

1

2 Origin of a folded repeat protein from an
3 intrinsically disordered ancestor

4

5 Hongbo Zhu[□], Edgardo Sepulveda, Marcus D. Hartmann, Manjunatha Kogenaru[†],

6 Astrid Ursinus, Eva Sulz, Reinhard Albrecht, Murray Coles, Jörg Martin,

7 Andrei N. Lupas[§]

8

9 Department of Protein Evolution,

10 Max Planck Institute for Developmental Biology,

11 Spemannstr. 35, D-72076 Tübingen, Germany

12

13

14

15

16 [□] hongbo.zhu@tuebingen.mpg.de

17 [†] present address: Department of Life Sciences, Imperial College London, London SW7 2AZ,

18 United Kingdom. E-mail: m.kogenaru@imperial.ac.uk

19 [§] andrei.lupas@tuebingen.mpg.de; Corresponding author

20 **Abstract**

21 **Repetitive proteins are thought to have arisen through the amplification of subdomain-sized**
22 **peptides. Many of these originated in a non-repetitive context as cofactors of RNA-based**
23 **replication and catalysis, and required the RNA to assume their active conformation. In search**
24 **of the origins of one of the most widespread repeat protein families, the tetratricopeptide**
25 **repeat (TPR), we identified several potential homologs of its repeated helical hairpin in non-**
26 **repetitive proteins, including the putatively ancient ribosomal protein S20 (RPS20), which only**
27 **becomes structured in the context of the ribosome. We evaluated the ability of the RPS20**
28 **hairpin to form a TPR fold by amplification and obtained structures identical to natural TPRs**
29 **for variants with 2-5 point mutations per repeat. The mutations were neutral in the parent**
30 **organism, suggesting that they could have been sampled in the course of evolution. TPRs**
31 **could thus have plausibly arisen by amplification from an ancestral helical hairpin.**

32

33 **1 Introduction**

34 Most present-day proteins arose through the combinatorial shuffling and differentiation of a set of
35 domain prototypes. In many cases, these prototypes can be traced back to the root of cellular life
36 and have since acted as the primary unit of protein evolution (Anantharaman et al., 2001; Apic et al.,
37 2001; Koonin, 2003; Kyripides et al., 1999; Orengo and Thornton, 2005; Ponting and Russell, 2002;
38 Ranea et al., 2006). The mechanisms by which they themselves arose are however still poorly
39 understood. We have proposed that the first folded domains emerged through the repetition, fusion,
40 recombination, and accretion of an ancestral set of peptides, which supported RNA-based replication
41 and catalysis (the RNA world (Bernhardt, 2012; Gilbert, 1986)) (Alva et al., 2015; Lupas et al., 2001;
42 Söding and Lupas, 2003). Repetition would have been a particularly prominent mechanism by which
43 these peptides yielded folds; six of the ten most populated folds in the Structural Classification of
44 Proteins (SCOP) (Murzin et al., 1995) – including the five most frequent ones – have repetitive
45 structures. In all cases, their amplification from subdomain-sized fragments can also be retraced at
46 the sequence level in at least some of their members.

47 One of these highly populated repetitive folds is the $\alpha\alpha$ -solenoid (SCOP a.118), whose most
48 widespread superfamily is the tetratricopeptide repeat (TPR; a.118.8). This was originally identified
49 as a repeating 34 amino-acid motif in Cdc23p of *Saccharomyces cerevisiae* (Sikorski et al., 1990) –
50 hence its name. Since then, TPR-containing proteins have been discovered in all kingdoms of life,
51 where they mediate protein-protein interactions in a broad range of biological processes, such as cell
52 cycle control, transcription, protein translocation, protein folding, signal transduction and innate
53 immunity (Cortajarena and Regan, 2006; Dunin-Horkawicz et al., 2014; Katibah et al., 2014; Keiski et
54 al., 2010; Kyripides and Woese, 1998; Lamb et al., 1995; Sikorski et al., 1990). The first crystal
55 structure of a TPR domain (Das et al., 1998) showed that the repeat units are helical hairpins, stacked
56 into a continuous, right-handed superhelical architecture with an inner groove that mediates the
57 interaction with target proteins (Forrer et al., 2004). The hairpins interact via a specific geometry

58 involving knobs-into-holes packing (Crick, 1953) and burying about 40% of their surface between
59 repeat units. This tightly packed, superhelical arrangement of a repeating structural unit is typical of
60 all α -solenoid proteins (Di Domenico et al., 2014; Kajava, 2012; Kobe and Kajava, 2000).

61 Comparison of TPRs from a variety of proteins reveals a high degree of sequence diversity, with
62 conservation observed mainly in the size of the repeating unit and the hydrophobicity of a few key
63 residues (D'Andrea and Regan, 2003; Magliery and Regan, 2004). Nevertheless, almost all known
64 TPR-containing proteins can be detected using a single sequence profile (Karpenahalli et al., 2007),
65 underscoring their homologous origin. As their name implies, TPR proteins generally contain at least
66 two unit hairpins in a repeated fashion. The few that have only one hairpin, notably the
67 mitochondrial import protein Tom20 (Abe et al., 2000), are clearly not ancestral based on their
68 phylogenetic distribution and functionality, implying that the ancestor of the superfamily already had
69 a repeated structure. In searching for the origin of TPRs, we hypothesized that the hairpin at the root
70 of the fold might either have been part of a different, non-repetitive fold or have given rise to both
71 repetitive and non-repetitive folds at the origin of folded domains. Either way we hoped that we
72 might find α -hairpins in non-repetitive proteins that are similar in both sequence and structure to the
73 TPR unit, suggesting a common origin. Here we show that such hairpins are detectable and that one
74 of them, from the ribosomal protein RPS20 (Schluenzen et al., 2000), can be customized to yield a
75 TPR fold by repetition, with only a small number of point mutations that are neutral for the parent
76 organism. Ribosomal proteins most likely constitute some of the oldest proteins observable today
77 and are still intimately involved in an RNA-driven process: translation (Fox, 2010; Hsiao et al., 2009).
78 They are mostly incapable of assuming their folds outside the ribosomal context (Peng et al., 2014)
79 and thus belong to a class of intrinsically disordered proteins that become structured upon binding to
80 a macromolecular scaffold (Dyson and Wright, 2005; Habchi et al., 2014; Oldfield and Dunker, 2014;
81 Peng et al., 2014; Varadi et al., 2014). This hairpin therefore plausibly retains today many of the
82 properties likely to have been present in the ancestral peptide that gave rise to the TPR fold.

83 **2 Results and Discussion**

84 **2.1 Recently amplified TPR arrays in present-day proteins**

85 Repetitive folds with variable numbers of repeats, such as HEAT, LRR, TPR or β -propellers, usually
86 have some members with a high level of sequence identity between their repeat units (Dunin-
87 Horkawicz et al., 2014). In these proteins, the units are more similar to each other than to any other
88 unit in the protein sequence database, showing that they were recently amplified. In a detailed study
89 of β -propellers (Chaudhuri et al., 2008), we found that this process of amplification and
90 differentiation has been ongoing since the origin of the fold. TPR proteins show a similar evolutionary
91 history. In some proteins, most of the repeats can be seen to have been amplified separately and to a
92 different extent in each ortholog, pointing to their recent origin (Figure 1a); in others, the
93 amplification must have occurred much earlier, as their ancestor already had fully differentiated
94 repeats (Figure 1b). In recently amplified proteins, such as the ones shown in Figure 1a, within which
95 repeats frequently have >80% pairwise sequence identity, tracking the probable α -hairpin at the root
96 of the amplification is a fairly straightforward proposition. We wondered, however, whether it might
97 be possible to go much further back in time and track the original α -hairpin from which the first TPR
98 protein was amplified. We therefore searched for TPR-like α -hairpins in non-repetitive proteins as
99 present-day descendants of the original hairpin.

100 **2.2 Identification of helical hairpins resembling the TPR unit**

101 We had previously developed a profile-based method, named TPRpred, specially designed for the
102 detection of TPRs and related repeat proteins with high sensitivity from sequence data (Karpenahalli
103 et al., 2007). Here, in a first step, we used TPRpred to scan protein sequences in the Protein Data
104 Bank (PDB) (Berman et al., 2000) for peptides that share statistically significant similarity to the TPR
105 sequence profile and yet have not been annotated as TPR in Pfam (Finn et al., 2014); we used a P -
106 value cutoff=1.0e-4, which leads to an estimated false discovery rate of 1.0%, see Materials and
107 Methods. We ignored tandem repeats in the hit list and focused only on the singleton cases.

108 Subsequently, we compared the structures of these helical hairpin singletons to the average TPR
109 hairpin and removed non-hairpin-like structures. This yielded 31 helical hairpins that are similar to
110 the TPR unit with respect to both sequence and structure. Among them, 21 are part of solenoid-like
111 structures and were discarded. The remaining nine hits belong to three families: (I) mitochondrial
112 import receptor subunit Tom20; (II) microtubule interacting and transport (MIT) domain including
113 katanin (Iwaya et al., 2010); and (III) 30S ribosomal protein S20 (RPS20) (Figure 2).

114 The similarity of Tom20 and MIT domains to TPR proteins has been noted before (Abe et al., 2000;
115 Iwaya et al., 2010; Scott et al., 2005), but the similarity of RPS20 was surprising and drew our
116 attention particularly due to the ancestry attributed to ribosomal proteins. To further explore the
117 similarity between the helical hairpin in RPS20 (in short, RPS20-hh) and TPR, we used TPRpred to
118 rank the RPS20 sequences in Pfam (Finn et al., 2014). The top-scoring hit was RPS20-hh from
119 *Thermus aquaticus* (NCBI accession number=WP_003044315.1, UniProt id=B7A5L8_THEAQ), which
120 matches the TPR unit sequence profile at a P -value of $5.4e-07$, almost an order of magnitude better
121 than the second hit (see supplementary file 1D). Furthermore, we examined the surface residues of
122 RPS20-hh fragments to assess their suitability to occur in a tandem repeat mode, as in TPRs. To this
123 end, we first defined five interface positions on the TPR helical hairpin and transferred the definition
124 to RPS20-hh according to their structure alignment (positions 3, 7, 10, 21 and 28 using TPR unit
125 numbering). Then, we searched for RPS20-hhs with as many hydrophobic residues as possible at
126 these interface positions. We found 42 RPS20-hhs that contain at least three hydrophobic residues
127 out of the five interface positions. Among them, the only RPS20-hh predicted to match the TPR unit
128 profile above a P -value of $1.0e-4$ was again the RPS20 from *T. aquaticus*, in which three of the five
129 interface residues are hydrophobic (L10, I21 and V28). We therefore chose this helical hairpin
130 (RPS20-hhta) to construct a TPR-like solenoid by amplification (Figure 3).

131 **2.3 Design of a TPR array from a RPS20**

132 We focused on the construction of three-repeat TPRs, which represent the most common form of
133 this fold (D'Andrea and Regan, 2003; Sawyer et al., 2013). For instance, 18 of the 54 non-identical

134 TPR domains in the extended Structural Classification of Proteins database (SCOPe v2.05) (Fox et al.,
135 2014) have three repeats. A previously designed three-repeat TPR protein, CTPR3, was also
136 demonstrated to be highly stable, even more so than natural three-repeat TPR proteins (Main et al.,
137 2003b). We concatenated three copies of RPS20-hhta as an initial construct, connected by the TPR
138 consensus loop sequence (DPNN). We annotate the two helices in each repeat unit as helix A_i and B_i ,
139 where i is the index of the repeat unit ($i=1, 2$ or 3) (Figure 3). Under the hypothesis of common
140 descent between TPR and RPS20 from the same ancestral peptide and retention of ancestral features
141 in RPS20, this basic construct would fold as a TPR solenoid with a minimal number of mutations,
142 ideally none.

143 When we experimentally made a construct containing no mutations (M0, Table 1), it was soluble but
144 remained unfolded under all conditions tested (see Section 2.4). We therefore introduced point
145 mutations into the sequence of RPS20-hhta, aimed at favoring the target structure. Here, we
146 followed the principle of consensus design (Forrer et al., 2004; Main et al., 2003a), which requires
147 the mutation positions to be occupied by the most commonly observed residues in homologous
148 proteins (Forrer et al., 2004). Consensus design methods have been successful in engineering several
149 different repeat proteins with solenoid folds, including ankyrin repeats (Binz et al., 2003; Kohl et al.,
150 2003; Mosavi et al., 2002), TPRs (Doyle et al., 2015; Kajander et al., 2007; Main et al., 2003b),
151 pentatricopeptide repeats (PPRs) (Coquille et al., 2014; Shen et al., 2016) and leucine rich repeats
152 (Rämisch et al., 2014; Stumpp et al., 2003). Following these principles, four different sites of
153 mutation (L4W, K7L/R, V9N, I23D/Y, see Figure 4) were considered to improve interface
154 hydrophobicity or preserve coevolved positions observed in TPRs (Sawyer et al., 2013) (see Materials
155 and Methods). Furthermore, as natural TPR proteins tend to exhibit zero net charge (Magliery and
156 Regan, 2004), four positively charged residues were also targeted (K2E, K6N, K22E, R25Q/E, see
157 Figure 4). This resulted in a set of eight candidate mutation sites. In order to preserve the character
158 of the RPS20-hhta sequence, we restricted the number of mutations in any repeat unit to be at most
159 five.

160 In most TPR proteins, there is an α -helix at the C-terminus, which interacts with the last TPR unit by
161 covering the hydrophobic surface. This so-called C-terminal “stop helix” had been observed in all
162 known TPR structures and was considered essential for the solubility of natural TPR proteins
163 (D'Andrea and Regan, 2003; Das et al., 1998; Main et al., 2003b). Most other designed TPRs employ
164 purpose-designed stop helix sequences. Here, we chose to use the RPS20 C-terminal helix to become
165 a natural stop helix, since it is already known to interact favorably with RPS20-hhta (Figure 3).
166 Further, we inserted two residues (Asn-Ser) before the first TPR unit as an N-terminal cap to the first
167 helix (Aurora and Rose, 1998; Kumar and Bansal, 1998), in analogy to a previously designed idealized
168 TPR protein, CTPR3 (Main et al., 2003b).

169 To model the structure of the designed proteins *in silico*, we fused two structures to create a hybrid
170 template: We used CTPR3 (PDB id: 1na0 chain A) as the structural template for the three RPS20-hhta
171 fragments, and the best-resolved RPS20 structure (PDB id: 2vqe chain T; 2.5 Å) for helix B3 and the
172 stop helix. We built structural models on this hybrid template and tested a variety of mutants using
173 the Rosetta programs *fixbb* and *relax*, which perform fixed-backbone design and structural
174 refinement (Das and Baker, 2008; Doyle et al., 2015; Park et al., 2015; Parmeggiani et al., 2015). The
175 Rosetta energy score of the models calculated for all mutants is depicted in a boxplot (Figure 4—
176 figure supplement 2). Among them, five were selected for further testing *in vitro* (see Materials and
177 Methods). These five tested mutants are termed M2, M4E, M4N, M4RD and M5. Their primary
178 structures are listed in Table 1.

179 **2.4 Biophysical characterization of designed TPRs and RPS20**

180 We cloned the five TPR designs plus the unmutated construct M0 into pET vectors for expression in
181 *Escherichia coli*. Three proteins (M0, M4RD and M5) could be purified from soluble extracts; the
182 other constructs were insoluble and were refolded from inclusion bodies. In far UV circular dichroism
183 (CD) spectra, all proteins displayed a strong alpha-helical pattern, except M0 and M4RD, which
184 appeared to be unfolded, but not prone to aggregation and precipitation, even at high
185 concentrations. When we studied the melting curves, M4N showed cooperative unfolding with a T_m

186 of 77°C (supplementary file 1F), while the unfolding of M2, M4E and M5 did not conform to a
187 classical two-state transition, consistent with an unstable molten globule-like state. On the other
188 hand, non-cooperative unfolding processes have been demonstrated for perfectly stable TPR repeats
189 and suggested to be common for various types of repeat proteins (Cortajarena and Regan, 2006;
190 Kajander et al., 2007; Stumpp et al., 2003). To clarify this point, urea-induced unfolding transitions
191 were monitored by CD. Like M4N, the three variants M2, M4E and M5 yielded typical cooperative
192 denaturation curves, indicative of folded polypeptides (Figure 5—figure supplement 2). The $\Delta G_{U-F}^{H_2O}$
193 values agree well with those reported for other designed TPRs (supplementary file 1F) (Main et al.,
194 2005). In line with these findings, M5, the only protein containing tryptophan residues, had a λ_{max} of
195 336 nm in fluorescence emission spectra, as expected for partially shielded aromatic residues. We
196 conclude that four of the five designed TPR variants, M2, M4E, M4N and M5, result in well-folded
197 repeat proteins. To determine the oligomeric state of our folded proteins, we performed static light
198 scattering experiments. Surprisingly, all four constructs were exclusively dimers (supplementary file
199 1F).

200 We also examined the ribosomal parent protein RPS20. Within the ribosome, RPS20 is partially
201 embedded in the 16S rRNA, making many nucleic acid contacts. Like many other ribosomal proteins,
202 it is not expected to adopt a stable structure in isolation. Indeed, it has a biased amino acid
203 composition and is predicted to be largely unstructured by many prediction programs (Figure 4—
204 figure supplement 1, see also Supplementary file 1J). It had been shown previously that isolated
205 RPS20 exhibits only one third helical content by CD (Paterakis et al., 1983). For *Thermus* RPS20
206 specifically, simulations predict a flexible conformation in solution (Burton et al., 2012). We cloned
207 RPS20 from *T. aquaticus* and its close relative *T. thermophilus*. Upon expression, both proteins were
208 insoluble and had to be refolded. In static light scattering measurements, both proteins behaved as
209 monomers (supplementary file 1F). Based on CD spectra, which showed a high proportion of random
210 structure, and the absence of defined melting and urea-denaturation curves (supplementary file 1F),
211 we conclude that RPS20 indeed exhibits considerable conformational variation in solution.

212 **2.5 Structure of a designed TPR**

213 To obtain high-resolution structural information on our designed proteins, we set up crystallization
214 trials for all four folded constructs. We obtained crystals and solved the structure of M4N to a
215 resolution of 2.2 Å (Figure 5a). The asymmetric unit (ASU) contains three polypeptide chains of
216 almost identical structure (all pairwise C_α RMSD values below 1.4 Å). Notably, all three chains exhibit
217 the desired TPR architecture with three repetitive hairpins, which interact via knobs-into-holes
218 packing between helices A_i and B(i-1), as is characteristic of TPR hairpins. A superposition to the
219 CTPR3 structure yields C_α RMSD values below 2.6 Å (supplementary file 1I). An unexpected difference
220 to canonical TPR structure is that the stop helix of M4N is not resolved in any of the three chains.
221 However, this missing helix is compensated for by a specific dimerization mode of two M4N
222 protomers. Therein, the C-terminal TPR units of the two protomers form a tight interface, in which
223 the B3 helix of each chain substitutes for the stop helix of the other, mimicking the capping effect of
224 the stop helix (Figure 6). A superposition of this mimicry to the last TPR unit and stop helix of CTPR3
225 yields C_α RMSD values as low as 1.2 Å over 44 residues. The third chain of the ASU, however, was
226 found as a monomer, capping its C-terminal TPR unit in a more unspecific manner by packing it
227 orthogonally against the two A1 helices of the dimer (Figure 5a).

228 Analysis by mass spectrometry revealed that the M4N stop helix had been partially proteolyzed upon
229 expression of the protein (Figure 5—figure supplement 3). Although we did not observe proteolysis
230 in the other folded constructs (M2, M4E and M5), which were also all dimeric, we analyzed whether
231 proteolysis might have favored the dimerization of M4N. Extending the stop helix with a C-terminal
232 His₆-tag prevented proteolysis, but did not affect stability or dimerization (M4N-His; supplementary
233 file 1F). We conclude that in the amplified constructs, the observed interactions are more favorable
234 than the interaction with the native stop helix, releasing it and rendering it prone to degradation.
235 This led us to ask whether this helix is in fact dispensable. Indeed, an M4NΔC construct, which
236 terminates with the B3 helix, showed the same stability and dimerization as M4N. We obtained two
237 structures for M4NΔC from different crystal forms at 2.0 Å and 1.7 Å resolution, respectively, the first

238 (CF I) with two dimers in the ASU and the second (CF II) with a single chain in the ASU, for which we
239 constructed the dimer by crystallographic symmetry. All three dimers superimpose to the M4N dimer
240 with C_α RMSD below 1.9 Å (Figure 7, supplementary file 1I). We conclude that the stop helix is
241 dispensable for folding, dimerization and stability of our designed constructs.

242 The geometry of dimerization in M4N has not been observed in TPR structures before. Although
243 there have been reports on the self-association of TPR-containing proteins involved in various
244 regulatory biological processes (Bansal et al., 2009a; Bansal et al., 2009b; Ramarao et al., 2001;
245 Serasinghe and Yoon, 2008), only a small number of oligomeric TPR structures have been determined
246 to date (Krachler et al., 2010; Lunelli et al., 2009; Zeytuni et al., 2012; Zeytuni et al., 2015; Zhang et
247 al., 2010). None of these resemble the ring-shaped dimer of M4N.

248 **2.6 Mutations introduced into RPS20-hhta are neutral to Thermus**

249 The results shown above suggest that the mutations we made to RPS20-hhta were crucial for
250 obtaining the TPR fold. If RPS20 and TPR proteins indeed share a common ancestor, such mutations
251 may have been sampled in the course of evolution. Since we cannot reconstruct the ancestor and do
252 not know what its function was beyond a general expectation of RNA binding, we decided to test
253 whether the mutations we introduced impaired the interaction between RPS20 and its cognate RNA,
254 as an indication of their compatibility with RNA interaction. Each mutation in M2 and M4N occurs in
255 natural RPS20 sequences (see supplementary file 1A), but no RPS20 sequence has all four mutations
256 simultaneously and we therefore tested if they can be tolerated *in vivo*. As genetic engineering in *T.*
257 *aquaticus* turned out to be unfeasible, we performed these tests in *T. thermophilus* HB8, which is a
258 well-established model organism. The RPS20 helical hairpins in *T. aquaticus* and *T. thermophilus*
259 differ only at four positions, of which two are highly conservative substitutions (Figure 8a).

260 We first attempted to substitute the chromosomal RPS20-encoding gene, *rpsT*, with a kanamycin
261 resistance cassette, to obtain *T. thermophilus* strain KM4 (Figure 8b). For complementation we
262 introduced plasmids bearing wild type *rpsT* from *T. thermophilus* (TT) or *T. aquaticus* (TA), *rpsT* from

263 *T. aquaticus* carrying the mutations from M2 (TA2) or M4N (TA4), or merely empty plasmids as
264 negative control (E). We monitored the substitution of *rpsT* by a PCR screening protocol, which will
265 amplify a 1500 bp region if WT *rpsT* is substituted and an 800 bp region otherwise (Figure 8b). Under
266 selection pressure from kanamycin, only the 1500 bp product was obtained in all cases where
267 plasmid-borne *rpsT* was introduced, whether in wild-type or mutated form (Figure 8c panels 1 and 2,
268 lanes TT, TA, TA2 and TA4), showing that the chromosomal gene had been fully substituted. In
269 contrast, PCR screening of strain KM4 complemented with an empty plasmid produced both 800 bp
270 and 1500 bp fragments (Figure 8c panels 1 and 2, lane E). Since *T. thermophilus* HB8 is a polyploid
271 organism (minimally tetraploid (Ohtani et al., 2010)), this result shows that *rpsT* can be reduced in
272 copy number, but not fully eliminated, suggesting that the gene is essential.

273 To assess the level of substitution achieved with the various plasmids, we designed a second PCR
274 screening protocol to specifically detect chromosomal *rpsT* via a 300 bp product. At low kanamycin
275 concentrations this protocol always generated a product (Figure 8d panel 1), but at increased
276 kanamycin concentration we did not obtain product for any *rpsT* allele (Figure 8d panel 2, lanes TT,
277 TA, TA2 and TA4). This demonstrates that plasmid-borne *rpsT* and its mutants were able to
278 complement the chromosomal *rpsT* and that the latter was displaced from the population to a level
279 that left it undetectable by PCR. In contrast, we could never completely suppress chromosomal *rpsT*
280 in strain KM4 complemented with an empty plasmid, even under high kanamycin conditions (120
281 $\mu\text{g/ml}$).

282 In *E. coli* and *Salmonella enterica*, *rpsT* has been reported to be non-essential, but its deletion
283 significantly lowers growth rates (Bubunencko et al., 2007; Tobin et al., 2010). We found that *rpsT* is
284 essential in *T. thermophilus*, but that its loss could be complemented by wild-type and mutant *T.*
285 *aquaticus rpsT*, and that this complementation restored wild-type levels of growth (Figure 8e).
286 Moreover, when the selection pressure from kanamycin was removed, no reversal in the PCR
287 products was detected for any strain (Figure 8c and d, panel 3), which confirms that chromosomal
288 *rpsT* was substantially displaced during kanamycin treatment. We conclude that *rpsT* from *T.*

289 *aquaticus* and its two mutated alleles are neutral with respect to survival and growth for *T.*
290 *thermophilus*. This demonstrates that the mutations we introduced do not affect negatively the
291 interaction between RPS20 and its cognate RNA, and that therefore such mutations could have been
292 sampled multiply and in a cumulative fashion by neutral drift in the course of evolution.

293 **2.7 Implications for the emergence of folded proteins**

294 Proteins are the most complex macromolecules synthesized in nature and by and large need to
295 assume defined structures for their activity. This folding process is complicated and easily disrupted,
296 witness the elaborate systems for protein folding, quality control and degradation universal to all
297 living beings. Despite widespread problems to reach and maintain the folded state, natural proteins
298 nevertheless form a best-case group, since the overwhelming majority of random polypeptides do
299 not appear to have a folded structure (Keefe and Szostak, 2001; Wei et al., 2003). It thus seems
300 impossible that, at the origin of life, the prototypes for the folded proteins we see today could have
301 arisen by random concatenation of amino acids. We have proposed that folding resulted from the
302 increasing complexity of peptides that supported RNA replication and catalysis, and that these
303 peptides assumed their structure through the interaction with the RNA scaffold (Lupas et al., 2001;
304 Söding and Lupas, 2003). In this view, protein folding was an emergent property of RNA-peptide
305 coevolution. We have recently described 40 such peptides whose conservation in diverse folds
306 suggests that they predated folded proteins (Alva et al., 2015). These peptides are enriched for
307 nucleic-acid binders, particularly in the context of the ribosome.

308 Due to its extremely slow rate of change, the ribosome essentially represents a living fossil, providing
309 the possibility to study the chronology of ancient events in molecular evolution (Hsiao et al., 2009).
310 Thus, core ribosomal proteins offer a window into the time when proteins were acquiring the ability
311 to fold. Those close to the catalytic center almost entirely lack secondary structure. Further away
312 from the center, their secondary structure content gradually increases and at the periphery, these
313 secondary structure elements become arranged into topologies that parallel those seen in cytosolic
314 proteins (Hsiao et al., 2009). Collectively, the structures of ribosomal proteins chart a path of

315 progressive emancipation from the RNA scaffold. Even the peripheral proteins, however, still mostly
316 assume their structure only in the context of the ribosomal RNA, as exemplified by RPS20 in our
317 study (supplementary file 1F, see also (Paterakis et al., 1983)).

318 The simplest mechanism to achieve an increase in complexity is the repetition of building blocks and
319 nature provides many examples for this, at all levels of organization. The dominant role of repetition
320 in the genesis of protein folds has been documented in many publications since the 1960s (Alva et
321 al., 2007; Blundell et al., 1979; Broom et al., 2012; Eck and Dayhoff, 1966; Kopec and Lupas, 2013;
322 Lee and Blaber, 2011; McLachlan, 1972; McLachlan, 1987; Remmert et al., 2010; Söding et al., 2006).
323 As a test of this mechanism, we explored whether a peptide originating from a ribosomal protein
324 that is disordered outside the context of the ribosome, could form a folded protein through an
325 increase in complexity afforded by repetition. For this, we chose a present-day representative of one
326 of the 40 fragments we reconstructed (Alva et al., 2015); this fragment is naturally found in single
327 copy in several different folds, including that of ribosomal protein RPS20, and repetitively in one fold,
328 TPR. Simple repetition was not sufficient in our case, but the repeat protein was so close to a folded
329 structure that only two point mutations per repeat were necessary to allow it to fold reliably. The
330 mutations needed for this transition did not appear to affect negatively the interaction with the RNA
331 scaffold, raising the possibility that they could have been among the variants sampled multiply in the
332 course of evolution.

333 Our experiments recapitulate a scenario for the emergence of a protein fold by a widespread and
334 well-documented mechanism, and show that this could have proceeded in a straightforward way.
335 These experiments represent a proof of concept, starting with a modern peptide likely to still retain
336 many features of an ancestral α -hairpin that gave rise to a number of folds, including TPR. Rather
337 than proposing proto-RPS20 as the parent of TPR domains, we see it as one of many proteins
338 emerging from this ancestral hairpin. Given the ease with which repetition of the RPS20 hairpin
339 yielded a TPR-like fold, we consider it likely that the hairpins belonging to the ancestral group were
340 amplified many times during the emergence of folded proteins to yield a range of TPR-like offspring,

341 of which only one may have survived to this day (but see also the figure legend to Figure 1). The
342 reason for this limited survival may lie in the fact that structure is a prerequisite for protein function,
343 but it is the function that is under biological selection. It could be that the newly emerged TPR-like
344 folds required many additional changes to achieve a useful activity and that therefore only very few –
345 possibly just one – survived. We consider a different scenario more probable, however. All present-
346 day TPR domains whose function has been characterized mediate protein-protein interactions by
347 binding to linear sequence motifs in unstructured polypeptide segments (D'Andrea and Regan, 2003;
348 Zeytuni and Zarivach, 2012). This activity would have been particularly relevant at a time of transition
349 from peptides dependent on RNA scaffolds for their structure, to autonomously folded polypeptides.
350 Functions relevant in this context would have been to prevent aggregation and increase the solubility
351 of newly emerging (poly)peptides, to promote autonomous folding, to serve as assembling factors
352 for RNA-protein and protein-protein complexes, and to recognize targeting sequences in the
353 emerging cellular networks. It therefore seems likely to us that many of the newly evolved TPR-like
354 folds became established in one or the other of these activities, only to be subsequently displaced by
355 folding becoming a general property of cellular polypeptides and by more advanced, energy-
356 dependent folding factors, which offered much better regulation. Exploring the extent to which our
357 new TPR protein could fulfill such functions represents the next frontier in our studies.

