

1 A study of the structural properties of sites
2 modified by the *O*-linked 6-N-
3 acetylglucosamine transferase
4

5
6 **Short title:** Structural characterisation of *O*-GlcNAc sites

7 **keywords:** post-translational modification, computational biology, structural analysis, *O*-
8 GlcNAc, protein disorder
9

10 Thiago Britto-Borges and Geoffrey J. Barton*

11
12 Division of Computational Biology
13 School of Life Sciences
14 University of Dundee
15 Dow Street, Dundee, DD1 5EH
16 UK`
17

18
19
20 *To whom correspondence should be addressed: g.j.barton@dundee.ac.uk
21
22

23

24 **Abstract**

25 Protein O-GlcNAcylation (O-GlcNAc) is an essential post-translational modification (PTM) in
26 higher eukaryotes. The O-linked β -N-acetylglucosamine transferase (OGT), targets specific
27 Serines and Threonines (S/T) in intracellular proteins. However, unlike phosphorylation, fewer
28 than 25% of known O-GlcNAc sites match a clear sequence pattern. Accordingly, the three-
29 dimensional structures of O-GlcNAc sites were characterised to investigate the role of structure
30 in molecular recognition. Of the 143/1,584 O-GlcNAc sites in 620 proteins were mapped to
31 protein X-ray structures. The modified S/T were 1.7x more likely to be annotated in the REM465
32 field which defines missing residues in a protein structure, while 7 O-GlcNAc sites were solvent
33 inaccessible and unlikely to be targeted by OGT. The 132/143 sites with complete backbone
34 atoms clustered into 10 groups, but these were indistinguishable from clusters from unmodified
35 S/T. This suggests there is no prevalent three-dimensional motif for OGT recognition. Predicted
36 features from the 620 proteins were compared to unmodified S/T in O-GlcNAcylated proteins
37 and globular proteins. The Jpred4 predicted secondary structure shows that modified S/T were
38 more likely to be coils. 5/6 methods to predict intrinsic disorder indicated O-GlcNAcylated S/T
39 to be significantly more disordered than unmodified S/T. Although the analysis did not find a
40 pattern in the site three-dimensional structure, it revealed the residues around the modification
41 site are likely to be disordered and suggests a potential role of secondary structure elements in
42 OGT site recognition.

43

44 I. Introduction

45 Protein *O*-GlcNAcylation, or *O*-GlcNAc, is a dynamic, intracellular glycosylation essential to
46 mammalian development^{1,2}. In animals, two enzymes mediate this post-translational
47 modification: the glycosyltransferase O-linked N-acetylglucosamine transferase (OGT), which
48 adds a single, non-extensible *O*-GlcNAc moiety to serine/threonine (S/T) in the target protein;
49 and the hexosaminidase *O*-GlcNAcase (OGA) that removes it. UDP-GlcNAc, the sugar donor to
50 the protein *O*-GlcNAcylation, is a product of the hexosamine pathway, hence the concentration
51 of intracellular glucose and the degree of protein *O*-GlcNAcylation levels are associated^{3, 4}. At
52 the physiological level, dysfunction of OGT activity has been linked to disease of the
53 cardiovascular system, diabetes, impaired development, cancer and neurodegeneration⁵⁻⁹. At the
54 cellular level, protein *O*-GlcNAcylation acts with phosphorylation, ubiquitylation and other
55 reversible post-translational modifications in a network of cell signalling events that promote
56 cellular adaptation to the viral infection process¹⁰, regulation of transcription¹¹ and metabolism¹²,
57 ¹³.

58 Technical advances in mass spectrometry have led to an increase in the number of
59 experimentally determined *O*-GlcNAc sites from 50 in the year 2000 to more than 1,000 today¹⁴.
60 However, there are still obstacles to mapping *O*-GlcNAc sites reliably. The modification has a
61 low abundance¹⁵ and is ten times less common than protein phosphorylation¹⁶. Thus, the
62 unmodified version of the peptide can suppress the *O*-GlcNAcylated peptide mass/charge signal.
63 In addition, methods to enrich *O*-GlcNAcylated peptides in samples have limited specificity^{16, 17},
64 and the β -glycosidic bond is labile under the peptide fragmentation step which determines the
65 modification's position within the peptide fragment.

66 Two machine learning methods have been used to detect patterns in the sequence of *O*-GlcNAc
67 sites^{18, 19} with limited success²⁰. One of the limiting factors was that, unlike phosphorylation
68 sites, *O*-GlcNAc sites lack a clear pattern in the primary structure. This is illustrated in Figure 1
69 which compares the relative sequence entropy for sites modified by OGT and three protein
70 kinases in the PhosphoSitePlus database¹⁴. The observed relative entropy for OGT sites shows
71 small signal, in contrast to protein kinase A (PKA; peak in -3 and +2), protein kinase C (PKC;
72 peak in -3) and casein kinase 2 (CK2, peak in +3) sites. This implies that the sequence in the
73 sites recognised by OGT carries less information than those recognised by PKA, PKC or CK2
74 and so are harder to distinguish from unmodified sites by sequence alone.

75 The crystal structure of OGT in a ternary complex with UDP-GlcNAc and a peptide substrate
76 revealed that the OGT and the peptides' residues predominantly make contact via the peptide
77 backbone^{21, 22}. This fact reduces the importance of the peptide side chain in the enzyme active
78 site, the cleft where the reaction occurs. A short structural motif, instead of sequence motif,
79 could work as a point of molecular recognition even with a degenerate sequence. Accordingly, in
80 this paper, the three-dimensional structures of S/T OGT substrates were examined to determine if
81 they have distinct structural motifs and patterns of secondary structure or solvent accessibility.
82 In addition, the predicted secondary structure and disorder were compared for known OGT
83 substrates and S/T unlikely to be modified.

84 **II. Methods**

85 **A. Data sources**

86 The data selection process is summarised in Figure 2. A total of 1,533 modified sites from 676
87 proteins were selected by combining proteins curated from the literature up until 2011¹⁸ and from
88 2011-2013²⁰. The sites were filtered to keep 7-residue long motifs with unique sequences. The

89 resulting dataset contained 1,385 sites in 620 proteins. This dataset is referred to hereafter as the
90 “modified sequence sites” (MSS). For comparison, 100,329 S/T from the same proteins, but not
91 annotated as OGT-modified, were selected as a background and are referred to here as the
92 “unmodified sequence sites” (USS).

