



Published in final edited form as:

ACS Catal. 2021 March 5; 11(5): 2977–2991. doi:10.1021/acscatal.0c04609.

## Structural basis for peptide substrate specificities of glycosyltransferase GalNAc-T2

Sai Pooja Mahajan<sup>1</sup>, Yashes Srinivasan<sup>2</sup>, Jason W. Labonte<sup>1,3</sup>, Matthew P. DeLisa<sup>4</sup>, Jeffrey J. Gray<sup>1,5,\*</sup>

<sup>1</sup>Department of Chemical and Biomolecular Engineering, Johns Hopkins University, Baltimore, Maryland 21218, United States

<sup>2</sup>Department of Bioengineering, University of California, Los Angeles, Los Angeles, California 90095, United States

<sup>3</sup>Department of Chemistry, Franklin & Marshall College, Lancaster, Pennsylvania 17604, United States

<sup>4</sup>Robert Frederick Smith School of Chemical and Biomolecular Engineering, Department of Microbiology, and Nancy E. and Peter C. Meinig School of Biomedical Engineering, Biochemistry, Molecular and Cell Biology, Cornell University, Ithaca, New York 14853, United States

<sup>5</sup>Sidney Kimmel Comprehensive Cancer Center, Johns Hopkins School of Medicine, Baltimore, Maryland, United States

### Abstract

The polypeptide *N*-acetylgalactosaminyl transferase (GalNAc-T) enzyme family initiates *O*-linked mucin-type glycosylation. The family constitutes 20 isoenzymes in humans. GalNAc-Ts exhibit both redundancy and finely tuned specificity for a wide range of peptide substrates. In this work, we deciphered the sequence and structural motifs that determine the peptide substrate preferences for the GalNAc-T2 isoform. Our approach involved sampling and characterization of peptide–enzyme conformations obtained from Rosetta Monte Carlo-minimization–based flexible docking. We computationally scanned 19 amino acid residues at positions –1 and +1 of an eight-residue peptide substrate, which comprised a dataset of 361 (19x19) peptides with previously characterized experimental GalNAc-T2 glycosylation efficiencies. The calculations recapitulated experimental specificity data, successfully discriminating between glycosylatable and non-glycosylatable peptides with a probability of 96.5% (ROC-AUC score), a balanced accuracy of 85.5% and a false positive rate of 7.3%. The glycosylatable peptide substrates *viz.* peptides with proline, serine, threonine, and alanine at the –1 position of the peptide preferentially exhibited cognate sequon-like conformations. The preference for specific residues at the –1 position of the peptide was regulated by enzyme residues R362, K363, Q364, H365 and W331, which modulate the pocket size and specific enzyme-peptide interactions. For the +1 position of the peptide,

\* Corresponding Author: Author email list, J.J.G: jgray@jhu.edu.

Supporting Information

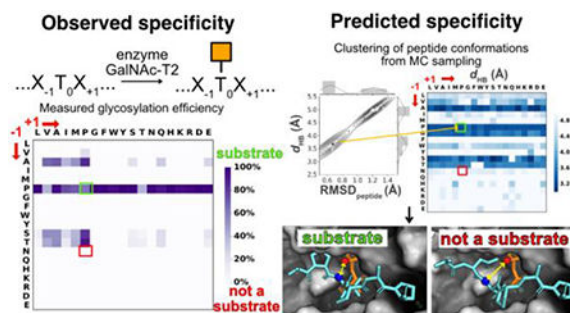
Supporting Figures S1–S24

Supporting Tables S1–4

Supporting methods and code snippets

enzyme residues K281 and K363 formed gating interactions with aromatics and glutamines at the +1 position of the peptide, leading to modes of peptide-binding sub-optimal for catalysis. Overall, our work revealed enzyme features that lead to the finely tuned specificity observed for a broad range of peptide substrates for the GalNAc-T2 enzyme. We anticipate that the key sequence and structural motifs can be extended to analyze specificities of other isoforms of the GalNAc-T family and can be used to guide design of variants with tailored specificity.

## Graphical Abstract



## Keywords

GalNAcT; glycosyltransferases; O-linked glycosylation; enzyme specificity; computational specificity prediction

## Introduction

In higher organisms, O-linked *N*-acetylgalactosamine (GalNAc) glycosylation (or mucin-type glycosylation) is an abundant and essential post-translational modification. This type of glycosylation is initiated by a family of glycosyltransferases (GTs) known as polypeptide *N*-acetylgalactosaminyltransferases or GalNAc-Ts. These enzymes transfer a GalNAc sugar from a donor uridine di-phosphate (UDP) nucleotide to the hydroxyl group of a threonine or serine residue of an acceptor peptide. This transfer is the first committed step of mucin-type O-glycosylation, and these enzymes, therefore, define the sites of O-glycosylation. The resulting O-linked GalNAc is further extended to one of the four common core structures, which can be subsequently extended to give mature glycans.<sup>1,2</sup> Aberrant O-glycosylation is a well-known marker of many cancers and has also been linked to developmental and metabolic disorders.<sup>3,4</sup>

In humans, the GalNAc-T family constitutes 20 isoforms. The unusually large number of isoforms for glycosylation is peculiar to O-glycosylation, and the multiplicity is conserved in mammalian evolution, suggesting that cell or tissue specific isoforms have specialized functions.<sup>5</sup> The isoforms exhibit specific substrate preferences that vary with isoenzyme surface charge, prior neighboring long-range and short-range glycosylation patterns and the sequence of the acceptor peptide substrate. Over the last two decades, the peptide substrate preferences for a large number of isoforms have been established by *in vitro* studies.<sup>6–8</sup> The peptide substrate is characterized by a sequence motif (or sequon), Thr/Ser-Pro-X-Pro (T/

SPXP), where T/S is the site of glycosylation (position 0). This sequon is the only conserved consensus motif modified by all isoforms except T7 and T10. The proline at the +3 position of the sequon is supported by a conserved structural motif, *viz.*, the “proline pocket” in the enzyme’s peptide binding groove in all isoforms that bind the T/SPXP motif.<sup>9–11</sup> While this may suggest that all peptides with this motif are valid substrates for the GalNAc-T enzymes, in practice, the T/SPXP motif was obtained by averaging preferences over all isoforms (with characterized specificity). The motif simply indicates preference (rather than a strict rule or constraint) for a specific amino acid residue at a given position on the sequon, such as proline at the +1 and +3 positions. For the remaining positions in the sequon, most isoforms exhibit overlapping yet selective preferences for different amino acid residues. For example, at the –1 position with respect to the glycosylation site, T1 favors aromatics<sup>12</sup> and T12 prefers bulky non-polar residues;<sup>13</sup> whereas T2 exhibits very little to no activity for these amino acids and instead prefers threonine, proline, alanine, and serine. Yet both T1 and T2 glycosylate the sequon T<sub>-1</sub>TP<sub>+1</sub><sup>12</sup> (with threonine at –1 and proline at +1 positions). Moreover, high-resolution glycosylation experiments from Kightlinger *et al.*<sup>12</sup> (explicit determination of glycosylation efficiency for a given sequon) demonstrate that the (T/S)PXP motif is neither sufficient nor necessary for glycosylation. For example, their data shows that while T2 readily glycosylated sequons T<sub>-1</sub>TP<sub>+1</sub>AP and S<sub>-1</sub>TP<sub>+1</sub>AP, it failed to modify sequons F<sub>-1</sub>TP<sub>+1</sub>AP, Q<sub>-1</sub>TP<sub>+1</sub>AP, K<sub>-1</sub>TP<sub>+1</sub>AP etc. (Figure 1 A). These observations have led to the hypothesis that GalNAc-Ts exhibit both redundancy and finely tuned specificity for a wide range of peptide substrates.

While there is ample experimental data on the peptide substrate specificities of various isoforms, the molecular basis for observed peptide substrate specificities is not well understood. Computational work, so far, has been focused on understanding the mechanism of sugar transfer,<sup>14,15</sup> conformational changes in the flexible loop in the catalytic domain,<sup>16,17</sup> and the effect of the flexible linker connecting the catalytic and lectin domains.<sup>11</sup> None of the computational studies so far have examined the amino acid preferences at different positions on the peptide. Computational studies can pinpoint key positions and structural motifs on an isoform that contribute to peptide substrate specificity. These sequence and structural motifs can be studied across isoforms to reveal more general patterns, to modulate enzyme specificity, and to gain insight into the consequences of enzyme and substrate mutations implicated in aberrant glycosylation, (*e.g.*, colorectal cancer associated mutations of GalNAc-T12<sup>18</sup>) paving the way for rational design of specific drugs/inhibitors<sup>19</sup>.

In this work, we seek to understand the sequence and structural motifs that determine the peptide substrate preferences for the GalNAc-T2 isoform. Our immediate goal is to recapitulate experimentally determined specificity in terms of glycosylation efficiency for sequon variations at positions –1 and +1 (19 amino acid residues tested for each position), as reported by Kightlinger *et al.*,<sup>12</sup> and to understand the structural motifs that best explain experimentally observed trends. To recapitulate experimentally observed specificities for a large dataset, we need an efficient, high-throughput computational method that can capture the key mechanisms of enzymatic catalysis.

Enzymatic catalysis relies primarily on selective transition state stabilization, ground state (reactants) destabilization, dynamics, and active-site gating.<sup>20,21</sup> In practice, these effects

occur at different length- and time-scales and therefore cannot be accurately captured by a single method, even when enzyme crystal structures are available.<sup>22</sup> Hybrid quantum mechanics/molecular mechanics (QM/MM) simulations have been able to recapitulate catalytic proficiency or mechanistic details for many enzymes such as Kemp eliminases<sup>23</sup> or glycoside hydrolases<sup>24</sup> as they are well-suited to characterize the transition state. QM/MM simulations, however, are not suitable to capture binding or dynamics over longer timescales and are prohibitively expensive for a larger dataset. Other factors that determine the stability of the transition state are electrostatic- and shape-complementarity at the peptide-enzyme interface. Electrostatic complementarity can be captured by various computational techniques (*e.g.*, Monte Carlo (MC) or molecular dynamics (MD)-based methods with Poisson-Boltzmann electrostatics or other continuum electrostatics models) at different length-scales. Other effects are determined by the thermodynamics of the enzyme-peptide interactions. To achieve a lower free energy of activation,<sup>20,25</sup> the enzyme must stabilize the transition state selectively relative to the reactants. Additionally, if the product is too stable in the enzyme's active site, product release becomes the catalytic rate-limiting step. This thermodynamic description demands the use of methods that capture multiple states (reactants, products and transition states).<sup>26,27</sup> Furthermore, dynamics is important in many catalytic mechanisms, from small vibrations that lead to rate-promoting motions<sup>28</sup> to large conformational changes and rearrangements in the molecular structure.<sup>29</sup> Active site gating is another important mechanism for catalysis by which key residues outside the active site regulate access to the active site.<sup>30,31</sup> These thermodynamic and kinetic effects, primarily in the nanosecond to microsecond timescales, can be captured faithfully by MD simulations though such simulations can be computationally prohibitive for comparing a large number of substrates. An alternative to MD simulations are Monte Carlo-minimization<sup>32</sup> (MCM) approaches, which are computationally faster and can be reliably used to determine thermodynamically stable native-like states.

Rosetta-based MCM computational protocols, notably, pepspec,<sup>22</sup> sequence-tolerance<sup>33,34</sup> and MFpred<sup>35</sup> have previously been used for predicting the sequence profiles of peptides recognized by various multi-specific protein recognition domains (PRDs) such as PDZ, SH2, SH3, kinases, and proteases. All protocols rely on MCM sampling and aim to approximate the stabilization of the substrate-bound state or transition state in the enzyme's peptide binding groove. The transition state is approximated by known cognate sequon conformations in the enzyme's active site (based on crystal structures and/or homology modeling) with additional constraints to preserve important structural motifs pertaining to the transition state, when available. In the absence of constraints, this approach is equivalent to evaluating the stabilization of the substrate-bound state.<sup>36</sup> MCM allows for faster sampling facilitating the scanning of a large number of amino acid residues at multiple positions of the peptide substrate. All three protocols achieve impressive accuracy in predicting experimentally observed profiles for many PRDs. However, since all three methods are developed with the broad goal of predicting sequence specificity profiles for a range of PRDs, the accuracy of prediction may not be sufficient to pinpoint subtle differences in specificity for a specific target of interest. For example, the sequence-tolerance protocol pre-calculates the interactions between all interacting residues ignoring changes in conformation of the peptide in the protein's binding pocket. All three methods

struggle to predict specificity for HIV-1 protease, which has a relaxed specificity profile and a preference for small hydrophobic residues, similar to GalNAc-T2. Additionally, all three protocols employ limited backbone sampling, prohibiting the free conformational sampling of the peptide in the binding groove. For the more targeted goal of designing a peptide inhibitor to discriminate between two similar PDZ domains, Zheng *et al.*<sup>1</sup> employ extensive conformational sampling using the full-fledged flexpepdock protocol<sup>38,39</sup> along with the CLASSY method to achieve a solution with desired specificity and affinity goals. Similarly, Pethe *et al.*<sup>40</sup> were able to obtain significantly improved prediction accuracies for proteases (including HIV-1 protease) compared to previous methods (MFPred, sequence tolerance and pepspec) by employing machine learning and a discriminatory score based on geometric features, interface score terms from Rosetta, and electrostatic score terms from Amber. In another work, Pethe *et al.*<sup>41</sup> used supervised learning on experimentally obtained deep-sequencing data and information from structure-based models to chart the specificity landscape of 3.2 million substrate variants of the viral protease HCV.

Here, we sought to understand specificity determinants for a specific isoform of the GalNAc-T family and to pinpoint sequence and structural motifs in the enzyme that explain fine-tuning of specificity. To this end, we developed a customized Rosetta-based protocol<sup>42,43</sup> that allowed us to model structures of all 361 peptide sequons (19×19) with the GalNAc-T2 enzyme and computationally determine the sequon preference for the GalNAc-T2 isoform. Our protocol was similar in spirit to earlier protocols in that it docks the peptide substrate into the enzyme's active site. However, unlike pepspec,<sup>22</sup> sequence-tolerance<sup>33,34</sup> and MFPred<sup>35</sup> and similar to the protocol of Zheng *et al.*,<sup>37</sup> we allowed fully flexible peptide sampling (as opposed to limited or no backbone sampling) followed by clustering and analysis of the sampled low energy decoys. Our strategy relied on characterizing the peptide binding to the enzyme with a range of structural features at the interface as a function of the amino acid residues at the +1 and -1 positions. Using our methodology, we were able to identify features that recapitulated high-quality experimental specificity data for GalNAc-T2. Extensive peptide backbone sampling revealed that the peptide binding groove of GalNAc-T2 stabilized multiple competing conformations/states – some leading to efficient glycosylation and others hampering it. Furthermore, multiple stable states suggested that kinetics might play an important role in determining specificity. Thus, finely-tuned specificity might be achieved by modulating the relative stability of these states to discriminate between peptide substrates for an isoform and across isoforms. Overall, our work reveals key residues on the enzyme that determine peptide substrate preferences at various sequon positions.