358 **3 Materials and Methods**

359 **3.1 Phylogeny for recently amplified TPR arrays**

360 All sequence similarity searches in this work were performed using the Web BLAST
361 (RRID:SCR_004870) from the National Institute for Biotechnology Information (NCBI;
362 <http://blast.ncbi.nlm.nih.gov>; (Boratyn et al., 2013)) and in the MPI Bioinformatics Toolkit
363 (RRID:SCR_010277, <https://toolkit.tuebingen.mpg.de/>; (Alva et al., 2016)). Examples of recently
364 amplified repeat units in TPR were taken from a previous investigation (Dunin-Horkawicz et al.,
365 2014). The TPR domain in serine/threonine-protein phosphatase 5 was chosen as a representative

366 three-repeat TPR, the most common TPR form in natural proteins (D'Andrea and Regan, 2003;
367 Sawyer et al., 2013), to study divergent evolution of TPR. We ran BLAST on the non-redundant
368 protein sequence database (nr) with an E-value threshold of 0.05 using the TPR domain of
369 serine/threonine-protein phosphatase 5 from *Homo sapiens* as query (Das et al., 1998). From the
370 results, we chose seven taxa to cover a diverse range of life.

371 TPRpred program (Karpenahalli et al., 2007) was used to help identify tandem repeats of TPR units.
372 The construction of multiple sequence alignments (MSAs) for TPR units was straightforward as all
373 TPR units are of the same size (34 aa) and no indels were allowed in the MSAs. We used Clustal X 2.1
374 (Larkin et al., 2007) to build phylogenetic trees using the neighbor-joining clustering algorithm and
375 1000 bootstrap trials (Bootstrap N-J Tree). SplitTree4 (Huson and Bryant, 2006) was used to render
376 the phylogenetic trees.

377 **3.2 Identification of helical hairpins resembling the TPR unit**

378 To find proteins homologous to the TPR unit, we first employed the TPRpred program (Karpenahalli
379 et al., 2007) to identify proteins that share significant sequence similarity to the TPR sequence
380 profile, then filtered them by comparing to the TPR structures.

381 First, TPRpred program with TPR profile tpr2.8 was used to identify TPR unit like sequences from all
382 protein sequences of known structures in the Protein Data Bank (PDB, RRID:SCR_012820) (Berman et
383 al., 2000). Protein sequences from the SEQRES record in PDB files were downloaded from the PDB.
384 We only considered sequences with at least 34 residues, which is the length of the TPR unit.
385 Redundancy was removed by keeping only non-identical sequences. In total, 68,197 sequences were
386 scanned by using TPRpred with default parameters. Only fragments predicted to be TPR with a *P*-
387 value lower than $1.0e-4$ were retained (646 hits). We estimated the false discovery rate (FDR)
388 (Noble, 2009) associated with this *P*-value cutoff using a simulated sequence dataset generated by
389 using the amino-acid composition derived from the PDB sequences. The dataset contains the same
390 number of sequences of the same length distribution as the PDB sequences. The FDR was estimated

391 to be the ratio of the number of hits in the simulated dataset to the number of detected hits in the
392 PDB sequences (Noble, 2009). We repeated the simulation 100 times and the FDR was estimated to
393 be $1.0 \pm 0.4\%$.

394 Within the 646 hits, we kept only TPR unit singletons, which are TPR units that have no other TPR
395 units close to them within a distance of 10 residues in the same sequence. TPR units of identical
396 sequences are considered only once. Subsequently, these TPR unit singletons were filtered by
397 removing those annotated to belong to clan *TPR* (CL0020) in Pfam 27.0 (RRID:SCR_004726).

398 The structures of the predicted TPR units obtained from the previous step were then compared to an
399 average TPR unit structure. A predicted TPR unit was discarded if the C_α RMSD of the 34 residues is
400 greater than 2.0 Å after superposition. The average TPR unit structure was generated by considering
401 all proteins belonging to family tetratricopeptide repeat (TPR) (a.118.8.1) in SCOP 1.75
402 (RRID:SCR_007039) (Murzin et al., 1995). TPR repeats in these proteins were again detected using
403 TPRpred and a per-repeat *P*-value cutoff of $1.0e-4$ was used. In total, 50 non-redundant TPR repeat
404 fragments were identified and superposed using a multiple structure alignment tool MultiProt
405 (Shatsky et al., 2004). The average C_α positions were calculated from the 50 structures after
406 superposition. We obtained 31 fragments after the structure filtering step (supplementary file 1C).
407 We then inspected the protein structures using PyMOL (RRID:SCR_000305) (Schrödinger, 2010).
408 Among them, 22 were observed to be involved in the formation of solenoid or tandem repeat
409 structures and were thus not further considered.

410 **3.3 Identification of TPR homologs in RPS20**

411 We applied TPRpred to scan all RPS20 sequences belonging to Pfam 27.0 family *Ribosomal S20p*
412 (PF01649), including sequences from both datasets “full” and “ncbi”. There are 4,402 and 2,284
413 sequences in the two sets. We merged the two sets and removed identical sequences to create a
414 dataset of 3,742 RPS20 sequences. TPRpred was used to detect TPR unit homologs in them. We

415 obtained 24 hits in these RPS20 sequences predicted by TPRpred to match TPR unit profile with a *P*-
416 value smaller than $1.0e-4$ (see supplementary file 1D).

417 We defined “interface positions” in the TPR unit and then transferred the definition to RPS20-hh
418 according to their structure superposition. We considered the residues on the outer side of the two
419 helices facing neighboring TPR units. Both helix A and helix B in the TPR unit are α -helices, which
420 have on average 3.6 residues per turn. Thus, every third or fourth residue always appears on the
421 same side of the helix. They are positions 3, 7 and 10 in helix A and positions 17, 21, 24 and 28 in
422 helix B. According to the TPR sequence profile compiled by Main et al. (Main et al., 2003b), the most
423 common residues at these positions are hydrophobic except for positions 17 and 24, where the most
424 common residues are both Tyr (see also Figure 4a). Therefore, positions 17 and 24 were not included
425 in the definition of interface positions. Furthermore, the residue at position equivalent to position 24
426 in RPS20 structure faces its C-terminal helix and is already an interface residue (Figure 4c). Thus, it
427 was not considered as an interface position to be checked in the study. In the end, only positions 3,
428 7, 10, 21 and 28 in RPS20-hh were defined to be interface positions to be examined, because they
429 are exposed to the solvent or interact with the RNA molecules in the ribosome, but would interact
430 with neighboring repeats in the TPR fold.

431 We searched all RPS20 sequences in Pfam 27.0 family *Ribosomal_S20p* (PF01649), including both
432 datasets “full” and “ncbi”, for candidates in which the interface positions are occupied by as many
433 hydrophobic residues as possible. In the MSA provided by Pfam, we extracted the 34 columns that
434 correspond to the sequence fragment of RPS20-hh from *Thermus aquaticus*, which was found by
435 TPRpred to be the hit with the best *P*-value and was thus used as the reference RPS20-hh. We
436 obtained 1,370 sequence fragments that do not contain any indels, in which the interface positions
437 were examined for hydrophobicity. Here, Ala, Ile, Leu, Met, Phe, Val were considered as hydrophobic
438 residues. Trp was not included as its side chain may be too large to be accommodated at the
439 interface.

440 We employed several low-complexity / intrinsically disordered region prediction methods (SEG
441 (Wootton, 1994), PONDR (Romero et al., 2001), DisEMBL (Linding et al., 2003), IUPred (Dosztányi et
442 al., 2005a; Dosztányi et al., 2005b)) to investigate putative intrinsically disordered regions in the
443 RPS20 of *Thermus aquaticus*. We ran SEG with three sets of recommended parameters (Wootton and
444 Federhen, 1996) and the other approaches with default parameters.

445 **3.4 Optimization of RPS20-hh in the designed TPRs**

446 We considered eight positions (2, 4, 6, 7, 9, 22, 23 and 25) in RPS20-hhta for optimization apart from
447 the four residues at the C-terminus.

448 Main et al. (Main et al., 2003b) discovered a set of eight “TPR signature residues” in the consensus
449 design: W4, L7, G8, Y11, A20, Y24, A27 and P32. Six of them are missing in RPS20-hhta except A20
450 and A27. Following the principle of consensus design, we introduced L4W and K7L into RPS20-hhta.
451 K7 is also one of the interface positions that ought to be mutated to hydrophobic residue for better
452 packing at interfaces. A8 and L11 were not optimized because they are the second and third most
453 common residues at positions 8 and 11 in the TPR profile, respectively. M24 was also retained
454 because it seems long hydrophobic side chains are favored at position 24 though Met is not one of
455 the three most common residues (YFL). P32 was introduced to replace D32 in RPS20-hhta as part of
456 the C-terminal consensus loop (DPNN) between repeats.

457 Co-evolution is commonly observed between physically interacting residues (de Juan et al., 2013).
458 We investigated if any positions we optimized are involved in co-evolution relationship so that we
459 can preserve such correlations. We performed direct coupling analysis (Morcos et al., 2011) and
460 computed the mutual information using MatrixPlot (Gorodkin et al., 1999) between all positions in
461 TPR repeat sequences. The results of both approaches revealed that the highest correlation occurs
462 between positions 7 and 23 (Figure 4—figure supplement 1), with the most commonly observed
463 combinations being R7-D23 and L7-Y23. Therefore, we always mutated I23 to the most commonly
464 observed residue tyrosine (I23Y) in the TPR consensus sequence together with aforementioned

465 mutation K7L. In addition, we considered combination K7R and I23D together. Combination K7-I23D
466 was also tested because of highly similar physicochemical properties between Lys and Arg side
467 chains.

468 The hydrophobic side chain of valine at position 9 in RPS20-hhta is buried between helices in RPS20,
469 but would be exposed on the surface of the designed protein except in the last repeat, in which V9
470 interacts with the stop helix. Therefore, it is considered to be mutated to the most common residue
471 asparagine (V9N) in the TPR repeat consensus except in the last repeat (Figure 4c).

472 RPS20-hhta sequence and surface is enriched with positively charged residues (Figure 4b). This would
473 lead to the exceedingly high theoretical iso-electric point (pI) of the designed proteins. Natural TPR
474 proteins tend to exhibit zero net charge (Magliery and Regan, 2004). Hence, we decided to randomly
475 mutate the positively charged residues (Lys and Arg) in the two helices of RPS20-hhta to the
476 corresponding most common residues in TPR sequence profile (K2E, K6N, K22E, R25Q/E). K26 was
477 not mutated as Lys is already the most common residue in the TPR profile.

478 At the C-terminus of the designed TPR, the last four residue of RPS20-hhta (IDKA) were replaced with
479 the TPR consensus loop sequence (DPNN) between repeat units. The reason is as follows. The
480 secondary structure of the TPR unit is helix (13 aa) – loop (3aa) – helix (14 aa) – loop (4aa), while the
481 secondary structure of the RPS20-hhta identified to be homologous to TPR unit is helix (13 aa) – loop
482 (3 aa) – helix (18 aa) (Figure 2 and 4). The last four residues may have been included in the prediction
483 by TPRpred merely to fulfill the size requirement of TPR repeat (34 aa). Indeed, when we scanned
484 RPS20-hhta sequence using the hidden Markov model constructed for Pfam family *TPR_1*, only
485 positions 2-28 were found to be similar to the *TPR_1* profile using HMMER 3.0 (RRID:SCR_005305)
486 (Eddy, 2009), even if all filters were switched off. So the four very C-terminal residues in RPS20-hhta
487 were not used in the designed TPR between repeat units. They were not replaced in the last repeat
488 unit (Figure 3).

489 **3.5 Structure modeling and refinement *in silico***

490 CTPR3 structure of an idealized TPR repeat (Main et al., 2003b) (PDB id: 1na0, chain A) was taken as
491 the main template to build an initial TPR structure model using RPS20-hhta. Helix B3 and the stop
492 helix in our designed protein are different from natural TPRs, but more similar to natural RPS20s. So
493 we also used a RPS20 protein as the structure template for the last repeat and the stop helix. The
494 structure of RPS20 from *Thermus thermophilus* HB8 (PDB id: 2vqe, chain T) was used because it was
495 the structure with the best resolution (2.5 Å). The C-terminal loop in 2vqeT was discarded. The two
496 structures 1na0A and 2vqeT were merged into a hybrid template based on the superposition of their
497 homologous helical hairpins: the third TPR unit in 1na0A and the RPS20-hh in 2vqeT (the very C-
498 terminal four residues were not used). We then modeled the designed TPR sequences using RPS20-
499 hhta onto the hybrid structure template using Rosetta programs *fixbb* and *relax* (Das and Baker,
500 2008). The Rosetta fixed backbone design application *fixbb* was used to make the initial model.
501 Subsequently, these models were relaxed using the Rosetta structure refinement application *relax*.
502 The two steps were iterated three times. See the supplementary file 1E for the command lines.
503 Rosetta 3.4 was used in the work.

504 We selected five constructs for further testing *in vitro* (Table 1). They are among the best-scoring
505 constructs according to the *in silico* simulation (Figure 4—figure supplement 2). If two constructs
506 have comparable scores (they are adjacent in the score ranking), the one with fewer mutations was
507 preferred. The selected constructs all differ at least at two positions in their sequences. When we
508 searched these optimized RPS20-hhta fragments in the NCBI *nr* database using BLAST (Camacho et
509 al., 2009), the top hits were still RPS20s.

510 **3.6 Cloning, protein expression and purification**

511 DNA sequences coding for the designed TPR repeats were gene-synthesized in codon-optimized form
512 (Eurofins) and cloned into vector pET-28b (Novagen) using NcoI/HindIII restriction sites, and into
513 pETHis_1a to generate proteins with an N-terminal cleavable His₆-tag. RPS20 *T. aquaticus* and *T.*

514 *thermophilus* genes were amplified from genomic DNA and cloned likewise. Recombinant plasmids were
515 transformed into *E. coli* strain BL21-Gold (DE3) grown on LB agar plates containing 50 µg/ml kanamycin.
516 For expression, cells were cultured at 25°C and induced with 1 mM isopropyl-D-thiogalactopyranoside
517 (IPTG) at an OD₆₀₀ of 0.6 for continued growth overnight.

518 Bacterial cell pellets were resuspended in buffer A (50 mM Tris pH 8, 150 mM NaCl), supplemented with
519 5 mM MgCl₂, DNaseI (Applichem) and protease inhibitor cocktail (cOmplete, Roche). After breaking the
520 cells in a French Press, the suspension was centrifuged twice at 37,000 g. Soluble His₆-tagged proteins
521 were purified by binding proteins to Ni-NTA columns (GE Healthcare) in buffer A (50 mM Tris pH 8.0, 300
522 mM NaCl) and elution with increasing concentrations of imidazole up to 0.6 M. Eluted proteins were
523 dialyzed against buffer A for cleavage by His₆-TEV-protease (50 U/mg protein). Cleavage leaves two
524 additional residues (Gly-Ala) as N-terminal extension to all proteins produced in this manner. After
525 incubation overnight, cleaved proteins were re-run on Ni-NTA columns and collected in the flow-
526 through. They were finally purified by gel size exclusion chromatography (Superdex G75, GE Healthcare)
527 in buffer A containing 0.5 mM EDTA. Insoluble proteins were dissolved in 6 M guanidinium chloride and
528 refolded by dialysis overnight against buffer A. Refolded proteins were further purified by sequential
529 anion-exchange (Q Sepharose HP) and cation-exchange (SP Sepharose HP) chromatography using 0-500
530 mM NaCl salt gradients in buffer D (20 mM Tris pH 8, 1 mM EDTA), and by gel size exclusion
531 chromatography (Superdex G75) in buffer A.

532 **3.7 Biophysical characterization**

533 To determine the native molecular mass of designed TPR repeats, static light scattering experiment
534 were performed by applying samples onto a superdex S200 gel size exclusion column to which a
535 miniDAWN Tristar Laser photometer (Wyatt) and an RI 2031 differential refractometer (JASCO) were
536 coupled. Runs were performed in buffer A. Data analysis and molecular mass calculations were
537 carried out with ASTRA V software (Wyatt). Tryptophan fluorescence spectra were recorded on a
538 Jasco FP-6500 spectrofluorometer at 23°C; excitation was at 280 nm, emission spectra were collected
539 from 300-400 nm. Circular dichroism (CD) spectra from 200-250 nm were recorded with a Jasco J-810

540 spectropolarimeter at 23°C in buffer E (30 mM MOPS pH 7.2, 150 mM NaCl). Cuvettes of 1 mm path
541 length were used in all measurements. For melting curves and determination of T_m, CD
542 measurements were recorded at 222 nm from 20-95°C, the temperature change was set to 1°C per
543 minute, using a Peltier-controlled sample holder unit. For equilibrium-unfolding experiments
544 performed at 23°C, native protein was mixed with different concentrations of urea in buffer A. After
545 equilibration, circular dichroism was monitored at 222 nm. The fraction of unfolded protein f_U was
546 determined based on $f_U = (y_F - y)/(y_F - y_U)$, where y_F and y_U are the values of y typical of the folded
547 and unfolded states. Data were fitted to a two-state model with the software ProFit (6.1) as
548 described (Grimsley et al., 2013), assuming a linear urea $[D]$ dependence of ΔG : $\Delta G_{U-F}^D = \Delta G_{U-F}^{H_2O} -$
549 $m[D]$, where ΔG_{U-F}^D is the free energy change at a given denaturant concentration, $\Delta G_{U-F}^{H_2O}$ the free
550 energy change in the absence of denaturant, and m the sensitivity of the transition to denaturant.
551 Fragment sizes of M4N were determined by ESI-microTOF mass spectrometry (Bruker Daltonics, Max
552 Planck Institute core facility Martinsried), followed by bioinformatic analysis using the Find-Pept tool
553 (ExpASY).

554 **3.8 Crystallization, structure solution and refinement**

555 For crystallization, the M4N and M4N Δ C protein solutions were concentrated to 70 and 30 mg/ml,
556 respectively, in buffer A. The buffer for M4N Δ C additionally contained 0.5 mM EDTA. Crystallization
557 trials were performed at 295 K in 96-well sitting-drop vapor-diffusion plates with 50 μ l of reservoir
558 solution and drops consisting of 300 nl protein solution and 300 nl reservoir solution in the case of
559 M4N, and 400 nl protein solution and 200 nl reservoir solution in the case of M4N Δ C. Crystallization
560 conditions for the crystals used in the diffraction experiments are listed in supplementary file 1H
561 together with the solutions used for cryo-protection. Single crystals were transferred into a droplet
562 of cryo-protectant before loop-mounting and flash-cooling in liquid nitrogen. For experimental
563 phasing, crystals of M4N were soaked overnight in a droplet containing reservoir solution
564 supplemented with 5 mM K₂PtCl₄ prior to cryo-protection and flash-cooling. All data were collected
565 at beamline X10SA (PXII) at the Swiss Light Source (Paul Scherrer Institute, Villigen, Switzerland) at

566 100 K using a PILATUS 6M detector (DECTRIS) at the wavelengths indicated in supplementary file 1H.
567 Diffraction images were processed and scaled using the XDS program suite (Kabsch, 1993). Using
568 SHELXD (Sheldrick, 2008), three strong Pt-sites were identified in the M4N derivative dataset. After
569 density modification with SHELXE, the resulting electron density map could be traced by Buccaneer
570 (Cowtan, 2006) to large extents, and revealed three chains of M4N in the asymmetric unit (ASU),
571 organized as one dimer and one monomer. Refinement was continued using the native dataset. The
572 two different crystal forms of M4N Δ C, CF I and CF II, were solved by molecular replacement on the
573 basis of the refined M4N coordinates. Using MOLREP (Vagin and Teplyakov, 2000), two copies of the
574 dimeric assembly of the M4N structure were located in the ASU of CF I, and one monomer in the ASU
575 of CF II. All models were completed by cyclic manual modeling with Coot (Emsley and Cowtan, 2004)
576 and refinement with PHENIX (RRID:SCR_014224) (Adams et al., 2010). Analysis with PROCHECK
577 (Laskowski et al., 1993) showed excellent geometries for all structures. Data collection and
578 refinement statistics are summarized in supplementary file 1H. The three structures are deposited in
579 the PDB (Berman et al., 2000) with accession codes: 5FZQ (M4N), 5FZR (M4N Δ C CF I), 5FZS (M4N Δ C
580 CF II).

581 **3.9 Testing mutations in *T. thermophilus***

582 *T. thermophilus* HB8 and *T. aquaticus* YT-1 were obtained from the German Collection of
583 Microorganisms and Cell Cultures (DSMZ). Growth in liquid medium was performed under mild
584 stirring (160 rpm) in long necked flasks at 68°C with DSMZ Medium 74 for *T. thermophilus* and DSMZ
585 Medium 878 for *T. aquaticus*. Agar (1.6% w/v) was added to the medium for growth on plates. When
586 required, kanamycin (30 μ g/ml) and bleocin (10 μ g/ml) were added to the media. For purification
587 experiments 25 ml cultures were grown to an optic density of 0.7 OD₆₀₀ (~12 hours) and then re-
588 inoculated in the same volume to an optical density of 0.035 OD₆₀₀. The process was repeated serially
589 three times and two 5 ml samples were taken in each step for glycerol stocks and DNA purification.
590 Transformation of *T. thermophilus* was performed as described previously (Nguyen and Silberg,

591 2010). Genomic and plasmid DNA from *Thermus* were purified from 5 ml cultures using the QIAamp
592 DNA Mini Kit and the QIAprep Spin Miniprep Kit, respectively.

593 *T. thermophilus* KM4 strain was generated by gene replacement as follows: two PCR products
594 comprising each one 1 Kb of DNA upstream and downstream of *rpsT* were amplified from *T.*
595 *thermophilus* HB8 genomic DNA and then fused by overlapping PCR. The resulting fragment, in which
596 *rpsT* is substituted by a PstI site, was cloned in the KpnI/XbaI sites of plasmid pBlueScript II SK (+).
597 Next, a fragment from plasmid pKT1 (Biotools, Spain), which contains the thermostable kanamycin
598 resistance *Kat* gene under the control of the constitutive PslpA promoter, was inserted into the new
599 PstI site. Direction of the *Kat* cassette insertion was selected, so transcription from the PslpA
600 promoter continues through *thx*, a gene that is located downstream and is predicted to form an
601 operon with *rpsT*. The 3 Kb final construct cloned in pBluescript was subsequently amplified by PCR
602 and the linear product was purified and transformed by electroporation in *T. thermophilus* HB8.
603 Integration of the *Kat* cassette was selected by growth in kanamycin.

604 For the complementation in trans of *rpsT* from *T. thermophilus*, a PCR product of *rpsT* was amplified
605 from genomic DNA and cloned in the SpeI/PstI sites of plasmid pJJSpro (Nguyen and Silberg, 2010)
606 generating plasmid pJJSpro-rps20Tt. The same approach was followed for *rpsT* in *T. aquaticus*
607 (pJJSpro-rpsTTa) and in *T. aquaticus* *rpsT* alleles with two (pJJSpro-rpsTTaM2) and four (pJJSpro-
608 rpsTTaM4N) amino-acid substitutions. The PCR product for the two later constructs was amplified
609 using the plasmids in which the synthesized genes were delivered as a template.

610 **4 Acknowledgments**

611 We thank Elisabeth Weyher from the Core Facility of the MPI for Biochemistry, Martinsried, for
612 analyzing proteins by mass spectrometry. We are grateful to the staff of beamline PXII/Swiss Light
613 Source for their technical support. We also thank Birte Höcker, Vikram Alva and Sergey Samsonov for
614 many helpful comments and discussions. This work was supported by institutional funds of the Max
615 Planck Society.

616 **5 References**

- 617 Abe Y, Shodai T, Muto T, Mihara K, Torii H, Nishikawa S, Endo T, Kohda D. 2000. Structural basis of
618 presequence recognition by the mitochondrial protein import receptor Tom20. *Cell* **100**(5):
619 551-560.
- 620 Adams PD, Afonine PV, Bunkóczi G, Chen VB, Davis IW, Echols N, Headd JJ, Hung L-W, Kapral GJ,
621 Grosse-Kunstleve RW, others. 2010. PHENIX: a comprehensive Python-based system for
622 macromolecular structure solution. *Acta Crystallogr Sect D Biol Crystallogr* **66**(2): 213-221.
- 623 Alva V, Ammelburg M, Söding J, Lupas AN. 2007. On the origin of the histone fold. *BMC Struct Biol*
624 **7**(17): 1-10.
- 625 Alva V, Söding J, Lupas AN. 2015. A vocabulary of ancient peptides at the origin of folded proteins.
626 *Elife* **4**(pii:e09410). doi: 10.7554/eLife.09410.
- 627 Alva V, Nam S-Z, Söding J, Lupas AN. 2016. The MPI bioinformatics Toolkit as an integrative platform
628 for advanced protein sequence and structure analysis. *Nucleic Acids Res* **44**(W1): W410-
629 W415.
- 630 Anantharaman V, Koonin EV, Aravind L. 2001. TRAM, a predicted RNA-binding domain, common to
631 tRNA uracil methylation and adenine thiolation enzymes. *FEMS Microbiol Lett* **197**(2): 215-
632 221.
- 633 Apic G, Gough J, Teichmann SA. 2001. Domain combinations in archaeal, eubacterial and eukaryotic
634 proteomes. *J Mol Biol* **310**(2): 311-325.
- 635 Aurora R, Rose GD. 1998. Helix capping. *Protein Sci* **7**(1): 21-38.
- 636 Bansal PK, Mishra A, High AA, Abdulle R, Kitagawa K. 2009a. Sgt1 dimerization is negatively regulated
637 by protein kinase CK2-mediated phosphorylation at Ser361. *J Biol Chem* **284**(28): 18692-
638 18698.
- 639 Bansal PK, Nourse A, Abdulle R, Kitagawa K. 2009b. Sgt1 dimerization is required for yeast
640 kinetochore assembly. *J Biol Chem* **284**(6): 3586-3592.

641 Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov IN, Bourne PE. 2000.
642 The Protein Data Bank. *Nucleic Acids Res* **28**(1): 235-242.

643 Bernhardt HS. 2012. The RNA world hypothesis: the worst theory of the early evolution of life (except
644 for all the others). *Biol Direct* **7**(23): 1-10.

645 Biegert A, Mayer C, Remmert M, Söding J, Lupas AN. 2006. The MPI Bioinformatics Toolkit for protein
646 sequence analysis. *Nucleic Acids Res* **34**(Web Server issue): W335-W339.

647 Binz HK, Stumpp MT, Forrer P, Amstutz P, Plückthun A. 2003. Designing repeat proteins: well-
648 expressed, soluble and stable proteins from combinatorial libraries of consensus ankyrin
649 repeat proteins. *J Mol Biol* **332**(2): 489-503.

650 Blundell TL, Sewell BT, McLachlan AD. 1979. Four-fold structural repeat in the acid proteases.
651 *Biochim Biophys Acta* **580**(1): 24-31.

652 Boratyn GM, Camacho C, Cooper PS, Coulouris G, Fong A, Ma N, Madden TL, Matten WT, McGinnis
653 SD, Merezuk Y, Raytselis Y, Sayers EW, Tao T, Ye J, Zaretskaya I. 2013. BLAST: a more
654 efficient report with usability improvements. *Nucleic Acids Res* **41**(Web Server issue): W29-
655 W33.

656 Broom A, Doxey AC, Lobsanov YD, Berthin LG, Rose DR, Howell PL, McConkey BJ, Meiering EM. 2012.
657 Modular evolution and the origins of symmetry: reconstruction of a three-fold symmetric
658 globular protein. *Structure* **20**(1): 161-171.

659 Bubunenko M, Baker T, Court DL. 2007. Essentiality of ribosomal and transcription antitermination
660 proteins analyzed by systematic gene replacement in *Escherichia coli*. *J Bacteriol* **189**(7):
661 2844-2853.

662 Burton B, Zimmermann MT, Jernigan RL, Wang Y. 2012. A computational investigation on the
663 connection between dynamics properties of ribosomal proteins and ribosome assembly.
664 *PLoS Comput Biol* **8**(5): e1002530.

665 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009. BLAST+:
666 architecture and applications. *BMC Bioinformatics* **10**(421): 1-9.

667 Chaudhuri I, Söding J, Lupas AN. 2008. Evolution of the beta-propeller fold. *Proteins* **71**(2): 795-803.

668 Chen C, Natale DA, Finn RD, Huang H, Zhang J, Wu CH, Mazumder R. 2011. Representative
669 proteomes: a stable, scalable and unbiased proteome set for sequence analysis and
670 functional annotation. *PLoS One* **6**(4): e18910.

671 Coquille S, Filipovska A, Chia T, Rajappa L, Lingford JP, Razif MFM, Thore S, Rackham O. 2014. An
672 artificial PPR scaffold for programmable RNA recognition. *Nat Commun* **5**: 5729.

673 Cortajarena AL, Regan L. 2006. Ligand binding by TPR domains. *Protein Sci* **15**(5): 1193-1198.

674 Cowtan K. 2006. The Buccaneer software for automated model building. 1. Tracing protein chains.
675 *Acta Crystallogr Sect D Biol Crystallogr* **62**(9): 1002-1011.

676 Crick FHC. 1953. The packing of α -helices: simple coiled-coils. *Acta Crystallogr* **6**(8-9): 689-697.

677 Crooks GE, Hon G, Chandonia J-M, Brenner SE. 2004. WebLogo: a sequence logo generator. *Genome*
678 *Res* **14**(6): 1188-1190.

679 D'Andrea LD, Regan L. 2003. TPR proteins: the versatile helix. *Trends Biochem Sci* **28**(12): 655-662.

680 Das AK, Cohen PW, Barford D. 1998. The structure of the tetratricopeptide repeats of protein
681 phosphatase 5: implications for TPR-mediated protein-protein interactions. *EMBO J* **17**(5):
682 1192-1199.

683 Das R, Baker D. 2008. Macromolecular modeling with rosetta. *Annu Rev Biochem* **77**: 363-382.

684 de Juan D, Pazos F, Valencia A. 2013. Emerging methods in protein co-evolution. *Nat Rev Genet* **14**(4):
685 249-261.

