

1 **Supplementary Information**

2 **Accurate Multi-Stage Prediction of Protein Crystallization**  
3 **Propensity Using Deep-Cascade Forest with Sequence-Based**  
4 **Features**

5 Yi-Heng Zhu<sup>1</sup>, Jun Hu<sup>2</sup>, Fang Ge<sup>1</sup>, Fuyi Li<sup>3,4</sup>, Jiangning Song<sup>3,4,\*</sup>, Yang Zhang<sup>5,\*</sup>, Dong-Jun Yu<sup>1,\*</sup>

6 <sup>1</sup>School of Computer Science and Engineering, Nanjing University of Science and Technology, 200  
7 Xiaolingwei, Nanjing, 210094, China;

8 <sup>2</sup>College of Information Engineering, Zhejiang University of Technology, Hangzhou, 310023, China;

9 <sup>3</sup>Biomedicine Discovery Institute and Department of Biochemistry and Molecular Biology, Monash  
10 University, Melbourne, VIC 3800, Australia

11 <sup>4</sup>Monash Data Futures Institute, Monash University, Melbourne, VIC 3800, Australia

12 <sup>5</sup>Department of Computational Medicine and Bioinformatics, University of Michigan, 100  
13 Washtenaw, Ann Arbor, Michigan, 48109-2218, United States

14  
15 \*To whom correspondence should be addressed.

## 16 **Supporting Texts**

### 17 **Text S1. How to extract proteins from TargetTrack database for constructing BD\_CRYS?**

18 We extracted 50275 proteins from TargetTrack database (<http://sbkb.org/>) [1, 2] by performing the  
19 following procedures:

20 (1) We extracted all the proteins, which were annotated with the most advanced experimental statuses  
21 by X-ray crystallography experiments, from TargetTrack; these statuses include “selected”, “cloned”,  
22 “expressed”, “soluble”, “purified”, “crystallized”, “diffraction”, “crystal structure” , “in PDB”, and  
23 “work stopped”;

24 (2) We removed all the extracted proteins, which were deposited in TargetTrack before 1 January,  
25 2012 or after 31 December, 2016;

26 (3) We removed all the extracted proteins, which were annotated with one of “selected”, “cloned”,  
27 “expressed”, “soluble”, “purified”, “crystallized” and “diffraction” and deposited in TargetTrack after  
28 31 December, 2014; this procedure could ensure that we did not select the proteins that are annotated  
29 with one of the above statuses and potentially experimented at present;

30 (4) We removed all the extracted proteins whose lengths were less than 30.

31

### 32 **Text S2. How to classify for the extracted proteins in BD\_CRYS?**

33 According to the annotation status, the extracted 50275 proteins can be classified into four classes,  
34 namely, production of protein material failed (MF), purification failed (PF), production of crystals  
35 failed (CF), and crystallizable (CRYS) [3]. MF, PF, and CF proteins are non-crystallizable proteins,  
36 while CRYS proteins are crystallizable proteins. More specifically, MF proteins fail in the first  
37 crystallization step (i.e., production of protein material); PF proteins succeed in the first step but fail  
38 in the second crystallization step (i.e., purification); CF proteins can pass through the previous two  
39 steps but fail in the last crystallization step (i.e., production of crystals); CRYS proteins can pass  
40 through all of three crystallization steps.

41 The classification standard for MF, PF, CF and CRYs proteins is summarized in Table S1. In addition,  
 42 for a protein annotated with “work stopped”, we classified it according to its last experimental status.  
 43 For example, if a protein was annotated with “work stopped” and its last experimental status was  
 44 “expressed”, we classified it as MF. Therefore, we can obtain 35102 MF proteins, 11315 PF proteins,  
 45 1207 CF proteins, and 2651 CRYs proteins. For each class, we used CD-HIT software [4] to remove  
 46 the redundant sequences and kept the proteins below 40% sequence identity. Then, the numbers of  
 47 four classes were 18523, 7164, 815 and 2106, respectively, after removing the redundant sequences.

48 **Table S1.** The classification standard of proteins according  
 49 to the corresponding annotation statuses.

Class deduced from protein annotation	Annotation status
	Selected
Production of protein material failed (MF)	Cloned
	Expressed
Purification failed (PF)	Soluble
	Purified
Production of crystals failed (CF)	Crystallized
	Diffraction
Crystallizable (CRYs)	Crystal structure
	In PDB

50

### 51 **Text S3. The construction procedures of BD\_MCRYs**

52 BD\_MCRYs is constructed as following. First, we downloaded all of 11269 proteins from PDBTM  
 53 database [5], a comprehensive and up-to-date membrane protein database (<http://pdbtm.enzim.hu>), to  
 54 form a dataset, represented as PDBTM\_Set. Then, we extracted 10036 proteins from TargetTrack  
 55 database [1, 2] to form another dataset, represented as Target\_Set, by performing the following  
 56 procedures:

57 (1) we extracted all the proteins, provided by the following membrane research centers, i.e., CSMP,  
 58 MPID, MPSBC, MPSbyNMR, NYCOMPS, TEMIMPS and TMPC, from TargetTrack;

59 (2) We removed all the extracted proteins, which were not annotated by X-ray crystallography  
 60 experiments;

61 (3) We removed all the extracted proteins, which were deposited in TargetTrack before 1 January

62 2010 or after 31 December 2016;

63 (4) We removed all the extracted proteins, which were annotated with one of “selected”, “cloned”,  
64 “expressed”, “soluble”, “purified”, “crystallized” and “diffraction” and deposited in TargetTrack after  
65 31 December, 2014;

66 (5) We removed all the extracted proteins whose lengths were less than 30.

67 Subsequently, we combined PDBTM\_Set with Target\_Set to form a new dataset, denoted as  
68 RBD\_MCRYs, and labeled the proteins by the following criteria:

69 (1) We labeled the samples, which were originated from PDBTM\_Set, as positive samples, i.e.,  
70 crystallizable proteins, because the 3D structures of proteins in PDBTM\_Set have been determined  
71 by X-ray crystallography experiments;

72 (2) The proteins annotated with one of “selected”, “cloned”, “expressed”, “soluble”, “purified”,  
73 “crystallized” and “diffraction” in Target\_Set were labeled as negative samples, i.e., non-  
74 crystallizable proteins, while the proteins annotated with one of “crystal structure” and “in PDB”  
75 were labeled as positives;

76 (3) For a protein annotated with “work stopped” in Target\_Set, we labeled it based on its last status.

77 Noted that the proteins in RBD\_MCRYs were divided into two classes (i.e., crystallizable and non-  
78 crystallizable proteins) rather than four classes (i.e., MF, PF, CF and CRYs proteins). The underlying  
79 reason is that the number of proteins belonging to CF class in RBD\_MCRYs is very limited.

80 Finally, a non-redundant dataset, denoted as BD\_MCRYs, can be generated by using the CD-HIT [4]  
81 with a threshold of 40% to remove the redundant sequences in RBD\_MCRYs. In this work, we  
82 randomly selected 20% sequences from BD\_MCRYs to form a test subset, denoted as MC\_TE, and  
83 the remaining sequences were formed a training subset, denoted as MC\_TR. MC\_TR includes 511  
84 crystallizable and 3569 non-crystallizable proteins, and MC\_TE contains 129 crystallizable and 891  
85 non-crystallizable proteins.

86

## 87 Text S4. Four existing sequence-based features

### 88 Amino acid composition

89 Amino acid composition (AAC) is one of the mostly used features in the protein crystallization  
90 propensity prediction [6-8]. Let  $A_1, A_2, \dots, A_{20}$  be the 20 ordered native amino acid types,  $N_i$  be the  
91 occurrence number of  $A_i$  in a given protein, and  $L$  be the length of the protein. Then, the AAC  
92 feature of a protein, called  $\mathbf{F}_{AAC}$ , is a 20-D vector and can be represented as follows:

$$93 \quad \mathbf{F}_{AAC} = \left( \frac{N_1}{L}, \frac{N_2}{L}, \dots, \frac{N_{20}}{L} \right)^T \quad (1)$$

94 where  $T$  represents the transpose of the vector.

### 95 Dipeptide composition

96 Dipeptide composition (DPC) reflects the frequency of two adjacent amino acids in a protein [9]. Let  
97  $A_1A_1, A_1A_2, \dots, A_{20}A_{19}, A_{20}A_{20}$  be the 400 potential amino acid pairs,  $N_{i,j}$  be the occurrence number  
98 of  $A_iA_j$  in a given protein, and  $L$  be the length of the protein. Then, the DPC feature of a protein,  
99 called  $\mathbf{F}_{DPC}$ , is a 400-D vector and can be represented as follows:

$$100 \quad \mathbf{F}_{DPC} = \left( \frac{N_{1,1}}{L-1}, \frac{N_{1,2}}{L-1}, \dots, \frac{N_{20,19}}{L-1}, \frac{N_{20,20}}{L-1} \right)^T \quad (2)$$

101 where  $T$  represents the transpose of the vector.

### 102 Pseudo-amino acid composition

103 Pseudo-amino acid composition (PseAAC) encodes both the composition information and the  
104 sequence-order information of a protein sequence [10, 11]. In this work, we used Type2 PseAAC [11],  
105 denoted as  $\mathbf{F}_{PseAAC}$ , to encode a protein sequence as a  $(20 \times \zeta \cdot \lambda)$ -D vector, where  $\zeta$  and  $\lambda$  are  
106 the number of amino acid physiochemical characteristics and the rank of correlation along the protein  
107 sequence, respectively. The first 20 components of  $\mathbf{F}_{PseAAC}$  are the traditional AAC, and the  
108 remaining  $\zeta \cdot \lambda$  components are scalar quantities which reflect the sequence-order information of  
109 the protein. The details of PseAAC can be found in [11]. In this study,  $\zeta$  and  $\lambda$  are separately set  
110 to be 6 and 8. Thus, the dimensionality of  $\mathbf{F}_{PseAAC}$  is  $20 \times 6 \times 8 = 68$ .

## 111 Pseudo-position specific scoring matrix

112 Pseudo-position specific scoring matrix (PsePSSM) [12] is the extension of the classical position  
113 specific scoring matrix (PSSM), and has been widely used in many protein attribute prediction tasks  
114 [13-15].

115 For a protein sequence with  $L$  amino acid residues, we first generate its PSSM feature by using the  
116 PSI-BLAST software [16] to search the SWISS-PROT database [17] via three iterations with 0.001  
117 as  $E$ -value cutoff. Then, each element of the generated PSSM is normalized through the logistic  
118 function  $f(x) = 1/(1 + e^{-x})$ , where  $x$  is the original element of PSSM. Let  $\mathbf{F}_{pssm} = (p_{i,j})_{L \times 20}$  be  
119 the normalized PSSM, the PsePSSM feature of a protein, represented as  $\mathbf{F}_{PsePSSM}$ , can be generated  
120 by the following two steps.

### 121 Step I. Calculate the PSSM composition

122 The PSSM composition, denoted as  $\mathbf{v}_{pssm}$ , is a 20-D vector and can be formulated as follows:

$$123 \quad \mathbf{v}_{pssm} = (v_1, v_2, \dots, v_{20})^T \quad (3)$$

124 where  $v_j = \sum_{i=1}^L p_{i,j} / L$ , and  $T$  represents the transpose of the vector.

### 125 Step II. Calculate the correlation factors

126 We calculate the  $g$ -tier correlation factor, denoted as  $\xi_j^g$ , for the  $j$ -th column of  $\mathbf{F}_{pssm}$  via coupling  
127 the  $g$ -most contiguous PSSM scores along the protein sequence as follows:

$$128 \quad \xi_j^g = \sum_{i=1}^{L-g} (p_{i,j} - p_{i+g,j})^2 / (L-g) \quad (4)$$

129 Let  $\boldsymbol{\xi}^g = (\xi_1^g, \xi_2^g, \dots, \xi_{20}^g)^T$  be the 20-D  $g$ -tier correlation factor vector and  $G (G \leq L)$  be the  
130 maximum value of  $g (g = 1, 2, \dots, G)$ ; then  $\mathbf{F}_{PsePSSM}$  can be generated by serially combining  $\mathbf{v}_{pssm}$   
131 with  $G$  correlation factor vectors as follows:

132

$$F_{PsePSSM} = \begin{pmatrix} \mathbf{v}_{pssm} \\ \xi^1 \\ \xi^2 \\ \vdots \\ \xi^G \end{pmatrix} \quad (5)$$

133

134

135

136

In this work, the value of  $G$  is set to be 8. Therefore, the dimensionality of  $F_{PsePSSM}$  is  $20 \times 20 \times 8 = 180$ .

137

138

139

140

141

142

143

144

145

146

### **Text S5. The performances of DCF, SVM, RF and CRTF with two types of feature combinations.**

Table S2 illustrates the performances of DCF models with ADPP and ADPPP on seven training datasets over five-fold cross-validation and seven test datasets over independent-validation. From Table S2, it can be concluded that PsePHSA helps improve the prediction accuracy of crystallization propensity. Specifically, over five-fold cross-validation, the  $Acc$ ,  $MCC$  and  $AUC$  of DCF-ADPPP (i.e., the DCF model using ADPPP as input) are separately 7.3%, 8.5% and 2.1% higher than those of DCF-ADPP (i.e., the DCF model using ADPP as input) on average on seven training datasets. Moreover, among all training datasets, the DCF-ADPPP achieves the maximal enhancements of  $Acc$ ,  $MCC$  and  $AUC$ , which are 20.9% ( $= (0.768 - 0.635) / 0.635 \times 100\%$ ), 15.5%, and 3.8%, on PF\_TR, CF\_TR and PF\_TR, respectively. In addition, on TRAIN3587, all five indices of DCF-ADPPP are increased in comparison with DCF-ADPP.

147

148

149

150

151

152

153

In independent-validation, the  $Acc$ ,  $MCC$  and  $AUC$  of DCF-ADPPP are also higher than the corresponding values measured for DCF-ADPP on each test dataset. Taking CF\_TE as an example, DCF-ADPPP gains 22.9%, 32.4%, and 3.8% increases of  $Acc$ ,  $MCC$  and  $AUC$ , respectively, than DCF-ADPP. Moreover, on three datasets, i.e., MC\_TE, TEST3585 and TEST500, all of five indices of DCF-ADPPP are better than the values yielded by DCF-ADPP. For example, on MC\_TE, compared with DCF-ADPP, DCF-ADPPP obtains 6.3%, 0.1%, 0.8%, 5.2% and 0.4% improvements of  $Sen$ ,  $Spe$ ,  $Acc$ ,  $MCC$  and  $AUC$ , respectively.

154

155

In addition, the performances of SVM, RF and CRTF models with ADPP and ADPPP on seven training datasets over five-fold cross-validation and seven test datasets over independent-validation

156 are summarized in Table S3, Table S4 and Table S5, respectively.

157 **Table S2.** The performances of DCF models with two types of feature combinations on seven training datasets over five-fold cross-  
 158 validation and seven test datasets over independent-validation.

Five-Fold Cross-Validation							Independent-Validation						
Dataset	Feature Combination	<i>Sen</i> (%)	<i>Spe</i> (%)	<i>Acc</i> (%)	<i>MCC</i>	<i>AUC</i>	Dataset	Feature Combination	<i>Sen</i> (%)	<i>Spe</i> (%)	<i>Acc</i> (%)	<i>MCC</i>	<i>AUC</i>
MF_TR	ADPP	<b>68.8</b>	66.9	67.5	0.327	0.734	MF_TE	ADPP	<b>70.4</b>	67.3	68.2	0.343	0.752
	ADPPP	62.0	<b>73.6</b>	<b>70.2</b>	<b>0.335</b>	<b>0.743</b>		ADPPP	63.6	<b>74.2</b>	<b>71.2</b>	<b>0.354</b>	<b>0.757</b>
PF_TR	ADPP	<b>71.1</b>	61.1	63.5	0.278	0.717	PF_TE	ADPP	<b>74.0</b>	60.9	64.1	0.302	0.741
	ADPPP	37.6	<b>89.8</b>	<b>76.8</b>	<b>0.315</b>	<b>0.744</b>		ADPPP	40.4	<b>89.3</b>	<b>77.2</b>	<b>0.333</b>	<b>0.762</b>
CF_TR	ADPP	54.5	<b>80.8</b>	61.7	0.317	0.734	CF_TE	ADPP	55.3	<b>79.7</b>	61.7	0.309	0.754
	ADPPP	<b>78.6</b>	59.5	<b>73.4</b>	<b>0.366</b>	<b>0.753</b>		ADPPP	<b>80.6</b>	62.2	<b>75.8</b>	<b>0.409</b>	<b>0.783</b>
CRYS_TR	ADPP	<b>58.1</b>	85.8	84.0	0.283	0.827	CRYS_TE	ADPP	<b>61.7</b>	86.1	84.5	0.314	0.844
	ADPPP	56.8	<b>88.7</b>	<b>86.7</b>	<b>0.316</b>	<b>0.843</b>		ADPPP	60.4	<b>88.4</b>	<b>86.6</b>	<b>0.339</b>	<b>0.863</b>
MC_TR	ADPP	66.1	<b>96.2</b>	92.4	0.643	0.916	MC_TE	ADPP	72.9	96.1	93.1	0.689	0.936
	ADPPP	<b>72.2</b>	95.7	<b>92.7</b>	<b>0.671</b>	<b>0.925</b>		ADPPP	<b>77.5</b>	<b>96.2</b>	<b>93.8</b>	<b>0.725</b>	<b>0.940</b>
TRAIN3587	ADPP	65.0	87.6	80.0	0.542	0.840	TEST3585	ADPP	61.8	90.5	80.8	0.555	0.851
	ADPPP	<b>69.7</b>	<b>88.9</b>	<b>82.4</b>	<b>0.599</b>	<b>0.870</b>		ADPPP	<b>65.4</b>	<b>91.1</b>	<b>82.5</b>	<b>0.595</b>	<b>0.870</b>
TRAIN1500	ADPP	88.1	<b>80.6</b>	84.4	0.690	0.915	TEST500	ADPP	86.5	80.1	83.2	0.666	0.913
	ADPPP	<b>91.0</b>	78.4	<b>84.7</b>	<b>0.700</b>	<b>0.921</b>		ADPPP	<b>89.8</b>	<b>80.5</b>	<b>85.0</b>	<b>0.704</b>	<b>0.923</b>

159

160 **Table S3.** The performances of SVM models with two types of feature combinations on seven training datasets over five-fold cross-  
 161 validation and seven test datasets over independent-validation.

