2

# 1    Accurate contact-based modelling of repeat proteins predicts the structure of

# 2    Curlin and SPW repeats.

3

4    Claudio Bassot[*] and Arne Elofsson[*]

5

6    [*] Science for Life Laboratory and Dep of Biochemistry and Biophysics, Stockholm

7    University.

8

## 9    Abstract

10    Repeat proteins are an abundant class in eukaryotic proteomes. They are involved in many

11    eukaryotic specific functions, including signalling. For many of these families, the structure is

12    not known. Recently, it has been shown that the structure of many protein families can be

13    predicted by using contact predictions from direct coupling analysis and deep learning.

14    However, their unique sequence features present in repeat proteins is a challenge for

15    contact predictions DCA-methods. Here, we show that using the deep learning-based

16    PconsC4 is more effective for predicting both intra and interunit contacts among a

17    comprehensive set of repeat proteins. In a benchmark dataset of 819 repeat proteins about

18    one third can be correctly modelled and among 51 PFAM families lacking a protein structure,

19    we produce models of five families with estimated high accuracy.

20

## 21    Author Summary

22    Repeat proteins are widespread among organisms and particularly abundant in eukaryotic

3                                                                    1

4 Structure Prediction of Repeats

5

23 proteomes. Their primary sequence present repetition in the amino acid sequences that

24 origin structures with repeated folds/domains. Although the repeated units are easy to be

25 recognized in primary sequence, often structure information are missing. Here we used

26 contact prediction for predicting the structure of repeats protein directly from their primary

27 sequences. We benchmark our method on a dataset comprehensive of all the known

28 repeated structures. We evaluate the contact predictions and the obtained models set for

29 different classes of proteins and different lengths of the target, and we benchmark the quality

30 assessment of the models on repeats proteins. Finally, we applied the methods on the

31 repeat PFAM families missing of resolved structures, five of them modelled with high

32 accuracy.

33

## **Introduction**

35 Repeat proteins contain periodic units in the primary sequence that are likely the result of

36 duplication event at the genetic level [1]. Repeat proteins emerge through replication

37 slippage [2] and double-strand break repair [3]. This protein class is present in all genomes

38 but is more frequent in eukaryotic organisms [4–6]  where they are involved in a wide range

39 of functions [7]. In particular, due to their extended structures repeat proteins often behave

40 as molecular scaffolds in protein signalling or for protein complexes as WD40 domain [8], or

41 ankyrin repeats [9,10].

42 Repeat proteins are often conserved among orthologs [4,11] while exhibiting a more

43 accelerated evolution and divergence among paralogs [11].

44 A classification of repeat proteins was proposed by Kajava [12,13] based on the length of

45 the repeat units and the tertiary structure of the repeat units. According to Kajava's

46 classification, there are five classes of repeat proteins. However, in this study, we ignore

47 class I and II because there are no available structures for class I, and class II structures are

48 folded in a coiled-coil structure easy to be predicted. Moreover, the extreme amino acid

7    Structure Prediction of Repeats

8

49    compositional bias of many of these proteins makes it very hard to find the coevolving

50    residues in these classes.

51    The dataset used in our study contains three classes of proteins divided into 20 subclasses

52    divided by their secondary structure, according to RepeatsDB [14] Fig 1. The three classes

53    are class III containing extended repeats (e.g. α and β solenoids), class IV containing closed

54    repeats structures (e.g. TIM and β barrels and β-propeller), class V where the units appear

55    as separate domains on a string. The units are also longer in class V than in other classes.

56

57    **Figure1. Repeats proteins classification.** *Representation of the repeats classes and*

58    *subclasses as classified in repeatsDB 2.0 [14]*

59

60    Class III is dominated by solenoid structures (Figure1 III.1, III.2 III.3) [13], and there is a wide

61    range in the numbers of units (from 4 to 38). Also, the length of the individual unit is widely

62    variable, e.g. β-solenoid have significantly shorter repeats compared with α and α/β solenoid

63    [13]. Two subclasses: β-trefoil/β-hairpins, anti-parallel and β-layer/β-hairpins form extended

64    beta strands without the bend typical of the solenoid.

65    Members of class IV are constrained in variability by the closed fold. Indeed despite ten

66    subclasses of different units fold the number of units go from 3 to 16, and the proteins with

67    more than ten units are rare. The length of the units is in between class III and V [13].

68    Class V has the longest units, which fold into proper domains and also a low number of units

69    with few interactions between them.

70    However, many repeat proteins lack a resolved structure or a template to perform homology

71    modelling. Residue-residue contact prediction is the most promising template free method

72    [15]. Contact prediction methods identify residues co-evolution from multiple sequence

9                                                                                    3

10    Structure Prediction of Repeats

11

73    alignment and identify the evolutionary constraints of the residues imposed by the tertiary

74    protein structure [16]. Nevertheless, repeat proteins are a difficult target for contact

75    prediction; the internal symmetry introduces artefacts in the contact map at a distance

76    corresponding to the repeated units [17].

77    Here, we benchmark the deep-learning-based contacts prediction program PconsC4 [18]

78    against the GaussDCA [19] on a comprehensive dataset generated from RepeatsDB [14].

79    The predicted contacts were then used as constraints to generate proteins model, and their

80    quality was then tested by Pcons [20]. On the base of the benchmark, we propose models

81    for the protein structures of PFAM protein families missing of resolved structures.

82

83    **Results and Discussion**

84    **General contact prediction analysis in repeat proteins**

85    To assess the quality of the contacts predictions among repeat protein classes, we generate

86    a dataset of proteins, clustering at 40% of identity, the reviewed entries of RepeatsDB [14].

87    For each repeats region present in the dataset we extract the sequence of a representative

88    repeat unit and a pair of repeats, obtaining in this way three datasets: i) a single unit

89    datasets; ii) a double unit datasets; iii) complete region datasets.

90    For all the three sets of proteins, multiple sequence alignments (MSA) and secondary

91    structure predictions were generated. Subsequently using the MSA as input for PconsC4

92    and GaussDCA [19] contacts were predicted for each family. The performance of the contact

93    predictions was evaluated for each subclass separately. As expected, PconsC4 over-

94    perform GaussDCA in all the three sets and all the classes of repeat proteins, Figure 2.

95

96    *Figure2.* **Precision of contact predictions.** *Positive Predicted Value (PPV) for the*

13 Structure Prediction of Repeats

14

97 *GaussDCA (red) and Pconsc4 (Blue) contact prediction for each subclass. In light colour the*

98 *single unit dataset, in medium colour the double units dataset, and in dark colour the*

99 *complete region dataset.*

100

101 Here, it should be remembered that PconsC4 use the GaussDCA prediction as an input for

102 the U-net [32] that learn to recognize specific contacts patterns [18].

