



University of Dundee

Ankyrin repeats in context with human population variation

Utgés, Javier S.; Tsenkov, Maxim I.; Dietrich, Noah J. M.; MacGowan, Stuart A.; Barton, Geoffrey J.

DOI:

[10.1101/2021.05.28.445974](https://doi.org/10.1101/2021.05.28.445974)

Publication date:

2021

Document Version

Other version

[Link to publication in Discovery Research Portal](#)

Citation for published version (APA):

Utgés, J. S., Tsenkov, M. I., Dietrich, N. J. M., MacGowan, S. A., & Barton, G. J. (2021). *Ankyrin repeats in context with human population variation*. BioRxiv. <https://doi.org/10.1101/2021.05.28.445974>

General rights

Copyright and moral rights for the publications made accessible in Discovery Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from Discovery Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain.
- You may freely distribute the URL identifying the publication in the public portal.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

1 **Ankyrin repeats in context with** 2 **human population variation**

3 **Javier S. Utgés^{1,2}, Maxim I. Tsenkov¹, Noah J. M. Dietrich¹,**

4 **Stuart A. MacGowan¹ and Geoffrey J. Barton^{1*}**

5 ¹Division of Computational Biology, School of Life Sciences, University
6 of Dundee, Scotland, UK

7 ²Universitat Pompeu Fabra (UPF), Barcelona, Spain

8
9 *Correspondence to: gjbarton@dundee.ac.uk

10

11 **Abstract**

12 Ankyrin protein repeats bind to a wide range of substrates and are one of the most common protein
13 motifs in nature. Here, we collate a high-quality alignment of 7,407 ankyrin repeats and examine
14 for the first time, the distribution of human population variants from large-scale sequencing of
15 healthy individuals across this family. Population variants are not randomly distributed across the
16 genome but are constrained by gene essentiality and function. Accordingly, we interpret the
17 population variants in context with evolutionary constraint and structural features including
18 secondary structure, accessibility and protein-protein interactions across 383 three-dimensional
19 structures of ankyrin repeats. We find five positions that are highly conserved across homologs
20 and also depleted in missense variants within the human population. These positions are
21 significantly enriched in intra-domain contacts and so likely to be key for repeat packing. In
22 contrast, a group of evolutionarily divergent positions are found to be depleted in missense
23 variants in human but significantly enriched in protein-protein interactions. Our analysis also
24 suggests the domain has three, not two surfaces, each with different patterns of enrichment in
25 protein-substrate interactions and missense variants. Our findings will be of interest to those
26 studying or engineering ankyrin-repeat containing proteins as well as those interpreting the
27 significance of disease variants.

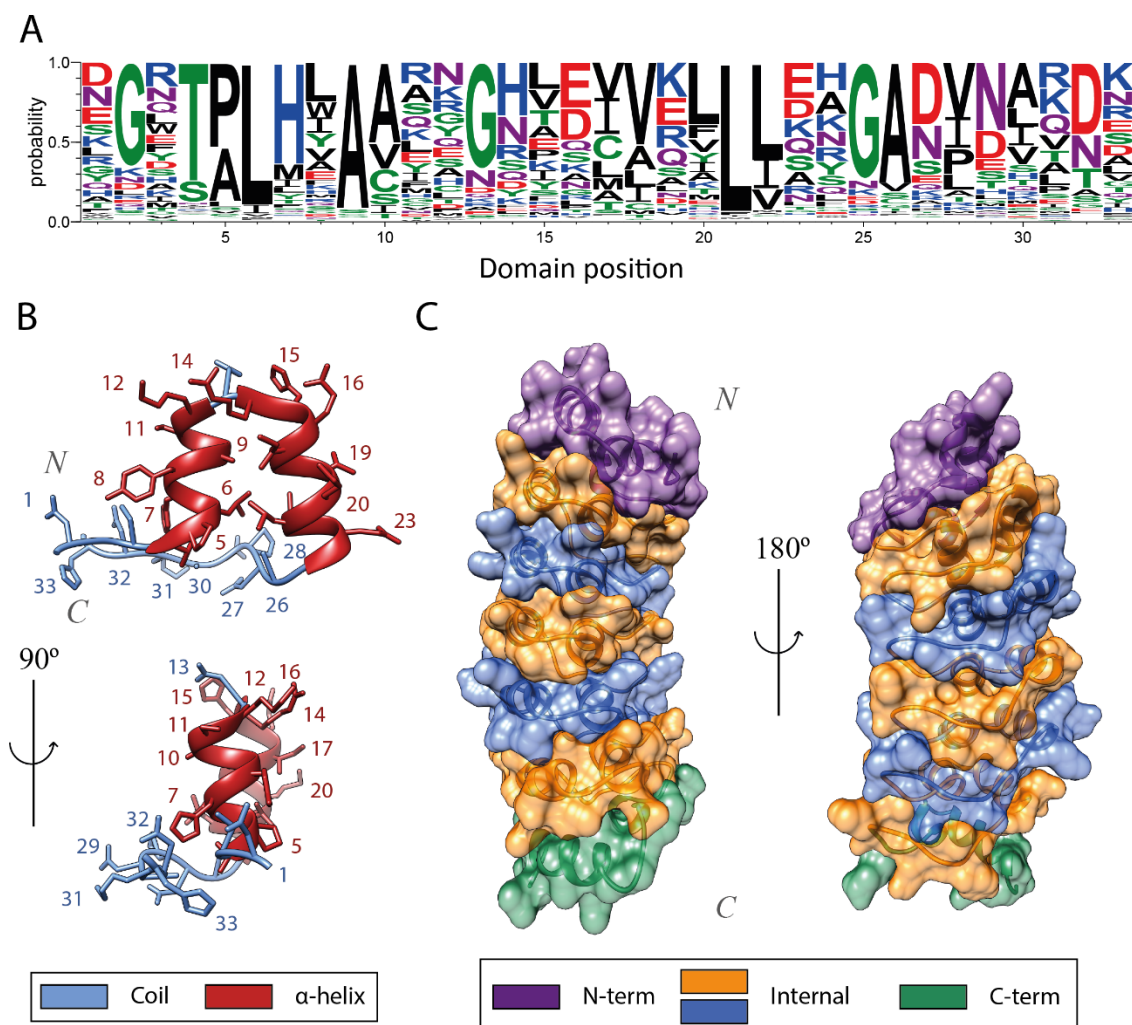
28 **Author Summary**

29 Comparison of variation at each position of the amino acid sequence for a protein across different
30 species is a powerful way to identify parts of the protein that are important for its structure and
31 function. Large-scale DNA sequencing of healthy people has recently made it possible to study
32 normal genetic variation within just one species. Our work combines information on genetic
33 differences between over 100,000 people with in-depth analysis of all available three-dimensional
34 structures for Ankyrin repeats which are a widespread family of binding proteins formed by units
35 with similar amino acid sequence that are found in tandem. Our combined analysis identifies sites
36 critical for ankyrin stability as well as the positions most important for substrate interactions and
37 hence function. Although focused only on the Ankyrins, the principles developed in our work are
38 general and can be applied to any protein family.

39 **Introduction**

40 The ankyrin repeat motif (ANK) is one of the most commonly observed protein motifs in nature,
41 with proteins containing this motif found in practically all phyla (1). Ankyrin repeats (AR) are
42 specialised in protein binding and take part in many processes including transcription initiation,
43 cell cycle regulation and cell signalling (2). ANK is 33 residues long (Fig 1A) and has a helix-
44 turn-helix conformation, with short loops at the N and C termini (Fig 1B). The last and first two
45 residues of adjacent repeats form a β -turn. These β -turns project outward at an angle of $\approx 90^\circ$ to
46 the antiparallel α -helices, yielding the characteristic L-shaped cross section of ankyrin repeats.
47 Ankyrin repeats are usually found in tandem with two or more forming an ankyrin repeat domain
48 (ARD). The stacking of repeats is mediated by the conserved hydrophobic faces of the helices as
49 well as the complementarity of repeat surfaces that assemble to form an extended helical bundle
50 (Fig 1C) (3). Less conserved positions in the motif, i.e., positions that present residues with many
51 different physicochemical properties, are located on the surface, and are likely to interact with
52 ligands. More conserved positions, on the other hand, tend to be buried in the structure and are
53 responsible for the correct packing of the domain. They do this by forming both intra- and inter-
54 repeat contacts, such as hydrophobic and hydrogen bond interactions (4).

55



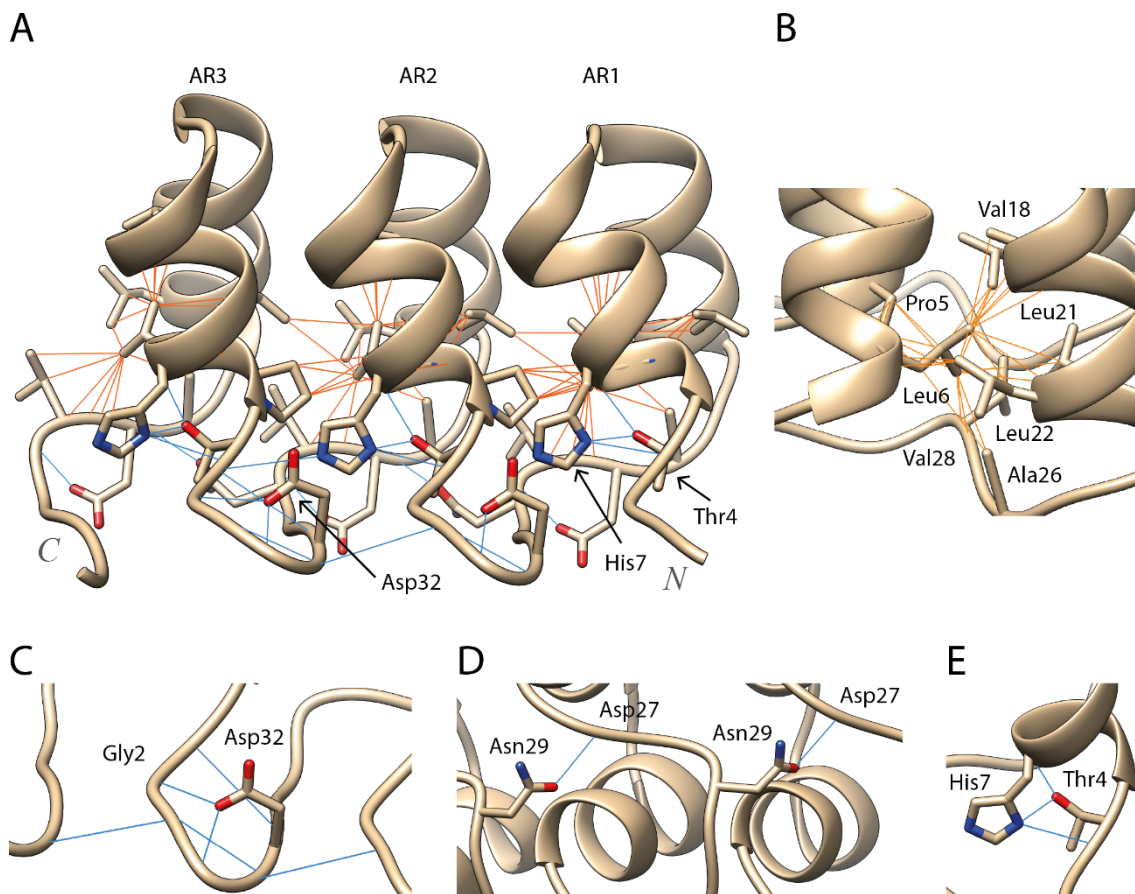
56

57 Figure 1. (A) Sequence logo of the ANK obtained with WebLogo (5) derived from the MSA generated in
 58 this work. The Y axis indicates the probability of observing an amino acid at any position within the motif;
 59 (B) Tertiary structure of an ankyrin repeat, coloured by secondary structure class: helices in red and coil in
 60 blue; (C) Representation of the complementary surfaces of individual ARs that form the human gankyrin
 61 ARD surface. N- and C-capping AR surfaces are coloured in purple and green respectively, whereas
 62 modular ones are coloured in blue and orange. (PDB ID: 1UOH) (6). Structure visualization with UCSF
 63 Chimera (7).

64 Proteins containing ankyrin repeats are known to bind many different protein and small
 65 molecule substrates. The concave face of an ARD, comprising the β -turn/loop region and the first
 66 α -helix is often associated with substrate binding. (8). Recent evidence suggests that ARDs might
 67 not only be able to bind small ligands or proteins, but also a range of sugars and lipids, thus
 68 extending their versatility and flexibility in substrate binding (9). This, coupled with the success
 69 that designed ankyrin repeat proteins (DARPin) are having in the clinical field (10) make
 70 ankyrins an extremely interesting target to study.

