI non-redundant protein sequence database filtered to 90% sequence identity for the calculation of position specific scoring matrices (PSSMs). Subsequently, the 22,510 structures of the positive and negative examples sets were merged, shuffled, and randomly assigned to the training and test sets comprising 20,295 and 2,215 sequences, respectively. Importantly, we ensured that all sequences of the test set show no more than 30% similarity to any sequence of the training set. Furthermore, we also ensured that the two datasets contained an equal percentage of the  $\pi$ -helical residues, which amounted to 0.46% of all residues in each set. Detailed statistics for all data sets are shown in Supplementary Table 1.

**Sequence encoding.** For encoding sequences and their corresponding PSSM profiles, we used a procedure in which a sequence is encoded as a  $700 \times 40$  matrix, where 700 and 40 corresponds to the maximal sequence length and the number of features associated with every residue, respectively. Out of 40 features, 20 denoted “one-hot” encoded amino acid and another 20 were the PSSM probabilities transformed by the sigmoid function. Finally, if the sequence was shorter than 700 residues it was padded with zeros randomly at the C- or N- terminal ends to match the  $700 \times 40$  matrix size.

**Deep-learning model architecture and training.** The cascaded convolutional and recurrent network architecture, based on the DCRNN secondary structure predictor<sup>16</sup>, with minor modifications, was implemented in Keras<sup>30</sup> using Tensorflow<sup>31</sup> backend (Fig. 6). Sequences and PSSMs encoded as  $700 \times 40$  matrices (see “Sequence encoding” section above for details) were independently introduced into three 1D convolutional layers (window lengths 3, 5 and 7), each of which contains 64 filters and is activated with the *tanh* function. The output of the three convolutional layers, i.e. three  $700 \times 64$  matrices, were passed through batch normalization layers and concatenated to yield a  $700 \times 192$  matrix. Next, the concatenated matrix was used as an input to a fully-connected layer containing 200 neurons and activated with the ReLU function. To detect the dependencies between distant residues based on the local features extracted by convolutional layers two bidirectional LSTM layers were used. Each consisted 200 neurons and their dropout and recurrent dropout parameters were set to 0.5, and activation and recurrent activation functions were set to *tanh* and sigmoid, respectively. Finally, a dense layer with 200 neurons and the ReLU activation function was used to connect the LSTM with an output layer containing 4 neurons and the softmax activation function. The final output of the network is a vector  $700 \times 4$  indicating the residue-wise probabilities of the 4 secondary structure classes: I –  $\pi$ -helix, H –  $\alpha$ -helix, E –  $\beta$ -strand, C – unstructured region. The implementation in Python and Keras is available at: <https://github.com/labstructbioinf/PiPred>.

The training process involved optimization of the network’s parameters using pairs of encoded sequences and their corresponding correct secondary structure labels. The training process was performed in a 10-fold cross-validation (CV) framework: the training set was randomly divided into 10 equally-sized parts, each containing approximately the same number of  $\pi$ -helical residues. In each CV round, one part served as validation set, whereas the remaining nine together as training set. The training was performed for 50 epochs with the ‘Adam’ optimizer<sup>32</sup> (the learning rate was set to 0.0003, whereas the remaining parameters were set to their default values) with categorical cross-entropy as the loss function (to account for significant underrepresentation of  $\pi$ -helices, the  $\pi$ -helix class was weighted by a factor of five). From each CV round, the best model (according to the F1-score of  $\pi$ -helix classification) was selected and the resulting 10 models were used to build the final ensemble predictor



**Figure 6.** Schematic representation of the prediction model architecture. Numbers below matrices indicate their dimensionality. The input sequence and the corresponding PSSM are encoded as  $700 \times 40$  matrix, which is introduced into three independent convolutional layers with window lengths of 3, 5, and 7 (for clarity only windows of the length of 3 and 5 are shown in green and yellow, respectively). The output of each convolutional layer is normalized and subsequently, all outputs are merged and passed through a dense layer, two LSTM layers, another dense layer, and an output layer. The output is a  $700 \times 4$  matrix where each row denotes probabilities of E (strand), H ( $\alpha$ -helix), I ( $\pi$ -helix), and C (unstructured coil) occurring at the given position of the input sequence. For details refer to “Deep-learning model architecture and training” section of Methods.