103 In general, the predictions for the full length regions (darker colors in Fig. 2) give better

104 results than split the proteins into units but with some exceptions. In particular in class V,

105 that is composed by bigger units forming repeats of the *"beds on a string"* type, the splitting

106 in units may help, especially in some subclasses, to reach better contacts prediction

107 performance as discussed later.

108 Furthermore, PconsC4 appear efficient in removing the DCA repeats artefacts compared

109 with GaussDCA. In Fig. 3 are shown some contact maps examples. In the GaussDCA

110 predictions are evident the periodic artefacts of wrong predictions (red dots) forming

111 perpendicular lines. These appear to be contacts between equivalent positions in the repeat

112 unit.

113

114 *Figure 3.* **GaussDCA and PconsC4 contact maps.** Co*ntact map for a prediction with a)*

115 *GaussDCA b) PconsC4. In grey, the real contacts from the structure, in green the corrected*

116 *predicted value, in red the false predicted value.*

117

118 Finally, it is well known that the quality of the prediction is directly correlated with the number

119 of sequences in the starting MSA, especially for the DCA methods [18]. The same trend is

120 confirmed among protein repeats, Fig. 4, where the repeats with a smaller MSA are

15 5

16    Structure Prediction of Repeats

17

121   predicted with lower PPV. PconsC4 and GaussDCA show the same pattern in the average

122   PPV except for an increase of the PPV for PconsC4 with MSA with a Neff Higher than 12.

123

124   *Figure 4.* **PPV versus Neff.** *Positively Predicted Value for GaussDCA in red and PconsC4 in*

125   *Blue on the Neff value (number of effective sequences length weighted with length).*

126

127   **Differences among repeat classes in contacts prediction**

128   Fig. 2 shows variations in the percentage of correct contacts among different protein repeat

129   classes and subclasses, to clarify the origin of these differences, we investigated more in-

130   depth the origin of the predicted contacts.

131   One central aspect that affects the difficulty of prediction is due to the pattern of contacts

132   [33]. In general, contacts that are parts of larger interaction areas are better predicted as well

133   as interactions between residues that are close in the sequence. A comparison between the

134   intra-unit and inter-unit contacts are shown in Fig 5a. Here, we obtained the number of intra

135   and inter-unit contacts from the PDB structures and we selected the same number of intra

136   and inter units from the contact predictions. The PPV was finally calculated as the number of

137   correct contacts over the number of the selected contacts.

138   On average the intra-units contacts are predicted with higher accuracy than the inter-unit

139   one, but this is not true for all protein classes. This behaviour is due to significant differences

140   of the units structures among the classes: in class III the unit are short, and the residues

141   form contacts mostly with the neighbour units; in class V, on the contrary, the units are long,

142   folded in independent domains and the contacts are predominantly inside the units with few

143   inter-unit contacts; class IV is halfway between class III and V. The inter units contacts of

144   class III and partially of class IV results easier do be predicted then class V ones, because

145   they form clearer patterns in contact maps. On the contrary, the intra-unit contacts of class V

18                         6

19    Structure Prediction of Repeats
20

146    are predicted better than class III and IV for the same reason. We plot the PPV versus the

147    ratio of the inter-unit contacts over the total number of contacts of each protein. The PPV

148    show an inverse relation with the ratio of inter-unit contact of the protein (Figure 5.b,c,d).


149    The inter-units PPV is low for the proteins with an inter-contact ratio lower than 20%

150    constituting class V. Figure 5c shows the lowest inter-units contacts PPV while the PPV

151    between inter- and intra- contacts invert the trends at a ratio of 80%. This switch

152    corresponds to solenoid structures and TIM Barrel that have a ratio between 80%-100%

153    larger interaction surfaces between different units than inside a single unit.


154


155    **Figure 5. Predicted contacts analysis.** *Positive Predicted Value (PPV) obtained by*

156    *PconsC4 for different types of contacts. a) Examples of inter- and intra- unit contacts. b) In*

157    *red, the PPV for intra-units contacts in blue PPV for inter-units contact. c) Repeats*

158    *subclasses. In red, the PPV for intra-units contacts in blue PPV for inter-units contact, colors*

159    *and shapes in the scatter plot indicate different protein subclasses. d) Secondary structure.*

160    *In red, the overall PPV, in blue, the α-helical subclasses, in green, the α-helix/β-strand*

161    *subclasses, and in orange, the β-strand subclasses.*


162


163    In Figure 5d, we divided the proteins into their secondary structure class. Proteins

164    subclasses containing only β-strand or α-helix/β-strand appear easier to predict. The plots

165    show a steep decrease in the PPV values around the ratio of 50% helix for both intra- and

166    inter-units contacts. This is due to the α-helix subclasses component. α-helix are harder to

167    predict because they produce a less clear contact pattern compared with β-strand.


168


169    **Protein model generation and quality assessment**

22    Structure Prediction of Repeats

23

170    Proteins models were generated using CONFOLD [28] starting from the contact predictions

171    of PconsC4 and the PSIpred secondary structure as constraints. In Fig. 6 we compare the

172    TM-score between the first model ranked by CONFOLD and the corresponding PDB protein

173    structure.

174

175    **Figure 6. Protein model quality.** *a) TM-score in all the subfamilies. In sea-green the single*

176    *unit prediction, in blue the double units prediction, in red the complete region prediction. b)*

177    *TM-score of the  subfamilies of class V. In green the single unit prediction and in brown the*

178    *prediction of that unit when the entire region is modelled.*

179

180    Although the best contact predictions were, on average, obtained with the complete regions,

181    still splitting the structure lead in some cases to a better model; this is true in particular for

182    the "Beads on a string" class V,   but single-unit models are also useful in the "propeller"

183    subclasses class IV: IV.4 β propeller,   IV.8 α/β propeller and   IV.5 α/β prism. Moreover

184    modelling a couple of units lead to the best models in two subclass III.3  α solenoid  IV.10

185    and aligned prism. All these subclasses except α-solenoid have a low ratio of inter-units

186    contacts (below 50%) Fig. 5b, however α-solenoid where the complete protein reaches in

187    some cases a length of 1000 residues. Moreover, the bend of the protein is very difficult to

188    predict, and the models result in a series of straight helices.

189    It is questionable if the lower quality of the models of the complete region is due to a general

190    decrease in the performance or only to the impossibility to model the correct interaction

191    among different domains. To answer this question, we analysed more in deeper class V,

192    where the decrease in the performance is most evident. We extract from the "complete

193    region model" the same units and the units previously modelled as single and double units

194    Fig. 6b. Interestingly, even the single units and the double units extracted from the complete

195    region modelling have a lower or similar accuracy compared with the single and double units

24                                                          8

25 Structure Prediction of Repeats

26

196 modelled alone, Figure 6b. This is observed is regardless of the quality of the prediction of

197 the contacts in the complete region prediction, Fig. 2, suggesting that the poor performance

198 is not only due to the more difficult prediction of the interdomain contacts but also due to a

199 limitation of the modelling of longer proteins.