71 Within the sequence variability found in protein repeats, ankyrins are relatively conserved
 72 and multiple amino acid patterns can be observed within the ANK (11) (Fig 1A). The TPLH
 73 motif, positions 4-7, is highly conserved across all ankyrin repeats. It is found at the beginning of
 74 the first α -helix. Thr4 establishes three hydrogen bonds with His7 (Fig 2E), Pro5 starts the helix
 75 with a tight turn and Leu6 forms multiple hydrophobic interactions both within and between

76 repeats (Fig 2B). The loops are more diverse in sequence, yet certain patterns are apparent as well.
77 The subsequence GADVN, 25-29, can be observed in the loop connecting the second α -helix with
78 the β -turn. Gly25 breaks the second helix. Ala26 and Val28 form intra- and inter-repeat
79 interactions (Fig 2B), whereas Asp27 and Asn29 form hydrogen bonds with adjacent repeats (Fig
80 2D). Gly13 is found in the loop between the antiparallel helices. Asp32 and Gly2 are highly
81 conserved at the β -turn that connects repeats (Fig 2C). A total of six hydrogen bonds take place
82 in this turn. This explains why Asp32 is conserved, and Gly2 would be similar to Gly13 and
83 Gly25, and is conserved due to its flexibility and special structural features (12). In the second α -
84 helix, the [I/V]VXLLL hydrophobic motif is observed. The residues on this motif, except X19,
85 which is usually hydrophilic, form intra- and inter-repeat hydrophobic networks that are thought
86 to help keep together the ARD structure (13) (Fig 2B).



87

88 Figure 2. (A) Trio of ARs from a designed ankyrin repeat protein (14) (PDB: 5MA3). These three ARs
89 display the main known interactions responsible for the correct packing of the ARD. Hydrogen bond
90 interactions are depicted by blue lines whereas hydrophobic ones are depicted by orange ones; (B)
91 Hydrophobic network formed by Leu6 in the hydrophobic core of the domain; (C) Hydrogen bonding
92 network at the β -turn between positions Asp32-Gly2; (D) Inter-repeat hydrogen bonds between conserved
93 Asn29 and Asp27; (E) Thr4 forms three hydrogen bonds with His7. Structure visualization with UCSF
94 Chimera (7).

95 In protein sequences, functional or structurally important residues are constrained in
96 evolution, resulting in amino acid conservation between homologues. In a similar way, the
97 genomic distribution of genetic variation within a species is affected by factors such as gene
98 essentiality (15), and protein domain architecture (16). Synonymous variants do not change the
99 protein sequence, and as a result they appear randomly distributed in structure. Missense variants,
100 on the other hand, change the protein sequence and are consequently constrained in space.

101 Pathogenic missense variants cluster in three-dimensional structure around functionally important
102 regions, such as catalytic sites, whereas neutral or non-pathogenic missense variants tend to
103 aggregate on regions which are tolerant to amino acid substitutions (17). Recently, in a study of
104 all Pfam domain families, MacGowan, Madeira (18) found that positions conserved across
105 homologues but also depleted in missense variants within the human population were of particular
106 functional and/or structural relevance since they are heavily enriched in disease-associated
107 variants. They also found a subset of evolutionary unconserved positions that were missense
108 depleted. These positions were enriched in ligand, DNA and protein interactions as well as in
109 pathogenic variants, suggesting their functional importance within the protein domain.

110 In this paper we perform a novel analysis that combines human population genetic
111 variation from gnomAD (19) across ankyrin repeats in context with evolutionary variation and all
112 available ankyrin protein structures. This is the first in-depth application to a repeat family of the
113 concepts developed in our earlier work across all Pfam domains (18). Application to a repeat
114 family boosts the statistical power of the method and highlights the positions in the ANK most
115 likely to be important for structural stability as well as those relevant to substrate specificity. We
116 anticipate that this work will be of value to those interested in understanding the function of ANK
117 containing proteins as well as those aiming to engineer novel AR specificity.

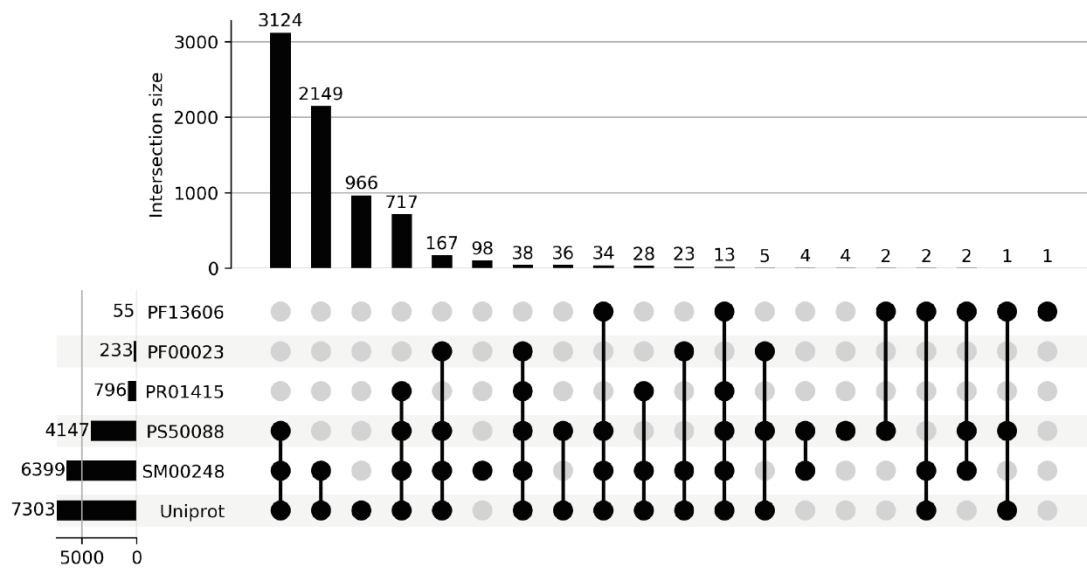
118 **Methods**

119 **Sequence extraction and database integration**

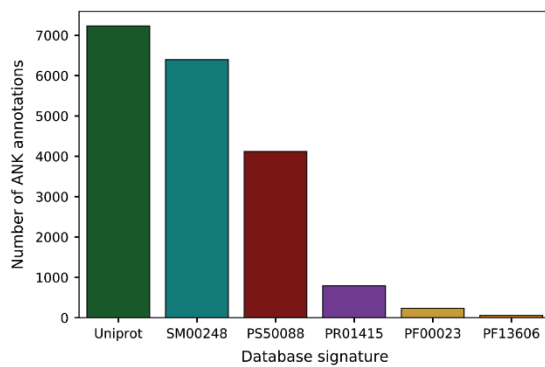
120 The UniProt (20), SMART (SM00248) (21), ProSite (PS50088) (22), PRINTS (PR01415) (23)
121 and PFAM (PF13606, PF00023) (24) were scanned for ankyrin repeat motif (ANK) definitions
122 in all species. The databases use slightly different algorithms resulting in variation in the number
123 as well as the length and coordinates of annotations between them (Fig 3A). Accordingly, we
124 retrieved all ankyrin repeat sequences found in Swiss-Prot reviewed proteins from the following
125 databases: UniProt (7,230 ankyrin repeats), SMART (6,396), ProSite (4,119), PRINTS (796) and
126 PFAM (288) (Fig 3B) resulting in a total of 18,825 ANK annotations. After redundancy filtering,
127 we established a high-quality set of 7,407 ankyrin repeat sequences: 4,109 (ProSite), 2,313
128 (SMART), 972 (UniProt) and 10 (PFAM) (Fig 3C) for analysis.

129

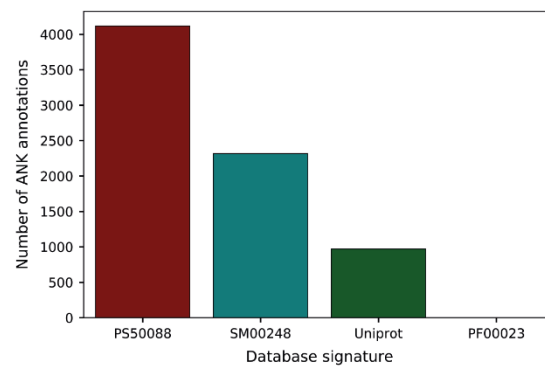
A



B



C



130

131 Figure 3. (A) Upset plot showing the distribution of ANK annotations and the overlap between different
 132 database signatures. Most of the annotations are shared between UniProt, SM00248 and PS50088. UniProt
 133 presents ≈ 1000 unique annotations which are not present in any other database; (B) This bar plot indicates
 134 the number of ANK annotations per database signature: 7,230, 6,396, 4,119, 796, 233 and 55 from left to
 135 right; (C) This bar plot shows the composition of the dataset resulting from the database merging, with
 136 ProSite accounting for $\approx 55\%$ of the annotations, SMART for $\approx 30\%$ and UniProt for the last $\approx 15\%$.

137 Multiple sequence alignment

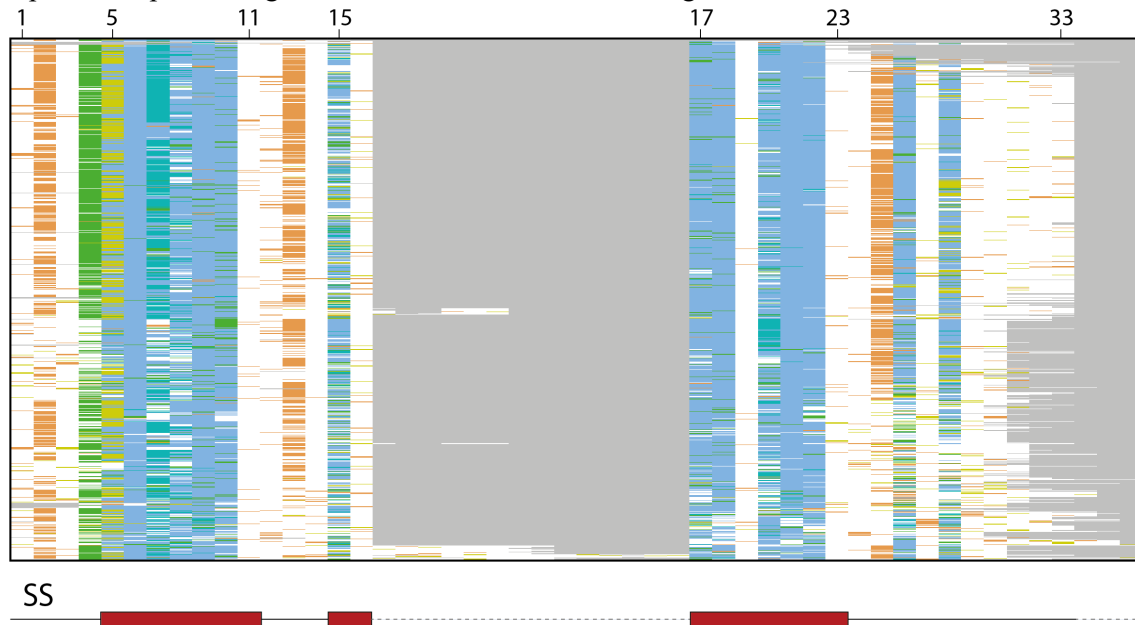
138 Several approaches were tried to align the 7,407 ankyrin repeat (AR) sequences, both sequence
 139 and structure-based. These included Clustal Ω (25), HMMER (26), T-Coffee (27), AMPS (28),
 140 Muscle (29) and STAMP (30). When applied to all 7,407 sequences, these aligners introduced
 141 many gaps and a high proportion of misaligned residues which were inconsistent with known key
 142 residues in the ankyrin repeat. Accordingly, the final multiple sequence alignment (MSA) was
 143 obtained by carrying out a series of sequences-to-profile multiple sequence alignments with
 144 Clustal Ω (Fig 4) as follows.

145 First, the sequences were divided into different groups according to their length and
 146 database of origin. Then, sequences that had the most common length, 33 residues, coming from
 147 the highest confidence database, ProSite, were aligned using Clustal Ω version 1.2.2 with defaults.
 148 Sequences introducing gaps in the 33 high-occupancy columns were removed and re-aligned with

149 a ClustalΩ sequences-to-profile alignment. Sequences inserting gaps yet again were removed
150 from the alignment.

151 For the rest of sequence groups, defined by sequence length and database of origin were
152 aligned to this growing alignment by consecutive sequence-to-profile alignments. As with the
153 first alignment, gap-introducing sequences were re-aligned and removed if necessary, for each
154 group.