93

94 **B. Mapping O-GlcNAc sites to protein structures**

95 Protein chains > 30 residues long from structures determined by X-ray crystallography to \leq
96 2.50 Å resolution were selected from the Protein Data Bank²³ (PDB: 2nd August of 2015).
97 Mapping the 1,385 OGT sites from 620 proteins to PDB structures by SIFTS²⁴ located 45 OGT
98 sites in 24 proteins of known structure. The structures of a further 107 sites were identified by
99 searching the sequences of O-GlcNAcylated proteins against the PDB chains with BLAST and
100 filtering by a conservative E-value $\leq 10^{-25}$ to minimise erroneous matches. The cutoff of $\leq 10^{-25}$
101 was found empirically to ensure the reliability of the match in the region of each site by
102 inspecting all alignments between query and PDB sequence at different thresholds. Selecting the
103 protein chain with highest coverage (SIFTS) or E-value (BLAST) left 143 sites in 107 proteins
104 for further analysis, referred to hereafter as the “143 Structural Sites” (SS143).

105

106 **C. Site definition and clustering**

107 The three-dimensional structure of OGT with its substrates suggests the region of contact
108 between OGT and a modifiable S/T includes the residues and +/- 3 amino acids either side^{21, 22}.
109 From the structural sites returned in Mapping O-GlcNAc-sites to protein structures, “132
110 Structural Sites” (hereafter SS132) had at least one match with all backbone atoms for the 7-
111 residue long site and were retained for further analysis. C α atoms of each residue and the C α and

112 the C β for the central S/T were superimposed for all pairs of sites. Hierarchical clustering by
113 complete linkage was applied on the resulting matrix of root-mean-square deviation (RMSD)
114 values and clusters selected where all pairs of peptides were within 3 Å RMSD of each other.

115

116 **D. Structural properties of sites**

117 Protein secondary structure assignments were obtained from DSSP²⁵. DSSP annotates 7 different
118 secondary structure states: 3_{10} helix (G), α helix (H), π helix (I), bends (S), turns (T), isolated (B)
119 and extended (E) β -bridge. These assignments were reduced to three states: G and H to helices
120 (H); I, B and E to strands (E); and all other, including residues with no assignment, to coils (C)²⁶.
121 The solvent accessible area from DSSP was normalised by the residue's maximum accessible
122 area²⁷. A S/T was considered exposed if its relative solvent accessibility (RSA) was > 25%;
123 partially buried if the RSA > 5% and \leq 25%, and buried if RSA \leq 5%. C α B-factors were
124 standardised (Z-score normalised) over the B-factors for all C α in the same chain.

125

126 **E. Prediction of protein disorder and secondary structure**

127 Protein secondary structure predictions for the proteins in the MSS dataset were performed by
128 JPred4²⁸. Since JPred4 limits sequence longer than 800 residues, 300 of the sequences in the
129 MSS dataset sequences were trimmed while ensuring the modified S/T was at least 100 residues
130 away from the N- and C-termini to avoid edge effects. The intrinsic disorder was predicted by
131 JRonn (Java implementation of Ronn²⁹), IUPred³⁰ and DisEMBL³¹ through the JABAWS³²
132 command line application. Between them, these methods provide 6 different disorder prediction
133 scores: DisEMBL-REM465 (0.6), DisEMBL-COILS (0.516), DisEMBL-HOTLOOPS (0.1204),
134 IUPred-Long (0.5), IUPred-Short (0.5) and JRonn (0.5). The score ordered/disordered classes

135 were defined by the cut-offs (in parenthesis) defined by the methods' authors. Disorder
136 predictions were also performed on a background set of 1,164 S/T selected at random from
137 globular proteins in the Astral dataset³³ version 2.04, referred to hereafter as the “Globular Set”
138 (GS).

139

140 **F. Statistical analysis and code**

141 The data collection, processing, analysis and the Ca clustering steps, were written in the Python
142 programming language (Python Software Foundation, version 2.7 <http://www.python.org>) with
143 the libraries Pandas (version 0.17)³⁴ and Biopython (version 1.65)³⁵. Statistical tests were
144 performed with the StatsModels (version 0.6) and Scipy (version 0.16) libraries. A *p* value (*p*)
145 threshold was set to 0.05.

146

147 **III. Results and Discussion**

148 **A. Analysis of O-GlcNAc sites in proteins of known three-dimensional** 149 **structure**

150 Previous reports have suggested that O-GlcNAc sites, like phosphorylation sites, are
151 predominantly present in disordered regions of proteins³⁶. One indication of structural disorder is
152 the crystallographic B-factor which indicates regions of the protein that lack crystallographic
153 contacts. However, the standardised B-factor distribution on the SS143 dataset is the same for
154 modified and unmodified S/T (Kruskal-Wallis two-sample test $p = 0.12$).

155 In X-ray crystal structures, the REM465 residue annotation indicates residues that are missing
156 from the protein structure model and has previously been used as an indicator of structural

157 disorder³¹. Of the 143 S/T in the SS143 dataset, 26 are in regions of the protein structure labelled
158 as REM465. In comparison, 553 of 4,811 unmodified S/T from the same protein structures are
159 also found in REM465 regions. Accordingly, O-GlcNAcylated S/T in these proteins are 1.7
160 times more likely to be in REM465 regions (Fisher's exact test $p = 0.02$). This finding is
161 consistent with O-GlcNAcylated S/T occurring more frequently in disordered or highly flexible
162 regions.

163 Table 2 summarises the DSSP assigned secondary structure for the SS143 compared to the 4,811
164 unmodified S/T in the same proteins. The proportions of H, E and C are equivalent for the two
165 groups implying that there is no preference in the secondary structure for modified S/T in this
166 dataset.

167 Residues that are buried in the protein structure are not thought to be targeted by protein kinases,
168 due to structural constraints. Figure 3 illustrates that there is no difference between modified and
169 unmodified S/T with respect to relative solvent accessibility (RSA). 45% (65) of S/T in the O-
170 GlcNAcylated proteins are exposed to solvent (RSA > 25%). Surprisingly, 7 O-GlcNAc sites,
171 listed in Table 2, have an RSA < 5%, suggesting they are inaccessible to OGT in the natively
172 folded protein.

