

RESEARCH ARTICLE

# A study of the structural properties of sites modified by the O-linked 6-N-acetylglucosamine transferase

Thiago Britto-Borges, Geoffrey J. Barton\*

Division of Computational Biology, School of Life Sciences, University of Dundee, Dundee, United Kingdom

\* [g.j.barton@dundee.ac.uk](mailto:g.j.barton@dundee.ac.uk)



## Abstract

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

## OPEN ACCESS

**Citation:** Britto-Borges T, Barton GJ (2017) A study of the structural properties of sites modified by the O-linked 6-N-acetylglucosamine transferase. PLoS ONE 12(9): e0184405. <https://doi.org/10.1371/journal.pone.0184405>

**Editor:** Iddo Friedberg, Iowa State University College of Veterinary Medicine, UNITED STATES

**Received:** March 21, 2017

**Accepted:** August 23, 2017

**Published:** September 8, 2017

**Copyright:** © 2017 Britto-Borges, Barton. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** All relevant data are within the paper and its Supporting Information files. Protein structure files were obtained from Protein Data Bank. Protein structure accession and site mapping are listed in S1 Table (<https://doi.org/10.6084/m9.figshare.4910141.v2>).

**Funding:** This work was supported by Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES process 1529/12-9; studentship to T.B.B). Website: <http://www.capes.gov.br/>. GJB Acknowledges the support of Wellcome Trust Strategic Awards: WT097945,

## Introduction

Protein O-GlcNAcylation, or O-GlcNAc, is a dynamic, intracellular glycosylation essential to mammalian development [1,2]. In animals, two enzymes mediate this post-translational modification: the glycosyltransferase O-linked 6-N-acetylglucosamine transferase (OGT), which adds a single, non-extensible O-GlcNAc moiety to serine/threonine (S/T) in the target protein; and the hexosaminidase O-GlcNAcase (OGA) that removes it. UDP-GlcNAc, the sugar donor to the protein O-GlcNAcylation, is a product of the hexosamine pathway, hence the concentration of intracellular glucose and the degree of protein O-GlcNAcylation levels are associated [3,4]. At the physiological level, dysfunction of OGT activity has been linked to disease of the

WT092340 and 098439/Z/12/Z. Website <https://wellcome.ac.uk/>. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing interests:** The authors have declared that no competing interests exist.

cardiovascular system, diabetes, impaired development, cancer and neurodegeneration [5–9]. At the cellular level, protein O-GlcNAcylation acts with phosphorylation, ubiquitylation and other reversible post-translational modifications in a network of cell signalling events that promote cellular adaptation to the viral infection process [10], regulation of transcription [11] and metabolism [12,13].

Technical advances in mass spectrometry have led to an increase in the number of experimentally determined O-GlcNAc sites from 50 in the year 2000 to more than 1,000 today [14]. However, there are still obstacles to mapping O-GlcNAc sites reliably. The modification has a low abundance [15] and is ten times less common than protein phosphorylation [16]. Thus, the unmodified version of the peptide can suppress the O-GlcNAcylated peptide mass/charge signal. In addition, methods to enrich O-GlcNAcylated peptides in samples have limited specificity [16,17], and the  $\beta$ -glycosidic bond is labile under the peptide fragmentation step which determines the modification's position within the peptide fragment.

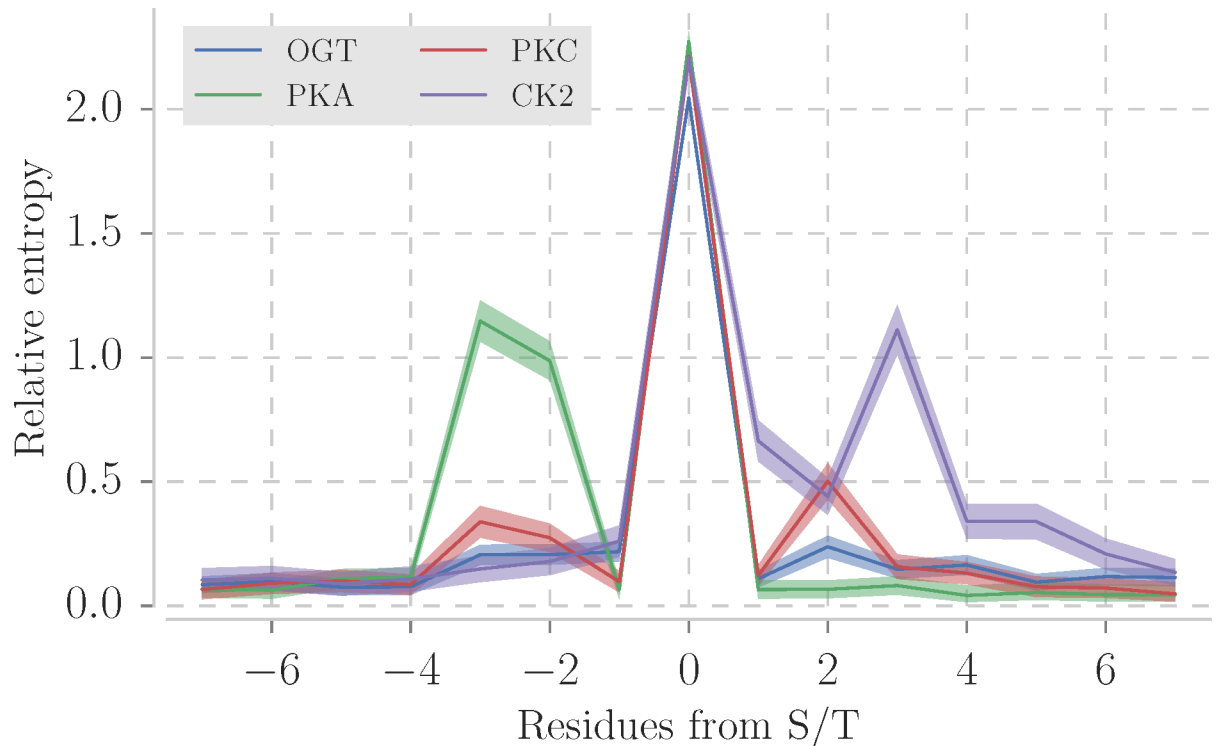
Two machine learning methods have been used to detect patterns in the sequence of O-GlcNAc sites [18–20] with limited success [21]. Newer predictors have exploited more complex machine learning approaches to classify potential novel sites [22–24] but to date have only been applied in a few studies. One of the limiting factors for site prediction is that, unlike phosphorylation sites, O-GlcNAc sites lack a clear pattern in the primary structure. This is illustrated in Fig 1 which compares the relative sequence entropy for sites modified by OGT and three protein kinases in the PhosphoSitePlus database [14]. The relative sequence entropy, calculated with the WebLogo library [25], describes the amount of information carried per position compared to the background amino acid distribution. OGT sites shows no peaks other than the modified S/T, in contrast to protein kinase A (PKA; peak in -3 and +2), protein kinase C (PKC; peak in -3) and casein kinase 2 (CK2, peak in +3) sites. This implies that the sequence in the sites recognised by OGT carries less information than those recognised by PKA, PKC or CK2 and so are harder to distinguish from unmodified sites by sequence alone.

