

1 Decoding glycosylation potential from protein 2 structure across human glycoproteins with a multi- 3 view recurrent neural network

4
5 Benjamin P. Kellman,^{1,2,3,4,5,*} Julien Mariethoz,⁷ Yujie Zhang,¹ Sigal Shaul,^{1,2} Mia Alteri,^{1,2} Daniel
6 Sandoval,⁶ Mia Jeffris,^{1,2} Erick Armingol,^{1,2,3} Bokan Bao,^{1,3} Frederique Lisacek,^{7,8} Daniel
7 Bojar,^{9,10,*} Nathan E. Lewis^{1,2,3,5,*}
8 * Corresponding authors

¹ Department of Pediatrics, University of California, San Diego, La Jolla, CA 92093, USA

² Department of Bioengineering, University of California, San Diego, La Jolla, CA 92093, USA

³ Bioinformatics and Systems Biology Graduate Program, University of California, San Diego, La Jolla, CA 92093, USA

⁴ Augment Biologics, La Jolla, CA 92092

⁵ Ragon Institute of MGH, MIT, and Harvard, Cambridge, MA, USA

⁶ Department of Cellular and Molecular Medicine, University of California, San Diego, La Jolla, CA 92093, USA

⁷ Proteome Informatics Group, Swiss Institute of Bioinformatics, CH-1227 Geneva, Switzerland

⁸ Computer Science Department & Section of Biology, University of Geneva, route de Drize 7, CH-1227, Geneva, Switzerland

⁹ Wallenberg Centre for Molecular and Translational Medicine, University of Gothenburg, Gothenburg 41390, Sweden

¹⁰ Department of Chemistry and Molecular Biology, University of Gothenburg, Gothenburg 41390, Sweden

Correspondence to:

Benjamin P. Kellman,

Daniel Bojar,

Nathan E. Lewis, 9500 Gilman Dr. MC 0760, La Jolla, CA 92093, email: nlewisres@ucsd.edu

9 Abstract

10 Glycosylation is described as a non-templated biosynthesis. Yet, the template-free premise is
11 antithetical to the observation that different N-glycans are consistently placed at specific sites. It
12 has been proposed that glycosite-proximal protein structures could constrain glycosylation and
13 explain the observed microheterogeneity. Using site-specific glycosylation data, we trained a
14 hybrid neural network to parse glycosites (recurrent neural network) and match them to feasible
15 N-glycosylation events (graph neural network). From glycosite-flanking sequences, the
16 algorithm predicts most human N-glycosylation events documented in the GlyConnect database
17 and proposed structures corresponding to observed monosaccharide composition of the
18 glycans at these sites. The algorithm also recapitulated glycosylation in Enhanced Aromatic
19 Sequons, SARS-CoV-2 spike, and IgG3 variants, thus demonstrating the ability of the algorithm
20 to predict both glycan structure and abundance. Thus, protein structure constrains glycosylation,
21 and the neural network enables predictive *in silico* glycosylation of uncharacterized or novel
22 protein sequences and genetic variants.

23
24
25
26

27 Introduction

28 Glycosylation is difficult to study as the one supposedly non-templated biopolymer.¹ Unlike
29 RNA, DNA, and proteins, glycan sequences are understood to be determined by local metabolic
30 and enzymatic conditions, including the availability of charged nucleotide sugars, enzyme
31 availability, Golgi localization, and substrate competition.² These well-supported claims do not
32 explain how different glycosylation sites within one protein are consistently differentially
33 glycosylated; a phenomenon called “microheterogeneity.”³

34

35 Indications of protein structure bounded biosynthesis for glycans has existed for decades. After
36 the N-glycosylation sequon (NX[S/T]) was defined, proximal-amino acid variation was found to
37 impact glycosylation complexity,^{4–6} occupancy,⁷ efficiency,⁸ and glycan class.⁹ Conversely,
38 amino acid sequence alignments of similarly glycosylated glycosites suggest the presence of
39 glycosite-flanking sequence conservation.¹⁰ In influenza and HIV, variation in glycosylation and

40 genetic variation proximal to glycosites can facilitate immune evasion.^{11,12} Examples of how the
41 protein context can constrain glycosylation include observations of higher-order structures such
42 as β -sheets and α -helices,¹³ accessibility,^{14–16} and glycosylation kinetics,^{17–20} all of which impact
43 glycan structure. We quantified associations between glycan substructures and local protein
44 structure, showing that protein structural constraints can predict glycosylation. Together, these
45 protein-glycan relations form a more comprehensive framework we call bounded biosynthesis,
46 wherein glycosylation is bounded by both metabolic conditions and genome-encoded protein
47 structural constraints.²¹ That study describes protein structure as a major determinant of
48 glycosylation, but there is a need to functionalize the proteomic bounds on glycosylation such
49 that it can be leveraged with ease to predict glycosylation from protein structure.

50
51 Machine learning can be applied to the complex structures of glycans for the analysis of glycan
52 structure, function, and classification. For example, natural language processing can encode
53 glycans longitudinally from the reducing end.^{22,23} The SweetTalk glycan embedding
54 recapitulated both antigenic glycans and microbial pathogenicity and phylogeny. Another study
55 leveraged the branched nonlinear glycan structure to scaffold graph convolutional neural
56 networks.²⁴ SweetNet identified glycan targets of viral lectins. Beyond glycan embedding,
57 biosynthetic constraints and outcomes have been modeled using neural networks.²⁵ Previous
58 attempts have been made to relate glycan branching with glycosite-proximal protein structure.²⁶
59 In the absence of meaningful embeddings and biosynthetic-substructure decomposition like
60 SweetNet and GlyCompare,²⁷ previous observations were limited to the association between
61 surface accessibility and glycan complexity. With these new embeddings and the knowledge
62 that glycan biosynthesis is a protein structure guided process, we can now functionalize protein-
63 based glycan predictions.

64
65 Here we present the Interloping Saccharide Neural Network Extrapolation (InSaNNE) model,
66 which predicts N-glycosylation from glycosite-proximal protein features. Using long short-term
67 memory (LSTM) units,²⁸ a type of recurrent neural network, we analyze glycosite-proximal
68 amino acids and leverage the functional and biosynthetic glycan encodings of SweetTalk,
69 SweetNet, and GlyCompare to generate an accurate mapping of glycan structure to protein
70 sequence and structure. We train and validate our glycosite-glycan pairing model on empirically
71 observed site-specific glycosylation. The model is trained using data from UniCarbKB²⁹ and
72 validated using more extensively curated data from GlyConnect³⁰. We further validate our
73 predictions on important glycosylation events on the coronavirus spike protein, immunoglobulin,

74 and the enhanced aromatic sequon. All N-glycan predictions are integrated in GlyConnect for
75 easy access. With InSaNNE, we leverage the new bounded biosynthesis paradigm to open
76 glycobiology to everyone by predicting expected and differential glycosylation onto their proteins
77 of interest.

78 Results

79 Graph convolutional neural networks accurately predict glycan- 80 glycosite pairs

81 We developed a model to predict the presence of specific glycans given the flanking amino acid
82 sequence at N-linked glycosylation sites. Specifically, glycan structures can be ranked to
83 indicate the most feasible glycosylation events at a glycosite of interest. To train, validate, and
84 test the model, we collected and annotated 1,721 unique glycosylation events across 75 human
85 glycoproteins from UniCarbKB²⁹ wherein glycan structure was previously fully determined (see
86 Methods). The model incorporates modules that analyzed both glycan structures (**Figure 1a**)
87 and the protein sequences (**Figure 1b**). To analyze the protein sequences, we used long short-
88 term memory (LSTM) units,²⁸ a recurrent neural network module effective at modeling protein
89 structure by asserting language-like processivity³¹ (**Figure 1b**). Both sequence-proximal
90 (glycosite-flanking) and spatially proximal (within n-Angstroms) protein features are important for
91 predicting feasible glycosylation. We examined two separate LSTM-based modules into our
92 model for analyzing the sequence-proximal and spatially proximal amino acids, separately. For
93 the analysis of the glycan component, we tested three glycan embeddings: (1) a fully connected
94 neural network using GlyCompare glycan substructure features²⁷ as input, (2) a glycan-based
95 language model in the style of SweetTalk,²³ and (3) a graph convolutional neural network based
96 on SweetNet.²⁴

97

98 On average, the model based on GlyCompare glycan substructure features achieved a 76.3%
99 accuracy in predicting which glycans have been observed at specific glycosites (**Table 1**). The
100 recurrent neural network (SweetTalk; 79.9%) or graph convolutional neural network (SweetNet;
101 83.1%) models further improved the performance, demonstrating that optimizing the glycan
102 analysis modules increases prediction performance. Choosing the SweetNet-based model as

103 our best-in-class performer, we used stochastic weight averaging (SWA; Izmailov et al., 2019)
104 to further optimize performance. SWA improved SweetNet-based model accuracy to 87.5%
105 (**Table 1**) and was therefore selected as our final model and used for all downstream analyses.

106 *Table 1 – A model for glycan-glycosite matching was developed to predict permissible glycans on a glycosylation site.*
107 *Modules analyzing the glycosite-flanking protein sequence and additional spatially proximal amino acids consisted of*
108 *recurrent neural networks, while the module analyzing glycans was either a fully connected neural network using*
109 *GlyCompare substructure features as input (GlyCompare), a glycan-based language model (SweetTalk), or a graph*
110 *convolutional neural network (SweetNet). We further tested the effect of stochastic weight averaging (SWA) on model*
111 *performance. Removing the information about spatially proximal amino acids from the model input is denoted by “-*
112 *Spatial” while the addition of the whole protein sequence as an additional input for the model is indicated by*
113 *“+Whole”. Results represent the mean values for accuracy and area under the curve (AUC) for the receiver-operator*
114 *curve (ROC) on our test set after five independent training runs.*

Metric	GlyCompare	SweetTalk	SweetNet	SweetNet SWA	SweetNet SWA -Spatial	SweetNet SWA +Whole
Accuracy	0.763	0.799	0.831	0.875	0.861	0.879
ROC AUC	0.823	0.871	0.894	0.929	0.920	0.930

115

116 After optimizing the glycan analysis module, we analyzed the role of protein sequences on
117 prediction performance. We trained a model that only had access to the glycan structure and
118 the glycosite-flanking sequence, without additional spatially (3D) proximal amino acids.
119 Compared to the full InSaNNE model (87.5% accuracy), the model without spatially proximal
120 amino acids achieved a slightly worse performance (86.1%, **Table 1**). The marginal
121 performance loss suggests that, while spatially proximal information helps, the glycosite-flanking
122 residues are most important.