173

174 **B. Groups of sites with similar local structure**

175 Since the secondary structure and relative accessibility of modified S/T were indistinguishable
176 from unmodified S/T, the local structure of the 7 residue peptides centred on S/T was
177 investigated by pairwise superposition and clustering (see Methods). 36 sites produce singlet
178 clusters, while the remaining 96 sites fall into 10 clusters. Sites in clusters had less than 3 Å
179 RMSD from each other. Figure 4 illustrates the superimposed structures for sites in clusters,

180 where green, yellow and grey represent residues in H, E, C secondary structures, respectively.
181 The clusters show that sites are found in a wide range of secondary structure states as
182 summarised in Appendix 1. The sites in Clusters E, G and J, have consistent consensus
183 secondary structures. Clusters A–D, F, H and I are all variants on coil-helix or coil-strand
184 transitions.
185 The buried sites, which are listed in Table 3, group in clusters D and G. The 3 sites in cluster D
186 are unlikely to be targeted by OGT because they are buried in the protein core. In contrast, the
187 2/4 sites in cluster G (structures 3abm and 4y7y) might be modified since are located at a dimer
188 interface, and so the monomer could be modified. The remaining two sites in cluster G
189 (structures 2zxe and 4l3j) lie on a loop that could potentially move to expose them to OGT.
190 To see if the clusters found for the SS132 dataset are features of *O*-GlcNAc modification or just
191 reflect the composition of the protein structures, 132 sites, centred on unmodified S/T, were
192 randomly sampled with replacement from the same proteins and clustered. The process was
193 repeated 1,000 times and the resulting clusters compared to those clusters in the SS132 dataset.
194 The number of clusters identified in each sample ranged from 10-14 (95% CI), which is
195 consistent with the SS132 dataset. Furthermore, the structural clusters identified for the random
196 sampling included structural clusters similar to those for the modified sites, suggesting there are
197 no dominant secondary structural or conformational patterns indicative of *O*-GlcNAc modified
198 sites in the SS132 dataset.
199

200 **C. Analysis of features predicted for the “Modified Sequence Sites”** 201 **dataset (MSS)**

202 Since the structural analysis of *O*-GlcNAc sites is limited by the number of sites in proteins of
203 known three-dimensional structure, prediction algorithms were applied to the sequences in the
204 MSS and USS datasets, as detailed in Methods. The proportions of S/T in the levels of solvent
205 accessibility predicted by JPred are equivalent in the MSS and USS datasets, as shown in
206 Table 4. 1% of the S/T are predicted to be buried in the MSS and USS datasets. Again, the result
207 is unexpected, since sites modified by PTM are thought to be accessible in the protein native
208 fold.

209 While the structural sites in the SS143 dataset have equal proportions of the secondary structure
210 states, the result from secondary structure predictions on the MSS set showed that *O*-GlcNAc
211 sites are likely to reside in coils, if compared to the USS dataset.

212 Table 5 shows an increase of the proportion of modified S/T in C ($p < 0.01$) and a corresponding
213 reduction in H ($p < 0.01$), but no change in E ($p = 0.6$). The enrichment of sites in C is consistent
214 with the need to place modified S/T in loops that are more likely to be mobile and so more
215 accessible to OGT. The proportions of secondary structure assigned by DSSP and predicted by
216 JPred4 differ. While secondary structure prediction has limited accuracy, the number of samples
217 in the SS143 dataset is limited and potentially biased toward structured regions in proteins. Also,
218 clustering sites in the SS132 dataset highlight groups that are more likely to occur near to the
219 transition between a secondary structure element and C, as observed in several members of
220 clusters A–D, F and H. The regions of transition between C and H/E are harder to predict than
221 contiguous secondary structure elements, and this may also contribute to the observed
222 enrichment in C.

223 The analysis of SS143 dataset showed an enrichment of S/T in REM465 regions likely to be
224 disordered or highly mobile. To explore this further, 3 disorder prediction algorithms, giving a
225 total of 6 disorder scores, were run on the MSS and USS datasets as detailed in Methods. Table 6
226 shows that, with the exception of DisEMBL-HOTLOOPS which is trained structural B-factors,
227 all methods report a small but significant increase in mean predicted disorder for the modified
228 S/T. To confirm this result, the MSS dataset was compared to the GS dataset, which was selected
229 from proteins known to be predominantly globular, and hence an ordered background. In Figure
230 5, DisEMBL-HOTLOOPS shows an increase in the ratio of disordered residues around the
231 modified S/T. DisEMBL-COILS and JRonn also indicate a small increase, not in a specific
232 region, but rather for 40 residues around the S/T. IUPred-Long, IUPred-Short and DisEMBL-
233 REM465 show a bigger increase of the ratio of disordered residues in the MSS dataset and
234 IUPred-Short and REM465 have a clearer peak within -15 to 15 residues from the modified S/T.
235 Overall, all methods indicate an increased proportion of predicted disorder in the MSS dataset
236 when compared to the GS dataset.

237 **IV. Conclusions and final remarks**

238 Despite the substantial evidence of protein structural disorder in the MSS and the SS143
239 datasets, the SS132 dataset clearly indicates that some of the examined sites appear within
240 ordered regions of the protein structure. Furthermore, InterproScan³⁷ analysis of *O*-GlcNAc sites
241 assigned 19 % of the sites to protein domains, this is similar to with the 25 % phosphoserines and
242 phosphothreonines in PFAM domains^{14, 38}, which are thought to be mostly ordered by definition.
243 So, like protein phosphorylation, *O*-GlcNAcylated S/T are found in both ordered and disordered
244 regions.

245 The local tertiary structure of *O*-GlcNAc sites is indistinguishable from unmodified sites. So,
246 how does OGT recognise the site it modifies? OGT may force the unfolding of the targeted
247 substrate³⁹. Moreover, OGT participates in macromolecular assemblies⁴⁰, and so adaptor proteins
248 should be important. Non-local interactions might also act in OGT substrate recognition. In
249 protein kinase C (PKC) substrate recognition, residues distant in the protein sequence but close
250 in its three-dimensional structure are critical⁴¹. Other components, such as UDP-GlcNAc
251 concentration and subcellular location-dependent interactions, modulate OGT activity⁴², but their
252 part in substrate recognition is still unknown. In conclusion, although no three-dimensional
253 fingerprint was detected during the structural characterisation of OGT-modified sites, the work
254 confirmed that S/T and surrounding residues are more disordered than the backgrounds tested
255 and that sites in transition between C to H/E might be involved, suggesting that the structural
256 flexibility has a role on OGT site recognition.

257 **Acknowledgements**

258 We would like to thank Dr. Tom Walsh and the University of Dundee IT department for
259 computing support; Prof. Daan van Aalten and DVA group for advice and discussions.

260 **Funding**

261 This work was supported by Coordenação de Aperfeiçoamento de Pessoal de Nível Superior
262 (CAPES process 1529/12-9; studentship to T.B.B).