## Results

### Clustering of low interaction energy decoys reveals that peptides exhibit multiple competing low-energy conformations

We studied all 361 (19x19) sequons obtained by scanning 19 amino acids (all amino acids except cysteine) at positions -1 and +1 with respect to the modified threonine (at position 0 or T<sub>0</sub>). The experimentally determined glycosylation efficiencies for all of these sequons was determined by Kightlinger *et al.*<sup>12</sup> and replotted in Figure 1A (Figure S1). For each

sequon, we started with the co-crystal structure of the peptide and UDP-sugar bound to the enzyme (pdb ids: 4d0z and 2ffu, respectively).<sup>16</sup> We mutated the residues at the  $-1$  and  $+1$  position of the peptide to the target sequon and repacked and minimized the nearby side chains to obtain a starting enzyme-peptide configuration for the target sequon. We then subjected this starting structure to MCM sampling of rigid-body displacements and peptide torsion angles in two stages - a low-resolution centroid stage with simulated annealing followed by a high-resolution all-atom stage to generate 2,000 structures (decoys) per sequon (see Methods). In this paper, we use the shorthand notation  $X_N$  for amino acid 'X' (denoted by 1-letter code) and sequon position 'N' and a sequon with the shorthand notation  $X_{-1}TX_{+1}$ . For example,  $P_{-1}$  denotes amino acid proline at the  $-1$  position,  $M_{+1}$  denotes a methionine at the  $+1$  position and  $P_{-1}TM_{+1}$  denotes a sequon or peptide with  $P_{-1}$  and  $M_{+1}$ . For brevity, we also refer to peptides or sequons containing residue X at position N as "XN peptides" or "XN sequons" respectively.

Preliminary analysis of the decoys showed that many peptides exhibited multiple stable states with comparable energies of interaction (interaction energy) between the peptide and enzyme. In Figure 1B, we show plots of interaction energy vs. distance to the reference crystal structure for four randomly chosen sequons obtained from MCM sampling of the peptide substrate with the respective sequons in the enzyme's peptide-binding groove. For all four sequons in Figure 1B, we observed multiple clusters of low interaction energy decoys, or "funnels." For example, for sequon  $T_{-1}TQ_{+1}$ , we observed two distinct funnels at  $RMSD_{peptide}$  (the root mean square deviation of  $C_\alpha$  carbons of the peptide backbone with respect to the peptide in the crystal structure) values of about  $\sim 0.65 \text{ \AA}$  and  $1.25 \text{ \AA}$  with comparable lowest interaction energies. Overall, 57% (205/361) of the sequons exhibited two significant clusters. 39/205 (19.0%) and 81/205 (39.5%) sequons exhibited lowest energy states for each cluster within 0.5 and 1.0 Rosetta energy units (or REUs) respectively, of each other, underscoring the importance of considering both states (Figure S2).

To characterize multiple low-energy conformations and to construct a more complete picture of the landscape of structural conformations sampled by the peptide substrate in the enzyme cavity, we developed the computational flow summarized in Figure 1C. For each sequon, we selected the top-10%-scoring decoys (by interaction energy) from MCM sampling and then clustered them using three features. The first feature,  $RMSD_{peptide}$ , characterized decoys on the basis of the similarity of the peptide backbone conformation and position of the peptide in the crystal structure. Next,  $doc$ , the distance between the hydroxyl group of  $T_0$  and the anomeric carbon ( $C_1$ ) on the sugar tracked the distance for the new glycosidic linkage. Finally,  $d_{HB}$ , the distance between the amide group of  $T_0$  on the peptide and the oxygen of the  $\beta$ -phosphate group of UDP ( $O_{\beta-PO4}$ ) was a reaction coordinate characterizing a transition-state-stabilizing hydrogen bond between the backbone amide of  $T_0$  and UDP.<sup>14</sup>

We characterized the lowest-energy decoy for the largest and second-largest clusters obtained for each sequon and plotted heatmaps to show the distribution of the lowest interaction energy (Figure 2A), normalized cluster size (Figure 2B),  $RMSD_{peptide}$  (Figure 2C) and  $d_{HB}$  (Figure 2D) for all 361 sequons (see also Figure S3 and S4). We also characterized the clusters by the decoy representing the center of the cluster, the average over all decoys with interaction energies within 1 REU and the average over the five decoys



in the cluster with the lowest interaction energies. All strategies resulted in similar heatmaps (Figure S5) and hence, going forward, we represented a cluster by the lowest energy decoy belonging to that cluster.

Horizontal stripes emerging across the  $\text{RMSD}_{\text{peptide}}$  and  $d_{\text{HB}}$  heatmaps (Figure 2C–D) suggested that sequons with the same amino acid at the  $-1$  position (horizontal axis) exhibited similar  $\text{RMSD}_{\text{peptide}}$  and  $d_{\text{HB}}$  values. To probe whether the low-energy conformations exhibited by various peptides depends on the identity of the residue in a position-specific manner, we plotted the  $\text{RMSD}_{\text{peptide}}$  and  $d_{\text{HB}}$  of the lowest energy decoy for the two largest clusters for each sequon colored by the amino acid residue at the  $-1$  (Figure 3A) and  $+1$  (Figure 3B) positions. It is apparent in Figure 3A that sequons with the same amino acid residue at the  $-1$  position, especially A, G, T, P, S and V, were grouped or clustered together. The clustering or grouping suggests that sequons with the same amino acid at the  $-1$  position exhibited similar conformations or low-energy states. Similar grouping was not observed for sequons with the same residue at the  $+1$  position (Figure 3B). This high-level analysis of low-energy conformations for the entire dataset suggested that the  $-1$  position plays a dominant role in determining the low-energy conformation(s) exhibited by a sequon and that the  $+1$  position contributed in a secondary capacity.

In the following sections, we present hypotheses to explain how each position ( $-1$ ,  $+1$ ) contributed in a characteristic manner to determine the low-energy conformations exhibited by a sequon and how these conformations, in turn, related to experimentally determined specificities. We characterized the low-energy conformations by a selected set of relevant features.

To compare our predictions with experiments, we employed logistic regression and, unless otherwise indicated, we labeled all sequons with experimental glycosylation efficiencies greater than 10% (efficiency threshold) as glycosylatable and those with efficiencies less than 10% as unglycosylatable, in line with previous work<sup>35</sup>. In Table 1 (and Table S1), for a chosen value of efficiency threshold, we have tabulated the area-under the curve (AUC) of the receiver operating curve (ROC). Since the dataset of glycosylation efficiency measurements is highly imbalanced with only 46/361 (12.7%) glycosylatable sequons (positive samples) and 315/361 (87.3%) unglycosylatable sequons (negative samples), in Table 1 (and Table S2), for each metric, we also report the true positives (TP), true negatives (TN), false positives (FP) and false negatives (FN) and the balanced accuracy ( $\text{BA} = (\text{TPR} + \text{TNR}) / 2$ , where the true positive rate  $\text{TPR} = \text{TP} / (\text{TP} + \text{FN})$  and the true negative rate  $\text{TNR} = \text{TN} / (\text{TN} + \text{FP})$  and false positive rate ( $\text{FPR} = \text{FP} / (\text{FP} + \text{TN})$ ). Note that a naive classifier that classifies all 361 sequons as unglycosylatable has an accuracy ( $\text{ACC} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN})$ ) of 87.3% ( $= 315/361$ ), a BA of 50% ( $= (0/46 + 315/315)/2$ ) and an FPR of 0%. Since accuracy (as opposed to balanced accuracy) higher than such a naive classifier may come at the expense of FPs, we report balanced accuracy instead of accuracy for our dataset.

Energy is a commonly used metric in determining specificity of peptide substrates (e.g.  $\text{pepspec}$ <sup>22</sup>,  $\text{sequence\_tolerance}$ <sup>33</sup> and  $\text{MFPred}$ <sup>35</sup>). However, for the purpose of prediction of specificity trends for the T2 enzyme, interaction energy by itself was a weak predictor of

specificity (AUC = 0.566, Table 1). The significantly lower AUC values based on interaction energy were due to the fact that P<sub>-1</sub> and S<sub>-1</sub> peptides bind the enzyme with significantly lower interaction energies than A<sub>-1</sub>, G<sub>-1</sub> or T<sub>-1</sub> peptides (Figure 2A). For comparison with other energy based approaches, we applied MFpred<sup>35</sup> to obtain the specificity profile for GalNAc-T2 (Figure S6). MFpred uses a mean-field approach that assumes each residue position is independent. MFpred obtains an AUC score of 0.68 at the -1 position and 0.50 at the +1 position (Table S3). For comparison of experimental and MFpred specificity logos, see Figure S6. In the MFpred study, the addition of structural-motif-preserving constraints improved predictions for some PBDs. For the MFpred scores reported here, we did not add any constraints.

We omitted  $d_{OC}$  from future characterization due to its low AUC score in classifying glycosylatable and non-glycosylatable sequons (AUC score 0.43; Table S1). However, we retained clustering in the 3d feature space formed by RMSD<sub>peptide</sub>,  $d_{OC}$  and  $d_{HB}$  since  $d_{OC}$  improves resolution of conformations for some sequons (Figure S7).

### Recapitulation of amino acid specificity trends for the -1 position

#### Sampling of TS-critical hydrogen bond recapitulates specificity for 84% of the sequons with a false positive rate of 14%—

In QM/MM simulations of the glycosylation of the EA2 peptide by GalNAc-T2, Gomez *et al.* characterized a hydrogen bond between the backbone amide of Thr<sub>0</sub> and the  $\beta$ -phosphate group on UDP<sup>14</sup>. They proposed that the hydrogen bond stabilizes the transition state (TS) in “a general catalytic strategy used in peptide O-glycosylation by retaining glycosyltransferases”. Hence, our first hypothesis was that successful glycosylation requires a peptide to exhibit a low-energy conformation with dim distances compatible with the proposed hydrogen bond. Thus, in Figure 4A, we show the heatmap of  $d_{HB}$  for the lowest interaction energy decoys. Applying a 4.0 Å threshold to the 19×19 grid of sequons splits the sequons into those that do not meet this condition (*i.e.*, > 4.0 Å,) and those that exhibited a representative low-energy conformation compatible with hydrogen bonding between the peptide and UDP. When compared with experimental results (Figure 1 A), this criterion discriminated well between substrate peptides and non-substrates of the enzyme (Figure 4B) with 44 TPs, 259 TNs, 56 FPs and 2 FNs (Table 1, Table S2).

A ROC analysis (Figure 4C) showed that the  $d_{HB}$  value of the lowest-energy decoy of the largest cluster correctly distinguished the glycosylatable sequons from the non-glycosylatable sequons with a probability of 92.3% (ROC-AUC value of 0.923; Figure 4C). Setting a threshold of  $d_{HB} < 4$  Å,  $d_{HB}$  classified with a balanced accuracy of 88.9% (TPR = 44/46, TNR = 259/316), accuracy of 84% (=303/316) (Table S2) and a false positive rate of 17.8% (56/(56+259)) (Table 1). If instead of  $d_{HB}$ , we used the fraction of the decoys ( $N_d$ ) that satisfied the  $d_{HB} < 4$  Å criterion, the probability of correctly classifying a sequon was 90.6% (ROC-AUC was 0.906; Figure 4C). Setting an arbitrary threshold of  $N_d > 0.70$ , this feature classified with a balanced accuracy of 85.2% (TPR=39/46, TNR=270/316) of the sequons with a false positive rate of 14.3% (45/(45+270)) (Table 1). The large number of FPs for both  $d_{HB}$  and  $N_d$  indicates that the of  $d_{HB} < 4$  Å criterion has low precision.



**Large amino acid residues are excluded from GalNAc-T2's "-1 pocket"**—Figure 4A reveals that peptides preferentially exhibited low-energy, highly populated states (higher fraction of decoys) with  $d_{\text{HB}}$  distances compatible with hydrogen-bonding when amino acid residues with smaller side chains such as proline, alanine, glycine, serine, or threonine were present at the -1 position. To understand the structural basis for the observed  $d_{\text{HB}}$  trends, we considered specific sequons and their lowest-energy conformations. In Figure 4D, we show the  $d_{\text{HB}}$  distances sampled by the top 10% low-energy decoys for four representative sequons –  $P_{-1}TP_{+1}$  and  $T_{-1}TP_{+1}$  (preferentially sampled conformations with  $d_{\text{HB}} < 4.0 \text{ \AA}$ ), and  $N_{-1}TP_{+1}$  and  $R_{-1}TP_{+1}$  (higher  $d_{\text{HB}}$  distances). Figure 4E shows the structures of the lowest energy conformation (largest cluster) for these four sequons. The sequons with  $P_{-1}$  and  $T_{-1}$  fit in the pocket-like cavity in the enzyme's peptide binding groove (Figure 4E, top panels), whereas sequons with  $N_{-1}$  or  $R_{-1}$  were excluded from this cavity due to steric hinderance thereby resulting in larger distances (Figure 4E, bottom panels). Hence, the structural basis for peptides to preferentially sample low-energy conformations compatible with sampling of the proposed TS stabilizing hydrogen bond was the *relative size of the side chain of the amino acid at the -1 position and fit into the "-1 pocket" on the enzyme* (highlighted in Figure 4F; discussed in more detail later).

The case of sequon  $N_{-1}TP_{+1}$  (Also  $V_{-1}TP_{+1}$ ; Figure S8) is also notable because it exhibited two significant clusters, the smaller one (normalized cluster size  $\sim 10\%$ ) exhibiting distances compatible with hydrogen bonding (Figure 4D) and the larger one (normalized cluster size  $\sim 90\%$ ) with comparable interaction energy exhibiting larger distances. Experimentally this sequon was non-glycosylatable. Sequon  $T_{-1}TP_{+1}$ , which is experimentally glycosylatable, also exhibits two significant clusters. However, the larger cluster exhibits a  $d_{\text{HB}}$  distance compatible with hydrogen bonding (Figure 4D). This suggests that a larger fraction of decoys exhibiting  $d_{\text{HB}}$  compatible with TS-stabilizing hydrogen bond renders a sequon more glycosylatable.