123 We next trained a model with access to the whole sequence of each protein, in addition to
124 glycosite-proximal amino acids, and glycan structures. The additional information from the
125 whole protein slowed training and inference, while providing a limited performance improvement
126 (87.9% accuracy, **Table 1**). We concluded that distant amino acids carry limited relevant
127 information for predicting permissible glycan structure that is not already captured in the nearby
128 sequence and spatially proximal amino acids.

129 Different glycosites prefer specific glycan features

130 True negatives, infeasible glycans, are hard to obtain experimentally, so we focused on recall
131 (True Positive Rate). InSaNNE achieved a recall of 84.8% for *N*-linked glycosylation events in
132 our dataset. The notable performance in these glycan-type-specific models, suggests that
133 InSaNNE performs with exceptional recall – recovering most permissible glycans at a given
134 glycosite.

135 Next, we examined which *N*-glycan motifs were more difficult for InSaNNE to predict. For this,
136 we calculated the average prediction accuracy for each glycan feature in the validation set.
137 Several rare glycan motifs (<10 observations) were more difficult to predict (**Figure 2a**).
138 However, InSaNNE exhibited a predictive accuracy of >80% for most motifs (**Figure 2b**). Since
139 glycan features represent a hierarchical feature set, rare motifs with low prediction accuracy are
140 not independent from each other and formed clusters based on glycan structure similarity
141 (**Figure 2c**). For example, glycan features with lower predictive performance were enriched for
142 oligomannose. Analogous to the glycan features, most glycosites exhibited an aggregate
143 predictive accuracy >90% (**Figure 2d**) and we found prediction performance correlated with the
144 number of observed glycans for similar glycosites (close in the embedding manifold;
145 **Supplementary Figure 1**). Predictions were robust to the removal of single amino acids or
146 short motifs, suggesting redundancy within glycosite-flanking sequences and soft boundaries on
147 the flanking window size (**Supplementary Figure 2**). Furthermore, the flanking residues, rather
148 than the central sequon-proximal residues, informed model predictions the most; ablation of
149 upstream residues was most impactful on performance (**Supplementary Figure 2**). In general,
150 given the consensus sequence of *N*-linked glycosylation, flanking residues are more variable,
151 and may carry more information for deep learning models, than more conserved sequon-
152 adjacent residues.

153 To illustrate the capabilities of InSaNNE, we used the model to predict the feasibility of all
154 glycans in our dataset at the glycosite GTVLTRNETHATYS (P07911:N396) from human
155 uromodulin – the most abundant protein in human urine and relevant for chronic kidney
156 disease.³² Notably, 58 of 61 experimentally observed glycans were placed in the top 80
157 predicted glycans (**Figure 2e**). Additionally, top glycans that were not previously reported at this
158 glycosite shared features with the observed glycans, such as a strong negative charge via
159 sialylation and/or sulfation. These results further demonstrate protein-sequence-based glycan
160 prediction and emphasize the value and relevance of our model.

161 Single amino acid changes modulate specific glycan features

162 While the ablation of individual glycosite-flanking amino acids does not substantially diminish
163 model performance (**Supplementary Figure 2**), glycosylation efficiency and range can be
164 impacted by glycosite-flanking mutations.^{5,6,9,11} Therefore, we tested if InSaNNE can predict how
165 changes to the glycosite-flanking sequence will impact glycosylation. This could facilitate
166 glycoengineering and elucidate structural interactions between protein and glycan structures at
167 the glycosylation site. We performed a deep mutational scan *in silico* (replacing each of the 14
168 glycosite-flanking amino acids with all amino acids) on every N-glycosite in our dataset. Using
169 the modified glycosite sequences as inputs for InSaNNE, we analyzed the changes in predicted
170 glycans compared to the wild-type sequence. To focus interpretation, we grouped glycans into
171 “sialylated” and “fucosylated.” This allowed us to track the changes in predicted probability for
172 each of these features following specific glycosite-flanking mutation (**Figure 3, Supplementary**
173 **Figure 3**). However, while these reflect general trends of individual glycosites across all
174 proteins, amino acid substitutions may have effects that deviate from these general trends.

175
176
177 For multiple amino acid substitutions, we observed distinct changes in the predicted
178 glycosylation of modified glycosites, with clear differences between changes to upstream and
179 downstream regions. The introduction of some amino acids (e.g., tyrosine; **Figure 3a**) had the
180 same qualitative effect regardless of where they were introduced. Meanwhile, other amino acids
181 (e.g., cysteine; **Figure 3b**) have diverging effects, with a decrease in predicted complex glycans
182 when introduced upstream and an increase when it is present downstream. We also observed
183 that predicted changes in glycosylation were impacted more strongly by mutations in the distal
184 parts of the glycosite-flanking sequence (e.g., glutamate; **Figure 3c**). These general trends of
185 amino acid-glycan associations could be useful for glycosite-specific glycoengineering.

186 Uncharacterized glycoproteins and glycan compositions can be 187 annotated with candidate glycan structures

188 Computational prediction to annotate protein features and functions is done routinely for newly
189 discovered proteins, yet limited *in silico* characterizations exist for glycosylation. However, the
190 relative speed of predicting glycosylation would make it invaluable for new, existing, or poorly
191 characterized proteins; typical glycoprofiling approaches can otherwise take several months.

192 Even many well-characterized glycoproteins have only compositional measurements
193 (unstructured monosaccharide counts) since glycan structure measurement and
194 characterization are resource and expertise-intensive processes. Thus, InSaNNE could be
195 invaluable for annotating glycosylation sites.

196 Predicting glycosite location is one of the few high-confidence bioinformatic predictions involving
197 glycosylation.^{33–37} To extend this capability, we predict the feasible glycan structures of 2,763
198 human N-linked glycosites in the GlyConnect database.³⁰ For this, we used InSaNNE to analyze
199 the annotated glycosylation sites together with the six upstream and seven downstream amino
200 acids. For each glycosite, we predicted the likelihood of 199 N-linked glycans (**Supplementary**
201 **Dataset 1**). Using our independent test set, we ascertain a threshold with an acceptable false-
202 positive rate (AUC 0.92, **Figure 4a**). A threshold of 0.6 (predicted presence) corresponded to a
203 false-positive rate <10% while maintaining a true positive rate >85%. This allowed us to assess
204 the recall or sensitivity of our predictions within GlyConnect by quantifying known glycan
205 structures that were successfully predicted (**Figure 4b**). Thus, InSaNNE could inform future
206 experiments and comparative analyses of structure-based constraints in glycosylation and
207 functional impacts.

208 InSaNNE predicts complex glycans in the enhanced 209 aromatic sequon and the SARS-CoV-2 Spike

210 N-glycans are commonly grouped into categories, such as highly processed complex glycans,
211 hybrid glycans, and immature oligomannose glycans.³⁸ Previous work showed that an aromatic
212 residue located two-positions N-terminal from a glycosylation site results in less complex N-
213 glycosylation at the site, termed the enhanced aromatic sequon.⁶ In this case, an L to F
214 substitution two residues upstream of the CD2 glycosylation site transformed the site from
215 predominantly complex (sialylated) and hybrid structures to low complexity (oligomannose)
216 structures. When InSaNNE evaluates the same sequences, the F allele sequence shows
217 significantly higher predicted presence for higher-mannose structures. We predict an
218 enrichment for 7-mannose structures (One-sided Mann-Whitney-Wilcoxon, $p=0.017$) and predict
219 an overall increase in oligomannose structure for the F allele (Linear model; Wald, $p<0.001$; F-
220 statistic, $p=7.44\times 10^{-5}$; **Figure 5a**). We see a corresponding decrease in sialylated structures in
221 the F allele (One-sided Mann-Whitney-Wilcoxon, $p<1e-4$; **Figure 5b**).

222

223 InSaNNE also recapitulates glycan types of SARS-CoV-2. These sites have been extensively
224 characterized throughout the pandemic.^{15,39-41} N234, N717, and N801 are highly reproducible
225 oligomannose sites.¹⁵ Oligomannose at N234 is consistently high (80-100%)¹⁵ and appears
226 necessary to support the open ACE2-binding spike conformation.⁴² Our predictions show strong
227 preference for Man5 and Man9 structures and a strong anticorrelation with sialylation (**Figure**
228 **5c-d**). Sites N717 and N801¹⁵) are predicted here to have almost no sialylation (**Figure 5c-d**).
229 Predictions for all glycosylation sites were mostly consistent with empirical observations
230 (**Supplementary Figure 4**).

231
232 We wondered if the spike protein of new strains shows predictable changes in glycosylation. We
233 examined InSaNNE predictions at site N616 in a simulated D614G variant (**Supplementary**
234 **Figure 6**) and N717 in a T716I variant (**Figure 5e-f**). We found distinct changes in predicted
235 glycosylation. T716I, between the furin cleavage site and the fusion peptide, is within the more
236 conserved S2 sequence and retains moderate antibody accessibility regardless of RBD
237 conformation.⁴³ To focus on relevant changes, we examined those with non-negligible ancestral
238 predicted-presence (>0.1) and substantial fold change ($|\log_{2}FC|>1$) relative to the ancestral
239 spike. At site N717 in the T716I variant, many asialylated sugars with one to three galactose
240 residues decrease relative to ancestral (**Figure 5f**, blue points). Additionally, a small number of
241 sugars with zero to two sialic acids and one to four galactose residues increase. Though
242 InSaNNE predicts that site N717 becomes variably permissible to mono-, di-, tri- and tetra-
243 antennary sialylated and asialylated structures, empirically, it is an oligomannose site,
244 suggesting these terminal galactoses may not be visible without additional mutations to the site.
245 Distinctly, InSaNNE reveals few confident changes at site N616 in the D614G variant
246 (**Supplementary Figure 6**). If glycan structure can be predicted from primary sequence, site
247 occupancy may also be bound by these constraints.