686 Di Domenico T, Potenza E, Walsh I, Parra RG, Giollo M, Minervini G, Piovesan D, Ihsan A, Ferrari C,
687 Kajava AV, Tosatto SCE. 2014. RepeatsDB: a database of tandem repeat protein structures.
688 *Nucleic Acids Res* **42**(Database issue): D352-D357.

689 Dosztányi Z, Csizmók V, Tompa P, Simon I. 2005a. The pairwise energy content estimated from amino
690 acid composition discriminates between folded and intrinsically unstructured proteins. *J Mol*
691 *Biol* **347**(4): 827-839.

692 Dosztányi Z, Csizmók V, Tompa P, Simon I. 2005b. IUPred: web server for the prediction of
693 intrinsically unstructured regions of proteins based on estimated energy content.
694 *Bioinformatics* **21**(16): 3433-3434.

695 Doyle L, Hallinan J, Bolduc J, Parmeggiani F, Baker D, Stoddard BL, Bradley P. 2015. Rational design of
696 α -helical tandem repeat proteins with closed architectures. *Nature* **528**(7583): 585-588.

697 Dunin-Horkawicz S, Kopec KO, Lupas AN. 2014. Prokaryotic ancestry of eukaryotic protein networks
698 mediating innate immunity and apoptosis. *J Mol Biol* **426**(7): 1568-1582.

699 Dyson HJ, Wright PE. 2005. Intrinsically unstructured proteins and their functions. *Nat Rev Mol Cell*
700 *Biol* **6**(3): 197-208.

701 Eck RV, Dayhoff MO. 1966. Evolution of the structure of ferredoxin based on living relics of primitive
702 amino acid sequences. *Science* **152**(3720): 363-366.

703 Eddy SR. 2009. A new generation of homology search tools based on probabilistic inference. *Genome*
704 *Inform* **23**(1): 205-211.

705 Emsley P, Cowtan K. 2004. Coot: model-building tools for molecular graphics. *Acta Crystallogr D Biol*
706 *Crystallogr* **60**(Pt 12 Pt 1): 2126-2132.

707 Finn RD, Bateman A, Clements J, Coghill P, Eberhardt RY, Eddy SR, Heger A, Hetherington K, Holm L,
708 Mistry J, Sonnhammer ELL, Tate J, Punta M. 2014. Pfam: the protein families database.
709 *Nucleic Acids Res* **42**(Database issue): D222-D230.

710 Forrer P, Binz HK, Stumpp MT, Plückthun A. 2004. Consensus design of repeat proteins.
711 *Chembiochem* **5**(2): 183-189.

712 Fox GE. 2010. Origin and evolution of the ribosome. *Cold Spring Harb Perspect Biol* **2**(9): a003483.

713 Fox NK, Brenner SE, Chandonia J-M. 2014. SCOPe: Structural Classification of Proteins--extended,
714 integrating SCOP and ASTRAL data and classification of new structures. *Nucleic Acids Res*
715 **42**(Database issue): D304-D309.

716 Gilbert W. 1986. Origin of life: The RNA world. *Nature* **319**: 618.

717 Gorodkin J, Staerfeldt HH, Lund O, Brunak S. 1999. MatrixPlot: visualizing sequence constraints.
718 *Bioinformatics* **15**(9): 769-770.

719 Grimsley GR, Trevino SR, Thurlkill RL, Scholtz JM. 2013. Determining the conformational stability of a
720 protein from urea and thermal unfolding curves. *Curr Protoc Protein Sci* **Chapter**
721 **28**(Unit28.4): 28.24.21-28.24.14.

722 Habchi J, Tompa P, Longhi S, Uversky VN. 2014. Introducing protein intrinsic disorder. *Chem Rev*
723 **114**(13): 6561-6588.

724 Hsiao C, Mohan S, Kalahar BK, Williams LD. 2009. Peeling the onion: ribosomes are ancient molecular
725 fossils. *Mol Biol Evol* **26**(11): 2415-2425.

726 Huson DH, Bryant D. 2006. Application of phylogenetic networks in evolutionary studies. *Mol Biol*
727 *Evol* **23**(2): 254-267.

728 Iwaya N, Kuwahara Y, Fujiwara Y, Goda N, Tenno T, Akiyama K, Mase S, Tochio H, Ikegami T,
729 Shirakawa M, Hiroaki H. 2010. A common substrate recognition mode conserved between
730 katanin p60 and VPS4 governs microtubule severing and membrane skeleton reorganization.
731 *J Biol Chem* **285**(22): 16822-16829.

732 Kabsch W, Sander C. 1983. Dictionary of protein secondary structure: pattern recognition of
733 hydrogen-bonded and geometrical features. *Biopolymers* **22**(12): 2577-2637.

734 Kabsch W. 1993. Automatic processing of rotation diffraction data from crystals of initially unknown
735 symmetry and cell constants. *J Appl Crystallogr* **26**(6): 795-800.

736 Kajander T, Cortajarena AL, Mochrie S, Regan L. 2007. Structure and stability of designed TPR protein
737 superhelices: unusual crystal packing and implications for natural TPR proteins. *Acta*
738 *Crystallogr D Biol Crystallogr* **63**(Pt 7): 800-811.

739 Kajava AV. 2012. Tandem repeats in proteins: from sequence to structure. *J Struct Biol* **179**(3): 279-
740 288.

741 Karpenahalli MR, Lupas AN, Söding J. 2007. TPRpred: a tool for prediction of TPR-, PPR- and SEL1-like
742 repeats from protein sequences. *BMC Bioinformatics* **8**(2): 1-8.

743 Katibah GE, Qin Y, Sidote DJ, Yao J, Lambowitz AM, Collins K. 2014. Broad and adaptable RNA
744 structure recognition by the human interferon-induced tetratricopeptide repeat protein
745 IFIT5. *Proc Natl Acad Sci U S A* **111**(33): 12025-12030.

746 Keefe AD, Szostak JW. 2001. Functional proteins from a random-sequence library. *Nature* **410**(6829):
747 715-718.

748 Keiski C-L, Harwich M, Jain S, Neculai AM, Yip P, Robinson H, Whitney JC, Riley L, Burrows LL, Ohman
749 DE, Howell PL. 2010. AlgK is a TPR-containing protein and the periplasmic component of a
750 novel exopolysaccharide secretin. *Structure* **18**(2): 265-273.

751 Kobe B, Kajava AV. 2000. When protein folding is simplified to protein coiling: the continuum of
752 solenoid protein structures. *Trends Biochem Sci* **25**(10): 509-515.

753 Kohl A, Binz HK, Forrer P, Stumpp MT, Plückthun A, Grütter MG. 2003. Designed to be stable: crystal
754 structure of a consensus ankyrin repeat protein. *Proc Natl Acad Sci U S A* **100**(4): 1700-1705.

755 Koonin EV. 2003. Comparative genomics, minimal gene-sets and the last universal common ancestor.
756 *Nat Rev Microbiol* **1**(2): 127-136.

757 Kopec KO, Lupas AN. 2013. β -Propeller blades as ancestral peptides in protein evolution. *PLoS One*
758 **8**(10): e77074.

759 Krachler AM, Sharma A, Kleanthous C. 2010. Self-association of TPR domains: Lessons learned from a
760 designed, consensus-based TPR oligomer. *Proteins* **78**(9): 2131-2143.

761 Kumar S, Bansal M. 1998. Dissecting alpha-helices: position-specific analysis of alpha-helices in
762 globular proteins. *Proteins* **31**(4): 460-476.

763 Kyrpides N, Overbeek R, Ouzounis C. 1999. Universal protein families and the functional content of
764 the last universal common ancestor. *J Mol Evol* **49**(4): 413-423.

765 Kyrpides NC, Woese CR. 1998. Tetratricopeptide-repeat proteins in the archaeon *Methanococcus*
766 *jannaschii*. *Trends Biochem Sci* **23**(7): 245-247.

767 Lamb JR, Tugendreich S, Hieter P. 1995. Tetratricopeptide repeat interactions: to TPR or not to TPR?
768 *Trends Biochem Sci* **20**(7): 257-259.

769 Larkin MA, Blackshields G, Brown NP, Chenna R, McGettigan PA, McWilliam H, Valentin F, Wallace
770 IM, Wilm A, Lopez R, Thompson JD, Gibson TJ, Higgins DG. 2007. Clustal W and Clustal X
771 version 2.0. *Bioinformatics* **23**(21): 2947-2948.

772 Laskowski RA, MacArthur MW, Moss DS, Thornton JM. 1993. PROCHECK: a program to check the
773 stereochemical quality of protein structures. *J Appl Crystallogr* **26**(2): 283-291.

774 Lee J, Blaber M. 2011. Experimental support for the evolution of symmetric protein architecture from
775 a simple peptide motif. *Proc Natl Acad Sci U S A* **108**(1): 126-130.

776 Linding R, Jensen LJ, Diella F, Bork P, Gibson TJ, Russell RB. 2003. Protein disorder prediction:
777 implications for structural proteomics. *Structure* **11**(11): 1453-1459.

778 Lunelli M, Lokareddy RK, Zychlinsky A, Kolbe M. 2009. IpaB-IpgC interaction defines binding motif for
779 type III secretion translocator. *Proc Natl Acad Sci U S A* **106**(24): 9661-9666.

780 Lupas AN, Ponting CP, Russell RB. 2001. On the evolution of protein folds: are similar motifs in
781 different protein folds the result of convergence, insertion, or relics of an ancient peptide
782 world? *J Struct Biol* **134**(2-3): 191-203.

783 Magliery TJ, Regan L. 2004. Beyond consensus: statistical free energies reveal hidden interactions in
784 the design of a TPR motif. *J Mol Biol* **343**(3): 731-745.

785 Main ERG, Jackson SE, Regan L. 2003a. The folding and design of repeat proteins: reaching a
786 consensus. *Curr Opin Struct Biol* **13**(4): 482-489.

787 Main ERG, Xiong Y, Cocco MJ, D'Andrea L, Regan L. 2003b. Design of stable alpha-helical arrays from
788 an idealized TPR motif. *Structure* **11**(5): 497-508.

789 Main ERG, Lowe AR, Mochrie SGJ, Jackson SE, Regan L. 2005. A recurring theme in protein
790 engineering: the design, stability and folding of repeat proteins. *Curr Opin Struct Biol* **15**(4):
791 464-471.

792 McLachlan AD. 1972. Repeating sequences and gene duplication in proteins. *J Mol Biol* **64**(2): 417-
793 437.

794 McLachlan AD. 1987. Gene duplication and the origin of repetitive protein structures. *Cold Spring*
795 *Harb Symp Quant Biol* **52**: 411-420.

796 Morcos F, Pagnani A, Lunt B, Bertolino A, Marks DS, Sander C, Zecchina R, Onuchic JN, Hwa T, Weigt
797 M. 2011. Direct-coupling analysis of residue coevolution captures native contacts across
798 many protein families. *Proc Natl Acad Sci U S A* **108**(49): E1293-E1301.

799 Mosavi LK, Minor J, Daniel L, Peng Z-Y. 2002. Consensus-derived structural determinants of the
800 ankyrin repeat motif. *Proc Natl Acad Sci U S A* **99**(25): 16029-16034.

801 Murzin AG, Brenner SE, Hubbard T, Chothia C. 1995. SCOP: a structural classification of proteins
802 database for the investigation of sequences and structures. *J Mol Biol* **247**(4): 536-540.

803 Nguyen PQ, Silberg JJ. 2010. A selection that reports on protein-protein interactions within a
804 thermophilic bacterium. *Protein Eng Des Sel* **23**(7): 529-536.

805 Noble WS. 2009. How does multiple testing correction work? *Nat Biotechnol* **27**(12): 1135-1137.

806 Ohtani N, Tomita M, Itaya M. 2010. An extreme thermophile, *Thermus thermophilus*, is a polyploid
807 bacterium. *J Bacteriol* **192**(20): 5499-5505.

808 Oldfield CJ, Dunker AK. 2014. Intrinsically disordered proteins and intrinsically disordered protein
809 regions. *Annu Rev Biochem* **83**: 553-584.

810 Orengo CA, Thornton JM. 2005. Protein families and their evolution-a structural perspective. *Annu*
811 *Rev Biochem* **74**: 867-900.

812 Park K, Shen BW, Parmeggiani F, Huang P-S, Stoddard BL, Baker D. 2015. Control of repeat-protein
813 curvature by computational protein design. *Nat Struct Mol Biol* **22**(2): 167-174.

814 Parmeggiani F, Huang P-S, Vorobiev S, Xiao R, Park K, Caprari S, Su M, Seetharaman J, Mao L, Janjua
815 H, Montelione GT, Hunt J, Baker D. 2015. A general computational approach for repeat
816 protein design. *J Mol Biol* **427**(2): 563-575.

817 Paterakis K, Littlechild J, Woolley P. 1983. Structural and functional studies on protein S20 from the
818 30-S subunit of the *Escherichia coli* ribosome. *Eur J Biochem* **129**(3): 543-548.

819 Peng Z, Oldfield CJ, Xue B, Mizianty MJ, Dunker AK, Kurgan L, Uversky VN. 2014. A creature with a
820 hundred waggly tails: intrinsically disordered proteins in the ribosome. *Cell Mol Life Sci* **71**(8):
821 1477-1504.

822 Ponting CP, Russell RR. 2002. The natural history of protein domains. *Annu Rev Biophys Biomol Struct*
823 **31**: 45-71.

824 Ramarao MK, Bianchetta MJ, Lanke J, Cohen JB. 2001. Role of rapsyn tetratricopeptide repeat and
825 coiled-coil domains in self-association and nicotinic acetylcholine receptor clustering. *J Biol*
826 *Chem* **276**(10): 7475-7483.

827 Rämisch S, Weininger U, Martinsson J, Akke M, André I. 2014. Computational design of a leucine-rich
828 repeat protein with a predefined geometry. *Proc Natl Acad Sci U S A* **111**(50): 17875-17880.

829 Ranea JAG, Sillero A, Thornton JM, Orengo CA. 2006. Protein superfamily evolution and the last
830 universal common ancestor (LUCA). *J Mol Evol* **63**(4): 513-525.

831 Remmert M, Biegert A, Linke D, Lupas AN, Söding J. 2010. Evolution of outer membrane beta-barrels
832 from an ancestral beta beta hairpin. *Mol Biol Evol* **27**(6): 1348-1358.

833 Romero P, Obradovic Z, Li X, Garner EC, Brown CJ, Dunker AK. 2001. Sequence complexity of
834 disordered protein. *Proteins* **42**(1): 38-48.

835 Sawyer N, Chen J, Regan L. 2013. All repeats are not equal: a module-based approach to guide repeat
836 protein design. *J Mol Biol* **425**(10): 1826-1838.

837 Schluenzen F, Tocilj A, Zarivach R, Harms J, Gluehmann M, Janell D, Bashan A, Bartels H, Agmon I,
838 Franceschi F, Yonath A. 2000. Structure of functionally activated small ribosomal subunit at
839 3.3 angstroms resolution. *Cell* **102**(5): 615-623.

840 Schrödinger L, 2010. The PyMOL Molecular Graphics System, Version 1.3r1.

841 Scott A, Gaspar J, Stuchell-Breerton MD, Alam SL, Skalicky JJ, Sundquist WI. 2005. Structure and
842 ESCRT-III protein interactions of the MIT domain of human VPS4A. *Proc Natl Acad Sci U S A*
843 **102**(39): 13813-13818.

844 Serasinghe MN, Yoon Y. 2008. The mitochondrial outer membrane protein hFis1 regulates
845 mitochondrial morphology and fission through self-interaction. *Exp Cell Res* **314**(19): 3494-
846 3507.

847 Shatsky M, Nussinov R, Wolfson HJ. 2004. A method for simultaneous alignment of multiple protein
848 structures. *Proteins* **56**(1): 143-156.

849 Sheldrick GM. 2008. A short history of SHELX. *Acta Crystallogr Sect A Found Crystallogr* **64**(1): 112-
850 122.

851 Shen C, Zhang D, Guan Z, Liu Y, Yang Z, Yang Y, Wang X, Wang Q, Zhang Q, Fan S, Zou T, Yin P. 2016.
852 Structural basis for specific single-stranded RNA recognition by designer pentatricopeptide
853 repeat proteins. *Nat Commun* **7**: 11285.

854 Sikorski RS, Boguski MS, Goebel M, Hieter P. 1990. A repeating amino acid motif in CDC23 defines a
855 family of proteins and a new relationship among genes required for mitosis and RNA
856 synthesis. *Cell* **60**(2): 307-317.

857 Söding J, Lupas AN. 2003. More than the sum of their parts: on the evolution of proteins from
858 peptides. *Bioessays* **25**(9): 837-846.

859 Söding J, Remmert M, Biegert A. 2006. HHrep: de novo protein repeat detection and the origin of
860 TIM barrels. *Nucleic Acids Res* **34**(Web Server issue): W137-W142.

861 Stumpp MT, Forrer P, Binz HK, Plückthun A. 2003. Designing repeat proteins: modular leucine-rich
862 repeat protein libraries based on the mammalian ribonuclease inhibitor family. *J Mol Biol*
863 **332**(2): 471-487.

864 Tobin C, Mandava CS, Ehrenberg M, Andersson DI, Sanyal S. 2010. Ribosomes lacking protein S20 are
865 defective in mRNA binding and subunit association. *J Mol Biol* **397**(3): 767-776.

866 Vagin A, Teplyakov A. 2000. An approach to multi-copy search in molecular replacement. *Acta*
867 *Crystallogr D Biol Crystallogr* **56**(Pt 12): 1622-1624.

868 Varadi M, Kosol S, Lebrun P, Valentini E, Blackledge M, Dunker AK, Felli IC, Forman-Kay JD, Kriwacki
869 RW, Pierattelli R, Sussman J, Svergun DI, Uversky VN, Vendruscolo M, Wishart D, Wright PE,

870 Tompa P. 2014. pE-DB: a database of structural ensembles of intrinsically disordered and of
871 unfolded proteins. *Nucleic Acids Res* **42**(Database issue): D326-D335.

872 Wei Y, Kim S, Fela D, Baum J, Hecht MH. 2003. Solution structure of a de novo protein from a
873 designed combinatorial library. *Proc Natl Acad Sci U S A* **100**(23): 13270-13273.

874 Wootton JC. 1994. Non-globular domains in protein sequences: automated segmentation using
875 complexity measures. *Comput Chem* **18**(3): 269-285.

876 Wootton JC, Federhen S. 1996. Analysis of compositionally biased regions in sequence databases.
877 *Methods Enzymol* **266**: 554-571.

878 Zeytuni N, Baran D, Davidov G, Zarivach R. 2012. Inter-phylum structural conservation of the
879 magnetosome-associated TPR-containing protein, MamA. *J Struct Biol* **180**(3): 479-487.

880 Zeytuni N, Zarivach R. 2012. Structural and functional discussion of the tetra-trico-peptide repeat, a
881 protein interaction module. *Structure* **20**(3): 397-405.

882 Zeytuni N, Cronin S, Lefèvre CT, Arnoux P, Baran D, Shtein Z, Davidov G, Zarivach R. 2015. MamA as a
883 Model Protein for Structure-Based Insight into the Evolutionary Origins of Magnetotactic
884 Bacteria. *PLoS One* **10**(6): e0130394.

885 Zhang Z, Kulkarni K, Hanrahan SJ, Thompson AJ, Barford D. 2010. The APC/C subunit Cdc16/Cut9 is a
886 contiguous tetratricopeptide repeat superhelix with a homo-dimer interface similar to Cdc27.
887 *EMBO J* **29**(21): 3733-3744.

888

889
890

Figure legends

901 one or more survivors among these would have become the ancestor(s) of today's TPR proteins. In
902 this more complex scenario, the homology of TPR proteins, which we trace through the comparison
903 of individual hairpins, is still given, but the TPR fold could have arisen from several independent
904 amplifications, and not just a single one.

905

a)

| | | | |
|-------|--|-------|---|
| Ana1 | VQQLFKQGETAESVGDNSQAETIWRKVLQVEPNP | Cal1 | TEQLFKQGEAAESVGNNSQAETIWRQVLQLEPSN |
| Ana2 | GKAYNNLGNALRRQGKLPALTAHQKALQLNPN | Cal2 | GKAYNNLGNALRRQGKLPALTAHQKALQLNPN |
| Ana3 | AEAYVIGNVLNAQGGKPEEVAAYRKAIEFDPKY | Cal3 | AEAYVIGNVLNAQGGKPEEVAAYRKAIEFDPKY |
| Ana4 | AKAYNSLGNALYDQKLEKAEVAAYRKAIEFDHXY | Cal4 | AAAYNGLGNALYDQKLEPAVAAYYQKAIQLDPNY |
| Ana5 | AAAYNGLGNVLYEQKELDEAVAAYRKAIELNPKY | Cal5 | AAAYNGLGNALRDQKLEPAIAAFQKAIQLNPNF |
| Ana6 | ATAYNNLGNALSDQKLEDEAVAAYQ EAIKLNPKD | Cal6 | AAAYNGLGNALYDQKLEPAVAAFQKAIQLNPNF |
| Ana7 | ATAYNNLGI ALS DQKLEDEAVAAYQKAIELDPKY | Cal7 | AAAYNGLGNALYDQKLEPAVAAYYQKAIQLDPNF |
| Ana8 | ATAYNNLGNALSDQKLEDEAVAAYQKAIELDPKY | Cal8 | AFAY YNL GNALYDQKLEPAIAAFQKAIQLNPNF |
| Ana9 | ATAYNNLGNALSDQKLEDEAVAAYQKAIELDPKY | Cal9 | AFAYNGLGNALYDQKLEPAIAAFQKAIQLNPNF |
| Ana10 | ATAYNNLGNAL RG QKLEDEAVAAYQKAIELNPKY | Cal10 | AFAYNGLGNALYDQKLEPAIAAFQKAIQLNPNF |
| Ana11 | ATAYNNLGI ALS DQKLEDEAVAAYQKAIELNPKD | Cal11 | AFAYNGLGNALYDQKLEPAIAAFQKAIQLNPNF |
| Ana12 | ATAYNNLGI ALS DQKLEDEAVAAYQKAIELDPKD | Cal12 | AFAYNGLGNALYDQKLEPAIAAFQKAIQLNPNF |
| Ana13 | AAVYNNLGNALSDQKLEAISNYKTALSLPEDT | Cal13 | AFAYNGLGNALYDQKLEPAIAAFQKAIQLNPNF |
| Ana14 | TLANNLGLALQDQEKFAEAIKYFDKAEELDPNF | Cal14 | AFAYNGLGNALYDQKLEPAIAAFQKAIQLNPNF |
| | | Cal15 | AFAYNGLGNALYDQKLEPAIAAFQKAIQLNPNF |
| | | Cal16 | AFAYNGLGNALYDQKLEPAIAAFQKAIQLNPNF |
| | | Cal17 | AFAYNGLGNALYDQKLEPAIAAFQKAIQL DPND |
| | | Cal18 | ANAYNNL GNALYDQKLEPAIAAYYQKAIQLNPNF |
| | | Cal19 | AAAYNNLGV ALS DQKLEPAIAAYQKVLTLPEDT |
| | | Cal20 | TAANNLGLLVFQEQGKLEKQAIIDYFDKSEALDPDY |
| Mic1 | LEQLWQQGETAQAQKKYPEAERIWRQIIQLDPNS | Cya1 | IDQLFQQGRTAGKMGKYTEAEIIFRRVIELDPNL |
| Mic2 | AVAFSNLCAALFRQNKLEAPIFCQKALALDPKL | Cya2 | ADAYNNLGNALYYQ GKLEDAIAAYYQKAIQLNPN D |
| Mic3 | PETY NNLGNVLYNQKLETEAEEMYRR TIEL DDK F | Cya3 | ADAYNNLGNALSDQ GKLEEAIAAYYQKAIQLNPN Y |
| Mic4 | VYAYNNLGLVLYDQKLEAEEMYRR TIEL DDK F | Cya4 | ADAYNNLGI ALSDQGKLEEAIAAYYQKAIQLNPN F |
| Mic5 | ALVY NNLGLVLYDQKLEAEEMYRR TIEL DDK F | Cya5 | TQAY NNLGI ALS DQGKLEEAIAAYYQKAIQLNPN Y |
| Mic6 | VYAYNGLGNV LR AQNKLEAEEMYRRALALDDQ Y | Cya6 | ADAY NNLGNAL FD QGKLEDAIAAYYQKAIQL DPND |
| Mic7 | VDAY NNLGNVLYDQKLEAEEMYRRALALDDQ Y | Cya7 | ANAY NNLGAAL YK QGLEEAIAAYYQKAIQLNPN L |
| Mic8 | VYAYNGLGNVLYDQKLEAEEMYRR TIEL DDK Y | Cya8 | AEAY NNLGV ALS DQGKRDEAIAAYYQKAIQLNPN L |
| Mic9 | VYAY KGL GNVLYNQKLEAEEMYRR IAL DDQ Y | Cya9 | AEAY NNLGV ALS DQGKRDEAIAAYYQKAIQLNPN F |
| Mic10 | VHAY NSLGNVLYNQKLEAEEMYRRALALDDQ Y | Cya10 | ALAYNNLGV ALS DQGKRDEAIAAYYQKAIQLNPN F |
| Mic11 | VHAY NNLGNVLYDQKLEAEEMYRRALALDDQ Y | Cya11 | ALAYNNLGV ALS DQGKRDEAIAAYYQKAIQLNPN F |
| Mic12 | VPAY HNLGNVLYNQKLEAEEMYRRALALDDK F | Cya12 | ALAYNNLGV ALS RNQGKRDEAIAAYYQKAIQL DPND |
| Mic13 | VYAYNNLGNVLYDQKLEAEEMYRRAL DL PD DT | Cya13 | ANAY NNLGLAL RN QGKRDEA IT AYYQKAIQLNPN F |
| Mic14 | TLAHNNLGLLLQEQQGLEAAIAEFAEKATKIDPQ Y | Cya14 | ALAYNNLGNAL YS QGKREEAIAAYYQKAIQLNPN F |
| | | Cya15 | ALAYNNLGNALSDQGKRDEAIAAYYQKAIQLNPN F |
| | | Cya16 | ALAYNNLGNALSDQGK LN EAIAT Y QKAIQLNPN F |
| | | Cya17 | ALAYNNLGNAL KD QGKLEAIAAYYQKALSLPEDT |
| | | Cya18 | TLAHNNLGLVYQ P QGKLEALREYEAALKIDPK F |
| fil1 | INQLFEQGNTAQQEGRYAEAEIWRQILEANPDN | | |
| fil2 | AGAY NNLGV ALS YLNGLQLEAVSAYQQAIALDPD Y | | |
| fil3 | AIAYNNLGI ALS RNGLQLEAV EAY QQAIALDPD F | | |
| fil4 | AIAY YNL GI ALS FDLQLEAVSAYQQAIALDPD D | | |
| fil5 | AIAYNNLGNALS NL GLQLEAVSAYQQAIALDPD D | | |
| fil6 | AIAY YNL GNALS NL GLQLEAVSAYQQA IT LDPD Y | | |
| fil7 | AIAYNNLGNALS NL GLQLEAV EAY QQAIALDPD D | | |
| fil8 | ADAY NNLGNAL RD LQLEAVSAYQQAIALDPD F | | |
| fil9 | ADAY NNLGI ALS RD L QLEAV EAY QQAIALDPD D | | |
| fil10 | AIAYNNLGV ALS YLNGLQLEAVSAYQQAIALDPD N | | |
| fil11 | AFAYNNLGYAY Q QGNLEAAIT EY KKAI AL APN Y | | |

b)

>gb|CP011382.1|c:3792015-3790486 Calothrix sp. 336/3, complete genome

GCT**G**CTGCTTACAATGGTCTCGGTAATGCCCTATATGACCAGAAGAACTAGAACCAGC**AG**TTGCCGCCT**AC**CAAAAAGCCAT**T**CAACT**C**ACCCTAAC**T**AT
GCC**G**CTGCTTAC**TACA**ATCTCGGCAATGCCCTAT**G**AGACGACGAGAAGAACTAGAACCAGCGATCGCCGCCTTCCAAAAAGCCATACA**ACT**CAAC**C**TAA**C**TTT
GCT**G**CTGCTTACAATGGTCTCGGTAATGCCCTATATGACCAGAAGAACTAGAACCAGC**AG**TTGCCGCCTTCCAAAAAGCCATACA**ACT**CAAC**C**TAA**C**TTT
GCT**G**CTGCTTACAATGGTCTCGGTAATGCCCTATATGACCAGAAGAACTAGAACCAGC**AG**TTGCCGCCT**AC**CAAAAAGCCAT**T**CAACT**C**ACCCTAAC**T**TT
GCTTTTGCTTAC**TACA**ATCTCGGCAATGCCCTATATGACCAGAAGAACTAGAACCAGCGATCGCCGCCTTCCAAAAAGCCATACA**ACT**CAAC**C**TAA**C**TTT
GCTTTTGCTTACAATGGTCTCGGTAATGCCCTATATGACCAGAAGAACTAGAACCAGCGATCGCCGCCTTCCAAAAAGCCATACA**ACT**CAAC**C**TAA**C**TTT
GCTTTTGCTTACAATGGTCTCGGTAATGCCCTATATGACCAGAAGAACTAGAACCAGCGATCGCCGCCTTCCAAAAAGCCATACA**ACT**CAAC**C**TAA**C**TTT
GCTTTTGCTTACAATGGTCTCGGTAATGCCCTATATGACCAGAAGAACTAGAACCAGCGATCGCCGCCTTCCAAAAAGCCATACA**ACT**CAAC**C**TAA**C**TTT
GCTTTTGCTTACAATGGTCTCGGTAATGCCCTATATGACCAGAAGAACTAGAACCAGCGATCGCCGCCTTCCAAAAAGCCATACA**ACT**CAAC**C**TAA**C**TTT
GCTTTTGCTTACAATGGTCTCGGTAATGCCCTATATGACCAGAAGAACTAGAACCAGCGATCGCCGCCTTCCAAAAAGCCATACA**ACT**CAAC**C**TAA**C**TTT
GCTTTTGCTTACAATGGTCTCGGTAATGCCCTATATGACCAGAAGAACTAGAACCAGCGATCGCCGCCTTCCAAAAAGCCATACA**ACT**CAAC**C**TAA**C**TTT
GCTTTTGCTTACAATGGTCTCGGTAATGCCCTATATGACCAGAAGAACTAGAACCAGCGATCGCCGCCTTCCAAAAAGCCATACA**ACT**CAAC**C**TAA**C**TTT
GCTTTTGCTTACAATGGTCTCGGTAATGCCCTATATGACCAGAAGAACTAGAACCAGCGATCGCCGCCTTCCAAAAAGCCATACA**ACT**CAAC**C**TAA**C**TTT
GCTTTTGCTTACAATGGTCTCGGTAATGCCCTATATGACCAGAAGAACTAGAACCAGCGATCGCCGCCT**AC**CAAAAAGCCATACA**G**CT**C**ACCCTAAC**G**AT
GCC**AA**C**G**CTTACA**CA**ATCTCGGCAATGCCCTATATGACCAGAAGAACTAGAACCAGCGATCGCCGCCT**AC**CAAAAAGCCAT**T**CAACT**C**ACCCTAAC**T**TT

906

907 **Figure 1—figure supplement 1.** Multiple sequence alignments of recently amplified TPR repeat units.