Five-Fold Cross-Validation							Independent-Validation						
Dataset	Feature combination	<i>Sen</i> (%)	<i>Spe</i> (%)	<i>Acc</i> (%)	<i>MCC</i>	<i>AUC</i>	Dataset	Feature combination	<i>Sen</i> (%)	<i>Spe</i> (%)	<i>Acc</i> (%)	<i>MCC</i>	<i>AUC</i>
MF_TR	ADPP	<b>68.4</b>	65.8	66.5	0.313	0.724	MF_TE	ADPP	<b>71.6</b>	65.9	67.5	0.339	0.746
	ADPPP	63.7	<b>71.5</b>	<b>69.2</b>	<b>0.327</b>	<b>0.736</b>		ADPPP	65.1	<b>71.8</b>	<b>69.9</b>	<b>0.341</b>	<b>0.749</b>
PF_TR	ADPP	<b>66.8</b>	64.0	64.7	0.269	0.705	PF_TE	ADPP	<b>66.6</b>	65.3	65.6	0.278	0.723
	ADPPP	57.0	<b>76.8</b>	<b>71.9</b>	<b>0.314</b>	<b>0.737</b>		ADPPP	57.6	<b>76.6</b>	<b>71.9</b>	<b>0.318</b>	<b>0.753</b>
CF_TR	ADPP	82.6	<b>47.9</b>	73.0	0.312	0.714	CF_TE	ADPP	83.4	43.4	72.9	0.277	0.747
	ADPPP	<b>85.6</b>	45.9	<b>74.6</b>	<b>0.334</b>	<b>0.740</b>		ADPPP	<b>86.4</b>	<b>44.1</b>	<b>75.3</b>	<b>0.325</b>	<b>0.775</b>
CRYS_TR	ADPP	56.3	86.6	84.7	0.283	0.817	CRYS_TE	ADPP	53.6	87.2	85.0	0.278	0.821
	ADPPP	<b>60.4</b>	<b>88.4</b>	<b>86.7</b>	<b>0.334</b>	<b>0.847</b>		ADPPP	<b>55.1</b>	<b>88.6</b>	<b>86.5</b>	<b>0.309</b>	<b>0.845</b>
MC_TR	ADPP	56.8	96.3	91.4	0.578	0.889	MC_TE	ADPP	<b>65.1</b>	96.1	92.2	0.634	0.919
	ADPPP	<b>59.3</b>	<b>97.2</b>	<b>92.5</b>	<b>0.628</b>	<b>0.905</b>		ADPPP	63.6	<b>97.2</b>	<b>92.9</b>	<b>0.659</b>	<b>0.928</b>
TRAIN3587	ADPP	60.9	82.0	74.9	0.433	0.794	TEST3585	ADPP	61.2	83.5	76.0	0.455	0.806
	ADPPP	<b>70.1</b>	<b>85.4</b>	<b>80.2</b>	<b>0.556</b>	<b>0.855</b>		ADPPP	<b>68.1</b>	<b>86.1</b>	<b>80.0</b>	<b>0.548</b>	<b>0.857</b>
TRAIN1500	ADPP	<b>90.6</b>	75.7	83.2	0.671	0.906	TEST500	ADPP	<b>93.4</b>	75.0	<b>84.0</b>	0.694	0.910
	ADPPP	89.3	<b>79.2</b>	<b>84.3</b>	<b>0.688</b>	<b>0.917</b>		ADPPP	89.8	<b>78.5</b>	<b>84.0</b>	<b>0.686</b>	<b>0.919</b>



162

163

164

**Table S4.** The performances of RF models with two types of feature combinations on seven training datasets over five-fold cross-validation and seven test datasets over independent-validation.

Five-Fold Cross-Validation							Independent-Validation						
Dataset	Feature combination	<i>Sen</i> (%)	<i>Spe</i> (%)	<i>Acc</i> (%)	<i>MCC</i>	<i>AUC</i>	Dataset	Feature combination	<i>Sen</i> (%)	<i>Spe</i> (%)	<i>Acc</i> (%)	<i>MCC</i>	<i>AUC</i>
MF_TR	ADPP	65.6	<b>67.7</b>	<b>67.1</b>	0.307	0.722	MF_TE	ADPP	66.6	<b>68.7</b>	<b>68.1</b>	0.323	0.738
	ADPPP	<b>70.5</b>	64.6	66.3	<b>0.320</b>	<b>0.732</b>		ADPPP	<b>70.8</b>	65.2	66.8	<b>0.326</b>	<b>0.740</b>
PF_TR	ADPP	<b>72.8</b>	57.7	61.5	0.264	0.706	PF_TE	ADPP	<b>74.9</b>	58.5	62.6	0.288	0.722
	ADPPP	59.5	<b>73.8</b>	<b>70.3</b>	<b>0.303</b>	<b>0.735</b>		ADPPP	58.7	<b>74.6</b>	<b>70.7</b>	<b>0.305</b>	<b>0.755</b>
CF_TR	ADPP	52.4	<b>81.6</b>	60.4	0.307	0.725	CF_TE	ADPP	52.4	<b>81.8</b>	60.1	0.303	0.744
	ADPPP	<b>75.1</b>	60.5	<b>71.1</b>	<b>0.334</b>	<b>0.745</b>		ADPPP	<b>77.2</b>	65.7	<b>74.2</b>	<b>0.398</b>	<b>0.777</b>
CRYS_TR	ADPP	<b>63.5</b>	80.7	79.6	0.258	0.810	CRYS_TE	ADPP	<b>67.3</b>	80.7	79.8	0.284	0.833
	ADPPP	57.4	<b>87.5</b>	<b>85.6</b>	<b>0.302</b>	<b>0.832</b>		ADPPP	60.4	<b>87.3</b>	<b>85.5</b>	<b>0.322</b>	<b>0.856</b>
MC_TR	ADPP	<b>65.8</b>	94.2	90.6	0.584	0.898	MC_TE	ADPP	<b>75.2</b>	94.8	92.4	0.670	0.927
	ADPPP	<b>65.8</b>	<b>96.3</b>	<b>92.5</b>	<b>0.645</b>	<b>0.903</b>		ADPPP	72.1	<b>96.5</b>	<b>93.4</b>	<b>0.698</b>	<b>0.929</b>
TRAIN3587	ADPP	68.4	<b>82.7</b>	77.9	0.507	0.821	TEST3585	ADPP	67.0	<b>83.6</b>	78.0	0.507	0.830
	ADPPP	<b>75.7</b>	80.2	<b>78.7</b>	<b>0.543</b>	<b>0.853</b>		ADPPP	<b>73.0</b>	80.8	<b>78.2</b>	<b>0.526</b>	<b>0.848</b>
TRAIN1500	ADPP	<b>92.5</b>	68.4	80.5	0.628	0.898	TEST500	ADPP	<b>92.2</b>	67.6	79.6	0.615	0.905
	ADPPP	89.0	<b>75.7</b>	<b>82.4</b>	<b>0.653</b>	<b>0.902</b>		ADPPP	88.1	<b>77.3</b>	<b>82.6</b>	<b>0.657</b>	<b>0.910</b>

165

166

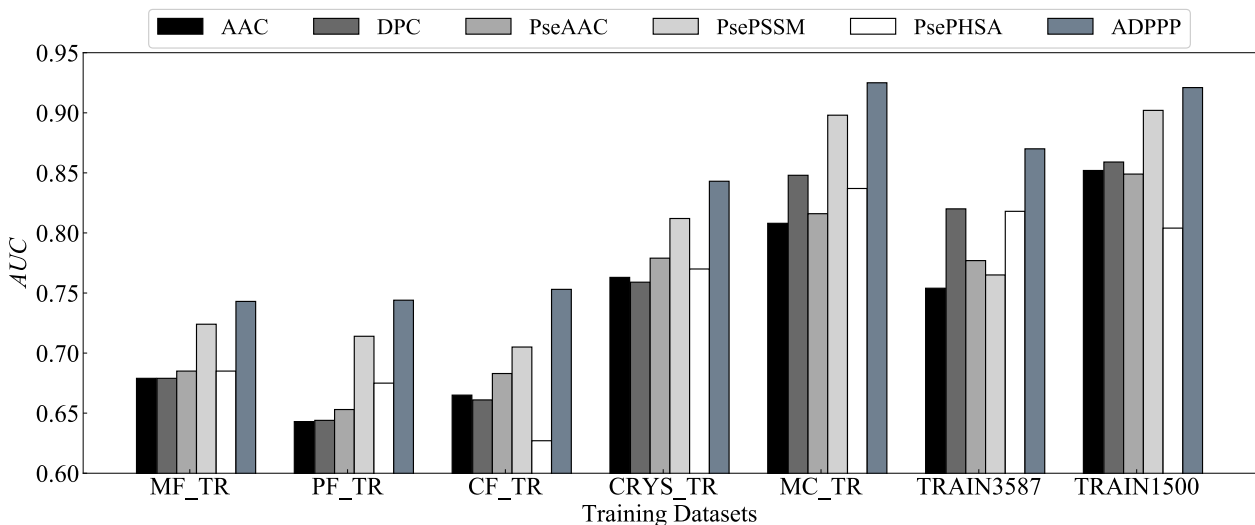
167

**Table S5.** The performances of CRTF models with two types of feature combinations on seven training datasets over five-fold cross-validation and seven test datasets over independent-validation.

Five-Fold Cross-Validation							Independent-Validation						
Dataset	Feature combination	<i>Sen</i> (%)	<i>Spe</i> (%)	<i>Acc</i> (%)	<i>MCC</i>	<i>AUC</i>	Dataset	Feature combination	<i>Sen</i> (%)	<i>Spe</i> (%)	<i>Acc</i> (%)	<i>MCC</i>	<i>AUC</i>
MF_TR	ADPP	69.5	64.2	65.8	0.308	0.721	MF_TE	ADPP	<b>71.2</b>	<b>65.0</b>	<b>66.8</b>	<b>0.328</b>	0.738
	ADPPP	<b>70.4</b>	<b>64.5</b>	<b>66.2</b>	<b>0.318</b>	<b>0.730</b>		ADPPP	70.4	64.6	66.2	0.316	<b>0.741</b>
PF_TR	ADPP	<b>77.1</b>	53.2	59.1	0.263	0.705	PF_TE	ADPP	<b>74.2</b>	56.6	61.0	0.267	0.728
	ADPPP	71.4	<b>63.5</b>	<b>65.5</b>	<b>0.303</b>	<b>0.734</b>		ADPPP	71.0	<b>66.0</b>	<b>67.2</b>	<b>0.322</b>	<b>0.753</b>
CF_TR	ADPP	<b>62.0</b>	72.1	<b>64.8</b>	0.306	0.723	CF_TE	ADPP	<b>62.8</b>	72.0	65.2	0.307	0.739
	ADPPP	57.9	<b>78.4</b>	63.6	<b>0.325</b>	<b>0.741</b>		ADPPP	59.8	<b>79.7</b>	<b>65.0</b>	<b>0.348</b>	<b>0.764</b>
CRYS_TR	ADPP	38.5	<b>92.1</b>	<b>88.8</b>	0.249	0.807	CRYS_TE	ADPP	39.9	<b>92.2</b>	<b>88.8</b>	0.264	0.832
	ADPPP	<b>55.8</b>	87.9	85.9	<b>0.297</b>	<b>0.831</b>		ADPPP	<b>59.5</b>	87.8	86.0	<b>0.324</b>	<b>0.857</b>
MC_TR	ADPP	53.2	<b>97.3</b>	91.8	0.585	0.894	MC_TE	ADPP	62.0	<b>97.4</b>	92.9	0.656	0.921
	ADPPP	<b>65.0</b>	96.4	<b>92.5</b>	<b>0.642</b>	<b>0.900</b>		ADPPP	<b>72.1</b>	96.1	<b>93.0</b>	<b>0.684</b>	<b>0.922</b>
TRAIN3587	ADPP	64.6	84.0	77.5	0.490	0.812	TEST3585	ADPP	64.0	83.9	77.2	0.484	0.818
	ADPPP	<b>69.4</b>	<b>84.2</b>	<b>79.2</b>	<b>0.534</b>	<b>0.848</b>		ADPPP	<b>67.4</b>	<b>84.5</b>	<b>78.7</b>	<b>0.521</b>	<b>0.847</b>
TRAIN1500	ADPP	<b>90.2</b>	73.3	81.8	0.645	0.906	TEST500	ADPP	<b>90.2</b>	72.7	81.2	0.636	0.912
	ADPPP	84.9	<b>80.8</b>	<b>82.9</b>	<b>0.658</b>	<b>0.910</b>		ADPPP	84.4	<b>83.2</b>	<b>83.8</b>	<b>0.676</b>	<b>0.915</b>

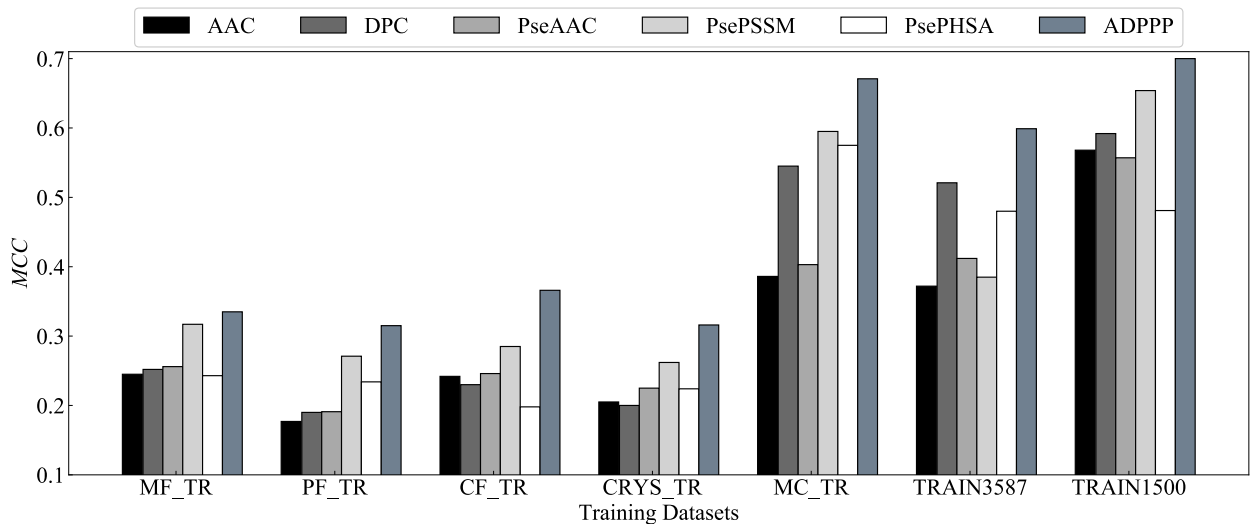
168 **Text S6. Analysis of the contributions of different types of features.**

169 The contributions of different types of features are carefully analyzed. Specifically, we separately use  
170 five individual features, including AAC, DPC, PseAAC, PsePSSM and PsePHSA, and their serial  
171 combination, i.e., ADPPP, as the inputs of the DCF models and evaluate the performances of these  
172 models. Figure S1 summarizes the *AUC* and *MCC* values of the DCF models with different types of  
173 features on seven training datasets over five-fold cross-validation.



174  
175

(A)



176  
177

(B)

178 **Figure S1.** The *AUC* and *MCC* values of the DCF models with different types of features on seven training datasets.