263 **Conflict of Interest:**

264 None declared.

265

266

267 **References**

- 268 1. Shafi R, Iyer SP, Ellies LG, O'Donnell N, Marek KW, Chui D, Hart GW, Marth JD. The
269 O-GlcNAc transferase gene resides on the X chromosome and is essential for embryonic
270 stem cell viability and mouse ontogeny. *Proc Natl Acad Sci U S A* 2000;97(11):5735–
271 5739.
- 272 2. O'Donnell N, Zachara NE, Hart GW, Marth JD. Ogt-Dependent X-Chromosome-Linked
273 Protein Glycosylation Is a Requisite Modification in Somatic Cell Function and Embryo
274 Viability. *Mol Cell Biol* 2004;24(4):1680–1690.
- 275 3. Buse MG. Hexosamines, insulin resistance, and the complications of diabetes: current
276 status. *Am J Physiol Endocrinol Metab* 2006;290(1):E1–E8.
- 277 4. Abdel Rahman AM, Ryczko M, Pawling J, Dennis JW. Probing the hexosamine
278 biosynthetic pathway in human tumor cells by multitargeted tandem mass spectrometry.
279 *ACS Chem Biol* 2013;8(9):2053–2062.
- 280 5. Liu J, Marchase RB, Chatham JC. Increased O-GlcNAc levels during reperfusion lead to
281 improved functional recovery and reduced calpain proteolysis. *Am J Physiol Heart Circ*
282 *Physiol* 2007;293(3):H1391-9.
- 283 6. McClain DA, Lubas WA, Cooksey RC, Hazel M, Parker GJ, Love DC, Hanover JA.
284 Altered glycan-dependent signaling induces insulin resistance and hyperleptinemia. *Proc*
285 *Natl Acad Sci U S A* 2002;99(16):10695–10699.
- 286 7. Mariappa D, Zheng X, Schimpl M, Raimi O, Ferenbach AT, Müller H-AJ, Aalten DMF
287 van. Dual functionality of O -GlcNAc transferase is required for *Drosophila* development.
288 *Open Biol* 2015;5(12):150234.
- 289 8. Lynch TP, Ferrer CM, Jackson SR, Shahriari KS, Vosseller K, Reginato MJ. Critical Role
290 of O-Linked -N-Acetylglucosamine Transferase in Prostate Cancer Invasion,
291 Angiogenesis, and Metastasis. *J Biol Chem* 2012;287(14):11070–11081.
- 292 9. Liu F, Iqbal K, Grundke-Iqbal I, Hart GW, Gong C-X. O-GlcNAcylation regulates
293 phosphorylation of tau: A mechanism involved in Alzheimer's disease. *Proc Natl Acad*
294 *Sci* 2004;101(29):10804–10809.
- 295 10. Chen D, Juárez S, Hartweck L, Alamillo JM, Simón-Mateo C, Pérez JJ, Fernández-
296 Fernández MR, Olszewski NE, García JA. Identification of secret agent as the O-GlcNAc
297 transferase that participates in Plum pox virus infection. *J Virol* 2005;79(15):9381–9387.
- 298 11. Kuo M, Zilberfarb V, Gangneux N, Christeff N, Issad T. O-glycosylation of FoxO1
299 increases its transcriptional activity towards the glucose 6-phosphatase gene. *FEBS Lett*
300 2008;582(5):829–834.
- 301 12. Wells L. Mapping Sites of O-GlcNAc Modification Using Affinity Tags for Serine and
302 Threonine Post-translational Modifications. *Mol Cell Proteomics* 2002;1(10):791–804.
- 303 13. Parker GJ, Lund KC, Taylor RP, McClain DA. Insulin resistance of glycogen synthase
304 mediated by o-linked N-acetylglucosamine. *J Biol Chem* 2003;278(12):10022–10027.
- 305 14. Hornbeck P V, Zhang B, Murray B, Kornhauser JM, Latham V, Skrzypek E.
306 PhosphoSitePlus, 2014: mutations, PTMs and recalibrations. *Nucleic Acids Res*
307 2015;43(Database issue):D512-20.
- 308 15. Roquemore EP, Dell A, Morris HR, Panico M, Reason AJ, Savoy LA, Wistow GJ, Zigler
309 JS, Earles BJ, Hart GW. Vertebrate lens alpha-crystallins are modified by O-linked N-
310 acetylglucosamine. *J Biol Chem* 1992;267(1):555–563.
- 311 16. Hahne H, Gholami AM, Kuster B, Moghaddas Gholami A, Kuster B. Discovery of O-

- 312 GlcNAc-modified proteins in published large-scale proteome data. *Mol Cell Proteomics*
313 2012;11(10):843–850.
- 314 17. Ma J, Hart GW. O-GlcNAc profiling: from proteins to proteomes. *Clin Proteomics*
315 2014;11(1):8.
- 316 18. Wang J, Torii M, Liu H, Hart GW, Hu Z. dbOGAP - An Integrated Bioinformatics
317 Resource for Protein O-GlcNAcylation. *BMC Bioinformatics* 2011;12(1):91.
- 318 19. Jia C-Z, Liu T, Wang Z-P. O-GlcNAcPRED: a sensitive predictor to capture protein O-
319 GlcNAcylation sites. *Mol Biosyst* 2013;9(11):2909–2913.
- 320 20. Jochmann R, Holz P, Sticht H, Stürzl M. Validation of the reliability of computational O-
321 GlcNAc prediction. *Biochim Biophys Acta - Proteins Proteomics* 2014;1844(2):416–421.
- 322 21. Lazarus MB, Nam Y, Jiang J, Sliz P, Walker S. Structure of human O-GlcNAc transferase
323 and its complex with a peptide substrate. In: *Nature*. Volume 469. Nature Publishing
324 Group; 2011. p 564–567.
- 325 22. Schimpl M, Zheng X, Borodkin VS, Blair DE, Ferenbach AT, Schüttelkopf AW,
326 Navratilova I, Aristotelous T, Albarbarawi O, Robinson D a, Macnaughtan M a, Aalten
327 DMF van. O-GlcNAc transferase invokes nucleotide sugar pyrophosphate participation in
328 catalysis. *Nat Chem Biol* 2012;8(12):969–974.
- 329 23. Velankar S, Alhroub Y, Alili A, Best C, Boutselakis HC, Caboche S, Conroy MJ, Dana
330 JM, Ginkel G van, Golovin A, Gore SP, Gutmanas A, Haslam P, Hirshberg M, John M,
331 Lagerstedt I, Mir S, Newman LE, Oldfield TJ, Penkett CJ, Pineda-Castillo J, Rinaldi L,
332 Sahni G, Sawka G, Sen S, Slowley R, Sousa da Silva AW, Suarez-Uruena A,
333 Swaminathan GJ, Symmons MF, Vranken WF, Wainwright M, Kleywegt GJ. PDBe:
334 Protein Data Bank in Europe. *Nucleic Acids Res* 2011;39(Database issue):D402-10.
- 335 24. Velankar S, Dana JM, Jacobsen J, Ginkel G van, Gane PJ, Luo J, Oldfield TJ, O'Donovan
336 C, Martin M-J, Kleywegt GJ. SIFTS: Structure Integration with Function, Taxonomy and
337 Sequences resource. *Nucleic Acids Res* 2013;41(Database issue):D483-9.
- 338 25. Kabsch W, Sander C. Dictionary of protein secondary structure: pattern recognition of
339 hydrogen-bonded and geometrical features. *Biopolymers* 1983;22(12):2577–2637.
- 340 26. Cuff J a, Barton GJ. Evaluation and improvement of multiple sequence methods for
341 protein secondary structure prediction. In: *Proteins*. Volume 34. 1999. p 508–519.
- 342 27. Cuff JA, Barton GJ. Application of multiple sequence alignment profiles to improve
343 protein secondary structure prediction. In: *Proteins*. Volume 40. John Wiley & Sons, Inc.;
344 2000. p 502–511.
- 345 28. Drozdetskiy A, Cole C, Procter J, Barton GJ. JPred4: A protein secondary structure
346 prediction server. *Nucleic Acids Res* 2015;43(W1):W389-94.
- 347 29. Yang ZR, Thomson R, McNeil P, Esnouf RM. RONN: the bio-basis function neural
348 network technique applied to the detection of natively disordered regions in proteins.
349 *Bioinformatics* 2005;21(16):3369–3376.
- 350 30. Dosztányi Z, Csizmók V, Tompa P, Simon I. The pairwise energy content estimated from
351 amino acid composition discriminates between folded and intrinsically unstructured
352 proteins. *J Mol Biol* 2005;347(4):827–839.
- 353 31. Linding R, Jensen LJ, Diella F, Bork P, Gibson TJ, Russell RB. Protein Disorder
354 Prediction. *Structure* 2003;11(11):1453–1459.
- 355 32. Troshin P V, Procter JB, Barton GJ. Java bioinformatics analysis web services for
356 multiple sequence alignment--JABAWS:MSA. *Bioinformatics* 2011;27(14):2001–2002.
- 357 33. Fox NK, Brenner SE, Chandonia J-M. SCOPe: Structural Classification of Proteins--

