



RESEARCH ARTICLE

NetSurfP-2.0: Improved prediction of protein structural features by integrated deep learning

Michael Schantz Klausen^{1†} | Martin Closter Jespersen² | Henrik Nielsen² |
Kamilla Kjærgaard Jensen² | Vanessa Isabell Jurtz² | Casper Kaae Sønderby³ |
Morten Otto Alexander Sommer¹ | Ole Winther^{3,4} | Morten Nielsen^{2,5} |
Bent Petersen^{2,6}  | Paolo Marcatili² 

¹Novo Nordisk Foundation Center for Biosustainability, Technical University of Denmark, Kongens Lyngby, Denmark

²Department of Bio and Health Informatics, Technical University of Denmark, Kongens Lyngby, Denmark

³The Bioinformatics Centre, Department of Biology, University of Copenhagen, Copenhagen, Denmark

⁴Department of Applied Mathematics and Computer Science, Technical University of Denmark, Kongens Lyngby, Denmark

⁵Instituto de Investigaciones Biotecnológicas, Universidad Nacional de San Martín, Buenos Aires, Argentina

⁶Faculty of Applied Sciences, Centre of Excellence for Omics-Driven Computational Biodiscovery (COMBio), AIMST University, Kedah, Malaysia

Correspondence

Paolo Marcatili and Bent Petersen, Department of Bio and Health Informatics, Technical University of Denmark, Kongens Lyngby, Denmark.

Emails: pamar@bioinformatics.dtu.dk; bentpetersenwork@gmail.com

Present address

Bent Petersen, Natural History Museum of Denmark, University of Copenhagen, 2200 Copenhagen N, Denmark.

Funding information

Novo Nordisk Fonden

Abstract

The ability to predict local structural features of a protein from the primary sequence is of paramount importance for unraveling its function in absence of experimental structural information. Two main factors affect the utility of potential prediction tools: their accuracy must enable extraction of reliable structural information on the proteins of interest, and their runtime must be low to keep pace with sequencing data being generated at a constantly increasing speed. Here, we present NetSurfP-2.0, a novel tool that can predict the most important local structural features with unprecedented accuracy and runtime. NetSurfP-2.0 is sequence-based and uses an architecture composed of convolutional and long short-term memory neural networks trained on solved protein structures. Using a single integrated model, NetSurfP-2.0 predicts solvent accessibility, secondary structure, structural disorder, and backbone dihedral angles for each residue of the input sequences. We assessed the accuracy of NetSurfP-2.0 on several independent test datasets and found it to consistently produce state-of-the-art predictions for each of its output features. We observe a correlation of 80% between predictions and experimental data for solvent accessibility, and a precision of 85% on secondary structure 3-class predictions. In addition to improved accuracy, the processing time has been optimized to allow predicting more than 1000 proteins in less than 2 hours, and complete proteomes in less than 1 day.

KEYWORDS

deep learning, disorder, local structure prediction, secondary structure, solvent accessibility

1 | INTRODUCTION

The Anfinsen experiment, showing that the structural characteristics of a protein are encoded in its primary sequence alone, is more than

50 years old.¹ As a practical application of it, several methods have been developed over the last decades to predict from sequence only several protein structural features, including solvent accessibility, secondary structure, backbone geometry, and disorder.^{2–7} These tools have tremendously impacted biology and chemistry, and some are among the most cited works in the field. They have been extensively used to annotate novel sequences, thus facilitating their

[†]Michael Schantz Klausen and Martin Closter Jespersen contributed equally to the work.

characterization. The accuracy of said methods plays a central role here: the rate of errors in many computationally-generated annotations is a well-known and unresolved problem affecting public databases.⁸ Such errors often propagate through databases and sequence annotations, and the availability of high-quality predictions is hence of primary importance to limit their occurrence.

On the other hand, the amount of novel sequences has been steadily increasing over the last years,⁹ and not only experimental methods, but also computational predictions of structural and functional features have a hard time keeping up with it. This creates a conflict between the need for accurate predictions, and the pace at which we can generate them.

NetSurfP-1.0¹⁰ is a tool published in 2009 for prediction of solvent accessibility and secondary structure using a feed-forward neural network architecture. Since then, deep learning techniques have affected the application of machine learning in biology expanding the ability of prediction tools to produce more accurate results on complex datasets.^{11–16} Here, we present NetSurfP-2.0, a new extended version of NetSurfP, that uses a deep neural network approach to accurately predict absolute and relative solvent accessibility, secondary structure using both 3- and 8-class definitions,¹⁷ φ and ψ dihedral angles, and structural disorder,¹⁸ of any given protein from its primary sequence only. By having an integrated deep model with several outputs, NetSurfP-2.0 can not only significantly reduce the computational time, but also achieve an improved accuracy that could not be reached by having separate models for each feature. In fact, when assessed on various test sets with less than 25% sequence identity to any protein used in the training, its accuracy was consistently on par with or better than that of other state-of-the-art tools.^{3,4,19–21} In particular, we observed a significant increase in the accuracy of solvent accessibility, secondary structure, and disorder over all the other tested methods.