OGT activity measured on peptide libraries demonstrate the enzyme substrate specificity and that point mutations near to the targeted S/T abolish peptide modification [26–28]. The crystal structure of OGT in a ternary complex with UDP-GlcNAc and a peptide substrate revealed that the OGT and the peptides' residues predominantly make contact via the peptide backbone [29,30]. This fact reduces the importance of the peptide side chain in the enzyme active site, the cleft where the reaction occurs. A short structural motif, instead of sequence motif, could work as a point of molecular recognition even with a degenerate sequence. Accordingly, in this paper, the three-dimensional structures of S/T OGT substrates were examined to determine if they have distinct structural motifs and patterns of secondary structure or solvent accessibility. In addition, the predicted secondary structure and disorder were compared for known OGT substrates and S/T unlikely to be modified.

## Methods

### Data sources

The data selection process is summarised in Table 1 and Fig 2. A total of 1,533 modified sites from 676 proteins were selected by combining proteins curated from the literature up until 2011 [18] and from 2011–2013 [21]. The majority of the sites were obtained from high-throughput experiments in mammals. The sites were filtered to keep 7-residue long motifs with unique sequences. The resulting dataset contained 1,385 sites in 620 proteins. This dataset is referred to hereafter as the “modified sequence sites” (MSS). For comparison, 100,329 S/T from the same proteins, but not annotated as OGT-modified, were selected as a background and are referred to here as the “unmodified sequence sites” (USS). S1 Table [<https://doi.org/>



**Fig 1. Sequence relative entropy of sites (+/- 7 residues) from 4 posttranslational modifications.** Three kinases with most sites in PhosphoSitePlus database [14] protein kinase A (PKA with 1285 sites), protein kinase C (PKC with 930 sites) and casein kinase 2 (CK2 with 742 sites). 1530 OGT sites were compiled from the same database. The sequence relative entropy was calculated with the WebLogo library [25]. Lines show mean relative entropy and the semi-transparent area represents 95% confidence intervals.

<https://doi.org/10.1371/journal.pone.0184405.g001>

[10.6084/m9.figshare.4910141.v2](https://doi.org/10.6084/m9.figshare.4910141.v2)] includes PDB accession code, chain identifier, position and summary data for each protein structure used in this study.

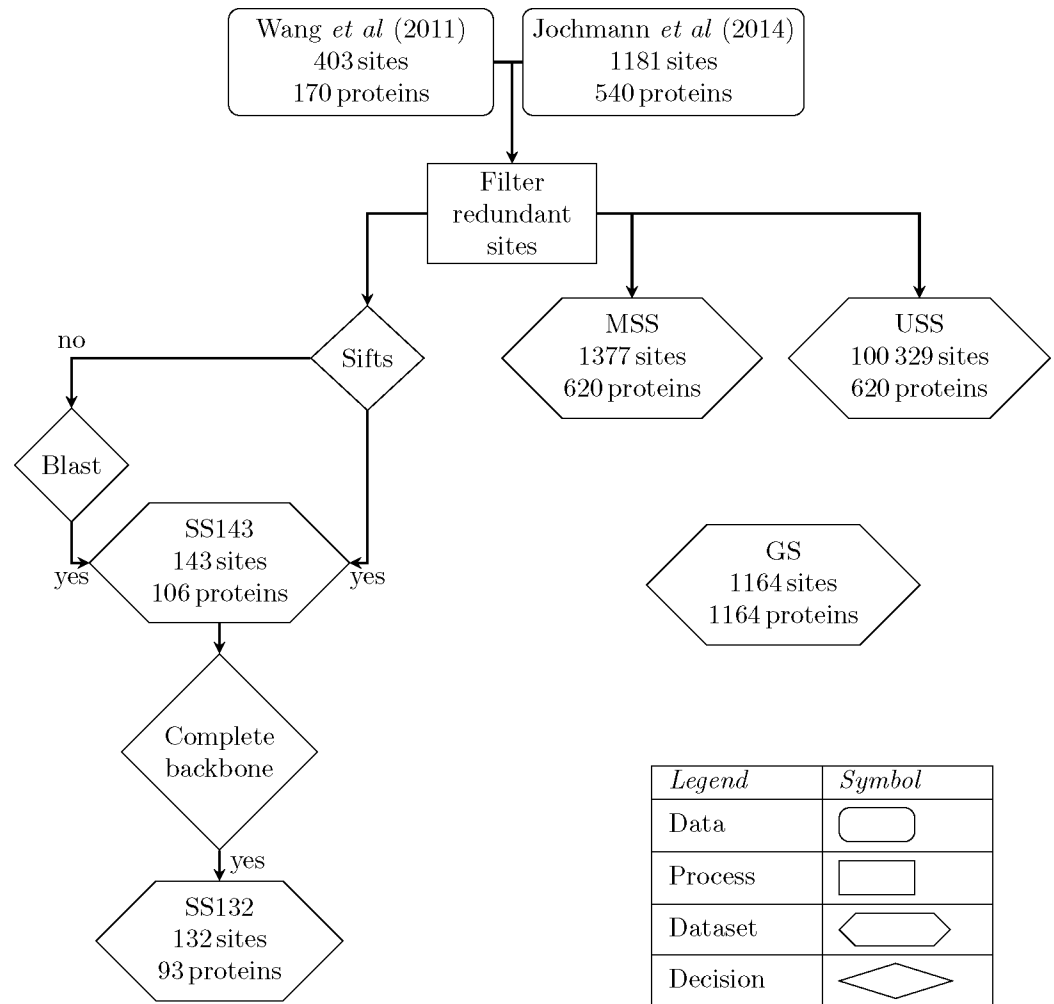
### Mapping O-GlcNAc sites to protein structures

Protein chains > 30 residues long from structures determined by X-ray crystallography to ≤ 2.50 Å resolution were selected from the Protein Data Bank [31] (PDB: 2<sup>nd</sup> August of 2015). Mapping the 1,385 OGT sites from 620 proteins to PDB structures by SIFTS [32] located 45 OGT sites in 24 proteins of known structure. The structures of a further 107 sites were identified by searching the sequences of O-GlcNAcylated proteins against the PDB chains with BLAST and filtering by a conservative E-value ≤ 10<sup>-25</sup> to minimise erroneous matches. The cutoff of ≤ 10<sup>-25</sup> was found empirically to ensure the reliability of the match in the region of each site by inspecting all alignments between query and PDB sequence at different thresholds.