For sequons with larger amino acid residues at the -1 position, characterization of the  $d_{\text{HB}}$  distance correlated with undetectable glycosylation in experimental assays as peptides/sequons with larger amino acids did not meet the hydrogen-bonding criteria and assumed conformations at distances farther from the UDP-GalNAc donor, making the reaction less likely. However, this  $d_{\text{HB}}$  criterion incorrectly classified all  $G_{-1}$  sequons. (Figure 4B; blue arrow). Since  $d_{\text{HB}}$  generated some false positives, especially  $G_{-1}$  peptides, the ability of the peptide to assume conformations amenable to the formation of the TS-stabilizing hydrogen bonding is a necessary but not sufficient condition to determine specificity.

**$G_{-1}$  results in distinct low-energy states characterized by higher  $\text{RMSD}_{\text{peptide}}$  values**—To probe why  $G_{-1}$  peptides may be non-glycosylatable even though they satisfy the  $d_{\text{HB}}$  metric, we examined the joint distribution of the  $\text{RMSD}_{\text{peptide}}$  and  $d_{\text{HB}}$  sampled by the top 10% of decoys for all  $G_{-1}$  and  $P_{-1}$  peptides (*i.e.*, averaged over all 19 amino acid at the +1 position) (Figures S9–S13 plots all 19 sequons for  $G_{-1}$ ,  $S_{-1}$ ,  $A_{-1}$ ,  $T_{-1}$  and  $P_{-1}$ .)  $G_{-1}$  peptides exhibited two low-energy states (Figure 5A), but  $P_{-1}$  peptides primarily exhibited a single, low  $\text{RMSD}_{\text{peptide}}$  state (Figure 5B).

**Amino acid residues with smaller side chains are sub-optimal for -1 pocket of the enzyme**—In Figure 5C, we have superposed the lowest-energy decoys for sequons P<sub>-1</sub>TP<sub>+1</sub> and G<sub>-1</sub>TP<sub>+1</sub>, representative of P<sub>-1</sub> and G<sub>-1</sub> peptides, respectively. For G<sub>-1</sub>TP<sub>+1</sub>, the backbone was shifted “up” with respect to that of the P<sub>-1</sub>TP<sub>+1</sub> backbone (Figure 5C; green arrow). For G<sub>-1</sub> peptides, the small size of the glycine residue allowed multiple configurations in the -1 pocket of the enzyme, all of which still made the TS stabilizing hydrogen bond (i.e., < 4.0 Å). Consequently, we also observed the “GTX-like state” (defined as  $\text{RMSD}_{\text{peptide}} \geq 1.0 \text{ \AA}$  and < 4.0 Å and marked in Figure 5 with black arrows) for A<sub>-1</sub> and S<sub>-1</sub> (shorter side chains) peptides (Figure 5D) but not for T<sub>-1</sub> peptides. Instead, T<sub>-1</sub> peptides exhibited a third state (Figure 5D; blue arrow) which we discuss later. Thus, while the  $d_{\text{HB}}$  metric explained why sequons with larger side chains at the -1 position were non-glycosylatable, the  $\text{RMSD}_{\text{peptide}}$  metric explains why certain sequons with smaller side chains at the -1 position may not be suitable for glycosylation.

**RMSD<sub>peptide</sub> metric improves sequon specificity predictions for G<sub>-1</sub> peptides and recapitulates specificity for 90% of the sequons**—P<sub>-1</sub> peptides, irrespective of the amino acid at the +1 position, experimentally exhibited high glycosylation efficiencies and also primarily exhibited the low-RMSD<sub>peptide</sub> or the PTX-like state (defined as  $\text{RMSD}_{\text{peptide}} < 1.0 \text{ \AA}$  and < 4.0 Å and marked in Figure 5 with red arrows). This leads us to hypothesize that besides the TS-stabilizing hydrogen bond (characterized by  $d_{\text{HB}}$ ), the second factor that determined the glycosylatability of a sequon was the precise positioning of the peptide in the enzyme’s peptide binding groove, i.e., how close the peptide backbone was, spatially and conformationally, to the cognate sequon peptide conformation in the crystal structure. We postulated that the “PTX-like state” (red arrows in Figure 5D) with  $\text{RMSD}_{\text{peptide}} < 1.0 \text{ \AA}$  leads to successful glycosylation (reactive state) whereas all other conformations or states with  $\text{RMSD}_{\text{peptide}} \geq 1.0 \text{ \AA}$  (e.g. the GTX-like state) did not lead to glycosylation (non-reactive).

Hence, we used the sampling of the PTX-like state by the top-scoring decoys, quantified by the  $\text{RMSD}_{\text{peptide}}$  of the lowest energy decoy of the largest cluster and the normalized size of the largest cluster, as the second criterion for successful glycosylation that includes multiple low-energy conformations. This criterion improved prediction for sequons that exhibited low-energy conformations in the GTX-like state, (A<sub>-1</sub>TH<sub>+1</sub>, A<sub>-1</sub>TG<sub>+1</sub>, S<sub>-1</sub>TG<sub>+1</sub>, G<sub>-1</sub>TG<sub>+1</sub>, etc.), including G<sub>-1</sub> peptides and was able to correctly classify many such peptides as non-glycosylatable (Figure 5E, F). When compared with experimental results, this criterion, based on the  $\text{RMSD}_{\text{peptide}}$  value of the lowest-energy decoy of the largest cluster as a classification metric for glycosylation gives a ROC-AUC value of 0.959 (Figure 5G, Table 1). Setting a threshold of  $\text{RMSD}_{\text{peptide}} < 1.0 \text{ \AA}$ , the criterion correctly classifies 90.3% ((39+287)/361) of sequons with a BA of 87.9% and an FPR of 8.9% (28/(28+287)) (Figure 5F, Table 1). When the fraction of decoys ( $N_r$ ) that satisfy the  $\text{RMSD}_{\text{peptide}} < 1 \text{ \AA}$  condition is the classification metric, the ROC-AUC value was 0.965 (Figure 5G). Setting an arbitrary threshold value of  $N_r > 0.53$ ,  $N_r$  with  $\text{RMSD}_{\text{peptide}} < 1 \text{ \AA}$  correctly classified 90.8% ((36+292)/361) of the sequons with a BA of 85.5% and an FPR rate of 7.3% (23/(23+293)) (Table 1).

Note that for most conformations that satisfy  $\text{RMSD}_{\text{peptide}} < 1.0 \text{ \AA}$ , the condition  $d_{\text{HB}} < 4.0 \text{ \AA}$  is also satisfied (Figures 5A, B, D).

However, since both classes of sequons, those that are glycosylatable experimentally (e.g.  $A_{-1}TA_{+1}$  and  $S_{-1}TA_{+1}$ , Figure 5F black boxes) and those that are non-glycosylatable experimentally (e.g.  $G_{-1}TA_{+1}$  and  $A_{-1}TH_{+1}$  Figure 5F brown boxes), exhibited non-reactive states, the criterion based on  $\text{RMSD}_{\text{peptide}}$  was not sufficient to correctly classify all peptides, especially for sequons that exhibited both reactive and non-reactive states with similar interaction energies and/or similar fraction of decoys.

**Amino acid residue at the -1 position dictates the low-energy conformations and glycosylatability for the majority of the sequons**—The analysis of the low-energy conformations characterized by  $d_{\text{HB}}$  and  $\text{RMSD}_{\text{peptide}}$  lead to the following observations. For the majority of sequons, those with  $K_{-1}$ ,  $R_{-1}$ ,  $F_{-1}$ ,  $Y_{-1}$ ,  $W_{-1}$ ,  $D_{-1}$ ,  $E_{-1}$ ,  $Q_{-1}$ ,  $N_{-1}$ ,  $H_{-1}$ ,  $I_{-1}$ ,  $M_{-1}$ ,  $L_{-1}$ , or  $V_{-1}$ , the peptide primarily sampled non-reactive low-energy conformations with  $d_{\text{HB}} > 4.0 \text{ \AA}$  and  $\text{RMSD}_{\text{peptide}} > 1.0 \text{ \AA}$ . For a small fraction of sequons ( $P_{-1}$  peptides), the peptide primarily sampled a reactive, cognate-sequon like state (or PTX-like state) with  $d_{\text{HB}} < 4.0 \text{ \AA}$  and  $\text{RMSD}_{\text{peptide}} < 1.0 \text{ \AA}$ . For both of these categories that primarily sample one state—either the non-reactive state or the reactive-state—the computational predictions based on either hypothesis ( $\text{RMSD}_{\text{peptide}}$  or  $d_{\text{HB}}$ ), agreed quite well with experimental data. These observations underscore the importance of the residue at the -1 position in determining the low-energy conformations and, consequently, the glycosylatability for the majority of the sequons ( $\sim 15 \times 19 = 285$  out of 361 peptides). However, four amino acids (G, A, S, T) at the -1 position are yet ambiguous, showing two states. Sometimes they are classified correctly, but sometimes not.

### Recapitulation of amino acid specificity trends for the +1 position

For  $G_{-1}$ ,  $A_{-1}$ ,  $S_{-1}$ ,  $T_{-1}$  sequons ( $4 \times 19 = 76$  out of 361), the peptide sampled both reactive and non-reactive states with comparable interaction energies. For many of these sequons, the computational predictions based on the effect of the -1 position did not accurately recapitulate experimental observations. Hence, for  $G_{-1}$ ,  $A_{-1}$ ,  $S_{-1}$ ,  $T_{-1}$ , to recapitulate experimental glycosylation trends, we must consider the effect of the +1 position.

**Amino acid at the +1 position confers secondary effects that modulate effects of the -1 position**—To investigate the effect of the +1 position for  $G_{-1}$ ,  $A_{-1}$ ,  $S_{-1}$  peptides, all of which exhibit the GTX-like state, we considered the variation in sampling of the GTX-like state as a function of the amino acid at the +1 position. Figure 6A shows these fractions for a subset of sequons, viz. the  $G_{-1}$ ,  $A_{-1}$ ,  $S_{-1}$ , and  $T_{-1}$  peptides. For  $T_{-1}$  peptides, no sequon exhibited the GTX-like state for a significant fraction of the decoys. For  $A_{-1}$ ,  $S_{-1}$  and  $G_{-1}$  peptides,  $G_{+1}$  and  $D_{+1}$  significantly increased the propensity to sample (indicated by a large fraction of decoys) the GTX-like state. Furthermore, for  $A_{-1}$  peptides,  $H_{+1}$ ,  $K_{+1}$ ,  $R_{+1}$ , and  $S_{+1}$  also resulted in a large fraction of decoys exhibiting the GTX-like state. The interaction energies of the lowest-energy decoys of the GTX-like state are comparable to those of the PTX-like state (Figure S14), suggesting that such a state could dominate or compete with the PTX-like. Hence, the +1 position, in these specific cases, enhanced the

sampling of the non-reactive, GTX-like state, modulating the glycosylatability of a peptide in a capacity secondary to the  $-1$  position.

We used these observations to test a classification based on the sampling of the GTX-like state. While the changes in classification accuracy are negligible, the prediction improves for sequons  $A_{-1}TW_{+1}$ ,  $A_{-1}TY_{+1}$ ,  $A_{-1}TT_{+1}$  and  $T_{-1}TG_{+1}$ , summarized in Table S4 and Figure S15.

### **Residues glutamine, glutamate, aspartate and the aromatics at the +1 position interact with residues K363/K281 on the enzyme to form competing states—**

To understand the variation in glycosylatability with the +1 position for  $T_{-1}$  peptides, we examined the sequons  $T_{-1}TQ_{+1}$ ,  $T_{-1}TF_{+1}$ ,  $T_{-1}TY_{+1}$  and  $T_{-1}TW_{+1}$ . Experimentally,  $T_{-1}TQ_{+1}$  was glycosylatable with  $\sim 20\%$  activity, whereas  $T_{-1}TF_{+1}$ ,  $T_{-1}TY_{+1}$  and  $T_{-1}TW_{+1}$  were non-glycosylatable. All four sequons exhibited the PTX-like state (red arrow in Figure 6B and Figure S16A). These sequons additionally exhibited a second, low-energy state with  $RMSD_{peptide} > 1.0$  and  $> 4.0$  Å (blue arrow in Figure 6B and Figure S16A). In this state, the residue  $T_{-1}$  occupied the  $-1$  pocket similar to the PTX-like state, while the residue  $Q_{+1}$  interacted with residue K281 on the enzyme, which lies at the rim of the peptide-binding groove (Figure 6B). The interaction between the residues  $Q_{+1}$  and K281 pulled the peptide backbone away from the catalysis site (Figure 6B and Figure S16B), resulting in a non-reactive state that competed with the reactive PTX-like state. We observed a similar interaction for residue  $D_{+1}$  (sequons  $T_{-1}TD_{+1}$  and  $P_{-1}TD_{+1}$ ), however, due to a shorter side chain compared to  $Q_{+1}$ , it was in a better position to interact with K363 residue (Figure 6C, Figure S16C, D).

In Figure 6D, we show the sampling of the “TTQ-like state” ( $d_{HB} > 4.0$  and  $RMSD_{peptide} > 1$  Å) for  $A_{-1}$ ,  $G_{-1}$ ,  $S_{-1}$ , and  $T_{-1}$  peptides. The TTQ-like state was observed primarily in  $S_{-1}$  and  $T_{-1}$  peptides.  $F_{+1}$ ,  $Y_{+1}$ ,  $W_{+1}$ ,  $M_{+1}$ ,  $I_{+1}$ ,  $E_{+1}$ ,  $Q_{+1}$ ,  $H_{+1}$ , and  $D_{+1}$  exhibited highly stabilized TTQ-like states. For non-polar residues such as  $M_{+1}$ , and  $I_{+1}$ , the stabilization arose from non-polar interactions of the +1 side chain with the K281 side chain.

For sequons that exhibited both states, we computed the energy difference between the lowest-energy decoys for the two states (Figure S17). Similar to the GTX-like state, for many sequons, the interaction energy of the lowest-energy decoys of the TTQ-like state is comparable to that of the PTX-like state. For sequons  $T_{-1}TD_{+1}$ ,  $T_{-1}TW_{+1}$  and  $T_{-1}TY_{+1}$ , the lowest interaction energy of the TTQ-like state was about 2 REU lower than that of the PTX-like state. For  $T_{-1}TQ_{+1}$ , the difference was small ( $-0.2$  REU), and for  $T_{-1}TE_{+1}$ , the PTX-like state was more stable by 2.4 REU. The relative stabilization of the PTX-like state over the TTQ-like state as measured by interaction energy ( $\Delta G_{PTX \rightarrow TTQ}^{int}$ ) correlated with higher experimental glycosylation efficiencies for sequons  $T_{-1}TQ_{+1}$  ( $\sim 22\%$ ) and  $T_{-1}TE_{+1}$  ( $\sim 13\%$ ) compared to  $T_{-1}TD_{+1}$  (3%) and  $T_{-1}TX_{+1}$  (0%), where X was an aromatic residue.