179 From Figure S1, the following two observations can be made:

180 First, PsePSSM is very useful for the prediction of protein crystallization propensity. Concretely, on

181 six out of seven datasets (i.e., MF\_TR, PF\_TR, CF\_TR, CRYST\_TR, MC\_TR and TRAIN1500),  
182 PsePSSM achieves the highest *AUC* and *MCC* values among all of five individual features. For  
183 examples, on PF\_TR, the *AUC* and *MCC* values of PsePSSM are 0.714 and 0.271, which are  
184 separately 5.8% ( $\square(0.714 - 0.675) / 0.675 \times 100\%$ ) and 15.8% higher than the corresponding values  
185 yielded by the second best individual feature, i.e., PsePHSA; on CF\_TR, PsePSSM generates 3.2%  
186 and 15.9% increases of *AUC* and *MCC*, respectively, in comparison with the second best performer,  
187 i.e., PseAAC. The good performance of PsePSSM indicates that crystallization may be closely related  
188 to evolutionary conservation information for a protein.

189 Second, the performance of the combination of five individual features (i.e., ADPPP) is significantly  
190 superior to that of each individual feature. Specifically, from the view of *AUC*, the corresponding  
191 values of ADPPP are 0.743, 0.744, 0.753, 0.843, 0.925, 0.870 and 0.921, which are separately 2.6%,  
192 4.2%, 6.8%, 3.8%, 3.0%, 6.1% and 2.1% higher than those of the optimal individual features, on  
193 MF\_TR, PF\_TR, CF\_TR, CRYST\_TR, MC\_TR, TRAIN3587 and TRAIN1500. With respect to *MCC*,  
194 ADPPP achieves 15.1% average increase in comparisons with the optimal individual features on  
195 seven training datasets.

196

#### 197 **Text S7. Performance comparisons between different prediction models.**

198 Table S6 displays the performances of four models on seven training datasets over five-fold cross-  
199 validation. From Table S6, we can conclude that DCF has the better performance than SVM, RF and  
200 CRTF. Over five-fold cross-validation, DCF shows the best performance on six out of seven training  
201 datasets (i.e., MF\_TR, PF\_TR, CF\_TR, MC\_TR, TRAIN3587 and TRAIN1500) with respect to the  
202 values of *MCC* and *AUC*. Specifically, compared with the second best model, i.e., SVM, DCF  
203 achieves 4.8% and 1.3% average improvements in *MCC* and *AUC*, respectively, on the mentioned-  
204 above six datasets. Moreover, the *Acc*, *MCC* and *AUC* values of DCF are higher than the  
205 corresponding values measured for RF and CRTF on each training dataset. Taking MC\_TR as an  
206 example, DCF gains 4.0% ( $\square(0.671 - 0.645) / 0.645 \times 100\%$ ) and 4.5% increases of *MCC* values in  
207 comparisons with RF and CRTF, respectively. As another example, the *AUC* value of DCF is 2.0%

208 and 2.6%, respectively, higher than that of RF and CRTF on TRAIN3587.

209 It cannot escape from our notice that the *MCC* and *AUC* values of DCF are lower than those of SVM  
 210 on CRYST\_TR over five-fold cross-validation. However, DCF gains 9.7% and 2.1% increases of *MCC*  
 211 and *AUC*, respectively, on the corresponding independent test dataset, i.e., CRYST\_TE (see details in  
 212 Table 3 in the Manuscript). The better performance of SVM on CRYST\_TR may be due to the over-  
 213 fitting in the training stage. As a result, SVM shows inferior generalization ability on CRYST\_TE.

214  
 215  
 216

**Table S6.** The performances of DCF, SVM, RF and CRTF on seven training datasets over five-fold cross-validation

Dataset	Model	<i>Sen</i> (%)	<i>Spe</i> (%)	<i>Acc</i> (%)	<i>MCC</i>	<i>AUC</i>
MF_TR	DCF	62.0	<b>73.6</b>	<b>70.2</b>	<b>0.335</b>	<b>0.743</b>
	SVM	63.7	71.5	69.2	0.327	0.736
	RF	<b>70.5</b>	64.6	66.3	0.320	0.732
	CRTF	70.4	64.5	66.2	0.318	0.730
PF_TR	DCF	37.6	<b>89.8</b>	<b>76.8</b>	<b>0.315</b>	<b>0.744</b>
	SVM	57.0	76.8	71.9	0.314	0.737
	RF	59.5	73.8	70.3	0.303	0.735
	CRTF	<b>71.4</b>	63.5	65.5	0.303	0.734
CF_TR	DCF	78.6	59.5	73.4	<b>0.366</b>	<b>0.753</b>
	SVM	<b>85.6</b>	45.9	<b>74.6</b>	0.334	0.740
	RF	75.1	60.5	71.1	0.334	0.745
	CRTF	57.9	<b>78.4</b>	63.6	0.325	0.741
CRYST_TR	DCF	56.8	<b>88.7</b>	<b>86.7</b>	0.316	0.843
	SVM	<b>60.4</b>	88.4	<b>86.7</b>	<b>0.334</b>	<b>0.847</b>
	RF	57.4	87.5	85.6	0.302	0.832
	CRTF	55.8	87.9	85.9	0.297	0.831
MC_TR	DCF	<b>72.2</b>	95.7	<b>92.7</b>	<b>0.671</b>	<b>0.925</b>
	SVM	59.3	<b>97.2</b>	92.5	0.628	0.905
	RF	65.8	96.3	92.5	0.645	0.903
	CRTF	65.0	96.4	92.5	0.642	0.900
TRAIN3587	DCF	69.7	<b>88.9</b>	<b>82.4</b>	<b>0.599</b>	<b>0.870</b>
	SVM	70.1	85.4	80.2	0.556	0.855
	RF	<b>75.7</b>	80.2	78.7	0.543	0.853
	CRTF	69.4	84.2	79.2	0.534	0.848
TRAIN1500	DCF	<b>91.0</b>	78.4	<b>84.7</b>	<b>0.700</b>	<b>0.921</b>
	SVM	89.3	79.2	84.3	0.688	0.917
	RF	89.0	75.7	82.4	0.653	0.902
	CRTF	84.9	<b>80.8</b>	82.9	0.658	0.910

217

218 **Text S8. How to generate CRYSTALP2, CRYSTALP2\_800, MC\_CRYSTALP2 and MC\_CRYSTALP2\_800?**

219 We remove the proteins, which cannot be accepted by the existing single-stage predictors, from  
220 CRYSTALP2 and MC\_CRYSTALP2 to form four new datasets, i.e., CRYSTALP2, CRYSTALP2\_800,  
221 MC\_CRYSTALP2 and MC\_CRYSTALP2\_800, as follows.

222 First, the proteins with a length of more than 1000 are removed from CRYSTALP2 and MC\_CRYSTALP2, and  
223 the remaining proteins in CRYSTALP2 and MC\_CRYSTALP2 form two new datasets, denoted as  
224 CRYSTALP2 and MC\_CRYSTALP2, respectively. CRYSTALP2 contains 320 crystallizable and  
225 4473 non-crystallizable proteins, and MC\_CRYSTALP2 includes 119 crystallizable and 874 non-  
226 crystallizable proteins. Then, the proteins with a length of more than 800 are further removed from  
227 CRYSTALP2 and MC\_CRYSTALP2, and the remaining proteins in CRYSTALP2 and  
228 MC\_CRYSTALP2 form two new datasets, denoted as CRYSTALP2\_800 and MC\_CRYSTALP2\_800, respectively.  
229 CRYSTALP2\_800 consists of 319 crystallizable and 4355 non-crystallizable proteins, and MC\_CRYSTALP2\_800  
230 is composed of 116 crystallizable and 855 non-crystallizable proteins.

231

232 **Text S9. The web servers of the existing predictors**

233 ParCrys and OB score servers are made freely available at [www.compbio.dundee.ac.uk/parcrys](http://www.compbio.dundee.ac.uk/parcrys).

234 CRYSTALP2 server is made freely available at <http://biomine.cs.vcu.edu/servers/CRYSTALP2/>.

235 SVMCRYSTALP2 software can be downloaded at [http://www3.ntu.edu.sg/home/EPNSugan/index\\_files/svmcrys.htm](http://www3.ntu.edu.sg/home/EPNSugan/index_files/svmcrys.htm).

236 TargetCrys server is made freely available at <http://csbio.njust.edu.cn/bioinf/TargetCrys/>.

237 fDETECT server is made freely available at <http://biomine.cs.vcu.edu/servers/fDETECT/>.

238 DeepCrystal server is made freely available at <https://deeplearning-protein.qcri.org>.

239 Crystalliz server is made freely available at <http://biotool.xmu.edu.cn/crystalliz/>.

240 TMCrys server is made freely available at <http://tmcrys.enzim.ttk.mta.hu>.

241

242

243

244 **Text S10. The performance comparisons between MDCFCrystal and the existing single-stage**  
 245 **protein crystallization propensity predictors.**

246 Table S7 displays the performance comparisons between ParCrys [7], OB-score [18], CRYSTALP2  
 247 [19], SVMCRY5 [20], TargetCrys [13], fDETECT [21], and MDCFCrystal on MC\_TER1000, which  
 248 consists of the proteins with a length of less than 1000. It is found that MDCFCrystal achieves the  
 249 best performances among all the predictors with respect to all of four indices, including *Sen*, *Spe*, *Acc*  
 250 and *MCC*. Taking TargetCrys as an example, which has the second highest *MCC* value, MDCFCrystal  
 251 obtains 237.0% ( $\square(0.765 - 0.227) / 0.227 \times 100\%$ ), 3.3%, 10.9%, and 295.0% increases in *Sen*, *Spe*,  
 252 *Acc* and *MCC*, respectively (*P*-value < 0.05 in student's t-test for the difference in *MCC* values). Table  
 253 S8 summarizes the performances of MDCFCrystal and DeepCrystal [22] on MC\_TER800, which is  
 254 consisted of the proteins with a length of less than 800. We observe that MDCFCrystal achieves the  
 255 better performances in terms of *Sen*, *Acc* and *MCC*, which are separately increased by 373.2%, 5.0%  
 256 and 124.0%.

257 **Table S7.** The performance comparisons between MDCFCrystal  
 258 and six single-stage predictors on MC\_TER1000

Predictor	<i>Sen</i> (%)	<i>Spe</i> (%)	<i>Acc</i> (%)	<i>MCC</i>	<i>P</i> -value <sup>f</sup>
ParCrys <sup>a</sup>	25.2	89.8	82.1	0.150	$7.0 \times 10^{-9}$
OB-Score <sup>a</sup>	28.6	89.4	82.1	0.174	$8.3 \times 10^{-9}$
CRYSTALP2 <sup>b</sup>	46.2	55.9	54.8	0.014	$2.9 \times 10^{-9}$
SVMCRY5 <sup>c</sup>	21.8	78.6	71.8	0.004	$2.8 \times 10^{-9}$
TargetCrys <sup>d</sup>	22.7	93.0	84.6	0.180	$8.7 \times 10^{-9}$
fDETECT <sup>e</sup>	21.8	91.4	83.1	0.143	$6.6 \times 10^{-9}$
<b>MDCFCrystal</b>	<b>76.5</b>	<b>96.1</b>	<b>93.8</b>	<b>0.711</b>	- <sup>g</sup>

259 <sup>a</sup> Results computed using ParCrys server at [www.compbio.dundee.ac.uk/parcrys](http://www.compbio.dundee.ac.uk/parcrys), which can  
 260 output both the ParCrys score and OB score for a protein.

261 <sup>b</sup> Results computed using CRYSTALP2 server at <http://biomine.cs.vcu.edu/servers/CRYSTALP2/>.

262 <sup>c</sup> Results computed using SVMCRY5 software downloaded at  
 263 [http://www3.ntu.edu.sg/home/EPNSugan/index\\_files/svmcrys.htm](http://www3.ntu.edu.sg/home/EPNSugan/index_files/svmcrys.htm).

264 <sup>d</sup> Results computed using TargetCrys server at <http://csbio.njust.edu.cn/bioinf/TargetCrys/>.

265 <sup>e</sup> Results computed using fDETECT server at <http://biomine.cs.vcu.edu/servers/fDETECT/>

266 <sup>f</sup> The *P*-values of student's t-test for the difference in *MCC* values between MDCFCrystal and  
 267 the existing single-stage predictors.

268 <sup>g</sup> '-' indicates that the corresponding value does not exist.

270 The reason for the poor performances of the above existing predictors is that they are not specially  
 271 designed for membrane proteins. More concretely, the training datasets of these predictors contain a

272 few (even no) membrane proteins, which leads that the corresponding prediction models learn very  
273 limited knowledge of crystallization of membrane proteins. As a result, these predictors show the  
274 poor performances in the prediction of crystallization propensity for membrane proteins.

275  
276 **Table S8.** The performance comparisons between MDCFCrystal  
277 and DeepCrystal on MC\_TER800

Predictor	<i>Sen</i> (%)	<i>Spe</i> (%)	<i>Acc</i> (%)	<i>MCC</i>	<i>P</i> -value <sup>b</sup>
DeepCrystal <sup>a</sup>	16.4	<b>99.3</b>	89.4	0.321	$2.4 \times 10^{-8}$
MDCFCrystal	<b>77.6</b>	96.1	<b>93.9</b>	<b>0.719</b>	- <sup>c</sup>

278 <sup>a</sup> Results computed using DeepCrystal server at <https://deeplearning-protein.qcri.org>.

279 <sup>b</sup> The *P*-value for the difference in *MCC* values between MDCFCrystal and DeepCrystal.

280 <sup>c</sup> ‘-’ indicates that the corresponding value does not exist.

281

## 282 **Text S11. Comparisons with predicted-structure-based predictors**

283 We will compare our predictors with the predicted-structure-based predictors, i.e., the predictors  
284 based on the predicted 3D structure information of proteins. However, to the best of our knowledge,  
285 there is no available predicted-structure-based predictors for estimating protein crystallization  
286 propensity. Therefore, we first design a predicted-structure-based predictor, denoted as PSTRCrystal,  
287 and then compare it with our sequence-based predictors.

288 PSTRCrystal is based on the fact that proteins with similar structures have similar functions and  
289 attributes. More specifically, if a query protein has a high structural similarity with the existing  
290 crystallizable proteins, PSTRCrystal will predict it as crystallizable protein; otherwise, it will be  
291 predicted as non-crystallizable protein. Nevertheless, the query proteins have no native structure  
292 information due to that they are candidates used for determining structures by X-ray crystallography  
293 experiments. Thus, we first use I-TASSER [23-25], one of the most powerful protein structure  
294 prediction tools, to predict the 3D structure of a query protein, and then measure the similarity  
295 between the predicted structure and the native structures of crystallizable proteins. Moreover, we use  
296 TM-align [26], one of the most popular structure alignment algorithms, to align two structures and  
297 output the structure similarity, which is measured by TM-score [27, 28]. The value of TM-score  
298 ranges from 0 to 1, and the higher TM-score means that the structures of two proteins are more similar.

299 In light of the above, the procedures of PSTRCrystal are described as follows.

300 In the training stage, given a training protein dataset, denoted as  $TPD$ , and a crystallizable protein  
301 database, denoted as  $CPD$ , we first use I-TASSER to predict the 3D structure of each protein in  $TPD$ .  
302 Then, for each protein in  $TPD$  (taking the  $i$ -th protein as an example), denoted as  $tpd_i$ , we calculate  
303 the average value of TM-scores between the predicted structure of  $tpd_i$  and the native structures of  
304 all proteins in  $CPD$ , denoted as  $AvgS_i = \sum_{j=1}^{N_{CPD}} ss_{i,j} / N_{CPD}$ , where  $N_{CPD}$  is the number of proteins in  
305  $CPD$ ,  $ss_{i,j}$  is the TM-score between the predicted structure of  $tpd_i$  and the native structure of the  $j$ -  
306 th protein in  $CPD$  by using TM-align. Subsequently, for a threshold  $T_m$ , if  $AvgS_i \geq T_m$ ,  $tpd_i$  is  
307 predicted as crystallizable protein; otherwise, it is predicted as non-crystallizable protein. Finally, we  
308 gradually increase the value of  $T_m$  from 0 to 1 with a step of 0.05 to search the optimal  $T_m$ , denoted  
309 as  $T_m^*$ , which maximizes the number of correctly predicted proteins in  $TPD$ .

310 In the prediction stage, for a query protein  $P_{query}$ , we first use I-TASSER to predict its 3D structure.  
311 Then, we calculate the average value of TM-scores between the predicted structure and the native  
312 structures of all proteins in  $CPD$ , denoted as  $AvgS_{query}$ . If  $AvgS_{query} \geq T_m^*$ ,  $P_{query}$  is predicted as  
313 crystallizable protein; otherwise, it is predicted as non-crystallizable protein.