**Table 1. Dataset summary.** See [Methods](#) for details.

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

<https://doi.org/10.1371/journal.pone.0184405.t001>



**Fig 2. Diagram of the relationships of among the 5 datasets used in this work.** See [Methods](#) for details.

<https://doi.org/10.1371/journal.pone.0184405.g002>

Selecting the protein chain with highest coverage (SIFTS) or E-value (BLAST) left 143 sites in 107 proteins for further analysis, referred to hereafter as the “143 Structural Sites” (SS143).

### Site definition and clustering

The three-dimensional structure of OGT with its substrates suggests the region of contact between OGT and a modifiable S/T includes the residues and +/- 3 amino acids either side [29,30]. From the structural sites returned in **Mapping O-GlcNAc-sites to protein structures**, “132 Structural Sites” (hereafter SS132) had at least one match with all backbone atoms for the 7-residue long site and were retained for further analysis. C $\alpha$  atoms of each residue and the C $\alpha$  and the C $\beta$  for the central S/T were superimposed for all pairs of sites. Hierarchical clustering by complete linkage was applied on the resulting matrix of root-mean-square deviation (RMSD) values and clusters selected where all pairs of peptides were within 3 Å RMSD of each other.

### Structural properties of sites

Protein secondary structure assignments were obtained from DSSP[33]. DSSP annotates 7 different secondary structure states:  $3_{10}$  helix (G),  $\alpha$  helix (H),  $\pi$  helix (I), bends (S), turns (T),

isolated (B) and extended (E)  $\beta$ -bridge. These assignments were reduced to three states: G and H to helices (H); I, B and E to strands (E); and all other, including residues with no assignment, to coils (C) [34]. The solvent accessible area from DSSP was normalised by the residue's maximum accessible area [35]. A S/T was considered exposed if its relative solvent accessibility (RSA) was  $> 25\%$ ; partially buried if the RSA  $> 5\%$  and  $\leq 25\%$ , and buried if RSA  $\leq 5\%$ .  $C\alpha$  B-factors were standardised (Z-score normalised) over the B-factors for all  $C\alpha$  in the same chain.

## Prediction of protein disorder and secondary structure

Protein secondary structure predictions for the proteins in the MSS dataset were performed by JPred4 [36]. Since JPred4 limits sequence longer than 800 residues, 300 of the sequences in the MSS dataset sequences were trimmed while ensuring the modified S/T was at least 100 residues away from the N- and C-termini to avoid edge effects. The intrinsic disorder was predicted by JRonn (Java implementation of Ronn [37]), IUPred [38] and DisEMBL [39] through the JABAWS [40] command line application. Between them, these methods provide 6 different disorder prediction scores: DisEMBL-REM465 (0.6), DisEMBL-COILS (0.516), DisEMBL-HOTLOOPS (0.1204), IUPred-Long (0.5), IUPred-Short (0.5) and JRonn (0.5). The score ordered/disordered classes were defined by the cut-offs (in parenthesis) defined by the methods' authors. Disorder predictions were also performed on a background set of 1,164 S/T selected at random from globular proteins in the Astral dataset [41] version 2.04, referred to hereafter as the "Globular Set" (GS).

## Statistical analysis and code

The data collection, processing, analysis and the  $C\alpha$  clustering steps, were written in the Python programming language (Python Software Foundation, version 2.7 <http://www.python.org>) with the libraries Pandas (version 0.17) [42] and Biopython (version 1.65) [43]. Statistical tests were performed with the StatsModels (version 0.6) and Scipy (version 0.16) libraries. A  $p$  value ( $p$ ) threshold was set to 0.05.

## Results and discussion

### Analysis of O-GlcNAc sites in proteins of known three-dimensional structure

Previous reports have suggested that O-GlcNAc sites, like phosphorylation sites, are predominantly present in disordered regions of proteins [44]. One indication of structural disorder is the crystallographic B-factor which indicates regions of the protein that lack crystallographic contacts. However, the standardised B-factor distribution on the SS143 dataset is the same for modified and unmodified S/T (Kruskal-Wallis two-sample test  $p = 0.12$ ).

In X-ray crystal structures, the REM465 residue annotation indicates residues that are missing from the protein structure model and has previously been used as an indicator of structural disorder [39]. Of the 143 S/T in the SS143 dataset, 26 are in regions of the protein structure labelled as REM465. In comparison, 553 of 4,811 unmodified S/T from the same protein structures are also found in REM465 regions. Accordingly, O-GlcNAcylated S/T in these proteins are 1.7 times more likely to be in REM465 regions (Fisher's exact test  $p = 0.02$ ). This finding is consistent with O-GlcNAcylated S/T occurring more frequently in disordered or highly flexible regions.

[Table 2](#) summarises the DSSP assigned secondary structure for the SS143 compared to the 4,811 unmodified S/T in the same proteins. The proportions of H, E and C are equivalent for

**Table 2. DSSP assigned secondary structure proportion of S/T in the SS143 dataset compared to unmodified S/T in same protein chains.**

Secondary structure	Modified		Unmodified		p 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		

95% CI— 95% confidence interval; n—number of S/T.

The p value refers to the two-tailed z-score test between the proportions of modified and unmodified groups.

<https://doi.org/10.1371/journal.pone.0184405.t002>

the two groups implying that there is no preference in the secondary structure for modified S/T in this dataset.