200

201 In order to evaluate the model quality, we plot the TM-scores of the models against the score

202 obtained from the quality assessment method Pcons [20], Fig. 7. In light of this result, we

203 consider the models of a complete repeats region reasonably correct when they reach a

204 Pcons score of 0.4. The complete dataset with Pcons prediction is reported in Table S1.

205

206 **Figure 7. TM-score versus Pcons-score.** *TM-score versus Pcons-score for complete*

207 *region models.*

208

209 **Modelling of repeat proteins without resolved structures**

210 In order to predict the structure of new repeats families, we selected 51 PFAM repeats family

211 without resolved structure. A representative sequence of each family was run against

212 Uniclust30 with HHpred, and the resulting MSA was used to predict the contact map that

213 was used together with the PSIpred prediction as constraints to generate the models.

214 All the models were evaluated with Pcons, but only five of them reach a Pcons score higher

215 than 0.4. These are the PFAM family; MORN 2, SPW, Curlin rpt, RTTN N, RHS repeat,

216 Table 1 (In Supplementary the target/template alignments).

217 In order to further prove the reliability of these models and perform a more comprehensive

218 protein modelling approach, we associated homology modelling and the contact-based

219 modelling approach. For three out of five proteins, HHsearch returned a highly reliable

27 9

28   Structure Prediction of Repeats

29

220   template, Table 1.

221

| PFAM Family | Rappresentative sequence (Uniprot ID) | Pcons score | Template with seq. coverage > 70% (PDB ID) | HHsearch probability | TM score contact model/homology model |
|---|---|---|---|---|---|
| **MORN 2 (PF07661)** | **Q8RH85** | **0.711** | **1MUF_A** | **99.37** | **0.5003** |
| SPW | A0A2A3HD64 | 0.674 | 5EQC_A | 29.35 | / |
| Curlin rpt | Q8EIH3 | 0.576 | 2N59_A | 1.97 | / |
| **RTTN N (PF14726)** | **W5P499** | **0.490** | **4U2X_E** | **94.03** | **0.3529** |
| **RHS repeat** | **A0A1G0MXS8** | **0.407** | **5KIS_B** | **99.46** | **0.5823** |

222

223   In Fig. 8, the superimposition between homology modelling and contact based model is

224   shown. In all three the protein family there is a substantial agreement between the two

225   approaches. MORN 2 family contact-based and homology model are in agreement except

226   for loops and the bend of the central beta-strand.

227

228   *Figure 8. High quality protein models. a) Superimposition between the contact-based*

229   *models and the Homology Model performed with Chimera [34] and their respective TM-*

230   *score. In red, the contact-based models and in light blue Homology models. b) Protein*

231   *model of SPW family in the membrane (light brown). On the left in blue and red the two*

232   *repeated units on the right in red the SPW motif. c ) Protein model of Curlin repeats, in blue*

30                                          10

31    Structure Prediction of Repeats

32

233    *and red the repeated units.*

234

235    The RTTN N is the family showing the lowest TM-score between the two models mostly due

236    to a different rearrangement of the firsts three alpha-helices. Has to be mentioned, however,

237    that despite a high probability score, the identity between the target and the best template is

238    only 7% (Figure S1b) making hard to determine which is the best model.

239    In RHS repeat family, the score between the contact-based models and the homology model

240    share a TM-score of 0.58. Only the N terminal is modelled in a different with an extra beta-

241    strand in the contact-based model and an alpha helix in the template-based modelling.

242    However, we argue that in this case, the contact-based modelling overperform the Homology

243    model; indeed the contact prediction mode is in agreement with the secondary structure

244    prediction that predicts an N-terminal Beta strand (Figure S1c).

245    The remaining two PFAM families do not have suitable templates, and contact-based

246    modelling is the best suitable method for model them.

247

248    **SPW family**

249    According to the PFAM database, the SPW family is present in Bacteria and Archaea in one

250    or two units, and in a few cases in association with a Vitamin K epoxide reductase or NAD-

251    dependent epimerase/dehydratase domain. Each repeated unit is formed by two

252    transmembrane alpha-helices and is characterized by an SPW motive [35]. According to our

253    model, the repeated motifs is buried in the membrane symmetrically located close to the

254    extracellular side, Fig. *8b*. PFAM architectures show many proteins with only a single SPW

255    motif however a more careful analysis of these sequences shows that in many cases they

256    contain a second degenerate SPW unit before or after the one identified where however the

257    proline residue is conserved (Figure S2).

33                                        11

34    Structure Prediction of Repeats

35

258    The Tryptophan is on the outer side of the protein facing the bilayer while the proline is on

259    the inner side of the protein promoting the formation of a kink in the transmembrane helix

260    [36]. The motif "SP" in particular, increase the bending effect of proline significantly due to

261    their hydrogen bond pattern [37], indeed due to the structural propriety, the motif is relatively

262    rare in membrane proteins [37].

263

264    **Curlin repeats family**

265    Our model results in a β-solenoid structure, Fig. *8c* DeBenedictis et al. in 2017 presented

266    and discussed ab initio models for the Curlin repeats family members CsgA and CsgB [38],

267    their best models is in agreement with our model (a direct comparison is difficult as the

268    coordinates is not available of their model). The model is furthermore confirmed by the

269    partial structure of the repeat units of CsgA published by Perov et al. [39] where they

270    crystallize in parallel β-sheets with individual units situated perpendicular to the fibril axis

271    (corresponding PDB IDs are 6G8C, 6G8D, 6G8E).

272

273    **Conclusion.**

274    The modelling of the unknown PFAM families was challenging. Only 10% of the datasets

275    had a Pcons score equal or higher to 0.4; compared to 21% in the benchmark dataset.

276    However, the differences between the two datasets have to be taken into account. It is

277    known that a smaller MSA affects the prediction of contacts and known structures are biased

278    towards the larger family [40] Indeed our "Unknown protein families" dataset shows a

279    significant lower Neff score compared with the PDB benchmark set, Figure 9a. Moreover, in

280    the Unknown protein set, there are more eukaryotic-specific protein families (Fig. 9b).

281

37    Structure Prediction of Repeats

38

282    ***Figure9. Datasets comparisons.*** *a) Neff score comparison between the two datasets. b)*

283    *The variation in the  membership to the three domains of life between the PFAM families of*

284    *the "Unresolved Proteins Dataset" and the "PDB dataset".*

285

286    Despite the significant improvement brought by deep-learning in contact prediction, there is

287    still room for improvement. The prediction of inter-domain contacts accuracy is often lower

288    than the intra-units one and the development of a model trained explicitly on repeats protein

289    datasets might improve the result. Furthermore, the folding part of the pipeline is a limiting

290    step, in particular for long proteins.