314 In this work, we can only implement PSTRCrystal on a small-scale benchmark dataset due to that it  
315 takes much computing time and resource to predict a 3D structure by I-TASSER. Concretely, we  
316 randomly select 200 crystallizable and 200 non-crystallizable proteins from CRYSDS to form a  
317 small-scale benchmark dataset, denoted as CRYSDS400; then, 20% samples in CRYSDS400 are randomly  
318 selected to form a test dataset, denoted as CRYSDS80, and the remaining samples are used as a training  
319 dataset, denoted as CRYSDS320. Moreover, to construct a crystallizable protein database, i.e.,  $CPD$ , we  
320 download all of 150247 proteins with the corresponding native structures from PDB database; then,  
321 the downloaded proteins whose sequence identities are more than 40% with the sequences in  
322 CRYSDS400 are removed by CD-HIT-2D software [4], and the remaining 144372 proteins are used as  
323  $CPD$ . In addition, because our predictors, i.e., DCFCrystal and MDCFCrystal, are trained on the  
324 large-scale datasets, it is unfair to directly compare them with PSTRCrystal trained on a small-scale



325 dataset (i.e., CRY320). In view of this, we retrain a single-stage predictor, denoted as SDCFCrystal,  
326 by using the proposed pipeline on CRY320, and then compare SDCFCrystal with PSTRCrystal on  
327 CRY80. The performance comparisons between PSTRCrystal and SDCFCrystal are summarized in  
328 Table S9.

329

330

331

**Table S9.** The performance comparisons between PSTRCrystal and SDCFCrystal on CRY80 over independent-validation

Predictor	<i>Sen</i> (%)	<i>Spe</i> (%)	<i>Acc</i> (%)	<i>MCC</i>
PSTRCrystal	<b>87.8</b>	54.8	75.0	0.458
SDCFCrystal	73.5	<b>87.1</b>	<b>78.7</b>	<b>0.590</b>

332 From Table S9, we find that the performance of SDCFCrystal is superior to that of PSTRCrystal.  
333 Concretely, SDCFCrystal achieves 58.9% ( $= (0.871 - 0.548) / 0.548 \times 100\%$ ), 4.9% and 28.8%  
334 improvements in *Spe*, *Acc* and *MCC*, respectively. In addition, we found that PSTRCrystal has the  
335 higher value of *Sen*, reaching 87.8%. The main reason is that PSTRCrystal learns a great deal of  
336 knowledge of positive samples (crystallizable proteins). More specifically, in the training stage,  
337 PSTRCrystal learns the knowledge of all crystallizable proteins deposited in the PDB database.  
338 Therefore, many crystallizable proteins in the test dataset are correctly predicted by PSTRCrystal.  
339 However, PSTRCrystal does not learn any knowledge of negative samples (non-crystallizable  
340 proteins) in the training stage, which leads that many negative samples cannot be correctly predicted.  
341 As a result, PSTRCrystal has a very low *Spe* value, only 54.8%, on the test dataset.

342 The performance of PSTRCrystal (*Sen*=87.8% and *Spe*=54.8% under the threshold  $T1=0.245$ ) further  
343 indicated that only 12.2% of the positive samples (crystallizable proteins) were mistakenly predicted  
344 as negatives (non-crystallizable proteins) but 45.2% of the negatives were mistakenly predicted as  
345 positives. Due to the fewer number of false negatives, it is relatively reliable for a protein to be  
346 predicted as non-crystallizable protein by PSTRCrystal. However, to further reinforce the confidence  
347 of the predicted negatives, the following suggestions are provided:

348 (1) Users can combine the prediction result of PSTRCrystal with the results of other effective protein  
349 crystallization propensity predictors (e.g. DCFCrystal) to obtain the final prediction result by voting

350 or weighted learning.

351 (2) Users can utilize BLAST software to search the query protein against the TargetTrack dataset. If  
352 there exists some non-crystallizable proteins that have a high sequence identity (i.e. more than 90%)  
353 with the query in TargetTrack, then the query is likely to be a non-crystallizable protein.

354 In addition, we have re-selected two new thresholds T2 and T3 which separately makes the rate of  
355 false positives and the rate of false negatives less than 10%. Specifically, T2=0.305 (Sen=38.8%,  
356 Spe=90.3%) and T3=0.235 (Sen=91.8%, Spe=45.2%). For a test protein, the corresponding predicted  
357 probability of belonging positive class by PSTRCrystal is represented as  $p$ . When  $p \leq 0.235$ , it can  
358 be predicted as a non-crystallizable protein (reliable); When  $p \geq 0.305$ , it can be predicted as a  
359 crystallizable protein (reliable). When  $0.235 < p < 0.305$ , the prediction results may be unreliable;  
360 accordingly, users can give up the current prediction result and use DCFCrystal to re-predict the test  
361 protein.

362

### 363 **Text S12. The construction procedures of CRY387**

364 CRY387 was consisted of 387 crystallizable protein sequences and constructed as follows: first, we  
365 downloaded all of 766 protein structures, which were deposited by X-ray crystallography experiments  
366 between October 1, 2019 and December 31, 2019, from PDB database; then, we extracted all of 2534  
367 sequences whose lengths range from 30 to 800 from these structures; subsequently, we used CD-HIT-  
368 2D software to remove the sequences which have more than 40% identity with the sequences in  
369 CRY387 (i.e., the training dataset in DCFCrystal); finally, the CD-HIT software was performed  
370 with a threshold of 40% on the remaining sequences to further remove redundant sequences.

371

372

373

374

375 **Text S13. The performances of DCFCrystal and the existing predictors for membrane proteins,**  
376 **multi-domain proteins and metal-binding proteins.**

377 We had a closer look at the details of proteins in CRYS387 and found that there are 10 membrane  
378 proteins, 27 multi-domain proteins and 63 metal-binding proteins. Moreover, for metal-binding  
379 proteins, only three types, including magnesium-binding proteins, calcium-binding proteins, and  
380 zinc-binding proteins, are considered in this work, and the corresponding numbers are 17, 21, and 31,  
381 respectively (there exists some proteins which can bind with multiple types of metal ligands). The  
382 IDs of the above proteins in PDB database are listed as follows.

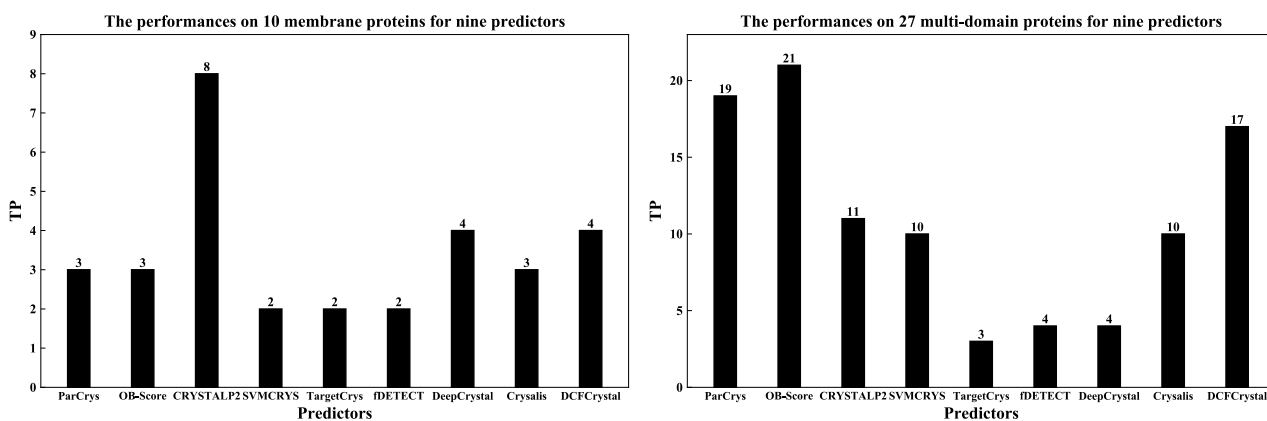
383 Membrane proteins: 6li0\_A, 6tpn\_A, 6tqj\_A, 6umg\_c, 6umg\_r, 6uz6\_B, 6uzr\_A, 6v3i\_B, 6v4l\_A,  
384 6v4l\_D.

385 Multi-domain proteins: 6l30\_A, 6l3f\_A, 6l3w\_A, 6lek\_A, 6ljc\_C, 6t1t\_A, 6t41\_A, 6t96\_A, 6t9i\_D,  
386 6tb2\_D, 6tsz\_U, 6uja\_B, 6ujd\_A, 6uke\_X, 6um4\_A, 6uqr\_A, 6uru\_A, 6uug\_A, 6uut\_A, 6uwm\_A,  
387 6v0k\_A, 6v1v\_A, 6v22\_E, 6v55\_A, 6vbu\_2, 6vbu\_5, 6vbu\_9.

388 Magnesium-binding proteins: 6ulg\_G, 6l3g\_A, 6syt\_C, 6ljc\_C, 6upp\_A, 6tfx\_B, 6syu\_A, 6tdz\_C,  
389 6v4p\_B, 5qtl\_B, 6uja\_B, 6vdd\_A, 6ul5\_A, 6ta4\_A, 6l7s\_B, 6ta4\_K, 6vdk\_A.

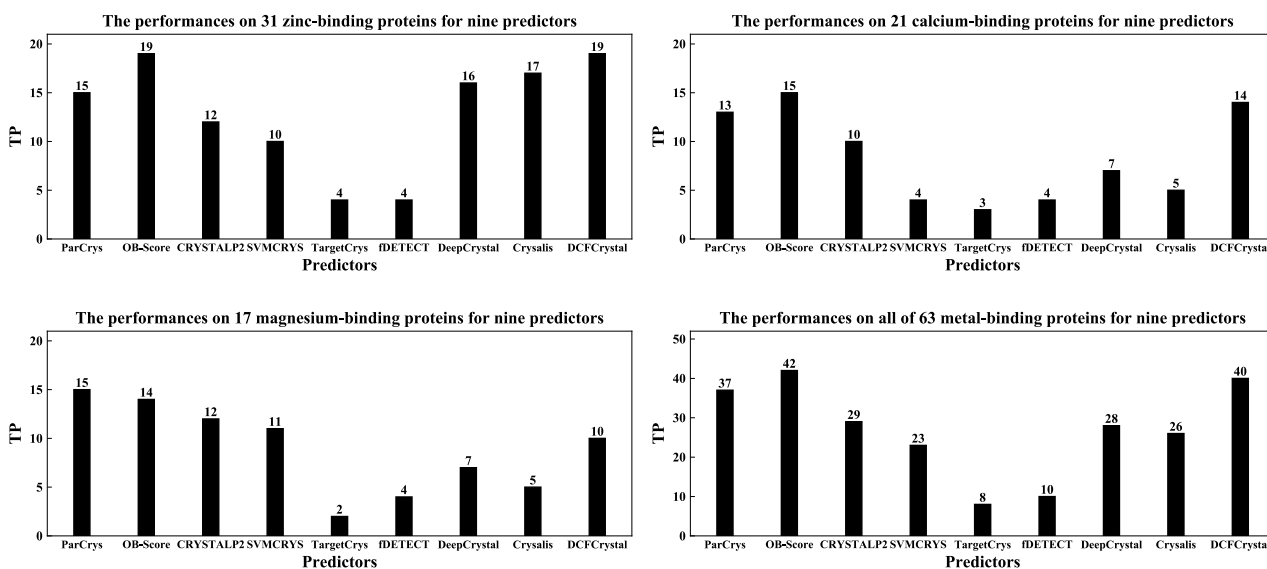
390 Calcium-binding proteins: 6uzb\_A, 6vbu\_2, 6uwg\_A, 6ljc\_C, 6uke\_X, 6v98\_A, 6th7\_A, 6l2j\_A,  
391 6tm6\_A, 6t72\_A, 6v55\_A, 6t9y\_A, 6usc\_A, 6uja\_B, 6uja\_A, 6v4p\_B, 6sz5\_A, 6t0q\_A, 6l2h\_A,  
392 6v0v\_A, 6uwr\_A.

393 Zinc-binding proteins: 6tjv\_P, 6t9l\_N, 6t9l\_M, 6t9l\_K, 6thp\_A, 6l7s\_B, 6t8h\_A, 6llb\_A, 6lhn\_A,  
394 6v4x\_H, 6ull\_A, 6uij\_A, 6v0v\_A, 6tmf\_Y, 6uro\_C, 6tmf\_W, 6t1b\_A, 6tmf\_Q, 6v77\_A, 6lai\_A,  
395 6tly\_A, 6tlx\_A, 6tm5\_B, 6uvn\_J, 6vdb\_A, 6thk\_A, 6ujd\_A, 6t9y\_A, 6tbz\_A, 6tld\_C, 6lae\_A.



396

397 **Figure S2.** The numbers of true positives for nine predictors on membrane proteins and multi-domain proteins.



398

399 **Figure S3.** The numbers of true positives for nine predictors on metal-binding proteins.

400 Figure S2 illustrates the numbers of true positives for DFCrystal and eight existing predictors,  
 401 including ParCrys [7], OB-score [18], CRYSTALP2 [19], SVMCRY [20], TargetCrys [13],  
 402 fDETECT [21], DeepCrystal [22] and Crystalis [29], on membrane proteins and multi-domain  
 403 proteins. It can be found that most of predictors cannot achieve the satisfactory performance. For  
 404 example, on multi-domain proteins, the number of true positives is less than half of all positives for  
 405 each of six predictors (i.e., CRYSTALP2, SVMCRY, TargetCrys, fDETECT, DeepCrystal, and  
 406 Crystalis). Figure S3 shows the numbers of true positives for nine predictors on metal-binding proteins.  
 407 For each of three types of metal-binding proteins, there exist at least four predictors which can only  
 408 predict less than half of all positives, correctly. For example, on 21 calcium-binding proteins, the

409 number of true positives is less than 10 for each of five predictors, including SVMCRYST, TargetCrys,  
410 fDETECT, DeepCrystal, and CrysaliS. The reason for the poor performances of these predictors has  
411 been explained in the Manuscript as follows: the numbers of membrane proteins, multi-domain  
412 proteins and metal-binding proteins are quite limited in the public databases; as a result, the existing  
413 machine-learning-based predictors could only learn very limited crystallization knowledge and show  
414 the inferior performance for these special proteins.

415 It cannot escape from our notice that DCFCrystal cannot achieve the best performance among all of  
416 nine predictors for several proteins in Figures 2 and 3. However, the performance of DCFCrystal still  
417 remains competitive. For example, on all of 63 metal-binding proteins, the number of true positives  
418 for DCFCrystal is reduced by 2 in comparison with OB-score. Nevertheless, DCFCrystal predicts  
419 more true positives than other seven predictors. Additionally, in the training datasets of the existing  
420 predictors, there may exist several proteins which have the high sequence identity with the proteins  
421 in CRY387. As a result, some of proteins in CRY387 may be accurately predicted by one existing  
422 predictor. This may explain that few predictors can achieve the better performance than DCFCrystal  
423 in Figures 2 and 3.

424

425 **Text S14. The performance comparisons between MDCFCrystal and the existing predictors on**  
426 **the membrane proteins recently deposited in PDB database.**

427 We compared MDCFCrystal with ParCrys [7], OB-score [18], CRYSTALP2 [19], SVMCRYST [20],  
428 TargetCrys [13], fDETECT [21], DeepCrystal [22] and CrysaliS [29] on a new constructed dataset,  
429 called CRY47, which contained 47 crystallizable membrane protein sequences. In CRY47, each  
430 protein were deposited in PDB database after July 1, 2019 by X-ray crystallography experiments and  
431 had less than 40% identity with the sequences in the MC\_TR (i.e., the training dataset of  
432 MDCFCrystal); moreover, the lengths of sequences in CRY47 were range from 30 to 800. Table  
433 S10 displays the performance comparison between MDCFCrystal and the existing predictors on  
434 CRY47. As illustrated in Table S10, MDCFCrystal correctly predicts the most (32) membrane

435 crystallizable proteins among all of 9 predictors. Compared with the second best performer, i.e.,  
 436 CRYSTALP2, MDCFCrystal achieves 14.3% increase with respect to the value of *Sensitivity*.

437 **Table S10.** Performance comparison of MDCFCrystal, and eight existing predictors on CRY547.

Predictor	TP	FN	<i>Sen</i> (%)
ParCrys <sup>a</sup>	18	29	38.3
OB-Score <sup>a</sup>	24	23	51.1
CRYSTALP2 <sup>a</sup>	28	19	59.6
SVMCRY5 <sup>a</sup>	15	32	31.9
TargetCrys <sup>a</sup>	16	31	34.0
fDETECT <sup>a</sup>	16	31	34.0
DeepCrystal <sup>a</sup>	24	23	51.1
Crysalis <sup>a,b</sup>	21	26	44.7
MDCFCrystal	32	15	68.1

438 <sup>a</sup>Results computed using the corresponding web servers, which are listed in Text S9.

439 <sup>b</sup>Results computed using CrysalisII, which is the sub-predictor of Crysalis.