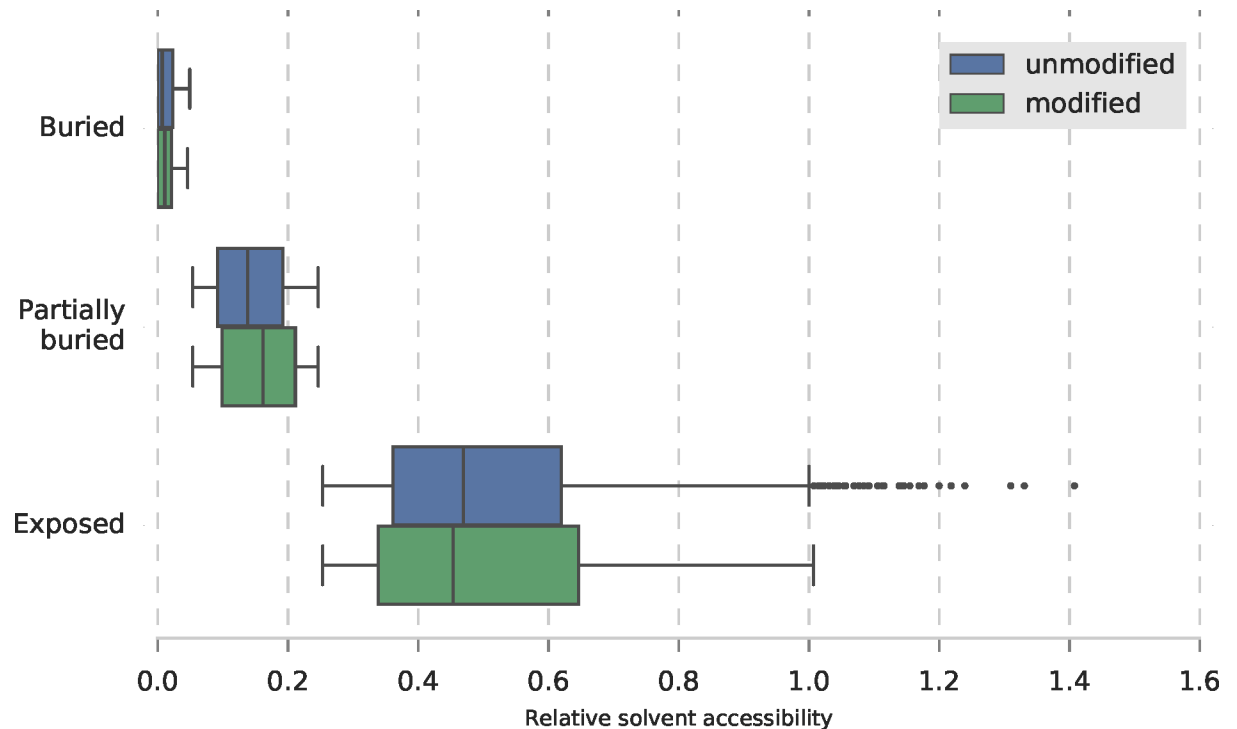
Residues that are buried in the protein structure are not thought to be targeted by protein kinases, due to structural constraints. Fig 3 illustrates that there is no difference between modified and unmodified S/T with respect to relative solvent accessibility (RSA). 45% of 65 S/T in the O-GlcNAcylated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

### Groups of sites with similar local structure

Since the secondary structure and relative accessibility of modified S/T were indistinguishable from unmodified S/T, the local structure of the 7 residue peptides centred on S/T was investigated by pairwise superposition and clustering (see Methods). 36 sites produce singlet clusters, where the majority of the residues are in C, while the remaining 96 sites fall into 10 clusters. Sites in clusters had less than 3 Å RMSD from each other. Fig 4 illustrates the superimposed structures for sites in clusters, where green, yellow and grey represent residues in H, E, C secondary structures, respectively. The clusters show that sites are found in a wide range of secondary structure states as summarised in S1 Table. The sites in Clusters E, G and J, have consistent consensus secondary structures. Clusters A–D, F, H and I are all variants on coil-helix or coil-strand transitions.

The buried sites, which are listed in Table 3, group in clusters D and G. The 3 sites in cluster D are unlikely to be targeted by OGT because they are buried in the protein core. In contrast, the 2/4 sites in cluster G (structures 3abm and 4y7y) might be modified since are located at a dimer interface, and so the monomer could be modified. The remaining two sites in cluster G (structures 2zxe and 4l3j) lie on a loop that could potentially move to expose them to OGT.

To see if the clusters found for the SS132 dataset are features of O-GlcNAc modification or just reflect the composition of the protein structures, 132 sites, centred on unmodified S/T, were randomly sampled with replacement from the same proteins and clustered. The process was repeated 1,000 times and the resulting clusters compared to those clusters in the SS132 dataset. The number of clusters identified in each sample ranged from 10–14 (95% CI), which is consistent with the SS132 dataset. Furthermore, the structural clusters identified for the random sampling included structural clusters similar to those for the modified sites, suggesting there are no dominant secondary structural or conformational patterns indicative of O-GlcNAc modified sites in the SS132 dataset. The analysis was also extended longer peptides with 20 residues either side of the modified S/T, but the structural clustering showed high heterogeneity for 41-residue peptides and no clear patterns were identified.



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

<https://doi.org/10.1371/journal.pone.0184405.g003>

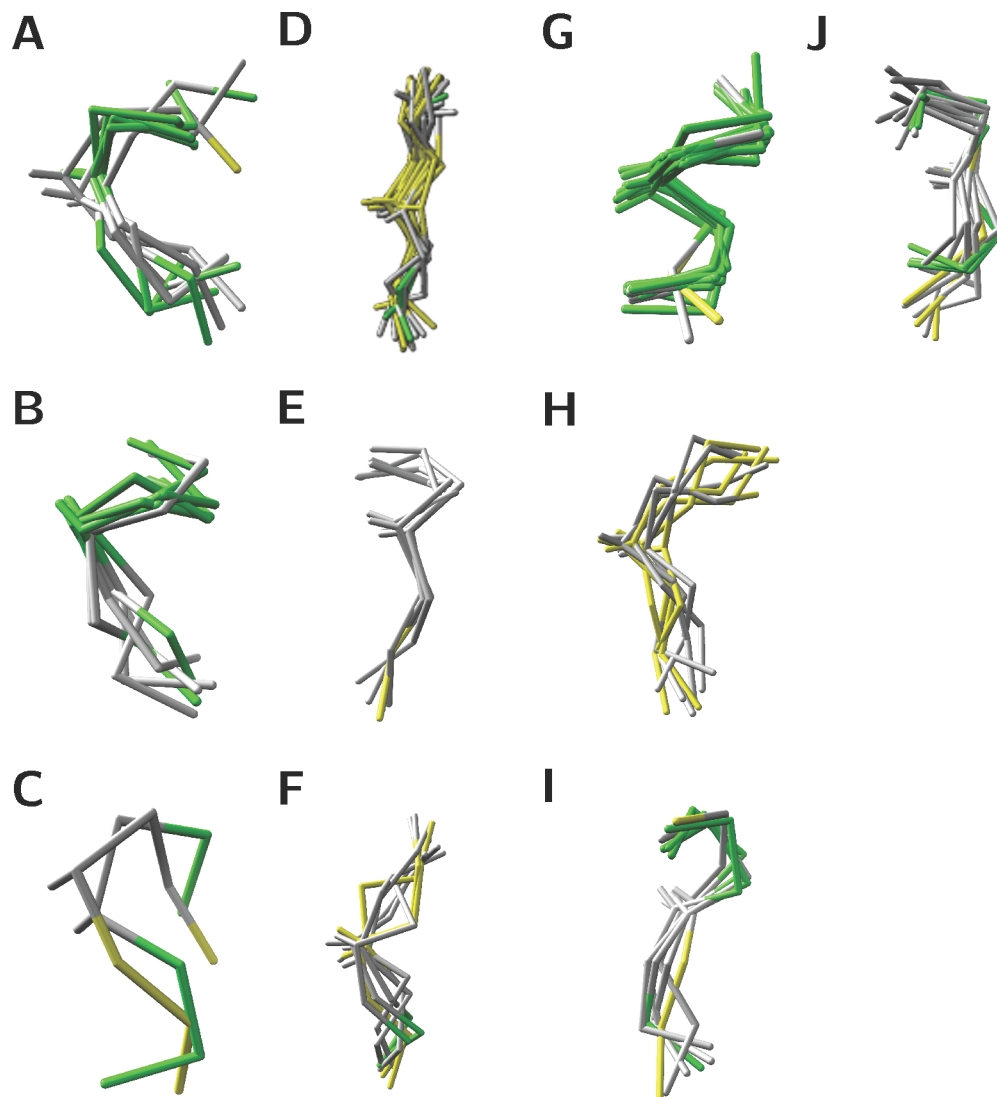
### Analysis of features predicted for the “modified sequence sites” dataset (MSS)