248
249 **InSaNNE predictions recapitulate biantennary abundance**
250 **on human IgG3**

251 Mutations can perturb glycosylation in IgG3.⁹ Eight complex biantennary structures in human
252 IgG3 were measured for wildtype (*wt*) and glycosite (N297; P01860:N227) proximal mutants.
253 While the *wt* IgG3 showed a preference for core-fucose and a1-6-branch galactose, R301A

254 increased all terminal galactose, and Y296A accepted no galactosylation (**Figure 6a**). Thus,
255 primary protein structure can profoundly influence glycosylation.

256 We compared InSaNNE predictions for the R301A and Y296A mutants and found that
257 predicted-presence and change in predicted-presence were correlated with empirical
258 occupancy. Abundance-prediction correlation was high for the R301A mutant ($R^2=0.876$; **Figure**
259 **6b**) and moderate for *wt* abundance ($R^2=0.25$; **Figure 6b**). Predicted presence was consistent
260 with measured abundance in the Y296A mutant ($R^2=0.33$; **Figure 6b**). Interestingly, prediction
261 performance increased when we compared changes relative to *wt*. The predicted presence log
262 fold-change in R301A relative to *wt* was highly correlated with measured abundance log fold-
263 change ($R^2=0.87$; **Figure 6c**). Yet, the consistency in predicted vs observed change for Y296A
264 decreased dramatically ($R<0$, $R^2=0.27$; **Figure 6c**). To further probe the prediction failure in
265 Y296A, we removed glycans with small predicted changes ($|\log FC|<1$). Without the low-
266 confidence changes, abundance prediction performance for *wt* ($R^2=0.52$), R301A ($R^2=0.99$),
267 and log fold-change (R301A vs. *wt*: $R^2=0.95$) improved (**Figure 6d-e**), while nearly all
268 predictions for Y296A dropped out. These results suggest that InSaNNE can predict occupancy
269 and occupancy change for non-small ($|\log \text{fold-change}|>1$) changes.

270 Accessing InSaNNE predictions and continuous comparison 271 through GlyConnect

272 We evaluated the agreement between InSaNNE predictions and GlyConnect data at the
273 compositional level. **Figure 7a** shows the protein-page d3 heatmap illustration comparing
274 GlyConnect-annotated glycosylation events for human coagulation factor XI (UniProt:P03951;
275 GlyConnect:818) with InSaNNE predictions; GlyConnect:818 is supported by four published
276 references. **Table 2** summarizes the comparison between GlyConnect annotation and InSaNNE
277 predictions for human coagulation factor XI.

278

	<i>reported structures</i>	<i>reported compositions</i>	<i>predicted structures</i>	<i>overlap</i>
Asn-90	7	0	2	2
Asn-126	5	4	4	2
Asn-163	2	1	9	4
Asn-450	4	4	4	2
Asn-491	10	5	6	6

Total glycans	42	27	16
Number of compositions	14	7	7

279 *Table 2 - Summary of knowledge of human coagulation factor XI (P03951) as stored in the*
280 *GlyConnect database at structural (first column) and compositional (second column) resolutions*
281 *along with predicted structures (third column). The overlap between stored and predicted*
282 *(predicted presence ≥ 0.8) structures is shown in the fourth column and the last row features the*
283 *overall number of compositions. Note that overlap refers to matches between predicted*
284 *structures and reported structures or compositions; one reported composition can map to*
285 *multiple predicted structures.*

286

287 The first composition, H5N4 (five hexoses and four hexosamines), matches three structures
288 with similar linkages recorded in GlyConnect. At site P03951:N491, in composition block H5N4,
289 we see InSaNNE correctly predicts the presence of GlyConnect glycan 3471; the dashed-line
290 compositional matches to glycans 2363 and 3233 are expected as all three glycans are
291 members of the same composition block. Additionally, glycan 2363 is highly predicted at N491
292 suggesting a partial linkage resolution for the incompletely determined structure stored in
293 GlyConnect. Likewise, structures matching the H5N4S2 (five hexoses, four hexosamines, and
294 two sialic acids) compositions contain glycan 3353 predicted and observed at all sites. Within
295 composition block H5N4S2, InSaNNE predicts a higher likelihood (>0.9) for glycan 1641 at
296 N163 (biantennary $\alpha 2,3$ -Neu5Ac). Glycan 1641 offers a complete resolution of structural
297 ambiguity for H5N4S2 at N163. Prediction and annotation both involve flexible linkage
298 definitions, particularly for non-core residues. In contrast, the prediction at site N163 is more
299 extensive than reported data. Interestingly, N163 is a rare NXC sequon, which may explain the
300 smaller number of reported structures and provides novel insights into the distinct preferences
301 of this rare sequon.

302

303 For human coagulation factor XI, GlyConnect contains site-specific observations of 42
304 structures and compositions, and 14 additional distinct but structurally related glycans (Table 2).
305 Compositional similarity was displayed using Compozitor (**Figure 7b**). The Compozitor graph
306 shows 14 compositional nodes connected through the addition of a single monosaccharide. Two
307 virtual nodes (green: H6N4S2 and H5N5S2) are needed to fully connect the graph.⁴⁴ All site-
308 specific InSaNNE-predicted structures correspond to previously annotated site-specific
309 compositions in GlyConnect (magenta). InSaNNE fails to predict structures corresponding to
310 three previously reported compositions the H6N5S2, H6N5F1S2, and H6N5F1S23.

311 Interestingly, the glycan property distribution (**Figure 7c**) is similar between reported and
312 predicted compositions, suggesting a lack of systematic bias that would diminish expected
313 performance for specific glycotypes. Other compositions were found in large scale
314 glycoproteomics experiments without any precise structural features and may be less reliable
315 annotations.

316 Discussion

317 Here we present InSaNNE, the Interloping Saccharide Neural Network Extrapolation, for
318 predicting glycans on membrane-bound and secreted proteins. This approach employs a
319 recurrent neural network and a graph convolutional neural network with stochastic weight
320 averaging to predict feasible glycan structures based on the underlying protein sequence.
321 InSaNNE successfully predicts known glycan structures on a wide range of proteins and
322 assesses the impact of single amino acid substitutions on resulting glycan structures. Beyond
323 initial cross-validation and test-set validation, we successfully predicted glycans on uromodulin,
324 SARS-CoV2, IgG3, and across the GlyConnect database. We have added the glycan
325 predictions to the glycome database GlyConnect, making them accessible for further study of
326 this discovery. Importantly, InSaNNE further questions the premise of template-free glycan
327 biosynthesis. Glycosylation through the bounded biosynthesis paradigm, and its accessibility
328 through the InSaNNE framework, will facilitate more accurate and accessible study of diverse
329 glycoproteins and glycoproteomic behaviors.

330
331 InSaNNE enables the draft annotation of glycosylation on novel proteins, glycoprotein
332 composition analyses, glycoinformatics, and whole proteomes. By increasing the predictability
333 of glycans, we have reduced the challenge of measuring glycans. Mass spectrometry is the gold
334 standard in glycan measurement today, but these measurements may produce partially
335 ambiguous structures and topologies. Consequently, the field is rich with datasets and
336 databases of partially or minimally assembled glycoprofiles.^{45–48} Combining measured glycan
337 compositions with site-specific predictions of feasible glycosylation should facilitate automated
338 glycoprofile assembly. These annotations can be completed for novel and existing glycoprofile
339 assemblies; because of the automated nature, structural glycoprofiles can be assembled for
340 single experiments or entire databases with comparable ease. The sequence-only nature of the
341 prediction is especially important, as many proteins lack experimental structural observations;
342 an algorithm that can operate on the primary sequence is considerably more portable than one

343 requiring structural information. A sequence-only prediction can even be used to quickly
344 compare different isoforms or predict glycans on newly discovered protein sequences.

345
346 We demonstrated our ability to glycosylate an entire proteome by predicting decoration
347 throughout GlyConnect. Newly glycosylated proteins can be used to identify lectin-binding,
348 glycan co-ligands, alternative charge, or steric conformations on proteins of interest, and
349 changes in protein dynamics. These predictions can be disseminated to enrich databases
350 detailing glycosylation^{30,49,50} and other post-translational modifications,⁵¹⁻⁵³ protein structure,^{54,55}
351 domains,^{56,57} and interactions.⁵⁸⁻⁶¹ Future work will extend this approach to O-linked glycans, an
352 even more challenging endeavor due to less available data for training and a seeming absence
353 of a clear consensus sequence on the protein side.⁶²

354
355 Predicted glycosylation can be used to inform large genetic and genome-wide studies. Genetic
356 variation can change protein function and resulting phenotype, but here we demonstrate that it
357 can impact glycosylation. InSaNNE can predict such changes and thus provide further
358 hypotheses for elucidating disease mechanisms. For example, adding predicted differential
359 glycosylation to a study of a high-heterogeneity critical immune gene like Human Leukocyte
360 Antigen (HLA) will be invaluable. This is because HLA has a functional binding-groove adjacent
361 glycosite^{63,64} that could contribute to the behavior, accessibility, and peptide presentation. Some
362 HLA molecules have already been observed to carry allotype-specific glycans.⁶⁵ Beyond HLA,
363 understanding differential glycosylation on reference and variant molecules can help distinguish
364 benign from pathogenic mutations: characterized (e.g., ClinVar) or uncharacterized (e.g.,
365 precision medicine). Additionally, certain glycoforms can modulate secretion.^{66,67} Because each
366 glycan may confer a change in behavior, phenotypes of highly diverse glycoproteins such as
367 secretion, protein-ligand interactions, cell-cell interactions, and extracellular protein complexes
368 can be enriched by knowledge of glycosylation. These are only a few of the studies that may
369 benefit from protein-predicted glycosylation potential.