155 At the end of this process, $\approx 98\%$ of the sequences were aligned in the resulting MSA.
156 The remaining 2% was formed by those gap-introducing sequences removed during the re-
157 alignment phase of the process. This 2% of sequences were re-aligned to the main alignment by
158 a profile-to-profile alignment shown as an overview in Figure 4.



159
160

161 Figure 4. Overview of the resulting MSA, including the 7,404 ankyrin repeat sequences. Only columns
162 presenting an occupancy $> 0.5\%$ are shown. Sequences are sorted by a tree generated in Jalview using the
163 average distance method and the BLOSUM62 matrix. Columns between 16-17 and after 33 represent
164 insertions in some ankyrin repeats. Red boxes below the overview indicate the location of the secondary
165 structure elements (SS), α -helices in this case, within the alignment. Grey dashed lines represent gaps and
166 are mostly found at low-occupancy columns. Columns are coloured according to the ClustalX colour
167 scheme (31). Hydrophobic residues are coloured in blue, glycines in orange, prolines in yellow, polar
168 residues in green and unconserved columns are coloured in white. Obtained with Jalview (32).

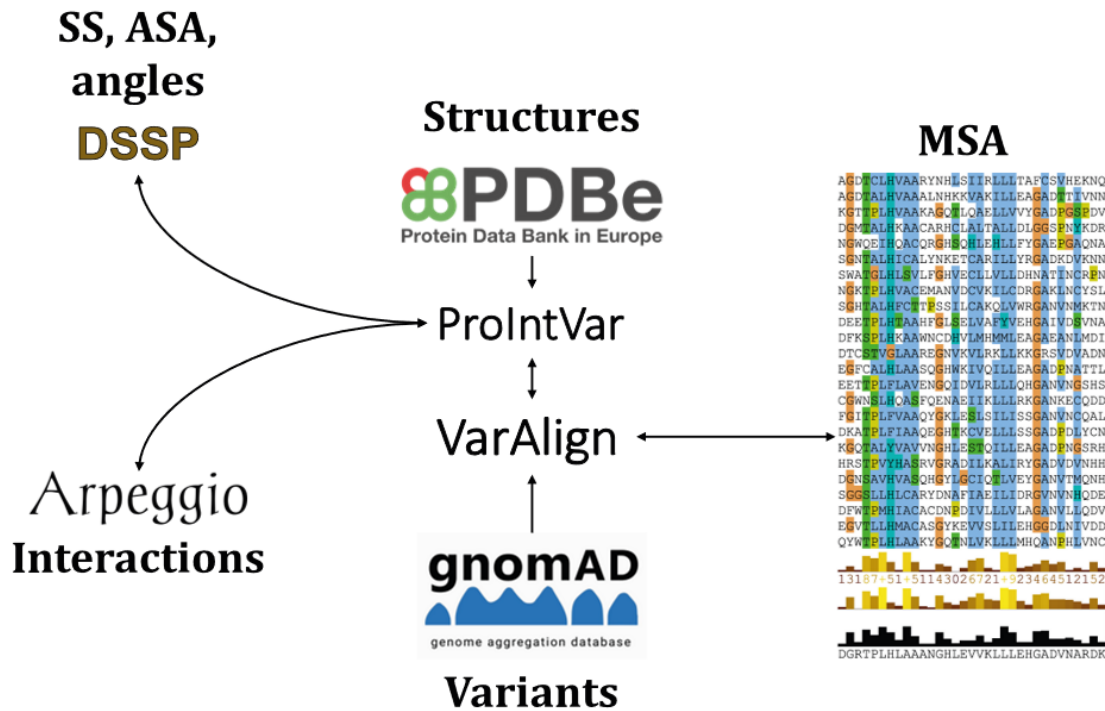
169 **VarAlign and ProIntVar**

170 A total of 35,691 variants found in the genome aggregation database (gnomAD) (19) coming from
171 1,435 human sequences were mapped to the MSA through VarAlign (18) (Fig 5).

172 419 sequences in the alignment were mapped to 209 different structures solved by X-ray
173 crystallography in the PDB (33-35) via SIFTS (36) through ProIntVar (37). These sequences
174 correspond to 419 unique ankyrin repeats, found in 80 proteins. The real-space R value (RSR)
175 and RSR-Z scores, as well as the real-space correlation coefficient (RSCC) quality metrics, as
176 calculated by (38), were retrieved by ProteoFAV (18) from the validation reports in PDBe. Only
177 residues with $RSCC > 0.85$ and $RSRZ < 2$ were considered for analysis. After this filtering step,
178 our structural dataset comprised 383/419 unique ARs coming from 176/209 PDBs, representing
179 73/80 proteins. This dataset included 11,186 of the 13,059 residues with structural coverage

180 before quality filtering. The average RSRZ per residue after filtering was -0.11 and the mean
181 RSCC had a value of $+0.95$.

182 DSSP (39) was run on all structures via ProIntVar and information from 381 ankyrin
183 repeat sequences was used to determine the consensus secondary structure as well as the relative
184 solvent accessibility (RSA) classification for all positions in the ANK, as described in
185 MacGowan, Madeira (18).



186

187 Figure 5. Diagram showing the main components of the pipeline. VarAlign retrieves variants found in
188 human sequences in the MSA from gnomAD. ProIntVar retrieves structures from the PDB and runs DSSP
189 and Arpeggio to get secondary structure, accessible surface area and inter-atomic contacts information.
190 Everything is mapped back to the residues and MSA columns (37).

191 Sequence divergence score

192 The Shenkin divergence score was used to characterise residue conservation at an alignment
193 position (40). This is a divergence score, based on Shannon's entropy (Equations 1 and 2).

194

$$V_{Shenkin} = 2^S \times 6 \quad (1)$$

195

$$S = - \sum_i^K p_i \log_2 p_i \quad (2)$$

196 Where S is Shannon's entropy and i is every one of the $K = 20$ different amino acid types. The
197 range of this diversity score is determined by Shannon's entropy. In a completely conserved
198 alignment column, one amino acid residue will be found with a frequency of 1.0, whereas the rest
199 will not be present, resulting in an entropy of 0.0, and a minimum $V_{Shenkin} = 2^0 \times 6 = 6$. At the
200 other extreme, an alignment column with all 20 amino acids at a frequency of $1/20 = 0.05$ would
201 give an entropy of $S \approx 4.32$, resulting in a maximum $V_{Shenkin} = 2^{4.32} \times 6 \approx 120$. Thus, low Shenkin
202 scores indicate higher conservation at a position and *vice versa*. To simplify the interpretation of
203 the score, we normalised the Shenkin score to 0-100 (Equation 3).

204
$$N_{Shenkin} = (V_{Shenkin} - V_{Shenkin_{min}}) / (V_{Shenkin_{max}} - V_{Shenkin_{min}}) \quad (3)$$

205 Where $V_{Shenkin_{min}}$ is the score of the most conserved column within the alignment, Position 9 with
206 a Shenkin score of 15.43 and $V_{Shenkin_{max}}$ is the score of the most diverse position, Position 3 with a
207 score of 103.96.

208 **Enrichment in variants**

209 The human genetic variants from gnomAD (19) were mapped to the MSA and missense variant
210 enrichment scores (MES) were calculated for the 33 positions of the ANK. MES is expressed as
211 the natural logarithm of an odds ratio (OR) and it represents the enrichment of variants in an
212 alignment column relative to the average for the other columns. Columns were classified as
213 depleted, enriched or neutral according to this MES (18). 95% confidence intervals and p-values
214 were calculated to assess the significance of these ratios (41).

215 **Enrichment in protein-substrate interactions**

216 For the structural analysis, the meaningful biological units were retrieved from PDBe. These are
217 the preferred assemblies for each structure, instead of the asymmetric units, which might not
218 reflect the packing of the protein observed in nature. All inter-atomic contacts were calculated by
219 Arpeggio (42). Atoms were considered to interact if they were within 5Å of each other.

220 We considered all interactions between an ankyrin repeat and any protein substrate
221 present in the preferred assembly as protein-protein interactions (PPIs). A log enrichment score
222 was calculated for PPIs per position in the motif in a similar manner to MES above. It is referred
223 as protein-protein interaction enrichment score (PPIES). The number of protein-protein
224 interactions per alignment column was normalised by the structural coverage of that column in
225 structures presenting an interaction between an ARD and a bound peptide substrate. We
226 considered that there was evidence of contact between an AR position and a bound peptide
227 substrate if there was at least one inter-atomic contact involving the repeat position and the
228 substrate in at least one of the structures representing the complex.

229 **Enrichment in intra-repeat contacts**

230 A contact map, shown as a 33×33 matrix, for the 33 positions in the ANK, was calculated to
231 show how often two positions interact within an AR. Each cell shows the proportion of repeats,
232 where evidence of contact between a given pair of residues has been observed. The absolute
233 frequency is normalised by the coverage of a given pair of residues within a repeat. This intra-
234 repeat contact map is symmetric. Thus, a given cell $c_{i,j} = c_{j,i}$. Contacts between adjacent residues
235 are not shown, resulting in a null diagonal.

236 Enrichment in intra-repeat contacts per position was calculated. Since the intra-repeat
237 contact matrix is symmetric, the total number of contacts per residue, C_i , was calculated using
238 Equation 4, where $c_{i,j}$ is the absolute frequency of contacts between any two amino acid residues
239 present at positions i and j within the $K = 33$ positions in the ANK. The same approach was used
240 to calculate the total structural coverage per ANK position, O_i , (Equation 5) where $o_{i,j}$ is the
241 absolute frequency of both positions i and j being present in the same repeat.

242
$$C_i = \sum_{j=1}^K c_{i,j} \quad (4)$$

243
$$O_i = \sum_j^K o_{i,j} \quad (5)$$

244 The total number of contacts and coverage of a position were calculated as the sum of all their
245 contacts and coverages, respectively (Equations 6 and 7).

$$246 \quad O_t = \sum_i^K O_i \quad (6)$$

$$247 \quad C_t = \sum_i^K C_i \quad (7)$$

248 Enrichment in these contacts was calculated per position in the same fashion as for variants and
249 PPIs. The same analysis was carried out on inter-repeat contacts but was of limited value and so
250 not included in this paper.

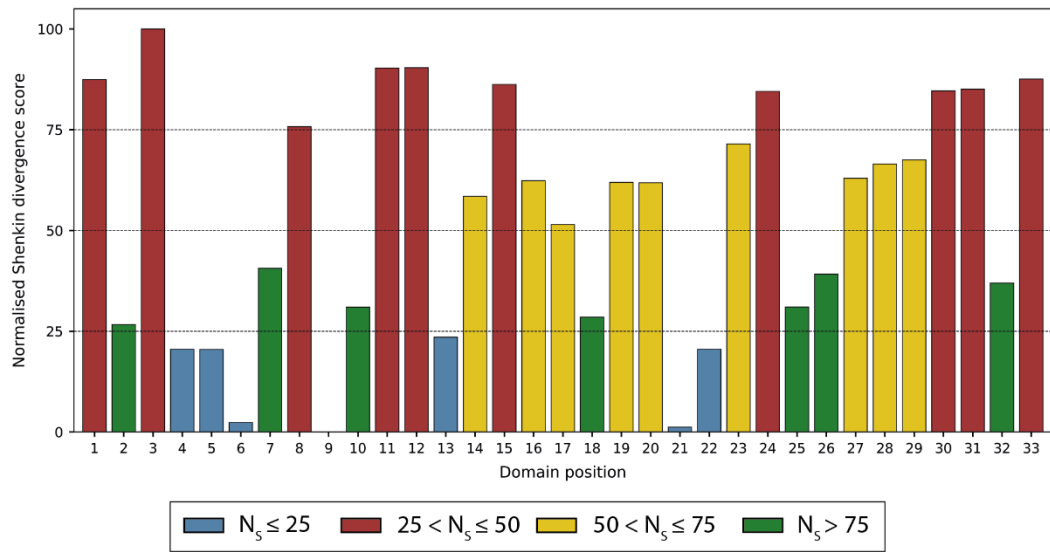
251 **Results and Discussion**

252 In this work, 7,407 ankyrin repeat sequences, including both human and other species, were used
253 to build a multiple sequence alignment and conservation profile of the motif. Human genetic
254 variation data coming from 1,435 human ankyrin repeats were used to study the distribution of
255 variation within the motif. Moreover, 176 three-dimensional structures, representing a total of
256 383 different ankyrin repeats were used to structurally characterise in detail this motif by
257 secondary structure, residue solvent accessibility, intra-domain contacts and protein-protein
258 interactions. For the first time, human population variation data was used to explain the
259 evolutionary constraint acting upon this family of protein repeats, integrating at the same time
260 these data with structure and sequence divergence.