- 358 extended, integrating SCOP and ASTRAL data and classification of new structures.
359 Nucleic Acids Res 2014;42(Database issue):D304-9.
- 360 34. McKinney W, Team PD. Pandas - Powerful Python Data Analysis Toolkit. In: Pandas -
361 Powerful Python Data Analysis Toolkit. 2015. p 1625.
- 362 35. Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, Friedberg I, Hamelryck
363 T, Kauff F, Wilczynski B, Hoon MJL de. Biopython: freely available Python tools for
364 computational molecular biology and bioinformatics. In: Bioinformatics. Volume 25.
365 2009. p 1422–1423.
- 366 36. Trinidad JC, Barkan DT, Gulledge BF, Thalhammer a., Sali a., Schoepfer R, Burlingame
367 a. L. Global Identification and Characterization of Both O-GlcNAcylation and
368 Phosphorylation at the Murine Synapse. Mol Cell Proteomics 2012;11(8):215–229.
- 369 37. Zdobnov EM, Apweiler R. InterProScan - an integration platform for the signature-
370 recognition methods in InterPro. Bioinformatics 2001;17(9):847–848.
- 371 38. Beltrao P, Albanèse V, Kenner LR, Swaney DL, Burlingame A, Villén J, Lim WA, Fraser
372 JS, Frydman J, Krogan NJ. Systematic Functional Prioritization of Protein
373 Posttranslational Modifications. Cell 2012;150(2):413–425.
- 374 39. Pathak S, Alonso J, Schimpl M, Rafie K, Blair DE, Borodkin VS, Schüttelkopf AW,
375 Albarbarawi O, Aalten DMF van. The active site of O-GlcNAc transferase imposes
376 constraints on substrate sequence. Nat Struct Mol Biol 2015;22(9):744–750.
- 377 40. Wells L, Kreppel LK, Comer FI, Wadzinski BE, Hart GW. O-GlcNAc Transferase Is in a
378 Functional Complex with Protein Phosphatase 1 Catalytic Subunits. J Biol Chem
379 2004;279(37):38466–38470.
- 380 41. Duarte ML, Pena DA, Nunes Ferraz FA, Berti DA, Paschoal Sobreira TJ, Costa-Junior
381 HM, Abdel Baqui MM, Disatnik MM-H, Xavier-neto J, Lopes de Oliveira PS,
382 Schechtman D, Augusto F, Ferraz N, Berti DA, José T, Sobreira P, Costa-Junior HM,
383 Muhammad M, Baqui A, Disatnik MM-H, Xavier-neto J. Protein folding creates
384 structure-based, noncontiguous consensus phosphorylation motifs recognized by kinases.
385 Sci Signal 2014;7(350):ra105-ra105.
- 386 42. Nagel AK, Ball LE. O-GlcNAc transferase and O-GlcNAcase: achieving target substrate
387 specificity. Amino Acids 2014;46(10):2305–2316.
- 388 43. Crooks GE, Hon G, Chandonia J-M, Brenner SE. WebLogo: a sequence logo generator.
389 Genome Res 2004;14(6):1188–1190.
- 390
391

392 Tables

393 Table 1 Dataset summary. See Methods for details.

394

Dataset name	Number of sites	Number of proteins	Short name
Modified Sequence Sites	1,385	620	MSS
Unmodified Sequence Sites	100,329	620	USS
Structural Sites	143	106	SS143
Structural Sites with backbone	132	93	SS132
Globular Set	1,164	1,164	GS

395

396 Table 2 DSSP assigned secondary structure proportion of S/T in the SS143 dataset compared to
 397 unmodified S/T in same protein chains. 95% CI – 95% confidence interval; n – number of S/T.
 398 The *p* value refers to the two-tailed z-score test between the proportions of modified and
 399 unmodified groups.

400

Secondary structure	Modified		Unmodified		<i>p</i> value
	Proportion (n)	95% CI [lower, upper]	Proportion (n)	95% CI [lower, upper]	
C	0.55 (78)	[0.46, 0.63]	0.51 (2475)	[0.50, 0.53]	0.36
H	0.25 (36)	[0.18, 0.32]	0.32 (1525)	[0.31, 0.33]	0.06
E	0.20 (29)	[0.13, 0.27]	0.17 (811)	[0.16, 0.18]	0.27
Total	143		4,811		

401

402

403 Table 3 Structural evidence of buried *O*-GlcNAc sites in the SS143 dataset. RSA – site mean
 404 relative solvent accessibility; Cluster id – Clusters in Figure 4.