To quantify the interaction energy of different amino acid residues at the +1 position to specific residues on the enzyme, we computed the pairwise energies of interaction between the residue at the +1 position and the enzyme (Figure S18) and, as expected, found that the residues that exhibit the TTQ-like state interact favorably with residues K281 or K363 on the

enzyme. On the other hand, residues  $P_{+1}$  and  $A_{+1}$  did not interact with K281 or K363 residues on the enzyme. The lack of interaction with K281 or K363 residues on the enzyme suggested that sequons  $T_{-1}TP_{+1}$ ,  $T_{-1}TA_{+1}$ ,  $S_{-1}TP_{+1}$  and  $S_{-1}TP_{+1}$  and  $S_{-1}TA_{+1}$  had no propensity for the TTQ-like state and may explain the high glycosylation efficiencies observed for these sequons.

With these observations, we tested using TTQ-like state populations for classifying the substrate and non-substrate sequons (classification accuracy in Table S4 and Figure S19). While the change in classification accuracy is negligible, the consideration of the TTQ-like state improves predictions for sequons  $T_{-1}Y_{+1}$ ,  $T_{-1}TW_{+1}$  and  $T_{-1}TD_{+1}$  (Figure S19).

## Characterization of the peptide-enzyme interface

**Shape complementarity and hydrogen bonding contribute to the finely tuned specificities at the  $-1$  and  $+1$  positions**—So far, our analysis focused on analyzing the landscape of low-energy conformations exhibited by the peptides and on recapitulating the experimentally observed specificity trends as a function of the amino acid at the  $-1$  and  $+1$  position of the sequon. In this process, we discovered the dominant modes of interaction between the peptide and the enzyme that lead to reactive (PTX-like) and non-reactive (GTX-like and TTQ-like) conformations. Comparison between experimental data and computational predictions also revealed that a majority of sequons that were glycosylatable exhibited a PTX-like conformation. Next, we characterize the PTX-like state to decipher the structural basis for the variation of specificity within the subset of peptides that exhibited this state.

First, we calculated the shape complementarity,<sup>44</sup>  $S_c$ , for the enzyme-peptide interface for all sequons for the ten lowest-interaction-energy decoys that satisfied the  $RMSD_{peptide} < 1.0 \text{ \AA}$  criterion (Figure 7A).  $P_{-1}$  peptides exhibited the highest shape complementarity at the peptide-enzyme interface. We further characterized the residue-wise and pairwise interaction energies at the interface (Figure S20). The  $P_{-1}$  residue exhibited generally higher attractive van der Waals energies with all enzyme residues at the interface (Figure S20) and especially with H365 of the enzyme (Figure 7B). The planar interface formed by a histidine residue at position 365 on the enzyme packs well against Pi residue (Figure 7C).  $T_{-1}$ ,  $S_{-1}$ , and  $A_{-1}$  residues shape complementarities and exhibited energies that varied to a significant extent with the residue at the  $+1$  position (Figure 7A, B and Figure S20, S21). Thus for these sequons, the  $+1$  position may additionally contribute to anchoring the peptide in the binding cavity.

In the  $+1$  position, proline exhibited the highest shape complementarity (Figure 7A and Figure S21) in the “ $+1$  pocket” formed by three aromatics F280, W282 and F361, stabilized by favorable interactions between the partially positively charged proline ring and the partially negatively charged  $\pi$  faces of aromatic side chains (Figure 7D, Figure S22). Also, similar to the TTQ-like state, sequons with aromatics, glutamine, glutamate and non-polar residues other than alanine, proline and glycine at the  $+1$  position interacted with K281 on the enzyme (Figure 7E).

Surprisingly, the median  $S_c$  of the top ten decoys that satisfied the  $\text{RMSD}_{\text{peptide}} < 1.0 \text{ \AA}$  criterion is a reasonably good classifier of sequon glycosylatability with a ROC-AUC score of 0.944.

For  $T_{-1}$  and  $S_{-1}$  residues, the PTX-like state was additionally stabilized by a hydrogen bond between the hydroxyl side chain and the backbone carboxyl of R362 on GalNAc-T2 (Figure S23).

In summary, the shape complementarity and pairwise-energies describe a  $-1$  pocket that is highly specific for the  $P_{-1}$  residue, underlying the high experimental glycosylation efficiencies measured for  $P_{-1}$  peptides.

### **Sequence motifs at the $-1$ pocket hint at modes of specificity modulation across isoforms T2, T14 and T16**

The  $-1$  pocket on the enzyme plays an important role in screening for optimally-sized side chains at the  $-1$  position of the sequon. This pocket is primarily formed by residues R362, K363, Q364, H365 and W331. These residues determine the size and chemical composition of the  $-1$  pocket. The residues R362, K363, Q364, and H365 reside on the flexible, semi-conserved catalytic loop<sup>45</sup> of GalNAc-T2. This flap-like loop can additionally contribute to the variability of the  $-1$  pocket size across the GalNAc-T isoforms.<sup>17</sup> Among the three isoforms of the GalNAc-T family that show a strong preference for Pi (T2, T14, T16), the H365 residue is conserved (Figure 7F). Residues K363 and Q364 reside at the point of entry for the  $-1$  residue on the peptide. Variation of amino acids at these positions could allow for variation in the size of the amino acid preferred at the  $-1$  position of the sequon. For example, isoforms T14 and T16, which are evolutionary most proximal to T2, have residues lysine or arginine at position 364. Unlike T2, both T14 and T16 prefer  $G_{-1}$ ,<sup>46</sup> indicative of a  $-1$  pocket suitable for smaller sidechains. In fact, when we repeated MCM sampling of the  $G_{-1}$  sequons for the T2 isoform with the Q364R mutation, we observed a complete shift towards conformations with  $\text{RMSD}_{\text{peptide}} < 1.0 \text{ \AA}$  (PTX-like state) and the elimination of the GTX-like state (Figure 7G), suggesting a possible strategy for varying the peptide substrate preference of various isoforms.

## **Discussion**

In this work, we attempted to understand the structural basis for the peptide substrate preferences of the T2 isoform of the GalNAc-T family. We expect this work to be useful in understanding how the preference for different peptide substrates is modulated across the 20 isoenzymes of this family.

We used a flexible backbone protocol with MCM sampling, resulting in more than one low-energy peptide conformation/state in the vicinity of the starting peptide conformation. Most existing protocols for determining peptide specificity for peptide binding domains employ limited backbone sampling, generating ensembles close to the starting structures (pepspec, MFPred) and usually employing additional constraints to sample TS-like conformations. While these studies have been successful at predicting specificity trends, a wealth of information can be garnered from sampling the peptide landscape without imposed constraints. Our work benefitted from the availability of crystal structures for the peptide-



enzyme complex but may be less accurate in the absence of crystal structures. Our approach also suffered from inaccuracies in the Rosetta energy function, the limitations of MCM sampling, and the use of implicit solvation models to name a few. We further note that an MD-based simulation, though computationally prohibitive for a large dataset, may be better suited for generating thermodynamically accurate ensembles and for characterizing the density of multiple stable states.

We investigated a range of features to predict the glycosylation efficiency of GalNAc-T2 (Table 1) and found that the features  $d_{\text{HB}}$  and  $\text{RMSD}_{\text{peptide}}$  are able to recapitulate binarized glycosylation specificity with a balanced accuracy of 88.9% and 87.9% and a false positive rate of 17.8% and 8.9% respectively. Alternatively, the fraction of decoys,  $N_d$  and  $N_r$ , that satisfy criterion  $d_{\text{HB}} < 4.0 \text{ \AA}$  and  $\text{RMSD}_{\text{peptide}} < 1.0 \text{ \AA}$  respectively, recapitulated specificity with a balanced accuracy of 85.2% and 85.6% and a false positive rate of 14.3% and 7.3% respectively.

Additionally, we found energy-based predictors (based on MFPred and interaction energy in this work) to be poor predictors of specificity, especially in the absence of structural-motif-preserving constraints. While MFPred is able to predict a preference for  $T_{-1}$  and  $S_{-1}$  residues, it fails to identify  $P_{-1}$  and  $A_{-1}$ . For the +1 position, MFPred performs much worse. Both of these AUC values as well as the average AUC over all positions predicted by MFPred are lower than those obtained with the  $d_{\text{HB}}$  and  $\text{RMSD}_{\text{peptide}}$  criteria (Table 1, Table S3). These results suggest that the stability of the peptide-enzyme complex or the interaction-energy at the interface, by itself, is a weak indicator of efficient catalysis by GalNAc-T2. In fact, since selective stabilization of the transition state over the reactants is important for catalysis, the over-stabilization of the reactant state (indicated by higher interaction energies) may increase the free energy of activation (difference between the energy of reactants and the transition state) thereby slowing or preventing the reaction<sup>20,25</sup>. The addition of constraints could partially alleviate this issue by restraining the enzyme-peptide complex in a configuration mimicking the transition state. However, the addition of constraints will omit the sampling of potential low energy states that may compete with or hinder the formation of the transition state. Such states can only be identified in protocols that allow flexible sampling of the backbone without constraints.

The  $-1$  position on the peptide strongly determined the glycosylation efficiency. Residues R362, K363, Q364 and H365 on the catalytic loop and residue W331 on the enzyme form the  $-1$  pocket and select for amino acids threonine, proline, serine, alanine or glycine at the sequon's  $-1$  position. For sequons with residues that did not fit this pocket, the peptide was not able to form a hydrogen bond with UDP that has been proposed to stabilize the TS. We further found that this pocket was especially favorable for recognizing peptides with proline at the  $-1$  position, as demonstrated by highly favorable interactions between H365 and proline and a high degree of shape complementarity at the +1 position irrespective of the amino acid. The flexible catalytic loop is especially suitable for modulating of the size of the peptide binding pocket.<sup>17</sup> Hence, by changing the size and other biophysical aspects of this pocket, the specificity for the  $-1$  position can be potentially modulated. These structural and sequence features are especially relevant for specificity modulation across isoforms, as the GalNAc-T family can glycosylate a wide range of amino acids at the  $-1$  position.

We additionally found that residues K281 and K363 acted as gating residues by interacting with peptide amino acid residues, such as Q<sub>+1</sub> and D<sub>+1</sub>, leading to low-energy states that compete with the reactive state. Hence, the specificity for the +1 position may be modulated by altering the lysine residues at positions 281 and 363 on the enzyme. Similar to the -1 position, such variation in specificity for the +1 position is already observed in the GalNAc-T family as certain isoforms (GalNAc-T1<sup>12</sup>, GalNAc-T14<sup>46</sup>) are capable of efficiently glycosylating D<sub>+1</sub>.

Key structural motifs identified in this work may be important for designing more promiscuous forms of the enzyme or tailored forms with specificities different from those seen in the 20 naturally occurring isoforms. Furthermore, since many members of the GalNAc-T family have been associated with various cancers, the sequence and structural motifs identified in this work may help decipher mutations that cause aberrant glycosylation.

## Methods

### Starting structure for enzyme-peptide complex

The primary starting structure of the enzyme-peptide complex was obtained from the crystal structure of the active conformation of GalNAc-T2 from the crystal structure of the complex (pdb id: 2ffu). Since the sugar is absent from that structure, we used a second GalNAc-T2 structure (pdb id: 4d0z) with bound peptide (mEA2), manganese and UDP-GalNAc-5S. While the sugar bound to UDP in 4d0z has a modification (sulfur instead of oxygen in the ring), it aligns exactly with 2ffu with the additional sugar (Figure S24). To generate the starting structure for each sequon, we used the crystal structure of the complex replacing the mEA2 peptide with A<sub>-2</sub>X<sub>-1</sub>T<sub>0</sub>X<sub>+1</sub>A<sub>+2</sub>P<sub>+3</sub>R<sub>+4</sub>C<sub>+5</sub>, where X is any amino acid residue except cysteine. Residues at positions -1 and +1 (denoted by Xs) were mutated to the target sequon for all 361 sequons studied in the work by Kightlinger *et al.*<sup>12</sup> using Rosetta's MutateResidue mover followed by side chain repacking and minimization using the PackRotamersMover. No backbone motion is allowed at this stage.

### Rosetta protocol for generating decoys

The glycosylation protocol is based on the flexpepdock<sup>38,39</sup> protocol with a few modifications. There are two main stages- 1) Low-resolution sampling with the centroid score (united atom) 2) High-resolution refinement with the all-atom ref2015 score function.<sup>47</sup> In the low-resolution phase, we use simulated annealing for enhanced sampling of the peptide. We vary the temperature from 2.0 to 0.6 in Rosetta temperature units (kT) over 30 Monte Carlo (MC) cycles. For each temperature cycle of simulated annealing, we use 50 inner MC cycles are used for perturbation followed by minimization in rigid body (across enzyme-peptide interface) and torsional (peptide) space. "Small" and "shear" movers from Rosetta are used for torsional sampling of the peptide<sup>48</sup> with rigid body perturbations using the RigidBodyPerturbMover. The final pose from the low-resolution stage is passed to the high-resolution stage. In the high-resolution stage, the attractive and repulsive potential weights are ramped down and up respectively over 10 outer cycles. Similar to the low-resolution stage, we apply rigid body sampling across the enzyme-peptide interface and torsional sampling of the peptide backbone followed by minimization and Metropolis

criterion. Additionally, both rigid body moves (30 cycles) and torsional moves (30 cycles) are accompanied by peptide side chain repacking every cycle and the interface side chains every 3<sup>rd</sup> cycle<sup>48</sup>. We used the default distance of 8 Å to define the interface. Additionally, the run was terminated if the peptide moved more than 8 Å away from the enzyme-peptide interface. The backbone of the enzyme is fixed throughout sampling. We generated 2000 decoys per sequon. Larger number of decoys (8000) were not found to alter the results.