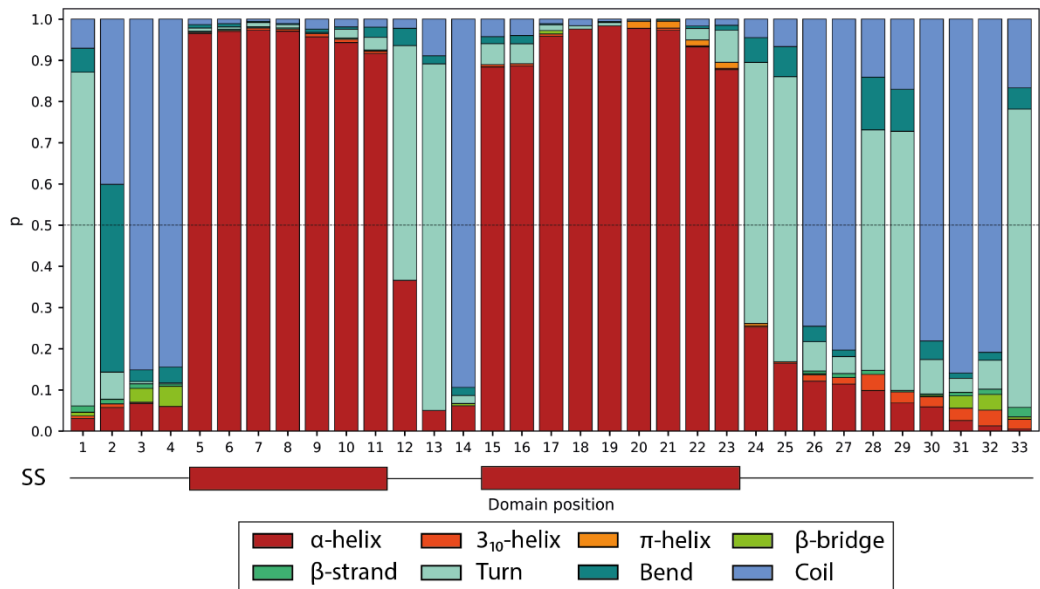
261

262

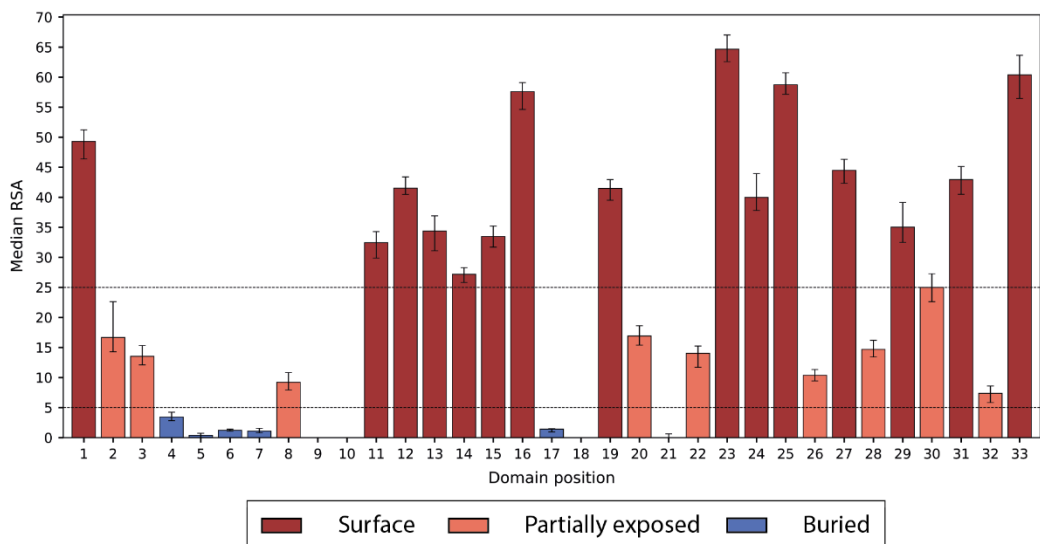
A



B



C



263 Figure 6. (A) Normalised Shenkin divergence score per domain position (Eq. 3) calculated from the MSA
264 containing 7,404 sequences. Positions are coloured according to their normalised Shenkin score as the
265 legend indicates; (B) Secondary structure assignment per position. Within each position, each coloured bar
266 represents the frequency of the eight states defined by DSSP: α -helix, 3_{10} -helix, π -helix, β -bridge, β -strand,
267 turn, bend and coil, observed for the residues with structural coverage at that column in the MSA. Most
268 helices range from 5-11 and 15-23 and finish in 5-turns, usually at positions 12-13 and 23-24. Two β -turns
269 are observed at positions 28-29 and 33-1; (C) Median residue relative surface accessibility per position,
270 calculated from DSSP's accessible surface area (39) as described in Tien, Meyer (43). Error bars indicate
271 95% CI of the median. Positions were classified according to the specified thresholds: surface (RSA \geq
272 25%), partially exposed (5% < RSA < 25%) or buried (RSA \leq 5%) (44).

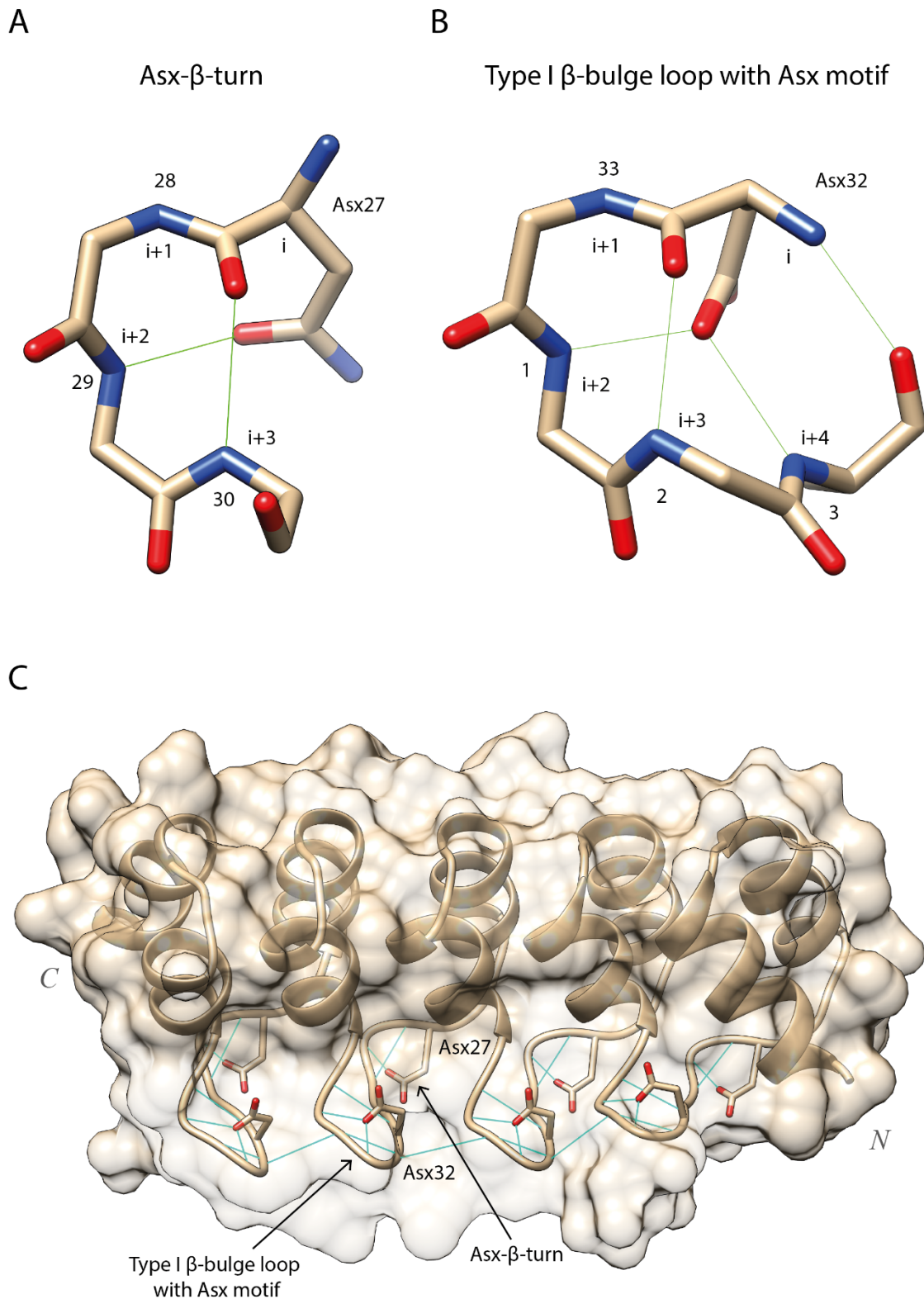
273 Conservation profile

274 The conservation profile derived from our MSA agrees with previous work (13). Figure 6A
275 shows the normalised Shenkin divergence score per position in the motif. As described in
276 Methods, this score goes from 0-100. Among the most conserved positions ($N_{Shenkin} < 25$) we find
277 Thr4, Pro5, Leu6, belonging in the TPLH motif, 4-7 as well as Ala9, Gly13, Leu21 and Leu22.
278 Some of the most evolutionary diverse positions, on the other hand, include positions 1, 3, 11, 12
279 and 33 among others, all presenting $N_{Shenkin} > 75$. Most of the highly diverse positions are found
280 on the concave surface and contribute to the variable interface where most of the substrate binding
281 takes place.

282 Secondary structure

283 Figure 6B shows the secondary structure assignment of the 33 positions in the ANK. Most of the
284 repeats present a seven-residue long first helix ranging from the fifth to the 11th position and a
285 second helix that in most cases is nine-residues long and extends from the 15th to the 23rd position.
286 Our results also show four turns along the ANK. Two of these turns, found at positions 12-13 and
287 24-25, are 5-turns and simply indicate the end of the α -helices, whereas the other two, positions
288 28-29 and 33-1, are β -turns. These two β -turns were classified as type I β -turns according to the
289 ϕ and ψ dihedral angles distribution of consensus columns 33, 1 and 28, 29 (45). Positions 27 and
290 32 in the alignment present either Asn or Asp with a high frequency of 44% and 58% respectively.
291 Consequently, we classified the turns they initiate as Asx motifs (Fig 7A-C). The turn at positions
292 27-30 was classified as an Asx- β -turn and the one at 32-2 as a type 1 β -bulge loop with an Asx
293 motif (46). Repeats that do not have Asx at positions 27 and 32, form a simple β -turn, instead of
294 an Asx motif since they lack the extra hydrogen bonds that this secondary structure motif requires.
295 The conservation of these Asx residues on both turns, suggests a structural relevance and role of
296 these Asx motifs on the correct packing of the ARD.