291    In our study, we performed a comprehensive coevolution analysis on repeat protein families,

292    and we show that PconsC4 contact-predictions method overcomes the traditional difficulties

293    of DCA methods for this class of proteins. We investigated the modelling of repeat units, and

294    we provided a "titration curve" for Pcons score for repeat proteins. Finally, we test our

295    pipeline on PFAM families without protein structures showing its usefulness in providing new

296    structural information.

297

298    **Materials and Methods**

299    **Datasets generation**

300    The repeat protein dataset was generated starting from the 3585 reviewed entries in

301    RepeatsDB [14,21], http://protein.bio.unipd.it/repeatsdb-lite/dataset. The proteins of class I

302    and II were removed, and then the dataset was homology reduced using CD-HIT [22] at 40%

303    identity resulting in 819 repeats regions. From this "complete region dataset" two others

304    datasets were generated: I) A "single unit" dataset with one repeat unit for each region; II) A

305    "double unit" dataset with a pair of units per each repeat region. In the two derived datasets,

306    the representative units were selected, avoiding or at least minimizing, the presence of

39                                                          13

40    Structure Prediction of Repeats

41

307    insertions.

308    The non-resolved repeats protein family dataset was generated, collecting all the repeat

309    proteins families with missing structural information present in PFAM [23] in May 2019 and

310    removing the domains with a significant overlap with the disorder prediction. It results in 51

311    protein families. The representative sequence for each family of repeat was chosen for

312    matching these criteria: 1) select the most common architecture; 2) Include when possible at

313    least three repeat units.

314

315    **Multiple sequence alignment (MSA)**

316    The multiple sequence alignments (MSA) were carried out using HHblits [24] with an E-value

317    cutoff of 0.001 against the Uniclust30_2017_04 database [25]. The number of effective

318    sequences of the alignment, expressed as Neff-score, was calculated by HHblits and used

319    for subsequent analysis.

320

321    **Contact prediction and models generation**

322    The protein models were generated following the PconsFold2 protocol of [26]. The

323    secondary structure of the repeat regions was predicted by PSIpred [27]. Protein contacts

324    were calculated with PconsC4 [18] and together with the secondary structure predictions

325    were used as input for Confold [28]. The modelling was run using the top scoring 1.5 L

326    contracts where L is the length of the modelled regions and the two-stage modelling.

327

328    **Contacts analysis**

329    A protein contact was defined as two residues having a beta carbon distance equal or lower

330    than 8Å in the PDB structure and farther than 5 residues in the sequence. Using this

42                                        14

43   Structure Prediction of Repeats

44

331   definition, we assess the number of correctly predicted contacts (the Positively Predicted

332   value (PPV)) taking into account the top-scoring 1.5 L contracts.

333   In the intra/inter unit contacts analysis, the predicted contacts of each protein were divided

334   between i) intra-unit contacts, if between residues inside the same unit; ii) inter-units if the

335   residues are in different repeat units. The units mapping was taken from the RepeatsDB

336   database [14]. In this analysis, we calculate the number of intra- and inter-unit contacts

337   existing in the PDB structure, and we selected the same number of intra- and inter-units

338   predictions. The PPV was then calculated as the number of correct contacts over the

339   number of the selected contacts.

340   **Homology modelling**

341   Templates for homology modelling were searched by HHsearch [29] using the HHpred web-

342   server with default settings on PDB_mmCIF70_3_Aug database. Subsequently, the models

343   were generated by HHpred [30].

344   **Protein models analysis**

345   The model quality was assessed using Pcons [20]. We download and installed Pcons. With

346   the option -d we predicted the quality among the model in the stage2 folder generated by

347   Confold. Pcons uses a clustering method, and the score is simply the average structural

348   similarity to all models, as measured by the S-score.

349   The TM-score was calculated using TMalign [31]. To ensure that the protein structure and

350   the model were properly aligned the option -I was used, providing a local protein alignment

351   for the two sequences.

352

353   **Bibliography**

354   1.   Heringa J. Detection of internal repeats: how common are they? Curr Opin Struct Biol. 1998;8:

45                                                              15

46  Structure Prediction of Repeats
47

355      338–345.

356  2.  Strand M, Prolla TA, Liskay RM, Petes TD. Destabilization of tracts of simple repetitive DNA in
357      yeast by mutations affecting DNA mismatch repair. Nature. 1993;365: 274–276.

358  3.  Pâques F, Leung W-Y, Haber JE. Expansions and Contractions in a Tandem Repeat Induced by
359      Double-Strand Break Repair. Molecular and Cellular Biology. 1998. pp. 2045–2054. doi:10.1128/
360      mcb.18.4.2045

361  4.  Schaper E, Gascuel O, Anisimova M. Deep conservation of human protein tandem repeats within
362      the eukaryotes. Mol Biol Evol. 2014;31: 1132–1148.

363  5.  E.M. Marcotte, M. Pellegrini, T.O. Yeates, D. Eisenberg. A census of protein repeats. J Mol Biol.
364      1999;293: 151–160.

365  6.  Björklund AK, Ekman D, Elofsson A. Expansion of protein domain repeats. PLoS Comput Biol.
366      2006;2: e114.

367  7.  Andrade MA, Perez-Iratxeta C, Ponting CP. Protein Repeats: Structures, Functions, and
368      Evolution. Journal of Structural Biology. 2001. pp. 117–131. doi:10.1006/jsbi.2001.4392

369  8.  Stirnimann CU, Petsalaki E, Russell RB, Müller CW. WD40 proteins propel cellular networks.
370      Trends Biochem Sci. 2010;35: 565–574.

371  9.  Li J, Mahajan A, Tsai M-D. Ankyrin repeat: a unique motif mediating protein-protein interactions.
372      Biochemistry. 2006;45: 15168–15178.

373  10. Mosavi LK, Cammett TJ, Desrosiers DC, Peng Z-Y. The ankyrin repeat as molecular architecture
374      for protein recognition. Protein Sci. 2004;13: 1435–1448.

375  11. Persi E, Wolf YI, Koonin EV. Positive and strongly relaxed purifying selection drive the evolution
376      of repeats in proteins. Nat Commun. 2016;7: 13570.

377  12. Kajava AV. Review: Proteins with Repeated Sequence—Structural Prediction and Modeling.
378      Journal of Structural Biology. 2001. pp. 132–144. doi:10.1006/jsbi.2000.4328

379  13. Kajava AV. Tandem repeats in proteins: From sequence to structure. Journal of Structural
380      Biology. 2012. pp. 279–288. doi:10.1016/j.jsb.2011.08.009

381  14. Paladin L, Hirsh L, Piovesan D, Andrade-Navarro MA, Kajava AV, Tosatto SCE. RepeatsDB 2.0:
382      improved annotation, classification, search and visualization of repeat protein structures. Nucleic
383      Acids Res. 2017;45: 3613.