NetSurfP-2.0 uses different approaches to make predictions for small and large sets of sequences, thus improving its time efficiency without compromising its accuracy. It has a user-friendly interface allowing non-expert users to obtain and analyze their results via a browser, thanks to its graphical output, or to download them in several common formats for further analysis.

2 | MATERIALS AND METHODS

We describe briefly the dataset and method used for training NetSurfP-2.0, and the validations performed.

2.1 | Structural dataset

A structural dataset consisting of 12 185 crystal structures was obtained from the Protein Data Bank (PDB),²² culled and selected by the PISCES server²³ with 25% sequence similarity clustering threshold and a resolution of 2.5 Å or better. To avoid over fitting, any sequence that had more than 25% identity to any sequences in the test datasets (see “Evaluation” section for details) was removed, as well as peptide chains with less than 20 residues, leaving 10 837 sequences. Finally, we randomly selected 500 sequences (validation

set) left out for early stopping and parameter optimization, leaving 10 337 sequences for training.

2.2 | Structural features

For all residues in each chain in the training dataset, we calculated its absolute and relative solvent accessibility (ASA and RSA, respectively), 3- and 8-class secondary structure classification (SS3 and SS8, respectively), and the backbone dihedral angles φ and ψ using the DSSP software.¹⁷ Finally, each residue that was present in the chain refseq sequence, but not in the solved structure, was defined as disordered. It is important to mention that disordered residues cannot be annotated with any of the other features, since no atomic coordinates are available for these residues.

2.3 | Protein sequence profiles

NetSurfP-2.0, like its predecessor, exploits sequence profiles of the target protein for its predictions. We used 2 different ways for generating such profiles. The first exploits the HH-suite software,²⁴ that runs quickly on individual sequences, while the second uses the MMseqs2 software,²⁵ that is optimized for searches on large data sets. In both cases, the profile-generation tools were run with default parameters, except MMseqs2 which used 2 iterations with the “-max-seqs” parameter set to 2000.

2.4 | Deep network architecture

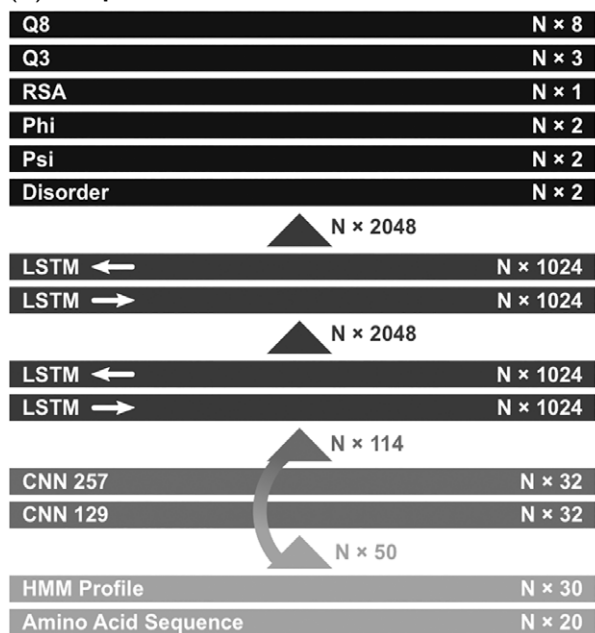
The model was implemented using the Keras library. The input layer of the model consists of the one-hot (sparse) encoded sequences (20 features) plus the full HMM profiles from HH-suite (30 features in total, comprising 20 features for the amino acid profile, 7 features for state transition probabilities, and 3 features for local alignment diversity), giving a total of 50 input features. This input is then connected to 2 convolutional neural network (CNN) layers, consisting of 32 filters each with size 129 and 257, respectively. The CNN output is concatenated with the initial 50 input features and connected to 2 bidirectional long short-term memory (LSTM) layers with 1024 nodes (Figure 1, panel A).

Each output (RSA, SS8, SS3, φ , ψ , and disorder) is calculated with a fully connected (FC) layer using the outputs from the final LSTM layer. RSA is encoded as a single output between 0 and 1. ASA output is not directly predicted, but calculated by multiplying RSA and ASAmx.²⁶ SS8, SS3, and disorder, are encoded as 8, 3, or 2 outputs with the target encoded as a sparse vector (target is set to 1, while rest of the elements are 0). φ and ψ are each encoded as a vector of length 2, where the first element is the sine of the angle and the second element is the cosine. This encoding reduces the effect of the periodicity of the angles,²⁷ and the predicted angle can be calculated with the arctan2 function.