297



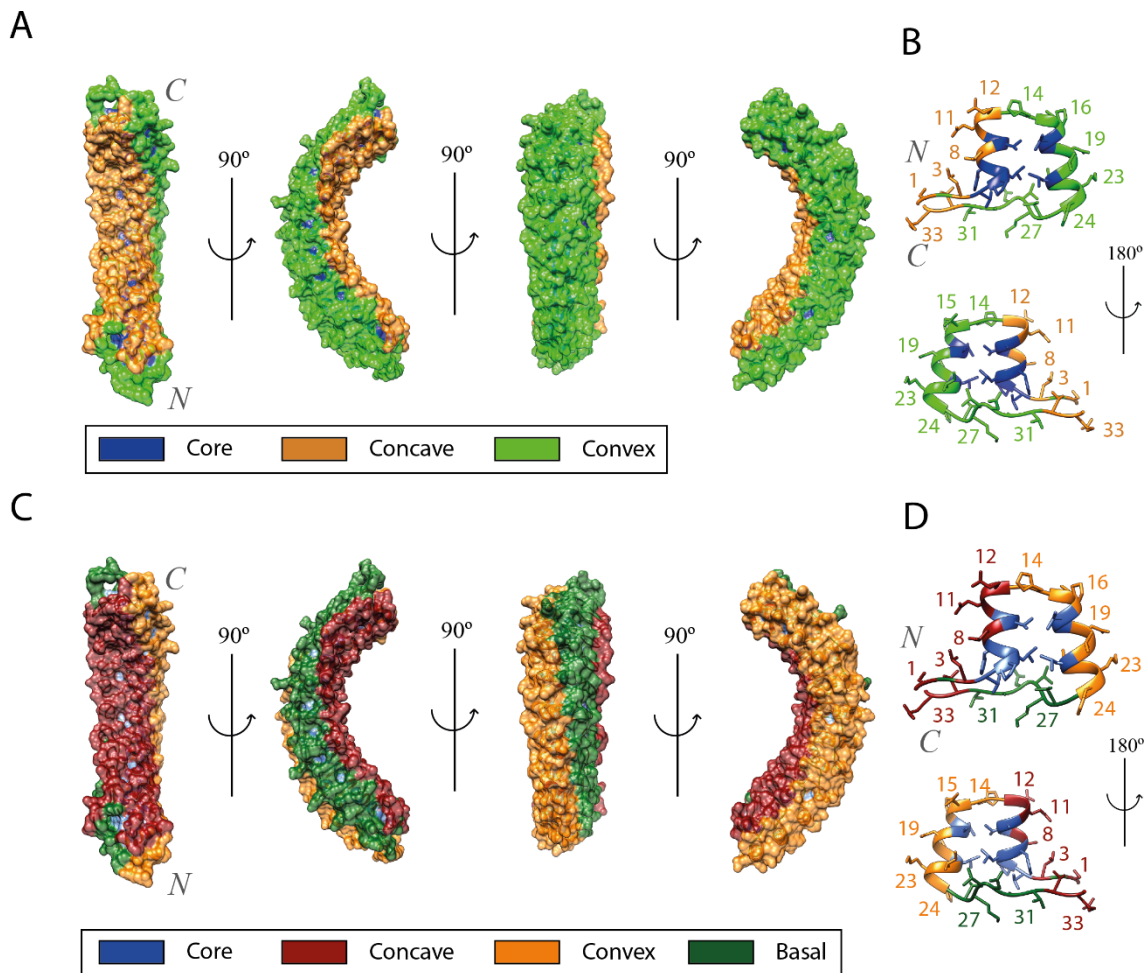
298

299 Figure 7. Hydrogen bonding patterns of the two Asx motifs found in the ANK and their location within the
300 ARD. Only repeats with either Asn or Asp at these positions will present this hydrogen bonding pattern.
301 (A) Asx- β -turn at positions 27-30. Conserved Asx, i.e., Asp/Asn, side chain at position $i = 27$ forms extra
302 hydrogen bond with backbone N at position $i + 2$; (B) Type I β -bulge loop with Asx motif at positions 32-
303 3. Conserved Asx side chain at domain position $i = 32$ forms two hydrogen bonds with backbone N of
304 residues $i + 2$ and $i + 4$. The rest of the hydrogen bonds originate from the backbone of the residues and are
305 not specific of Asx motifs. PDB: 5MA3 (14); (C) DARPin-8.4 (Barandun J, Schroeder T, Mittl PRE,

306 Grutter MG). Light blue lines represent the hydrogen bonds that determine these secondary structure motifs.
307 The conservation of these Asx residues at positions 27 and 32, together with their depletion in missense
308 variants (Fig 9) and the hydrogen bonding network they facilitate, suggests that these Asx motifs are one
309 of the most structurally important components of the ankyrin repeat domain structure. Figure obtained with
310 UCSF Chimera (7).

311 Relative solvent accessibility and surface classification

312 The surface of the ankyrin repeat domain has previously been divided into two faces: concave
313 (positions 32-12) and convex (positions 13-31) (Fig 8A, B) (47). Positions with high RSA (RSA
314 $\approx 50\%$), such as 1, 12 and 33 are found near positions 13 and 32. Due to their high solvent
315 accessibility, these positions were used to define ridges at the limits of the concave and convex
316 surface. However, our analysis of all available structures also showed positions 23 and 25 to have
317 a high RSA (Fig 6C). In addition, in the same fashion as positions 1, 33 and 12, positions 23 and
318 25 from different repeats form a ridge on the domain structure. This ridge suggests the definition
319 of a third surface of the domain or basal surface as shown in Figure 8C and enabled the
320 classification of all positions that were not buried into one of the three defined surfaces, (Table
321 1). This classification is shown in Figure 8C for an ARD containing 12 repeats (48).



322

323 Figure 8. Comparison of the original definition of the ARD surfaces (A, B) with the new definitions derived
324 from the results of this study (C, D). All panels refer to the D34 region of ANK1 ARD, PDB accession:
325 1N11 (48). This structure shows 12 out of the 23 ARs found on this ARD; (A) Surface of an ARD. Residues
326 conforming the concave surface are coloured in orange, residues on the convex surface in green and buried
327 residues in blue; (B) Surface of an individual repeat. The first α -helix and the β -turn region form the concave
328 surface, whereas the second helix and the loop form the convex one; (C) Residues conforming the concave

329 surface are coloured in dark red; residues on the convex surface in orange; the basal surface is coloured on
330 dark green and buried residues in blue. Figure obtained with UCSF Chimera (7).

331 Some blue-coloured regions can be observed on the ARD surface on Figure 8C. These
332 are the side chains of buried residues within the motif dominated by Thr4 and His7. The correct
333 classification of the positions in the ANK as either buried or any of the defined surfaces is critical
334 to calculate accurate enrichment scores in missense variants and protein-protein interactions on a
335 surface basis later in the analysis.

336 Table 1: Classification of the 33 positions within the ANK in the different surfaces.

Surface	Consensus residue positions
Core	4, 5, 6, 7, 9, 10, 17, 18, 21
Concave	1, 2, 3, 8, 11, 12, 32, 33
Convex	13, 14, 15, 16, 19, 20, 22, 23, 24
Basal	25, 26, 27, 28, 29, 30, 31

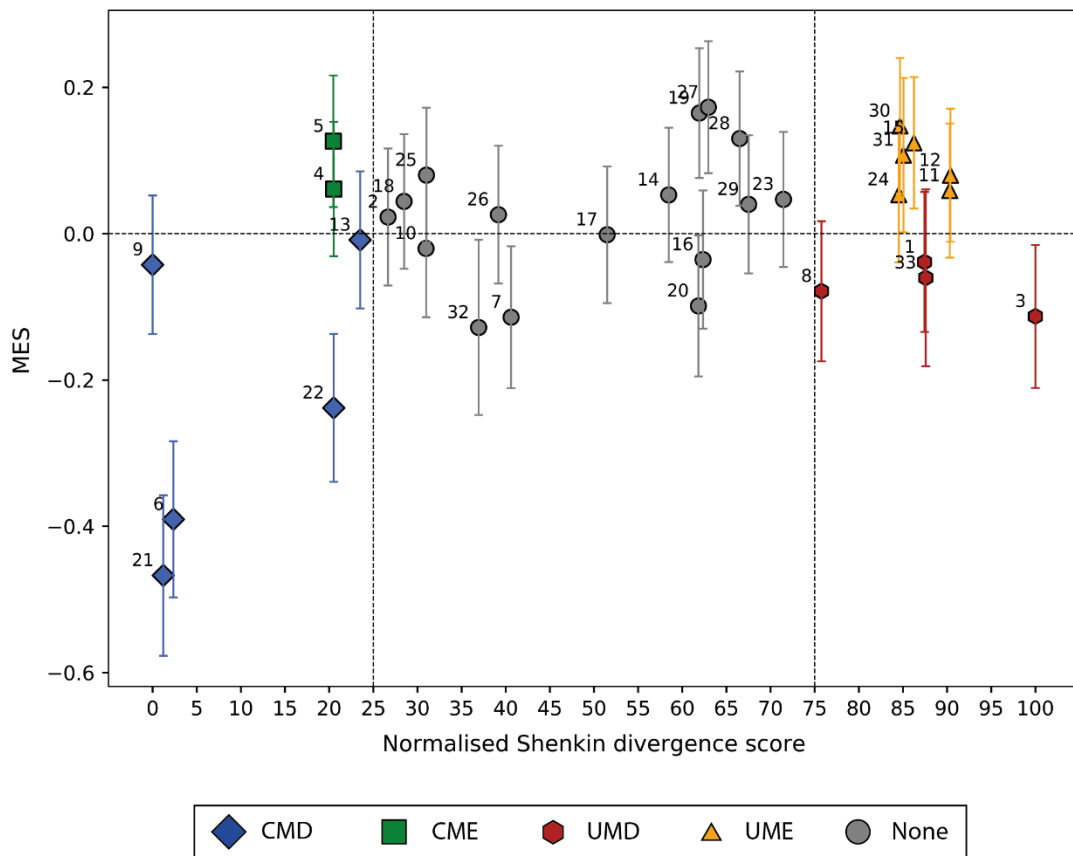
337

338 Missense variants enrichment analysis

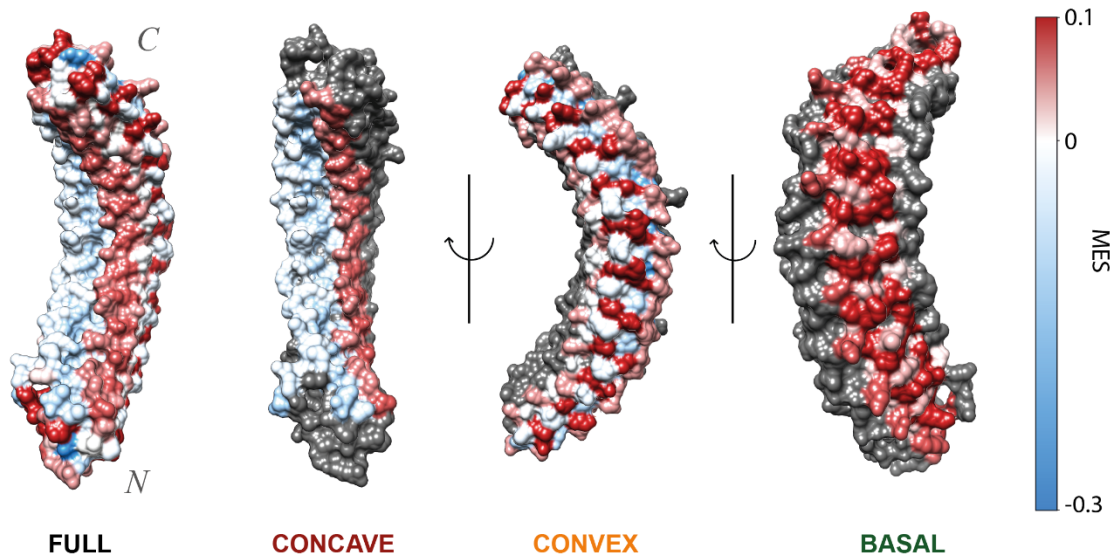
339 21,338 missense variants from 1,435 human ankyrin repeat sequences were used to calculate
340 column-specific missense enrichment scores (MES). The MES measures how enriched in
341 missense variants an alignment column is compared to the average of the other columns in the
342 alignment (18). The 33 columns of the motif were classified into four categories according to
343 their normalised Shenkin divergence score ($N_{Shenkin}$) and MES. Columns with $0 \leq N_{Shenkin} \leq 25$ and
344 $MES < 0$ were classified as conserved and missense depleted (CMD), whereas columns satisfying
345 $0 \leq N_{Shenkin} \leq 25$ and $MES > 0$ were called conserved and missense enriched (CME). We also
346 classified those columns with $75 \leq N_{Shenkin} \leq 100$ as unconserved and either missense depleted
347 (UMD) if $MES < 0$ or enriched if $MES > 0$ (UME). $N_{Shenkin}$ ranges from 0 for the most conserved
348 (Position 9) to 100 for the most divergent (Position 3) column within the ANK alignment. As a
349 consequence, positions with $N_{Shenkin} < 25$ will be those with divergence scores between the
350 minimum and $\frac{1}{4}$ of the maximum score and positions with $N_{Shenkin} > 75$ will be those with
351 divergence scores on the fourth quartile of the range. Figure 9A shows the enrichment in human
352 population missense variants per position in the ANK relative to their Shenkin divergence score.
353 Positions that are depleted in missense variants relative to the rest of positions within the ANK
354 are the most interesting and are likely to be functionally important. Depletion in missense
355 variation directly results from evolutionary constraint within the human population and is
356 therefore indicative of functional and/or structural relevance (18).

357

A



B



358

359 Figure 9. (A) Relative Missense Enrichment Score (MES) against normalised Shenkin divergence score for
 360 the 33 positions of the domain. Blue diamonds: CMD positions (6, 9, 13, 21, 22); Green squares: CME
 361 positions (4, 5), UMDs are coloured in red hexagons (1, 3, 8, 33) and UMEs in orange triangles (11, 12,
 362 15, 23, 24, 30, 31). Error bars represent 95% CI of the MES, i.e., $\ln(OR)$. Error bars for the Shenkin score
 363 are not shown as the uncertainty associated with it is negligible. Positions coloured in grey circles are
 364 classed as “None”, for they do not meet our divergence score thresholds; (B) D34 region of ANK1 ARD,
 365 PDB accession: 1N11 (48) This structure shows 12 out of the 23 ARs found on this ARD. Residues are

366 coloured according to the missense enrichment score of the alignment column they align to in the MSA.
367 The colour scale goes from blue (missense-depleted) to red (missense-enriched) going through white
368 (neutral). From left to right, the full domain, then concave, convex and basal surface are coloured. On each
369 of the last three representations, only one surface is coloured. Residues not belonging in said surface are
370 coloured in grey. Overall, the concave surface is coloured in a light blue colour (except positions 11 and
371 12), indicating its depletion in missense variants, relative to the other positions within the ANK. Figure
372 obtained with UCSF Chimera (7).