440

441 **Text S15. The performance comparisons between CDCFCrystal and the existing predictors**

442 Table S11 summarizes the performance comparisons between CDCFCrystal and TargetCrys [13] on  
 443 TRAIN3587 (i.e., the training subset of CRY57172) over five-fold cross-validation. From Table S11,  
 444 we can see that CDCFCrystal achieves the better performance. Concretely, the *Sen*, *Acc* and *MCC* of  
 445 CDCFCrystal are separately 18.1% ( $(0.697 - 0.590) / 0.590 \times 100\%$ ), 2.1% and 8.9% higher than  
 446 the corresponding values yielded by TargetCrys.

447 **Table S11.** The performance comparisons between CDCFCrystal and  
 448 TargetCrys on TRAIN3587 over five-fold cross-validation.

Predictor	<i>Sen</i> (%)	<i>Spe</i> (%)	<i>Acc</i> (%)	<i>MCC</i>
TargetCrys <sup>a</sup>	59.0	<b>91.7</b>	80.7	0.550
CDCFCrystal	<b>69.7</b>	88.9	<b>82.4</b>	<b>0.599</b>

449 <sup>a</sup> Data excerpted from [13]

450 Further, we compare our predictor with ParCrys [7], OB-score [18], CRYSTALP2 [19], MetaPPCP  
 451 [30], SVMCRY5 [20], XtalPred [31], SCMCRY5 [32], PPCpred [3], RFCRY5 [8], PPCinter [33] and  
 452 TargetCrys [13] on TEST3585 (i.e., the test subset of CRY57172) over independent-validation, as  
 453 show in Table S12. It is found that CDCFCrystal achieves the best performance in terms of *ACC* and

454 *MCC* among all predictors. Compared with the second best predictor, i.e., TargetCrys, CDCFCrystal  
 455 achieves 1.9% and 6.3% improvements of *Acc* and *MCC*, respectively. More importantly, the *ACC*  
 456 and *MCC* values of CDCFCrystal are obviously higher than those of ParCrys, OB-score,  
 457 CRYSTALP2, MetaPPCP, SVMCRY5 and XtalPred. For example, there are 46.5% and 183.3%  
 458 enhancements of *Acc* and *MCC*, respectively, between CDCFCrystal and SVMCRY5. Moreover, all  
 459 of the four indices of CDCFCrystal are highest in comparisons with PPCinter, PPCpred, SCMCRY5  
 460 and MetaPPCP. Taking PPCpred as an example, CDCFCrystal gains 6.9%, 7.4%, 7.4%, and 26.6%  
 461 increases of *Sen*, *Spe*, *Acc* and *MCC*, respectively. It cannot escape from our notice that OB-score and  
 462 ParCrys have the highest values of *Sen*, but with the lowest values of *Spe*. The reason is that too many  
 463 negative samples are predicted as positives by these two predictors. Together with the fact that the  
 464 number of negatives is larger than that of positives, thus this makes their *MCC* performances lowest.

465

466

467

**Table S12.** The performance comparisons between CDCFCrystal and eleven existing predictors on TEST3585 over independent-validation.

Predictor	<i>Sen</i> (%)	<i>Spe</i> (%)	<i>Acc</i> (%)	<i>MCC</i>
ParCrys <sup>a</sup>	78.6	31.8	47.5	0.110
OB-score <sup>a</sup>	<b>80.3</b>	31.4	47.8	0.120
CRYSTALP2 <sup>a</sup>	74.4	45.7	55.3	0.200
MetaPPCP <sup>a</sup>	61.7	59.0	59.9	0.200
SVMCRY5 <sup>a</sup>	75.2	46.7	56.3	0.210
XtalPred <sup>a</sup>	67.0	62.3	63.9	0.280
SCMCRY5 <sup>a</sup>	46.0	91.0	76.1	0.440
PPCpred <sup>a</sup>	61.2	84.8	76.8	0.470
RFCRY5 <sup>a</sup>	51.0	<b>95.0</b>	80.0	0.530
PPCinter <sup>a</sup>	61.7	89.8	80.4	0.550
TargetCrys <sup>a</sup>	58.0	92.7	81.0	0.560
CDCFCrystal	65.4	91.1	<b>82.5</b>	<b>0.595</b>

468

<sup>a</sup> Data excerpted from [13]

469

470 **Text S16. The details of proteins selected from four families.**

471 We separately selected 18, 12, 38 and 32 proteins from four protein families, including PF13419,  
 472 PF00583, PF13649, and PF03061 for case studies. The details of selected proteins in each family are

473 described as follows.

474

475 **PF13419**

476 >371319\_JCSG (crystallizable)

477 MTPIAQRDGGAIQLVGFDDTLWKSSEDYYRTAEADFEAILSGYLDLGDSRMQQHLLAVERRNLKIFGYGAKGMTLSMIETAIELTEARIEARDIQRIV

478 EIGRATLQHPVEVIAGVREAVAAIAADYAVVLITKGDLFHQEQKIEQSGLSDFPPIEVVSEKDPQTYARVLEFDLPAERFVMIGNSLRSDVEPVLAIGG

479 WGIYTPYAVTWAHEQDHGVAADPEPRLREVPDPSPGWPAAVRALDAQAGRQQ

480 >508504\_EFI (non-crystallizable)

481 MATAHPIKAAIFDMDGLLIDSEPLWLQAELDIFTALGLDTSRDSLPTLGLRIDLVKLVYQTMPWQGPSQEEVCNRIIARAIDLVEDTRPVLPGIEYAL

482 ALCRQQGLKIGLASASPLHMQERVLAMLGVEKYFDCLVSAEYLPYSKPHPEVYLNAAAQLDVPDLPQCVTLEDSVNGMIATKAARMRSIVIPSVEYRA

483 DPRWALADIQLES LDQLRKDDIS

484 >508213\_EFI (non-crystallizable)

485 MIKNLIFDFGKVLVNHDLQPLLERHFGDDEVSLIGFHKILSDPEFINMCDRGIIPFEKMIDGAIKRYPEYSDAFTFFKDNYLEEITGEIEGMRVLLKLLKQ

486 SGFKLYGLTNWSDTIYRVMEKFDIFQLLDGTVISCEEHFHFKPEKEIYLRCDKYGLKPSECLFTDDRMVNVLGAKAIGMEAVLFTTPKEYIFEIERIFGIKI

487 EQ

488 >508460\_EFI (non-crystallizable)

489 MKQPTAVIFDWNLTIDTSINIDRTTFYQVLDQMGYKNIDLSIPNSTIPKYLITLLGKRWKEATILYENSLEKSQKSDNFMLNDGAIELLDTLKKNITM

490 AIVSNKNGERLRSEIHHKLNTHYFDSIIGSGDTGTIKPSPEPVLAALTNINIEPSKEVFFIGDSISDIQSAIEAGCLPIKYGSTNIKDILSFKNFYDIRNFICQLI

491 NI

492 >508136\_EFI (non-crystallizable)

493 MVADIELDAAQSIAAVLFDKDGTLLGYDASWGPVNRELASIAAKGDAALADRLLAACGMDPVTGHVVPDLSLLAAGNTAEIAAGLVAAGSSCDVVEL

494 TQRDLRLFTEAADKSVPVTDLKAFFARLKARGYKLGIASSDNENSIRQTAIRFGFEEDIDFVAGYDSGYGTPQPGMVLGFCEAIGFPPERVAVVGDNN

495 HDLHMAKNAGAGLRIAVLTGTGSRESLGADAHYCFDDITGLEALLPERAV

496 >508399\_EFI (non-crystallizable)

497 MLLFDIDGTLIRTRGFRQTEAAALSEWLGRPVTTTEGVDFAGRTDPAILLDILKASGLPEHTARHLLPEALEVYSRAMIRRLRPEHLEVLPGVVMLEEE

498 LSEWPDVYLGLVTGNLRPVAFHKLAMAGLAGYFEGEAFGCDHANRNLPLAIERIREATGYPFTGADAVIIGDTPHDVACARHAGASVAVVCTGGYS

499 RDALEACRPDLLEDLSDPEPLFKLLTQQALS RKAS

500 >508064\_EFI (non-crystallizable)

501 MVEEREFDSIIFDL DGTMW DSTENAAIVWKEIAKKDSRITDEVTGPKLKALYGLPLEDIARGFLSVPEDVAIETMEKCVVAQCYPYLAHEHGILLGKIEE

502 TLKELSKKYRLFIVSNCKSGYIEAFLEAHKLGQYFDDFECPPGGTGKLLADNIRIVMKRNQLRNPIYVGD TGGDGDAAHQAKIPFVYARYGFGEATEYE

503 YVIDSFDQLTTLRMTE

504 >508516\_EFI (non-crystallizable)

505 MIEAFIFDL DGVITDTAYHYMAWRKLAHKV GIDIDTKFNESLKGISRME SLDRILEFGNKKYSFSEEEKVRMAEEKNNYVSLIDEITSNDILPGIESLL

506 IDVKSNNIKIGLSSASKNAINVLNHLGISDKFDIADAGCKNNKPHPEIFLMSAKGLNVNPQNCIGIEDASAGIDAINSANMFSVGVGNENLKKANL



507 VVDSTNQLKFEYIQEKYNEYIVRRII  
508 >508355\_EFI (non-crystallizable)  
509 MPENVDLNALLFDMGVLVLDVSRSYRAIETVEHFTGRQIGENAIQRYKNYGGFEDDWKLTHAIVTDTAMEVPISRVIDEFQDRYRGDDWDGFIITEE  
510 PPLIDDQTLDRNLNQSILGIVTGRPEEEAQWTLDHQNWTDYFPLLVGKEKQGDRAKPNPFLEHSLTMLAAAGCPIDPEEAVYIGDSVDDMDAAREAG  
511 MWRIGVVPPYVETDEHKPLLEEHAHVVIDDLNTPDVLSTLDERTPARSTR  
512 >508560\_EFI (non-crystallizable)  
513 MLTGIGAIVFDLDGTYQSESLGGQIAACADRYLADLLSVSPEEAGEIVRRVRRELTARFGREASLSDACRELGGDLRELHRRFAAEVAPEPHLRDRSRV  
514 VQLRLTLGANRELYLYTNNNRALSGRIMDAIGVTGLFRRVVTIEDSWRPKPDQLALEALFAALGRKPSECLFVGDYRVIDLRLPAELGCSVYLSRTVDE  
515 LLGLTLPLSEEQQ  
516 >508267\_EFI (non-crystallizable)  
517 MSIATPTGYPRAVFDLLTGLLDSWTAWNAAAGSEPAGRAWRAEYLRLAYGCGRYVPYEQLVREAAARATGLPESAPAALEAGWHELPVWDDARALL  
518 RALRPHCKLAVTNCSDRLGRQAAGLLGVDWDAIVTAEAGFYKPDPRPYRMALQALQVPADAAAFVAGSGHDLFGTAAVGLRTCWHNRLGLARP  
519 QGAPEPELQSATLAVALPWLQAFRPAVR  
520 >508330\_EFI (non-crystallizable)  
521 MSLPHAPRTPAPAAAASSAVPLRLAVVDMSGTSIVEHGLQDTAFARTLDQHGVPAGTPEHDDAARRFRALRPTSRTAVFPRVFADRAVAAAATRTFEATF  
522 DALLGQHGVQAVPGAEEALVRLRALGLHVCLCTGYARHTQNMILESGLWMGLDLSLSPDDAGRGVPPYPMILTALLGLDLDVRSVLVVGDTAED  
523 MTAGRRAGAGLVVGVRTGRDADDVLLAAGADRVVPLADVPDLVARTR  
524 >508151\_EFI (non-crystallizable)  
525 MITRLTDIELEQIRGVIFDLDTLAHSNPDFKGLRAALGIGSGTDILEHIHSLETTVAKMQALEIVHDYELESSRQASWIEGAQALIAFLKTRQLPLAILTR  
526 NMPEAAKITIEKLGIDIPVLTRYDAEPKHPQGIYLICEQWQLNPADILYVGDYFLDLQTAQNAAGSRCALYCPEDVPDYAQAADLLVSCYHSLIQAWPK  
527 >508552\_EFI (non-crystallizable)  
528 MDAVFFDFDGVLTDDKYGSDTTNRYLGEATGLGFDRIDQALERYNDDLLLGRLGHPDVWSALCCELGLEMDYNLLDAAFRSTPMNEGVLALARRLQ  
529 GRFRLGIITDNKSDRMDCLRAMHELDALFDPIVSAAVGANKSGGEIFQHALALCGLRPERSLFIDNSRRNLEVAAGLGMATLFHDDVRNDVLALRQA  
530 LERILGISLA  
531 >508141\_EFI (non-crystallizable)  
532 MFQKFYPTHEYVESSYEIDYEKLYKNGYRGIIFDIDNTLVEHGADASERAVALIKRLKKIGFEVCLISNNKEDRVKRFNQDIKIKYIFNAHKPSIKNYLKAM  
533 EYMNTNKSNTIFVGDQIFTDVYGANRAGITSYLVKPIGKKEEQIVIKRLLERIVLSFYRRKQAKQKKRS  
534 >508225\_EFI (non-crystallizable)  
535 MMPLLVEELGLSCFAADLLTHYDALTALDHAQAMPHAAEVLTELRRRGVNIQVVTNGWEDAQTRCLAGCDLSLADDDVISEAVGLSKPDCIYHLALK  
536 RLGVTTAHSWVFGDSDPRNDVWGPQVGMRAAYLPTGHPLNGERPDVTLTDLRDVLNLP  
537 >508190\_EFI (non-crystallizable)  
538 MAHAAPKPKVIIFDVNETLLDLETMRYSVGKALDQGEELTTLWFSTMLHHSVTTVTGQDFGKIGVAALMMVAQNNNIDITEEQAVTAIKTPLLSLP  
539 AHPDVKAGLTALKAQGFKIVSLTNSSNKGVETQFKNAGLTDYFDKMSIEDIKVYKPDLSRYAWALEQLNIKPEEALMVAAHGWDVAGAKAAGLQTA  
540 FVARPGKALYPLAQEPDYIVKDLSELVEILK  
541 >508424\_EFI (non-crystallizable)

542 MRVTMRRLLLWDIDGTLSTDGIAANAMRTALRQLVGPVRIERTSYAGKTDWQIVRESLPSVDEATIQSRLQEFIALYTAELTAQREALIARSTVFAGV  
543 VEALHALSTHAYQAPLTGNVAAAARIKLECTGLLRWLEVEAGAYGDDHFDRLALPPIAAGRARERYRYAFTPADVVIGDTPRDIACGRAFGARTVAVA  
544 TGPFSMAELAEYNPDVLLPDLRDTVAVVEAVLGN  
545  
546 **PF00583**  
547 >BhR182\_NESG (crystallizable)  
548 MIIREATVQDYEEVARLHTQVHEAHVKERGDIFRSNEPTLNPSFFQAAVQGEKSTVLVFDEREKIGAYSVIHLVQTPLLPTMQQRKTVYISDLCVDETR  
549 RGGGIGRLIFEAIISYGKAHQVDAIELDVYDFNDRAKAFYHSLGMRCQKQTMELPL  
550 >GilaA\_00357\_a\_SSGCID (non-crystallizable)  
551 MAKYLGRYSLRSIQERDLSRLTTLLEQLSVVGEVPREKLVSFYKSVSTNPSHDVTVVVDETDTVCCATLIIEPKLLHAGRSVGHIEDVVVDLTLRNQGI  
552 GRFLITSLIERARNNDYKVIDTDPDTAEFYKCKGMKQKGLMMAIYF  
553 >033720\_NYSGRC (non-crystallizable)  
554 MKPDETPMFDPSLLKEVDWSQNTATFSPAISPTHPEGLVLRPLCTADLNRGFFKVLGQLTETGVVSPEQFMKSFEHMKKS GDYVYVTVVEDVTLGQIV  
555 ATATLIIEHKFIHSCAKRGRVEDVVVSDECRGKQLGKLLLSTLTLSSKLNKYKITLECLPQNVGFYKFGYTVSEENYMCRRFLK  
556 >APC103096\_MCSG (non-crystallizable)  
557 MTTYVWRGAVDDRALGELHAEAFEHAYVDIGWSAQLKGHSLGWVTAHDDGHDDPVGFVNVAWDGGVHAFVLDTMVARAVRGRGIGRGLVARAA  
558 SGARHAHCEWLHVDYEPLEPFYAACGFEPAGLVRLR  
559 >358652\_JCSG (crystallizable)  
560 MRTLNKDEHNYIKQIANIHETLLSQVESNYKCTKLSIALRYEMICSRLHTNDKIYIYENEGQLIAFIWGHFSNEKSMVNIELLYVEPQFRKLG IATQLKI  
561 ALEKWAKTMNAKRISNTIHKNNLPMISLNKDLGYQVSHVKMYKDID  
562 >021319\_NYSGRC (non-crystallizable)  
563 METVRIDEGFDRYEELLSLIRASFAYMDGRIDPPSSAHALTAASLNKRARDEIAFAAVAGRELLGCIFCKPEADCLYIGKLAVAPGRQKGVGRMLIAAA  
564 EETARDLGLPALRLQTRIELAGNQATFAAWGFVETARTAHPGFTRPTSVEMRKVLS  
565 >IDP92632\_CSGID (non-crystallizable)  
566 MFGYRSNVPKVRLLTDRLLVRLVHDRDAWRLADYYAENRHFLKPWEPVRDESHCYPSGWQARLGMINEFHKQGSIFYGLFDPDEKEIIGVANFNSV  
567 VRGSFHACYLGYSIGQKWQKGLMFEALTAAIRYMQRTQHIHRIMANYMPHNKRSGDLLARLGFKEGYAKDYLLIDGQWRDHVLTALTPDWTPG  
568 R  
569 >APC103015\_MCSG (non-crystallizable)  
570 MASVLVAEELRPDHAWRVDAEVDGVVGAIAIVERNGPEAIVRAIVTMPGFERLGLAGELLEACAELARRHGATVLSASCAPGDRAKALFEEAGLRTV  
571 ERLARSLA  
572 >APC103097\_MCSG (non-crystallizable)  
573 MSHLRLHEVTDQNLRALTDQFKPGQERFVAPVVLISAEAYVTPTAWPRAILEQDKIVGFVMANFDPDNEIEAFRCGIWRLNIAAHAQGRGVGRFAVE  
574 EVASEARKRGQDRMTVLWAEEGGPEFYLRCGFEPGERIFDQTLGVRTTAASAPRT  
575 >022418\_NYSGRC (non-crystallizable)  
576 MKKHDLVYLTEDASHDAAIEIINEEAFGPGRFTRAAARIREQPHDRALSFCADNGETIASVRMTPVTAGSVKGHLLGLAVRPSHKNQGIGRELVRIA