405

PDB id	Chain	Position	Cluster id	RSA
1f4j	B	114	D	0.05
3cb2	B	170	D	0.02
4qvp	T	131	D	0.01
2zxe	A	366	G	0.02
3abm	R	63	G	0.01
4l3j	A	180	G	0.01
4y7y	Z	190	G	0.04

406
 407 Table 4 JPred4 predicted solvent accessibility for S/T in the MSS and USS datasets. The
 408 proportions of buried S/T as predicted by the Jnetsol method in JPred4. The proportions of
 409 buried S/T are significantly smaller for modified group. 95% CI – 95% confidence interval; n –
 410 number of S/T predicted to be buried. The *p* value refers to the two-tailed z-score test between
 411 the modified and unmodified groups.
 412

Buried at	Modified (MSS)		Unmodified (USS)		<i>p</i> value
	Proportion (n)	95% CI [lower, upper]	Proportion (n)	95% CI [lower, upper]	
0%	0.01 (7)	[0.00, 0.01]	0.01 (836)	[0.008, 0.009]	0.18
5%	0.04 (55)	[0.03, 0.05]	0.04 (3,917)	[0.038, 0.040]	0.86
25%	0.29 (403)	[0.27, 0.31]	0.35 (28,044)	[0.27, 0.28]	0.31

413
 414
 415 Table 5 JPred4 predicted secondary structure proportions for S/T in the MSS and USS datasets.
 416 95% CI – 95% confidence interval; n – the number of S/T; the *p* value refers to the two-tailed z-
 417 score test between the modified and unmodified groups.
 418

Secondary structure	Modified (MSS)		Unmodified (USS)		<i>p</i> value
	Proportion (n)	95% CI [lower, upper]	Proportion (n)	95% CI [lower, upper]	
C	0.88 (1,205)	[0.86, 0.90]	0.829 (83,150)	[0.826, 0.831]	<0.01
H	0.08 (106)	[0.07, 0.09]	0.126 (12,684)	[0.124, 0.128]	<0.01
E	0.05 (66)	[0.04, 0.06]	0.045 (4,495)	[0.044, 0.046]	0.6

419
 420
 421
 422 Table 6 Predicted disorder between modified and unmodified S/T. All disorder prediction
 423 methods, excepting DisEMBL-HOTLOOPS, reveal a small but significant increase of mean
 424 disorder score for modified S/T over unmodified ones. The *p* value refers to the two-tailed t-test
 425 between the modified and unmodified groups. SE – standard error.
 426

Method	Mean score modified (MSS) ± SE	Mean score unmodified (USS) ± SE	<i>p</i> value
DisEMBL-REM465	0.48 ± 0.004	0.47 ± 0.001	0.01
DisEMBL-COILS	0.60 ± 0.004	0.58 ± 0.001	<0.01

DisEMBL- HOTLOOPS	0.10 ± 0.001	0.10 ± 0.001	0.45
IUpred-Long	0.59 ± 0.006	0.55 ± 0.001	<0.01
IUpred-Short	0.48 ± 0.005	0.45 ± 0.001	<0.01
JRonn	0.62 ± 0.004	0.61 ± 0.001	0.02

427

428

429

430 **Legend to figures**

431 **Figure 1:** Sequence relative entropy of sites (+/- 7 residues) from 4 post-translational
432 modifications. Three kinases with most sites in PhosphoSitePlus database¹⁴: protein kinase A
433 (PKA with 1285 sites), protein kinase C (PKC with 930 sites) and casein kinase 2 (CK2 with 742
434 sites). 1530 OGT sites were compiled from the same database. The sequence relative entropy
435 was calculated with the WebLogo library⁴³. Lines show mean relative entropy and the semi-
436 transparent area represents 95% confidence intervals. Unlikely the PKA, PKC and CK2, known
437 OGT sites have no clear sequence consensus.

438 **Figure 2:** Diagram of the relationships among the 5 datasets used in this work. See Methods for
439 details.

440 **Figure 3:** RSA of modified S/T in the SS143 dataset and unmodified S/T in same proteins.
441 DSSP calculated solvent accessibility was normalised by the residue theoretical maximum
442 accessibility, and the derived scores were reduced to three levels: buried ($RSA \leq 0.05$), partially
443 buried ($0.05 < RSA \leq 0.25$) and exposed ($RSA > 0.25$). The y-axis and x-axis carry the RSA
444 levels and the RSA distribution for each level, respectively. The mean RSA is equivalent
445 between modified and unmodified residues, at all three levels.

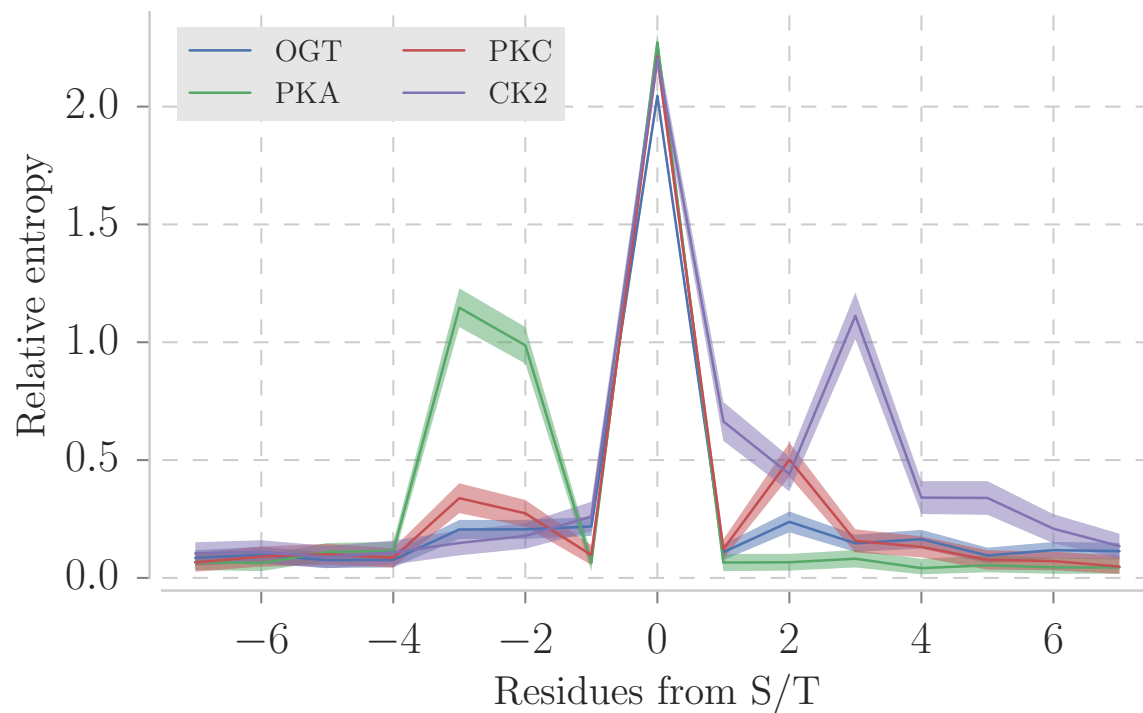
446 **Figure 4:** Structural superimpositions for the 10 clusters comprising 96 sites in the SS132
447 dataset. Pairs of sites were superimposed on their 7 C α atoms and the C β of the central S/T.
448 Their pairwise RMSD were clustered with complete linkage and Euclidean distance. Clusters
449 were defined by a 3 Å RMSD threshold. Green, yellow and grey represent residues in H, E, C
450 secondary structures elements, respectively.