2.5 | Training

The training was performed using mini-batches of size 15. The individual learning rate of each neuron was optimized using the Adam function.²⁸ Early stopping was performed on the validation set. Since the different

(A) Deep model architecture overview



(B) Computational time per sequence

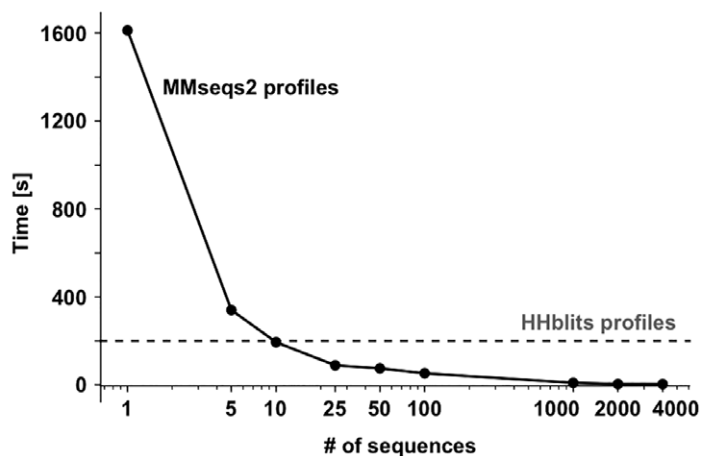


FIGURE 1 Network architecture and computation time plot. In panel A the network architecture is shown. N is the number of amino acids in the target protein sequence. Each box represents a different layer of the network, from the input (bottom) to the output (top), and the corresponding number of neurons/filters. The arrows represent the features that are passed between consecutive layers. The computation time per sequence of NetSurfP-2.0 is reported in panel B. The x-axis represents the number of input sequences (logarithmic scale), the y-axis the average computation time in seconds per sequence. The method implementation using HH-suite profiles is reported as a gray dashed line, and the one using MMSeqs2 profiles is reported as a solid black line

target values for each output have different distributions, a weighted sum of different loss functions was used: SS8, SS3, and disorder use cross entropy loss, while RSA, φ and ψ use mean squared error loss. Such weights were adjusted so each loss contribution was approximately equal and then fine-tuned for maximum overall performance. When the target value for a feature of a given residue was missing, for example, for secondary structure of disordered residues, or φ angles of N-terminal residues, the loss for that output was set to 0.

2.6 | Evaluation

The final 2 models, one trained with the HH-suite and one with MMseqs2 profiles, were tested on 3 independent datasets: the TS115 dataset (115 proteins),²⁹ the CB513 dataset (513 protein regions from 434 proteins),³⁰ and a dataset consisting of all the free-modeling targets (21 proteins) at the CASP 12 experiment.³¹ No protein with more than 25% sequence identity to the proteins in these datasets was present in the training. Disorder prediction was not performed on the CB513 dataset since it contains very few disordered regions.

We used different metrics to evaluate each feature: Pearson's correlation coefficient (PCC) for solvent accessibilities, Q3 and Q8 accuracy for SS3 and SS8 respectively,¹⁷ mean absolute error in degrees for φ and ψ angles (MAE), Matthews correlation coefficient (MCC) and false positive rate (FPR) for disorder. In each dataset, the performance was calculated both as the average over all the residues in the dataset (*per residue*) and as the average of the performances per structure, the latter being defined as the average of the metric on all the residue of each given structure. The P values between the top-

scoring method and all the other methods on a given feature in a dataset were calculated using a 2-tailed paired Student's t test on the corresponding performances per structure.

2.7 | DisProt dataset

We retrieved all the entries in DisProt,³² a database containing proteins annotated with several experimentally validated disorder types. All proteins with an available solved structure were removed to avoid overlaps with the training set. For each residue of each protein, we compare its experimental disorder annotation to the disorder prediction from NetSurfP-2.0. We classified proteins with more than 75% of their residues being disordered as completely disordered proteins.

3 | RESULTS

We have compared the performance of NetSurfP-2.0 to other state-of-the-art tools with similar functionality: NetSurfP-1.0,¹⁰ Spider3,⁴ SPOT-Disorder,³ RaptorX,^{20,21} and JPred4.¹⁹ In order to check whether the results of the methods were significantly different, we calculated a P value for each feature by using a pairwise Student's t test on the results of the 2 methods. Results on the independent test sets CASP12,³¹ TS115,²⁹ and CB513³⁰ are presented in Table 1, and in an extended version in Supporting Information Table S1. They match to a very high degree with the results obtained by using a 4-fold cross validation on the training set (Supporting Information Table S2).