577 VEAARRKGSEAVILVGDPPYYQPLGFQVVRHGALQFPGPVDPARVLVVPVALDVHARLEGMIAWRDDGATCLTARAEAQGAAA  
578 >APC102186\_MCSG (non-crystallizable)  
579 MQLRSVIAADHPALFALWSRTPGIRLRAEDAYPFFLAYLQRNPGLSLLVETEVEVIACLMAGHDGRRGYLQHLVDPGYRGLGLARRMLDEVLARLA  
580 REGIGKSHVFLDAAEEAQAFWRAQSDWERRKDIQVFSTREGHA  
581 >026752\_NYSGRC (non-crystallizable)  
582 MTSELSKIKKTYEYKKFIRCLIDELLINKEKIDELDAKRKKQRVEDIKAKCLRKFNLNIGFPPNSDVLALATDEEKKKLSLIRKKPIRTLSGVAVVAVMT  
583 SPEPCPHGKCAFPCPGKESVFGDVPQSYTGKEPATMRGIMYNFDPYVQTSERLKQLENVGHPTDKVELIIMGGTFPARDTSYQENFIKGLDAMNGVI  
584 SETLEEAQKINETASHRCVALTIETRPDYCKEEHVNEMLKLGATRVELGIQSTYDEILDFVKRGHSVSESIKATSRKNSGLKVSYHIIPGLPHTTEEMDK  
585 ENIRRVFNPEFKPDLIKFYPCLVIEGTELYDLWKKGEYKIPITDEEAVELITYGKSIMPKWIRTSRIQRDIPATVIDEGVKKSNLDELVYNNLEKKGKCKC  
586 IRCREVGHVYKKGKIPDNSSIKLLIEEYASGGKEFFITYEDIKNDLLIGYLRRLIPDMNSVFRPEIDENTALIRQVHVCQQQELGSKIEDTKNWQHK  
587 GYGKMLLEEAELAKSLGSKILITSGIGVREYYKKQGYDRIGPYMGKLN  
588  
589 **PF13649**  
590 >029991\_NYSGRC (non-crystallizable)  
591 MALRAGAGKVGTLMASDEFNRWEGRAVEDYVFGTAPNAFLSSCRDILPKQGRALSIADGEGRNGVFLAECGLSVLSVDFSPAAQGKAQRLAAARG  
592 VTIETQTADLLTWDWPGDFDVIAGIFFQFVEADERPRIFQAIRDALKPGGLLLIEGYRPKQLIYKTTGGPSRADNLYTHELLEAAFQDFDNLISIREHDEISA  
593 EGSQHVGMALIDLIGWPK  
594 >030055\_NYSGRC (non-crystallizable)  
595 MNIGYKVGELIYDANIYDGLNTFLSDLQFYKKWLPKNKEAEILELCCGTGRLTIPIAKDGYSICGVDTYTPSMLEQAKMKAIEAELVIDFIEADIRMLDLQ  
596 EKFDLIFIPNSIHHLRNEDLNFALGCVRNHLKAGGLFLLDCFNPNQIYVESEKQVAVIAEYTTDDGRDVLIKQTMRYESTTQINRIEWHYFINGEFHS  
597 IQNLDMRMFFPQELDSYLERAGFDIIHKFGSFEFEEAFNDNSEKQIYVLTLDNKNVLYEKIQNR  
598 >023596\_NYSGRC (non-crystallizable)  
599 MSDVATRPNYDLRDEIKAYWSERAATFDLSPGHEIFSEEERAAWHRLILRHLGEGAGRSALDLASGTGVVSHLLDDLGRVAGMDWSEPMLEARQK  
600 AKSRGRDISFRMGDAENTMEPDDHYDVVNRHLVWTLVDPAAAAREWLRVLKPGGRVLIVDGDVFNATRLERFFSSLSVWGQVRVGLLRPDAPSQPR  
601 EMLETHRSILARVHFSQARAEAVVGLLRAAGFADITVDTDLGEIHRMQAKNWNLFKGLARRSQHRAIRASKPVA  
602 >030692\_NYSGRC (non-crystallizable)  
603 MNNEINNNYSTFSEKYDQLFDSELYQEWFVDFVTKNSKATSIMDLGGGAGRLAVLLAQLGYTVDVLDLSPEMPLSLAQKHANEANVDLSLLQADMRDF  
604 SDWKKEYPIIVSFADALNYLPLNSDFKLAIQVYDHLAVGGQFLFDVITPYQVNVLYDNYYYNNDDDDENIFMWTSYSGEQENSVDHDLKFFVYDEA  
605 IDAFKIMREIHHEQTYDLKTYQETLRSAGFHNIIEVFANFGQNNIDENTERWFFRAVK  
606 >029936\_NYSGRC (non-crystallizable)  
607 MSRPEPPVSRDPWLERWLPPLREAGGQGPVLEIGCGEGEDSRALAEAGVRLIAFDLSADAVAAAASARAPGARFVCQDVRQAFPLGGERAGAVVASLS  
608 LHYFPWDETVEIVERIRDCLPPGGKLLCRLNATDDHHYGASGHPEIAPDFYLVGDGEPKRFFDERSARALFADGWRILSLEHRVSGKYALPKALWEAAL  
609 EKTG  
610 >029963\_NYSGRC (non-crystallizable)  
611 MSDYVKLNKTNWDERAPLHAASADYAVQRFVDEAGYLSDVVQFDRPLLDIRGLRGVHLQCHIGTDTLSLARLGAQMSGVDFSPASLAEARTLARR

612 CNTPIDYHESDVFLAAEVLPGNFDLVYTGIGALCWLPISIERWAQTVGALLKPGGRLLFIREGHPMLWAINEDHDDSLRVELPYFETREPLVWDDENTY  
613 VETDSPLKATMTHEWNHGLGEIISALLAQGLDITGLVEHQSIPEALPGQMVVDERGEWRLKEAPWRLPLSYTLQAVKRGG  
614 >029945\_NYSGRC (non-crystallizable)  
615 MASKQTLFRIFYRIGFTFPWDGHPLAQSLRDLVEGTGDAALPAGKALEIGCGTGDCAIYLAQHGWNVTAVDVFAKPLERARAKAGAAGAAVDFVQA  
616 DVTRLGQAGIGTGFELIVDNGCLHNMSDADRDAYVREVTGVAAPQARLLIVAFVPGGRFGVGVVEDAEMQRRFTADWTLAAGPERELDGAERTPA  
617 RYYLFQRR  
618 >030033\_NYSGRC (non-crystallizable)  
619 MTTPPFDTEDLFDEDYLHFAARPLEETSDAAGLVERLLELRPGERVLDLACGHGRIANRLAARGLEVSGLDITPVFLDRARQDAHERGVEVDYVRG  
620 DMRELPWTGHFDAIVNWFTAFGYFDDAGNRRVLDQARQALRPGRLLLELNNQAHVLRAFQQAIVHEADGDLLVDRHRLDPLTSRNIVERTIVRAG  
621 RVRRFHYFTRMFAPPELRDWLLAAGFDEVEGFGEDGGPLTADSRMLVLARR  
622 >029937\_NYSGRC (non-crystallizable)  
623 MHPMTLDAIDFNLLYRLQKHSSTYKKSQEEWDGKAWINEKIHEGFYNDEMERRIDLGVQSLLDVCGPGTFALRFAPRLKQVYALDYSPKMLEV  
624 LEHNAQKRAISNICPLCLDLEESWEGVAPCDVVIASRCLEVEDMRAVLQKLHEKAKKAVYITYKNGGSFLEAEVLEAMGRKITPKPDYIYLLNILYQM  
625 GIRASLDFISPGEPYGTDFSFSYHRSTLWSIGEMTPQEEAGLKDYEACQKKGKTPAHKNSSWAFISWRK  
626 >030020\_NYSGRC (non-crystallizable)  
627 MRDKYKYIGPVYDFLSNLYSGKNIHRCKTAMLDVETVKPGDRILFAGVGHGRDAIRAAELGAEVTVVDLSETMLRKFADAQHKEAPHLTIRRIHSDIM  
628 KVDEFEQYDMVVANFFLNVFDEDMVKVLEHLIRLGKADARVVVGDFCYPTGNILSRMFKKLYWYMAVFIFWLFANNAFHKIYNYPEHMQRGLGQ  
629 VTEKKHFKLLNMDCYWSILGRKQA  
630 >APC103455\_MCSG (non-crystallizable)  
631 MELYQRLVPWYRLLTPPSEYAEAEASCYRAAFERAVPDARTLLELGAGAGHNAYHLKQRFACLTLDIADEMLDLSRALNPACEHLPGDMRTLRLERQFD  
632 LVFVHDVAVGYMRTPEELAAAIGTAFATRPGGAALFAPDHLRETFTTETHVHEGGDAQRALRCMEWTWDPDPEDTECVVDFVFALRENGRTEVVHD  
633 QHRVGLFARATWQKLLSQAGYQVEIVARPLTEEYSDEIFLCRRPAA  
634 >029943\_NYSGRC (non-crystallizable)  
635 MKDNTLYGPIAHLYESFSDATDHIKVEIRTIFNLAGDIHGKSVLDLACGYGLFSREYRNRGASKVIGVDISENMIAIAKSKSQYGDIEFHVNRICKME  
636 SESFGKFDIVNAAWLFCHAESLEDETMRVIAAHLKPAKGLIAYTFEPDYRLEKGNENYCIKILSEEPVKDITTLVKAELTTPPSPFTMYRWSREQYQ  
637 TAIQKAGFKQFKWQKPMLLERDIEAHPPGFWDDFQRNCLDTALVCQI  
638 >029924\_NYSGRC (non-crystallizable)  
639 MLPATAVLRKKRRRSKMTMIHMLLTNSRSAAKLYRVWKEKKEMDIHEIDWNEVWKDLHEQNLRRRKGECAIWSRESALEFLERSNKNPQRVA  
640 KVFSDLGVGPASRVLDVAGPGTLAVPLASRCAHVTAVEPAAGMVEVMKEFAQKEGVENLEIVSKRWEDIDPAELSGPYDVVFASYSLGMPDIRAAVE  
641 KMCKLATKRVCYWLFGSSPWEQWMIDLWPALHGQEYRSGPKADVLFHVLYDMGIYPNMETLQLLYTRTFPDFDAAVENFKREYHVETDAQEKILR  
642 EYLSSVLKKEGGEFVLSSENSMRVKLWWEVDQCA  
643 >030758\_NYSGRC (non-crystallizable)  
644 MERQLISHIAHYDHPIAAPVSEQNLERLLTRAKLAPGARILDGCGEAPWVLRALHELHPEAVADGVDISEHALTAAQKAADQRGLSDRLGLHHVPAAD  
645 FTGTEPYDLVLCVGSTHAFDGLTATMQDIRRHLRPGGLALVGEFVETPPTPEALTKLGANLDDYGDLSATVAQAEDAGYATVYGHTSDLAEWHEYE  
646 WSWIGTLTNWALDHPGPDGDAALAAARDHRDMWLNNGYRDILGFVTLRLRRTD

647 >029247\_NYSGRC (non-crystallizable)  
648 MLYGEEFHVVDLACGPGSFSMRLNRFPAIRVTAIDLPLLLTLAKEALSEYKDRIQFFSGDIATADCFAAITDKPQAVVSSTAIHWLLPEQQVALYRNIF  
649 NLLDEHGLFMNADHQRFNDRNPCQKRIAQLHDEDTQKKAWTAGVQDWDWSWFASATRHHELADLMDARTAIFKDRPTPLPTTVEFQLTALRQAGFSE  
650 TGTLWQFLDDYVIAGWK  
651 >030103\_NYSGRC (non-crystallizable)  
652 MFLYQYFKNPKQTGAFCASSKLSKLITSHVQHAKNIVEIGPTGSGFTKYILKQKSHNASFFAVEINPHMAKKLEQNIKNIDIEISSAEFLPNILEKRAINT  
653 VDLIISGIPWALLNSSEQDLLLSIHEALEENGCFATFAYILPTPKGRVFKKLFATFSKVEISPIWRNLPPAFVYFCTK  
654 >029958\_NYSGRC (non-crystallizable)  
655 MNQDELKEAFDQATGYEERQMKLAPVYEGIYFQLQWVFSGLPDNARILCVGSGVGTSEISYLAKRFPNWRFIAVEPSGGMLDICRKRAEQEGFSSRC  
656 VFHEGYLDSLGLPECDGATCLMVSQFFLDKEDRVSFQSIASLKPGGILVSSDLSEMVGSNSEYSTLINLWAKMLHGSNVSSDTVDKIHSAWLKDVAI  
657 LPPEEIKHLIHRGGFELAVQTYQAQLVRAWVARRL  
658 >029890\_NYSGRC (non-crystallizable)  
659 MRLPGMLRPTAERHFHSIFYLRHNARRQEHLATLGLDLGNKSVLEVAGIGDHTQFFLDRGCKVLCTEPRGENLDVIRQRFSGSNPNVTVDHLLDLGDG  
660 LPAEAHQYDVVYCYGVLYHLSRPAEALAWMCDRAVDLLLLLETCVSYSGEDEPFLVSRASSPSQAITGTGCRPSRVWVMNRLREKMPHVYV TATQPR  
661 HRQFPLDWRANGPIASTGLARAVFVASRAPLNLPTLVEELPMVQRRC  
662 >030143\_NYSGRC (non-crystallizable)  
663 MTLDENRLNELLGRFVDFGGSYQALSAVLGDRLGLYRALQMTMPATPEEVAEERAAERYVREWLAGQAAGGYVTYDPASGRYSLTEEQAFLLD  
664 TGGLQAAAAFHIPVAVAKNIERITEAVRTGGGFPWHDHDELPEFEGVERFFRPGYAMNLVSSWIPALDGVAKLRAGARVADVCGGHGSSTLLADAFP  
665 ASEVIGFDYHPASIDLARKRAVEAGISDRVSFEVASAADFPDGTDYDLVAIFDALHDMPLGAAKHIREVLEEDGTLLEVPDAGDRVEDNLNPVGRLY  
666 YGASTVICVAHMSAEPRTALGAQAGEATLLELHDAGFGSVRRAAETPFNIVLEARI  
667 >020588\_NYSGRC (non-crystallizable)  
668 MALDRADFYDAELARHNRQLRVAADFGADDRVLDIGCGAGQTTREARAAPQGEAIGVDISAEMLEEARRRSAAEGLRNAMFEQGDAQFHGFPTGS  
669 FDLCISRFVGMFFADPAAAFANIGRAMRPGARLVVMVWQSRERNEWSRAIRQALAPAIASAGAANPFSGLDPPVATDLLSAAGFTSIDFADVQEPVF  
670 YGSDVDAAFDALTSLYLVQDALASTNEPPDKPLQRLRDLLEGHMTPEGVFFDSRAWIITARRAGGGG  
671 >029948\_NYSGRC (non-crystallizable)  
672 MSVTPGADPYALSAEFYEVMAIPHWDMKRQVLVSALTARGPVKDHVLDIGAGTGLSTVTVADTIADVPIHAVEPSAAMRAALVSRILSRPDLIDRVTV  
673 HPVNLEELDLPERLGAVVFLGVIGYMDKQARQHFWAALRPRLTTPRAPVIVEVMALDQMPVPENTIAQQRIGVRHNEVRISGQPAGSDAEHWTMRYV  
674 VSEGDKVTREFTAETWHTVGLAELAHEAEAHDMTFEQLHPIIGVLHPR  
675 >030036\_NYSGRC (non-crystallizable)  
676 MVEAAFDVIGERYVKETQELRAQLAGQWVIDRLPARGARVLDLGCCTGFPTAEQFDGAGVEVGVDESPRMLELAARRVPGARLVRGDMRALDA  
677 GLGDFDAVTAFFSLLMLSRAEVS AVLESVRDRLRGPRLALAMVQGDSDAERMSFLGADLTVTAYSPRALGEVVSAGAGFVVEELREVEVVCESDRPF  
678 AVEPQVFVYARAAG  
679 >030069\_NYSGRC (non-crystallizable)  
680 MTRRTGVAPADETLYTDARLVAVYDLFNAGDHDFAFYAARIGAAPQRILDIGCGTGFARRLAAAGHDVVAIDPASAMMDYARRQPGADAVRWIACD  
681 LRDLPFGAPFDAAMVTGHAFQCLLSDDAIDSTLHGVRRLVTSGGRFLFETRNPRLEPWRAWTPQSSARRVDSAPFGAVELQHVSHAVEGPIVSFDTHY