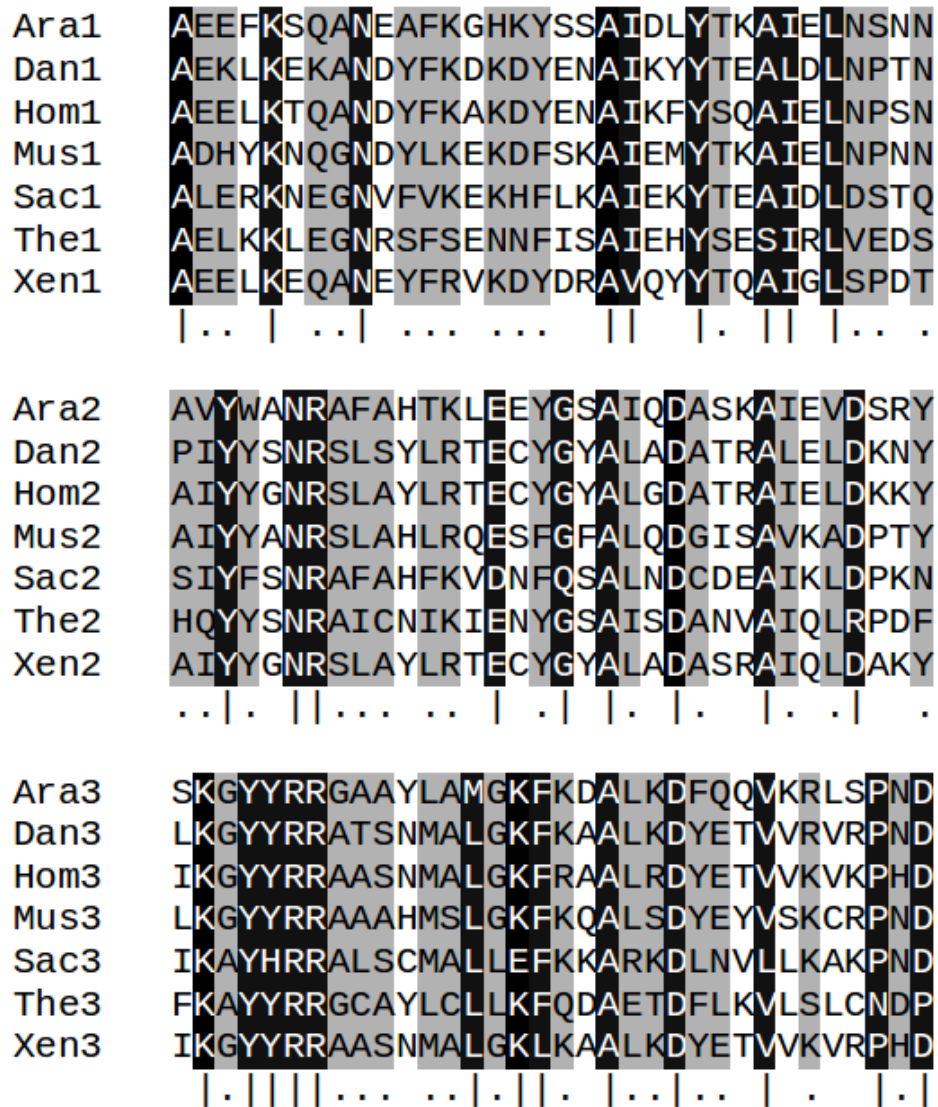
908 a) Alignments of the TPR units used for the phylogeny in Figure 1a. Residues different from the most

909 common one in each column are shown in bold face and highlighted in yellow. Abbreviations: Ana:

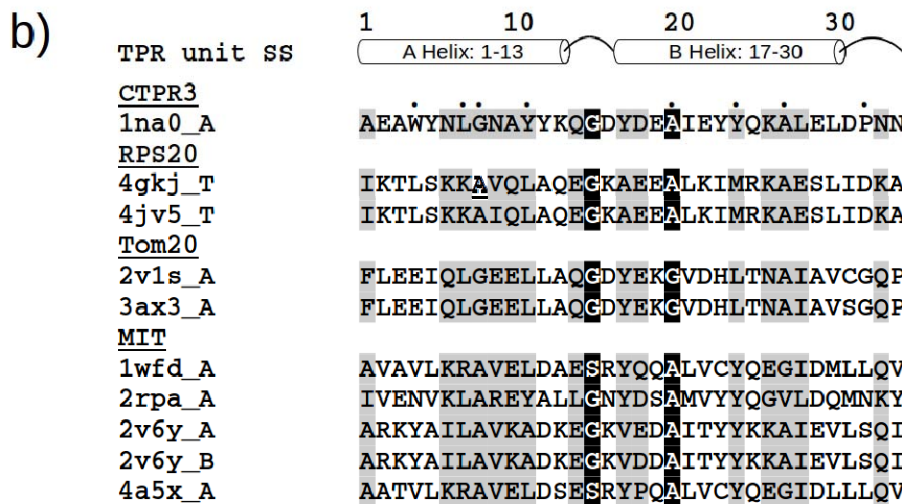
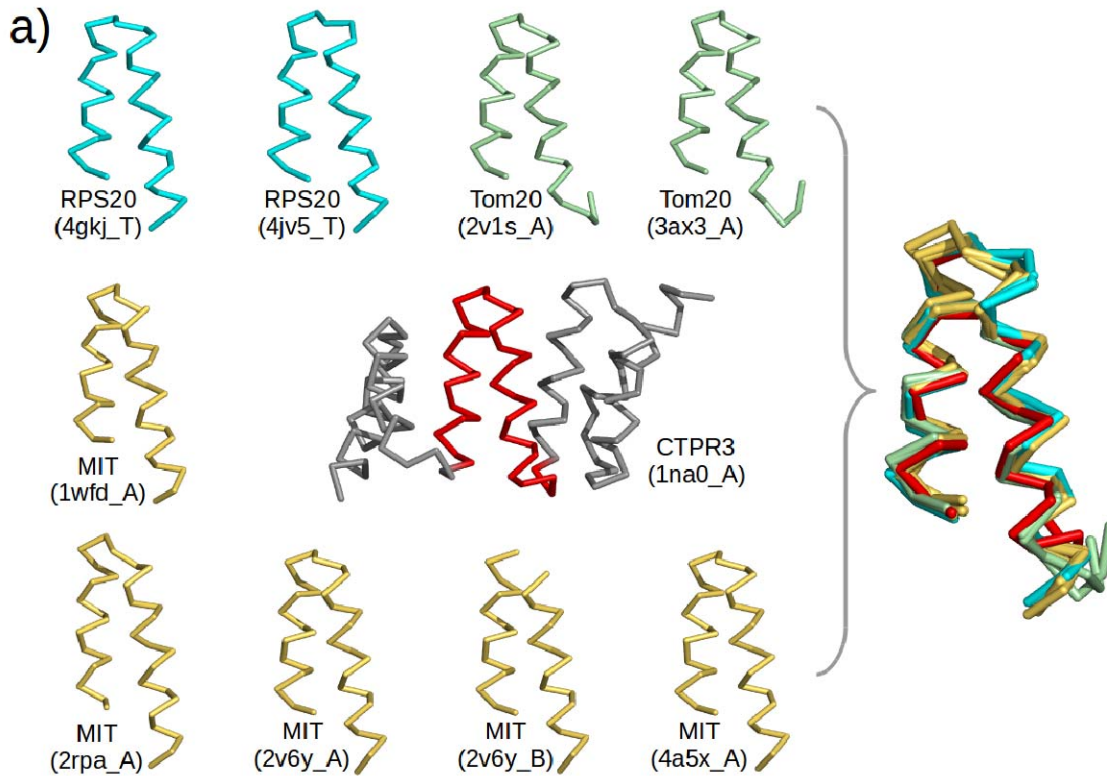
910 Anabaena sp. 90 (gi: 752818954, accession: WP_041458168.1); Cal: Calothrix sp. 336/3 (gi:

911 821031795, accession: WP_046815017.1); Cya: *Cyanothece* sp. PCC 8801 (gi: 501590504, accession:
912 WP_012594639.1); fil: filamentous cyanobacterium ESFC-1 (gi: 740500649, accession:
913 WP_038331513.1); Mic: *Microcystis aeruginosa* SPC777 (gi: 513477764, accession: EPF24195.1). b)
914 The corresponding alignment of the DNA sequences for the most recently amplified TPR units, Cal4-
915 Cal18, of which the central repeats, Cal9-Cal16, are fully identical. Synonymous mutations
916 (highlighted in gray) are found at less than 1% of the nucleotides, illustrating the recent time point of
917 the amplification. Non-synonymous mutations (highlighted in yellow) are about 2.5 times as frequent
918 as synonymous ones.

919

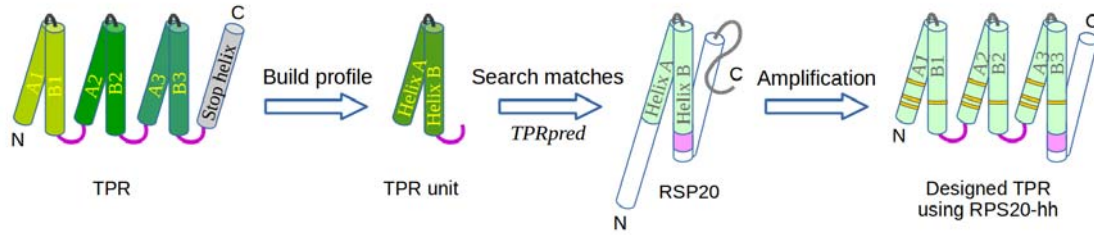


920
921 **Figure 1—figure supplement 2.** Multiple sequence alignment of the three TPR repeat units in
922 serine/threonine-protein phosphatase 5 from seven taxa. Columns with identify $\geq 80\%$ are
923 highlighted in black and marked by vertical bars (|); column with identify $< 80\%$ but $\geq 50\%$ are
924 highlighted in gray and marked by dots (.). Abbreviations: Ara: *Arabidopsis thaliana* (gi: 18406066,
925 accession: NP_565985.1); Dan: *Danio rerio* (gi: 126158897, accession: NP_001014372.2); Hom: *Homo*
926 *sapiens* (gi: 5453958, accession: NP_006238.1) ; Mus: *Musca domestica* (gi: 557765703, accession:
927 XP_005182549.1); Sac: *Saccharomyces cerevisiae* S288c (gi: 398365781, accession: NP_011639.3);
928 The: *Theileria annulata* strain Ankara (gi: 84994100, accession: XP_951772.1); Xen: *Xenopus tropicalis*
929 (gi: 56118654, accession: NP_001007891.1).



930 **Figure 2. a)** Structure gallery of non-repetitive helical hairpins in the PDB that share both sequence
 931 and structure similarity to TPR unit hairpin. Only the 34 amino-acid helical hairpins are shown. The
 932 helical hairpins in 30S ribosomal protein s20 (RPS20), mitochondrial import receptor subunit
 933 (Tom20), and microtubule interacting and transport domain (MIT) are depicted in cyan, green, and
 934 yellow, respectively. The structure of a TPR with a consensus sequence, CTPR3, is shown in the
 935 center with the middle TPR unit highlighted in red. PDB IDs and chain names of the proteins are given
 936

937 in parentheses. In the superposition, all helical hairpins are superimposed onto the middle TPR unit
938 of CTPR3. **b)** Multiple sequence alignment of the helical hairpin sequences listed in a). The eight TPR
939 signature positions are marked by dots in CTPR3. Columns with sequence identity $\geq 80\%$ are in black,
940 and columns with sequence identity $\geq 50\%$ are in gray.



941

942

Figure 3. The design of TPR using RPS20. RPS20-hh is identified by TPRpred to match the sequence

943

profile of TPR units. Their structures are also very similar (helices are shown as cylinders), except for

944

the last four residues (colored in light and dark magenta). We designed a TPR protein using a RPS20-

945

hh with up to five mutations (yellow strips) in each repeat unit. The C-terminal loop in the TPR unit

946

(dark magenta loop) is used to replace the corresponding C-terminus (light magenta cylinder) of

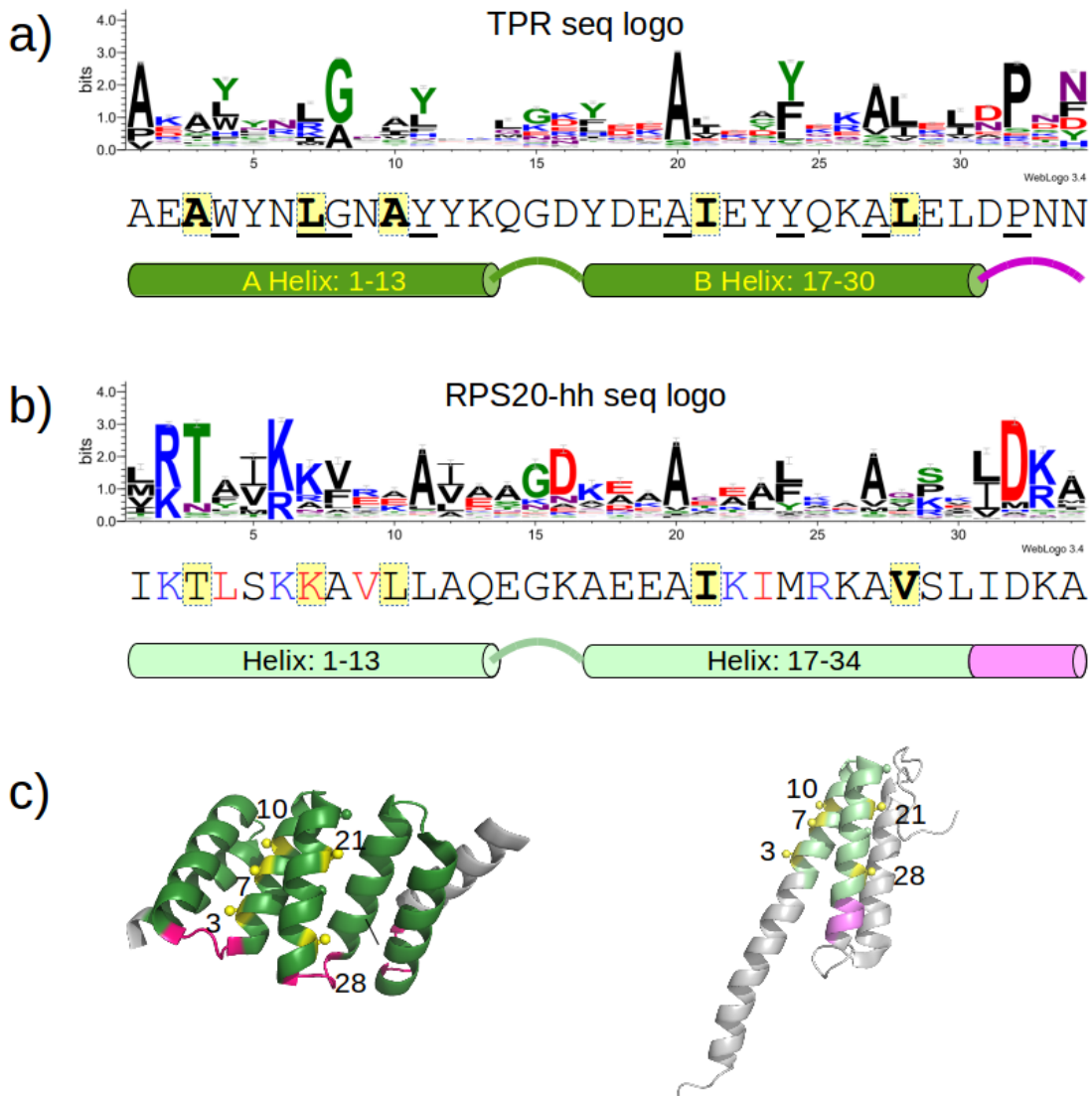
947

RPS20-hh to connect adjacent repeats. The C-terminal helix in RPS20 (white cylinder) was used as the

948

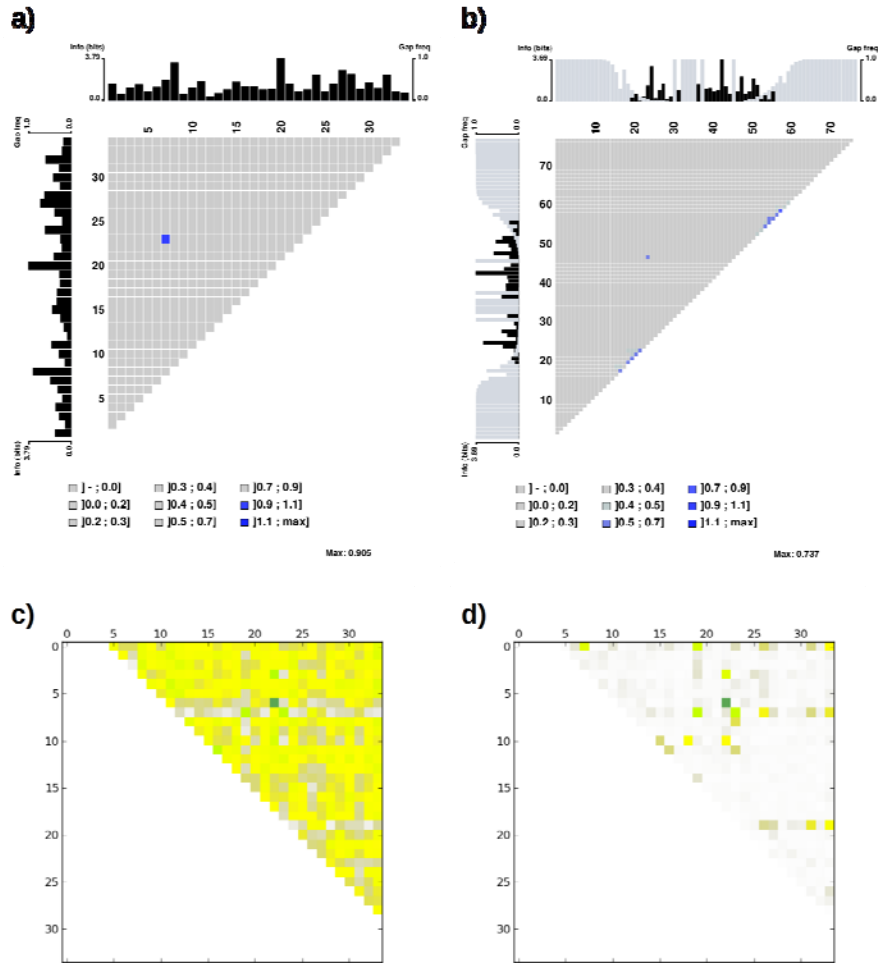
stop helix in the design.

949



950 **Figure 4.** a) Sequence logo of the TPR motif. A TPR consensus sequence (Main et al., 2003b) (PDB:
 951 1na0, chain A) and its secondary structure determined by DSSP (Kabsch and Sander, 1983) are
 952 aligned below the sequence logo. The eight TPR signature positions are underscored in the consensus
 953 sequence. The five interface positions are highlighted in yellow. **b)** Sequence logo of RPS20-hh. The
 954 RPS20-hhta sequence and its predicted secondary structure using Quick2D (Biegert et al., 2006) is
 955 aligned below the sequence logo. The derived interface positions are highlighted in yellow. The four
 956 residues subjected to mutations are colored in red. The four positively charged residues selected for
 957 mutation to lower the surface charge are in blue. **c)** The locations of the interface positions displayed
 958 on a TPR (left) and a RPS20 structure (right). In both structures, the interface positions are labeled
 959 and highlighted as yellow spheres. The TPR structure is CTPR3 (PDB: 1na0, chain A), which is shown
 960

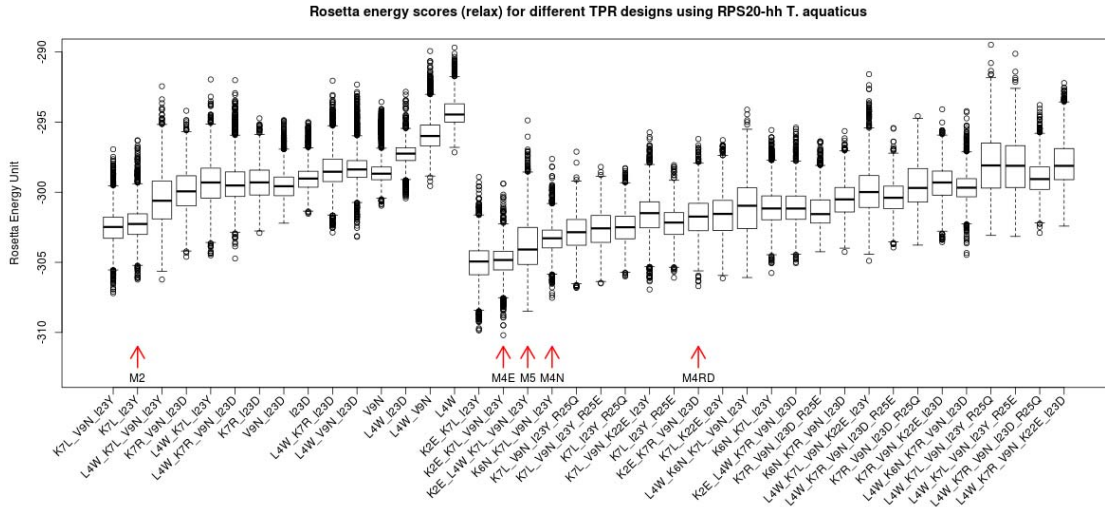
961 as a cartoon and is colored using the same scheme as the secondary structure representation in a).
962 The stop helix is in gray. The RPS20 structure is from *T. thermophilus* (PDB: 4gkj, chain T), in which
963 the RPS20-hh fragment is colored using the same scheme as the secondary structure representation
964 in b).
965 The sequence logos were generated using WebLogo (Crooks et al., 2004). Sequences from
966 representative proteome 75% (Chen et al., 2011) downloaded from Pfam families *TPR_1* and
967 *Ribosomal_S20p* were used as input to WebLogo (9338 and 972 sequences, respectively). The
968 structures were rendered using PyMOL (Schrödinger, 2010).
969



970

971 **Figure 4—figure supplement 1.** Mutual information plot (a and b) and direct coupling analysis plot (c
 972 and d) for TPR repeat sequences. The subfigures a) and c) were generated using the seed alignment
 973 sequences from Pfam family *TPR_1* (558 sequences. Sequence Q29585_PIG/28-61 was removed as it
 974 contains unknown residue X). The largest mutual information value is observed between position 7
 975 and 23. The subfigures b) and d) were generated using the multiple alignment of representative
 976 proteomes rp75 sequences from Pfam family *TPR_1* (9338 sequences). The largest non-local mutual
 977 information value was observed between position 24 and 47, corresponding to position 7 and 23
 978 using TPR repeat numbering. Alignments were taken from Pfam 27.0. Subfigures a) and b) were
 979 generated using MatrixPlot. Subfigures c) and d) were generated using DCA Workbench
 980 (<http://dca.rice.edu/portal/dca/workbench>).

981

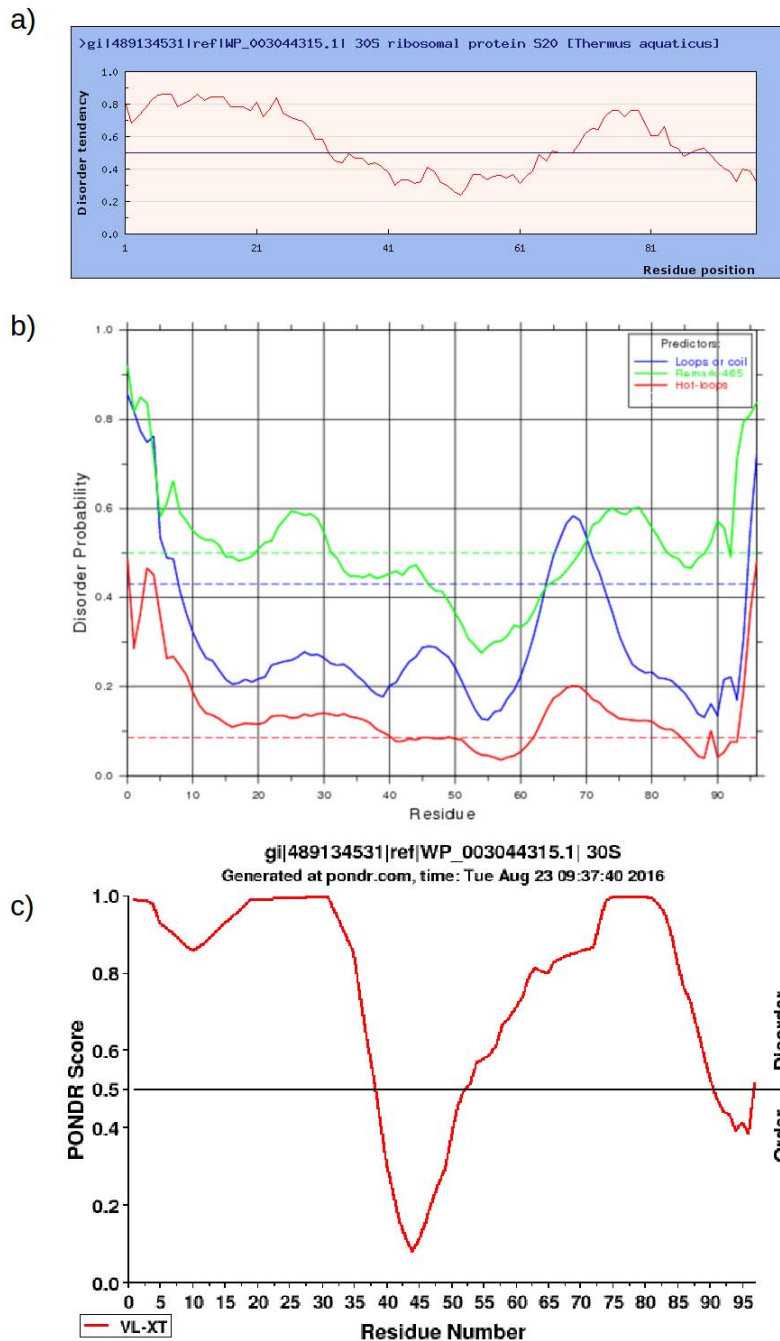


982

983 **Figure 4—figure supplement 2.** Rosetta energy scores (*fixbb+relax*) for TPR designs based on RPS20-
 984 hhta sequence and various sets of mutations. The scores for the designs are shown in two groups:
 985 the group to the left are combinations involving only primary mutations (see supplementary file 1E).
 986 The group to the right are designs involving both primary and secondary mutations (supplementary
 987 file 1E). The design variants are sorted by the average of the lowest 10% scores. The designs tested in
 988 the lab are marked by red arrows (M2, M4E, M5, M4N, M4RD). The *in silico* simulation was
 989 performed using Rosetta 3.4.