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

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

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





**Fig 4. Structural superimpositions for the 10 clusters comprising 96 sites in the SS132 dataset.** Pairs of sites were superimposed on their 7 C $\alpha$  atoms and the C $\beta$  of the central S/T. Their pairwise RMSD were clustered with complete linkage and Euclidean distance. Clusters were defined by a 3 Å threshold. Green, yellow and grey represent residues in H, E, C secondary structures respectively.

<https://doi.org/10.1371/journal.pone.0184405.g004>

**Table 3. Structural evidence of buried O-GlcNAc sites in the SS143 dataset.**

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

RSA—site mean relative solvent accessibility; Cluster id—Clusters in Fig 4.

<https://doi.org/10.1371/journal.pone.0184405.t003>



**Table 4. JPred4 predicted solvent accessibility for S/T in the MSS and USS datasets.** The proportions of buried S/T as predicted by the Jnetsol method in JPred4. The proportions of buried S/T are significantly smaller for modified group.

Buried at	Modified (MSS)		Unmodified (USS)		p 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

95% CI— 95% confidence interval; n—number of S/T predicted to be buried. The p value refers to the two-tailed z-score test between the modified and unmodified groups.

<https://doi.org/10.1371/journal.pone.0184405.t004>

harder to predict than contiguous secondary structure elements, and this may also contribute to the observed enrichment in C.

The analysis of SS143 dataset showed an enrichment of S/T in REM465 regions likely to be disordered or highly mobile. To explore this further, 3 disorder prediction algorithms, giving a total of 6 disorder scores, were run on the MSS and USS datasets as detailed in Methods.

Table 6 shows that, with the exception of DisEMBL-HOTLOOPS which is trained structural B-factors, all methods report a small but significant increase in mean predicted disorder for the modified S/T. To confirm this result, the MSS dataset was compared to the GS dataset, which was selected from proteins known to be predominantly globular, and hence an ordered background. In Fig 5, DisEMBL-HOTLOOPS shows an increase in the ratio of disordered residues around the modified S/T. DisEMBL-COILS and JRonn also indicate a small increase, not in a specific region, but rather for 40 residues around the S/T. IUPred-Long, IUPred-Short and DisEMBL-REM465 show a bigger increase of the ratio of disordered residues in the MSS dataset and IUPred-Short and REM465 have a clearer peak within -15 to 15 residues from the modified S/T. Overall, all methods indicate an increased proportion of predicted disorder in the MSS dataset when compared to the GS dataset.

### Conclusions and final remarks

Despite the substantial evidence of protein structural disorder in the MSS and the SS143 datasets, the SS132 dataset clearly indicates that some of the examined sites appear within ordered regions of the protein structure. Furthermore, InterproScan [45] analysis of O-GlcNAc sites assigned 19% of the sites to protein domains, this is similar to with the 25% phosphoserines and phosphothreonines in PFAM domains [14,46], which are thought to be mostly ordered by definition. So, like protein phosphorylation, O-GlcNAcylated S/T are found in both ordered and disordered regions.

**Table 5. JPred4 predicted secondary structure proportions for S/T in the MSS and USS datasets.**

Secondary structure	Modified (MSS)		Unmodified (USS)		p 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

95% CI— 95% confidence interval; n—the number of S/T; the p value refers to the two-tailed z-score test between the modified and unmodified groups.

<https://doi.org/10.1371/journal.pone.0184405.t005>

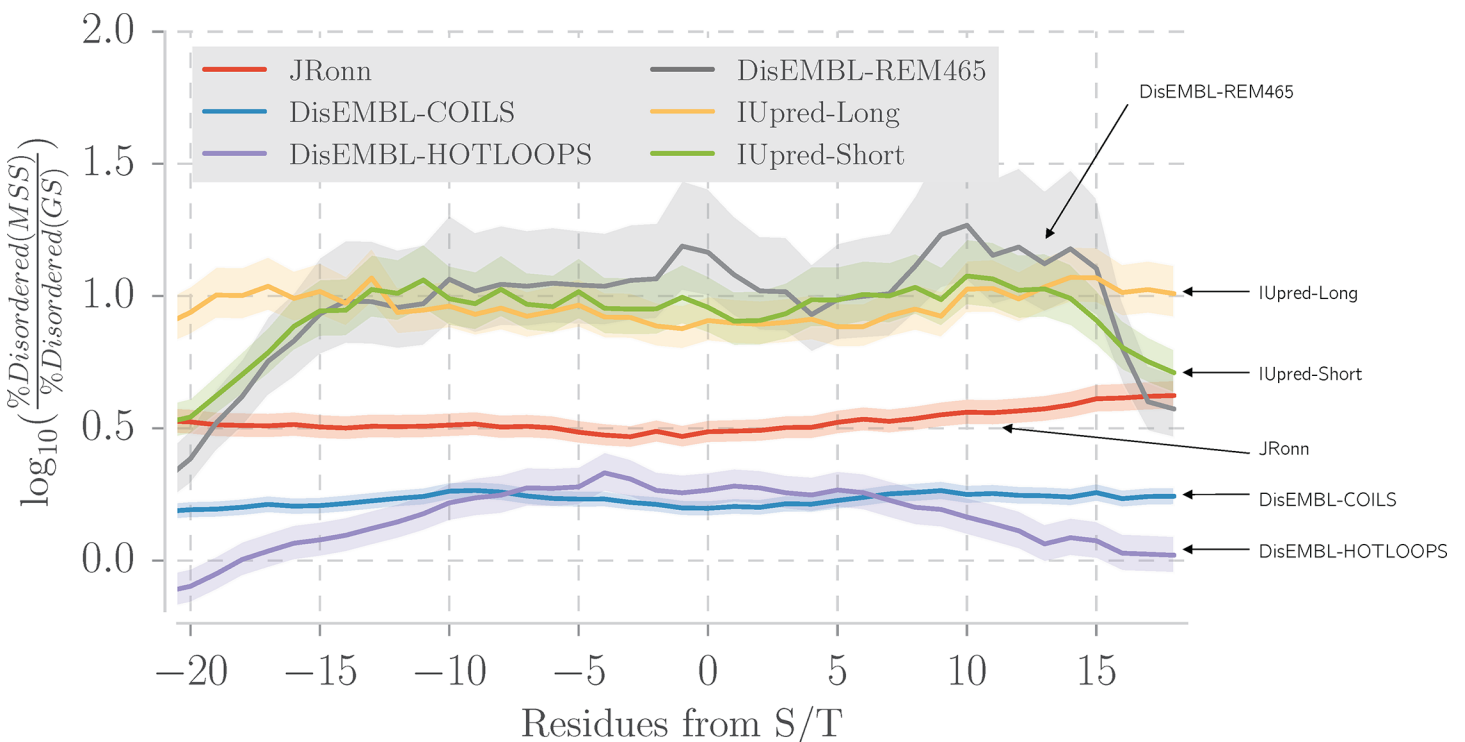
**Table 6. Predicted disorder between modified and unmodified S/T.** All disorder prediction methods, excepting DisEMBL-HOTLOOPS, reveal a small but significant increase of mean disorder score for modified S/T over unmodified ones.