370
371 Bounded biosynthesis provides a more complete picture of immune evasion by evolving
372 pathogens. Glycan-coated viruses have been responsible for many pandemics, while nearly
373 every decade has seen epidemic strains of viruses, such as influenza. Recent work has
374 highlighted the alignment of these fluctuations with changes in glycans decorating these
375 viruses.¹² Without specific glycoforms, it is not possible to determine which of these viruses
376 successfully disguised critical immune epitopes and which viruses created or maintained new

377 lectin-targeted epitopes. With specific glycan prediction, we may predict the most concerning
378 mutations, those that may reinforce a glycan shield,^{11,68–70} stabilize virulence factors,⁴² or
379 occlude immunogenic antigens.⁷¹ Glycoform predictions can provide these missing data along
380 with previously inaccessible insight into the history and future of viral evolution.

381
382 In summary, bounded glycan biosynthesis, as functionalized by InSaNNE and made accessible
383 through GlyConnect, will enable investigators to easily consider glycosylation across many
384 areas of biological study. InSaNNE will thereby sharpen our understanding of the extracellular
385 space and innumerable intercellular phenotypes.

386

387 Acknowledgements

388 This work was supported by NIGMS (R35 GM119850, NEL), the Novo Nordisk Foundation
389 (NNF20SA0066621, NEL), and a Branco Weiss Fellowship – Society in Science awarded to
390 D.B.

391

392 Conflicts

393 This work is associated with a provisional patent filed by the authors, and Augment Biologics,
394 founded by BK and NEL.

395

396 Methods

397 Site-specific glycosylation training set construction

398 Empirical site-specific glycosylation data from humans was obtained from UnicarbKB²⁹ and
399 Glyconnect⁷² with supplemental information from GlyGen.⁷³ The protein structure annotation
400 was done using the Structural Systems Biology (ssbio) package in python.⁷⁴ Protein structure
401 analysis was performed in Python v2.7.15 using ssbio v0.9.9.8 to retrieve and calculate: existing
402 empirical and homology models from PDB and SWISSMOD (PDBe SIFTS),⁷⁵ *de novo*
403 homology models (I-TASSER v5.1), sequence properties (EMBOS v6.6.0.0 pepstats), sequence
404 alignment (EMBOS v6.6.0.0 needle), secondary structure (DSSP v3.0.0, SCRATCHv1.1::sspro
405 and SCRATCHv1.1::sspro8), solvent accessibility (DSSPv3.0.0 and FreeSASAv2.0.2), and
406 residue depth (MSMSv2.2.6.1). Additional amino acid aggregate features were calculated using

407 R::seqinr. Glycan structures were annotated using a combination of glypy⁷⁶ and GlyCompare²⁷
408 for structure parsing and comparison, respectively. All glycan substructures, a connected subset
409 of monosaccharides with and without linkage information, were extracted from each glycan,
410 merged to make a superset of substructures, then mapped to each glycan. This resulted in a
411 mapping from every glycan in the input database to shared substructures.

412
413 For the dataset used to train InSaNNE, we extracted 1,721 unique glycosylation events from
414 UniCarbKB.²⁹ This included the glycan structure that was observed and the glycosite-flanking
415 sequence (14 amino acids, with the glycosylated amino acid in the center) and structural
416 information in the form of additional amino acids within 6Å if structural simulations converged.
417 As negative examples, we generated the same number of combinations of glycosites and
418 glycans that have not been observed.

419 Model construction

420 All glycan-glycosite matching models comprised (1) a recurrent neural network that analyzed
421 the amino acid sequence of the glycosite, (2) another recurrent neural network analyzing the
422 amino acids of the three-dimensional glycosite surroundings, (3) a model analyzing the glycan
423 structure, described below, and (4) a part consisting of fully connected layers to use the
424 concatenated features generated by the previous modules to predict whether a glycan is
425 permissible at a glycosite. The recurrent neural networks consisted of a 128-dimensional
426 embedding layer followed by two bidirectional long short-term memory (LSTM) layers. The fully
427 connected model part consisted of a linear layer, a leaky ReLU (rectified linear unit) activation
428 function, a batch normalization layer, and a multi-sample dropout scheme⁷⁷ followed by a
429 sigmoid function.

430 We compared three different model architectures for the glycan analysis module. For assessing
431 GlyCompare,²⁷ the glycan analysis module comprised a fully connected neural network using
432 the 12,259 GlyCompare features as inputs for two linear layers interspersed with dropout, leaky
433 ReLU, and batch normalization layers. For the model containing a SweetTalk-based language
434 model for glycan analysis,²² we converted glycans to glycowords and used a bidirectional
435 recurrent neural network for protein sequences. For the SweetNet-based model,²⁴ we converted
436 glycans to graphs by constructing a list of nodes (representing monosaccharides or linkages)
437 and edges to denote graph connectivity. All glycan processing for SweetTalk and SweetNet was

438 done using glycowork version 0.5.⁷⁸ The corresponding model contained an embedding layer
439 and three graph convolutional layers, interspersed by leaky ReLUs, Top-K pooling layers, and
440 both global mean and global maximum pooling operations. Model architectures and
441 hyperparameters were optimized using cross-validation.

442 Model training and prediction

443 All models were trained with an NVIDIA[®] Tesla[®] K80 GPU using PyTorch version 1.11.0.⁷⁹ We
444 split the data on a protein level into 80% for training and 20% for testing. For the RNNs, all
445 glycosite-flanking protein sequence and glycan structure were brought to the same length by
446 padding. Linear layers and RNNs were initialized using Xavier initialization⁸⁰ while SweetNet-
447 type models were initialized using a sparse initialization scheme with a sparsity of 10%.

448 We used a batch size of 64 for all models. As an optimizer, we used ADAM (adaptive moment
449 estimation) with a weight decay value of 0.00001 and a starting learning rate of 0.00001, which
450 was decayed according to a cosine function over 170 epochs. We trained models for a
451 maximum of 250 epochs, with an early stopping criterion of 25 epochs without a decrease in
452 validation loss. As a loss function, we used binary cross-entropy. Beginning from epoch 150, we
453 additionally employed stochastic weight averaging⁸¹ with a learning rate of 0.0001.

454 The presence or absence of each glycan can be predicted from the trained InSaNNE model by
455 inputting a glycosite and glycans to predict whether these glycans could occur on this glycosite.
456 To heuristically boost signal for glycans with limited representation in the training set, we
457 generated a naturalistic background of predicted presence for each glycan. Predictions were
458 generated from all training-set glycosites to capture the biases and variation of the dataset as a
459 background predicted-presence distribution for each glycan. The background-adjusted
460 predicted-presence is the product of predicted presence and the predicted-presence cumulative
461 probability (statsmodels::ECDF v0.12.2) relative to the naturalistic background for that glycan.

462

463 Integration and display of predictions in GlyConnect

464 Using InSaNNe, we calculate the predicted presence of 512 N-linked glycans for each N-linked
465 glycosite in the GlyConnect dataset. Prediction data were processed to fit the requirements of

466 the GlyConnect database format, mainly storing association between glycans, glycoproteins and
467 glycosites.³⁰

468
469 IUPAC-represented glycans,^{82,83} output by InSaNNe, were transformed to GlycoCT⁸⁴ using the
470 GlyConnect API function, `convertIupacToGlycoct` ([https://bitbucket.org/sib-pig/sugar-
471 converter/downloads/](https://bitbucket.org/sib-pig/sugar-converter/downloads/)). Transformed prediction data was integrated in the database to enable
472 dynamic mapping through predefined queries for glycan structures and glycoprotein sites. Once
473 transformed, any update of the InSaNNe prediction will easily be reflected in the database.

474
475 JSON files resulting from querying GlyConnect REST API are used for data export and display.
476 A d3.js heatmap (<https://d3-graph-gallery.com/heatmap>) was selected as an appropriate data
477 visualizer. The dimensions are defined as glycan structures/compositions and glycoprotein sites
478 (designated by UniProt accession numbers and glycosylated amino acid sequence position).
479 Heatmaps are created in three types of pages: (1) protein page featuring all glycan structures
480 and compositions found attached to that protein, (2) structure page, featuring one structure and
481 the many proteins on which they are found attached, and (3) composition page, featuring all
482 matching glycan structures and the many proteins on which they are found attached. This data
483 can be exported as csv files. Prediction data can also be visualized and compared using
484 GlyConnect Compozitor.⁴⁴

485
486
487

489 Figure Captions

490 **Figure 1 InSaNNE model architecture.** **A)** Three model architectures were used to embed glycan structures in
491 meaningful manifolds,^{23,24,27} given a glycan, these models output glycan-specific coordinates within the embeddings.
492 To analyze the GlyCompare features of glycans, we used a fully connected neural network, while a SweetTalk-based
493 language model used linear glycan sequences and a SweetNet-based graph convolutional neural network relied on
494 glycan connectivity (see Methods for details). **B)** Full model architecture of InSaNNE. The results of one of the glycan
495 embedding modules (A) is concatenated with protein-structure and protein-sequence embeddings output by the two
496 protein-language models. These outputs were analyzed by a fully connected neural network and yielded the
497 predicted probability of a glycan-glycosite match. Specifically, InSaNNE takes in a 14 amino acid glycosite-flanking
498 sequence, optional spatially proximal amino acids, and a comprehensive library of 700 representative glycans on
499 which InSaNNE was trained. Glycan libraries containing non-represented glycans can be used following additional
500 training.