This protocol is available in the Rosetta software suite (revision $\geq$ 275). See Supporting Information for the complete list of steps to run the protocol. The protocol is run from commandline as follows:

```
>mucintypeglycosylation.<system><compiler><mode> @flags
```

```
System=linux; compiler=gcc; mode=release
```

Where “flags” is a plain text file and contains the following options:

```
-in:file:s <input pdb file>
-in:fde:native <input pdb file>
-nstruct 2000 #no. of decoys
-residue_to_glycosylate 3P #Threonine on peptide chain P
-substrate_type peptide
-low_res true #enables low resolution stage
-tree_type docking
-sugardonor_residue 495 #residue number of UDP-5SGalNac
-enable_backbone_moves_pp #enables peptide backbone moves in high resolution
-ex1
-ex2aro
-nevery_interface 3 # pack enzyme peptide interface every 3 MC cycles in high resolution
-ntotal_backbone 30 # mn 30 MC cycles in High resolution
-output_distance_metrics true #output rmsd, distance, interaction energies to score file
```

### Clustering and analysis of decoys

The top 10% decoys (200/2000) were clustered using the dbscan clustering algorithm<sup>49,50</sup> in sklearn<sup>51</sup> with parameters set to  $\text{eps} = 0.3 \text{ \AA}$  (maximum distance between samples for one to be considered in the neighborhood of the other) and  $\text{min\_samples} = 10$  (number of samples in the neighborhood of a point to be considered a core point).

## Calculation of features

We report two RMSD metrics in this work – RMSD<sub>peptide</sub> and RMSD<sub>sequon</sub>. Both metrics are calculated over backbone C<sub>α</sub> atoms only with respect to the backbone of the peptide in the starting structure. For RMSD<sub>peptide</sub>, RMSD is calculated over all peptide positions (8). For RMSD<sub>sequon</sub>, RMSD is calculated for positions -1 to +3 (XTXAP). Shape complementarity<sup>44</sup> is calculated using PyRosetta<sup>52</sup> as described in Supporting Information. The interaction energy at the enzyme-peptide interface is calculated as the difference between the ref2015 score for the bound complex and the ref2015 score for the enzyme (includes the UDP-sugar molecule) and the peptide, separated from the complex without relaxing or repacking side chains.

## Specificity prediction with MFPred

We used MFPred as described:<sup>35</sup> 1) The starting structure was relaxed. 2) The lowest energy decoy from relax step was used as the starting structure for the FastRelax protocol for each sequon. 3) The lowest energy decoy for each sequon from the FastRelax protocol was processed by the GenMeanFieldMover. All calculations were performed as described in study<sup>35</sup> with Rosetta software suite<sup>43</sup> (revision 226).

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgements

We thank Tyler J. Stewart for critically reading the manuscript. We thank Nadine L. Samara at National Institutes of Health for helpful discussions and for critically reading and editing the manuscript. We also thank Weston Kightlinger, Liang Lin and Michael C. Jewett at Northwestern University for sharing data and helpful discussions.

## Funding

This work was supported by National Science Foundation Grant DBI-1659649 (to Y.S.), National Institutes of Health Grants (R01-GM127578 to J.J.G and M.P.D, R01-GM078221 to J.J.G and R01-GM137314 to M.P.D), Defense Threat Reduction Agency (HDTRA1-15-10052/P00001 and HDTRA1-20-1-0004 to M.P.D.), and the National Science Foundation Grant CBET-1605242 (to M.P.D.).

## References

- (1). Steen P. Van den; Rudd PM; Dwek RA; Opdenakker G Concepts and Principles of O-Linked Glycosylation. Crit. Rev. Biochem. Mol. Biol 1998, 33, 151–208. 10.1080/10409239891204198. [PubMed: 9673446]
- (2). Brockhausen I; Stanley P Chapter 10 O-GalNAc Glycans. Essentials Glycobiol. 2017, 1, 1–9. 10.1101/glycobiology.3e010.
- (3). Sletmoen M; Gerken TA; Stokke BT; Burchell J; Brewer CF Tn and STn Are Members of a Family of Carbohydrate Tumor Antigens That Possess Carbohydrate-Carbohydrate Interactions. Glycobiology 2018, 28, 437–442. 10.1093/glycob/cwy032. [PubMed: 29618060]
- (4). Kudelka MR; Ju T; Heimburg-Molinaro J; Cummings RD Simple Sugars to Complex Disease—Mucin-Type O-Glycans in Cancer. In Advances in Cancer Research; Academic Press Inc., 2015; Vol. 126, pp 53–135. 10.1016/bs.acr.2014.11.002. [PubMed: 25727146]
- (5). Ten Hagen KG; Fritz TA; Tabak LA All in the Family: The UDP-GalNAc:Polypeptide N-Acetylgalactosaminyltransferases. Glycobiology 2003, 13, 1R–16. 10.1093/glycob/cwg007. [PubMed: 12634318]

- (6). Gerken TA; Raman J; Fritz TA; Jamison O Identification of Common and Unique Peptide Substrate Preferences for the UDP-GalNAc:Polypeptide  $\alpha$ -N-Acetylgalactosaminyltransferases T1 and T2 Derived from Oriented Random Peptide Substrates. *J. Biol. Chem* 2006, 281, 32403–32416. 10.1074/jbc.M605149200. [PubMed: 16912039]
- (7). Gerken TA; Ten Hagen KG; Jamison O Conservation of Peptide Acceptor Preferences between *Drosophila* and Mammalian Polypeptide-GalNAc Transferase Ortholog Pairs. *Glycobiology* 2008, 18, 861–870. 10.1093/glycob/cwn073. [PubMed: 18669915]
- (8). Perrine CL; Ganguli A; Wu P; Bertozzi CR; Fritz TA; Raman J; Tabak LA; Gerken TA Glycopeptide-Preferring Polypeptide GalNAc Transferase 10 (PpGalNAc T10), Involved in Mucin-Type O-Glycosylation, Has a Unique GalNAc-O-Ser/Thr-Binding Site in Its Catalytic Domain Not Found in PpGalNAc T1 or T2. *J. Biol. Chem* 2009, 284, 20387–20397. 10.1074/jbc.M109.017236. [PubMed: 19460755]
- (9). Bennett EP; Mandel U; Clausen H; Gerken TA; Fritz TA; Tabak LA Control of Mucin-Type O-Glycosylation: A Classification of the Polypeptide GalNAc-Transferase Gene Family. *Glycobiology* 2012, 22, 736–756. 10.1093/glycob/cwr182. [PubMed: 22183981]
- (10). Fritz TA; Raman J; Tabak LA Dynamic Association between the Catalytic and Lectin Domains of Human UDP-GalNAc:Polypeptide  $\alpha$ -N-Acetylgalactosaminyltransferase-2. *J. Biol. Chem* 2006, 281, 8613–8619. 10.1074/JBC.M513590200. [PubMed: 16434399]
- (11). Lira-Navarrete E; De Las Rivas M; Compañón T; Pallarés MC; Kong Y; Iglesias-Fernández J; Bernardes GJL; Peregrina JM; Rovira C; Bernadó P; Bruscolini P; Clausen H; Lostao A; Corzana F; Hurtado-Guerrero R Dynamic Interplay between Catalytic and Lectin Domains of GalNAc-Transferases Modulates Protein O-Glycosylation. *Nat. Commun* 2015, 6. 10.1038/ncomms7937.
- (12). Kightlinger W; Lin L; Rosztoczy M; Li W; Delisa MP; Mrksich M; Jewett MC Design of Glycosylation Sites by Rapid Synthesis and Analysis of Glycosyltransferases Article. *Nat. Chem. Biol* 2018, 14, 627–635. 10.1038/s41589-018-0051-2. [PubMed: 29736039]
- (13). Gerken TA; Jamison O; Perrine CL; Collette JC; Moinova H; Ravi L; Markowitz SD; Shen W; Patel H; Tabak LA Emerging Paradigms for the Initiation of Mucin-Type Protein O-Glycosylation by the Polypeptide GalNAc Transferase Family of Glycosyltransferases. *J. Biol. Chem* 2011, 286, 14493–14507. 10.1074/jbc.M111.218701. [PubMed: 21349845]
- (14). Gómez H; Rojas R; Patel D; Tabak LA; Lluch JM; Masgrau L A Computational and Experimental Study of O-Glycosylation. Catalysis by Human UDP-GalNAc Polypeptide:GalNAc Transferase-T2. *Org. Biomol. Chem* 2014, 12, 2645–2655. 10.1039/c3ob42569j. [PubMed: 24643241]
- (15). Trnka T; Kozmon S; Tvaroška I; Koča J Stepwise Catalytic Mechanism via Short-Lived Intermediate Inferred from Combined QM/MM MERP and PES Calculations on Retaining Glycosyltransferase PpGalNAcT2. *PLoS Comput. Biol* 2015, 11. 10.1371/journal.pcbi.1004061.
- (16). Lira-Navarrete E; Iglesias-Fernández J; Zandberg WF; Compañón L; Kong Y; Corzana F; Pinto BM; Clausen H; Peregrina JM; Vocadlo DJ; Rovira C; Hurtado-Guerrero R Substrate-Guided Front-Face Reaction Revealed by Combined Structural Snapshots and Metadynamics for the Polypeptide *N*-Acetylgalactosaminyltransferase 2. *Angew. Chemie Int. Ed* 2014, 53, 8206–8210. 10.1002/anie.201402781.
- (17). de las Rivas M; Paul Daniel EJ; Narimatsu Y; Compañón F; Kato K; Hermosilla P; Thureau A; Ceballos-Laita L; Coelho H; Bernadó P; Marcelo F; Hansen L; Maeda R; Lostao A; Corzana F; Clausen H; Gerken TA; Hurtado-Guerrero R Molecular Basis for Fibroblast Growth Factor 23 O-Glycosylation by GalNAc-T3. *Nat. Chem. Biol* 2020, 16, 351–360. 10.1038/s41589-019-0444-x. [PubMed: 31932717]
- (18). Fernandez AJ; Daniel EJP; Mahajan SP; Gray JJ; Gerken TA; Tabak LA; Samara NL The Structure of the Colorectal Cancer-Associated Enzyme GalNAc-T12 Reveals How Nonconserved Residues Dictate Its Function. *Proc. Natl. Acad. Sci* 2019, 116, 20404–20410. 10.1073/pnas.1902211116. [PubMed: 31548401]
- (19). Tarp MA; Clausen H Mucin-Type O-Glycosylation and Its Potential Use in Drug and Vaccine Development. *Biochim. Biophys. Acta - Gen. Subj* 2008, 1780, 546–563. 10.1016/j.bbagen.2007.09.010.