384  15. Abriata LA, Tamò GE, Monastyrskyy B, Kryshtafovych A, Dal Peraro M. Assessment of hard
385      target modeling in CASP12 reveals an emerging role of alignment-based contact prediction
386      methods. Proteins. 2018;86 Suppl 1: 97–112.

387  16. Pazos F, Helmer-Citterich M, Ausiello G, Valencia A. Correlated mutations contain information
388      about protein-protein interaction. J Mol Biol. 1997;271: 511–523.

389  17. Espada R, Parra RG, Mora T, Walczak AM, Ferreiro DU. Capturing coevolutionary signals
390      inrepeat proteins. BMC Bioinformatics. 2015;16: 207.

391  18. Michel M, Hurtado DM, Elofsson A. PconsC4: fast, accurate, and hassle-free contact predictions.
392      Bioinformatics. 2018. doi:10.1093/bioinformatics/bty1036

393  19. Baldassi C, Zamparo M, Feinauer C, Procaccini A, Zecchina R, Weigt M, et al. Fast and accurate
394      multivariate Gaussian modeling of protein families: predicting residue contacts and protein-
395      interaction partners. PLoS One. 2014;9: e92721.

396  20. Lundström J, Rychlewski L, Bujnicki J, Elofsson A. Pcons: a neural-network-based consensus
397      predictor that improves fold recognition. Protein Sci. 2001;10: 2354–2362.

49    Structure Prediction of Repeats

50

398   21.   Hirsh L, Paladin L, Piovesan D, Tosatto SCE. RepeatsDB-lite: a web server for unit annotation of
399        tandem repeat proteins. Nucleic Acids Res. 2018;46: W402–W407.

400   22.   Li W, Godzik A. Cd-hit: a fast program for clustering and comparing large sets of protein or
401        nucleotide sequences. Bioinformatics. 2006. pp. 1658–1659. doi:10.1093/bioinformatics/btl158

402   23.   El-Gebali S, Mistry J, Bateman A, Eddy SR, Luciani A, Potter SC, et al. The Pfam protein families
403        database in 2019. Nucleic Acids Res. 2019;47: D427–D432.

404   24.   Remmert M, Biegert A, Hauser A, Söding J. HHblits: lightning-fast iterative protein sequence
405        searching by HMM-HMM alignment. Nature Methods. 2012. pp. 173–175.
406        doi:10.1038/nmeth.1818

407   25.   Mirdita M, von den Driesch L, Galiez C, Martin MJ, Söding J, Steinegger M. Uniclust databases of
408        clustered and deeply annotated protein sequences and alignments. Nucleic Acids Res. 2017;45:
409        D170–D176.

410   26.   Bassot C, Menendez Hurtado D, Elofsson A. Using PconsC4 and PconsFold2 to Predict Protein
411        Structure. Curr Protoc Bioinformatics. 2019; e75.

412   27.   McGuffin LJ, Bryson K, Jones DT. The PSIPRED protein structure prediction server.
413        Bioinformatics. 2000. pp. 404–405. doi:10.1093/bioinformatics/16.4.404

414   28.   Adhikari B, Bhattacharya D, Cao R, Cheng J. CONFOLD: Residue-residue contact-guided ab
415        initio protein folding. Proteins. 2015;83: 1436–1449.

416   29.   Soding J, Biegert A, Lupas AN. The HHpred interactive server for protein homology detection and
417        structure prediction. Nucleic Acids Research. 2005. pp. W244–W248. doi:10.1093/nar/gki408

418   30.   Zimmermann L, Stephens A, Nam S-Z, Rau D, Kübler J, Lozajic M, et al. A Completely
419        Reimplemented MPI Bioinformatics Toolkit with a New HHpred Server at its Core. J Mol Biol.
420        2018;430: 2237–2243.

421   31.   Zhang Y, Skolnick J. TM-align: a protein structure alignment algorithm based on the TM-score.
422        Nucleic Acids Res. 2005;33: 2302–2309.

423   32.   Ronneberger O, Fischer P, Brox T. U-Net: Convolutional Networks for Biomedical Image
424        Segmentation. Lecture Notes in Computer Science. 2015. pp. 234–241. doi:10.1007/978-3-319-
425        24574-4_28

426   33.   Skwark MJ, Raimondi D, Michel M, Elofsson A. Improved contact predictions using the
427        recognition of protein like contact patterns. PLoS Comput Biol. 2014;10: e1003889.

428   34.   Pettersen EF, Goddard TD, Huang CC, Couch GS, Greenblatt DM, Meng EC, et al. UCSF
429        Chimera--a visualization system for exploratory research and analysis. J Comput Chem. 2004;25:
430        1605–1612.

431   35.   Yeats C, Bentley S, Bateman A. New knowledge from old: in silico discovery of novel protein
432        domains in Streptomyces coelicolor. BMC Microbiol. 2003;3: 3.

433   36.   von Heijne G. Proline kinks in transmembrane alpha-helices. J Mol Biol. 1991;218: 499–503.

434   37.   Deupi X, Olivella M, Govaerts C, Ballesteros JA, Campillo M, Pardo L. Ser and Thr Residues
435        Modulate the Conformation of Pro-Kinked Transmembrane α-Helices. Biophysical Journal. 2004.
436        pp. 105–115. doi:10.1016/s0006-3495(04)74088-6

437   38.   DeBenedictis EP, Ma D, Keten S. Structural predictions for curli amyloid fibril subunits CsgA and
438        CsgB. RSC Adv. 2017;7: 48102–48112.

439   39.   Perov S, Lidor O, Salinas N, Golan N, Tayeb-Fligelman E, Deshmukh M, et al. Structural Insights
440        into Curli CsgA Cross-β Fibril Architecture Inspired Repurposing of Anti-amyloid Compounds as
441        Anti-biofilm Agents. doi:10.1101/493668

52    Structure Prediction of Repeats

53

442    40. Orlando G, Raimondi D, Vranken WF. Observation selection bias in contact prediction and its
443         implications for structural bioinformatics. Sci Rep. 2016;6: 36679.

444

445

54                                              18

55   Structure Prediction of Repeats

56

446  **Supporting Information**

447

448  **Table S1. Unknown protein family dataset.** In the columns are reported respectively: the UniProt ID

449  of the modelled sequence, the PFAM family, the Pcons score.

450

451  **Figure S1 Target/template alignments.** Target/template alignments for the homology modelling.

452

453  **Figure S2 Amino Acid frequency of the single domain architecture sequences**. From the logo is

454  possible recognize two SPW domains, one of them degenerated (in particular the first Serine in the

455  second motif) that is not recognized by PFAM.