682 RFLRDDTRVTHASRLRFIAQRELQARVAAAGFSAVEWYGDWQGASFDDATSVEIIACRV  
683 >030714\_NYSGRC (non-crystallizable)  
684 MTQNNWRYFFDEYAEKYDNEIFTKNTKAEIDFIEQELNIPAGSFILDVGCCTGRHSIELAKRGYSVTGIDISERMLSIARKKCEKESVSVDFIQANAVDFK  
685 VNKLYDACICLCEGAFGLLSEGEDPFDRDIMILKNINKTLKTGSKFIFTALNGLRMIRLFNDEDVSKGKFDQLAIVESSPMSDYLENAPDNIFLREKGFIA  
686 SELVYMLKIAGFLVENIWGGTAGSWNRKPLRMDEIELMLVSKKERKC  
687 >029974\_NYSGRC (non-crystallizable)  
688 MNDDSERWDERYRSEEFLLGEKPSRFLAERIEEVKCLCPGRKALDIACGEGRNSIFLARHGYSVTGLDISPVAVEKARRWAGREGLACDFRLADLETYA  
689 FDERFDLIINFNLLRDLPQEVAAALTPGGVVIFDTILESPTAPVPHRKEFLLQPGELARFFAPYPGTILFCGEYPDSATPTAKLIYRHSK  
690 >030643\_NYSGRC (non-crystallizable)  
691 MTRTPKTQWSPQEYSRFGDERSRPFELLARVQAVDPRTVVDLGCASGVLTLELARRWPNASVLGLDSSAELLATAPADLPANVRLEQGDIAFRADG  
692 VDVVFTNAALQWLPQHRDLISAWAHQLNPGGWLALQVPGNFGAPSHALMRQVAESPRWAARLAGVLRGTSTDGAEDYARLAISSGLVPDAWETT  
693 YVHLLGGDDPVLRWVHGTGLRPVISALTADEFAEFESY GALLRRAYPRSGDITPFGFRIFCVASKPDGAR  
694 >BuceA\_17257\_a\_SSGCID (non-crystallizable)  
695 MLLKNLRPANDYDRFATETLEPWDLLFISIRQLARGMTAGTIADIGTATGVVVPRLATDPAMRGWRYVIGIDLPAMLDEGRPRIHELGLDDTIEMRVG  
696 DALALPFDDGTLTMAVGRATLHHLDPKALSLETEMYRVLAPGGIALVHDMRRDAPQHLLDRFTAMRAAADYPPTHVEEKITLDEAHALVAEAGLAEVA  
697 SIYSPSMGLGALGFEILLKCPALA  
698 >030085\_NYSGRC (non-crystallizable)  
699 MVTDYNQGEIAERYKKTKAIPVRTRIEAYSFLKHIGDVTGEKVVDIACGAGDYARVMRRAGAARVVGFDISEKMIGLAREQEAHEPLGIEYFVADASQ  
700 EVAQQDYDLAISAYLLVYARDRDELARMCRGVACRVRPGGRFVTLTNPGLYTFGRVPDYRKYGFQIKLADAAFEAGAPIELTAFVDGAPLVIENYYLPI  
701 AAYEAALRQAGFHDFVHVPELAAAPQGEDEGDFWDDYLNYPPIAIVIECVRD  
702 >029904\_NYSGRC (non-crystallizable)  
703 MVQLIKKYEDLLCLLDELIKNESTFRWDEFYLERERDVPFFVLAPDEQLVEYVCTGLIESGKVLLELGC GPRNAIYLAENHFEVDVVDLSQK AIDWAM  
704 DRANERKASIRFIRENIFNLNVNKASYDLVYDSGCFHHIPHRMDYIHLVTTALKPGGHFGLTCFIENGLGAAISDLDVYRQQSLQGG LGFTQKLI  
705 QIFDDFAVIEIRKMKKEYPNDSRFGVNGLLTALFQKKKVEKVK  
706 >030006\_NYSGRC (non-crystallizable)  
707 MARQPQRAAAGGASGSRSLIPRVSLAADDGHLRVAVFYGENPDVAVDPSIASQFKLPFIGSQT LQRVVVMKISELARRCGLARSTLLY YEKLGVIAGT  
708 RAANGYRHYDDEDLQRLLMVQALQAGGLSLKQCLACLAGLEQATLLARVRELDEKLAQMQRARDLLADLAGLRAHSGDEFKAWQRQLQHQPQ  
709 AYFAWVMKQGFSEKERYHLQWLSKDMNEHERYIRDFKLLLDGMSYWGPGDSLFTQQQFAALSSQPRRIFDMGCGRGAATLALA QVTDATIVAILDDE  
710 EALA AVARSASAVGLGQVTTLCANMAALPADLAPADLIWAEGSAYTIGVANALQAWRPYLAGPAACLVLSDLVWLTDTPPEEALAFWQRDYPAMQTL  
711 AGLLKTQVEAGYRCLSHTPLPQRAWHNYLDPIERNLARHRAELGDSPA WQDLSREVAIHRQYLGSYGYVICCLQAA  
712 >030718\_NYSGRC (non-crystallizable)  
713 MQAYTGFAEVYDTFMDNVPYEEWSEYLAGLLKEYGVKDGLVLELGC GTGSTRRLFERGYDMIGIDLSEDMLEIAREKMDVGYSFDDILYNQDM  
714 REFELYGTVS AVVSICDSMN YITKPEELKQVFRLVNNYLD PQGIFIDMNTIHKYRDILGETTIAENREDCSFIWENFYHDKEAINQYDITIYKKTEIELED  
715 EELESTTSLYERNLYQRWEETHYQRAYELEIEIKALLIEAGMEYVAAYDAMTKNAPSEESERIYIAREKKQENKFYL  
716 >029919\_NYSGRC (non-crystallizable)

717 MFCLMSKEYLSWDKFIVTKIPSSVKLDPIIYEYIKKDHLILDIGCGVGKVSQQLAFQGFYVEGIDINETGILAAQDSARKLNLADKAHFRVGDADKLPY  
718 MDDKFDIVIMHGLLTIIVDNSDRNKIIQEAYRVLNPEGHLYIVDFGQTWHSIDIYRERYLKDFPITKEEGSFLVYNKDTGEIEFISHHYTEKELIFLLVNGF  
719 KIDYFRDFLLFSGFSFLYIVQNNKILILIIYFFYLWEKETLQ  
720 >029900\_NYSGRC (non-crystallizable)  
721 MRAVGLLTLYPCGRGPAEGRGEGFLGQMLYEWGSKGRCLCFDSLTPYPPVNYDDLAPIYDQQYDSYRDDLHFYAGLAERAGGRVLEIGAGTGRVTA  
722 FLTRRGA AVLGVPESGEMIVGAQARAAREGLTLELVQATAQTFASDERFGLIIPFNALMHLYTPAEQLAALQNI RAHLAPGGQFVFDLYVPHFGAMNT  
723 LRHEGETFHVPDGSRTDLFLLQRHDAPRQVITTEYFADTTAPD GALRR AHHTLTQRYYTRFEME WLLRCAGFEAPRV TGSFQGGPLVETSEVMVFQAR  
724 GA  
725 >030062\_NYSGRC (non-crystallizable)  
726 MGRDSTDRTTTRTSLEERVAKERGRPRSTSTVLWEGVQRL LADAAAAPGAGELTVVDLGGGTGGLAVRVAALGHRVVVVDPSDALAALERTVEAG  
727 LSDRVRALQGDATDLAAVLEPGRADVLLCHGVLEVVD DPRAALRAAH DALRPGGRSL LVAQWPASVLARVLSGHL DQALHVLSDADHRWSGHDP  
728 LRRRFDRASATALAQGAGFTVTAVEGTRTFSDVVPSTRTESDADVELLRLEALASTSPELLGLAGHLHLHATR  
729 >030101\_NYSGRC (non-crystallizable)  
730 MSRDMLEQASKYDHWAWLYNRTLGPYRGAYKIGPIERVV LPHVPAGGAILDLCCGTGQLAAALSERGFNVIGLDGSADMLRYARENAPS VTFTEGDA  
731 CNFTFDTPFDVAVLCTASLNHMQLNDLVFVSSVS RALKPGGIFVFDVNHPAQMSRYWRGHPTGEINTDFAWLITPQYD SAANRGAFTVDIYRRPD  
732 AHPVSM DLRFVRLAQFRIRLALLSRFSRLRPHWEHHSV VNRWGHNL DAMSRALHESGFSVELRSTQGGPVDDSHAAYFFCRKAPTAEKQAETAK  
733 ETAL  
734 >030088\_NYSGRC (non-crystallizable)  
735 MPNPAKPEKTGWRKYMLPRLIRLSRSAPKDRPLAWDRY WAGITATGPGGEVLWDAGSDHEFLGYRDLILRHLDPALPVVDVCGGHGSFTRALAAHF  
736 PQAIGVDISAHAAALAAEPGNPGNVSEVRDMTAPGAGAGLVAGPANVFVRGVLHVLSPADQAALAE NLRVLAGGRGTVFLAETNFQGNPVEYVT  
737 HLGATQRSIPAPLERAIRGLPMPGHFGPKERSRALPPASWELLEEGAAA IETNPVTGVEGQSRVPGYVAALRPRRPSGETPKHAGEDAPLTGSS  
738 >030131\_NYSGRC (non-crystallizable)  
739 MAPQVTD FSWNEMWKQSYDREK GIRANPLLEYW DKRANDFSLMRKSN DYDFGRKVYAAALSSV LTPDSSMLDIGAGPGSFTIPFAQH IKS VTAIEPS  
740 KGMVA VFKENAKELGVENFNIEEMVQDL PQDGSFDSQFDVVAISLVLMFPD VWPRI LQMEQYSKGYCAIVAGIPDWKNPRAASKSDVEEFQILYNM  
741 LLSQGRFPNVSVIDYK CERMVEDEIECRKIIYEQYEGDLTPEAEEQIRKEVIARSKDNKCLISSRS AVIWWNPKEIV  
742 >029999\_NYSGRC (non-crystallizable)  
743 MPQLQKEWFDTWFDTPY YHILYSNRDESEAEVFLTNLMNHMAV PKGASILD LPCGKGRHTLFLAEKGYTLTGADLSVASIALAQSHAPAGV TFLVHDL  
744 RKP AWNESFDYVLNLF TSGYFETEAE DRAAFTLSKALKSGGSLVIDFMNVTRAVNLLKEETKVMEGIEFQLKRYVKDGYIHK EIRFEADNTPFFFT  
745 ERVKALTLDDFKEFFTFAGLTLVDTFGSYQLDAYDAAESDRLIMI AKK  
746  
747 **PF03061**  
748 >212066\_NYSGRC (non-crystallizable)  
749 MIPTPANPRILETMSQLFVSFPHCATLGF EYVGTGRKPTLKLQWREDLVGNPATGILHGGVITSLVDTCSAIAVTAHLPELETIATLDLRIDY LKSATPGK  
750 AIHCTAE CYRLASQIAFTRAVCYHDNPADPIAHGVATFMRESSRTPMLQEDGR  
751 >212054\_NYSGRC (non-crystallizable)

752 MDSNPESTGGRPAIPPGMDFDPSRFGSFMRRHGHTGFIGMQYRDHGDNWIELALPWREDLVGDPETGVLASGPIISLLDNATSMSVWALRGGFRPQV  
753 TLDLRVDYVRAATPGKTIVAWAECYQLKRSMFVVRGIAHDGDISDPVAHAAGIFIQVDADGWSRDSGAKA  
754 >211853\_NYSGRC (non-crystallizable)  
755 MENSALLTELEERSFTEKNYQTWLGVRLVRHKPGFVHLELPVRPEFLNTLGTVHGGFLANLADSALCSAILSELPPGITCSSIEIKVNYLLPVRGNILRAD  
756 ASVIRRGKNIGVSRAELFAPDGALAAVATGTFMIQSLSSFCLPRVD  
757 >211881\_NYSGRC (non-crystallizable)  
758 MTETYSINEEEIARRWQKFAHVSPYNRELGLLPHVVRPDWCVLKVEYQDALVGDQPTRVLHGGVVTTALLDAAFGEAIFVKLPAFRPMATLDRIDYLK  
759 PATPGRAVLGGAVCYKLTPELAFVRGCAYHESLEDPIATAVGIYMFTEGRPVISNEEVPR  
760 >212127\_NYSGRC (non-crystallizable)  
761 MPSSPEADPAVAIPPTVSAPFDVELGLEFTELTADGARAQLEVKPKHLQPMGLVHGGVYCSMVESMASMAAFTWLSTRGGGGVVGNNSTDFLRAIS  
762 SGTVYGTAEPLHRGRRQQLWLVVITDDADHVIARGQVRLQNLEAPPSDG  
763 >212132\_NYSGRC (non-crystallizable)  
764 MTHYSSLGSEGAGEEVDPEYEHGGFPEYGPASPGPGFGRFVAAMRRLQDLAVSADPGDGVWDQAAEQADALASLLGPFQADEGQGPAGRTPDLPG  
765 MGSLLLPPWTLTRYAPDGVEMRGSFSRFHVGGNSAVHGGVLPPLFDHFMGMISHAAARPISRFLHVDYRKITPIDTPLLVRGRVTRTEGRKAFVAAE  
766 LVDADEMPLLAEANGLMVRLLPDQP  
767 >021665\_NYSGRC (non-crystallizable)  
768 MDGTERENVTEIRIPFRDIMHGHMHNAAYYAHAEAALANLWRHRPAIAEPPAYLVRRSACIFHRGLRFDEPARFTVTVAKIGGSSIGFAVRVETGDKL  
769 AAEVEIVWVAVDRARHQPVQLPGPTREWLAYAA  
770 >212237\_NYSGRC (non-crystallizable)  
771 MIGSMTETGAPNSVTLRFLAAPTVDGHSGSVDAGTVLEWVDKAAAYAAVWAKSYCVTAYVGNHIFADPVNSGDMVEVEATIVYTGRSSMHIRTVV  
772 SSSDPKGGPATMRSQCMVIFAVGEDGKPIPVKQFEPASDAEIEQRDHALARIKVREQIVEAMNHQEYTDAGTAERVTLRFMAAPTDVNWGGKVHGGI  
773 VMKWIDEAAYVCASRYCGKDTVAVFSGGVRFYRPLLIGHVVEVEARLVYTGTKGMHIAVHVRSGDPKGRELNLTYYCLTVMVARDGAGNSVPIPAW  
774 VPVSEDEKRLHAHARELLEIRGTAPGNRLPNHLLAGD  
775 >211994\_NYSGRC (non-crystallizable)  
776 MIYSMTQIEARYSETDQMGVIYHGNYPTWFEVARTDYISKLGFSYKDMEDSGIISPVIDLDIKYIKSIFYPEKVTIKTWVERYRSLRSIYKYEIYNEAGEL  
777 ATTGSTALTCIKKEDFKPIRLDKYFPEWHATYSEVDKRNKAGECLEVINGL  
778 >212075\_NYSGRC (non-crystallizable)  
779 MTTQHTPPPAPIPPGFVALADSGGPMHHIGPLYRLRQGLVKFGRVRRHVNPLDILHGGMMASFCMMLPLSVHDKSAEVADRFLPTISLQID  
780 YLAAVPLGAWVEGQAQPLRVTRSLVFAQGLVSADGIPCARTSGVFKIGPALSGLVAQ  
781 >212227\_NYSGRC (non-crystallizable)  
782 MSGTSAALKPPADATTPVRHPDAPAPGELLGAHYGHCFGCGEESHGLHLAARAGQGVSITAEFTVQPAHQGAPGLAHGGVLATALDETGLSLNWLL  
783 RTIAVTGRLEDFVRPVPVGTVLYLEAEVTAVAGRKIYSTATGRIGGPEGPVAVRADALFVEVKVDHFTDNGRQEEIRAAMNDPQLRRARAFEVNP  
784 >212240\_NYSGRC (non-crystallizable)  
785 MKYKLLVLDVDGTLLNDAKEISKRTLASLLKVQMGIRVALASGRPTYGLMPLAKTLELGNYGFFIISYNGGQIINAQNGEILFERRINPEMLPYLEKK  
786 ARKNNFAIFTYHDDTILTDSSDNEHVRAEANLNLKIIQEEEFSTAI DFAPCKCILVSNDEEALKDLEEHWKKRLDGLDVFCEPYFLEVPCGIDKAN