TABLE 1 Results of the method's validation on independent test datasets

	RSA [PCC]	ASA [PCC]	SS3 [Q3]	SS8 [Q8]	Disorder [MCC]	Disorder [FPR]	Phi [MAE]	Psi [MAE]
CASP12								
NetSurfP-2.0 (mmseqs)	0.726	0.735	0.820	0.703	0.660	0.015	20.3	31.8
NetSurfP-2.0 (hhblits)	0.725	0.737	0.824	0.711	0.604	0.011	20.0	31.2
NetSurfP-1.0	0.617	0.641	0.709					
Spider3		0.688	0.791		0.582	0.026	21.6	33.2
RaptorX	0.594		0.786	0.661	0.621	0.045		
Jpred4			0.760					
TS115								
NetSurfP-2.0 (mmseqs)	0.778	0.797	0.857	0.750	0.656	0.006	17.2	25.8
NetSurfP-2.0 (hhblits)	0.775	0.795	0.853	0.744	0.663	0.008	17.5	26.5
NetSurfP-1.0	0.661	0.691	0.779					
Spider3		0.769	0.839		0.575	0.027	18.5	27.3
RaptorX	0.651		0.822	0.716	0.567	0.044		
Jpred4			0.767					
CB513								
NetSurfP-2.0 (mmseqs)	0.794	0.807	0.854	0.723			20.1	28.0
NetSurfP-2.0 (hhblits)	0.788	0.803	0.853	0.720			20.2	28.6
NetSurfP-1.0	0.700	0.723	0.788					
Spider3		0.797	0.845				20.4	28.2
RaptorX	0.676		0.827	0.706				
Jpred4			0.779					

The performance of NetSurfP-2.0 (using HH-suite and MMSeqs2 profiles), NetSurfP-1.0, Spider3, SPOT-disorder, RaptorX, and JPred4, is displayed for the CASP12, TS115, and CB513 datasets. SPOT-disorder and Spider3 predictions are reported as a single row. The following performance metrics are used: Pearson correlation coefficient (PCC), Q3 and Q8 accuracy, Matthew's correlation coefficient (MCC), False Positive Rate (FPR), and mean absolute error (MAE) in degrees. RaptorX RSA prediction classes B (buried), M (medium), and E (exposed), are mapped to the middle of the corresponding interval used to define them, that is, to 0.05, 0.25, and 0.7, respectively. The different predicted features are reported in the column header, together with the corresponding performance metric. For each feature and each dataset, the best score is reported in bold. Scores in italics are the ones for which no significant difference with respect to the top scoring method is observed (P value > 0.05). Empty cells represent predictions that were not performed, either because not part of a method's output, or because the feature was not present in the corresponding dataset.

We give here the results for each individual feature, and a benchmark of the time performance of the tool. An example of the ASA and SS3 predictions for the human Orotate phosphoribosyltransferase (OPRTase) domain, displayed on its solved structure (PDB id 2WNS) is illustrated in Figure 2.

3.1 | Solvent accessibility

The main focus of the original NetSurfP-1.0 and of its updated version is to predict solvent accessibility for individual residues. Both tools predict RSA, and also calculate the corresponding ASA as described in the Methods section. We have compared the performance of the old and updated version of our tool, and compared to that of Spider3, a recent tool with a similar architecture that only predicts the ASA, and RaptorX, which performs a 3-class prediction of Buried ($RSA < 0.1$), Medium ($0.1 < RSA < 0.4$) and Exposed ($0.4 < RSA$) residues. RaptorX 3-class predictions were transformed to numeric predictions by assigning to each class the middle of its RSA interval. We have also calculated NetSurfP-2.0 accuracy by projecting its predictions to RaptorX 3-class RSA classes (Supporting Information Figure S1). The results shown in Table 1 demonstrate a consistently improved performance of NetSurfP-2.0, compared to the other methods, with a PCC of ~ 0.8 on the test datasets TS115 and CB513 and on the validation set. The PCC of NetSurfP-2.0 on the CASP12 dataset is ~ 0.72 , inferior to the other data sets but still significantly better than all other

methods. It has to be noted that many of the structures in the CASP12 dataset are not obtained through X-ray crystallography, and they contain a number of disordered regions (as defined in the Methods) substantially larger than both the other external dataset and the training data. We see in general (Supporting Information Table S3) that all the predictions are less accurate in the few residues before and after a disordered region. We will further discuss this later.