456

# CLASS III



III.1 β-solenoid

III.2 α/β solenoid

III.3 α-solenoid

III.4  β trefoil / β hairpins

III.5  anti-parallel β layer / β hairpins

# CLASS IV



IV.1 TIM-barrel

IV.2 β-barrel / β hairpins

IV.3 β-trefoil

IV.4 β-propeller

IV.5 α/β prism

IV.6 α-barrel

IV.7 α/β barrel

IV.8 α/β propeller

IV.9 α/β trefoil

IV.10 aligned prism

# CLASS V



V.1 α-beads

V.2 β-beads

V.3 α/β-beads

V.4 β sandwich beads

V.5 α/β sandwich beads

Legend:

Single unit Pconsc4

Single unit GDCA

Double units Pconsc4

Double units GDCA

Complete PconsC4

Complete GDCA

GaussDCA
1wg0_A.0 (PDB: 1wg0)
PPV = 0.37

PconsC4
1wg0_A.0 (PDB: 1wg0)
PPV = 0.51

GaussDCA
2prt_A.0 (PDB: 2prt)
PPV = 0.16

PconsC4
2prt_A.0 (PDB: 2prt)
PPV = 0.58

GaussDCA
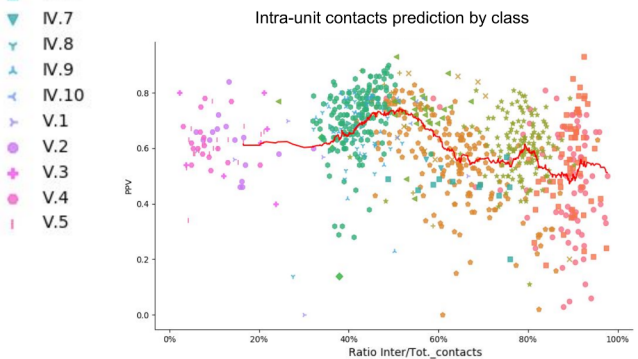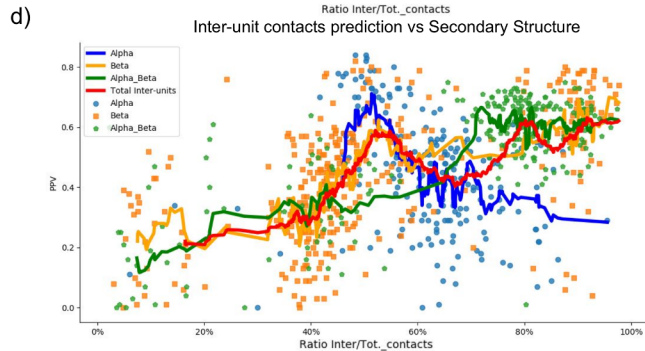4hbd_A.0 (PDB: 4hbd)
PPV = 0.34

PconsC4
4hbd_A.0 (PDB: 4hbd)
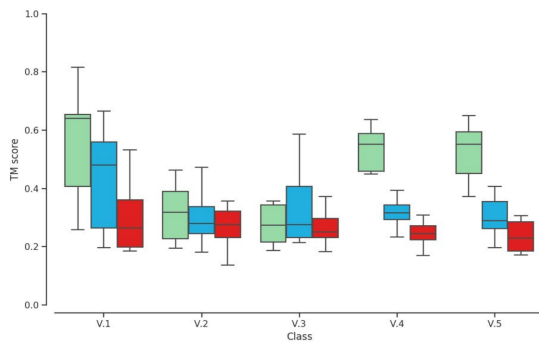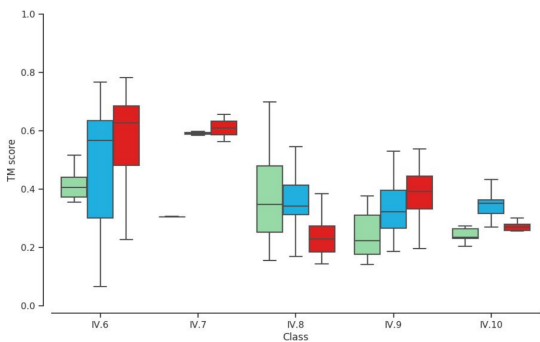PPV = 0.44

GaussDCA
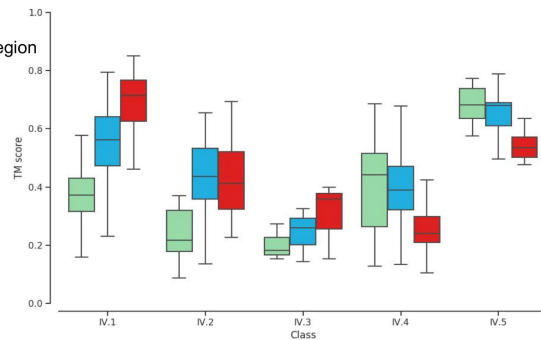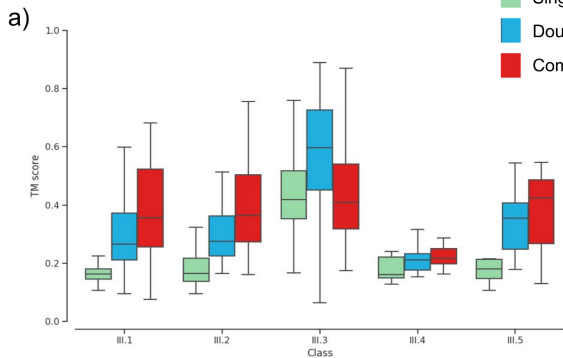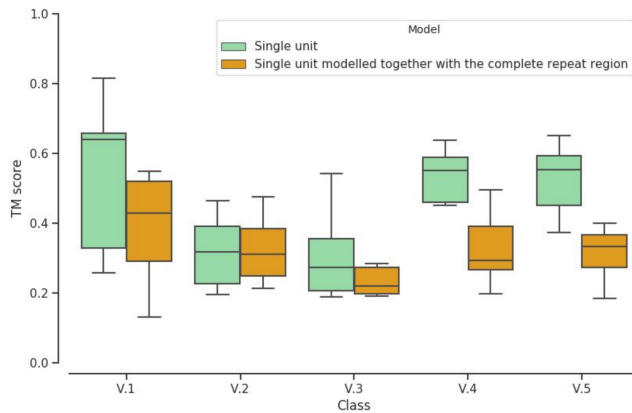4l3f_D.0 (PDB: 4l3f)
PPV = 0.51

PconsC4
4l3f_D.0 (PDB: 4l3f)
PPV = 0.69

a) Inter unit contact

Intra unit contact

b) Inter and intra unit contacts prediction vs Inter unit contacts ratio

- Inter
- Intra
- PPV
- PPV

PPV

Ratio Inter/Tot._contacts

c) Inter-unit contacts prediction by class

class

| | |
|---|---|
| ● | III.1 |
| ■ | III.2 |
| ⬠ | III.3 |
| ✕ | III.4 |
| ✚ | III.5 |
| ★ | IV.1 |
| ◄ | IV.2 |
| ◆ | IV.3 |
| ● | IV.4 |
| ▶ | IV.5 |
| ■ | IV.6 |
| ▼ | IV.7 |
| Y | IV.8 |
| ▲ | IV.9 |
| ◄ | IV.10 |
| ▶ | V.1 |
| ● | V.2 |
| ✚ | V.3 |
| ● | V.4 |
| ❙ | V.5 |