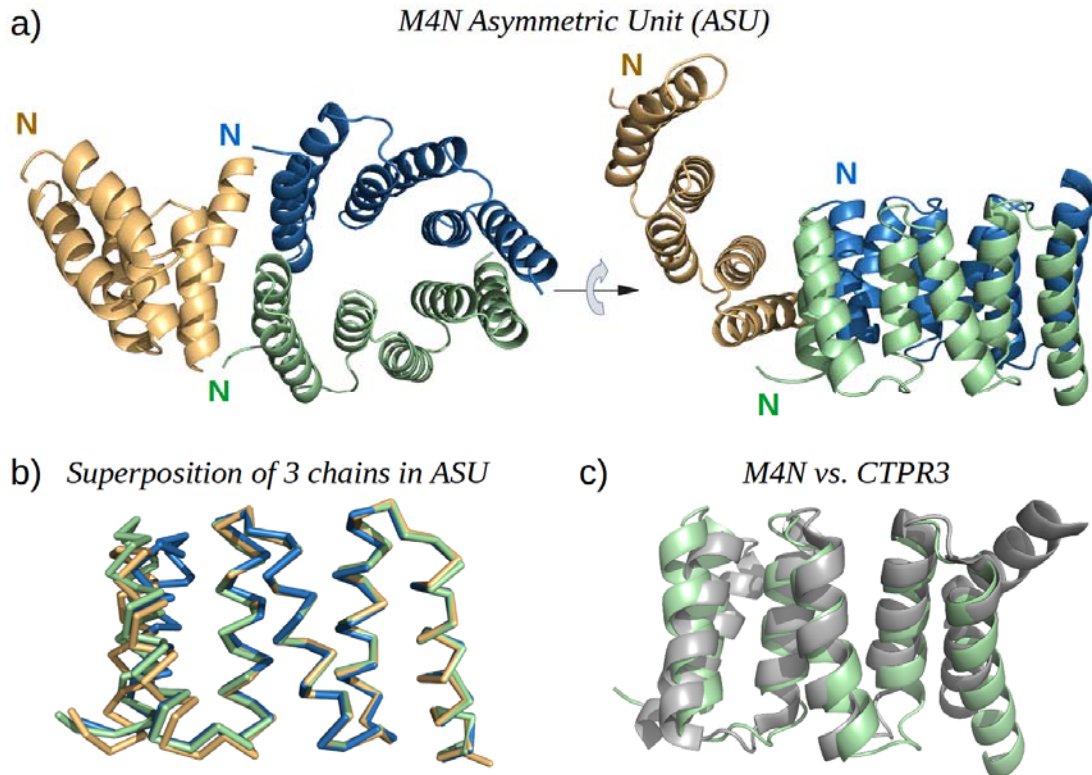
990

991



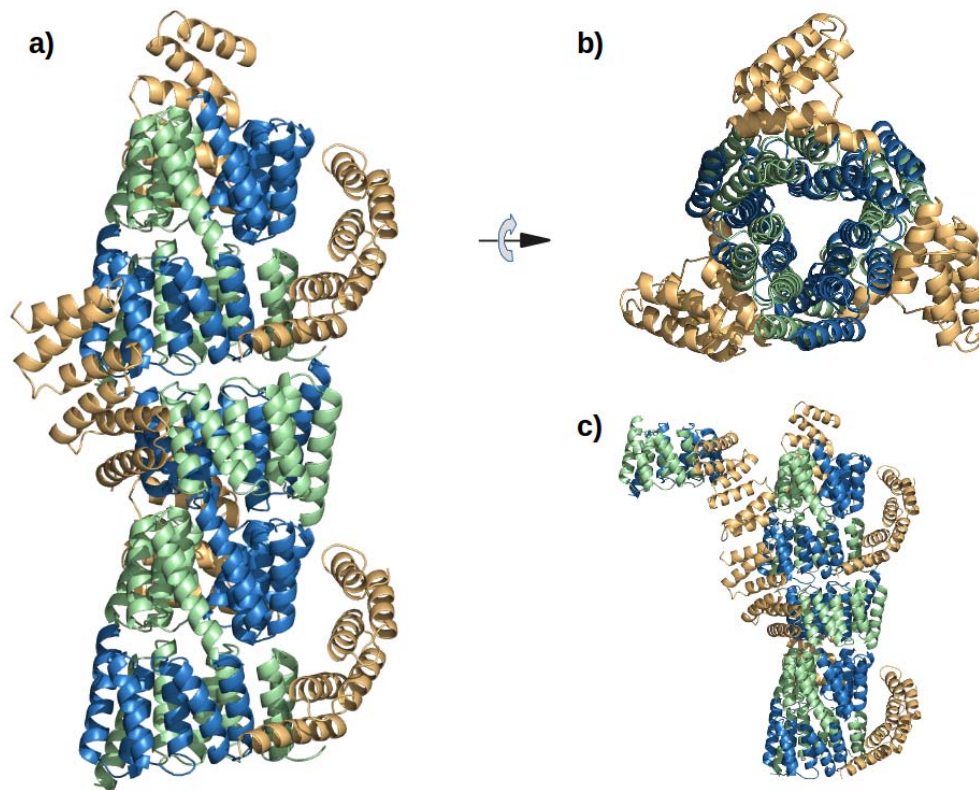
992

993 **Figure 4—figure supplement 3.** Prediction of intrinsically disordered regions in RPS20 of *Thermus*
 994 *aquaticus* (NCBI gi: 489134531, accession: WP_003044315.1) using a) IUPred
 995 (<http://iupred.enzim.hu/>) ; b) DisEMBL (<http://dis.embl.de/>) and c) PONDR
 996 (<http://www.pondr.com/>).



997
 998 **Figure 5.** a) The X-ray structure of M4N. The three chains A, B and C in the asymmetric unit are
 999 colored green, blue and yellow, respectively. Chains A and B form a dimer. b) Superposition of the
 1000 three chains. Only C α traces are shown for clarity. c) Superposition of M4N (chain A, green) and the
 1001 designed consensus TPR CTPR3 (PDB: 1na0, chain A, gray).

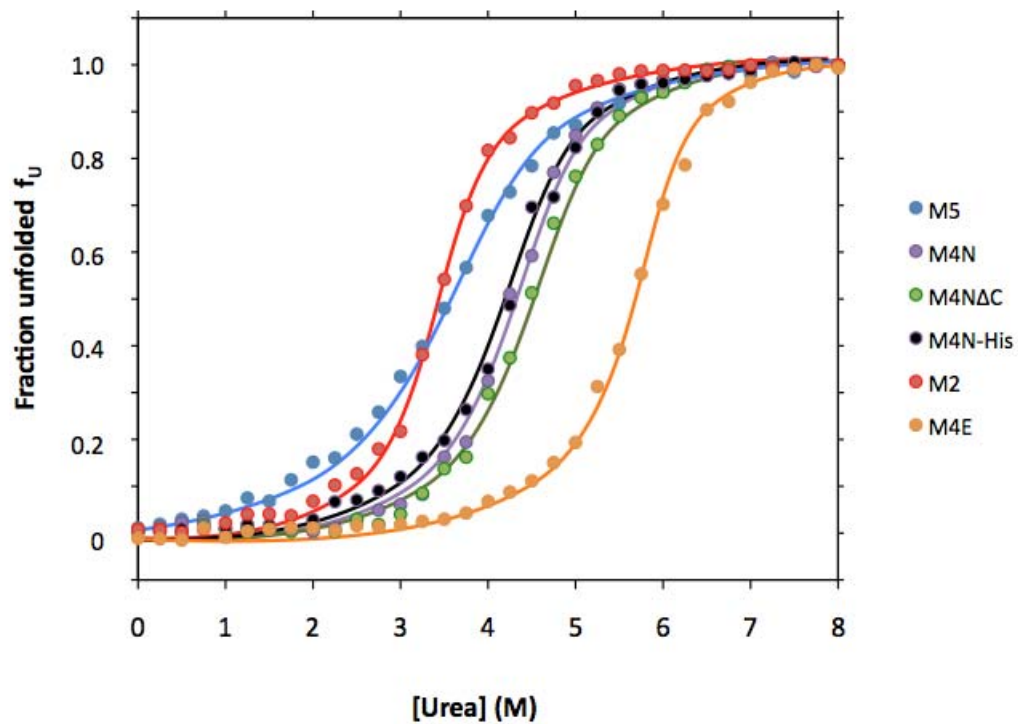
1002



1003

1004 **Figure 5—figure supplement 1.** The interaction of M4N molecules in the crystal. **a)** Five adjacent
1005 ASUs are depicted. Chain A (green) and B (blue) form a dimer, while chain C (yellow) packs its C-
1006 terminus to the N-termini of chains A and B. **b)** Top view. **c)** An additional ASU (top-left) is shown to
1007 demonstrate the packing of N-termini of chains C.

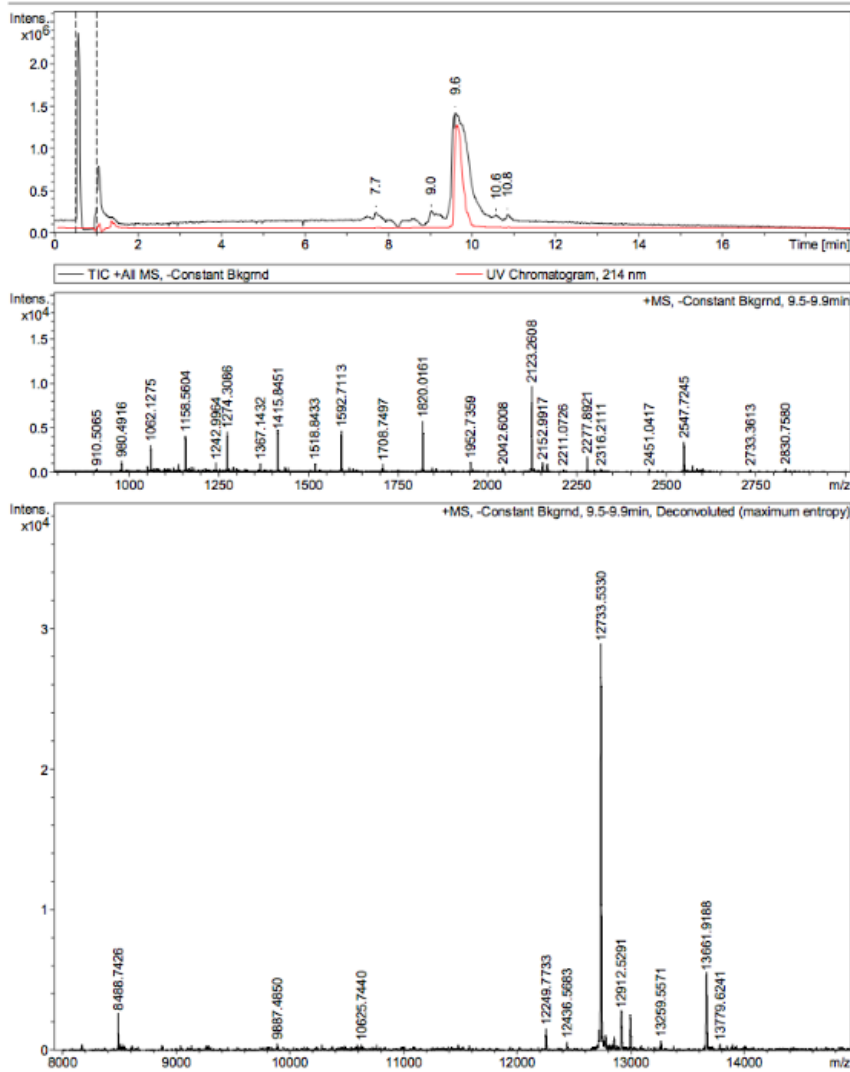
1008



1009

1010 **Figure 5—figure supplement 2.** Urea denaturation of designed TPR repeats. Urea-induced
 1011 equilibrium unfolding at 23°C was monitored by circular dichroism at 222 nm. Data were converted
 1012 to the fraction of unfolded protein f_U and fitted to a two-state model. The protein concentration was
 1013 15 μ M. See supplement file 1F for obtained parameters.

1014



>M4N
 MGNS
IKTLSNLANLLAQEGKAEAAIKYMRKAVSLDPNN
IKTLSNLANLLAQEGKAEAAIKYMRKAVSLDPNN
IKTLSNLAVLLAQEGKAEAAIKYMRKAVSLIDKA
AKGSTLHKNAARRKSRLMRKVQKL

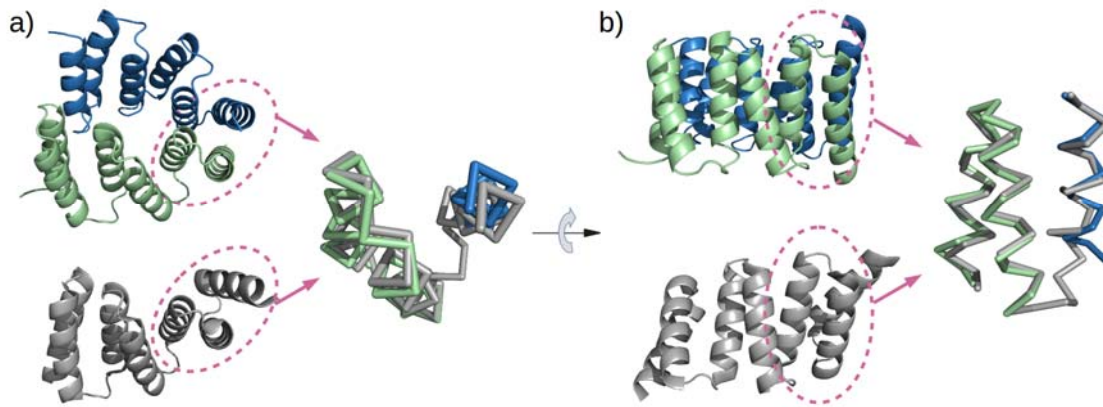


1015

1016 **Figure 5—figure supplement 3.** Mass spectrometry (MS) analysis of M4N. The M4N fragment with a
 1017 mass of 12733.533 Da in MS is underlined and highlighted in blue (theoretical mass 12733.77 Da).

1018 The C-terminus of M4N as observed in the crystal structure is marked by a red arrow.

1019



1020

1021

Figure 6. Mimicry of the stop helix in the M4N dimer. The C-terminal TPR unit in chain A (green) and

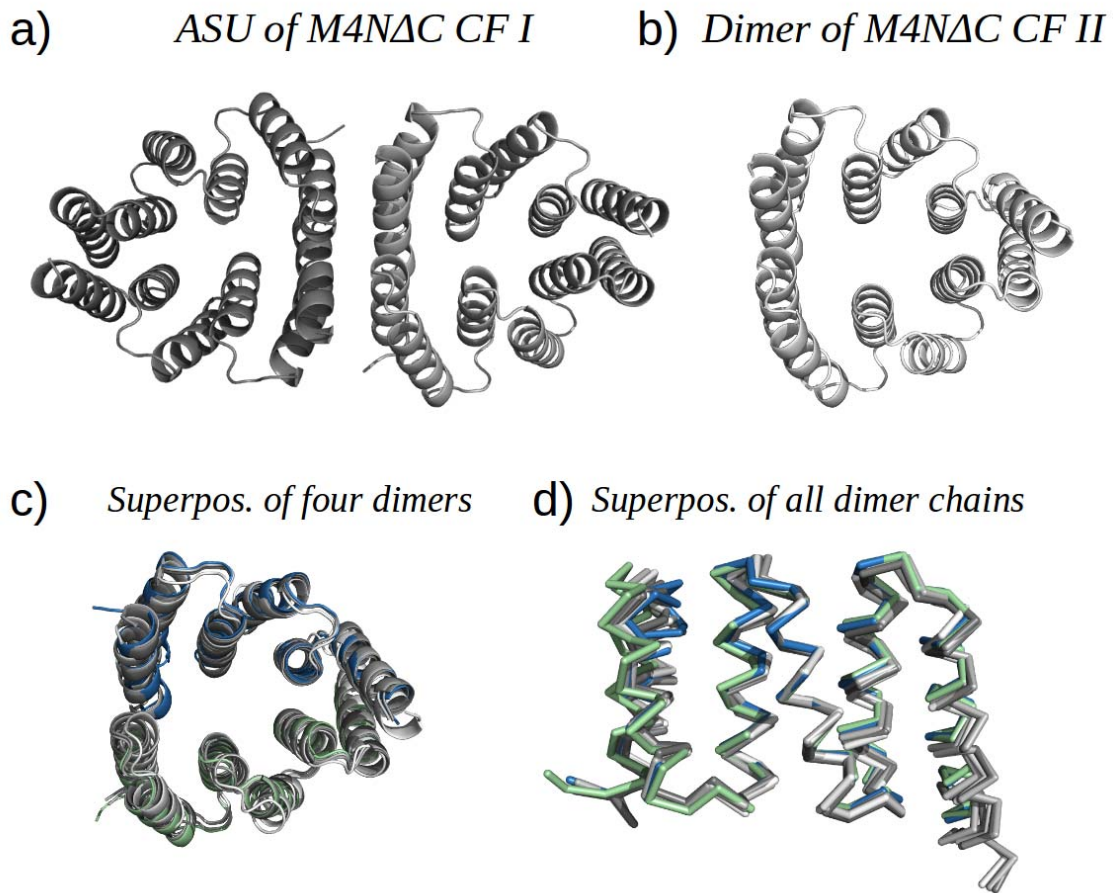
1022

the C-terminal helix B3 in chain B (blue) are superposed to the last TPR unit plus the stop helix in

1023

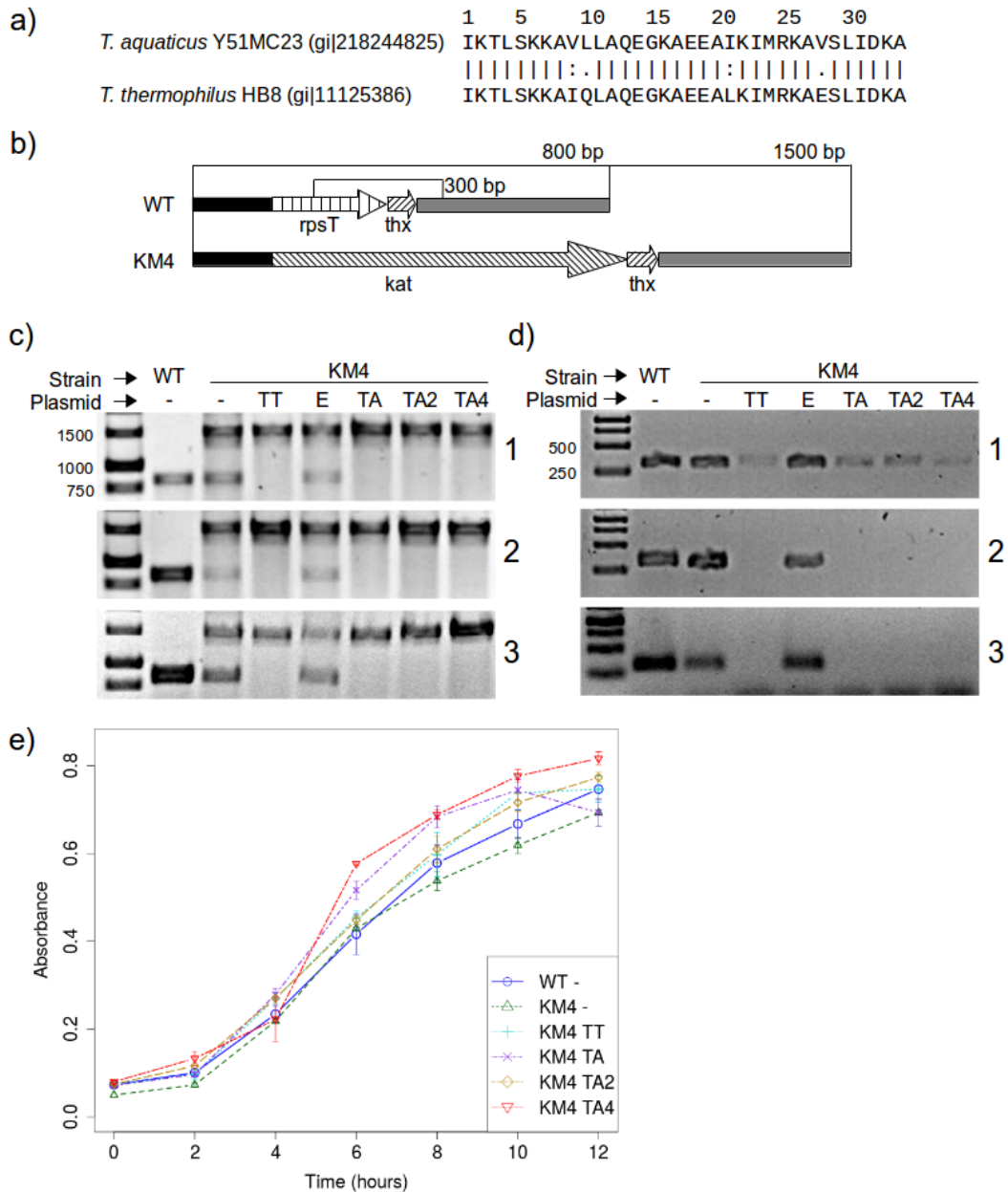
CTPR3 (gray).

1024



1025
 1026 **Figure 7.** M4NΔC structures of two different crystal forms and their comparison to the M4N dimer. **a)**
 1027 Two dimers in the ASU of M4NΔC CF I. **b)** Dimer constructed by applying the crystallographic
 1028 symmetry to the single chain in the ASU of M4NΔC CF II. **c)** Superposition of all the four M4N and
 1029 M4NΔC dimers. The M4N dimer is in green and blue. The three M4NΔC dimers are in different
 1030 shades of gray as in a) and b). **d)** Superposition of all the chains in the M4N and M4NΔC dimers (eight
 1031 chains in total). Only C α traces of proteins are shown for clarity.

1032



1033

1034

Figure 8. RPS20 variants M2 and M4N are functional proteins. **a)** The 34 amino-acid long RPS20-hh

1035

fragments in *T. aquaticus* and *T. thermophilus* differ only at four positions, including two conservative

1036

mutations (V9I and I21L). **b)** Scheme of the *rpsT* region before (upper) and after (lower) substitution

1037

of *rpsT* with the kanamycin resistance cassette (*kat*). Base pair (bp) values indicate the PCR products

1038

that can be amplified. Regions depicted with the same pattern are identical. Regions in solid black

1039

and gray also contain genes which are not marked for clarity. **c)** PCR to detect substitution of *rps20*

1040

by the *kat* gene and **d)** PCR to detect the presence of chromosomal *rpsT* in *T. thermophilus* strains

1041

(WT: *T. thermophilus* HB8; KM4: *T. thermophilus* KM4) carrying various plasmids (TT: pJJSpro-rps20Tt;

1042 E: pJJSpro; TA: pJJSpro-rpsTTa; TA2: pJJSpro-rpsTTaM2; TA4: pJJSpro-rpsTTaM4N; -: No plasmid) after
1043 sequential grow under different selective pressures (1: 30 µg/ml kanamycin; 2: 120 µg/ml kanamycin;
1044 3: 0 µg/ml kanamycin). e) Corresponding growth curves of the host bacteria with various
1045 substitutions and plasmids.

1046

1047 **Supplementary file 1:** Further supporting computational and experimental results.

1048 **Section A:** Sequence variation in RPS20-hh at positions 6, 7, 9 and 23 (TPR unit numbering) observed
1049 in RPS20 sequences.

1050 **Section B:** Most commonly observed amino acids in RPS20-hh.

1051 **Section C:** List of putative TPR homologs identified in the PDB by sequence and structure analysis.

1052 **Section D:** RPS20-hh sequences that resemble a TPR profile according to TPRpred.

1053 **Section E:** Mutations tested in silico on RPS20-hh for TPR design.

1054 **Section F:** Biophysical parameters of designed TPRs.

1055 **Section G:** Primary structures of M4N molecules observed in the crystal structures.

1056 **Section H:** Crystallization conditions, and data collection/refinement statistics.

1057 **Section I:** Detailed structure comparison results of different chains in M4N structures, and of M4N to
1058 CTPR3.

1059 **Section J:** SEG prediction of low-complexity regions in RPS20-hhta.

1060

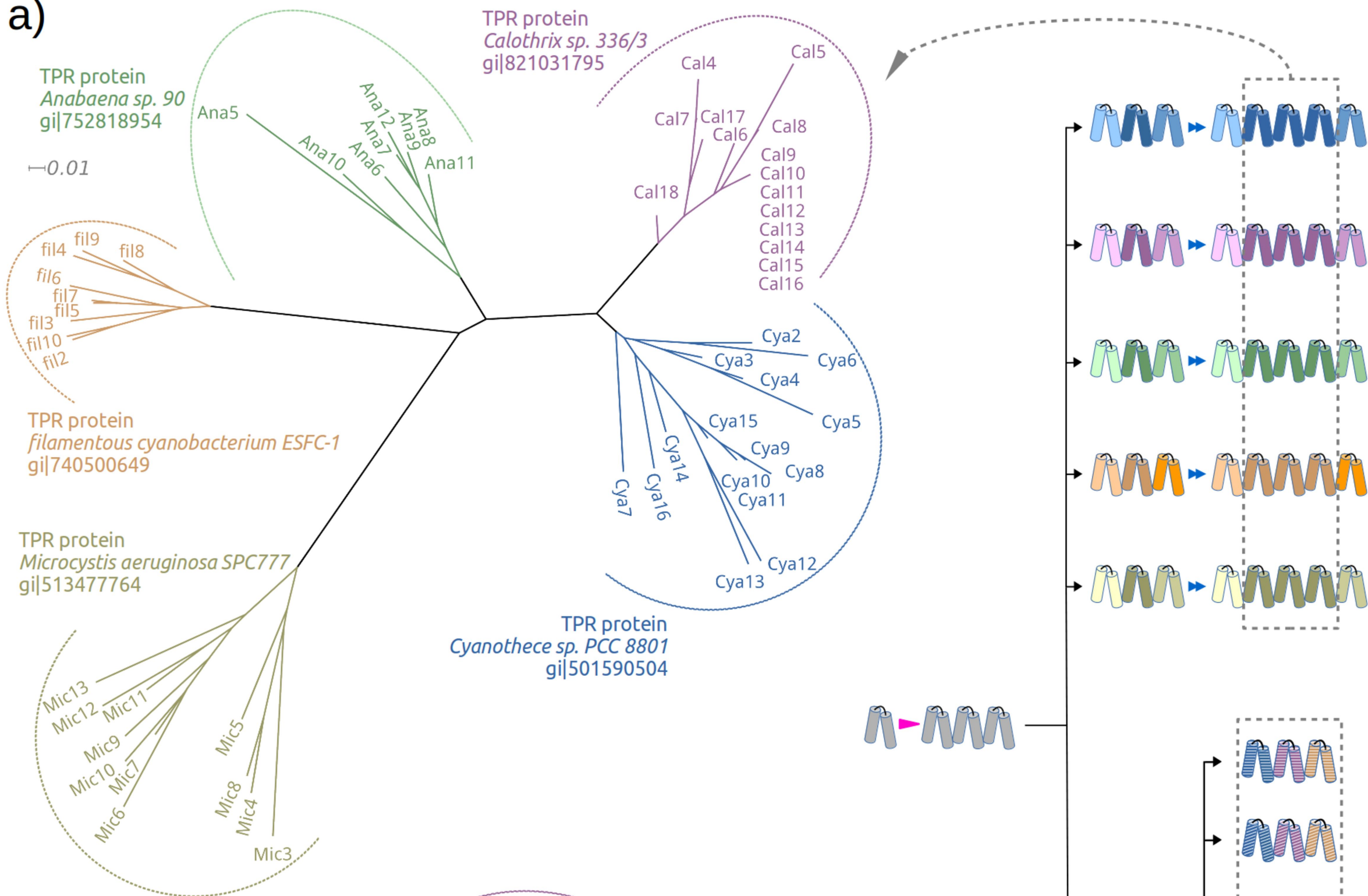
1061

| Name | Mutations | Sequence |
|-------|--------------------------|--|
| M0 | - | NS IKTLSK K AVLLAQEGKAEAAIKIMRKA V SLDPNN IKTLSK K AVLLAQEGKAEAAIKIMRKA V SLDPNN IKTLSK K AVLLAQEGKAEAAIKIMRKA V SLIDKA <i>AKGSTLHKNAARRKSRLMRKVQKL</i> |
| M2 | K7L, I23Y | NS IKTLSK L AVLLAQEGKAEAAIK Y MRKA V SLDPNN IKTLSK L AVLLAQEGKAEAAIK Y MRKA V SLDPNN IKTLSK L AVLLAQEGKAEAAIK Y MRKA V SLIDKA <i>AKGSTLHKNAARRKSRLMRKVQKL</i> |
| M4E | K2E, K7L, V9N, I23Y | NS I <u>ETLSKLAN</u> LLAQEGKAEAAIK Y MRKA V SLDPNN I <u>ETLSKLAN</u> LLAQEGKAEAAIK Y MRKA V SLDPNN I <u>ETLSKLAV</u> LLAQEGKAEAAIK Y MRKA V SLIDKA <i>AKGSTLHKNAARRKSRLMRKVQKL</i> |
| M4N | K6N, K7L, V9N, I23Y | NS IKTLS NL ANLLAQEGKAEAAIK Y MRKA V SLDPNN IKTLS NL ANLLAQEGKAEAAIK Y MRKA V SLDPNN IKTLS NL AVLLAQEGKAEAAIK Y MRKA V SLIDKA <i>AKGSTLHKNAARRKSRLMRKVQKL</i> |
| M4RD | K2E, K7R, V9N, I23D | NS I <u>ETLSKRAN</u> LLAQEGKAEAAIK D MRKA V SLDPNN I <u>ETLSKRAN</u> LLAQEGKAEAAIK D MRKA V SLDPNN I <u>ETLSKRAV</u> LLAQEGKAEAAIK D MRKA V SLIDKA <i>AKGSTLHKNAARRKSRLMRKVQKL</i> |
| M5 | K2E, L4W, K7L, V9N, I23Y | NS I <u>ETLSKLAN</u> LLAQEGKAEAAIK Y MRKA V SLDPNN I <u>ETWSKLAN</u> LLAQEGKAEAAIK Y MRKA V SLDPNN I <u>ETWSKLAV</u> LLAQEGKAEAAIK Y MRKA V SLIDKA <i>AKGSTLHKNAARRKSRLMRKVQKL</i> |
| M4NΔC | K6N, K7L, V9N, I23Y | NS IKTLS NL ANLLAQEGKAEAAIK Y MRKA V SLDPNN IKTLS NL ANLLAQEGKAEAAIK Y MRKA V SLDPNN IKTLS NL AVLLAQEGKAEAAIK Y MRKA V SLIDKA <i>AK</i> |