Method	Mean score modified (MSS) ± SE	Mean score unmodified (USS) ± SE	p 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

The p value refers to the two-tailed t-test between the modified and unmodified groups. SE—standard error.

<https://doi.org/10.1371/journal.pone.0184405.t006>

The local tertiary structure of O-GlcNAc sites is indistinguishable from unmodified sites, and so how does OGT recognise the site it modifies? OGT may force the unfolding of the targeted substrate [26]. Moreover, OGT participates in macromolecular assemblies [47], and the role of adaptor proteins cannot be ignored. In protein kinase C (PKC) substrate recognition, residues distant in the protein sequence but close in its three-dimensional structure are critical [48] and non-local interactions might also act in OGT substrate recognition. Other components, such as UDP-GlcNAc concentration and subcellular location-dependent interactions, modulate OGT activity [49], but their part in substrate recognition is still unknown. In



**Fig 5. Predicted disorder around O-GlcNAc-sites in the MSS compared to randomly selected S/T in the GS-dataset.** The y-axis shows the  $\log_{10}$  odds ratio of the between the proportion of disordered residues in the MSS dataset and the proportion of disordered residues in the GS dataset. The semi-transparent area represents 95% confidence intervals. A residue was defined as disordered according to each method's threshold. The x-axis represents the distance in residues to the central residue which is always a S/T. DisEMBL-REM465, IUpred-short predict protein structural disorder specifically around the modification site, while the other methods predict intrinsic disorder over O-GlcNAcylated proteins. DisEMBL-REM465 shows a less pronounced increase in predicted disorder compared to the other methods.

<https://doi.org/10.1371/journal.pone.0184405.g005>

conclusion, although no three-dimensional fingerprint was detected during the structural characterisation of OGT-modified sites, the work confirmed that S/T and surrounding residues are more disordered than the backgrounds tested and that sites in transition between C to H/E might be involved, suggesting that the structural flexibility has a role on OGT site recognition.

## Supporting information

**S1 Table. Properties of sites in the SS132 dataset.** List of all entries in the SS132 dataset. PDB, PDB accession number; Chain, chain in the PDB file; Position, residue position within the chain; Cluster, cluster id. RSA, relative solvent accessibility; SS, secondary structure. (CSV)

## Acknowledgments

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

## Author Contributions

**Conceptualization:** Thiago Britto-Borges, Geoffrey J. Barton.

**Data curation:** Thiago Britto-Borges, Geoffrey J. Barton.

**Formal analysis:** Thiago Britto-Borges, Geoffrey J. Barton.

**Funding acquisition:** Geoffrey J. Barton.

**Investigation:** Thiago Britto-Borges, Geoffrey J. Barton.

**Methodology:** Thiago Britto-Borges, Geoffrey J. Barton.

**Resources:** Thiago Britto-Borges, Geoffrey J. Barton.

**Software:** Thiago Britto-Borges.

**Supervision:** Geoffrey J. Barton.

**Validation:** Thiago Britto-Borges.

**Visualization:** Thiago Britto-Borges.

**Writing – original draft:** Thiago Britto-Borges, Geoffrey J. Barton.

**Writing – review & editing:** Thiago Britto-Borges, Geoffrey J. Barton.

## References

1. Shafi R, Iyer SP, Ellies LG, O'Donnell N, Marek KW, Chui D, et al. The O-GlcNAc transferase gene resides on the X chromosome and is essential for embryonic stem cell viability and mouse ontogeny. *Proc Natl Acad Sci U S A*. 2000; 97: 5735–9. <https://doi.org/10.1073/pnas.100471497> PMID: 10801981
2. O'Donnell N, Zachara NE, Hart GW, Marth JD. Ogt-Dependent X-Chromosome-Linked Protein Glycosylation Is a Requisite Modification in Somatic Cell Function and Embryo Viability. *Mol Cell Biol*. 2004; 24: 1680–1690. <https://doi.org/10.1128/MCB.24.4.1680-1690.2004> PMID: 14749383
3. Buse MG. Hexosamines, insulin resistance, and the complications of diabetes: current status. *Am J Physiol Endocrinol Metab*. 2006; 290: E1–E8. <https://doi.org/10.1152/ajpendo.00329.2005> PMID: 16339923
4. Abdel Rahman AM, Ryczko M, Pawling J, Dennis JW. Probing the hexosamine biosynthetic pathway in human tumor cells by multitargeted tandem mass spectrometry. *ACS Chem Biol*. 2013; 8: 2053–62. <https://doi.org/10.1021/cb4004173> PMID: 23875632