PPV

Ratio Inter/Tot._contacts

Intra-unit contacts prediction by class

PPV

Ratio Inter/Tot._contacts

d) Inter-unit contacts prediction vs Secondary Structure

- Alpha
- Beta
- Alpha_Beta
- Total Inter-units
- Alpha
- Beta
- Alpha_Beta

PPV

Ratio Inter/Tot._contacts

Intra-unit contacts prediction vs Secondary Structure

- Alpha
- Beta
- Alpha_Beta
- Total Inta-units
- Alpha
- Beta
- Alpha_Beta

PPV

Ratio Inter/Tot._contacts

a)

b)

a)

MORN2    RTTN N    RHS repeat

**TMscore**  0.5003    0.3529    0.5823

Contact Based Model
Homology Model

b)

90°C

c)

90°C

a)

Neff

14

12

10

8

6

4

2

PDB dataset    Unresolved proteins

Set

b)

Bacteria

24.0%

46.0%    Bacteria and Eukaryota

30.0%

Eukaryota

Bacteria

10.5%

Eukaryota

17.3%

72.3%

Bacteria and Eukaryota

## a) MORN 2, PF07661 , target Q8RH85

**1MUF_A** SET9 (2.1.1.43); SET domain, histone lysine methyltransferase; HET: MSE; 2.26A {Homo sapiens} SCOP: b.76.2.1, b.85.7.1
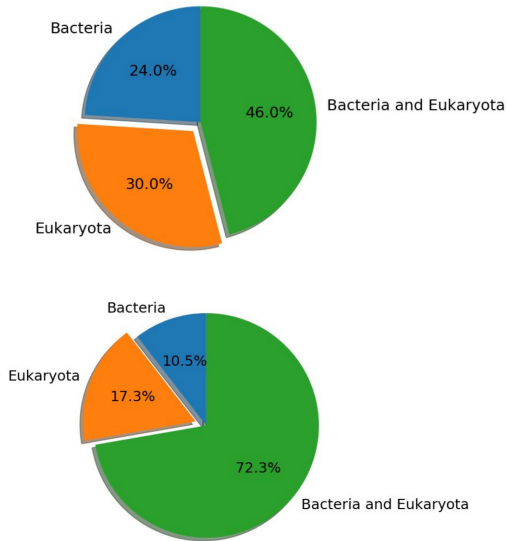
Probability: 99.37%,  E-value: 1.3e-13,  Score: 73.11,  Aligned cols: 70,  Identities: 21%,  Similarity: 0.357,

```
Q ss_pred         CCceEEEEcCCCcEEEEEEeeCCEEEEeEEEEeecCCeEEEEEeeeCCeEceeEEE-EECCCCCEEEEEEEeeC
Q Q8RH85       1  QVGVEKSYYESGELLSECSYKNGKMDGIAKIYYQNGQVEIEDPYKNGERNGVIK~VYDENGKLVRQATFKN   70 (70)
Q Consensus    1  ~~G~~~~~~g~~~~~~~~~~~g~~~~G~~~~~~~~~~~~~~~~~~~~~~~g~~~~~~~~~g~~~~~~~~~   70 (70)
                  ++|.+..|+++|.+..++.|.++.++|.++.|+++|.+.....|.++.++|.++ .|+++|.+.+..|++
Q Consensus   11  ~~G~~~~~~G~~~g~~~~g~~~G~~~~~~~~g~~~~~~~~~~~~~~~~g~~~~~~~~g~~~~~~~~~~~   81 (257)
Q Q8RH85      11  LNGPAQEYDTDGRLIFKGQYKDNIRHGVCWIYYPDGGSLVGEVNEDGEMTGEKIAYVYPDERTALYGKFID   81 (257)
T ss_dssp         EEEEEEEECTTCcEEEEEEEETTEEEEEEEECTTSCEEEEECCTTSCSCEEEEEEECTTSSEEEEEEEET
T ss_pred         EEeeEEEECCCCCEEEEEEEECCeEEeeEEEEeCCCCeEEEEEcCCCceeceEEEEEECCCCCEEEEEEEEC
```

## b) RTTN N, PF14726 , target W5P499

**4U2X_D** eVP24, KPNA5C; eVP24, importin alpha6, immune antagonist; 3.153A {Zaire ebolavirus}; Related PDB entries: 4U2X_F 4U2X_E

Probability: 94.16%,  E-value: 0.22,  Score: 31.77,  Aligned cols: 97,  Identities: 7%,  Similarity: 0.015,

```
Q ss_pred         CHHHHHHHHHHHHcccchHHHHhcCHHHHHHHHHhCCCCCCCHHHHHHHHHthcCchHHH------------HHHH
Q W5P499       1  EIRERALRSILCKLEHSLVCGADLASHRLLFLHLLEWFNFPSVPHKDEVLGLLSRLVKYPPAVQ------------HLVD   68 (97)
Q Consensus    1  EIR~RAL~nI~sKL~~gL~~~~dl~~~~~Ll~~Ll~WFn~~~~~~~~~~~VL~Ll~L~k~~~~~~~~~~~~~~~~l~~   68 (97)
                  ++|...|+.~+..=..++~~.....+....+..|++.+..++......++.+.++.+......          .+..
Q Consensus   65  ~~~~~a~~~L~~~~~~~~~~~~~~i~~l~~~l~~~~~~~~~~~~~~~~~~l~~~~~~~~~~~~~~~~~~~  144 (175)
T W5P499      65  RTRKEAAWAITNATSGGTPEQIRYLVALGCIKPLCDLLTVMDSKIVQVALNGLENILRLGEQESKQNGIGINPYCALIEE  144 (175)
T ss_dssp         HHHHHHHHHHCCCHHHHHHHHHHITCHHHHHGGGCSCHHHHHHHHHHHHHHHHHHHHHHHHHHHIC--CCSCHHHHHHHH
T ss_pred         HHHHHHHHHHHcCCCHHHHHHHHHHhcCCCHHHHHHHHHHHHHHHhcCCCHHHHHHHHHHHHHHHHHHhchHHHthhhCCCCchHHHHHHHHHH
```

```
Q ss_pred         cCHHHHHHHHHthhhcCCHHHHHHHHHHHHHhcC
Q W5P499      69  LGAVEFLSKLRPNVEPNLQAEIDGILDGL   97 (97)
Q Consensus   69  ~G~~~fL~~Lr~~i~~~~~~~id~I~~~l   97 (97)
                  .|.+.+.+.+.+-+++.+...+.++.++
Q Consensus  145  ~~~~~l~ll~~~~~~~v~~~a~~~l~~l  173 (175)
T W5P499     145  AYGLDKIEFLQSHENQEIYQKAFDLIEHY  173 (175)
T ss_dssp         TTHHHHHHHHTTCSSHHHHHHHHHHHHHH
T ss_pred         hChHHHHHHHthhhCCcHHHHHHHHHHHHH
```