3.2 | Secondary structure

Many tools for secondary structure prediction have been published over the last 20 years, with their reported accuracy improving over time.²⁹ In many cases, these tools have been subject of independent validation studies^{30,33} to get an independent assessment of their actual capabilities. We have decided to compare our tool with Spider3, RaptorX, and Jpred4, as these are among the most commonly used and accurate tools available. All the aforementioned tools perform 3-class prediction of secondary structure, while only NetSurfP-2.0 and RaptorX also provide an 8-class prediction. The results of the benchmark are given in Table 1. In all cases, NetSurfP-2.0 produce significantly more accurate predictions than all the other tools, with an average Q3 accuracy of $\sim 85\%$, and a Q8 accuracy between 72% and 75% (Supporting Information Table S1 and Figure S2). As for solvent accessibility prediction, the results on the CASP12 dataset are less accurate for all tested tools. A difference between the old version of

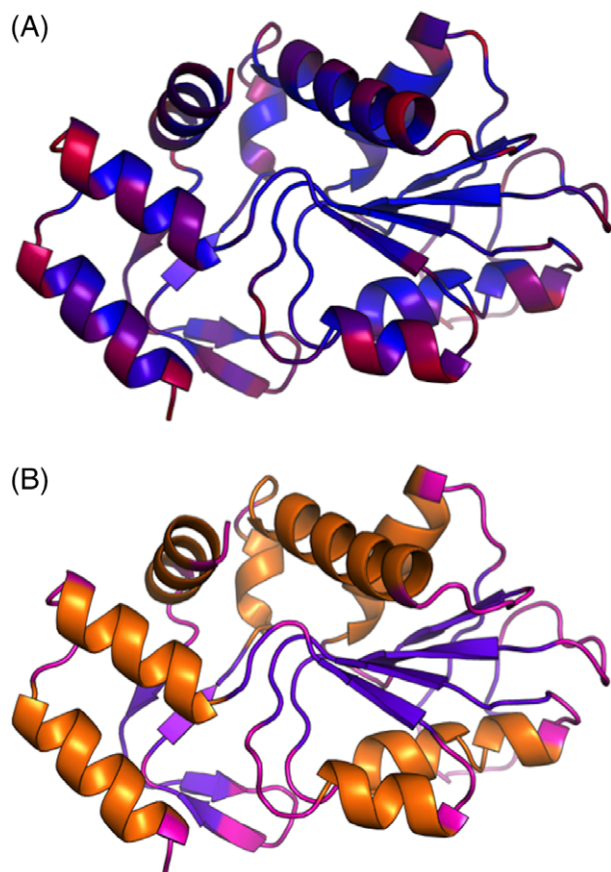


FIGURE 2 NetSurfP-2.0 predictions mapped on the OPRTase domain structure. Panel A represents the predicted ASA in a color gradient from blue (low) to red (high). Panel B represents SS3 helix, strand, and coil classes in orange, purple, and pink, respectively. The actual secondary structure of the protein is displayed in the cartoon representation of the structure. Both color codings are consistent with the web server graphical output

the tool and its successor is that the latter is trained including disordered regions. We have also tested whether this affects the accuracy of the prediction of features other than disorder in the proximity of the disordered regions. This is actually the case for Q3 and Q8, which are significantly higher for our tool when compared to a modified version of it in which the disordered regions have been completely removed from the training sequences (Supporting Information Table S4). This effect is more pronounced for the datasets that have more disordered regions (CASP12 and TS115) and less so for the CB513 dataset. Even though such residues constitute a very small portion of the total amount, these results suggest that including the

TABLE 2 Comparison of the disorder prediction on TS115, Disprot, and completely disordered proteins

Dataset	#Proteins	%Disordered (%)	Recall (%)
TS115	115	7.4	51.7
DisProt proteins	803	24.2	54.7
Completely disordered proteins	138	95.0	58.0

The #proteins, %disordered, and recall column report the number of proteins, the average content of disorder, and the recall per dataset. The recall is calculated on all the disordered residues using the default 0.5 threshold.

disordered regions in general help our model to achieve a better internal representation of the protein sequences.

3.3 | Disorder

There is nowadays a general support that disorder plays a fundamental role in the function and dynamics of proteins, and several different types of disorder have been described and annotated. Our tool is focused on protein regions with missing atomic coordinates in their solved structures, corresponding to the REMARK-465 regions in the DisEMBL annotation.³⁴ We have compared our prediction to both RaptorX and SPOT-disorder, a method developed by the same group that developed Spider-3. As customary in cases where the amount of negative data points greatly outnumbers the amount of positive data points,³⁵ we have decided to use the MCC to compare the tools. In all cases, NetSurfP-2.0 produces the most accurate results, with an average MCC of 0.65 (see Table 1). It has to be noted that this difference, though large, is not statistically significant. This might be partially due to the limited number of proteins that contain disordered regions among all the datasets, and on which the MCC and the corresponding statistical tests could be calculated. On the other hand, if we extend our evaluation of the disorder prediction by including the FPR, we see that our tool produces far less false positives than all the other tools, and in this case the difference is statistically significant.

This improved performance could be due to the specific training method we used, that includes many non-disordered proteins. We will show in the following that this, however, does not reflect in a general underprediction tendency of the tool, as exemplified in disorder-enriched proteins.