501
502 **Figure 2 – Characterizing the glycan-glycosite-matching model InSaNNE.** **A)** Dependence of glycan feature
503 prediction performance on occurrence. Using our trained InSaNNE model, we plotted the averaged prediction
504 performance of glycan features against their counts in our dataset. **B)** Glycan feature accuracy distribution. A
505 histogram of the prediction performance for all observed glycan features is shown. **C)** Clusters of difficult-to-predict
506 glycan features. We used t-SNE to visualize the glycan representation learned by InSaNNE for all glycan features.
507 Each feature was colored by its averaged prediction performance to identify structurally related clusters of glycan
508 features that are more difficult to predict for InSaNNE (shown in brighter colors). **D)** Prediction performance
509 depending on the glycosite was visualized using a t-SNE of the glycosite representations learned by InSaNNE. For all
510 glycosites in our dataset, we averaged prediction performance over all glycans and colored glycosites by prediction
511 performance to identify difficult to predict glycosite clusters. **E)** Experimentally observed and predicted glycans at a
512 glycosylation site of human uromodulin were compared. GTVLTR \underline{N} ETHATYS (P07911:N396) was used to predict
513 permissible glycans using the trained InSaNNE model, and the top 80 predicted glycans were analyzed and
514 compared to previously observed glycans at that site³².

515
516 **Figure 3 - Effects of amino acid substitutions on predicted glycosylation ranges.** A-C) For all N-linked
517 glycosites in our dataset, we substituted each amino acid with tyrosine (A), cysteine (B), or glutamate (C) and input
518 the modified glycosite-flanking sequences into our InSaNNE model and predicted feasible glycosylation. We then
519 calculated the average change (predicted presence difference) compared to the predicted wild-type glycosylation
520 glycosites; shown here with a 95% confidence interval. Lines for changes to fucosylated (red) and sialylated (purple)
521 glycans are shown. See **Supplementary Figure 3** for analogous plots for other amino acid substitutions.

522
523 **Figure 4 - Enriching GlyConnect with InSaNNE predictions.** **A)** For classification thresholds between 0 and 1, we
524 assessed true and false positive rates of InSaNNE predictions on the independent test set and compared it to a
525 random classifier baseline. **B)** We validated InSaNNE predictions with existing structures on GlyConnect by
526 investigating the influence of classification threshold on the hit rate (i.e., recall/sensitivity) of InSaNNE accurately
527 predicting known glycan structures in GlyConnect. The grey dotted line marks the 0.6 threshold used.

528
529 **Figure 5 InSaNNE predicts complex glycans around the enhanced aromatic sequon and the SARS-CoV-2**
530 **spike protein.** (A-B) Boxplot distributions of predicted-presence for the L and F variants at N-2 stratified by number
531 of (A) mannoses per glycan and (B) sialic acids per glycan. (C-D) Boxplots describing predicted glycosylation by (C)
532 mannose per glycan and (D) sialic acid per glycan for three oligomannose sites in the SARS-CoV-2 spike
533 glycoprotein. See **Supplementary Figure 4** for all SARS-CoV-2 spike glycosylation sites. (E-F) Fold changes of
534 predicted glycans at site N717, labeled by number of (E) galactose and (F) sialic acid units, between the wild-type
535 and B.1.1.7 spike protein. Predicted-presence fold-change (y-axis) is stratified by the basal predicted-presence for
536 each glycan in the wild-type (x-axis). Predicted-presence fold-change from wild-type by galactose, mannose, GlcNAc,
537 and sialic acid is provided for N717 and N616 in B.1.1.7 (**Supplementary Figure 5**) and D615G (**Supplementary**
538 **Figure 6**) variants respectively. ns: $p>0.05$, *: $p<0.05$, **: $p<0.01$, ***: $p<0.001$, ****: $p<1e-3$, *****: $p<1e-4$

539

540 **Figure 6 - InSaNNE predictions of relative abundance on IgG3.** **A)** Heatmap showing the log-scale abundance of
541 various glycan species observed in wt and mutant Fc on human IgG3.⁹ **B)** The background-adjusted InSaNNE
542 predicted-presence is compared with the empirical abundance in wild type (black), R301A mutant (blue), and the
543 Y296A mutant (teal). **C)** Log fold change between glycan abundance for mutants relative to wildtype were compared
544 between empirical and predicted abundance for all glycans. **D-E)** The bottom panels mirror panels **B-C** except
545 glycans with a predicted absolute log fold-change less than 1 were removed.

546

547 **Figure 7 Predicted glycosylation pattern of human coagulation factor XI (P03951).** H: hexose, N: hexosamine, F:
548 fucose, S: sialic acid. **(A)** The heatmap displays the predicted presence for glycan structures at each known N-
549 glycosite and indicates agreement with glycans previously observed at those sites retrieved from GlyConnect. The
550 structures in each row are ordered by glycan composition; columns represent the five annotated N-glycosites of
551 P03951. Site-specific glycan structure predictions are many-to-many relationships in the GlyConnect database since
552 the same structure may be associated with several sites and conversely a single site may be predicted to present
553 several similar yet non-mutually exclusive glycan structures. Composition blocks contain all structures matching a
554 specific composition. Color indicates the strength of the predicted presence from 0.8 (lower-bound cutoff) to 1
555 (predicted presence upper-bound). A solid-line borders indicate exact structural matches (identical precise
556 monosaccharides and identical linkages) while dashed lines indicate composition matches (monosaccharide
557 category, e.g., hexose) with at least one non-identical linkage; composition-equivalent blocks (e.g., H5N4) are
558 labelled. **(B)** A Compozitor graph representing compositional similarity between predicted and observed glycans.
559 Fourteen glycan compositions are reported in GlyConnect for human coagulation factor XI. Nodes are connected via
560 single monosaccharide additions represented as the edge label. Seven compositions are predicted and all included in
561 the fourteen previously observed compositions (magenta). Two virtual nodes (green) were added to connect the
562 graph. Numbers within the blue nodes express a correspondence in GlyConnect data between a composition and
563 structures. When the number is absent it means we only have compositional data. The size of the non-blue nodes
564 represents a comparison with the total content of GlyConnect to indicate the likelihood of the composition. For large
565 nodes, the composition occurs often, irrespective of the protein where it is seen. **(C)** The bar chart represents glycan
566 properties mapped in all subsets (database, predicted and virtual). It highlights the similarity across properties of
567 predicted and stored structures.

568

569 References

570

571

572 1. Pothukuchi, P., Agliarulo, I., Russo, D., Rizzo, R., Russo, F., and Parashuraman, S. (2019).
573 Translation of genome to glycome: role of the Golgi apparatus. *FEBS Lett.* 593, 2390–
574 2411.

575 2. Kellman, B.P., and Lewis, N.E. (2021). Big-Data Glycomics: Tools to Connect Glycan
576 Biosynthesis to Extracellular Communication. *Trends Biochem. Sci.* 46, 284–300.

577 3. Johnson, R.L., and Deutsch, H.F. (1970). Preparation and studies of myeloma Fab
578 subfractions. *Immunochemistry* 7, 207–215.

579 4. Petrescu, A.-J., Milac, A.-L., Petrescu, S.M., Dwek, R.A., and Wormald, M.R. (2004).
580 Statistical analysis of the protein environment of N-glycosylation sites: implications for
581 occupancy, structure, and folding. *Glycobiology* 14, 103–114.

- 582 5. Huang, Y.-W., Yang, H.-I., Wu, Y.-T., Hsu, T.-L., Lin, T.-W., Kelly, J.W., and Wong, C.-H.
583 (2017). Residues Comprising the Enhanced Aromatic Sequon Influence Protein N-
584 Glycosylation Efficiency. *J. Am. Chem. Soc.* *139*, 12947–12955.
- 585 6. Murray, A.N., Chen, W., Antonopoulos, A., Hanson, S.R., Wiseman, R.L., Dell, A., Haslam,
586 S.M., Powers, D.L., Powers, E.T., and Kelly, J.W. (2015). Enhanced Aromatic Sequons
587 Increase Oligosaccharyltransferase Glycosylation Efficiency and Glycan Homogeneity.
588 *Chem. Biol.* *22*, 1052–1062.
- 589 7. Shakin-Eshleman, S.H., Spitalnik, S.L., and Kasturi, L. (1996). The Amino Acid at the X
590 Position of an Asn-X-Ser Sequon Is an Important Determinant of N-Linked Core-
591 glycosylation Efficiency (*). *J. Biol. Chem.* *271*, 6363–6366.
- 592 8. Kasturi, L., Eshleman, J.R., Wunner, W.H., and Shakin-Eshleman, S.H. (1995). The
593 Hydroxy Amino Acid in an Asn-X-Ser/Thr Sequon Can Influence N-Linked Core
594 Glycosylation Efficiency and the Level of Expression of a Cell Surface Glycoprotein*. *J.*
595 *Biol. Chem.* *270*, 14756–14761.
- 596 9. Lund, J., Takahashi, N., Pound, J.D., Goodall, M., and Jefferis, R. (1996). Multiple
597 interactions of IgG with its core oligosaccharide can modulate recognition by complement
598 and human Fc gamma receptor I and influence the synthesis of its oligosaccharide chains.
599 *J. Immunol.* *157*, 4963–4969.
- 600 10. Gastaldello, A., Alocci, D., Baeriswyl, J.-L., Mariethoz, J., and Lisacek, F. (2016).
601 GlycoSiteAlign: Glycosite alignment based on glycan structure. *J. Proteome Res.* *15*, 3916–
602 3928.
- 603 11. Yu, W.-H., Zhao, P., Draghi, M., Arevalo, C., Karsten, C.B., Suscovich, T.J., Gunn, B.,
604 Streeck, H., Brass, A.L., Tiemeyer, M., et al. (2018). Exploiting glycan topography for
605 computational design of Env glycoprotein antigenicity. *PLoS Comput. Biol.* *14*, e1006093.
- 606 12. Altman, M.O., Angel, M., Košík, I., Trovão, N.S., Zost, S.J., Gibbs, J.S., Casalino, L.,
607 Amaro, R.E., Hensley, S.E., Nelson, M.I., et al. (2019). Human Influenza A Virus
608 Hemagglutinin Glycan Evolution Follows a Temporal Pattern to a Glycan Limit. *mBio* *10*.
609 10.1128/mbio.00204-19.
- 610 13. Silverman, J.M., and Imperiali, B. (2016). Bacterial N-Glycosylation Efficiency Is Dependent
611 on the Structural Context of Target Sequons. *J. Biol. Chem.* *291*, 22001–22010.
- 612 14. Thaysen-Andersen, M., and Packer, N.H. (2012). Site-specific glycoproteomics confirms
613 that protein structure dictates formation of N-glycan type, core fucosylation and branching.
614 *Glycobiology* *22*, 1440–1452. 10.1093/glycob/cws110.
- 615 15. Allen, J.D., Chawla, H., Samsudin, F., Zuzic, L., Shivgan, A.T., Watanabe, Y., He, W.-T.,
616 Callaghan, S., Song, G., Yong, P., et al. (2021). Site-specific steric control of SARS-CoV-2
617 spike glycosylation. Cold Spring Harbor Laboratory, 2021.03.08.433764.
618 10.1101/2021.03.08.433764.
- 619 16. García-García, A., Serna, S., Yang, Z., Delso, I., Taleb, V., Hicks, T., Artschwager, R.,
620 Vakhrushev, S.Y., Clausen, H., Angulo, J., et al. (2021). FUT8-directed core fucosylation of
621 N-glycans is regulated by the glycan structure and protein environment. *ACS Catal.* *11*,
622 9052–9065.