(PiPred). For each residue of a given sequence and PSSM, PiPred returns four probabilities corresponding to four secondary structure classes (I, H, E, and C). In each position of the input sequence, the predicted secondary structure is defined as the one with the highest probability.

**Function of  $\pi$ -helices.** 22,510 structures of the initial dataset (positive and negative examples) were used to identify associations between the presence of  $\pi$ -helices and biological functions. The PDB to Gene Ontology (GO) mappings were downloaded from [www.geneontology.org/gene-associations/goa\\_pdb.gaf](http://www.geneontology.org/gene-associations/goa_pdb.gaf)<sup>33</sup> and 11,066 out of 22,510 structures were found to have at least one associated GO term. These structures were analyzed with GOATOOLS<sup>34</sup> to identify the enrichment of GO terms in 1,773 structures containing one or more  $\pi$ -helices (p-values were adjusted with the Holm method). In addition to the GO enrichment statistics, we analyzed 2,555  $\pi$ -helices present in structures containing ligands using PLIP<sup>35</sup>.

**Differential geometry analyses.** A differential geometry representation of protein backbones was used to evaluate and compare the geometry of (i)  $\pi$ -helices correctly predicted as  $\pi$ -helices, (ii)  $\alpha$ -helices correctly predicted as  $\alpha$ -helices, and (iii)  $\alpha$ -helices mispredicted as  $\pi$ -helices (false positives). To this end, we used the FlexGeo<sup>36</sup> method, which implements an approach analogous to that used in CHORAL<sup>37</sup>, ARABESQUE<sup>38</sup>, and Polyphony<sup>39</sup>. The protein backbone is represented by a piecewise cubic spline interpolation using the  $C\alpha$  atoms as knots, i.e. as a regular smooth curve  $\vec{r}(t)$  parametrized by  $C\alpha$  residue number  $t$ . According to the Fundamental Theorem of Curves, any regular spatial curve can be fully characterized by its curvature,  $\kappa$ , and torsion,  $\tau$ , values as a function of arc length  $s$ . Given the allowable change of parameters,  $s$  and  $t$  are related by:

$$\frac{ds}{dt} = \left\| \frac{d\vec{r}}{dt} \right\| \quad (1)$$

Therefore, it is possible to calculate  $\kappa$  and  $\tau$  using the following equations:

$$\kappa = \bar{\kappa} = \left\| \frac{d\vec{T}}{ds} \right\| = \frac{\|\dot{r} \times \ddot{r}\|}{\|\dot{r}\|^3} \quad (2)$$

$$\tau = \|\bar{\tau}\| = \left\| \frac{d\vec{B}}{ds} \right\| = \frac{\|(\dot{r} \times \ddot{r}) \cdot \ddot{r}\|}{\|\dot{r} \times \ddot{r}\|^2} \quad (3)$$

where,  $\vec{T}$  is the tangent vector and  $\vec{B}$  is the binormal vector of the Frenet-Serret frame of the curve. Each residue can then be represented by its respective  $C\alpha$   $\kappa$  and  $\tau$  values. The curvature values express how much a given point

of a curve deviates from a straight line in comparison to the previous point, i. e.  $\kappa = 0$  only if the point does not change the tangent vector  $\vec{T}$  of the curve. Similarly, the torsion expresses how a given point deviates from a plane in comparison to the previous point, i. e.  $\tau = 0$  only if the point does not change the binormal vector  $\vec{B}$  of the curve. The units of  $\kappa$  and  $\tau$  are per  $\text{\AA}^{-1}$ .

Based on the per-segment classification (Fig. 2B), we defined 53  $\alpha$ -helices mispredicted as canonical  $\pi$ -helices (comprising seven or more residues), 127  $\pi$ -helices correctly predicted as  $\pi$ -helices (only predicted  $\pi$ -helical segments comprising seven or more residues were considered), and 127  $\alpha$ -helices correctly predicted as  $\alpha$ -helices (out of >12,000 correctly predicted  $\alpha$ -helices 127 were randomly selected to ensure the balance). These helical segments were analyzed using the differential geometry procedure and the resulting curvature as well as torsion values were averaged for each segment and plotted (Fig. 4A).