As we mentioned above, we focus on one particular definition of disorder, while many more are available. To check if our tool can produce meaningful predictions on different types of protein disorder (eg, intrinsically disordered proteins, functional disorder), we conducted a benchmark against the proteins in the DisProt database,³² a resource containing experimentally annotated disordered regions of different types. The results of this benchmark are reported in Table 2 and demonstrate that our tool performs well also on other types of disorder, and that it can also be used to identify completely disordered proteins, which were not used in its training. NetSurfP-2.0, when compared to the other tools, presents a much lower level of false positives, that is particularly evident in proteins that contain few or no disordered region at all. This, however, does not affect NetSurfP-2.0 accuracy on proteins that contain many disordered regions or that are intrinsically disordered. By comparing the recall on the TS115, the whole DisProt data set, and the subset of completely disordered proteins, we see that the tool not only performs well also on proteins enriched in disordered regions, but that the low level of false positives noticed before is not trivially linked to a general underprediction effect.

Another way to define disordered regions is by the so called “hot loops”,³⁴ that is loop regions that present high temperature factors (B factor) for their C α atoms. We tested if our prediction captures this definition by analyzing its correlation to the B factor of backbone atoms for the TS115 dataset. To compare B factors from different structures, each B factor was normalized by the average B factor of the protein chain it belongs to. We see a Spearman

TABLE 3 Accuracy of phi and psi prediction according to the secondary structure for the CASP12, TS115, and CB513 datasets

		Helices			Strands		Coils		
		G	H	I	B	E	S	T	C
CASP12	Phi	17.35	8.75	15.72	29.39	18.43	34.59	27.65	28.67
		9.69			18.89		29.91		
	Psi	30.18	15.51	16.11	33.73	20.95	60.73	33.54	52.54
		16.93			21.49		50.13		
TS115	Phi	17.74	6.85	15.49	23.69	17.60	31.85	23.33	27.72
		7.94			17.94		27.38		
	Psi	32.39	11.15	17.27	40.83	19.16	54.73	33.83	42.87
		13.13			20.35		42.89		
CB513	Phi	19.35	7.88	16.78	26.14	17.98	35.67	27.20	28.66
		9.24			18.46		29.89		
	Psi	32.66	10.90	18.14	40.22	20.18	56.00	36.35	41.87
		13.32			21.36		43.64		

In each cell we report the MAE for all residues with a specific secondary structure, either based on the 3-class (Empty cells) or 8-class (white cells) definition. For the 8 class definition: G = 3–10 helix, H = α helix, I = π helix, B = β bridge, E = extended strand, S = bend, T = h-bonded turn, C = coil.

correlation of 0.43 (Supporting Information Figure S3), confirming the ability of our model to produce a meaningful and consistent internal representation of the protein characteristics.

3.4 | Backbone dihedral angles

To complete the evaluation of our tool, we report its performance on the prediction of the backbone dihedral angles φ and ψ . We compared the results of NetSurfP-2.0 and Spider3. In this case, the 2 tools performed almost identically, with NetSurfP-2.0 producing only marginally better predictions, with no statistically significant difference. Both tools produced more accurate prediction of the φ angle compared to ψ . This is expected, given the much broader distribution of the ψ angle in the Ramachandran plot compared to the φ angle, that is almost always confined to values between -180° and -40° .

We also observed a very poor prediction accuracy for both angles in the proximity of disordered regions (Supporting Information Table S3), and, to a lower extent, in loop and coil regions (Table 3 and Supporting Information Figure S4).

3.5 | Individual versus integrated model

Thanks to the multi-task architecture and training strategy used, it is possible to predict all the features concurrently using a single model. Though this architecture improves the time optimization of our tool, this could potentially be sub-optimal with respect to accuracy. In order to test this, we trained single-output models for RSA, secondary structure (3- and 8-class), and disorder. We see (Supporting Information Table S4) that the performance of the integrated model is comparable or better than that of the individual models.

It is also interesting to notice that the hyperparameter optimization described in Methods plays an important role in the training of the integrated model: if no relative weight is assigned to the different output, we observe a small degradation of the performance of the RSA with respect to both the integrated model and the individual ones.

3.6 | Time performance

We have shown that NetSurfP-2.0 outperforms all other methods in all the tests. This is the case for both the NetSurfP-2.0 models trained using different profile generation tools, namely HH-suite and MMseqs2. Moreover, both the HH-suite and MMseqs2 models perform similarly on all datasets tested. However, they have very different running time: the runtime on a single protein sequence for the HH-suite model is approximately 2 minutes, but it scales linearly with the number of sequences. MMseqs2, conversely, is slower for small datasets, but on large datasets it provides a speed-up of up to 50 times and the ability to parallelize on multiple processors (Figure 1, panel B). Given this, NetSurfP-2.0 is implemented to use the HH-suite model for searches of less than 100 sequences, and the MMseqs2 model otherwise, thus offering a good trade-off between computation time and resource demand, without sacrificing the method's accuracy.