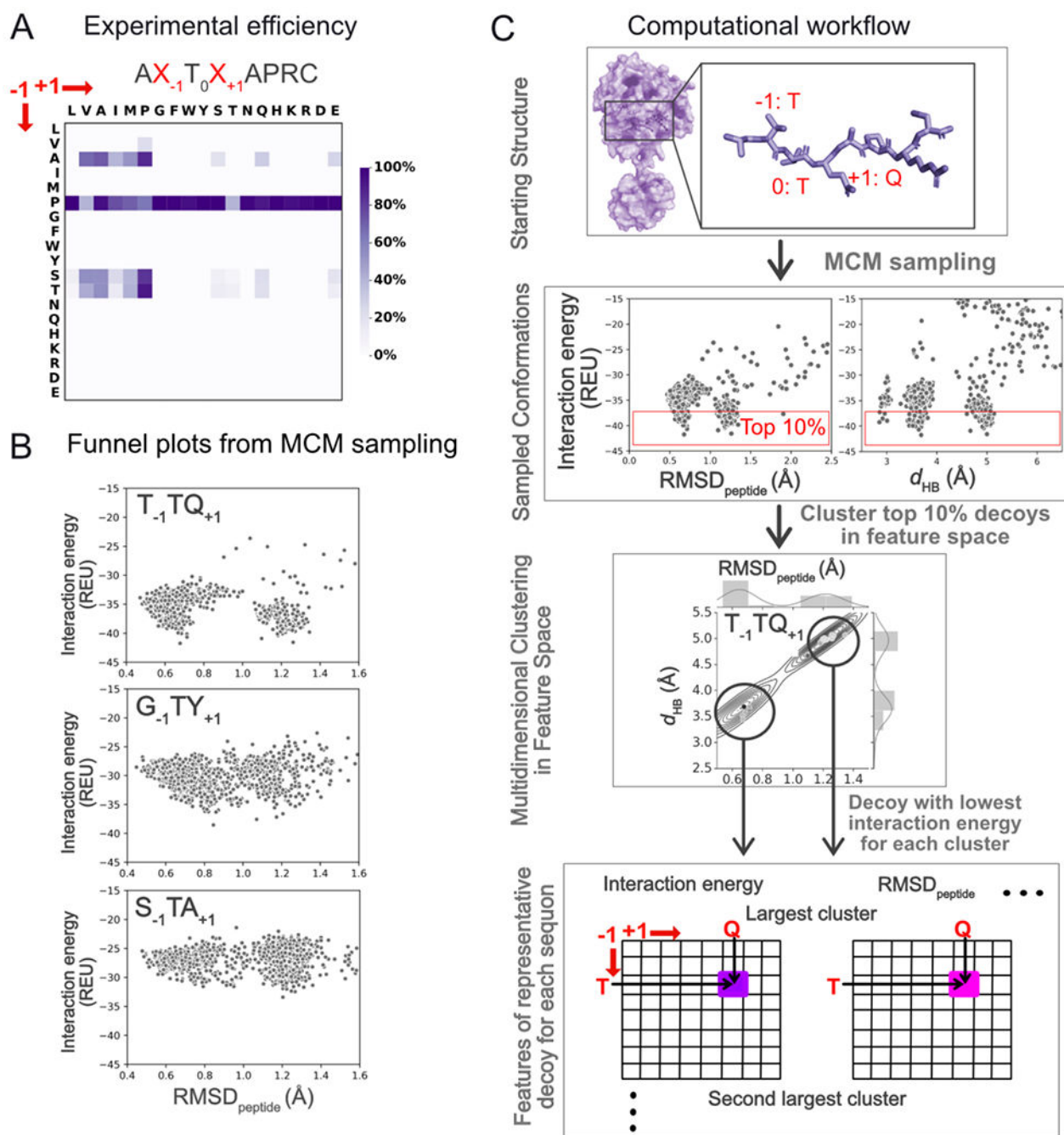
- (20). Mak WS; Siegel JB Computational Enzyme Design: Transitioning from Catalytic Proteins to Enzymes. *Curr. Opin. Struct. Biol* 2014, 27, 87–94. 10.1016/j.sbi.2014.05.010. [PubMed: 25005925]
- (21). Kundert K; Kortemme T Computational Design of Structured Loops for New Protein Functions. *Biol. Chem* 2019, 400, 275–288. 10.1515/hsz-2018-0348. [PubMed: 30676995]
- (22). King CA; Bradley P Structure-Based Prediction of Protein-Peptide Specificity in Rosetta. *Proteins Struct. Fund. Bioinforma* 2010, 78, 3437–3449. 10.1002/prot.22851.
- (23). Frushicheva MP; Cao J; Warshel A Challenges and Advances in Validating Enzyme Design Proposals: The Case of Kemp Eliminase Catalysis. *Biochemistry* 2011, 50, 3849–3858. 10.1021/bi200063a. [PubMed: 21443179]
- (24). Mayes HB; Knott BC; Crowley MF; Broadbelt LJ; Ståhlberg J; Beckham GT Who's on Base? Revealing the Catalytic Mechanism of Inverting Family 6 Glycoside Hydrolases. *Chem. Sci* 2016, 7, 5955–5968. 10.1039/c6sc00571c. [PubMed: 30155195]
- (25). Pauling L Molecular Architecture and Biological Reactions. *Chem. Eng. News* 1946, 24, 1375–1377. 10.1021/cen-v024n010.p1375.
- (26). Leaver-Fay A; Jacak R; Stranges PB; Kuhlman B A Generic Program for Multistate Protein Design. *PLoS One* 2011, 6, e20937. 10.1371/journal.pone.0020937. [PubMed: 21754981]
- (27). St-Jacques AD; Eyahpaise ve C.; Chica RA Computational Design of Multi substrate Enzyme Specificity. 2019, 14, 15. 10.1021/acscatal.9b01464.
- (28). Antoniou D; Schwartz SD Protein Dynamics and Enzymatic Chemical Barrier Passage. *J. Phys. Chem. B* 2011, 115, 15147–15158. 10.1021/jp207876k. [PubMed: 22031954]
- (29). Boehr DD; Nussinov R; Wright PE The Role of Dynamic Conformational Ensembles in Biomolecular Recognition. *Nat. Chem. Biol* 2009, 5, 789–796. 10.1038/nchembio.232.
- (30). Brouk M; Derry N-L; Shainsky J; Ben-Barak Zelas Z; Boyko Y; Dabush K; Fishman A The Influence of Key Residues in the Tunnel Entrance and the Active Site on Activity and Selectivity of Toluene-4-Monooxygenase. *J. Mol. Catal. B Enzym* 2010, 66, 72–80. 10.1016/j.molcatb.2010.03.006.
- (31). Lee H-L; Chang C-K; Jeng W-Y; Wang AH-J; Liang P-H Mutations in the Substrate Entrance Region of  $\beta$ -Glucosidase from *Trichoderma Reesei* Improve Enzyme Activity and Thermostability. *Protein Eng. Des. Sel* 2012, 25, 733–740. 10.1093/protein/gzs073. [PubMed: 23077275]
- (32). Li Z; Scheraga HA Monte Carlo-Minimization Approach to the Multiple-Minima Problem in Protein Folding. *Proc. Natl. Acad. Sci. U. S. A* 1987, 84, 6611–6615. 10.1073/pnas.84.19.6611. [PubMed: 3477791]
- (33). Smith CA; Kortemme T Predicting the Tolerated Sequences for Proteins and Protein Interfaces Using RosettaBackrub Flexible Backbone Design. *PLoS One* 2011, 6. 10.1371/journal.pone.0020451.
- (34). Smith CA; Kortemme T Structure-Based Prediction of the Peptide Sequence Space Recognized by Natural and Synthetic PDZ Domains. *J. Mol. Biol* 2010, 402, 460–474. 10.1016/J.JMB.2010.07.032. [PubMed: 20654621]
- (35). Rubenstein AB; Pethe MA; Khare SD MFpred: Rapid and Accurate Prediction of Protein-Peptide Recognition Multispecificity Using Self-Consistent Mean Field Theory. *PLoS Comput. Biol* 2017, 13, e1005614. 10.1371/journal.pcbi.1005614. [PubMed: 28650961]
- (36). Chaudhury S; Gray JJ Identification of Structural Mechanisms of HIV-1 Protease Specificity Using Computational Peptide Docking: Implications for Drug Resistance. *Structure* 2009, 77, 1636–1648. 10.1016/j.str.2009.10.008.
- (37). Zheng F; Jewell H; Fitzpatrick J; Zhang J; Mierke DF; Grigoryan G Computational Design of Selective Peptides to Discriminate between Similar PDZ Domains in an Oncogenic Pathway. *J. Mol. Biol* 2015, 427, 491–510. 10.1016/j.jmb.2014.10.014. [PubMed: 25451599]
- (38). Raveh B; London N; Zimmerman L; Schueler-Furman O Rosetta FlexPepDock Ab-Initio: Simultaneous Folding, Docking and Refinement of Peptides onto Their Receptors. *PLoS One* 2011, 6, e18934. 10.1371/journal.pone.0018934. [PubMed: 21572516]



- (39). Raveh B; London N; Schueler-Furman O Sub-Angstrom Modeling of Complexes between Flexible Peptides and Globular Proteins. *Proteins Struct. Funct. Bioinforma* 2010, 78, 2029–2040. 10.1002/prot.22716.
- (40). Pethe MA; Rubenstein AB; Khare SD Large-Scale Structure-Based Prediction and Identification of Novel Protease Substrates Using Computational Protein Design. *J. Mol. Biol* 2017, 429, 220–236. 10.1016/j.jmb.2016.11.031. [PubMed: 27932294]
- (41). Pethe MA; Rubenstein AB; Khare SD Data-Driven Supervised Learning of a Viral Protease Specificity Landscape from Deep Sequencing and Molecular Simulations. *Proc. Natl. Acad. Sci. U. S. A* 2019, 116, 168–176. 10.1073/pnas.1805256116. [PubMed: 30587591]
- (42). Leaver-Fay A; Tyka M; Lewis SM; Lange OF; Thompson J; Jacak R; Kaufman KW; Renfrew PD; Smith CA; Sheffler W; Davis IW; Cooper S; Treuille A; Mandell DJ; Richter F; Ban Y-EA; Fleishman SJ; Com JE; Kim DE; Lyskov S; Berrondo M; Mentzer S; Popović Z; Havranek JJ; Karanicolas J; Das R; Meiler J; Kortemme T; Gray JJ; Kuhlman B; Baker D; Bradley P Rosetta3. In *Methods in Enzymology*; Academic Press Inc., 2011; Vol. 487, pp 545–574. 10.1016/B978-0-12-381270-4.00019-6. [PubMed: 21187238]
- (43). Leman JK; Weitzner BD; Lewis SM; Adolf-Bryfogle J; Alam N; Alford RF; Aprahamian M; Baker D; Barlow KA; Barth P; Basanta B; Bender BJ; Blacklock K ; Bonet J; Boyken SE; Bradley P; Byströff C; Conway P; Cooper S; Correia BE; Coventry B; Das R; De Jong RM; DiMaio F; Dsilva L; Dunbrack R; Ford AS; Frenz B; Fu DY; Geniesse C; Goldschmidt L; Gowthaman R; Gray JJ; Gront D; Guffy S; Horowitz S; Huang PS; Huber T; Jacobs TM; Jeliakov JR; Johnson DK; Kappel K; Karanicolas J; Khakzad H; Khar KR; Khare SD; Khatib F; Khrumushin A; King IC; Kleffner R; Koepnick B; Kortemme T; Kuenze G; Kuhlman B; Kuroda D; Labonte JW; Lai JK; Lapidoth G; Leaver-Fay A; Lindert S; Linsky T; London N; Lubin JH; Lyskov S; Maguire J; Malmström L; Marcos E; Marcu O; Marze NA; Meiler J; Moretti R; Mulligan VK; Nerli S; Norm C; Ó'Conchúir S; Ollikainen N; Ovchinnikov S; Pacella MS; Pan X; Park H; Pavlovicz RE; Pethe M; Pierce BG; Pilla KB; Raveh B; Renfrew PD; Burman SSR; Rubenstein A; Sauer MF; Scheck A; Schief W; Schueler-Furman O; Sedan Y; Sevy AM; Sgourakis NG; Shi L; Siegel JB; Silva DA; Smith S; Song Y; Stein A; Szegedy M; Teets FD; Thyme SB; Wang RYR; Watkins A; Zimmerman L; Bonneau R Macromolecular Modeling and Design in Rosetta: Recent Methods and Frameworks. *Nat. Methods* 2020, 77, 665–680. 10.1038/s41592-020-0848-2.
- (44). Lawrence MC; Colman PM Shape Complementarity at Protein/Protein Interfaces. *J. Mol. Biol* 1993, 234, 946–950. 10.1006/jmbi.1993.1648. [PubMed: 8263940]
- (45). De Las Rivas M; Paul Daniel EJ; Coelho H; Lira-Navarrete E; Raich L; Compañón I; Diniz A; Lagartera L; Jiménez-Barbero J; Clausen H; Rovira C; Marcelo F; Corzana F; Gerken TA; Hurtado-Guerrero R Structural and Mechanistic Insights into the Catalytic-Domain-Mediated Short-Range Glycosylation Preferences of GalNAc-T4. *ACS Cent. Sci* 2018, 4, 1274–1290. 10.1021/acscentsci.8b00488. [PubMed: 30276263]
- (46). de las Rivas M; Lira-Navarrete E; Gerken TA; Hurtado-Guerrero R Polypeptide GalNAc-Ts: From Redundancy to Specificity. *Curr. Opin. Struct. Biol* 2019, 56, 87–96. 10.1016/J.SBI.2018.12.007. [PubMed: 30703750]
- (47). Alford RF; Leaver-Fay A; Jeliakov JR; O'Meara MJ; DiMaio FP; Park H; Shapovalov MV; Renfrew PD; Mulligan VK; Kappel K; Labonte JW; Pacella MS; Bonneau R; Bradley P; Dunbrack RL; Das R; Baker D; Kuhlman B; Kortemme T; Gray JJ The Rosetta All-Atom Energy Function for Macromolecular Modeling and Design. *J. Chem. Theory Comput* 2017, 13, 3031–3048. 10.1021/acs.jctc.7b00125. [PubMed: 28430426]
- (48). Rohl CA; Strauss CEM; Misura KMS; Baker D Protein Structure Prediction Using Rosetta. *Methods Enzymol.* 2004, 383, 66–93. 10.1016/S0076-6879(04)83004-0. [PubMed: 15063647]
- (49). Ester M; Kriegl H-P; Sander J; Xu X A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise. In *Proc. 2nd Int. Conf. on Knowledge Discovery and Data Mining*. Portland, OR; 1996; pp 226–231.
- (50). Schubert E; Sander J; Ester M; Kriegel HP; Xu X DBSCAN Revisited, Revisited: Why and How You Should (Still) Use DBSCAN. *ACM Trans. Database Syst* 2017, 42, 1–21. 10.1145/3068335.
- (51). Pedregosa F; Varoquaux G; Gramfort A; Michel V; Thirion B; Grisel O; Blondel M; Prettenhofer P; Weiss R; Dubourg V; Vanderplas J; Passos A; Cournapeau D; Brucher M; Perrot M;

Duchesnay E Scikit-Learn: Machine Learning in Python. *J. Mach. Learn. Res* 2011, 12, 2825–2830.

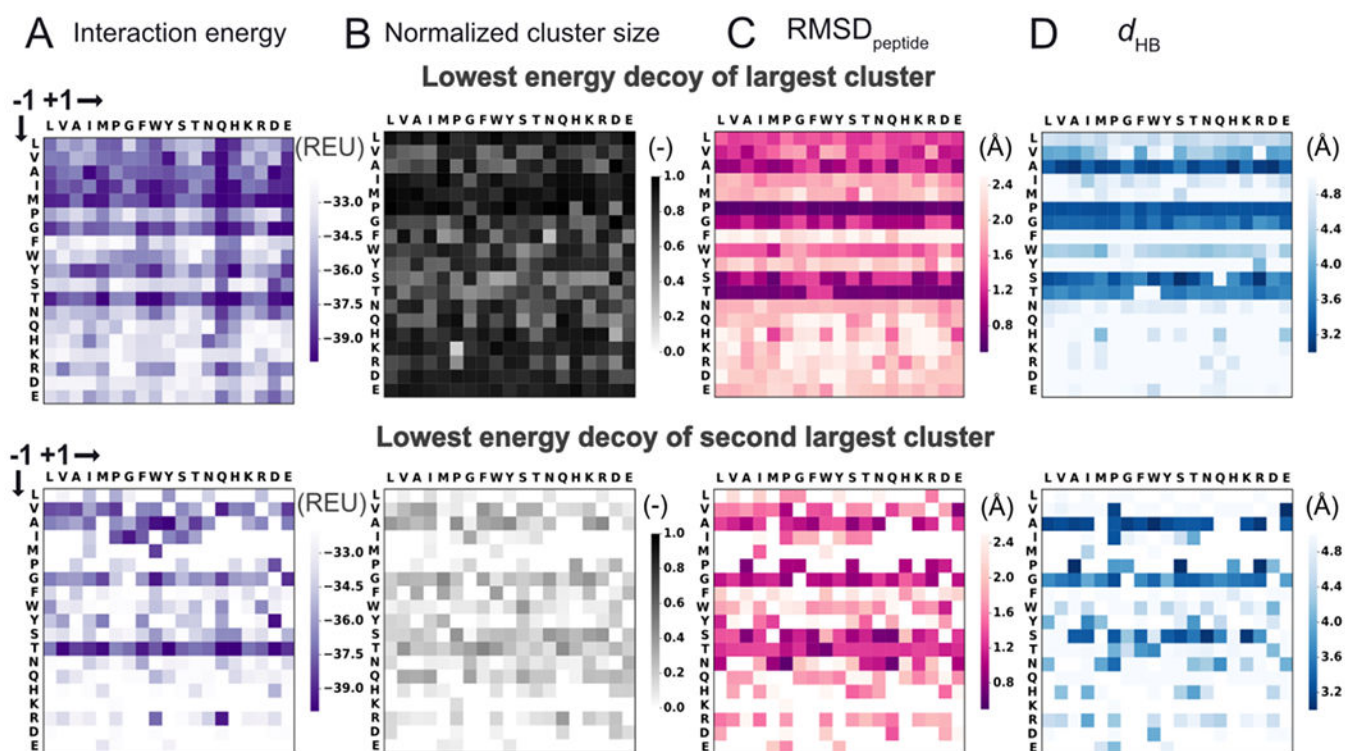
- (52). Chaudhury S; Lyskov S; Gray JJ PyRosetta: A Script-Based Interface for Implementing Molecular Modeling Algorithms Using Rosetta. *Bioinformatics* 2010, 26, 689–691. 10.1093/bioinformatics/btq007. [PubMed: 20061306]



**Figure 1. Computational workflow to determine glycosylation efficiency of GalNAc-T2 for peptide substrates for an experimentally characterized dataset obtained by scanning 19 amino acid residues (all except cysteine) at positions  $-1$  and  $+1$  of the peptide.**