5. Liu J, Marchase RB, Chatham JC. Increased O-GlcNAc levels during reperfusion lead to improved functional recovery and reduced calpain proteolysis. *Am J Physiol Heart Circ Physiol*. 2007; 293: H1391–9. <https://doi.org/10.1152/ajpheart.00285.2007> PMID: 17573462
6. McClain DA, Lubas WA, Cooksey RC, Hazel M, Parker GJ, Love DC, et al. Altered glycan-dependent signaling induces insulin resistance and hyperleptinemia. *Proc Natl Acad Sci U S A*. 2002; 99: 10695–10699. <https://doi.org/10.1073/pnas.152346899> [pii] PMID: 12136128
7. Mariappa D, Zheng X, Schimpl M, Raimi O, Ferenbach AT, Müller H-AJ, et al. Dual functionality of O-GlcNAc transferase is required for *Drosophila* development. *Open Biol*. 2015; 5: 150234. <https://doi.org/10.1098/rsob.150234> PMID: 26674417
8. Lynch TP, Ferrer CM, Jackson SR, Shahriari KS, Vosseller K, Reginato MJ. Critical Role of O-Linked N-Acetylglucosamine Transferase in Prostate Cancer Invasion, Angiogenesis, and Metastasis. *J Biol Chem*. 2012; 287: 11070–11081. <https://doi.org/10.1074/jbc.M111.302547> PMID: 22275356
9. Liu F, Iqbal K, Grundke-Iqbal I, Hart GW, Gong C-X. O-GlcNAcylation regulates phosphorylation of tau: A mechanism involved in Alzheimer's disease. *Proc Natl Acad Sci. National Academy of Sciences*; 2004; 101: 10804–10809. <https://doi.org/10.1073/pnas.0400348101> PMID: 15249677
10. Chen D, Juárez S, Hartweck L, Alamillo JM, Simón-Mateo C, Pérez JJ, et al. Identification of secret agent as the O-GlcNAc transferase that participates in Plum pox virus infection. *J Virol*. 2005; 79: 9381–7. <https://doi.org/10.1128/JVI.79.15.9381-9387.2005> PMID: 16014901
11. Kuo M, Zilberfarb V, Gangneux N, Christeff N, Issad T. O-glycosylation of FoxO1 increases its transcriptional activity towards the glucose 6-phosphatase gene. *FEBS Lett. Elsevier*; 2008; 582: 829–34. <https://doi.org/10.1016/j.febslet.2008.02.010> PMID: 18280254
12. Wells L. Mapping Sites of O-GlcNAc Modification Using Affinity Tags for Serine and Threonine Post-translational Modifications. *Mol Cell Proteomics*. 2002; 1: 791–804. <https://doi.org/10.1074/mcp.M200048-MCP200> PMID: 12438562
13. Parker GJ, Lund KC, Taylor RP, McClain DA. Insulin resistance of glycogen synthase mediated by o-linked N-acetylglucosamine. *J Biol Chem*. 2003; 278: 10022–7. <https://doi.org/10.1074/jbc.M207787200> PMID: 12510058
14. Hornbeck P V, Zhang B, Murray B, Kornhauser JM, Latham V, Skrzypek E. PhosphoSitePlus, 2014: mutations, PTMs and recalibrations. *Nucleic Acids Res*. 2015; 43: D512–20. <https://doi.org/10.1093/nar/gku1267> PMID: 25514926
15. Roquemore EP, Dell A, Morris HR, Panico M, Reason AJ, Savoy LA, et al. Vertebrate lens alpha-crystallins are modified by O-linked N-acetylglucosamine. *J Biol Chem*. 1992; 267: 555–63. Available: <http://www.ncbi.nlm.nih.gov/pubmed/1730617> PMID: 1730617
16. Hahne H, Gholami AM, Kuster B, Moghaddas Gholami A, Kuster B. Discovery of O-GlcNAc-modified proteins in published large-scale proteome data. *Mol Cell Proteomics*. 2012; 11: 843–50. <https://doi.org/10.1074/mcp.M112.019463> PMID: 22661428
17. Ma J, Hart GW. O-GlcNAc profiling: from proteins to proteomes. *Clin Proteomics*. 2014; 11: 8. <https://doi.org/10.1186/1559-0275-11-8> PMID: 24593906
18. Wang J, Torii M, Liu H, Hart GW, Hu Z. dbOGAP—An Integrated Bioinformatics Resource for Protein O-GlcNAcylation. *BMC Bioinformatics. BioMed Central Ltd*; 2011; 12: 91. <https://doi.org/10.1186/1471-2105-12-91> PMID: 21466708
19. Jia C-Z, Liu T, Wang Z-P. O-GlcNAcPRED: a sensitive predictor to capture protein O-GlcNAcylation sites. *Mol Biosyst. The Royal Society of Chemistry*; 2013; 9: 2909–13. <https://doi.org/10.1039/c3mb70326f> PMID: 24056994
20. Gupta R, Brunak S. Prediction of glycosylation across the human proteome and the correlation to protein function. [Internet]. Pacific Symposium on Biocomputing. Pacific Symposium on Biocomputing. 2002. pp. 310–22. [https://doi.org/10.1142/9789812799623\\_0029](https://doi.org/10.1142/9789812799623_0029) PMID: 11928486
21. Jochmann R, Holz P, Sticht H, Stürzl M. Validation of the reliability of computational O-GlcNAc prediction. *Biochim Biophys Acta—Proteins Proteomics. Elsevier B.V.*; 2014; 1844: 416–421. <https://doi.org/10.1016/j.bbapap.2013.12.002> PMID: 24332980
22. Zhao X, Ning Q, Chai H, Ai M, Ma Z. PGlcS: Prediction of protein O-GlcNAcylation sites with multiple features and analysis [Internet]. *Journal of Theoretical Biology. Elsevier*; 2015. pp. 524–529. <https://doi.org/10.1016/j.jtbi.2015.06.026> PMID: 26116363
23. Kao H-J, Huang C-H, Bretaña N, Lu C-T, Huang K-Y, Weng S-L, et al. A two-layered machine learning method to identify protein O-GlcNAcylation sites with O-GlcNAc transferase substrate motifs [Internet]. *BMC Bioinformatics. BioMed Central Ltd*; 2015. p. S10. <https://doi.org/10.1186/1471-2105-16-S18-S10> PMID: 26680539
24. Wu H-Y, Lu C-T, Kao H-J, Chen Y-JY-J, Chen Y-JY-J, Lee T-Y. Characterization and identification of protein O-GlcNAcylation sites with substrate specificity. [Internet]. *BMC bioinformatics. BioMed Central Ltd*; 2014. p. S1. <https://doi.org/10.1186/1471-2105-15-S16-S1> PMID: 25521204