3.7 | Tool availability

NetSurfP-2.0 is available both as a web-server, and as an independent software (<http://www.cbs.dtu.dk/services/NetSurfP-2.0/>). The web-server version accepts up to 4000 sequences or 4 000 000 residues per job.

4 | DISCUSSION

The NetSurfP-2.0 web server provides state-of-the-art sequence-based predictions for solvent accessibility, secondary structure, disorder, and backbone geometry. Its superior performance is likely due to a few different characteristics: its architecture, its training strategy, and its data representation. The core of NetSurfP-2.0 architecture is composed by the 2 bidirectional LSTM layers. Thanks to their memory, these elements can keep track of the context around a residue better than convolutional elements. In fact, in preliminary analyses, we observed that the CNN layers that precede the LSTMs only provide a minor improvement to the overall prediction accuracy. We have also tested the effect of including a conditional random field layer downstream of the LSTM layers to

enforce a more robust grammar of the different structural elements, but this did not bring any advantage (data not shown). On the other hand, by training a multi-task weight-sharing integrated model with several structural features, we improve the accuracy of disorder with respect to models trained on individual features. This improvement likely results from a more robust and informative internal state of the system, which is extremely valuable for features, such as disorder, where only a few positives are present on average.

This integration was enabled by using improved representations of the structural data. The previous version of NetSurfP, as well as other tools, are trained only on the residues that are observed in the solved structure. In this way, the models are presented with cases that are neither physically nor biologically meaningful, such as residues divided by a disordered region, that are far apart in primary and tertiary structure but presented to the model as consecutive. In contrast, by using a recent training procedure strategy,¹⁶ we can train the model on all residues, including the disordered ones, thus increasing the accuracy of annotated features in the data and reduce the frustration during training.

The NetSurfP-2.0 framework is extremely flexible and allows to potentially include many more structural features. We have shown that the disorder prediction of our model has a fair correlation with the residues' B factor. Given this result, we believe that including the latter as an additional output for the system might actually improve the disorder prediction itself. Other possible features to be added are proline cis/trans conformation, metal binding sites, phosphorylation, glycosylation, and many others.

Having an integrated model has an effect on the accuracy of the tool, but most importantly makes it much more time-efficient. On top of that, our software uses 2 different profile creation strategies in order to achieve an even better efficiency both for small sets of sequences, and for large batches of thousands of proteins. This allows the tool to annotate a single proteome in less than a day, a very important feature in present day biology.

Thanks to its accuracy, its fast computation time, and its easy and intuitive interface, we believe that NetSurfP-2.0 will become a valuable resource that will aid researchers both with and without extensive computational knowledge to analyze and understand protein structure and function.

ACKNOWLEDGMENT

MSK and MOAS acknowledge funding from the Novo Nordisk Foundation.