787 TLGVLLSYLNIAREEVIAIGDGVCDVNMLQVAGLGIAMGHAQDSVKVCADYVTASNEEDGVAQSVEKLILAEVHAAEIPDLLNERARHALMGNLGI  
788 QYTYASEERIEATMPVDHRTRQPFILHGGATLALAETVAGLGSMITCQPDEIVVGMQVSGNHSSAHEGDTVRAVATIVHKGRSSHVWNVDTVSTN  
789 KLVSSVRVNSVLKRR  
790 >211953\_NYSGRC (non-crystallizable)  
791 MSFSSRAESDGMRGNAAKHADNQRVRERRMAEQAAEYGLQEARQMLQEAFAFPWVLDLGLSIEALELDPAGSPPDWQPGATLRMAFSERLCRSG  
792 GVICGQALMALADTAMVFAVCAGYRGFRPMTTVDQTTFLKAVASTDVIADARLVRLGRTMSFGRVTLGASDRKPVAMVSSAFAML  
793 >212179\_NYSGRC (non-crystallizable)  
794 MSDPSASFELPVRVYIEDTDAGGIVFHAKYLHYMERARTEWVRSQGVGLRAGLEHNISYVVQKMNLHFRMPAKLDDQLLVTAELKAASRVWVGFR  
795 QCVYRADDRQLLCDADVACVALDTGKPRRLPENMQEILKNFV  
796 >212134\_NYSGRC (non-crystallizable)  
797 MDVPGGSAQAWPTRCPTPRLVMTQEPFSPGSTARVELVTVADTAQAIGSGDVPVLGTPRILALAEAATVAAIARQLPSGATTGVRVELDHQAATPVG  
798 RTVVARARLAEVDGRLLFAVSVTEDGSTVAEGRVERLLVDRQRFIERAGRSS  
799 >019613\_NYSGRC (non-crystallizable)  
800 MQPCNQSAPCWPFQTPSQTLQTLSPALMLRCGLKVVSAARQATTRMTDPDPNHTAATDGIPDGFVRHARSSPLTAPWEPLYAKTTADAVTLGLRI  
801 RECHTSNRGLAHGGLITALADNAMGYSCGLKLGQGLTSSLAIDFIGPAKIGQWLQIEPEVIKLGAKLCVAQCFVTADGTRCARANGTFSVVKAKE  
802 >212126\_NYSGRC (non-crystallizable)  
803 MLASTVEKLSGRITLVNAIEADSSRSKTVTWSPLVGAELAKTMSGLAYMQAMIDGKIPPPISGLMNMATAVSAETGLVTFACPTDESQYNPIGTVHGGL  
804 VCTLLDSVCGCAVQTTLPAGQSYTSLEIKINYLRPVLAHTGELIAVGRVTKPGSSAAFAEGEIRDQAGKLIATASSTLLVFPV  
805 >212124\_NYSGRC (non-crystallizable)  
806 MQEIHKDLLTFVEQSIPFHKLLGIRVEHAVPGFARVRLPYQDAFCGNMARGALHGGVTAVLVDICGAVALWTHFGPLDKTATIDMRVDYQRPAPFDDLL  
807 AEGEVVMGNRIASVHVRVTAAPDQLIAEGRCVYYVVRVPPQPETAGTQE  
808 >212129\_NYSGRC (non-crystallizable)  
809 MSSRHKVEAMTHQATAAEESEIPGAHQGAAPGEGPGEIPGKPTSASRTTLSSHIMTHNDTNLLGTVHGGVIMKLVVDAAGAVAGRHSQGPVAVTASMDM  
810 VFLEPVRVGDVHVKVQVNWGTGRSMEVGVVLAERWNEAPATQVGSAYLVFAAVDADGKPRRVPPVLPETERDKRRYQEAQIRRTHLARRRAIM  
811 DLREKRAAEGLDD  
812 >GO\_111497\_MPP (non-crystallizable)  
813 MLRSCAARLRTLALCLPPVGRRLPGSEPRPELRSFSSEEVILKDCSVNPSWNKDLRLFLDQFMKKCEDGSKRLPSYKRTPTEWIQDFKTHFLDPKL  
814 MKEEQMSQAQLFTRSFDDGLGFYVMFYNDIEKRMVCLFQGGPYLEGPPGFHGGAIATMIDATVGMCAMMAGGIVMTANLNINIKRPIPLCSVMI  
815 NSQLDKVEGRKFFVSCNVQSVDEKTYSEATSLFIKLNPAKSLT  
816 >212223\_NYSGRC (non-crystallizable)  
817 MNHSRCNAVERAMHSFHPQTPDWEPRVRSFARQGLMQALHAVIEHLKPGVAITMPADPTYSQQHGVIHGGAIASILDSACGYAALTLPVGREVLT  
818 VEFKVNFLSPARGRFLAVGRVVRAGKTVTVCAAGEAFTVDGDRRVPALMQATMMAVPEMPERAAH  
819 >212232\_NYSGRC (non-crystallizable)  
820 MVSGGGQEHHTELRVRYAETDAMAVAHATYPVWFVARTELMHALGLPYTEMETRGYYLMLSLGHVQYRRAARYDDRLDITRITRTEIRSRTLKFAY  
821 EVHRIGADGTRELLATGETHHIATDHQYRPSRMPDDVLALLAGEG

822 >211894\_NYSGRC (non-crystallizable)  
823 MHDMMHFTKQYPIRFSFCDPAGIVYFPQYLVLNSWLIEDWFNEGLGIDFASFIGHRRRLGLPIVKLNCEFISPSHHGDTLTLQLRVAKLGQRSITLDLEGHVD  
824 MIMRLRCQQVLVFTSLDTEKSTPMPDDVDGALRALLSQEDIA  
825 >212180\_NYSGRC (non-crystallizable)  
826 MSDQSIQDPVMSFDDKARMIQHRRSIHGAIIGLQLDRYAPAEAWSSLPYHPVFGDVSTGVIHGGVVTAMLEDESCGMAVQLALPGTTAIATLDRIDYL  
827 RPATPGQVMRAHAHCYHLTRSIAFVRATAYQDAEDVPIATATAMFMVGANRTDMLRQTPKVTMDSAPELVAPEDPDGGPLAISPYPRFLGIRVDGDAQ  
828 AMMPYAPKLVGNPILPALHGGVIGAFLETAIVSVREIGLATAPKPIGLTVNYLRSRPLDTFAKVSIVKQGRRVVAFEAQAYQRDPAEPIASCYGHFK  
829 LRSGPAE  
830 >211972\_NYSGRC (non-crystallizable)  
831 MPEPINPLTQIRAINETAPFNHFFGIEVKSAGVGVVELSMPWRPEAGQYSGFLHAGVIGALIDTACGFAAATLVGPVLASHYSVNCLRPVAVGESFLARAR  
832 VVKPGKSQVFTSCEVFALLDGSEKLVATGETLLSVVQDKT  
833 >366861\_JCSG (crystallizable)  
834 MSDDLTAQTAAIPEGFSQLNWSRGFGRQIGPLFEHREGPGQARLAFRVEEHHTNGLGNCHGGMLMSFADMAWGRIISLQKYSWVTVRLMCDFLS  
835 GAKLGDWVEGEGELISEEDMLFTVRGRIWAGERTLITGTGVFKALSARKPRPGELAYKEEA  
836 >211883\_NYSGRC (non-crystallizable)  
837 MEIGAMFERIPFAAELGIEFDEVADGHAEGRLPLREEHSSNPRQIAHGGVTFSLADTVGGAADVSKSESVSPTIDMRIDYLAPATADLRAVADVVRAGE  
838 SVTAVDIEVYDADDHHVASARGVYKTGGQGEETPWTDTGTDVEPASDGAEQRSEE  
839 >212139\_NYSGRC (non-crystallizable)  
840 MSTEEKPRVTDAEMLARFQNSKKRPPCSETLGMRLADLNQDKQWVKMEFDVSPSFANPTGAVQGGFIAAMLDEAMSTAVIIASNVMTAPTLEMKTS  
841 YLRRMLPGKASVEARILKLGKSAAFMEADCFDAEGKLVARATATAIPMAFKRL  
842 >212176\_NYSGRC (non-crystallizable)  
843 MENYTSLIRLIRISAHDAHAYAGGLVDGARMHLHFGDVATELLIRSDGDEGLFVAYDEVQFLAPVHAGDYIEASGRIVAMGKTSRKMVFARKVIVPAGM  
844 AGQPSAADVLAELPLVCKASGTCVVPLHCQRMARPPVRC  
845 >212050\_NYSGRC (non-crystallizable)  
846 MNAPEPAPSASYFHRPFQQLVGYDILHNERGLYFRLPIRRDHLNPHGVLHGGVPLTLLDAVGGRTLIDRRIPGSDQRILSSVTVTLTVDFMRAIGSGVLF  
847 ASATPDHIGKTLAYVSMKVTLDDLDGDIVSHGIGTYRIYTKSLVKHV  
848 >211971\_NYSGRC (non-crystallizable)  
849 MTAPTILTDRLPPYARAMGMRIEGLLDGAPLLAMDFSDRAMGRPGYLHGGAIAGMLEIAAIMALHADLGAEDASVRIKPVNISVEYLRGGITVETFAR  
850 GEVIRAGRRIANVRAEAWQADRDKPLASCWMNFLIKPKG  
851 >211895\_NYSGRC (non-crystallizable)  
852 MPIALAGGGFFRVSAPPARTPYPEAMSYAEVLGMILTLDASPDLTRVALTVTEAGLNMHGTAHGGILFSLADEAFAVISNLDAQAVAAETHMSFFRAA  
853 REGERLVAVATPERVGRTLATYRIEVRRGEEGEVLALFLGTVSRREKQS

## References

- 855 [1] M. J. Gabanyi, P. D. Adams, K. Arnold *et al.* The structural biology knowledgebase: a portal to protein structures,  
856 sequences, functions, and methods. *J Struct Funct Genomics* 2011;12:45-54.
- 857 [2] H. M. Berman, J. D. Westbrook, M. J. Gabanyi *et al.* The protein structure initiative structural genomics  
858 knowledgebase. *Nucleic Acids Res* 2008;37:D365-D368.
- 859 [3] M. J. Mizianty, and L. Kurgan. Sequence-based prediction of protein crystallization, purification and production  
860 propensity. *Bioinformatics* 2011;27:i24-i33.
- 861 [4] W. Li, and A. Godzik. CD-HIT: a fast program for clustering and comparing large sets of protein or nucleotide  
862 sequences. *Bioinformatics* 2006;22:1658-1659.
- 863 [5] K. Dániel, S. István, and G. E. Tusnády. PDBTM: protein data bank of transmembrane proteins after 8 years. *Nucleic  
864 Acids Res* 2013;41:524-529.
- 865 [6] K. Chen, L. Kurgan, and M. Rahbari. Prediction of protein crystallization using collocation of amino acid pairs.  
866 *Biochem Biophys Res Commun* 2007;355:764-769.
- 867 [7] I. Overton, G. Padovani, and M. Girolami, G. ParCrys: a Parzen window density estimation approach to protein  
868 crystallization propensity prediction. *Bioinformatics* 2008;24:901-907.
- 869 [8] S. Jahandideh, and A. Mahdavi. RFCRYS: Sequence-based protein crystallization propensity prediction by means  
870 of random forest. *J Theor Biol* 2012;306:115-119.
- 871 [9] C. Ding, L. F. Yuan, S. H. Guo *et al.* Identification of mycobacterial membrane proteins and their types using over-  
872 represented tripeptide compositions. *J Proteomics* 2012;77:321-328.
- 873 [10] K. Chou, and M. P. Com. Prediction of protein cellular attributes using pseudo-amino acid composition. *Proteins:  
874 Struct Funct Bioinf* 2010;43:246-255.
- 875 [11] K. C. Chou. Using amphiphilic pseudo amino acid composition to predict enzyme subfamily classes. *Bioinformatics*  
876 2005;21:10-19.
- 877 [12] K. C. Chou, and H. B. Shen. MemType-2L: a web server for predicting membrane proteins and their types by  
878 incorporating evolution information through Pse-PSSM. *Biochem Biophys Res Commun* 2007;360:339-345.
- 879 [13] J. Hu, K. Han, Y. Li *et al.* TargetCrys: protein crystallization prediction by fusing multi-view features with two-  
880 layered SVM. *Amino Acids* 2016;48:1-15.
- 881 [14] L. Wei, J. Tang, and Q. Zou. Local-DPP: an improved DNA-binding protein prediction method by exploring local  
882 evolutionary information. *Inf Sci* 2016;384:135-144.
- 883 [15] J. Hu, X. G. Zhou, Y. H. Zhu *et al.* TargetDBP: accurate DNA-binding protein prediction via sequence-based multi-  
884 view feature Learning. *IEEE/ACM Trans Comput Biol Bioinf* 2019;(In press).
- 885 [16] A. A. Schäffer, L. Aravind, T. L. Madden *et al.* Improving the accuracy of PSI-BLAST protein database searches  
886 with composition-based statistics and other refinements. *Nucleic Acids Res* 2001;29:2994-3005.
- 887 [17] A. Bairoch, and R. Apweiler. The SWISS-PROT protein sequence database and its supplement TrEMBL in 2000.  
888 *Nucleic Acids Res* 2000;28:45-48.
- 889 [18] I. M. Overton, and G. J. Barton. A normalised scale for structural genomics target ranking: the OB-score. *Febs Lett*  
890 2006;580:4005-4009.
- 891 [19] L. Kurgan, A. A. Razib, S. Aghakhani *et al.* CRYSTALP2: sequence-based protein crystallization propensity  
892 prediction. *BMC Struct Biol* 2009;9:50.
- 893 [20] K. Krishna Kumar, P. Ganesan, P. N. Suganthan *et al.* SVMCRYST: an SVM approach for the prediction of protein  
894 crystallization propensity from protein sequence. *Protein Pept Lett* 2010;17:423-430.
- 895 [21] F. Meng, C. Wang, and L. Kurgan. fDETECT webserver: fast predictor of propensity for protein production,  
896 purification, and crystallization. *BMC Bioinf* 2017;18:580.
- 897 [22] A. Elbasir, B. Moovarkumudalvan, K. Kunji *et al.* DeepCrystal: a deep learning framework for sequence-based

898 protein crystallization prediction. *Bioinformatics* 2018;35:2216-2225.

899 [23] Y. Zhang. I-TASSER server for protein 3D structure prediction. *BMC Bioinf* 2008;9:40.

900 [24] A. Roy, A. Kucukural, and Y. Zhang. I-TASSER: a unified platform for automated protein structure and function  
901 prediction. *Nat Protoc* 2010;5:725.

902 [25] J. Yang, R. Yan, A. Roy *et al.* The I-TASSER Suite: protein structure and function prediction. *Nat Methods* 2015;12:7.

903 [26] Y. Zhang, and J. Skolnick. TM-align: a protein structure alignment algorithm based on the TM-score. *Nucleic Acids*  
904 *Res* 2005;33:2302-2309.

905 [27] Y. Zhang, and J. Skolnick. Scoring function for automated assessment of protein structure template quality. *Proteins:*  
906 *Struct Funct Bioinf* 2004;57:702-710.

907 [28] J. Xu, and Y. Zhang. How significant is a protein structure similarity with TM-score=0.5? *Bioinformatics*  
908 2010;26:889-895.

909 [29] H. Wang, L. Feng, Z. Zhang *et al.* CrysAlis: an integrated server for computational analysis and design of protein  
910 crystallization. *Sci Rep* 2016;6:21383.

911 [30] M. J. Mizianty, and L. Kurgan. Meta prediction of protein crystallization propensity. *Biochem Biophys Res Commun*  
912 2009;390:10-15.

913 [31] L. Slabinski, L. Jaroszewski, L. Rychlewski, Ia *et al.* XtalPred: a web server for prediction of protein crystallizability.  
914 *Bioinformatics* 2007;23:3403-3405.

915 [32] P. Charoenkwan, W. Shoombuatong, H. C. Lee *et al.* SCMCRYIS: predicting protein crystallization using an  
916 ensemble scoring card method with estimating propensity scores of p-collocated amino acid pairs. *Plos One*  
917 2013;8:e72368.

918 [33] J. Gao, H. Gang, Z. Wu *et al.* Improved prediction of protein crystallization, purification and production propensity  
919 using hybrid sequence representation. *Curr Bioinf* 2014;9:57-64.

920