## c) RHS repeat, PF05593, target A0A1G0MXS8

**5KIS_B** YenB, RHS2; ABC toxin, RHS, TOXIN; 2.4A {Yersinia entomophaga}

Probability: 99.47%,  E-value: 4.2e-14,  Score: 98.25,  Aligned cols: 125,  Identities: 20%,  Similarity: 0.284,

```
Q ss_pred         CCccCCeEEEECCCCCEEEEEECCCCCEEEEEcCCCCCEEEEEECCCCCEEEEEeC------CCeEEEEEcCCCCeEEEe
Q A0A1G0MXS8    1  YDAAGRHTSSTDSNGRYLQVSYDTTGKKTKTIYPEGSVVSYSYDGTGRLATITNG------GGRTYGYSYDKLGRRSKLT   74 (127)
Q Consensus     1  yd~~g~~~~~~~~~~~~~~yd~~g~~~~~~~~~~~~~~yd~~g~~~~~~~~~~~~~~~~~~~~~~~yd~~g~~~~~   74 (127)
                  ++..+++.....+.+.+.+.|..|+++.+..+.+....|.||..|+++.+...       ...+.||+.|++++.
Q Consensus   214  ~~~~~~~~~~~~~~i~~l~~~~~~~~~~~i~~l~~~~~~~~~~~~~~~~~~~~~~~~~~YD~~g~l~~~  293 (965)
T A0A1G0MXS8  214  GEGASAWNDLLSGEEYVTLTTADATGTVLTTTDAKGNIQRVRYDVAGLLSGSWLTVRDRTEQVIVKSLTYSAAGQKQRED  293 (965)
T ss_dssp         SSSHHHHHTTBCSCCeEEEEEEECTTSCEEEEECTTSCEEEEEECTTSCEEEEEEEECTTSCEEEEE
T ss_pred         CCccchhhhcCCCceEEEEEECCCCCEEEEEecCCCCEEEEEeCCCCCeEEEEEEecCCCCceEEEEEEEECCCCCEEEEE
```

```
Q ss_pred         cCCCCEEEEEEECCC-CCEEEEEEeCC------CCcEEEEEEEEEEcCCCCCEEEEEEcCCC
Q A0A1G0MXS8   75  YPSGATANYAYDAA~GRLTSLEHKQS-----NGRILASFAYTHDNVGNRHTKTEPDG  125 (127)
Q Consensus    75  ~~~~~~~~yd~~g~~~~~~~~~~~~~~~~~~~~~~~~Yd~~g~~~~~~~~~~~~g  125 (127)
                  .+++....|.||..|++...........        ........|.||+.|++++....+
Q Consensus   294  ~~~G~~~y~YD~~grl~~~~~~~~~~~~~~~~~~~~~~Y~D~~G~l~~~~~~  350 (965)
T A0A1G0MXS8  294  HGNGVVITYTYEAETQRLTGIRTERPAGHASGAKVLQDLRYEYDPVGNVLKITNDAE  350 (965)
T ss_dssp         ETTSCEEEEEECTTTCCEEEEEEECTTCCCEEEEEEEEECTTSCEEEEEETTC
T ss_pred         eCCCCEEEEEEcCCCceEEEEEecCCCCCCCCCceeeEeEEEECCCCCEEEEEEeCcc
```

SPW
motif

SPW
motif

TableS1

| Uniprot entry | Pfam | Pcons_score |
| --- | --- | --- |
| S6TLB9 | PF14882 | 0.071 |
| W5U916 | PF15907 | 0.072 |
| D7MCA5 | PF07725 | 0.079 |
| A0A2A2LSA2 | PF14625 | 0.08 |
| R5P8A5 | PF07538 | 0.081 |
| G1XIQ8 | PF13446 | 0.083 |
| A0A0B0HSH2 | PF13753 | 0.094 |
| J1S4N0 | PF11966 | 0.094 |
| U5QIU9 | PF06739 | 0.096 |
| A0A1A9WU23 | PF02363 | 0.102 |
| A0A252E8A5 | PF17660 | 0.103 |
| L8TNF3 | PF08310 | 0.108 |
| W4LGN0 | PF14312 | 0.111 |
| U7Q0S5 | PF10281 | 0.116 |
| D3UXB8 | PF07634 | 0.117 |
| Q29AL9 | PF14939 | 0.121 |
| A0A1Z4C3E9 | PF17164 | 0.122 |
| Q9NZT2 | PF04680 | 0.123 |
| D5SU36 | PF07639 | 0.139 |
| G3VIY2 | PF00400 | 0.143 |
| D3BR65 | PF00526 | 0.15 |
| I3BT02 | PF03640 | 0.153 |
| A0A094KVK3 | PF00880 | 0.159 |
| R6YH89 | PF14903 | 0.16 |
| O64827 | PF18868 | 0.164 |
| A0A257INW4 | PF13573 | 0.165 |
| A7S4G3 | PF07016 | 0.167 |
| R2SEH8 | PF18780 | 0.168 |
| T0NQR8 | PF08043 | 0.176 |
| A0A1V1NWB1 | PF08309 | 0.179 |
| A0A0L7M9B8 | PF07981 | 0.188 |
| A7RAI0 | PF06598 | 0.2 |
| G1RYA9 | PF06049 | 0.209 |
| S7J9T7 | PF02415 | 0.215 |
| V5CQL0 | PF03406 | 0.232 |
| A0A1I7SWM5 | PF00839 | 0.251 |
| Q7RTC2 | PF12135 | 0.252 |
| Q8PXT0 | PF06848 | 0.253 |
| W5N853 | PF03128 | 0.253 |
| Q6YQH3 | PF11178 | 0.262 |
| R7MCC4 | PF13475 | 0.278 |
| Q6P6X2 | PF10578 | 0.295 |
| H9GJM4 | PF13330 | 0.319 |
| U2FCE1 | PF12779 | 0.355 |
| F3GDU0 | PF00818 | 0.392 |
| A0A1G0MXS8 | PF05593 | 0.407 |
| W5P499 | PF14726 | 0.49 |
| Q8EIH3 | PF07012 | 0.576 |
| A0A2A3HD64 | PF03779 | 0.674 |
| Q8RH85 | PF07661 | 0.711 |
| W5Q8K9 | PF15390 | 0.079 |