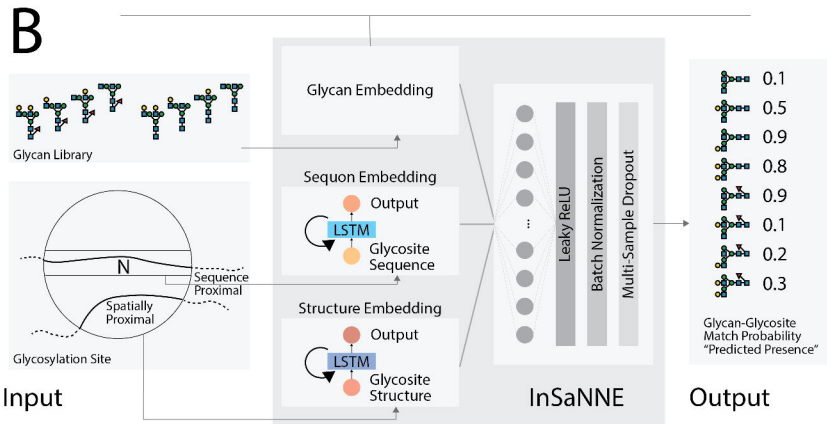
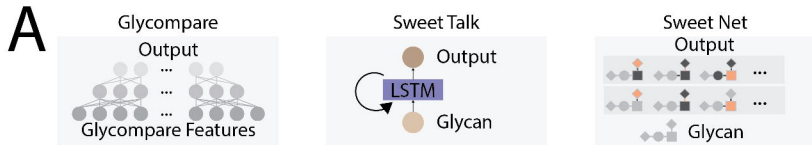
- 623 17. Losfeld, M.-E., Scibona, E., Lin, C.-W., Villiger, T.K., Gauss, R., Morbidelli, M., and Aebi, M.
624 (2017). Influence of protein/glycan interaction on site-specific glycan heterogeneity. *FASEB*
625 *J.* *31*, 4623–4635.
- 626 18. Losfeld, M.-E., Scibona, E., Lin, C.-W., and Aebi, M. (2022). Glycosylation network
627 mapping and site-specific glycan maturation in vivo. *iScience*, 105417.
- 628 19. Mathew, C., Weiß, R.G., Giese, C., Lin, C.-W., Losfeld, M.-E., Glockshuber, R., Riniker, S.,
629 and Aebi, M. (2021). Glycan-protein interactions determine kinetics of N-glycan remodeling.
630 *RSC Chem Biol* *2*, 917–931.
- 631 20. Adams, T.M., Zhao, P., Chapla, D., Moremen, K.W., and Wells, L. (2022). Sequential in
632 vitro enzymatic N-glycoprotein modification reveals site-specific rates of glycoenzyme
633 processing. *J. Biol. Chem.*, 102474.
- 634 21. Kellman Protein structure, a genetic encoding for glycosylation. Unpublished co-
635 submission.
- 636 22. Bojar, D., Camacho, D.M., and Collins, J.J. (2020). Using Natural Language Processing to
637 Learn the Grammar of Glycans. *bioRxiv*, 2020.01.10.902114. 10.1101/2020.01.10.902114.
- 638 23. Bojar, D., Powers, R.K., Camacho, D.M., and Collins, J.J. (2021). Deep-learning resources
639 for studying glycan-mediated host-microbe interactions. *Cell Host Microbe* *29*, 132-144.e3.
- 640 24. Burkholz, R., Quackenbush, J., and Bojar, D. (2021). Using graph convolutional neural
641 networks to learn a representation for glycans. *Cell Rep.* *35*, 109251.
- 642 25. Kotidis, P., and Kontoravdi, C. (2020). Harnessing the potential of artificial neural networks
643 for predicting protein glycosylation. *Metabolic Engineering Communications*, e00131.
- 644 26. Senger, R.S., and Karim, M.N. (2008). Prediction of N-linked glycan branching patterns
645 using artificial neural networks. *Math. Biosci.* *211*, 89–104.
- 646 27. Bao, B., Kellman, B.P., Chiang, A.W.T., Zhang, Y., Sorrentino, J.T., York, A.K.,
647 Mohammad, M.A., Haymond, M.W., Bode, L., and Lewis, N.E. (2021). Correcting for
648 sparsity and interdependence in glycomics by accounting for glycan biosynthesis. *Nat.*
649 *Commun.* *12*, 4988.
- 650 28. Sherstinsky, A. (2020). Fundamentals of recurrent neural network (RNN) and long short-
651 term memory (LSTM) network. *Physica D* *404*, 132306.
- 652 29. Campbell, M.P., Peterson, R., Mariethoz, J., Gasteiger, E., Akune, Y., Aoki-Kinoshita, K.F.,
653 Lisacek, F., and Packer, N.H. (2014). UniCarbKB: building a knowledge platform for
654 glycoproteomics. *Nucleic Acids Res.* *42*, D215-21.
- 655 30. Alocci, D., Mariethoz, J., Gastaldello, A., Gasteiger, E., Karlsson, N.G., Kolarich, D.,
656 Packer, N.H., and Lisacek, F. (2019). GlyConnect: Glycoproteomics Goes Visual,
657 Interactive, and Analytical. *J. Proteome Res.* *18*, 664–677.
- 658 31. Bepler, T., and Berger, B. (2021). Learning the protein language: Evolution, structure, and
659 function. *Cell Syst.* *12*, 654-669.e3.