ORCID

Bent Petersen  <https://orcid.org/0000-0002-2472-8317>

Paolo Marcatili  <https://orcid.org/0000-0003-2615-5695>

REFERENCES

- Anfinsen CB, Haber E, Sela M, White FH. The kinetics of formation of native ribonuclease during oxidation of the reduced polypeptide chain. *Proc Natl Acad Sci USA*. 1961;47:1309-1314.
- Cuff JA, Clamp ME, Siddiqui AS, Finlay M, Barton GJ. JPred: a consensus secondary structure prediction server. *Bioinformatics*. 1998;14(10):892-893.
- Hanson J, Yang Y, Paliwal K, Zhou Y. Improving protein disorder prediction by deep bidirectional long short-term memory recurrent neural networks. *Bioinformatics*. 2017;33(5):685-692.
- Heffernan R, Yang Y, Paliwal K, Zhou Y. Capturing non-local interactions by long short-term memory bidirectional recurrent neural networks for improving prediction of protein secondary structure, backbone angles, contact numbers and solvent accessibility. *Bioinformatics*. 2017;33(18):2842-2849.
- Jones DT. Protein secondary structure prediction based on position-specific scoring matrices. *J Mol Biol*. 1999;292(2):195-202.
- McGuffin LJ, Bryson K, Jones DT. The PSIPRED protein structure prediction server. *Bioinformatics*. 2000;16(4):404-405.
- Rost B. PHD: predicting one-dimensional protein structure by profile-based neural networks. *Methods Enzymol*. 1996;266:525-539.
- Schoes AM, Brown SD, Dodevski I, Babbitt PC. Annotation error in public databases: misannotation of molecular function in enzyme superfamilies. *PLoS Comput Biol*. 2009;5(12):e1000605.
- Kodama Y, Shumway M, Leinonen R, International Nucleotide Sequence Database Collaboration. The sequence read archive: explosive growth of sequencing data. *Nucleic Acids Res*. 2012;40(Database issue):D54-D56.
- Petersen B, Petersen T, Andersen P, Nielsen M, Lundegaard C. A generic method for assignment of reliability scores applied to solvent accessibility predictions. *BMC Struct Biol*. 2009;9(1):51.
- Almagro Armenteros JJ, Sønderby CK, Sønderby SK, Nielsen H, Winther O. DeepLoc: prediction of protein subcellular localization using deep learning. *Bioinformatics*. 2017;33(21):3387-3395.
- Esteve A, Kuprel B, Novoa RA, et al. Dermatologist-level classification of skin cancer with deep neural networks. *Nature*. 2017;542(7639):115-118.
- Min S, Lee B, Yoon S. Deep learning in bioinformatics. *Brief Bioinform*. 2017;18(5):851-869.
- Wang S, Sun S, Li Z, Zhang R, Xu J. Accurate De novo prediction of protein contact map by ultra-deep learning model. *PLoS Comput Biol*. 2017;13(1):e1005324.
- Jurtz VI, Johansen AR, Nielsen M, et al. An introduction to deep learning on biological sequence data: examples and solutions. *Bioinformatics*. 2017a;33(22):3685-3690.
- Jurtz V, Paul S, Andreatta M, Marcatili P, Peters B, Nielsen M. NetMHCpan-4.0: improved peptide-MHC class I interaction predictions integrating eluted ligand and peptide binding affinity data. *J Immunol*. 2017b;199(9):3360-3368.
- Kabsch W, Sander C. Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features. *Bio-polymers*. 1983;22(12):2577-2637.
- Linding R, Jensen LJ, Diella F, Bork P, Gibson TJ, Russell RB. Protein disorder prediction. *Structure*. 2003a;11(11):1453-1459.
- Drozdetskiy A, Cole C, Procter J, Barton GJ. JPred4: a protein secondary structure prediction server. *Nucleic Acids Res*. 2015;43(W1):W389-W394.
- Wang S, Li W, Liu S, Xu J. RaptorX-property: a web server for protein structure property prediction. *Nucleic Acids Res*. 2016a;44(W1):W430-W435.
- Wang S, Peng J, Ma J, Xu J. Protein secondary structure prediction using deep convolutional neural fields. *Sci Rep*. 2016b;6(18962):1-11.
- Berman HM, Westbrook J, Feng Z, et al. The Protein Data Bank. *Nucleic Acids Res*. 2000;28(1):235-242.
- Wang G, Dunbrack RL. PISCES: recent improvements to a PDB sequence culling server. *Nucleic Acids Res*. 2005;33(Web Server):W94-W98.
- Remmert M, Biegert A, Hauser A, Söding J. HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment. *Nat Methods*. 2011;9(2):173-175.
- Steinegger M, Söding J. MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nat Biotechnol*. 2017;35(11):1026-1028.
- Ahmad S, Gromiha MM, Sarai A. Real value prediction of solvent accessibility from amino acid sequence. *Proteins*. 2003;50(4):629-635.
- Lyons J, Dehzangi A, Heffernan R, et al. Predicting backbone C α angles and dihedrals from protein sequences by stacked sparse

- auto-encoder deep neural network. *J Comput Chem.* 2014;35(28):2040-2046.
28. Kingma DP, Ba J. (2014). Adam: a method for stochastic optimization. ArXiv:1412.6980 [Cs]. Retrieved from <http://arxiv.org/abs/1412.6980>.
 29. Yang Y, Gao J, Wang J, et al. Sixty-five years of the long march in protein secondary structure prediction: the final stretch? *Brief Bioinform.* 2018;19(3):482-494.
 30. Cuff JA, Barton GJ. Evaluation and improvement of multiple sequence methods for protein secondary structure prediction. *Proteins.* 1999;34(4):508-519.
 31. Abriata LA, Tamò GE, Monastyrskyy B, Kryshtafovych A, Dal Peraro M. Assessment of hard target modeling in CASP12 reveals an emerging role of alignment-based contact prediction methods. *Proteins Struct Funct Bioinform.* 2018;86:97-112.
 32. Piovesan D, Tabaro F, Mičetić I, et al. DisProt 7.0: a major update of the database of disordered proteins. *Nucleic Acids Res.* 2017;45(D1):D219-D227.
 33. Zhang H, Zhang T, Chen K, et al. Critical assessment of high-throughput standalone methods for secondary structure prediction. *Brief Bioinform.* 2011;12(6):672-688.
 34. Linding R, Jensen LJ, Diella F, Bork P, Gibson TJ, Russell RB. Protein disorder prediction: implications for structural proteomics. *Structure.* 2003b;11(11):1453-1459.
 35. Chicco D. Ten quick tips for machine learning in computational biology. *BioData Mining.* 2017;10(1):1-5.

SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of this article.

How to cite this article: Klausen MS, Jespersen MC, Nielsen H, et al. NetSurfP-2.0: Improved prediction of protein structural features by integrated deep learning. *Proteins.* 2019; 87:520-527. <https://doi.org/10.1002/prot.25674>