(A) For reference, a replot of the experimentally determined efficiencies (data from Kightlinger *et al.*<sup>12</sup>). (B) Monte Carlo minimization (MCM) sampling of peptides docked to GalNAc-T2 result in “funnel plots” like the three shown. Each point represents one structural model, or “decoy,” at its corresponding RMSD from the reference structure and the interaction energy calculated by Rosetta. (C) Computational workflow to characterize enzyme-peptide interactions for a representative sequon, T<sub>-1</sub>TQ<sub>+1</sub>, with T at the  $-1$  position

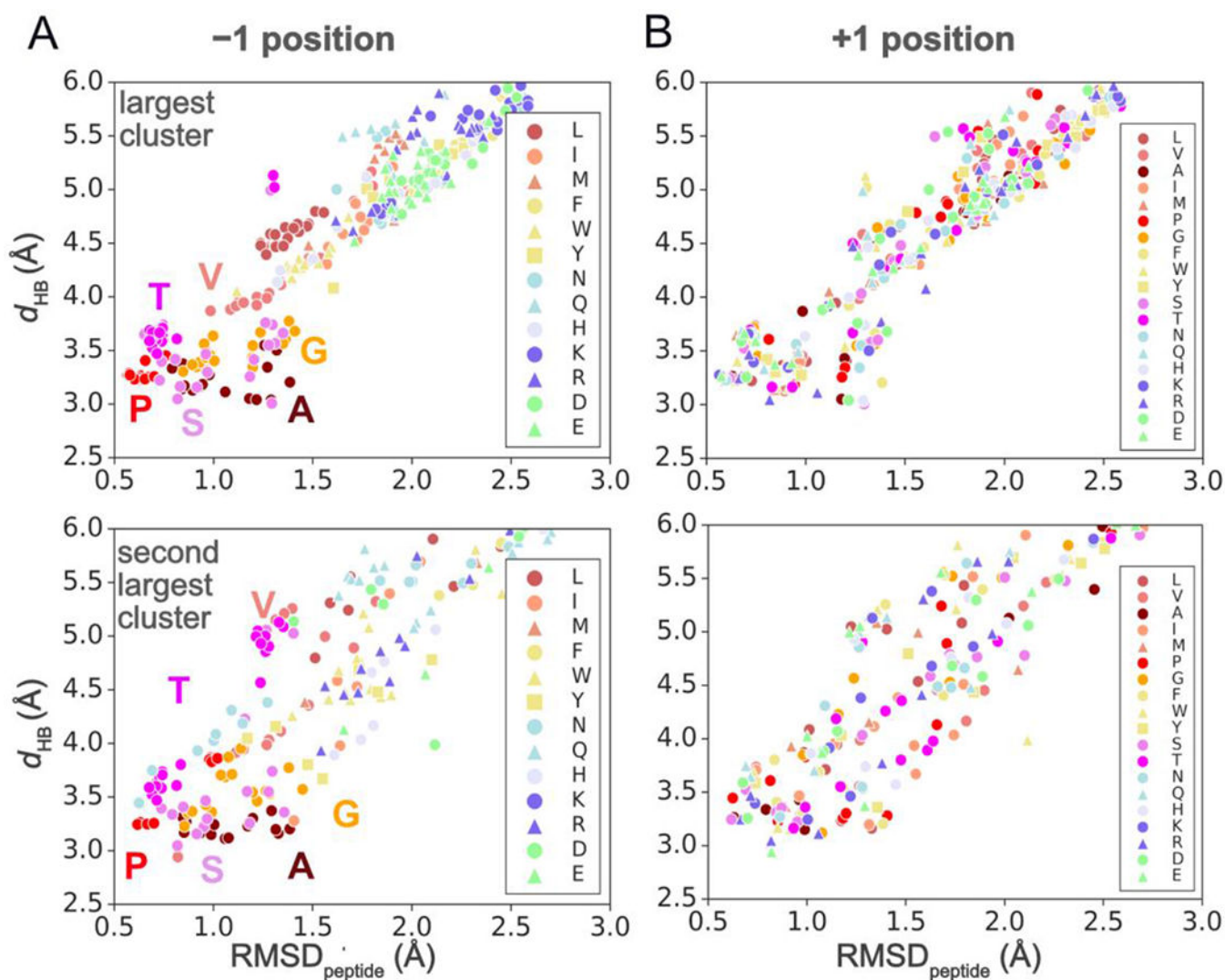
and Q at the +1 position. For each sequon, we selected the top-10%-scoring decoys (by interaction energy) from MCM sampling and clustered them using three features (see main text for features). For each sequon, we characterized the two largest clusters by the lowest interaction energy decoy belonging to that cluster and then examined which features (or combinations of features) could recapitulate the experimental glycosylation efficiencies (A).



**Figure 2. Characterization of the lowest-energy representative conformation for the top two clusters in Rosetta runs.**

As in Figure 1A, each heatmap shows results for the 19x19 peptide sequons on a color scale.

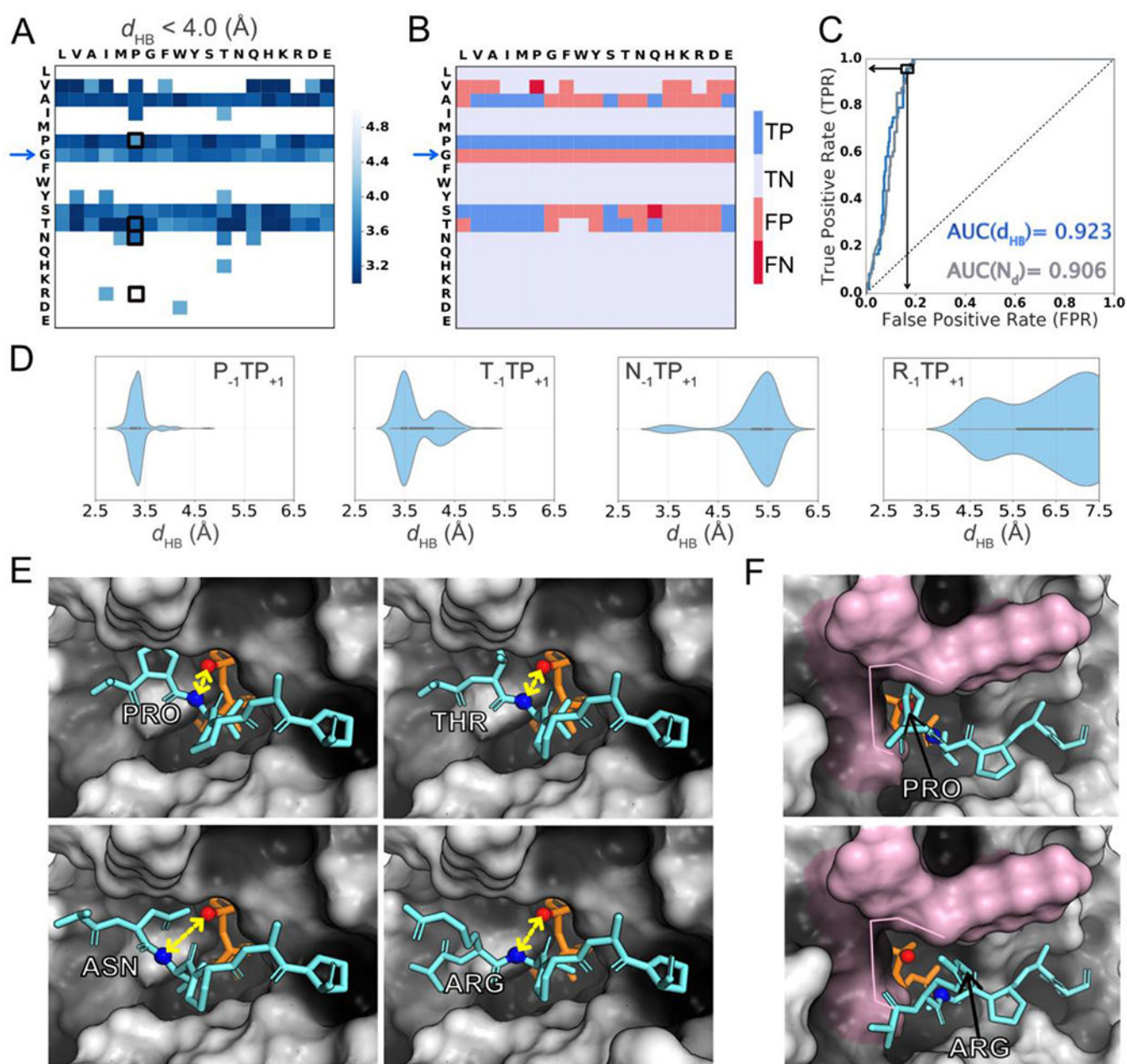
(A) Interaction energy, (B) normalized cluster size, (C)  $\text{RMSD}_{\text{peptide}}$ , and (D)  $d_{\text{HB}}$  of the largest (top) and second-largest (bottom) clusters characterized by the lowest interaction energy decoy for each cluster.



**Figure 3. The major determinant of  $d_{\text{HB}}$  and RMSD<sub>peptide</sub> sampled by lowest energy decoys is the amino acid residue at the -1 position.**

Lowest energy decoys belonging to the largest (top) and second largest clusters (bottom) for all sequons plotted as a function of the RMSD<sub>peptide</sub> and  $d_{\text{HB}}$  and colored by (A) the residue at the -1 position of the sequon and (B) the residue at the +1 position of the sequon.





**Figure 4. Substrate specificity based on TS stabilizing hydrogen bond criterion with  $d_{\text{HB}} < 4 \text{ \AA}$ .** (A) Heatmaps of (left panel)  $d_{\text{HB}}$  distances of the lowest-interaction-energy decoy belonging to a cluster with the cluster centroid satisfying the criterion. (B) True positives (TP, dark blue), true negatives (TN, light blue), false positives (FP, light red) and false negatives (FN, dark red) predicted based on  $d_{\text{HB}} < 4 \text{ \AA}$  threshold applied to the  $d_{\text{HB}}$  value of the lowest-interaction-energy decoy of the largest cluster. (C) ROC curve for  $d_{\text{HB}}$  distances and fraction of decoys ( $N_d$ ) satisfying criterion. (D) Violinplot of distribution of  $d_{\text{HB}}$  distances sampled by the top-scoring 10% decoys for four representative sequons ( $P_{-1}TP_{+1}$ ,  $T_{-1}TP_{+1}$ ,  $N_{-1}TP_{+1}$  and  $R_{-1}TP_{+1}$ ). (E) Lowest interaction energy decoys for four sequons ( $P_{-1}TP_{+1}$ ,  $T_{-1}TP_{+1}$ ,  $N_{-1}TP_{+1}$  and  $R_{-1}TP_{+1}$  – black boxes in the heatmap in (A)),  $d_{\text{HB}}$  is calculated between the

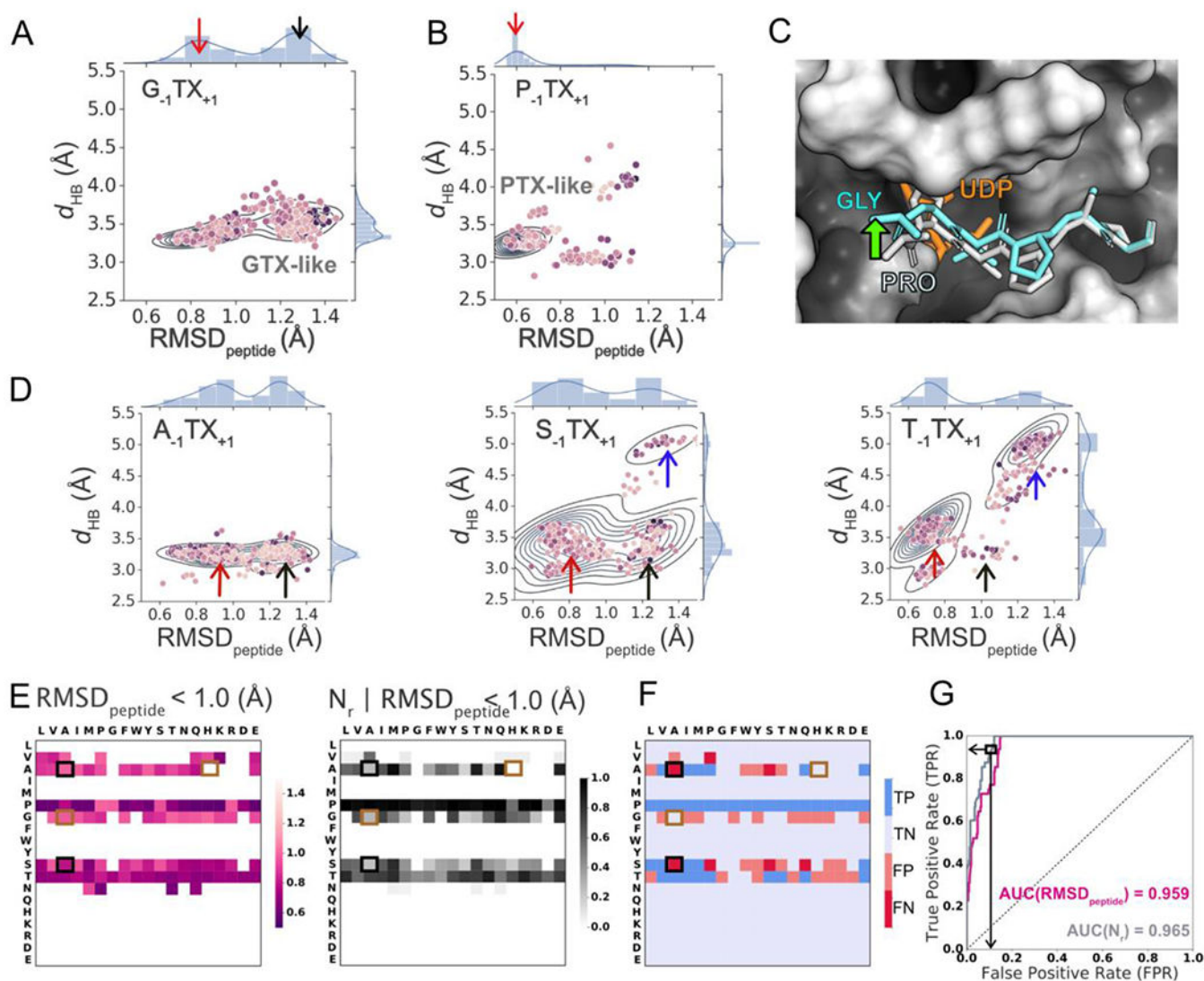
amide nitrogen (blue sphere) of  $T_0$  on peptide(aquamarine) and the  $O_{\beta-PO4}$  (red sphere) on UDP (orange),  $d_{HB}$  is shown with double-ended yellow arrows. (F) Pocket-like cavity formed by enzyme residues (pink surface) that contacts the amino acid at the -1 position on the peptide.