373 The relationship between residue solvent accessibility and enrichment in missense
374 variants was examined. As expected, on average, buried residues ($RSA \leq 5\%$) were depleted in
375 missense variants relative to residues present on the surface ($MES = -0.10, p = 1.9 \times 10^{-7}$).
376 Furthermore, residues present on the concave surface of the ankyrin repeat domain were
377 significantly depleted in missense variants relative to the other surfaces, ($MES = -0.08, p =$
378 4.4×10^{-4}). The convex surface was neither enriched nor depleted, whereas the basal surface
379 was significantly enriched in missense variants: ($MES = 0.09, p = 6.2 \times 10^{-6}$). Moreover, the
380 basal surface is significantly enriched in missense variation relative to the convex one ($MES =$
381 $0.08, p = 8.8 \times 10^{-4}$). These results can be observed in structure in Figure 9B and are further
382 discussed in “*Different surfaces of the ARD*” below.

383 **Ankyrin repeat contact maps and enrichment**

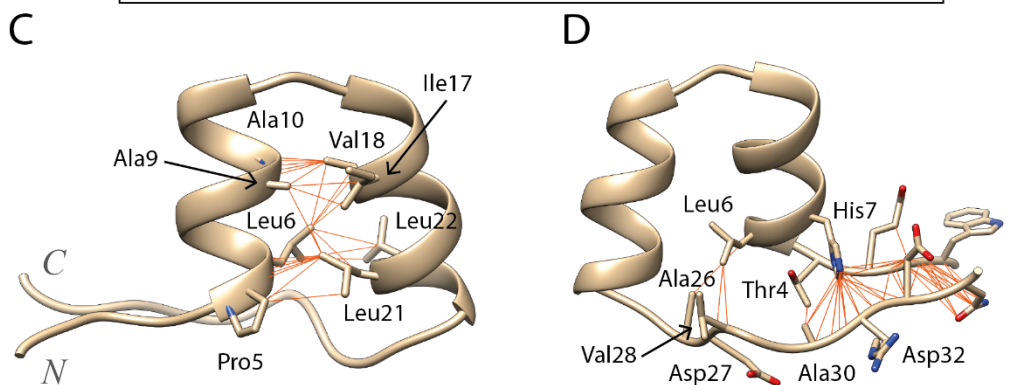
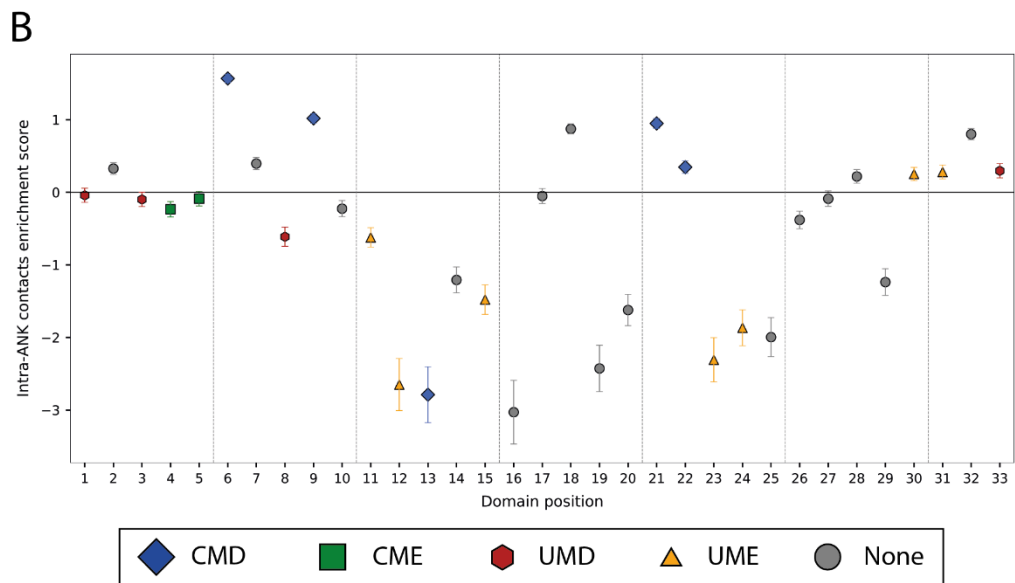
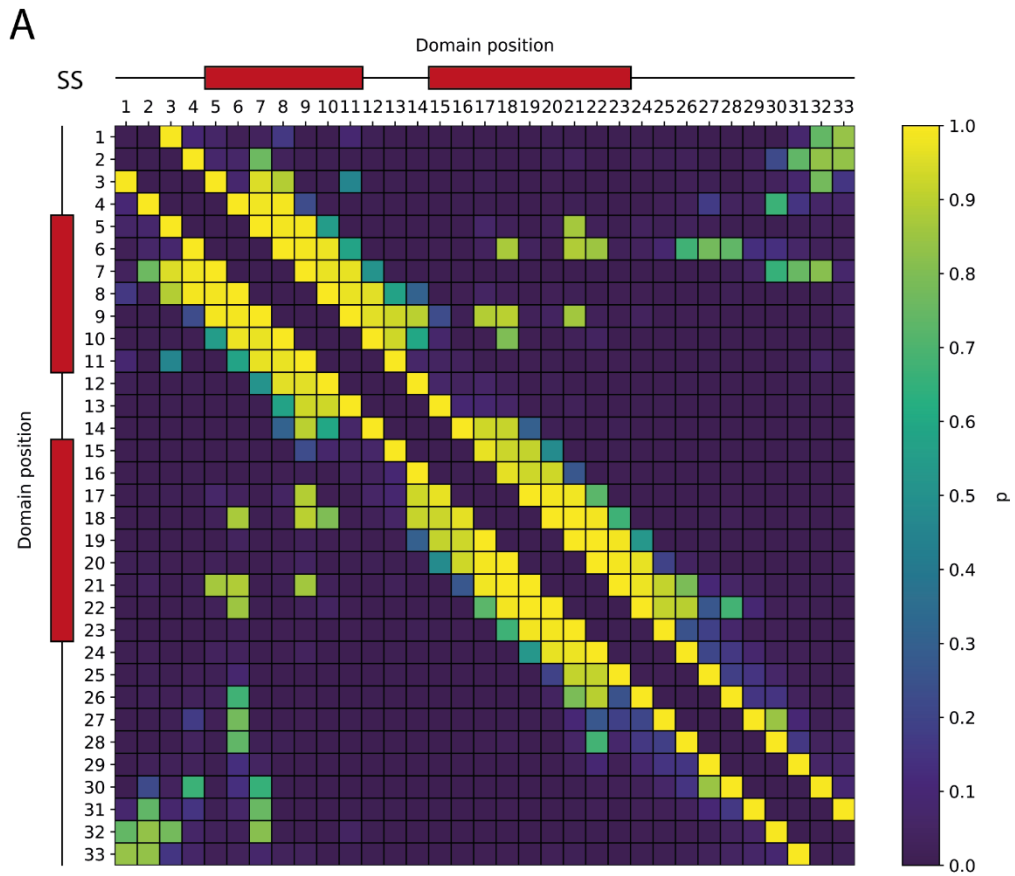
384 In this work, contacts across all known ankyrin repeat structures were considered instead of just
385 a single repeat or domain. This allowed a comprehensive contact map for the repeat motif to be
386 calculated as well as enrichment scores for each residue’s contacts within the repeat, thus
387 highlighting the most structurally important positions.

388 **Intra-repeat contacts**

389 Figure 10A shows the symmetric contact matrix that defines the ankyrin repeat motif. Contacts
390 between residues within 2-5 amino acids of each other are around the diagonal. Most other
391 contacts are between the residues along first and the second α -helices, from positions 5-11 and
392 15-23, respectively or contacts between residues close in sequence within the loops. This pattern
393 of contacts is typical of helical or turn secondary structures. Accordingly, we focused on contacts
394 most relevant to the ANK fold, i.e., helix-helix contacts, by filtering out contacts between
395 positions within ≤ 6 residues of each other. Figure 10B shows the enrichment in these intra-repeat
396 contacts for each position within the ANK. CMD positions are among the most enriched in intra-
397 repeat interactions which suggests an important role in ankyrin repeat packing and may explain
398 their conservation across homologs and depletion in missense variants within the human
399 population.

400

401

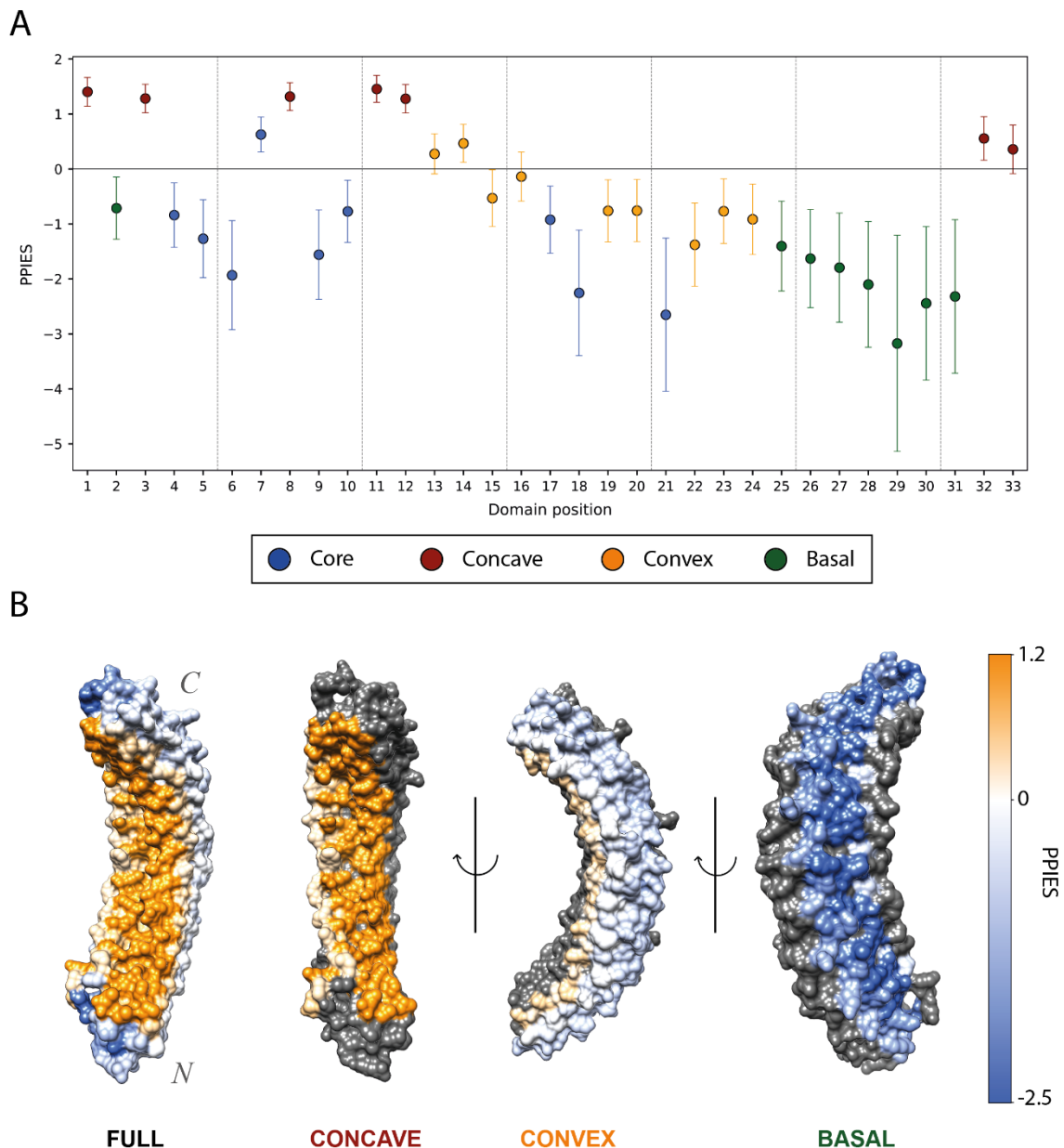


402 Figure 10. (A) Contact map for intra-repeat residue-residue interactions. Cells are coloured according to
403 the probability of observing contact between two positions with the viridis colour palette. Red boxes above
404 axis indicate the location of the secondary structure (SS) elements, α -helices, in the motif; (B) Intra-repeat
405 contacts enrichment plot. Error bars indicate 95% CI of the enrichment score, i.e., $\ln(\text{OR})$. Data points are
406 coloured according to their missense enrichment and residue conservation classification (Fig 9); (C) Cluster
407 of intra-repeat contacts between the first and second helices. Residues 5, 6, 9 and 10 in the first helix interact
408 with residues 17, 18, 21 and 22 by forming hydrophobic interactions. These positions are all buried and
409 conserved; (D) Cluster of intra-repeat contacts between the start and end residues of an AR. These
410 interactions are not as specific as the ones in the first cluster and they include diverse positions such as 1,
411 3, 31 or 33. These are the most frequently observed contacts across all structure displayed in an example
412 repeat. PDB: 5MA3 (14). Figure obtained with UCSF Chimera (7).