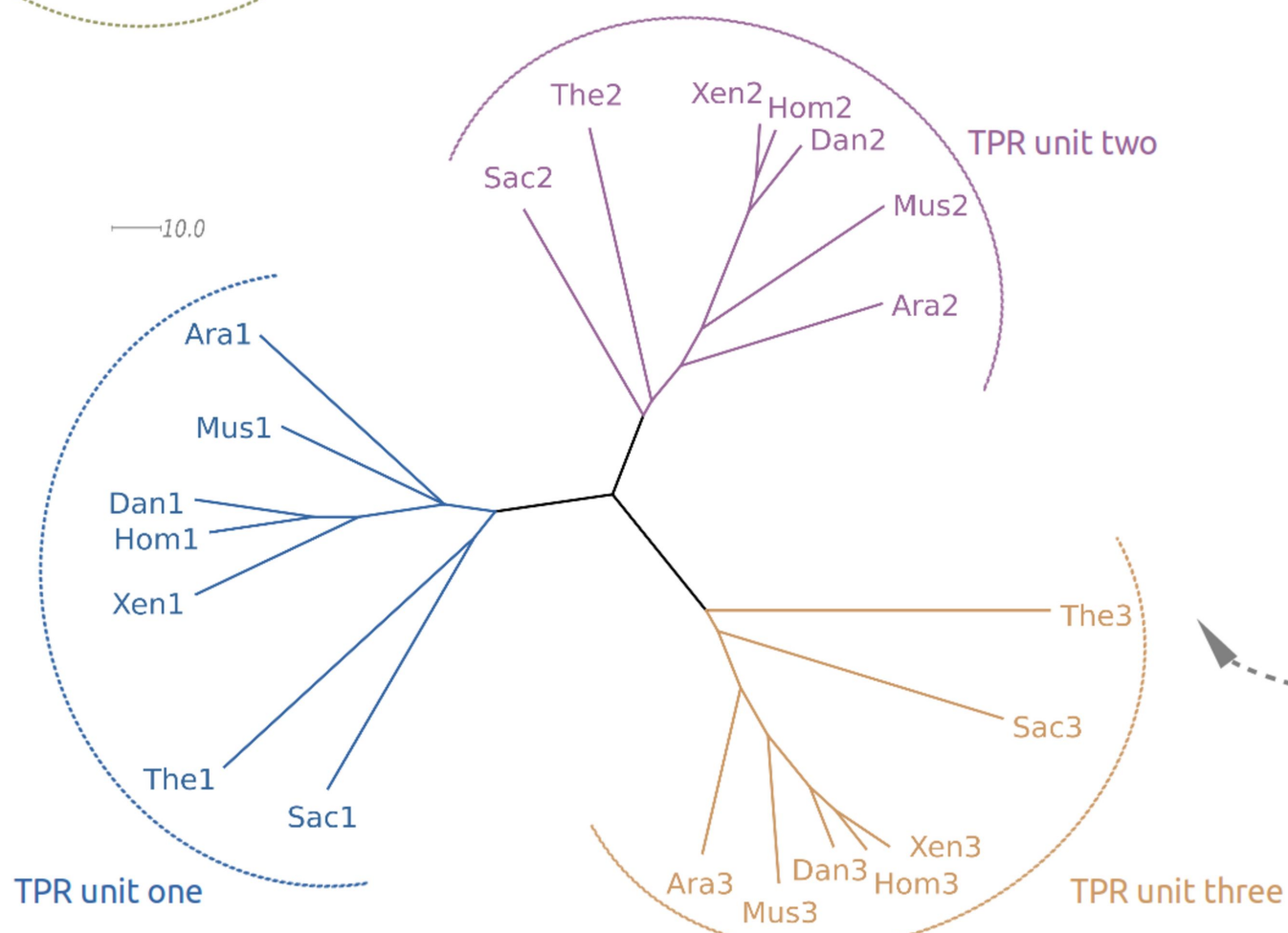
1062

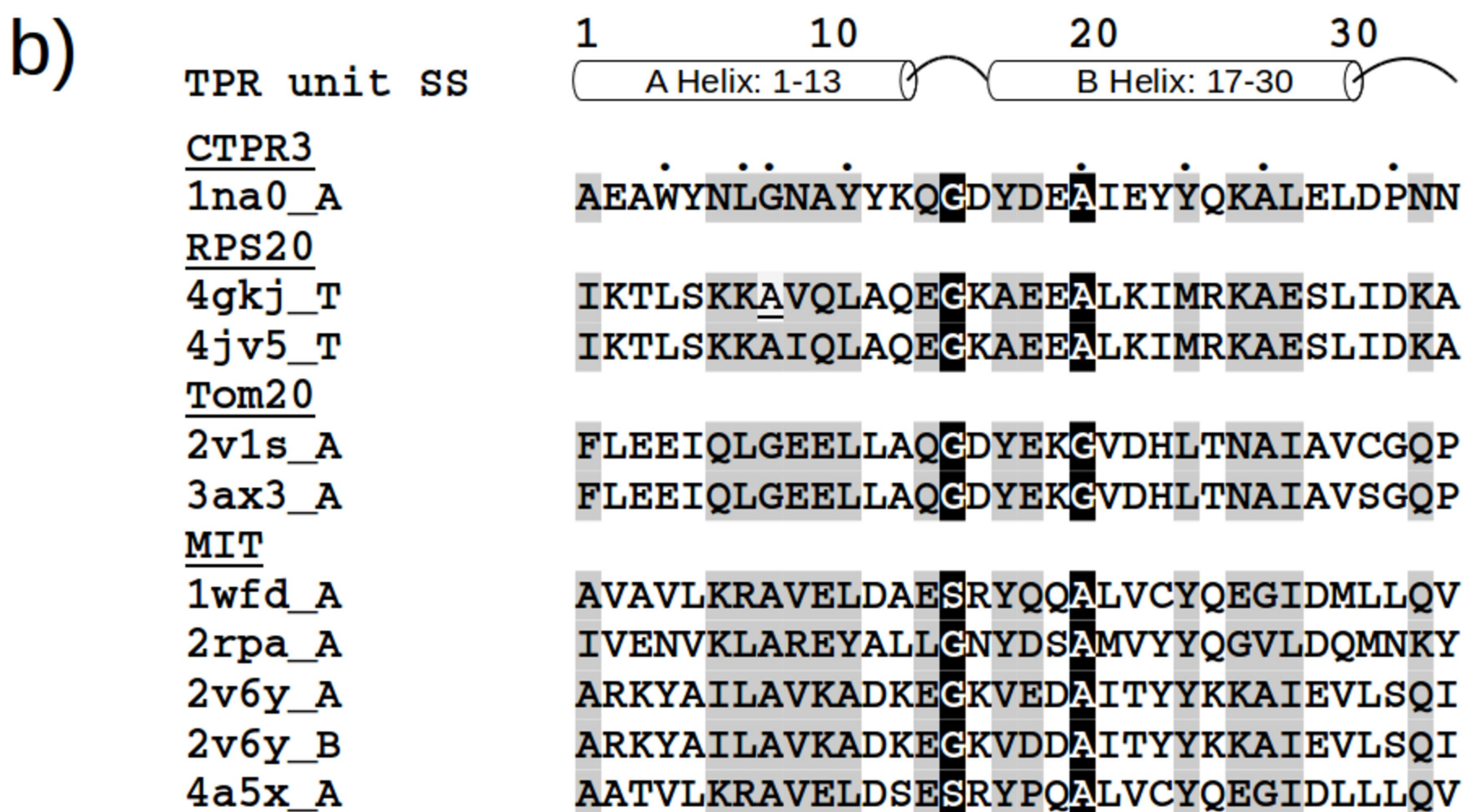
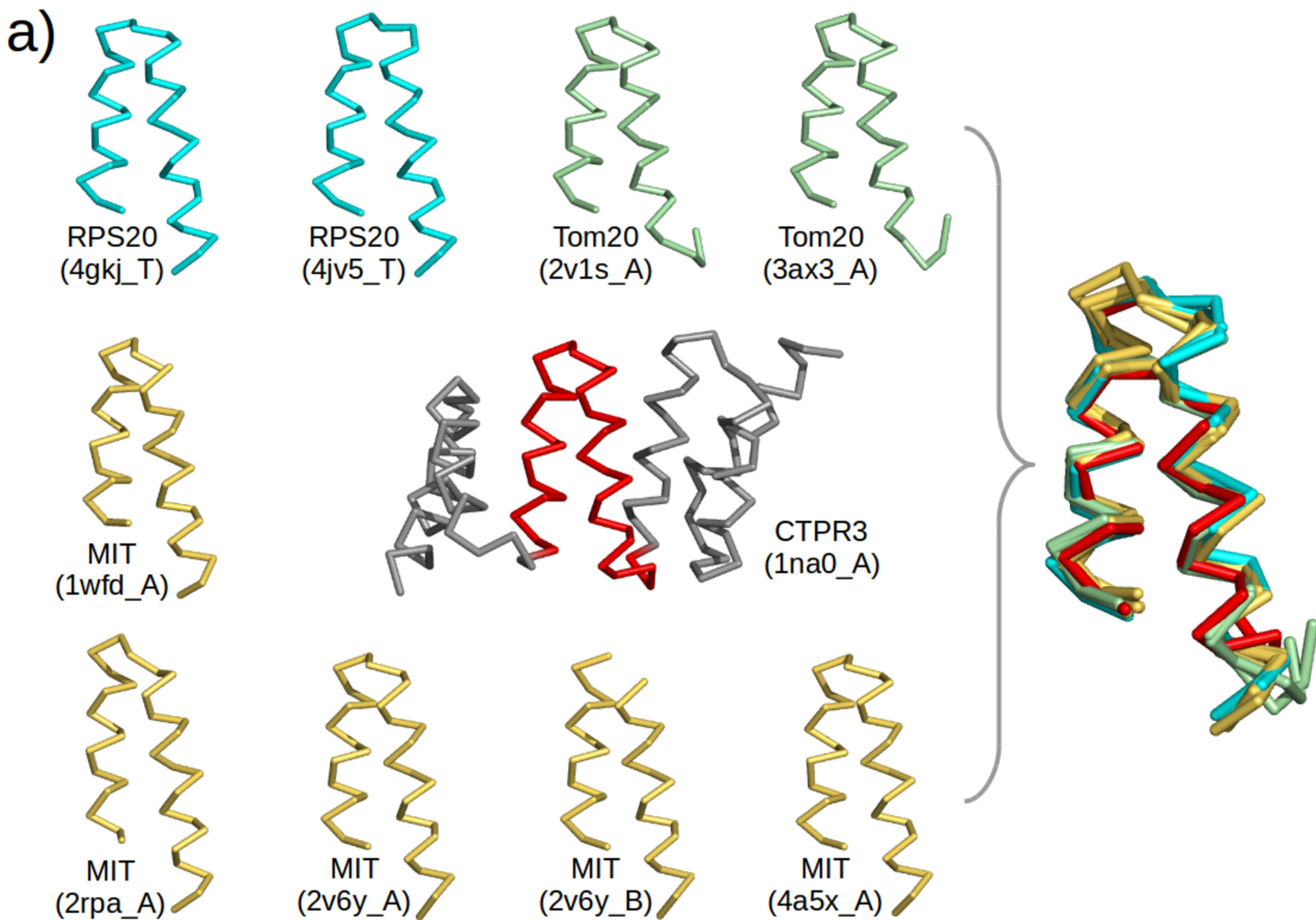
1063 **Table 1.** The primary structures of the six designed proteins using RPS20-hhta tested *in vitro*. Point
1064 mutations introduced into RPS20-hhta are shown in bold and underlined. The C-terminal four
1065 residues in RPS20-hhta were replaced by the consensus loop sequence DPNN in TPRs (underlined).
1066 The stop helix is in gray (italic). M4NΔC is M4N without stop helix.

a)



b)







TPR

Build profile

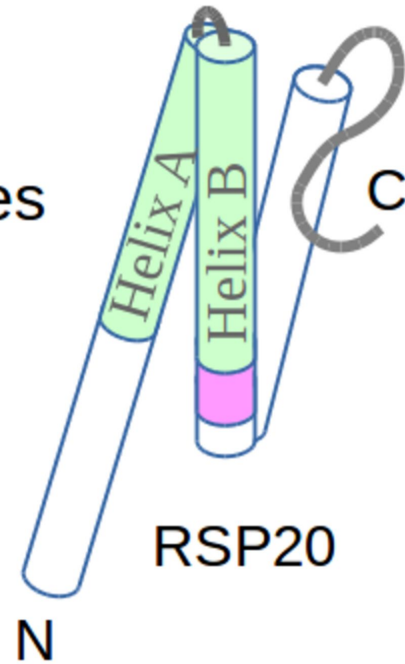


TPR unit

Search matches

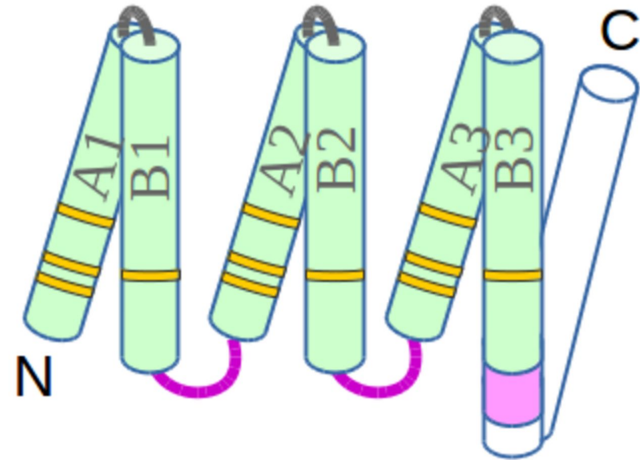


TPRpred

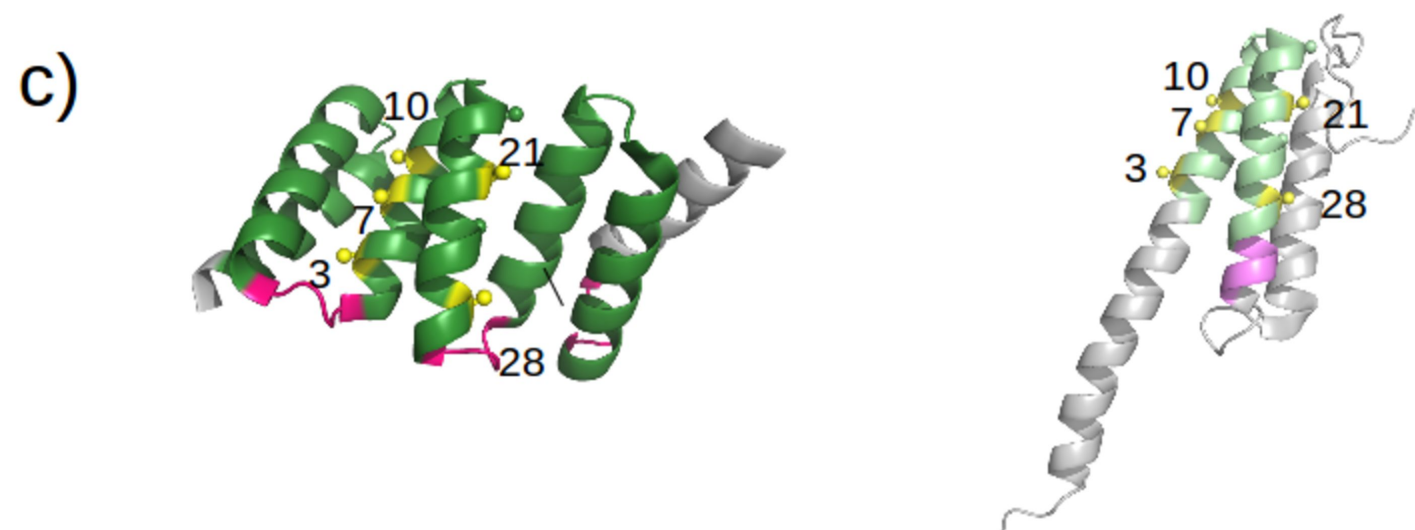
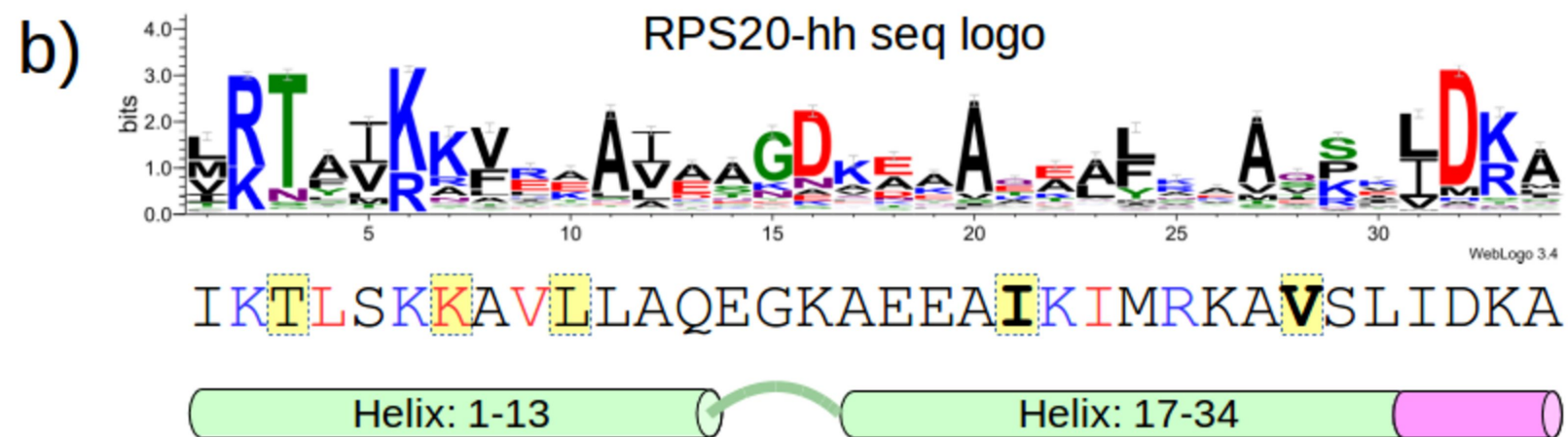
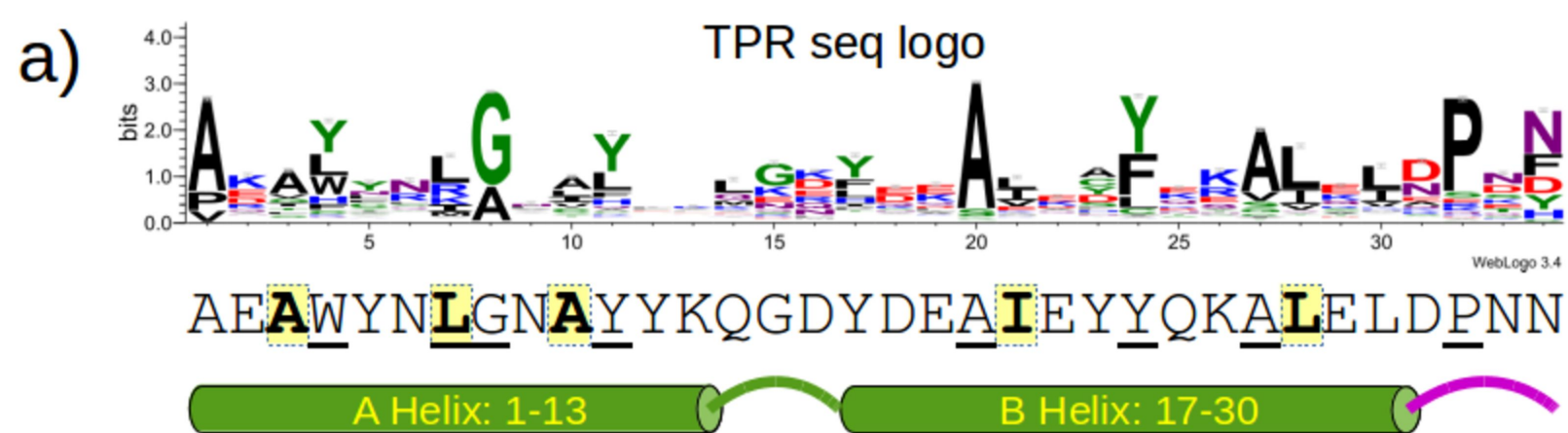


RSP20

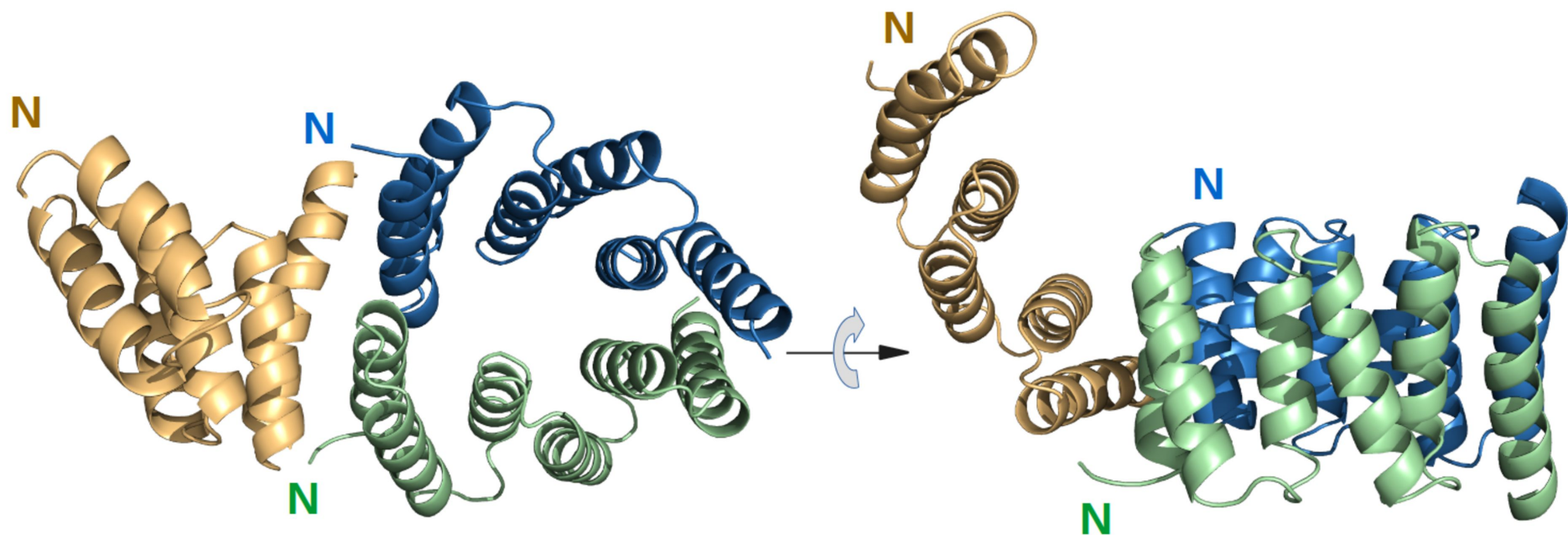
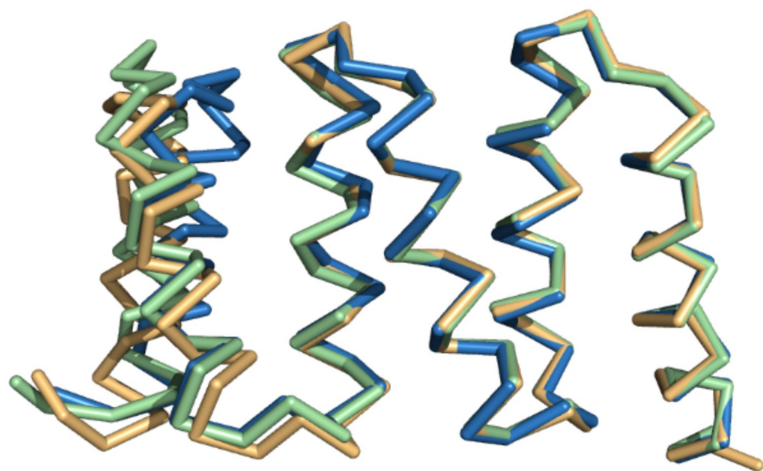
Amplification



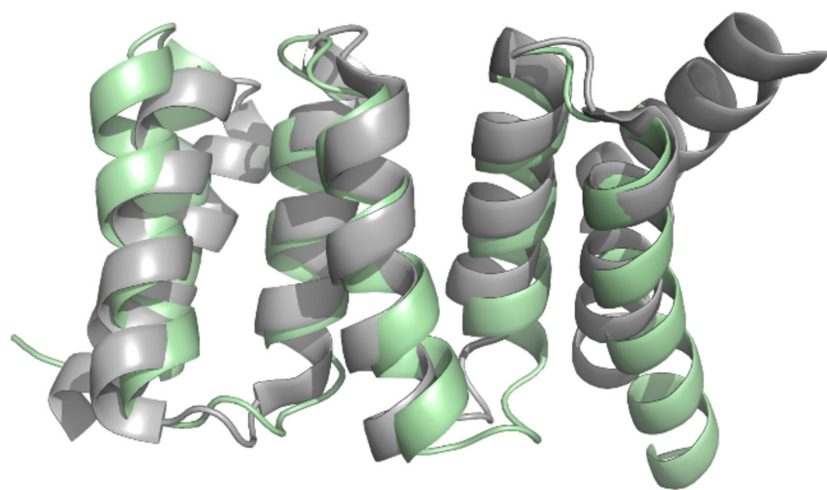
Designed TPR
using RPS20-hh



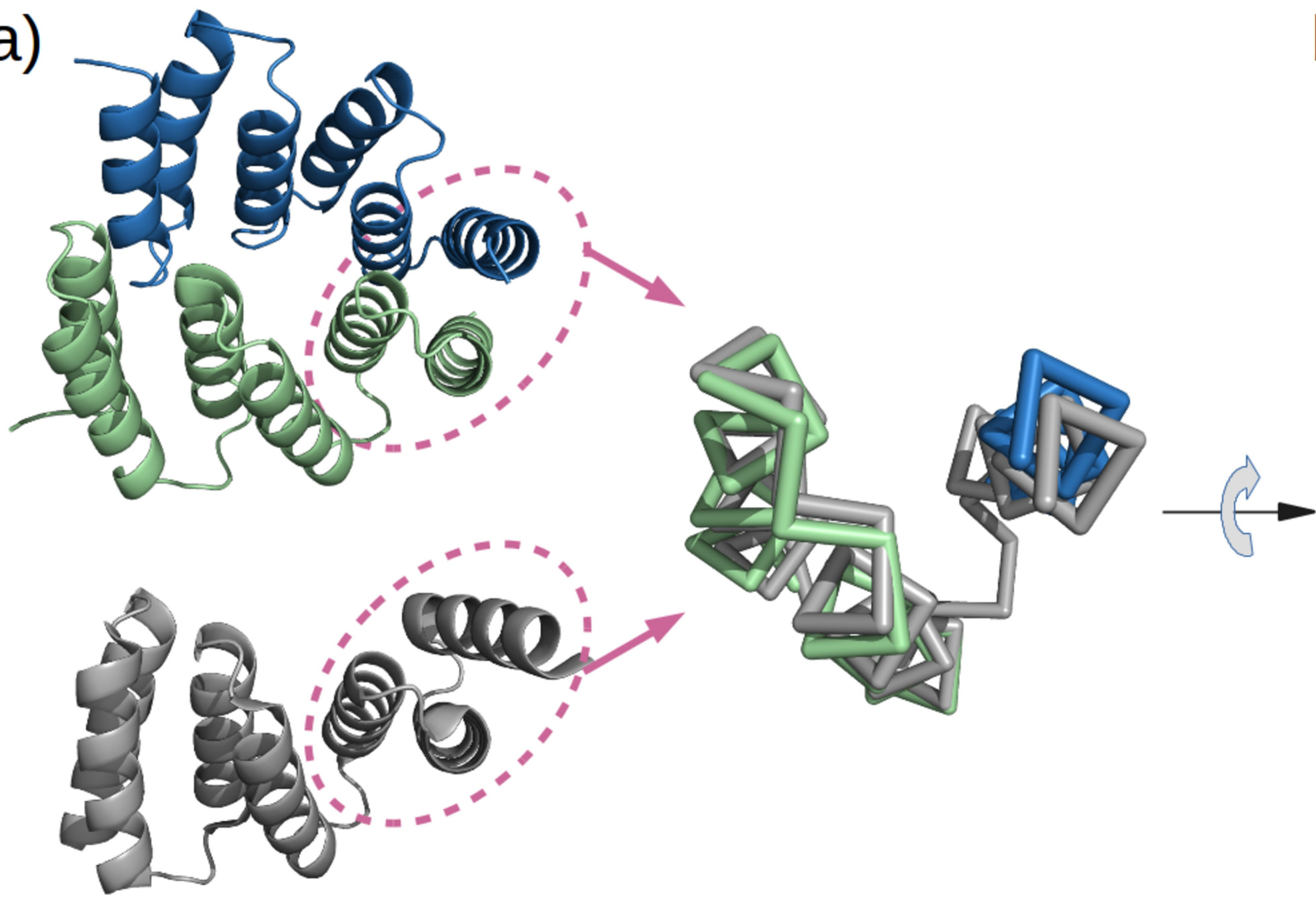
a)

M4N Asymmetric Unit (ASU)b) *Superposition of 3 chains in ASU*

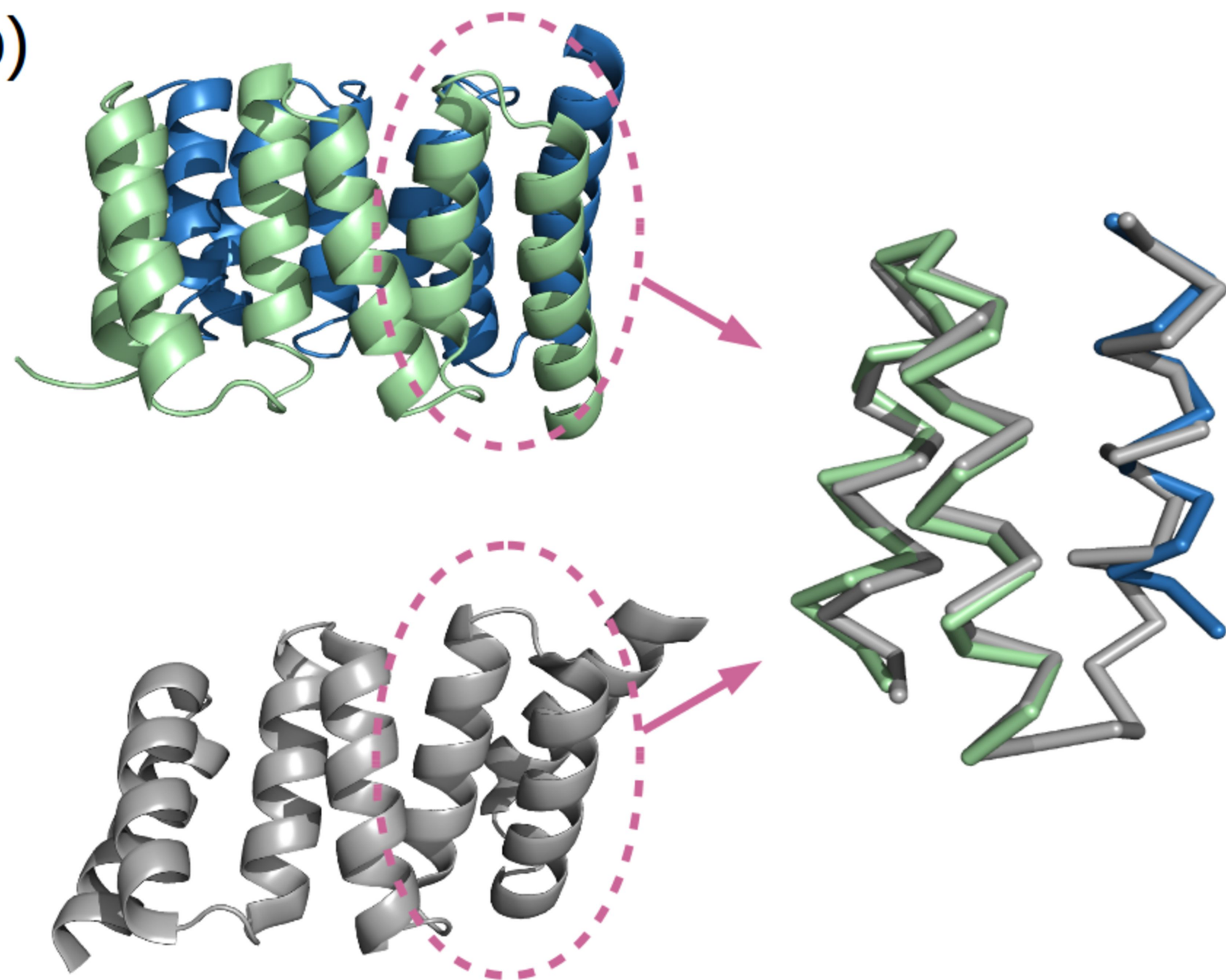
c)