**Correction of the CB6133, CB5926, CB513, CASP10, and CASP11 datasets.** CB6133 and CB513<sup>19</sup> are standard datasets used in the development of computational tools for protein secondary structure prediction. The CB6133 dataset is composed of training and test subsets, and therefore can be used on its own in the development process. Alternatively, a CB6133 dataset variant, CB6133\_filtered, obtained by removing sequences that exhibit >25% sequence identity to sequences in the CB513 dataset, can be used for training and the CB513 dataset for validation. Recently, after the discovery of duplicates in the original CB6133 dataset, an updated version of it (CB5926) as well as an updated version of its filtered dataset (CB5926\_filtered) were released.

As these datasets were constructed using older DSSP assignments, which did not annotate  $\pi$ -helices accurately, they show a low abundance of  $\pi$ -helices (Supplementary Table 2). Consequently, as most secondary structure prediction methods train their models using these datasets, they do not predict  $\pi$ -helices. We therefore decided to correct the annotation of  $\pi$ -helices within these datasets. To this end, we downloaded datasets CB5926, CB5926\_filtered, CB6133, CB6133\_filtered, and CB513 from <https://www.princeton.edu/~jzthree/datasets/ICML2014/> and extracted sequences for all entries. Since the sequences did not contain the corresponding PDB identifiers, we built a sequence database corresponding to crystallographic (resolution less than 2.5  $\text{\AA}$ ) and NMR structures (preference was given to X-ray structures if multiple entries were matched). Subsequently, each sequence from the five datasets was used to search our custom PDB database with BLAST and ideal matches, spanning whole query sequence range, were sorted according to their resolution. In each case, the match with the lowest resolution was used to generate updated secondary structure labels with the aid of the procedure described in “Construction of the training and test sets”.

To investigate the effects of these corrections, we re-trained two state-of-the-art secondary structure prediction methods, CNNH\_PSS<sup>17</sup> and DCRNN<sup>16</sup>. Both were trained and tested using the original CB5926\_filtered and CB513 datasets, respectively, and their updated variants with corrected secondary structure labels. For testing, in addition to the CB513 dataset, we also used the test set developed for the benchmarking of PiPred as well as CASP10 and CASP11 datasets. Testing for the significance of the difference between the performance of PiPred and the other methods was performed as in<sup>40</sup> using paired t-test; F1 scores were calculated for the individual sequences containing  $\pi$ -helices.

## Data Availability

PiPred is available as a web service (<https://toolkit.tuebingen.mpg.de/#/tools/quick2d>) and as a standalone software (<https://github.com/labstructbioinf/PiPred>). Results of the Pfam scan and corrected CB6133, CB5926, CB513, CASP10, and CASP11 datasets can be downloaded from <https://lbs.cent.uw.edu.pl/pipred>.