413 **Protein-substrate interaction enrichment**

414 Figure 11A shows the enrichment in Protein-Protein interactions (PPIs) per position in the ANK.
415 Out of the 176 protein structures that satisfied our quality thresholds, as described in Methods, 63
416 include protein substrates. These represent the interaction between 35 different ARDs and their
417 substrates, accounting for a total of 142 repeats. All the positions that are found on the concave
418 surface are enriched in PPIs. His7 is highly conserved and even though it is buried, part of its side
419 chain is accessible to the concave surface. This way, position 7 interacts with the substrate and
420 appears enriched in the analysis. Positions 13 and 14 define the beginning of the convex surface,
421 at the loop between the two helices. These positions, despite not forming part of the concave
422 surface, are very close to it. Thus, they are enriched in PPIs, though not as significantly as those
423 positions on the concave surface.

424



425

426 Figure 11. (A) Protein-substrate interactions enrichment plot. Error bars indicate 95% CI of the protein-
 427 protein interactions enrichment score (PPIES), i.e., $\ln(\text{OR})$. Data points are coloured according to their
 428 surface classification (Table 1); (B) D34 region of ANK1 ARD, PDB accession: 1N11 (48). This structure
 429 shows 12 out of the 23 ARs found on this ARD. Residues are coloured according to the PPIES of the
 430 alignment column they align to in the MSA. The colour scale goes from blue (depleted in PPIs) to orange
 431 (enriched in PPIs) going through white (neutral). From left to right, the whole domain, then the concave,
 432 convex and basal surface are coloured. On each of the last three representations, only one surface is
 433 coloured. Residues not belonging in that surface are coloured in grey. Overall, the concave surface is
 434 coloured in a strong orange colour, indicating its importance in protein binding, whereas the basal one
 435 presents a dark blue colour, indicative of its overall depletion in PPIs. Figure obtained with UCSF Chimera
 436 (7).

437 We also compared the enrichment in PPIs between different surfaces. As expected, buried
 438 residues were significantly depleted on average relative to surface residues, ($PPIES =$
 439 $-1.02, p < 10^{-16}$). Compared to residues belonging to other surfaces, concave residues are
 440 highly enriched in PPIs: ($PPIES = 1.86, p < 10^{-16}$). Conversely, convex and basal residues are
 441 both depleted in PPIs relative to residues present in the other surfaces: ($PPIES = -0.79, p <$

442 10^{-16}) and ($PPIES = -2.19, p < 10^{-16}$), respectively. In addition, the direct comparison
443 between basal and convex showed that there is a significant difference regarding their
444 involvement in substrate binding. Residues present on the convex surface were in average
445 enriched in PPIs relative those in the basal surface: ($PPIES = 1.31, p = 5.40 \times 10^{-15}$). All
446 these differences in enrichment in PPIs between different surfaces can be observed in Figure 11B,
447 which shows the different surfaces of an ARD, where residues are coloured according to the
448 PPIES of the column they align to. These results agree with (8) and show the prevalence of the
449 concave surface in substrate binding. They also illustrate the rare, though existing, convex binding
450 as well as the practically null contribution of the basal surface to substrate binding.

451 **Different surfaces of the ARD**

452 In this work, we define the binding mode of an ARD as given by the number of repeats and
453 residues that bind the substrate, as well as the surface the latter belong to. These modes can either
454 be absolute or combined/mixed. In the former, one surface dominates the binding, whereas in the
455 latter, a combination of different surfaces accounts for most of the substrate binding residues. For
456 this part of the analysis, only those proteins with a minimum of two repeats and four residues
457 binding the substrate were considered. Of the remaining 25 proteins, 21 (84%) presented a
458 concave binding mode, only one (4%) presented a convex binding mode and none presented a
459 basal mode. The other three (12%) proteins presented a mixed binding mode, where concave,
460 convex, basal and even buried residues participate in the substrate binding.

461 These binding surfaces present different patterns of enrichment in variation as well as
462 protein-substrate interactions. The concave surface is significantly depleted in missense variants
463 and enriched in PPIs, whereas the basal is the complete opposite and is enriched in missense
464 variants and depleted in PPIs. These results confirm the dominance of the concave binding mode.
465 In addition, we have observed that ARDs can also present a convex binding mode (49), whereas
466 no basal binding mode was observed in the dataset. The differential importance in substrate
467 binding seems to influence the distribution of missense variants within the motif.

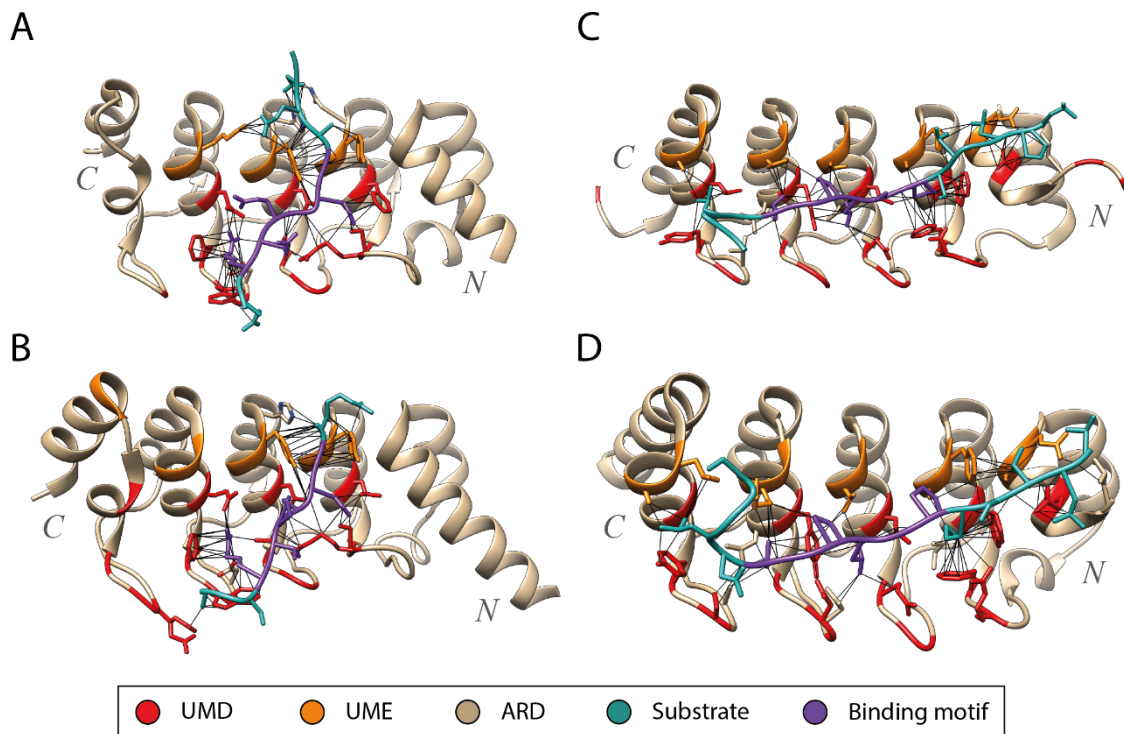
468 **Conserved and missense depleted positions**

469 Positions 6, 9, 21 and 22 were found to be highly conserved and depleted in missense variants
470 relative to the other motif positions (CMD). These positions are mostly buried, and present
471 hydrophobic residues. CMDs are enriched in intra-repeat contacts. This population-level
472 constraint agrees with the amino acid conservation and is proof of the structural relevance of these
473 residues. Positions 7 and 32 are not as conserved as the residues we have classified as CMD;
474 however, they are significantly depleted in missense variants as well. These two residues are
475 structurally relevant due to the hydrogen bonding networks they create, as can be seen in Figure
476 2C, E.

477 **Unconserved and missense depleted positions**

478 It is known that the concave surface allows high sequence variability, in order to accommodate
479 the diversity of protein substrates that ankyrins bind (8). Positions 1, 3, 8 and 33 are amongst the
480 most diverse positions within the ANK, though at the same time depleted in missense variants in
481 the human population. These positions are enriched in PPIs ($PPIES = 3.6, p < 10^{-16}$) and
482 constitute most of the concave surface of the ARD. Missense depletion at these sites show that
483 they are constrained at a population level, thus confirming the functional importance of these
484 residues.

485 Figure 12 illustrates how the ARDs of homolog pairs with gene names
486 ANKRA2/RFXANK and TNKS1/TNKS2 bind their protein substrates. ANKRA2 and RFXANK
487 are human proteins that are involved in the regulation of transcription by RNA polymerase II.
488 Both ankyrin repeat domains present five repeat units. The domains are very similar in sequence,
489 including UMD and unconserved positions 11 and 12, which do not vary across these proteins'
490 homologs. Multiple structures have been solved portraying the interaction between these ARDs
491 and more than five different protein substrates. All the substrates present the shared binding motif
492 PXLPX[I/L] (50) (51). A similar pattern can be observed with TNKS1 and TNKS2 and the
493 substrates they bind, which share the tankyrase binding motif RXXPDG (52).



501 Figure 12. ARDs in complex with substrates. (A) RFXANK and RFX5 (PDB ID: 3V30) (50); (B) ANKRA2
502 and HDAC4 (3V31) RN1262; (C) TNKS2 and ARPIN (4Z68) (53); (D) TNKS1 and USP25 (5G7) (54).
503 UMD positions (red) and UMEs 11, 12 (orange) are conserved across proteins that bind similar substrates
504 (dark cyan). For example, these positions are conserved across TNKS2 and TNKS1, which are known to
505 bind substrates with the motif RXXPDG (purple). Similarly, RFXANK and ANKRA2, bind substrates with
506 the motif PXLPX[I/L] (purple). Figure obtained with UCSF Chimera (7).

501 These examples show how ankyrin domains that present similar concave surfaces,
502 determined by their UMD positions (1, 3, 8 and 33) bind similar protein substrates, or at least,
503 substrates that share a binding motif. At the same time, it seems that all substrates binding these
504 domains share a binding motif. These findings further support the hypothesis presented by
505 MacGowan, Madeira (18), which states that UMDs are determinant for substrate binding
506 specificity.

507 Conclusions

508 The multiple sequence alignment of homologues and the aggregation of genetic variants, or other
509 features, over alignment columns, as described in MacGowan, Madeira (18), can provide insight
510 at the residue level on the evolutionary constraint acting on functional domains as well as
511 highlight structural or functionally relevant residues in protein domains. Overall, a clear variation
512 distribution pattern can be observed within the ankyrin repeat motif. There are five positions that

513 are conserved and depleted in missense variation due to their structural importance, e.g.,
514 enrichment in intra-repeat contacts. Four other positions are highly variable within the family and
515 overall depleted in missense variants, for they are the key for a specific and successful substrate
516 binding.

517 In this study, we use 7,407 ankyrin repeat sequences, 21,338 human missense variants and
518 160 3D structures to study the distribution of missense variants within the ankyrin repeat motif
519 and explain the observed patterns with structural data. The general conclusions are as follows.

- 520 1. Two of the turns found on the secondary structure of the ANK, positions 28-29 and 33-
521 1, are Asx motifs. Positions 27 and 32 present conserved Asx.
- 522 2. The surface of the ARD can be divided in three different surfaces using the RSA of the
523 repeat positions.
- 524 3. Positions that are conserved and depleted in missense variants (CMD) are significantly
525 enriched in intra-repeat contacts ($OR = 2.8, p \approx 0$) and are key for the correct packing
526 of the motif as well as the domain.
- 527 4. Positions that are unconserved yet depleted in missense variants (UMD) are heavily
528 enriched in protein-protein interactions ($OR = 3.6, p < 10^{-16}$) and might be responsible
529 for substrate binding specificity in the motif.
- 530 5. The concave surface of the ARD is significantly enriched in PPIs ($PPIES = 1.86, p <$
531 10^{-16}) and consequently depleted in missense variation ($MES = -0.08, p = 4.4 \times$
532 10^{-4}) whereas the other two surfaces are less constrained in line with their reduced
533 importance in substrate binding.