451 **Figure 5:** Predicted disorder around O-GlcNAc sites in the MSS compared randomly selected
452 S/T in the GS-dataset. The y-axis shows the log₁₀ odds ratio of the between the proportion of

453 disordered residues in the MSS dataset and the proportion of disordered residues in the GS
454 dataset. The semi-transparent area represents 95% confidence intervals. A residue was defined as
455 disordered given a method threshold. The x-axis represents the distance in residues of the central
456 residue, always a S/T. DisEMBL-REM465, IUpred-short predict protein structural disorder
457 specifically around the modification site, while the other methods detect intrinsic disorder over
458 O-GlcNAcylated proteins. DisEMBL-HOTLOOPS shows a less pronounced increase compared
459 to the other methods.
460

461 **Figures**

462 Figure 1



463

464

465

466

467

468

469

470

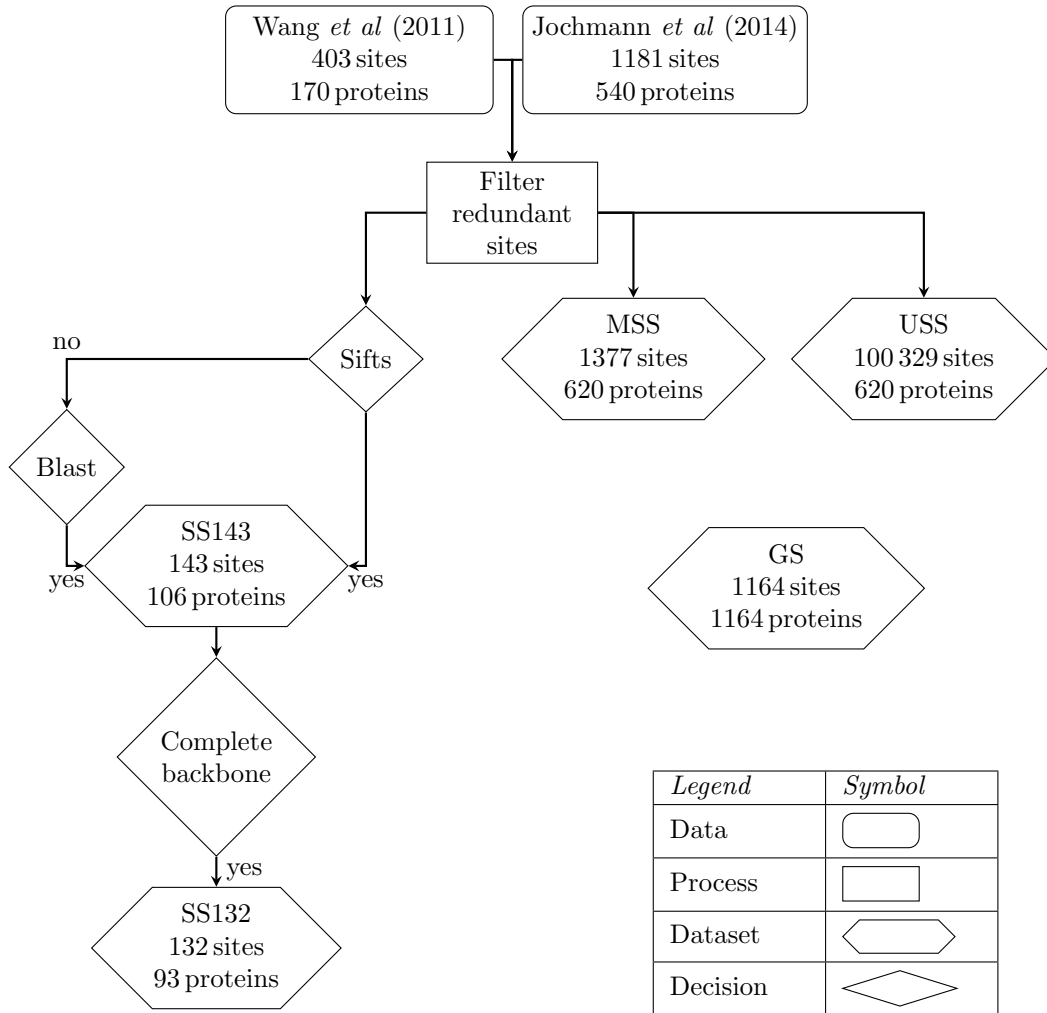
471

472

473

474

475 Figure 2

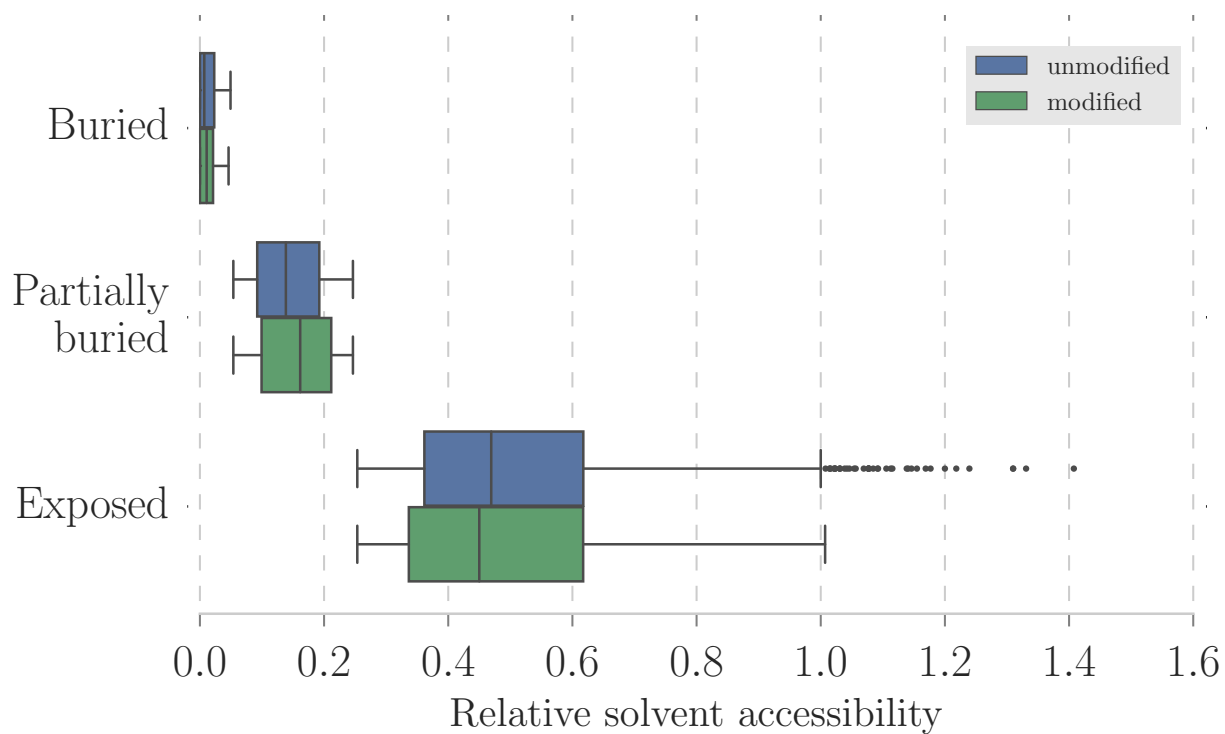


<i>Legend</i>	<i>Symbol</i>
Data	
Process	
Dataset	
Decision	

476

477

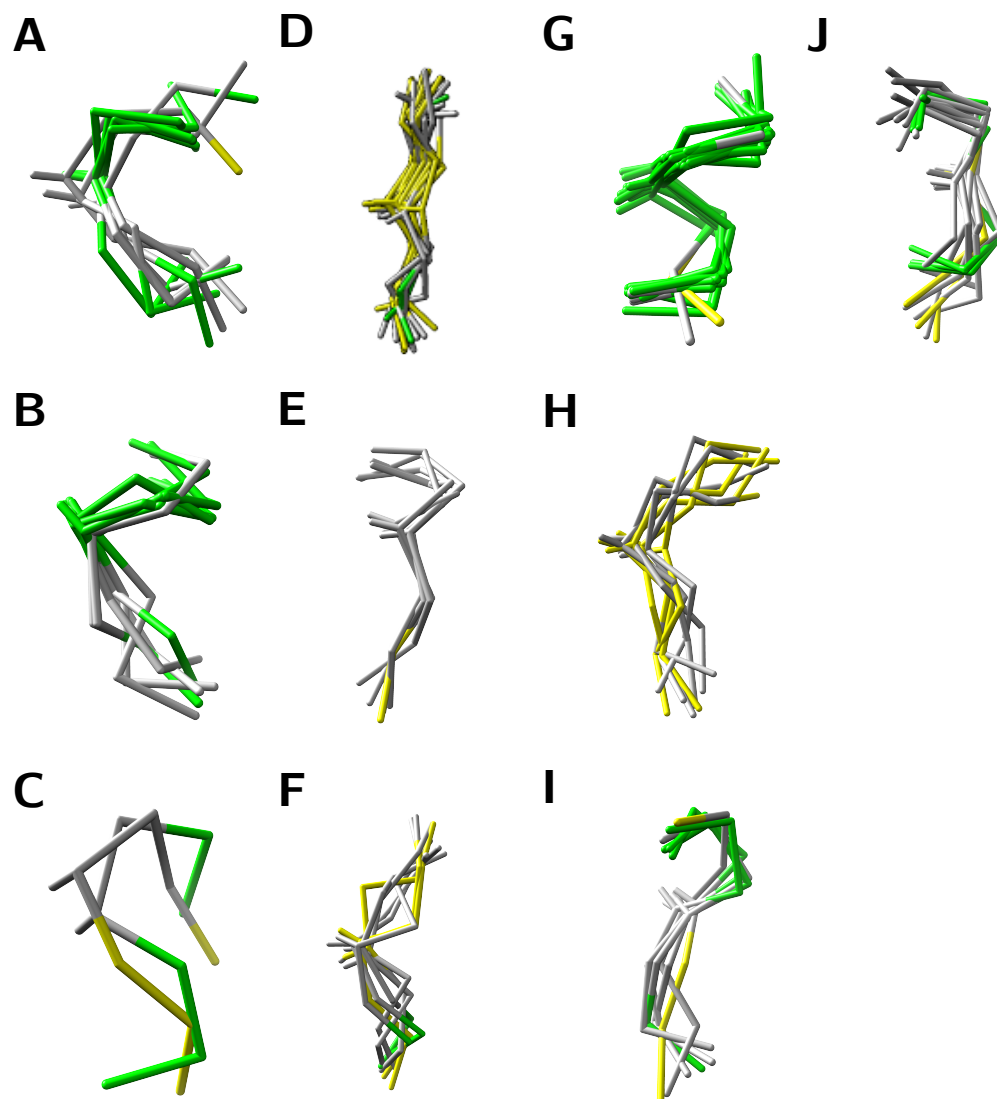
478 Figure 3



479

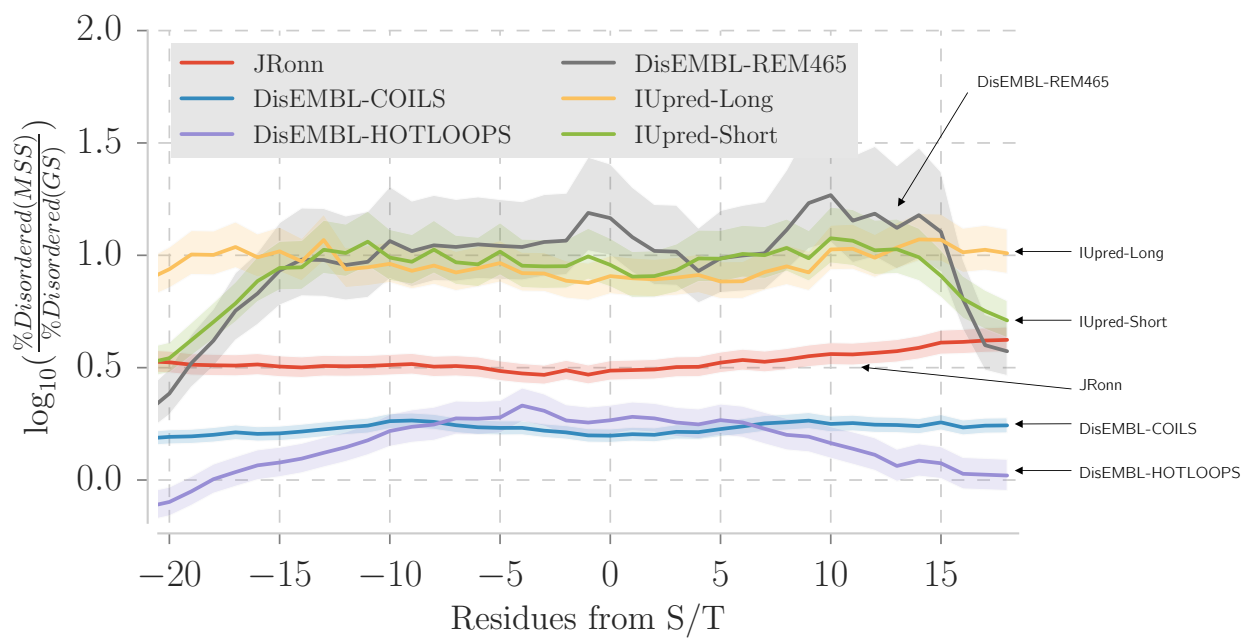
480

481 Figure 4



482

483 Figure 5



484