## References

- Cooley, R. B., Arp, D. J. & Karplus, P. A. Evolutionary origin of a secondary structure:  $\pi$ -helices as cryptic but widespread insertional variations of  $\alpha$ -helices that enhance protein functionality. *J. Mol. Biol.* **404**, 232–46 (2010).
- Ramachandran, G. N. & Sasisekharan, V. Conformation of polypeptides and proteins. *Adv. Protein Chem.* **23**, 283–438 (1968).
- Rohl, C. A. & Doig, A. J. Models for the 3(10)-helix/coil, pi-helix/coil, and alpha-helix/3(10)-helix/coil transitions in isolated peptides. *Protein Sci.* **5**, 1687–96 (1996).
- Cartailler, J.-P. & Luecke, H. Structural and functional characterization of pi bulges and other short intrahelical deformations. *Structure* **12**, 133–44 (2004).
- Ren, Z., Ren, P. X., Balusu, R. & Yang, X. Transmembrane Helices Tilt, Bend, Slide, Torque, and Unwind between Functional States of Rhodopsin. *Sci. Rep.* **6**, 34129 (2016).
- Riek, R. P., Rigoutsos, I., Novotny, J. & Graham, R. M. Non-alpha-helical elements modulate polytopic membrane protein architecture. *J. Mol. Biol.* **306**, 349–62 (2001).
- Kabsch, W. & Sander, C. Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers* **22**, 2577–637 (1983).
- Frishman, D. & Argos, P. Knowledge-based protein secondary structure assignment. *Proteins* **23**, 566–79 (1995).
- van der Kant, R. & Vriend, G. Alpha-bulges in G protein-coupled receptors. *Int. J. Mol. Sci.* **15**, 7841–64 (2014).
- Riek, R. P. & Graham, R. M. The elusive  $\pi$ -helix. *J. Struct. Biol.* **173**, 153–60 (2011).
- Fodje, M. N. & Al-Karadaghi, S. Occurrence, conformational features and amino acid propensities for the pi-helix. *Protein Eng.* **15**, 353–8 (2002).
- Kumar, P. & Bansal, M. Dissecting  $\pi$ -helices: sequence, structure and function. *FEBS J.* **282**, 4415–32 (2015).
- Jiang, Q., Jin, X., Lee, S.-J. & Yao, S. Protein secondary structure prediction: A survey of the state of the art. *J. Mol. Graph. Model.* **76**, 379–402 (2017).
- Yang, Y. *et al.* Sixty-five years of the long march in protein secondary structure prediction: the final stretch? *Brief. Bioinform.* **bbw129**, <https://doi.org/10.1093/bib/bbw129> (2016).
- Wang, S., Peng, J., Ma, J. & Xu, J. Protein Secondary Structure Prediction Using Deep Convolutional Neural Fields. *Sci. Rep.* **6**, 18962 (2016).
- Li, Z. & Yu, Y. Protein Secondary Structure Prediction Using Cascaded Convolutional and Recurrent Neural Networks. in *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence* 2560–2567 (AAAI Press, 2016).
- Zhou, J., Wang, H., Zhao, Z., Xu, R. & Lu, Q. CNNH\_PSS: protein 8-class secondary structure prediction by convolutional neural network with highway. *BMC Bioinformatics* **19**, 60 (2018).