534 Acknowledgements

535 We thank Drs Marek Gierlinski, Matthew Parker and Jim Procter for their insightful suggestions
536 during this research. We also thank Prof. Ulrich Zachariae for critical reading of the manuscript
537 and the IT service of the University of Dundee for their support of our HPC infrastructure. This
538 work was supported by Biotechnology and Biological Sciences Research Council Grants
539 (BB/J019364/1 and BB/R014752/1) and Wellcome Trust Biomedical Resources Grant
540 (101651/Z/13/Z).
541

542 References

- 543 1. Bork P. Hundreds of ankyrin-like repeats in functionally diverse proteins: mobile
544 modules that cross phyla horizontally? *Proteins*. 1993;17(4):363-74.
- 545 2. Andrade MA, Perez-Iratxeta C, Ponting CP. Protein repeats: structures, functions, and
546 evolution. *J Struct Biol*. 2001;134(2-3):117-31.
- 547 3. Sedgwick SG, Smerdon SJ. The ankyrin repeat: a diversity of interactions on a common
548 structural framework. *Trends Biochem Sci*. 1999;24(8):311-6.
- 549 4. Forrer P, Stumpp MT, Binz HK, Pluckthun A. A novel strategy to design binding
550 molecules harnessing the modular nature of repeat proteins. *FEBS Lett*. 2003;539(1-3):2-6.
- 551 5. Crooks GE, Hon G, Chandonia JM, Brenner SE. WebLogo: a sequence logo generator.
552 *Genome Res*. 2004;14(6):1188-90.
- 553 6. Krzywda S, Brzozowski AM, Higashitsuji H, Fujita J, Welchman R, Dawson S, et al. The
554 crystal structure of gankyrin, an oncoprotein found in complexes with cyclin-dependent kinase
555 4, a 19 S proteasomal ATPase regulator, and the tumor suppressors Rb and p53. *J Biol Chem*.
556 2004;279(2):1541-5.
- 557 7. Pettersen EF, Goddard TD, Huang CC, Couch GS, Greenblatt DM, Meng EC, et al. UCSF
558 Chimera--a visualization system for exploratory research and analysis. *J Comput Chem*.
559 2004;25(13):1605-12.

- 560 8. Mosavi LK, Minor DL, Jr., Peng ZY. Consensus-derived structural determinants of the
561 ankyrin repeat motif. *Proc Natl Acad Sci U S A*. 2002;99(25):16029-34.
- 562 9. Islam Z, Nagampalli RSK, Fatima MT, Ashraf GM. New paradigm in ankyrin repeats:
563 Beyond protein-protein interaction module. *Int J Biol Macromol*. 2018;109:1164-73.
- 564 10. Stumpp MT, Binz HK, Amstutz P. DARPins: a new generation of protein therapeutics.
565 *Drug Discov Today*. 2008;13(15-16):695-701.
- 566 11. Main ER, Jackson SE, Regan L. The folding and design of repeat proteins: reaching a
567 consensus. *Curr Opin Struct Biol*. 2003;13(4):482-9.
- 568 12. Kohl A, Binz HK, Forrer P, Stumpp MT, Pluckthun A, Grutter MG. Designed to be stable:
569 crystal structure of a consensus ankyrin repeat protein. *Proc Natl Acad Sci U S A*.
570 2003;100(4):1700-5.
- 571 13. Li J, Mahajan A, Tsai MD. Ankyrin repeat: a unique motif mediating protein-protein
572 interactions. *Biochemistry*. 2006;45(51):15168-78.
- 573 14. Hansen S, Stuber JC, Ernst P, Koch A, Bojar D, Batyuk A, et al. Design and applications
574 of a clamp for Green Fluorescent Protein with picomolar affinity. *Sci Rep*. 2017;7(1):16292.
- 575 15. Petrovski S, Wang Q, Heinzen EL, Allen AS, Goldstein DB. Genic intolerance to
576 functional variation and the interpretation of personal genomes. *PLoS Genet*.
577 2013;9(8):e1003709.
- 578 16. Gussow AB, Petrovski S, Wang Q, Allen AS, Goldstein DB. The intolerance to functional
579 genetic variation of protein domains predicts the localization of pathogenic mutations within
580 genes. *Genome Biol*. 2016;17:9.
- 581 17. Sivley RM, Dou X, Meiler J, Bush WS, Capra JA. Comprehensive Analysis of Constraint
582 on the Spatial Distribution of Missense Variants in Human Protein Structures. *Am J Hum Genet*.
583 2018;102(3):415-26.
- 584 18. MacGowan SA, Madeira F, Britto-Borges T, Schmittner MS, Cole C, Barton GJ. Human
585 Missense Variation is Constrained by Domain Structure and Highlights Functional and
586 Pathogenic Residues. *bioRxiv*. 2017:127050.
- 587 19. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alföldi J, Wang Q, et al. The
588 mutational constraint spectrum quantified from variation in 141,456 humans. *bioRxiv*.
589 2020:531210.
- 590 20. UniProt C. UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Res*.
591 2019;47(D1):D506-D15.
- 592 21. Schultz J, Milpetz F, Bork P, Ponting CP. SMART, a simple modular architecture
593 research tool: identification of signaling domains. *Proc Natl Acad Sci U S A*. 1998;95(11):5857-
594 64.
- 595 22. Sigrist CJ, de Castro E, Cerutti L, Cuče BA, Hulo N, Bridge A, et al. New and continuing
596 developments at PROSITE. *Nucleic Acids Res*. 2013;41(Database issue):D344-7.
- 597 23. Attwood TK. The PRINTS database: a resource for identification of protein families.
598 *Brief Bioinform*. 2002;3(3):252-63.
- 599 24. El-Gebali S, Mistry J, Bateman A, Eddy SR, Luciani A, Potter SC, et al. The Pfam protein
600 families database in 2019. *Nucleic Acids Res*. 2019;47(D1):D427-D32.
- 601 25. Sievers F, Higgins DG. Clustal Omega, accurate alignment of very large numbers of
602 sequences. *Methods Mol Biol*. 2014;1079:105-16.
- 603 26. Eddy SR. Multiple alignment using hidden Markov models. *Proc Int Conf Intell Syst Mol*
604 *Biol*. 1995;3:114-20.
- 605 27. Notredame C, Higgins DG, Heringa J. T-Coffee: A novel method for fast and accurate
606 multiple sequence alignment. *J Mol Biol*. 2000;302(1):205-17.
- 607 28. Barton GJ. The AMPS package for multiple protein sequence alignment. *Methods Mol*
608 *Biol*. 1994;25:327-47.
- 609 29. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high
610 throughput. *Nucleic Acids Res*. 2004;32(5):1792-7.

- 611 30. Russell RB, Barton GJ. Multiple protein sequence alignment from tertiary structure
612 comparison: assignment of global and residue confidence levels. *Proteins*. 1992;14(2):309-23.
- 613 31. Jeanmougin F, Thompson JD, Gouy M, Higgins DG, Gibson TJ. Multiple sequence
614 alignment with Clustal X. *Trends Biochem Sci*. 1998;23(10):403-5.
- 615 32. Waterhouse AM, Procter JB, Martin DM, Clamp M, Barton GJ. Jalview Version 2--a
616 multiple sequence alignment editor and analysis workbench. *Bioinformatics*. 2009;25(9):1189-
617 91.
- 618 33. Berman H, Henrick K, Nakamura H. Announcing the worldwide Protein Data Bank. *Nat*
619 *Struct Biol*. 2003;10(12):980.
- 620 34. Berman H, Henrick K, Nakamura H, Markley JL. The worldwide Protein Data Bank
621 (wwPDB): ensuring a single, uniform archive of PDB data. *Nucleic Acids Res*. 2007;35(Database
622 issue):D301-3.
- 623 35. wwPDBconsortium. Protein Data Bank: the single global archive for 3D
624 macromolecular structure data. *Nucleic Acids Res*. 2019;47(D1):D520-D8.
- 625 36. Velankar S, Dana JM, Jacobsen J, van Ginkel G, Gane PJ, Luo J, et al. SIFTS: Structure
626 Integration with Function, Taxonomy and Sequences resource. *Nucleic Acids Res*.
627 2013;41(Database issue):D483-9.
- 628 37. MacGowan SA, Madeira F, Britto-Borges T, Warowny M, Drozdetskiy A, Procter JB, et
629 al. The Dundee Resource for Sequence Analysis and Structure Prediction. *Protein Sci*.
630 2020;29(1):277-97.
- 631 38. Kleywegt GJ, Harris MR, Zou JY, Taylor TC, Wahlby A, Jones TA. The Uppsala Electron-
632 Density Server. *Acta Crystallogr D Biol Crystallogr*. 2004;60(Pt 12 Pt 1):2240-9.
- 633 39. Kabsch W, Sander C. Dictionary of protein secondary structure: pattern recognition of
634 hydrogen-bonded and geometrical features. *Biopolymers*. 1983;22(12):2577-637.
- 635 40. Shenkin PS, Erman B, Mastrandrea LD. Information-theoretical entropy as a measure
636 of sequence variability. *Proteins*. 1991;11(4):297-313.
- 637 41. Szumilas M. Explaining odds ratios. *J Can Acad Child Adolesc Psychiatry*.
638 2010;19(3):227-9.
- 639 42. Jubb HC, Higuieruelo AP, Ochoa-Montano B, Pitt WR, Ascher DB, Blundell TL. Arpeggio:
640 A Web Server for Calculating and Visualising Interatomic Interactions in Protein Structures. *J*
641 *Mol Biol*. 2017;429(3):365-71.
- 642 43. Tien MZ, Meyer AG, Sydykova DK, Spielman SJ, Wilke CO. Maximum allowed solvent
643 accessibilities of residues in proteins. *Plos One*. 2013;8(11):e80635.
- 644 44. Miller S, Janin J, Lesk AM, Chothia C. Interior and surface of monomeric proteins. *J Mol*
645 *Biol*. 1987;196(3):641-56.
- 646 45. de Brevern AG. Extension of the classical classification of beta-turns. *Sci Rep*.
647 2016;6:33191.
- 648 46. Wan WY, Milner-White EJ. A natural grouping of motifs with an aspartate or
649 asparagine residue forming two hydrogen bonds to residues ahead in sequence: their
650 occurrence at alpha-helical N termini and in other situations. *J Mol Biol*. 1999;286(5):1633-49.
- 651 47. Wang C, Wei Z, Chen K, Ye F, Yu C, Bennett V, et al. Structural basis of diverse
652 membrane target recognitions by ankyrins. *Elife*. 2014;3.
- 653 48. Michaely P, Tomchick DR, Machius M, Anderson RG. Crystal structure of a 12 ANK
654 repeat stack from human ankyrinR. *EMBO J*. 2002;21(23):6387-96.
- 655 49. Hesketh GG, Perez-Dorado I, Jackson LP, Wartosch L, Schafer IB, Gray SR, et al. VARP is
656 recruited on to endosomes by direct interaction with retromer, where together they function
657 in export to the cell surface. *Dev Cell*. 2014;29(5):591-606.
- 658 50. Xu C, Jin J, Bian C, Lam R, Tian R, Weist R, et al. Sequence-specific recognition of a
659 PxLPxL/L motif by an ankyrin repeat tumbler lock. *Sci Signal*. 2012;5(226):ra39.
- 660 51. Nie J, Xu C, Jin J, Aka JA, Tempel W, Nguyen V, et al. Ankyrin repeats of ANKRA2
661 recognize a PxLPxL motif on the 3M syndrome protein CCDC8. *Structure*. 2015;23(4):700-12.

- 662 52. Sbdio JI, Chi NW. Identification of a tankyrase-binding motif shared by IRAP, TAB182,
663 and human TRF1 but not mouse TRF1. NuMA contains this RXXPDG motif and is a novel
664 tankyrase partner. *J Biol Chem.* 2002;277(35):31887-92.
- 665 53. Fetics S, Thureau A, Campanacci V, Aumont-Nicaise M, Dang I, Gautreau A, et al.
666 Hybrid Structural Analysis of the Arp2/3 Regulator Arpin Identifies Its Acidic Tail as a Primary
667 Binding Epitope. *Structure.* 2016;24(2):252-60.
- 668 54. Xu W, Lau YH, Fischer G, Tan YS, Chattopadhyay A, de la Roche M, et al. Macrocyclized
669 Extended Peptides: Inhibiting the Substrate-Recognition Domain of Tankyrase. *J Am Chem Soc.*
670 2017;139(6):2245-56.
- 671