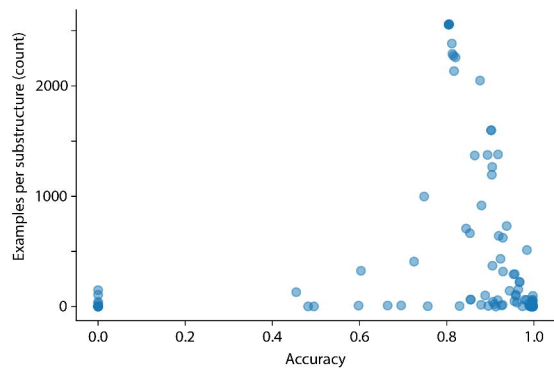
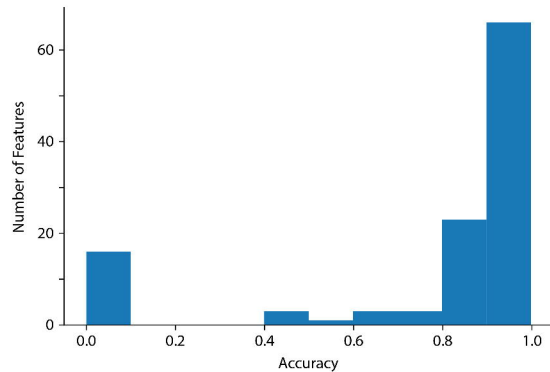
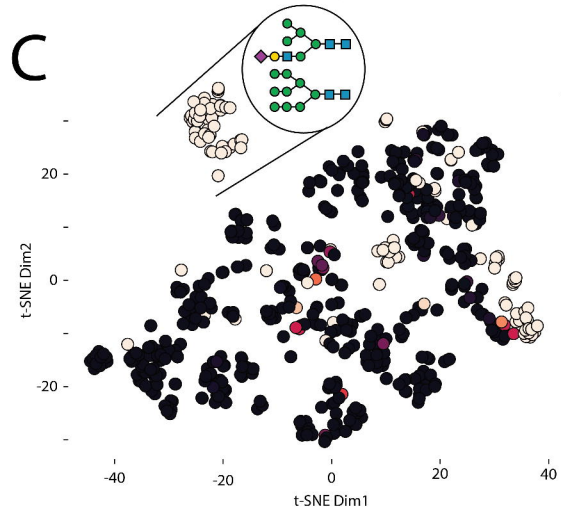
- 660 32. Devuyst, O., Olinger, E., and Rampoldi, L. (2017). Uromodulin: from physiology to rare and
661 complex kidney disorders. *Nat. Rev. Nephrol.* 13, 525–544.
- 662 33. Gupta, R., and Brunak, S. (2002). Prediction of glycosylation across the human proteome
663 and the correlation to protein function. *Pac. Symp. Biocomput.*, 310–322.
- 664 34. Steentoft, C., Vakhrushev, S.Y., Joshi, H.J., Kong, Y., Vester-Christensen, M.B.,
665 Schjoldager, K.T.-B.G., Lavrsen, K., Dabelsteen, S., Pedersen, N.B., Marcos-Silva, L., et al.
666 (2013). Precision mapping of the human O-GalNAc glycoproteome through SimpleCell
667 technology. *EMBO J.* 32, 1478–1488.
- 668 35. Pitti, T., Chen, C.-T., Lin, H.-N., Choong, W.-K., Hsu, W.-L., and Sung, T.-Y. (2019). N-
669 GlyDE: a two-stage N-linked glycosylation site prediction incorporating gapped dipeptides
670 and pattern-based encoding. *Sci. Rep.* 9, 15975.
- 671 36. Taherzadeh, G., Dehzangi, A., Golchin, M., Zhou, Y., and Campbell, M.P. (2019). SPRINT-
672 Gly: predicting N- and O-linked glycosylation sites of human and mouse proteins by using
673 sequence and predicted structural properties. *Bioinformatics* 35, 4140–4146.
- 674 37. Pakhrin, S.C., Aoki-Kinoshita, K.F., Caragea, D., and Kc, D.B. (2021). DeepNGlyPred: A
675 deep neural network-based approach for human N-linked glycosylation site prediction.
676 *Molecules* 26, 7314.
- 677 38. Varki, A., Cummings, R.D., Esko, J.D., Stanley, P., Hart, G.W., Aebi, M., Mohnen, D.,
678 Kinoshita, T., and Packer, N.H. eds. (2022). *Essentials of glycobiology*, fourth edition 4th
679 ed. (Cold Spring Harbor Laboratory Press).
- 680 39. Watanabe, Y., Allen, J.D., Wrapp, D., McLellan, J.S., and Crispin, M. (2020). Site-specific
681 glycan analysis of the SARS-CoV-2 spike. *Science* 369, 330–333.
- 682 40. Zhao, P., Praissman, J.L., Grant, O.C., Cai, Y., Xiao, T., Rosenbalm, K.E., Aoki, K.,
683 Kellman, B.P., Bridger, R., Barouch, D.H., et al. (2020). Virus-Receptor Interactions of
684 Glycosylated SARS-CoV-2 Spike and Human ACE2 Receptor. *Cell Host Microbe* 28, 586-
685 601.e6.
- 686 41. Shajahan, A., Supekar, N.T., Gleinich, A.S., and Azadi, P. (2020). Deducing the N- and O-
687 glycosylation profile of the spike protein of novel coronavirus SARS-CoV-2. *Glycobiology*
688 30, 981–988. 10.1093/glycob/cwaa042.
- 689 42. Casalino, L., Gaieb, Z., Goldsmith, J.A., Hjorth, C.K., Dommer, A.C., Harbison, A.M.,
690 Fogarty, C.A., Barros, E.P., Taylor, B.C., McLellan, J.S., et al. (2020). Beyond Shielding:
691 The Roles of Glycans in the SARS-CoV-2 Spike Protein. *ACS Central Science* 6, 1722–
692 1734. 10.1021/acscentsci.0c01056.
- 693 43. Harvey, W.T., Carabelli, A.M., Jackson, B., Gupta, R.K., Thomson, E.C., Harrison, E.M.,
694 Ludden, C., Reeve, R., Rambaut, A., COVID-19 Genomics UK (COG-UK) Consortium, et
695 al. (2021). SARS-CoV-2 variants, spike mutations and immune escape. *Nat. Rev. Microbiol.*
696 19, 409–424.
- 697 44. Robin, T., Mariethoz, J., and Lisacek, F. (2020). Examining and fine-tuning the selection of
698 glycan compositions with GlyConnect Compozitor. *Mol. Cell. Proteomics.*
699 10.1074/mcp.RA120.002041.

- 700 45. Liu, M.-Q., Zeng, W.-F., Fang, P., Cao, W.-Q., Liu, C., Yan, G.-Q., Zhang, Y., Peng, C.,
701 Wu, J.-Q., Zhang, X.-J., et al. (2017). pGlyco 2.0 enables precision N-glycoproteomics with
702 comprehensive quality control and one-step mass spectrometry for intact glycopeptide
703 identification. *Nat. Commun.* *8*, 438.
- 704 46. Rojas-Macias, M.A., Mariethoz, J., Andersson, P., Jin, C., Venkatakrisnan, V., Aoki, N.P.,
705 Shinmachi, D., Ashwood, C., Madunic, K., Zhang, T., et al. (2019). Towards a standardized
706 bioinformatics infrastructure for N- and O-glycomics. *Nat. Commun.* *10*, 3275.
- 707 47. Riley, N.M., Hebert, A.S., Westphall, M.S., and Coon, J.J. (2019). Capturing site-specific
708 heterogeneity with large-scale N-glycoproteome analysis. *Nat. Commun.* *10*, 1311.
- 709 48. Yang, Y., Yan, G., Kong, S., Wu, M., Yang, P., Cao, W., and Qiao, L. (2021). GproDIA
710 enables data-independent acquisition glycoproteomics with comprehensive statistical
711 control. *Nat. Commun.* *12*, 6073.
- 712 49. Kahsay, R., Vora, J., Navelkar, R., Mousavi, R., Fochtman, B.C., Holmes, X., Pattabiraman,
713 N., Ranzinger, R., Mahadik, R., Williamson, T., et al. (2020). GlyGen data model and
714 processing workflow. *Bioinformatics* *36*, 3941–3943.
- 715 50. Yamada, I., Shiota, M., Shinmachi, D., Ono, T., Tsuchiya, S., Hosoda, M., Fujita, A., Aoki,
716 N.P., Watanabe, Y., Fujita, N., et al. (2020). The GlyCosmos Portal: a unified and
717 comprehensive web resource for the glycosciences. *Nat. Methods.* [10.1038/s41592-020-0879-8](https://doi.org/10.1038/s41592-020-0879-8).
- 719 51. Minguéz, P., Letunic, I., Parca, L., Garcia-Alonso, L., Dopazo, J., Huerta-Cepas, J., and
720 Bork, P. (2015). PTMcode v2: a resource for functional associations of post-translational
721 modifications within and between proteins. *Nucleic Acids Res.* *43*, D494-502.
- 722 52. Craveur, P., Rebehmed, J., and de Brevern, A.G. (2014). PTM-SD: a database of
723 structurally resolved and annotated posttranslational modifications in proteins. *Database*
724 (Oxford) *2014*, bau041–bau041.
- 725 53. Li, Z., Li, S., Luo, M., Jhong, J.-H., Li, W., Yao, L., Pang, Y., Wang, Z., Wang, R., Ma, R., et
726 al. (2022). dbPTM in 2022: an updated database for exploring regulatory networks and
727 functional associations of protein post-translational modifications. *Nucleic Acids Res.* *50*,
728 D471–D479.
- 729 54. Goodsell, D.S. (2021). Fifty years of open access to PDB structures. *RCSB Protein Data*
730 *Bank.* [10.2210/rcsb_pdb/mom_2021_10](https://doi.org/10.2210/rcsb_pdb/mom_2021_10).
- 731 55. UniProt Consortium (2021). UniProt: the universal protein knowledgebase in 2021. *Nucleic*
732 *Acids Res.* *49*, D480–D489.
- 733 56. Mistry, J., Chuguransky, S., Williams, L., Qureshi, M., Salazar, G.A., Sonnhammer, E.L.L.,
734 Tosatto, S.C.E., Paladin, L., Raj, S., Richardson, L.J., et al. (2021). Pfam: The protein
735 families database in 2021. *Nucleic Acids Res.* *49*, D412–D419.
- 736 57. Sigrist, C.J.A., Cerutti, L., de Castro, E., Langendijk-Genevaux, P.S., Bulliard, V., Bairoch,
737 A., and Hulo, N. (2010). PROSITE, a protein domain database for functional
738 characterization and annotation. *Nucleic Acids Res.* *38*, D161–D166.

- 739 58. Clerc, O., Deniaud, M., Vallet, S.D., Naba, A., Rivet, A., Perez, S., Thierry-Mieg, N., and
740 Ricard-Blum, S. (2019). MatrixDB: integration of new data with a focus on
741 glycosaminoglycan interactions. *Nucleic Acids Res.* *47*, D376–D381.
- 742 59. Chatr-Aryamontri, A., Breitkreutz, B.-J., Oughtred, R., Boucher, L., Heinicke, S., Chen, D.,
743 Stark, C., Breitkreutz, A., Kolas, N., O'Donnell, L., et al. (2015). The BioGRID interaction
744 database: 2015 update. *Nucleic Acids Res.* *43*, D470-8.
- 745 60. Szklarczyk, D., Gable, A.L., Lyon, D., Junge, A., Wyder, S., Huerta-Cepas, J., Simonovic,
746 M., Doncheva, N.T., Morris, J.H., Bork, P., et al. (2019). STRING v11: protein-protein
747 association networks with increased coverage, supporting functional discovery in genome-
748 wide experimental datasets. *Nucleic Acids Res.* *47*, D607–D613.
- 749 61. Franz, M., Rodriguez, H., Lopes, C., Zuberi, K., Montojo, J., Bader, G.D., and Morris, Q.
750 (2018). GeneMANIA update 2018. *Nucleic Acids Res.* *46*, W60–W64.
- 751 62. Malaker, S.A., Riley, N.M., Shon, D.J., Pedram, K., Krishnan, V., Dorigo, O., and Bertozzi,
752 C.R. (2022). Revealing the human mucinome. *Nat. Commun.* *13*, 3542.
- 753 63. Ryan, S.O., and Cobb, B.A. (2012). Roles for major histocompatibility complex
754 glycosylation in immune function. *Semin. Immunopathol.* *34*, 425–441.
- 755 64. Ilca, F.T., and Boyle, L.H. (2021). The glycosylation status of MHC class I molecules
756 impacts their interactions with TAPBPR. *Mol. Immunol.* *139*, 168–176.
- 757 65. Hoek, M., Demmers, L.C., Wu, W., and Heck, A.J.R. (2021). Allotype-specific glycosylation
758 and cellular localization of human leukocyte antigen class I proteins. *J. Proteome Res.* *20*,
759 4518–4528.
- 760 66. Sagt, C.M., Kleizen, B., Verwaal, R., de Jong, M.D., Müller, W.H., Smits, A., Visser, C.,
761 Boonstra, J., Verkleij, A.J., and Verrips, C.T. (2000). Introduction of an N-glycosylation site
762 increases secretion of heterologous proteins in yeasts. *Appl. Environ. Microbiol.* *66*, 4940–
763 4944.
- 764 67. Olczak, M., and Szulc, B. (2021). Modified secreted alkaline phosphatase as an improved
765 reporter protein for N-glycosylation analysis. *PLoS One* *16*, e0251805.
- 766 68. Harbison, A.M., Fogarty, C.A., Phung, T.K., Satheesan, A., Schulz, B.L., and Fadda, E.
767 (2022). Fine-tuning the spike: role of the nature and topology of the glycan shield in the
768 structure and dynamics of the SARS-CoV-2 S. *Chem. Sci.* *13*, 386–395.
- 769 69. Wei, C.-J., Boyington, J.C., Dai, K., Houser, K.V., Pearce, M.B., Kong, W.-P., Yang, Z.-Y.,
770 Tumpey, T.M., and Nabel, G.J. (2010). Cross-neutralization of 1918 and 2009 influenza
771 viruses: role of glycans in viral evolution and vaccine design. *Sci. Transl. Med.* *2*, 24ra21.
- 772 70. Go, E.P., Ding, H., Zhang, S., Ringe, R.P., Nicely, N., Hua, D., Steinbock, R.T., Golabek,
773 M., Alin, J., Alam, S.M., et al. (2017). Glycosylation Benchmark Profile for HIV-1 Envelope
774 Glycoprotein Production Based on Eleven Env Trimers. *J. Virol.* *91*. 10.1128/JVI.02428-16.
- 775 71. Grant, O.C., Montgomery, D., Ito, K., and Woods, R.J. (2020). Analysis of the SARS-CoV-2
776 spike protein glycan shield reveals implications for immune recognition. *Sci. Rep.* *10*,
777 14991.