M4N vs. CTPR3

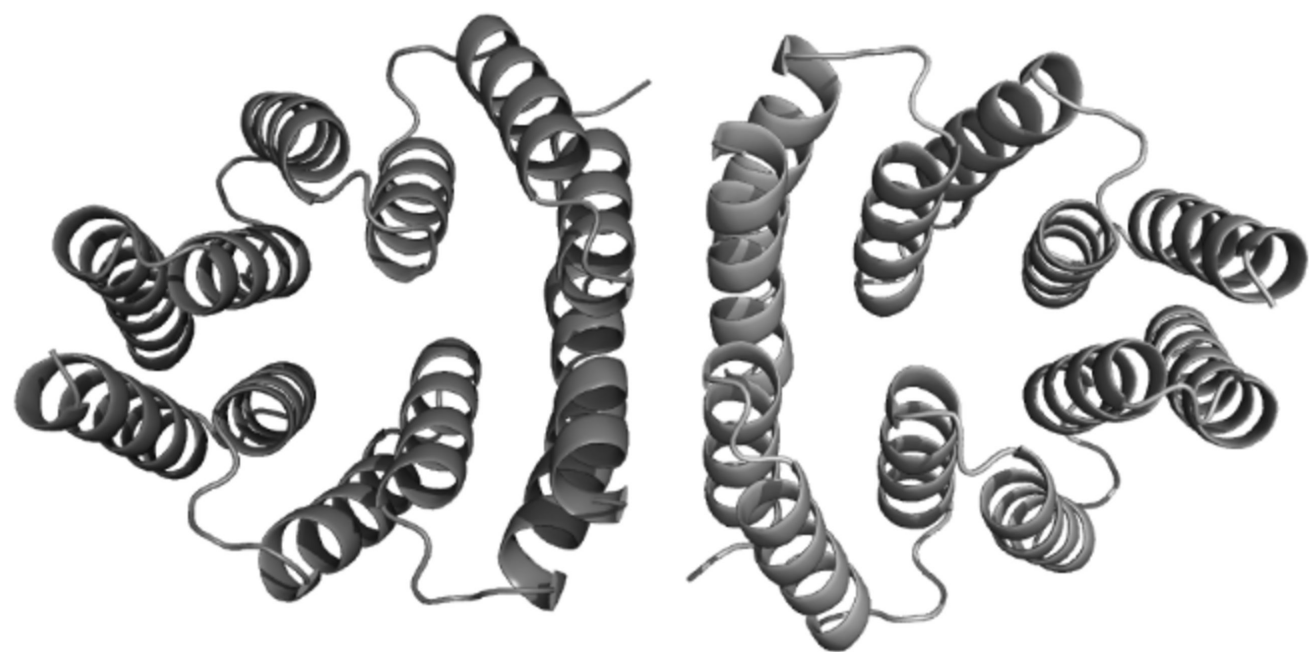
a)



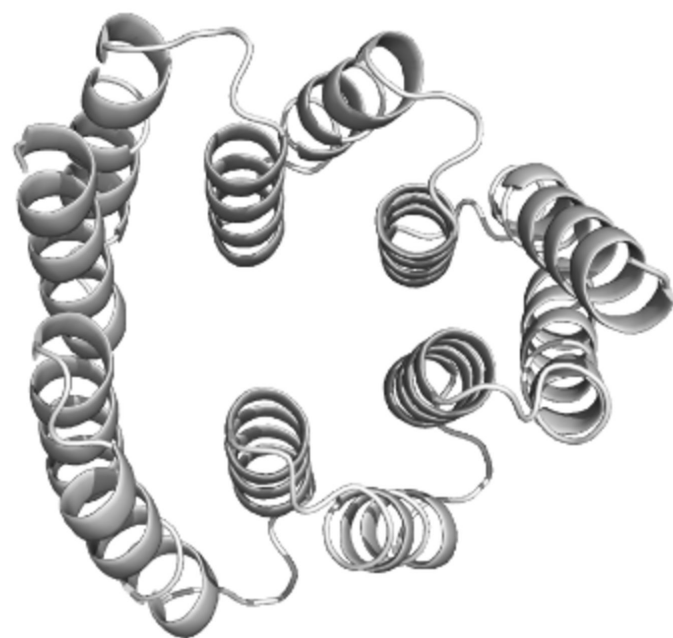
b)



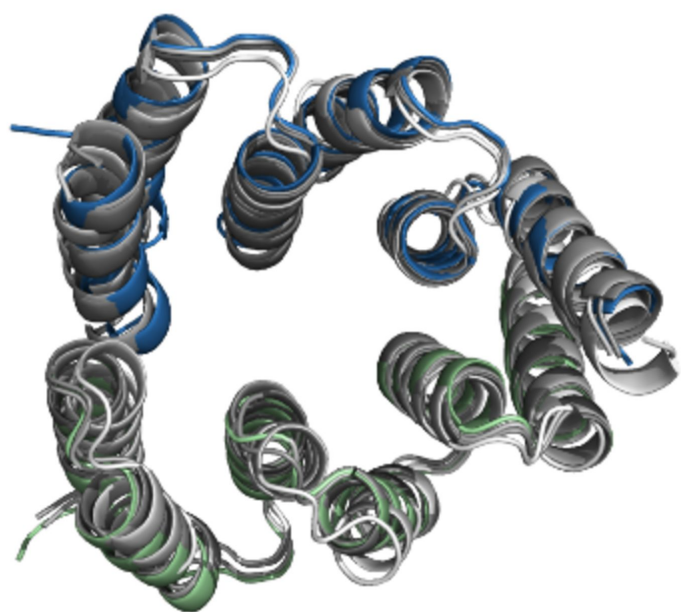
a) *ASU of M4NΔC CF I*



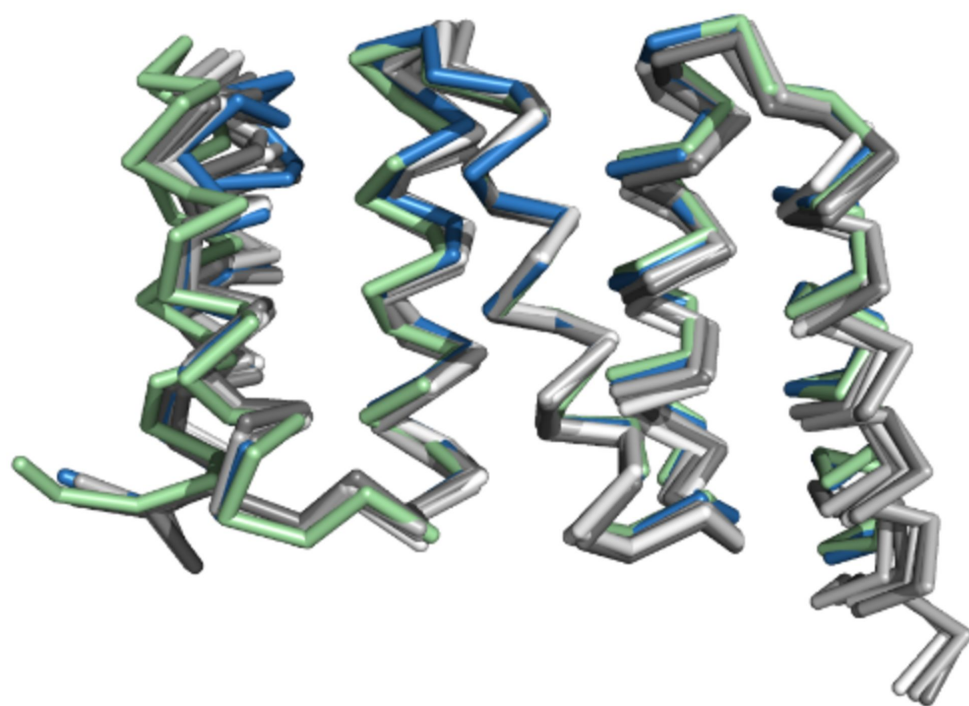
b) *Dimer of M4NΔC CF II*



c) *Superpos. of four dimers*

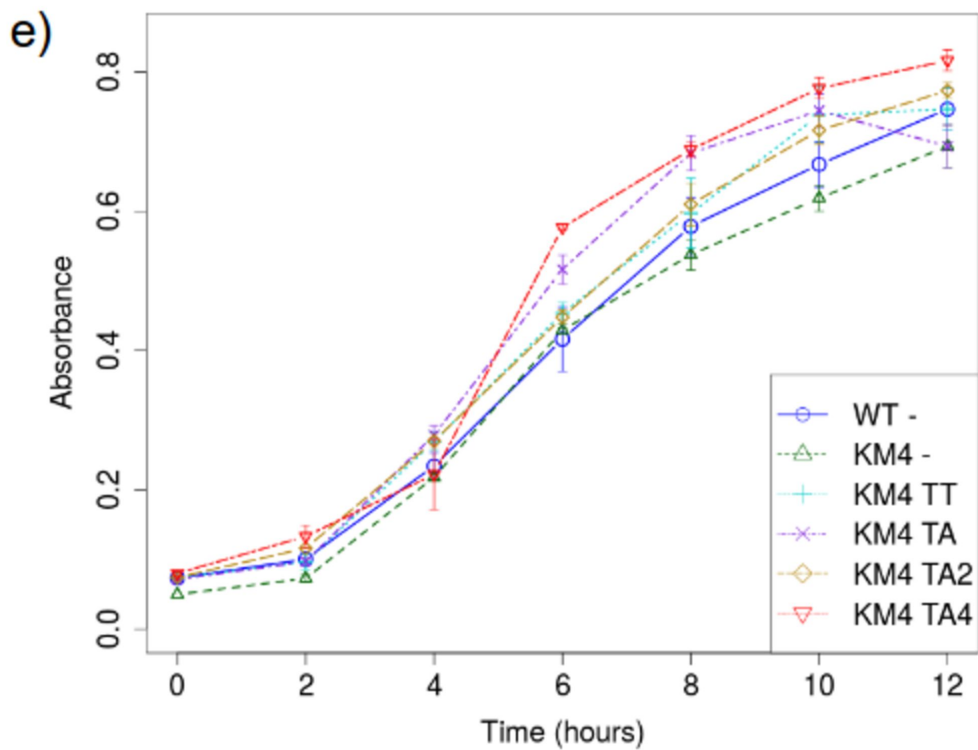
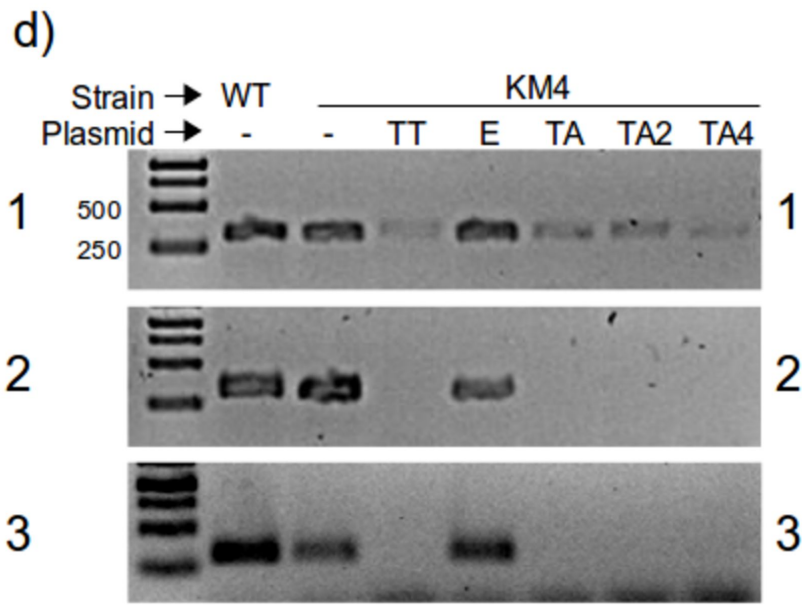
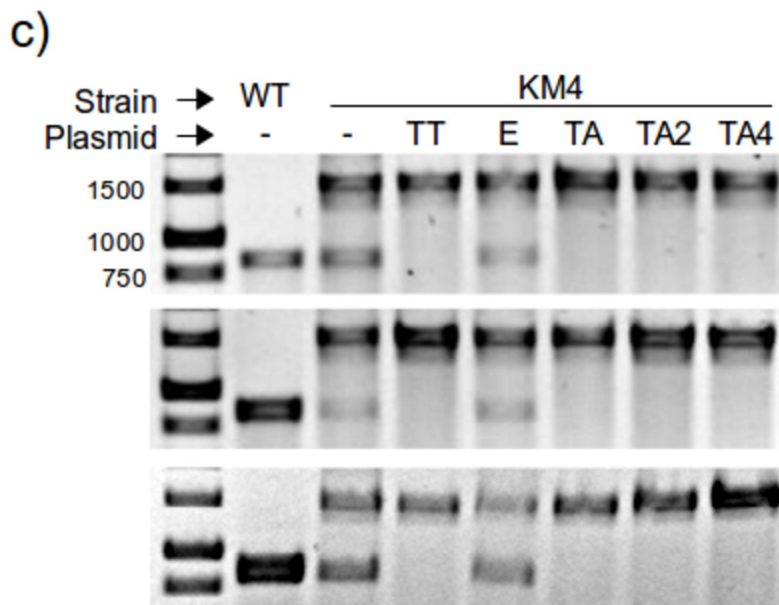
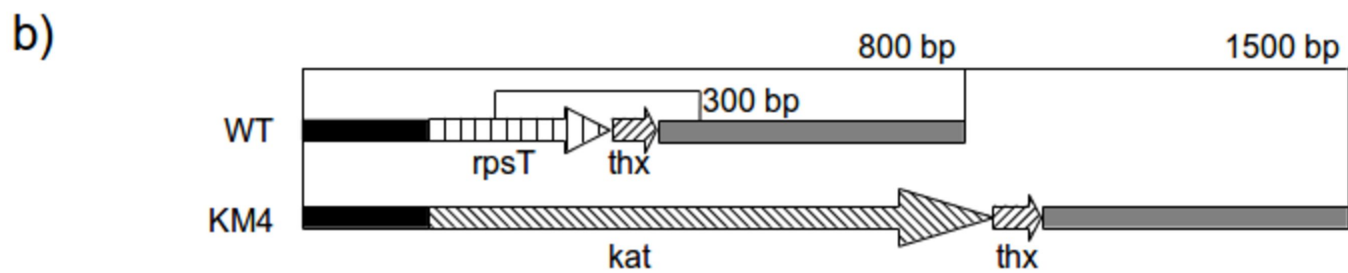


d) *Superpos. of all dimer chains*



a)

| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|--|---|---|----|----|----|----|----|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|--|--|
| | 1 | 5 | 10 | 15 | 20 | 25 | 30 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| <i>T. aquaticus</i> Y51MC23 (gi 218244825) | I | K | T | L | S | K | K | A | V | L | L | A | Q | E | G | K | A | E | E | A | I | K | I | M | R | K | A | V | S | L | I | D | K | A | | |
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| <i>T. thermophilus</i> HB8 (gi 11125386) | I | K | T | L | S | K | K | A | I | Q | L | A | Q | E | G | K | A | E | E | A | L | K | I | M | R | K | A | E | S | L | I | D | K | A | | |

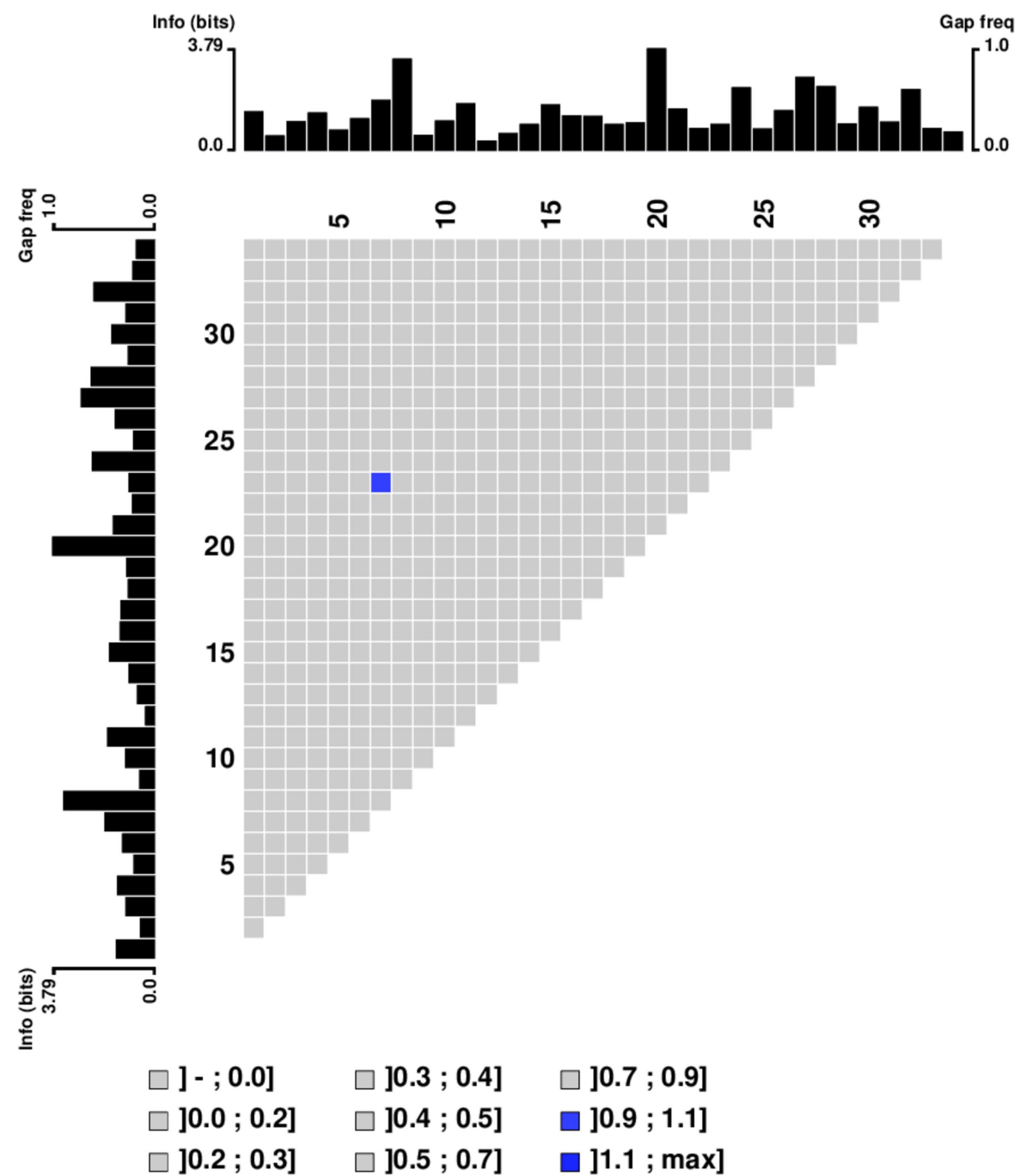


| | | |
|------|--------------------------------|----------------|
| Ara1 | AEEFKSQANEAFKGHKYSSAIDL | YTKAIELNSNN |
| Dan1 | AEKLKEKANDYFKDKDYENAIKYYTEALD | LNPTN |
| Hom1 | AEELKTQANDYFKAKDYENAIKFYSQAIEL | LNPSN |
| Mus1 | ADHYKNQGN DY LKEKDFSKAIE | MYTKAIELNPNN |
| Sac1 | ALERKNEGNV FVKEKHFLKAIE | KYTEAIDL DSTQ |
| The1 | AELKKLEGNRSFSENNFISAI | IEHYSESIRLVEDS |
| Xen1 | AEELKEQANEYFRVKDYDRAVQYYTQA | IGLSPDT |
| | | |

| | | |
|------|----------------------------------|------------------|
| Ara2 | AVYWANRAFAHTKLEEEYGS | AIQDASKAIEVDSRY |
| Dan2 | PIYYSNRSLSYLRTECYGYALAD | ATRALELDKNY |
| Hom2 | AIYYGNRSLAYLRTECYGYALGD | ATRAIELDKKY |
| Mus2 | AIYYANRSLAHLRQESFGFALQD | GISAVKADPTY |
| Sac2 | SIYFSNRAFAHFKVDNFQSALN | DCDEAIKLD PKN |
| The2 | HQYYSNRAICNIKIENYGS | AI SDANVAIQLRPDF |
| Xen2 | AIYYGNRSLAYLRTECYGYALAD | ASRAIQ L DAKY |
| | | |

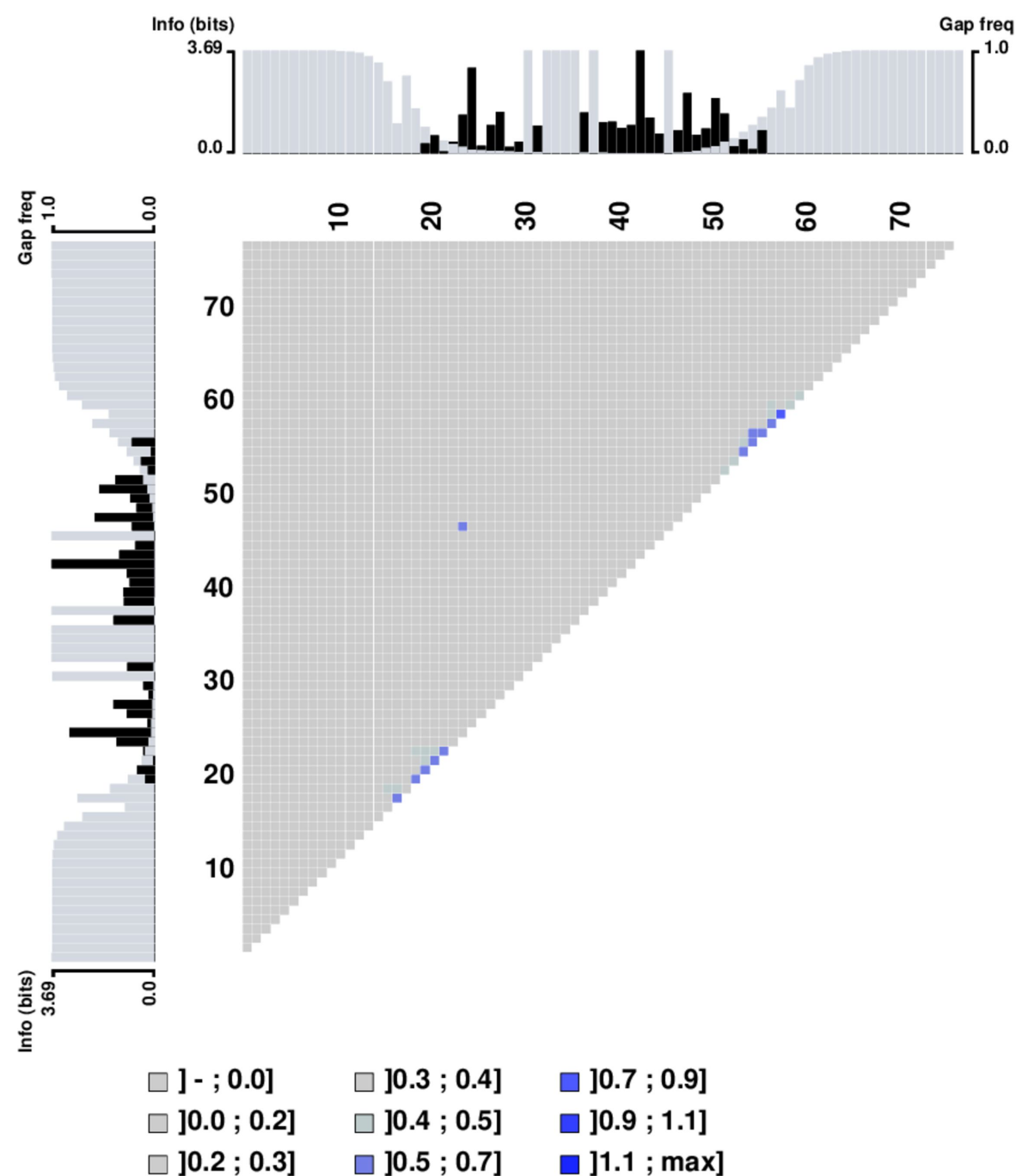
| | | |
|------|-----------------------------------|---------------|
| Ara3 | SKGYYRRGAAYLAMGKFKDAL | KDFQQVKRLSPND |
| Dan3 | LKGYYRRATSNMALGKFKAAL | KDYETVVRVRPND |
| Hom3 | IKGYYRRAASNMALGKFRAAL | RDYETVVKVPHD |
| Mus3 | LKGYYRRAAAHMSLGKFKQAL | SDYEYMSKCRPND |
| Sac3 | IKAYHRRALSCMALLEFKKARK | DLNVLLKAKPND |
| The3 | FKAYYRRGCAYLCLLKFQDAE | TDFLKVLSLCNDP |
| Xen3 | IKGYYRRAASNMALGKLLKAAL | KDYETVVKVRPHD |
| | | |

a)



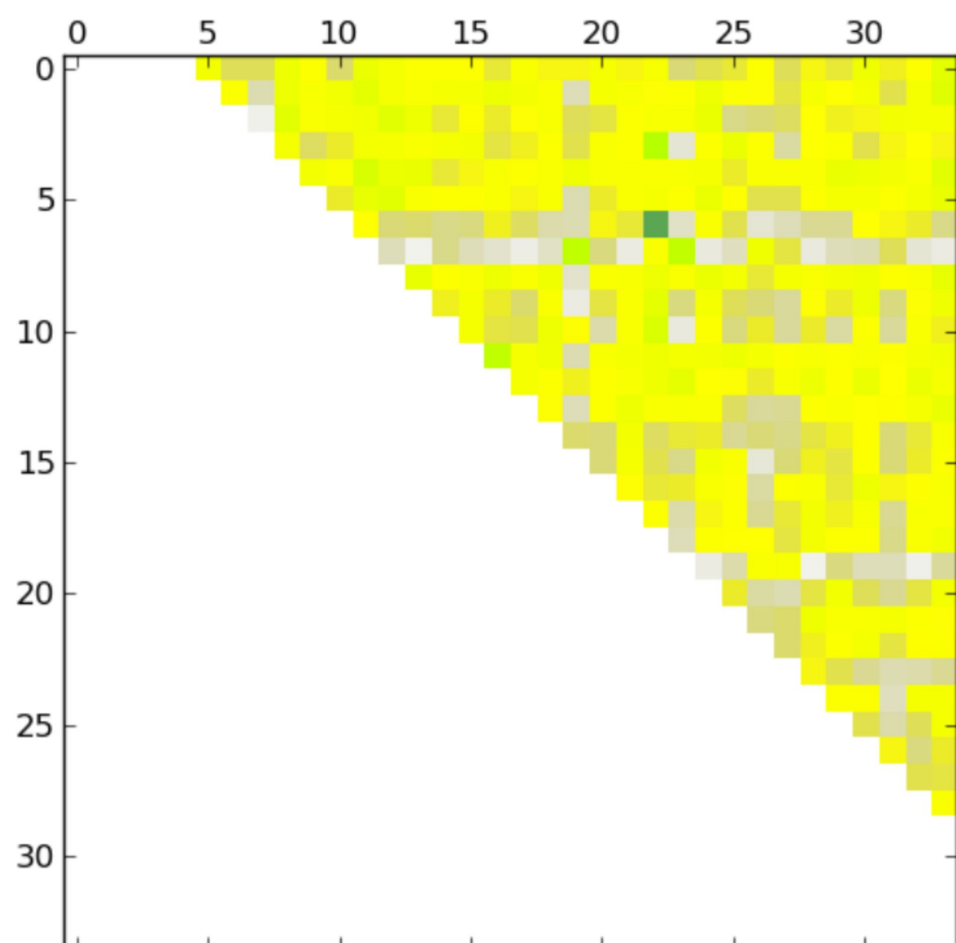
Max: 0.905

b)