Author Manuscript

Author Manuscript

Author Manuscript

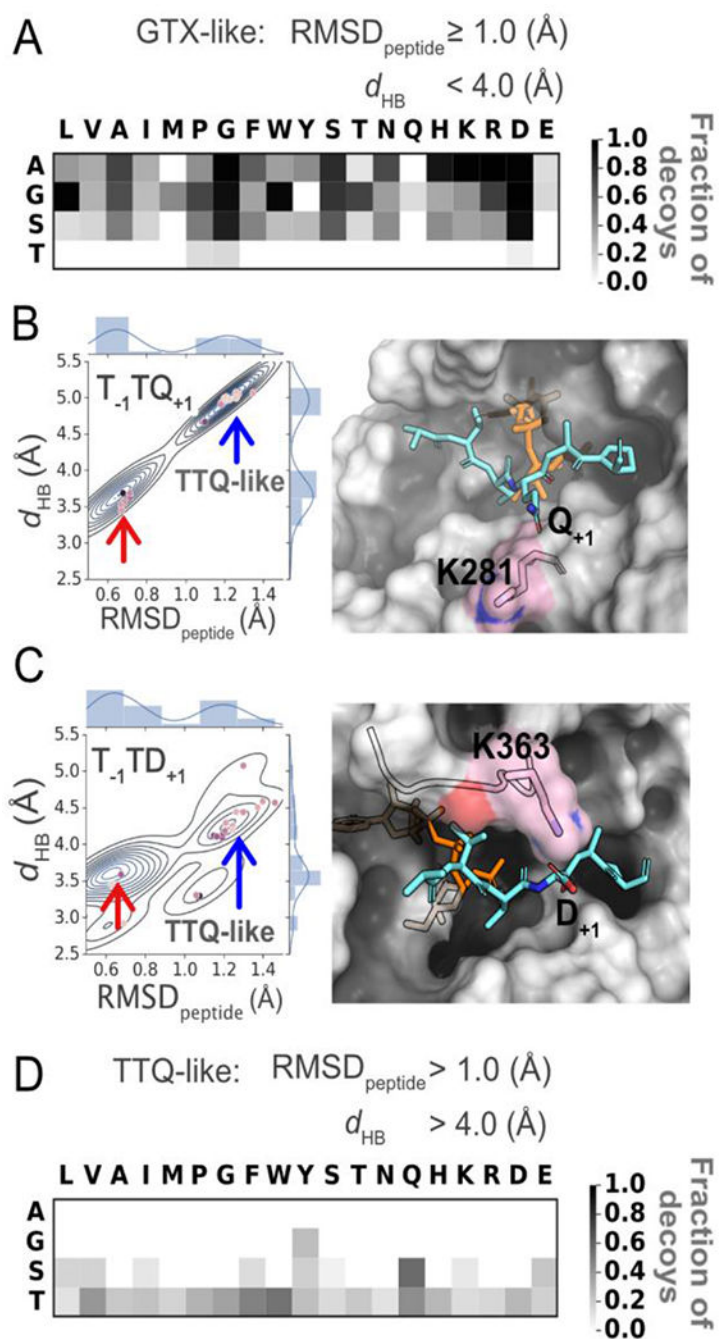
Author Manuscript



**Figure 5. Substrate specificity based on  $\text{RMSD}_{\text{peptide}} < 1.0 \text{ \AA}$  criterion.**

Joint and marginal probability densities for the top 10% of structures by score (200/2000) for a given amino acid at the  $-1$  position and aggregated over all amino acids at the  $+1$  position. (A)  $G_{-1}$  and (B)  $P_{-1}$  and all amino acid residues at  $X_{+1}$ ; Top 1% (20/2000) decoys per sequon shown as points where darker color indicates lower interaction energy. “GTX-like” state is marked with a black arrow and “PTX-like state” is marked with a red arrow. (C) Lowest interaction energy decoy in the enzyme’s peptide binding groove for representative sequon  $P_{-1}TP_{+1}$  (white;  $\text{RMSD}_{\text{peptide}} < 1.0 \text{ \AA}$ ) superposed with that for  $G_{-1}TP_{+1}$  (aquamarine;  $\text{RMSD}_{\text{peptide}} > 1.0 \text{ \AA}$ ). (D) Joint and marginal probability densities of Top 10%(200/2000) sequons for all peptides with fixed amino acids  $A_{-1}$ ,  $S_{-1}$ ,  $T_{-1}$  and all amino acid residues at  $X_{+1}$ ; Top 1% (20/2000) decoys per sequon shown as points where darker color indicates lower interaction energy. The blue arrow indicates a third state distinct from PTX- and GTX-like states. (E) Heatmap of  $\text{RMSD}_{\text{peptide}}$  (left panel) of the lowest energy decoy per sequon, and fraction of decoys ( $N_r$ ) satisfying RMSD criterion (right panel) for  $\text{RMSD}_{\text{peptide}} < 1.0 \text{ \AA}$ . (F) True positives (TP), true negatives (TN), false positives (FP), and false negatives (FN). (G) True Positive Rate (TPR) vs False Positive Rate (FPR) curve showing  $\text{AUC}(\text{RMSD}_{\text{peptide}}) = 0.959$  and  $\text{AUC}(N_r) = 0.965$ .

(FP) and false negatives (FN) predicted based on  $\text{RMSD}_{\text{peptide}} < 1 \text{ \AA}$  threshold applied to the lowest-interaction-energy decoy of the largest cluster. (G) ROC curve for  $\text{RMSD}_{\text{peptide}}$  (magenta) and  $N_r$  (grey) satisfying  $\text{RMSD}_{\text{peptide}} < 1.0 \text{ \AA}$ . Black and brown boxes in (E) and (F) indicate examples of glycosylatable and non-glycosylatable sequons, respectively, that also exhibit the GTX-like states.



**Figure 6. Secondary effects of the amino acid at the +1 position.**

(A) Fraction of decoys sampling the GTX-like state for sequons with A, G, S, or T at the  $-1$  position. (B) Joint and marginal probability densities of top 10% (200/2000) sequons for  $T_{-1}TQ_{+1}$  (left panel) and lowest energy decoy for TTQ-like state for sequon  $T_{-1}TQ_{+1}$  state, where  $Q_{+1}$  position interacts with K281. (C) Joint and marginal probability densities of top 10% (200/2000) sequons for  $T_{-1}TD_{+1}$  (left panel) and lowest energy decoy for TTQ-like state for sequon  $T_{-1}TD_{+1}$  state, where  $D_{+1}$  position interacts with K363. (D) Fraction of low-energy decoys in the TTQ-like state. Top 1% (20/2000) decoys per sequon shown as

points in (B) and (C) where darker color indicates lower interaction energy. “TTX-like” state is marked with a black arrow and “PTX-state” is marked with a red arrow.

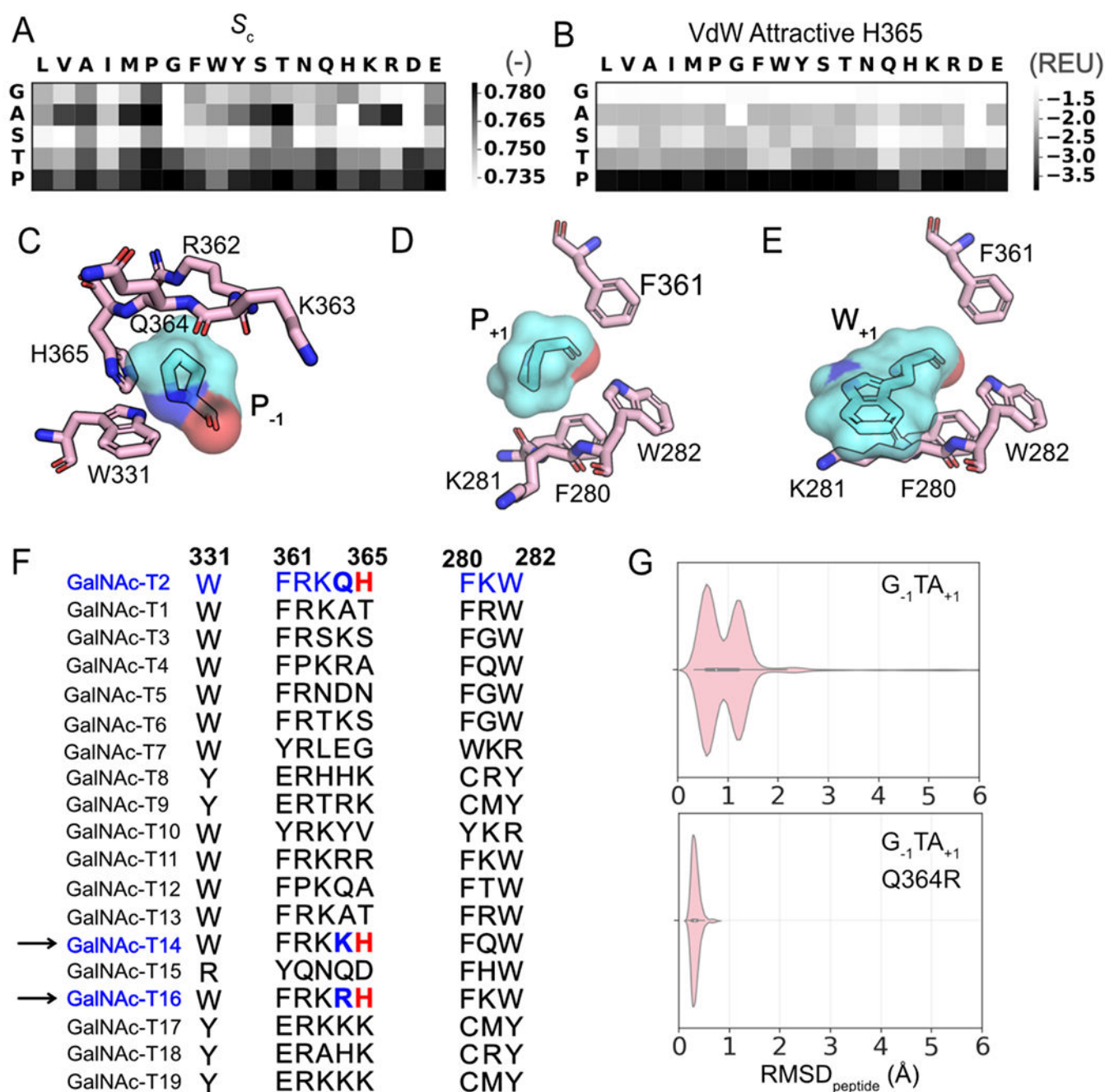
Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript





**Figure 7. Characterization of enzyme–peptide interactions for top 10 decoys  $G_{-1}$ ,  $A_{-1}$ ,  $S_{-1}$ ,  $T_{-1}$ ,  $P_{-1}$  peptides.**

(A) Median shape complementarity ( $S_c$ ). (B) Attractive component of the van der Waals (VdW) potential in Rosetta score function between the residue at the  $-1$  position and H365 on the enzyme. (C)  $-1$  pocket of at the enzyme peptide interface with H365 on the enzyme (pink) interacting with the proline at the  $-1$  position on the peptide (aquamarine). (D) Residues 280, 281, 282 and 361 on the enzyme (pink) interacting with proline at the  $+1$  position on the peptide (aquamarine). (E) Residues 280, 281, 282 and 361 on the enzyme (pink) interacting with tryptophan at the  $+1$  position on the peptide (aquamarine). (F)

Multiple sequence alignment of isoform T2 with other isoforms for the residues at the enzyme-peptide interface for +1 and -1 positions on the peptide. (G) Violinplots for  $\text{RMSD}_{\text{peptide}}$  distributions sampled for sequon  $\text{G}_{-1}\text{TA}_{+1}$  for isoform T2 (top) and a variant T2-Q364R (bottom). Residue numbering based on GalNAc-T2 Uniprot entry Q10471.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

**Table 1.**

Summary of AUC scores, true positives (TP), true negatives (TN), false positives (FP) and false negatives (FN), and balanced accuracy (BA=(TPR+TNR)/2; TPR=TP/(TP+FN) and TNR=TN/(TN+FP)) and false positive rate (FPR=FP/(FP+TN)) for predictions at an experimental glycosylation efficiency threshold of 10%.

Feature	AUC	Feature threshold	TP	TN	BA (%)	FPR (%)
			FP	FN		
Interaction Energy ( $d_{\text{HB}} < 4.0 \text{ \AA}$ )	0.875	$< -34 \text{ (REU)}$	42	255	86.1	19.0
			60	4		
$d_{\text{HB}}$ (largest cluster)	0.923	$< 4.0 \text{ (\AA)}$	44	259	88.9	17.8
			56	2		
Fraction of decoys ( $d_{\text{HB}} < 4.0 \text{ \AA}$ )	0.906	$> 0.70 \text{ (-)}$	39	270	85.2	14.3
			45	7		
RMSD <sub>peptide</sub> (largest cluster)	0.959	$< 1.0 \text{ (\AA)}$	39	287	87.9	8.9
			28	7		
Fraction of decoys (RMSD <sub>peptide</sub> $< 1.0 \text{ \AA}$ )	0.965	$> 0.53 \text{ (-)}$	36	292	85.5	7.3
			23	10		
$S_c$ (median top 10; RMSD <sub>peptide</sub> $< 1.0 \text{ \AA}$ )	0.944	$> 0.735$	42	265	87.7	15.9
			50	4		
Interaction Energy (largest cluster)	0.564	$< -34 \text{ (REU)}$	36	98	54.7	68.9
			217	10		
Interaction Energy (RMSD <sub>peptide</sub> $< 1.0 \text{ \AA}$ )	0.908	$< -34 \text{ (REU)}$	39	271	85.4	14.0
			44	7		

The calculation of TPs, TNs, FPs, FNs, BA and FPR requires a threshold. Feature thresholds were chosen in two cases ( $d_{\text{HB}}$  (largest cluster)  $< 4.0 \text{ \AA}$  and RMSD<sub>peptide</sub> (largest cluster)  $< 1.0 \text{ \AA}$ ) to match criteria discussed in the main text. For all other cases, thresholds were chosen arbitrarily. Also see Table S2 for precision and F1 score.