18. Yaseen, A. & Li, Y. Template-based C8-SCORPION: a protein 8-state secondary structure prediction method using structural information and context-based features. *BMC Bioinformatics* **15**, S3 (2014).
19. Zhou, J. & Troyanskaya, O. G. Deep Supervised and Convolutional Generative Stochastic Network for Protein Secondary Structure Prediction. In *Proceedings of the 31st International Conference on International Conference on Machine Learning - Volume 32* I-745–I-753 (JMLR.org, 2014).
20. Wang, Z., Zhao, F., Peng, J. & Xu, J. Protein 8-class secondary structure prediction using conditional neural fields. *Proteomics* **11**, 3786–3792 (2011).
21. Magnan, C. N. & Baldi, P. SSpro/ACCpro 5: almost perfect prediction of protein secondary structure and relative solvent accessibility using profiles, machine learning and structural similarity. *Bioinformatics* **30**, 2592–2597 (2014).
22. Wang, G. & Dunbrack, R. L. PISCES: a protein sequence culling server. *Bioinformatics* **19**, 1589–91 (2003).
23. Cuff, J. A. & Barton, G. J. Evaluation and improvement of multiple sequence methods for protein secondary structure prediction. *Proteins* **34**, 508–19 (1999).
24. Rigoutsos, I., Riek, P., Graham, R. M. & Novotny, J. Structural details (kinks and non-alpha conformations) in transmembrane helices are intrahelically determined and can be predicted by sequence pattern descriptors. *Nucleic Acids Res.* **31**, 4625–31 (2003).
25. Zimmermann, L. *et al.* A Completely Reimplemented MPI Bioinformatics Toolkit with a New HHpred Server at its Core. *J. Mol. Biol.* **430**, 2237–2243 (2018).
26. Kozma, D., Simon, I. & Tusnády, G. E. PDBTM: Protein Data Bank of transmembrane proteins after 8 years. *Nucleic Acids Res.* **41**, D524–9 (2013).
27. Cao, Z. & Bowie, J. U. Shifting hydrogen bonds may produce flexible transmembrane helices. *Proc. Natl. Acad. Sci. USA* **109**, 8121–6 (2012).
28. Chen, K.-Y. M., Sun, J., Salvo, J. S., Baker, D. & Barth, P. High-resolution modeling of transmembrane helical protein structures from distant homologues. *PLoS Comput. Biol.* **10**, e1003636 (2014).
29. Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.* **25**, 3389–402 (1997).
30. Chollet, F. & others. Keras. <https://keras.io> (2015).
31. Abadi, M. *et al.* TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed Systems. *CoRR* **abs/1603.0** (2016).
32. Kingma, D. P. & Ba, J. L. Adam: a Method for Stochastic Optimization. *Int. Conf. Learn. Represent. 2015*. <http://doi.acm.org.ezproxy.lib.ucf.edu/10.1145/1830483.1830503> (2015).
33. Ashburner, M. *et al.* Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat. Genet.* **25**, 25–9 (2000).
34. Klopfenstein, D. V. *et al.* GOATOOLS: A Python library for Gene Ontology analyses. *Sci. Rep.* **8**, 10872 (2018).
35. Salentin, S., Schreiber, S., Haupt, V. J., Adasme, M. F. & Schroeder, M. PLIP: fully automated protein-ligand interaction profiler. *Nucleic Acids Res.* **43**, W443–7 (2015).
36. Silva Neto, A. M., Silva, S. R., Vendruscolo, M., Camilloni, C. & Montalvão, R. W. A superposition free method for protein conformational ensemble analyses and local clustering based on a differential geometry representation of backbone. *Proteins Struct. Funct. Bioinforma. prot.* 25652, <https://doi.org/10.1002/prot.25652> (2019).
37. Montalvão, R. W., Smith, R. E., Lovell, S. C. & Blundell, T. L. CHORAL: A differential geometry approach to the prediction of the cores of protein structures. *Bioinformatics* **21**, 3719–3725 (2005).
38. Leung, H. T. A., Montaña, B. O., Blundell, T., Vendruscolo, M. & Montalvão, R. W. Arabesque: a Tool for Protein Structural Comparison Using Differential Geometry and Knot Theory. *World Res. J. Pept. Protein* **1**, 33–40 (2012).
39. Pitt, W. R., Montalvão, R. W. & Blundell, T. L. Polyphony: superposition independent methods for ensemble-based drug discovery. *BMC Bioinformatics* **15**, 324 (2014).
40. Hanson, J., Paliwal, K., Litfin, T., Yang, Y. & Zhou, Y. Improving prediction of protein secondary structure, backbone angles, solvent accessibility and contact numbers by using predicted contact maps and an ensemble of recurrent and residual convolutional neural networks. *Bioinformatics*, <https://doi.org/10.1093/bioinformatics/bty1006> (2018).

## Acknowledgements

This work was supported by the Polish National Science Centre (grant number 2015/18/E/NZ1/00689 to SDH). VA was supported by institutional funds from the Max Planck Society. Computations were carried out with the support of the Interdisciplinary Centre for Mathematical and Computational Modeling (ICM) University of Warsaw (grant number GA67-18 to SDH).

## Author Contributions

S.D.-H., J.L. and V.A. wrote the main manuscript text and S.D.-H. prepared Figures 1, 2, 4, and 5. A.W. prepared Figures 3 and 6. All authors reviewed the manuscript. J.L. performed pi-helix annotation and prepared datasets used in the study. J.L. and A.W. designed and implemented machine learning framework, A.M. performed differential geometry analysis, whereas K.S. performed protein-ligand interactions analysis. V.A. incorporated PiPred to the MPI Bioinformatics Toolkit.

## Additional Information

**Supplementary information** accompanies this paper at <https://doi.org/10.1038/s41598-019-43189-4>.

**Competing Interests:** The authors declare no competing interests.

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2019