25. Crooks GE, Hon G, Chandonia J-M, Brenner SE. WebLogo: a sequence logo generator. *Genome Res.* Cold Spring Harbor Laboratory Press; 2004; 14: 1188–90. <https://doi.org/10.1101/gr.849004> PMID: 15173120
26. Pathak S, Alonso J, Schimpl M, Rafie K, Blair DE, Borodkin VS, et al. The active site of O-GlcNAc transferase imposes constraints on substrate sequence. *Nat Struct Mol Biol.* Nature Publishing Group, a division of Macmillan Publishers Limited. All Rights Reserved.; 2015; 22: 744–750. <https://doi.org/10.1038/nsmb.3063> PMID: 26237509
27. Shi J, Sharif S, Ruijtenbeek R, Pieters RJ. Activity Based High-Throughput Screening for Novel O-GlcNAc Transferase Substrates Using a Dynamic Peptide Microarray. *PLoS One.* 2016; 11: e0151085. <https://doi.org/10.1371/journal.pone.0151085> PMID: 26960196
28. Liu X, Li L, Wang Y, Yan H, Ma X, Wang PG, et al. A peptide panel investigation reveals the acceptor specificity of O-GlcNAc transferase. *FASEB J.* 2014; 28: 3362–3372. <https://doi.org/10.1096/fj.13-246850> PMID: 24760753
29. Lazarus MB, Nam Y, Jiang J, Sliz P, Walker S. Structure of human O-GlcNAc transferase and its complex with a peptide substrate [Internet]. *Nature.* Nature Publishing Group; 2011. pp. 564–567. <https://doi.org/10.1038/nature09638> PMID: 21240259
30. Schimpl M, Zheng X, Borodkin VS, Blair DE, Ferenbach AT, Schüttelkopf AW, et al. O-GlcNAc transferase invokes nucleotide sugar pyrophosphate participation in catalysis. *Nat Chem Biol.* 2012; 8: 969–74. <https://doi.org/10.1038/nchembio.1108> PMID: 23103942
31. Velankar S, Alhroub Y, Alili A, Best C, Boutselakis HC, Caboche S, et al. PDBE: Protein Data Bank in Europe. *Nucleic Acids Res.* 2011; 39: D402–10. <https://doi.org/10.1093/nar/gkq985> PMID: 21045060
32. Velankar S, Dana JM, Jacobsen J, van Ginkel G, Gane PJ, Luo J, et al. SIFTS: Structure Integration with Function, Taxonomy and Sequences resource. *Nucleic Acids Res.* 2013; 41: D483–9. <https://doi.org/10.1093/nar/gks1258> PMID: 23203869
33. Kabsch W, Sander C. Dictionary of protein secondary structure: pattern recognition of hydrogen-bonded and geometrical features. *Biopolymers.* 1983; 22: 2577–2637. <https://doi.org/10.1002/bip.360221211> PMID: 6667333
34. Cuff J a, Barton GJ. Evaluation and improvement of multiple sequence methods for protein secondary structure prediction. *Proteins.* 1999. pp. 508–519. [https://doi.org/10.1002/\(SICI\)1097-0134\(19990301\)34:4<508::AID-PROT10>3.0.CO;2-4](https://doi.org/10.1002/(SICI)1097-0134(19990301)34:4<508::AID-PROT10>3.0.CO;2-4) [pii] PMID: 10081963
35. Cuff JA, Barton GJ. Application of multiple sequence alignment profiles to improve protein secondary structure prediction. [Internet]. *Proteins.* John Wiley & Sons, Inc.; 2000. pp. 502–511. [https://doi.org/10.1002/1097-0134\(20000815\)40:3<502::AID-PROT170>3.0.CO;2-Q](https://doi.org/10.1002/1097-0134(20000815)40:3<502::AID-PROT170>3.0.CO;2-Q) PMID: 10861942
36. Drozdetskiy A, Cole C, Procter J, Barton GJ. JPred4: A protein secondary structure prediction server. *Nucleic Acids Res.* 2015; 43: W389–94. <https://doi.org/10.1093/nar/gkv332> PMID: 25883141
37. Yang ZR, Thomson R, McNeil P, Esnouf RM. RONN: the bio-basis function neural network technique applied to the detection of natively disordered regions in proteins. *Bioinformatics.* 2005; 21: 3369–76. <https://doi.org/10.1093/bioinformatics/bti534> PMID: 15947016
38. Dosztányi Z, Csizsók V, Tompa P, Simon I. The pairwise energy content estimated from amino acid composition discriminates between folded and intrinsically unstructured proteins. *J Mol Biol.* 2005; 347: 827–39. <https://doi.org/10.1016/j.jmb.2005.01.071> PMID: 15769473
39. Linding R, Jensen LJ, Diella F, Bork P, Gibson TJ, Russell RB. Protein Disorder Prediction. *Structure.* 2003; 11: 1453–1459. <https://doi.org/10.1016/j.str.2003.10.002> PMID: 14604535
40. Troshin P V, Procter JB, Barton GJ. Java bioinformatics analysis web services for multiple sequence alignment—JABAWS:MSA. *Bioinformatics.* Oxford University Press; 2011; 27: 2001–2002. <https://doi.org/10.1093/bioinformatics/btr304> PMID: 21593132
41. Fox NK, Brenner SE, Chandonia J-M. SCOPe: Structural Classification of Proteins—extended, integrating SCOP and ASTRAL data and classification of new structures. *Nucleic Acids Res.* 2014; 42: D304–9. <https://doi.org/10.1093/nar/gkt1240> PMID: 24304899
42. McKinney W, Team PD. Pandas—Powerful Python Data Analysis Toolkit. *Pandas—Powerful Python Data Analysis Toolkit.* 2015. p. 1625.
43. Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, et al. Biopython: freely available Python tools for computational molecular biology and bioinformatics [Internet]. *Bioinformatics.* 2009. pp. 1422–1423. <https://doi.org/10.1093/bioinformatics/btp163> PMID: 19304878
44. Trinidad JC, Barkan DT, Gullede BF, Thalhammer A, Sali A, Schoepfer R, et al. Global Identification and Characterization of Both O-GlcNAcylation and Phosphorylation at the Murine Synapse. *Mol Cell Proteomics.* 2012; 11: 215–229. <https://doi.org/10.1074/mcp.O112.018366> PMID: 22645316

45. Zdobnov EM, Apweiler R. InterProScan—an integration platform for the signature-recognition methods in InterPro. *Bioinformatics*. 2001; 17: 847–848. <https://doi.org/10.1093/bioinformatics/17.9.847> PMID: [11590104](https://pubmed.ncbi.nlm.nih.gov/11590104/)
46. Beltrao P, Albanèse V, Kenner LR, Swaney DL, Burlingame A, Villén J, et al. Systematic Functional Prioritization of Protein Posttranslational Modifications. *Cell*. 2012; 150: 413–425. <https://doi.org/10.1016/j.cell.2012.05.036> PMID: [22817900](https://pubmed.ncbi.nlm.nih.gov/22817900/)
47. Wells L, Kreppel LK, Comer FI, Wadzinski BE, Hart GW. O-GlcNAc Transferase Is in a Functional Complex with Protein Phosphatase 1 Catalytic Subunits. *J Biol Chem*. 2004; 279: 38466–38470. <https://doi.org/10.1074/jbc.M406481200> PMID: [15247246](https://pubmed.ncbi.nlm.nih.gov/15247246/)
48. Duarte ML, Pena DA, Nunes Ferraz FA, Berti DA, Paschoal Sobreira TJ, Costa-Junior HM, et al. Protein folding creates structure-based, noncontiguous consensus phosphorylation motifs recognized by kinases. *Sci Signal*. 2014; 7: ra105–ra105. <https://doi.org/10.1126/scisignal.2005412> PMID: [25372052](https://pubmed.ncbi.nlm.nih.gov/25372052/)
49. Nagel AK, Ball LE. O-GlcNAc transferase and O-GlcNAcase: achieving target substrate specificity. *Amino Acids*. 2014; 46: 2305–2316. <https://doi.org/10.1007/s00726-014-1827-7> PMID: [25173736](https://pubmed.ncbi.nlm.nih.gov/25173736/)