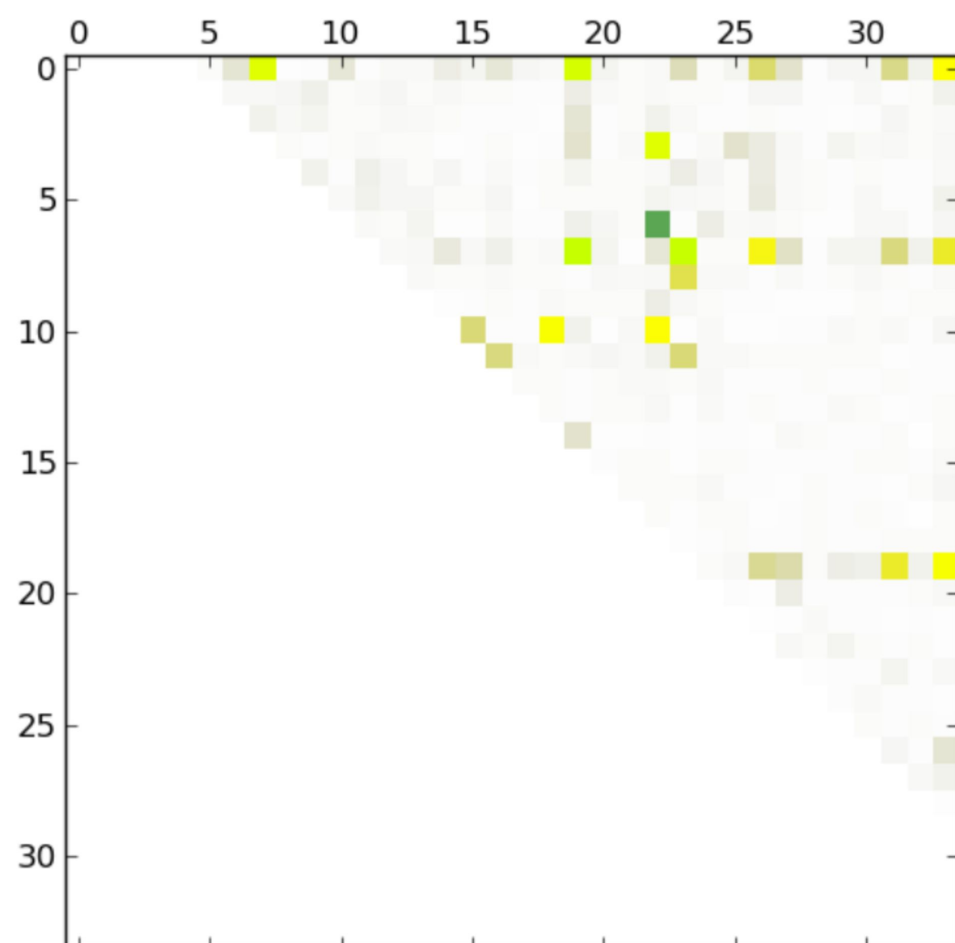


Max: 0.737

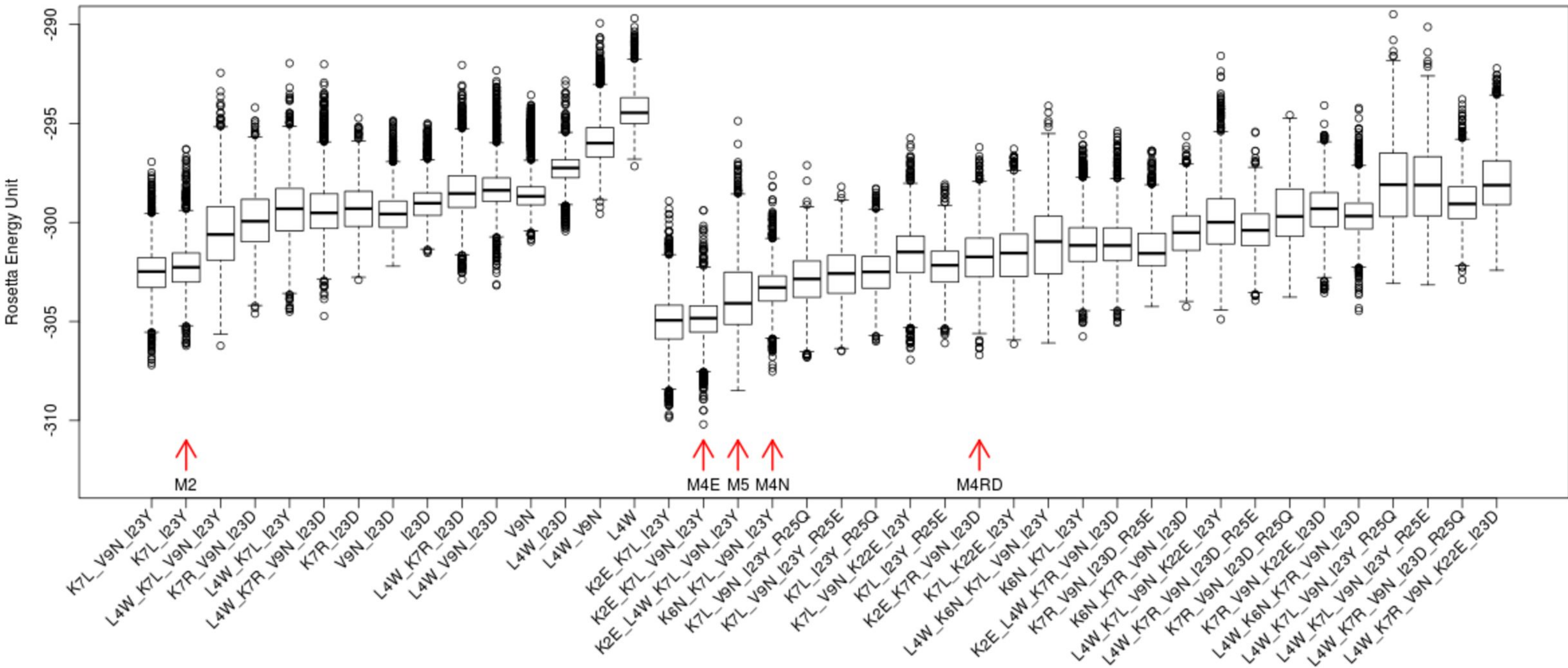
c)



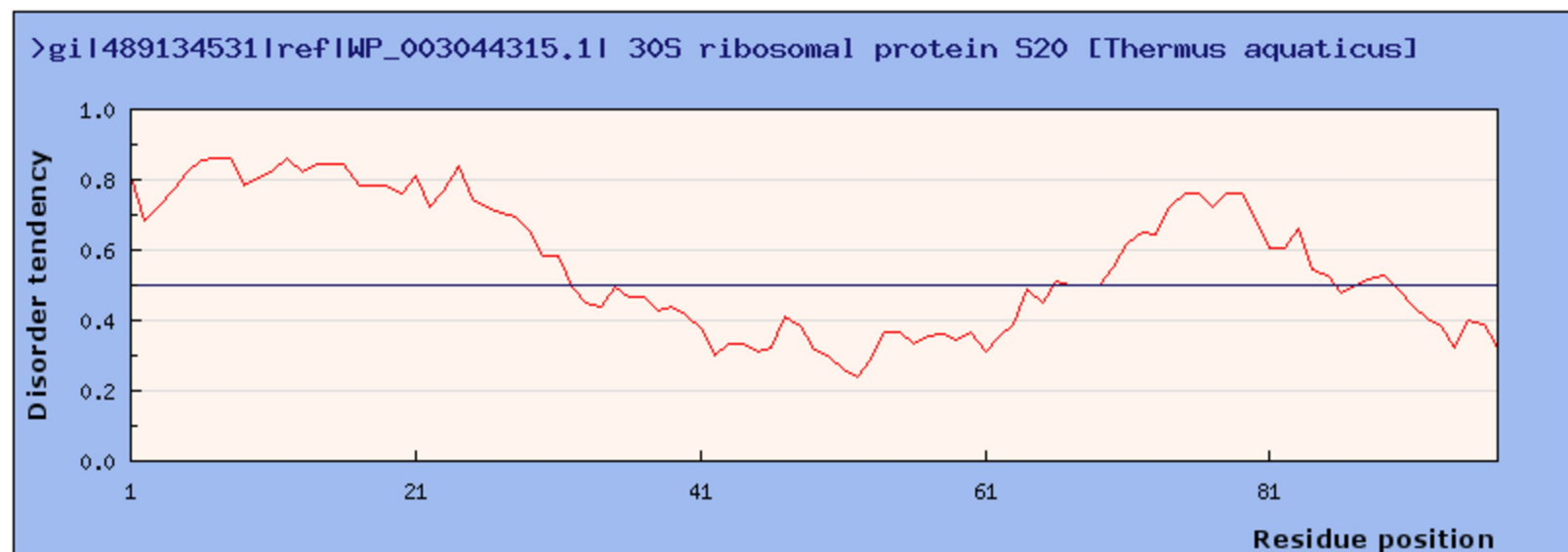
d)



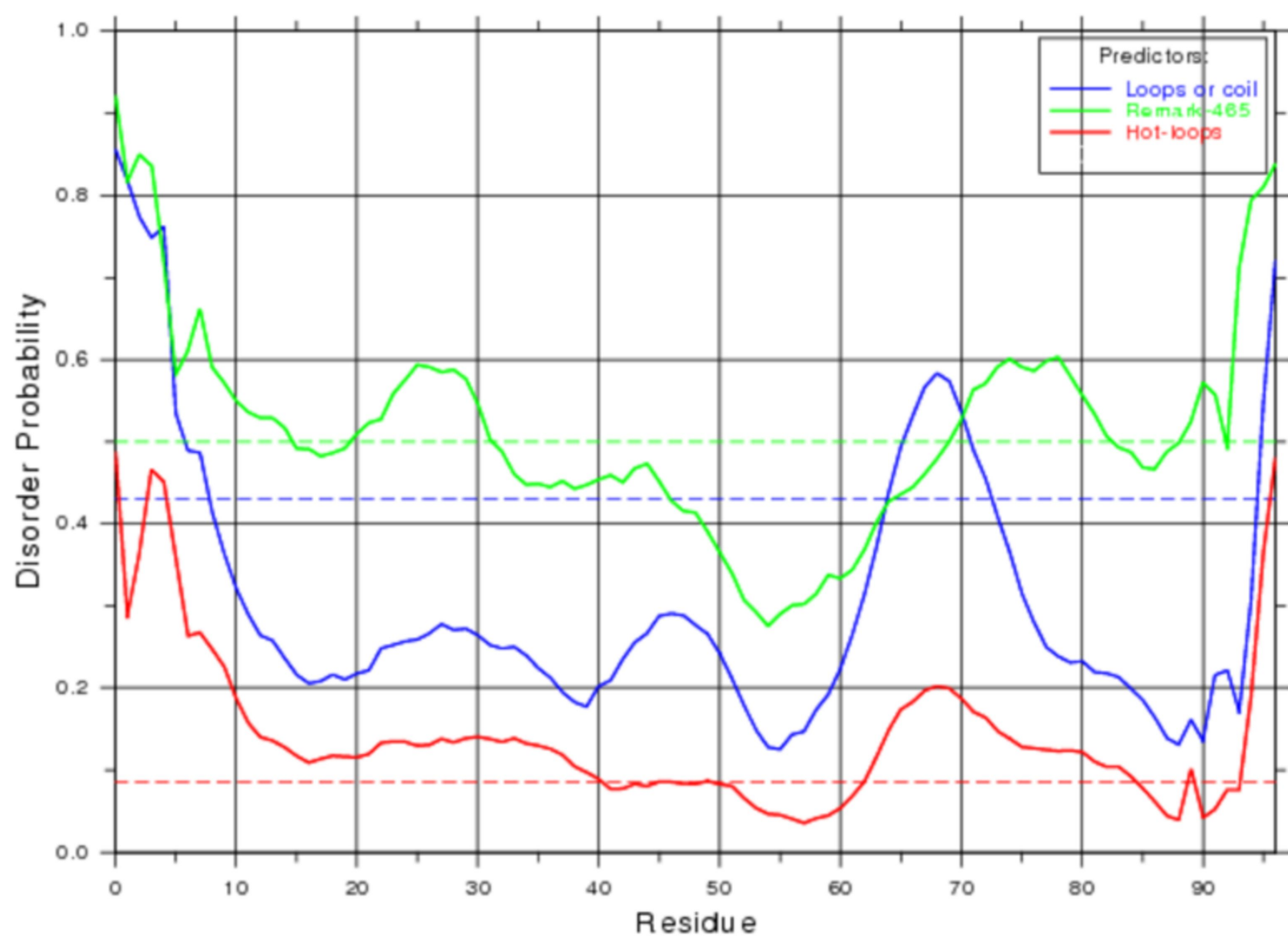
Rosetta energy scores (relax) for different TPR designs using RPS20-hh *T. aquaticus*



a)



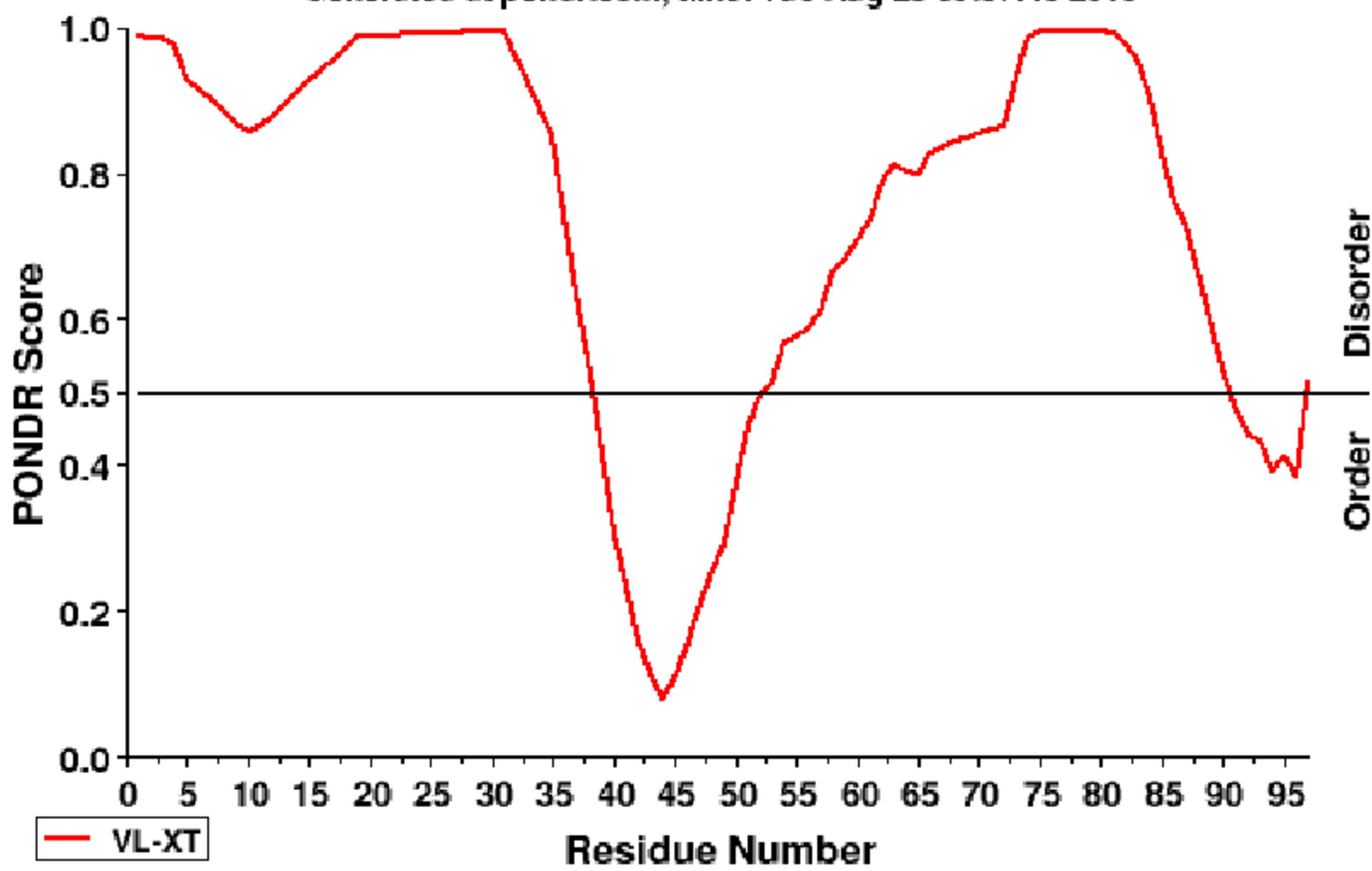
b)



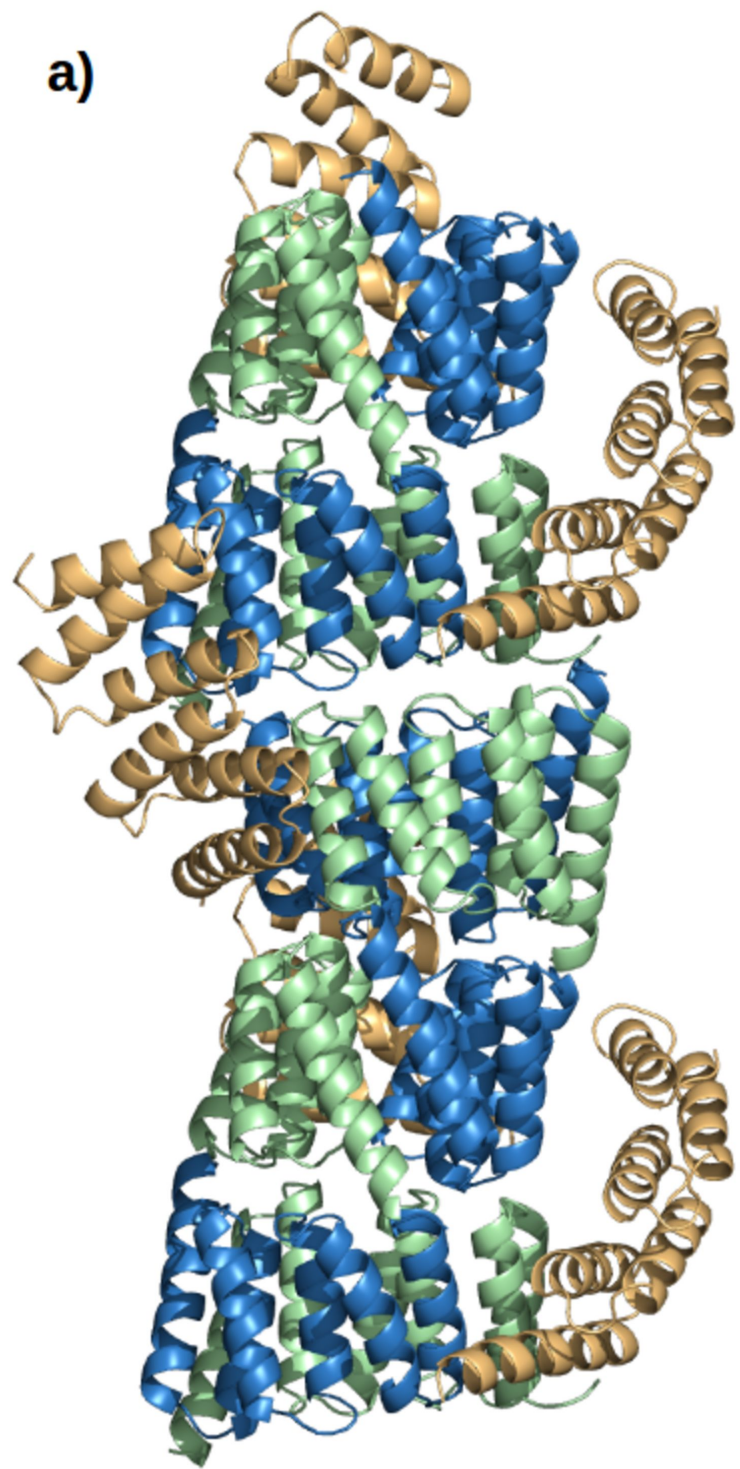
gi|489134531|ref|WP_003044315.1| 30S

Generated at pondr.com, time: Tue Aug 23 09:37:40 2016

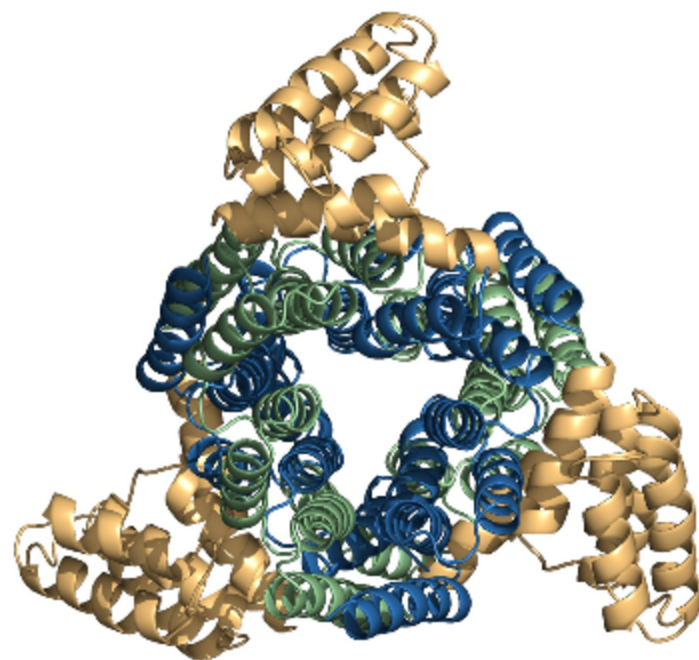
c)



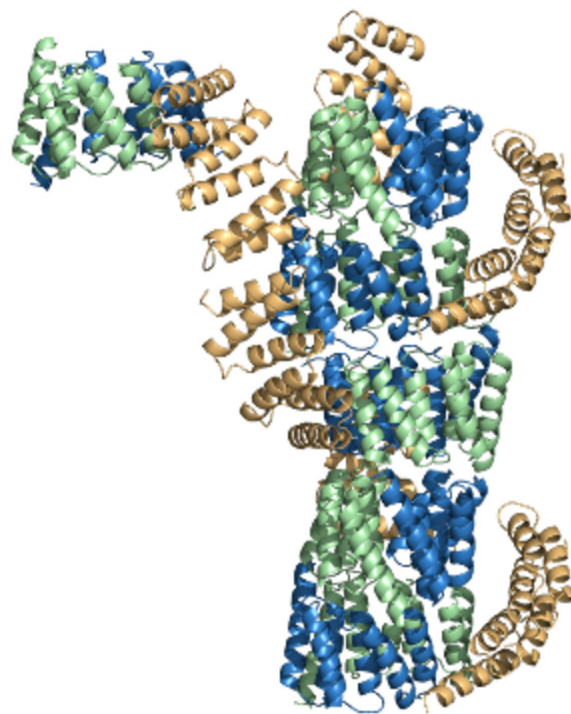
a)

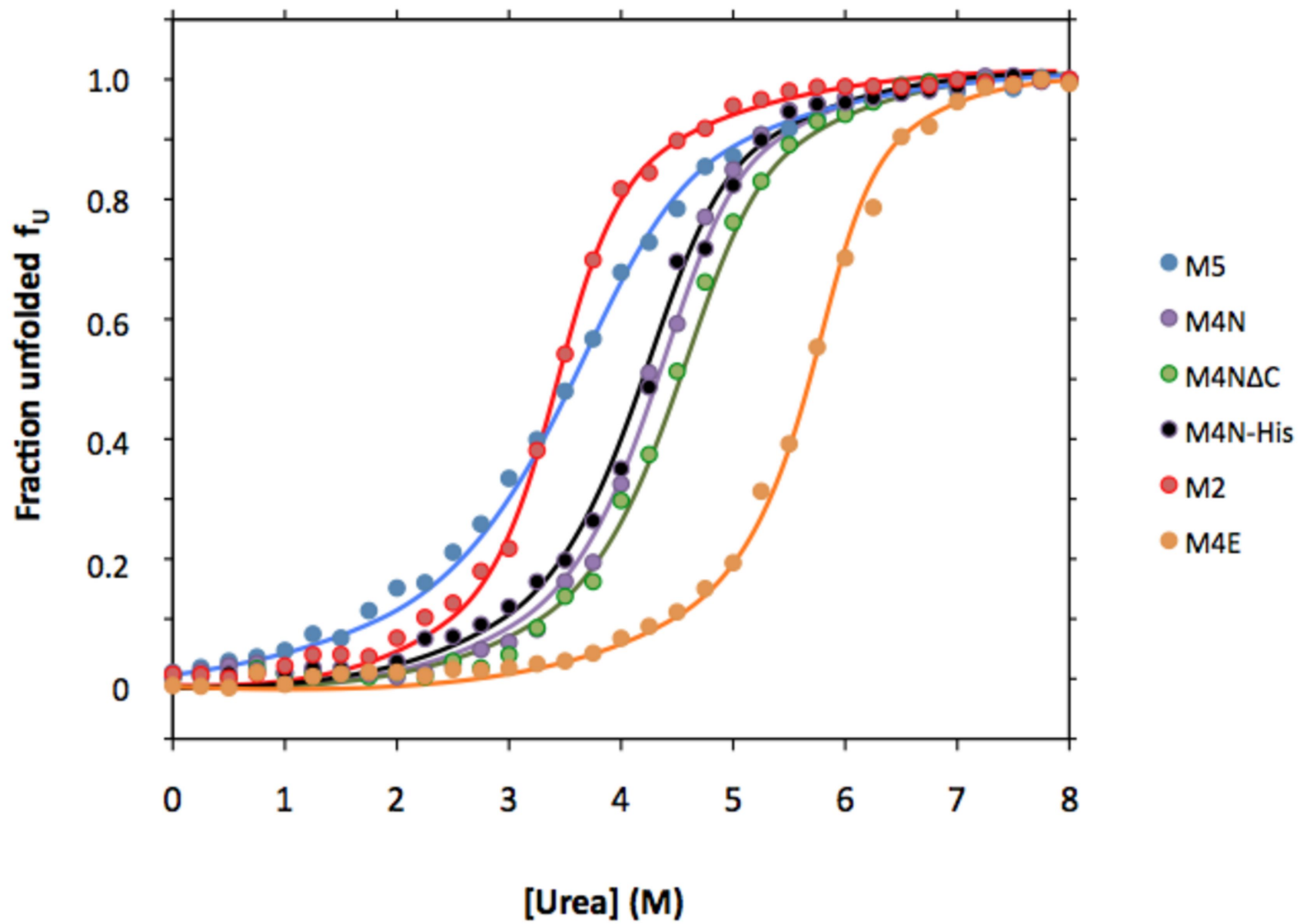


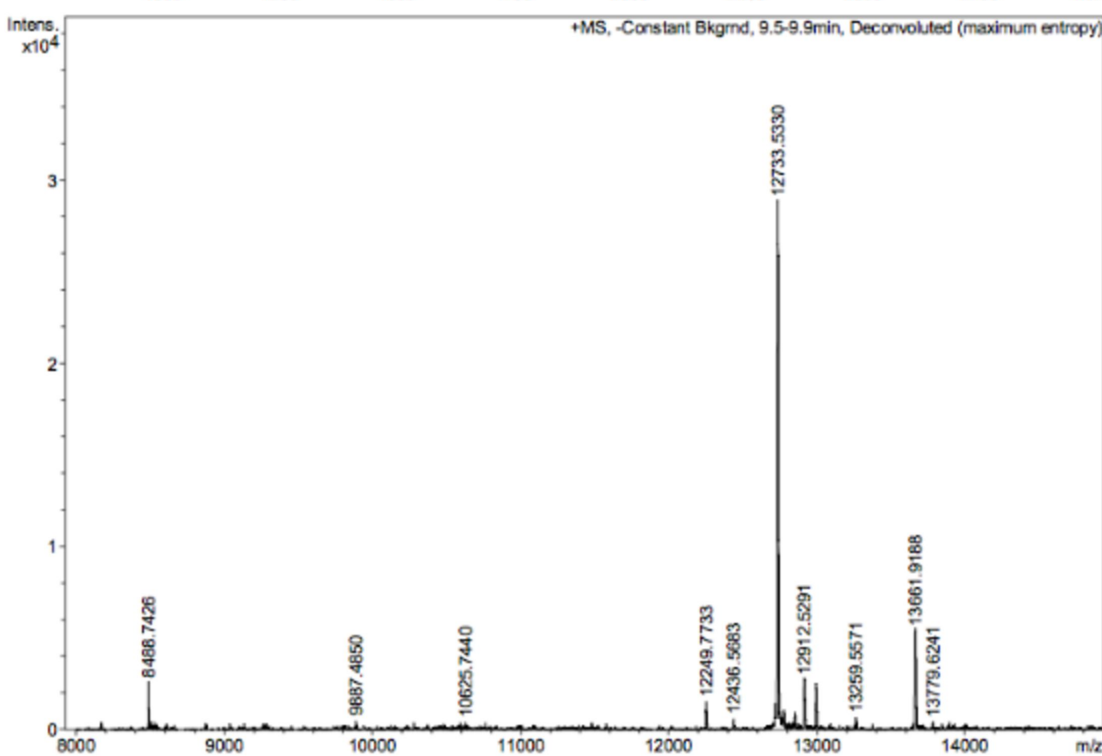
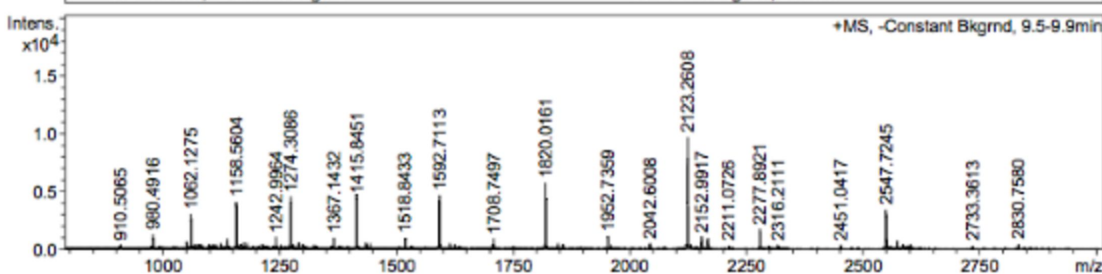
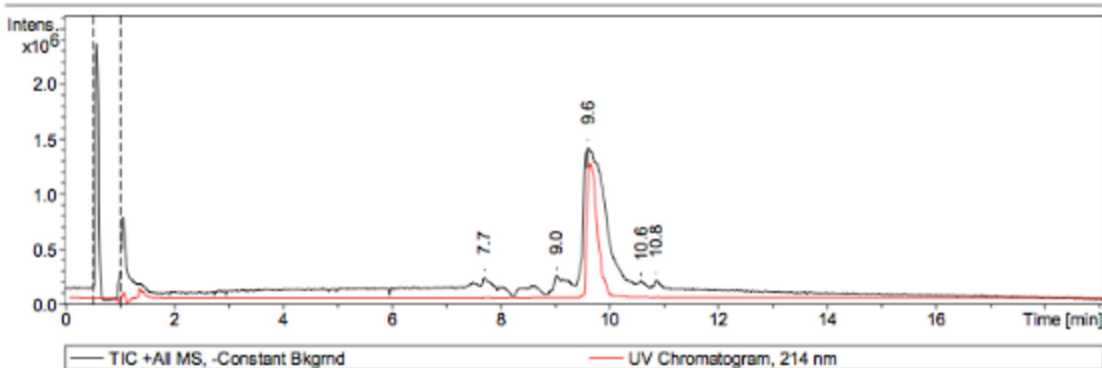
b)



c)







>M4N

MGNS

IKTLSNLANLLAQEGKAE EAIK YMRKAVSLDPNN

IKTLSNLANLLAQEGKAE EAIK YMRKAVSLDPNN

IKTLSNLAVLLAQEGKAE EAIK YMRKAVSLIDKA

AKGSTLHKNA AARRKSRLMRKVQKL