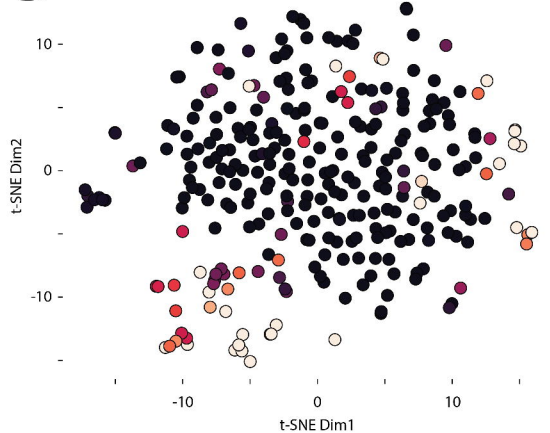
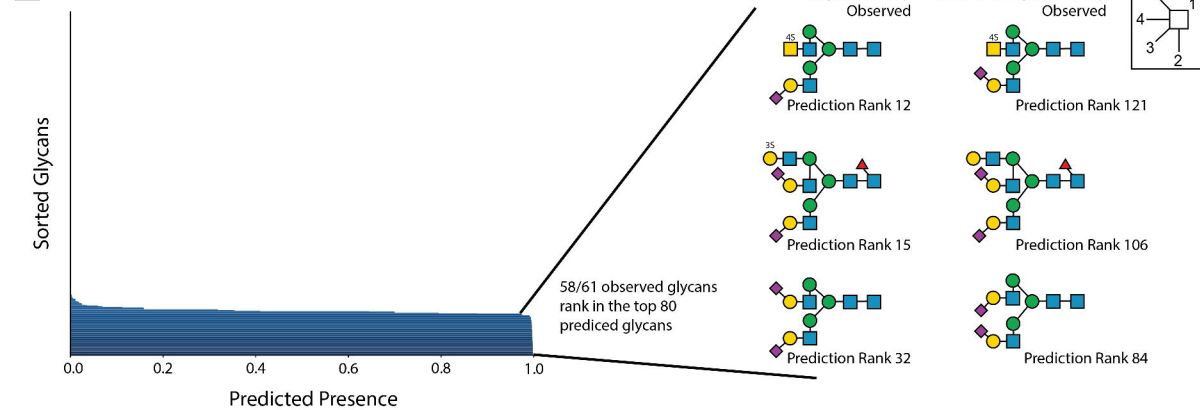
- 778 72. Mariethoz, J., Alocci, D., Gastaldello, A., Horlacher, O., Gasteiger, E., Rojas-Macias, M.,
779 Karlsson, N.G., Packer, N.H., and Lisacek, F. (2018). Glycomics@ExPASy: Bridging the
780 Gap. *Molecular & Cellular Proteomics* 17, 2164–2176. 10.1074/mcp.ra118.000799.
- 781 73. York, W.S., Mazumder, R., Ranzinger, R., Edwards, N., Kahsay, R., Aoki-Kinoshita, K.F.,
782 Campbell, M.P., Cummings, R.D., Feizi, T., Martin, M., et al. (2020). GlyGen:
783 Computational and Informatics Resources for Glycoscience. *Glycobiology* 30, 72–73.
- 784 74. Mih, N., Brunk, E., Chen, K., Catoi, E., Sastry, A., Kavvas, E., Monk, J.M., Zhang, Z., and
785 Palsson, B.O. (2018). ssbio: a Python framework for structural systems biology.
786 *Bioinformatics* 34, 2155–2157.
- 787 75. Velankar, S., Dana, J.M., Jacobsen, J., van Ginkel, G., Gane, P.J., Luo, J., Oldfield, T.J.,
788 O'Donovan, C., Martin, M.-J., and Kleywegt, G.J. (2013). SIFTS: Structure Integration with
789 Function, Taxonomy and Sequences resource. *Nucleic Acids Res.* 41, D483-9.
- 790 76. Klein, J., and Zaia, J. (2019). glypy: An Open Source Glycoinformatics Library. *J. Proteome*
791 *Res.* 18, 3532–3537.
- 792 77. Inoue, H. (2019). Multi-sample dropout for accelerated training and better generalization.
793 arXiv [cs.NE].
- 794 78. Thomès, L., Burkholz, R., and Bojar, D. (2021). Glycowork: A Python package for glycan
795 data science and machine learning. *Glycobiology* 31, 1240–1244.
- 796 79. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z.,
797 Gimelshein, N., Antiga, L., et al. (2019). PyTorch: An Imperative Style, High-Performance
798 Deep Learning Library. In *Advances in Neural Information Processing Systems*, H. Wallach,
799 H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, eds.
800 (Curran Associates, Inc.).
- 801 80. Glorot, X., and Bengio, Y. (2010). Understanding the difficulty of training deep feedforward
802 neural networks. In *Proceedings of the Thirteenth International Conference on Artificial*
803 *Intelligence and Statistics Proceedings of Machine Learning Research.*, Y. W. Teh and M.
804 Titterton, eds. (PMLR), pp. 249–256.
- 805 81. Izmailov, P., Podoprikin, D., Garipov, T., Vetrov, D., and Wilson, A.G. (2018). Averaging
806 weights leads to wider optima and better generalization. arXiv [cs.LG].
- 807 82. Sharon, N. (1986). IUPAC-IUB Joint Commission on Biochemical Nomenclature (JCBN).
808 Nomenclature of glycoproteins, glycopeptides and peptidoglycans. *Glycoconj. J.* 3, 123–
809 133.
- 810 83. McNaught, A.D. (1997). Nomenclature of carbohydrates. *Carbohydr. Res.* 297, 1–92.
- 811 84. Herget, S., Ranzinger, R., Maass, K., and Lieth, C.-W.V.D. (2008). GlycoCT-a unifying
812 sequence format for carbohydrates. *Carbohydr. Res.* 343, 2162–2171.

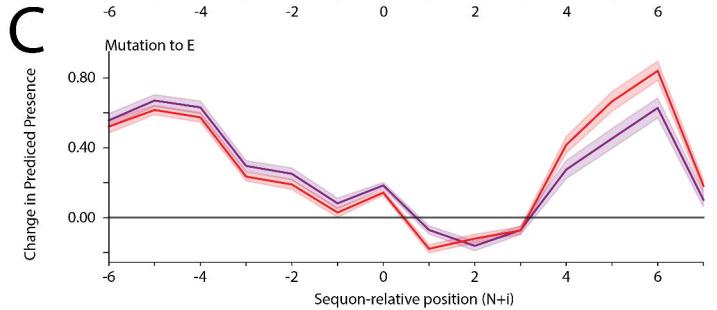
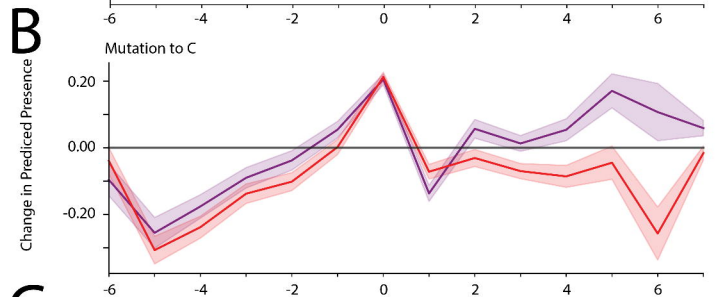
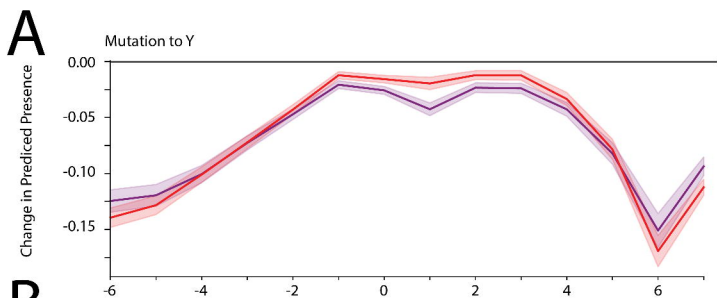


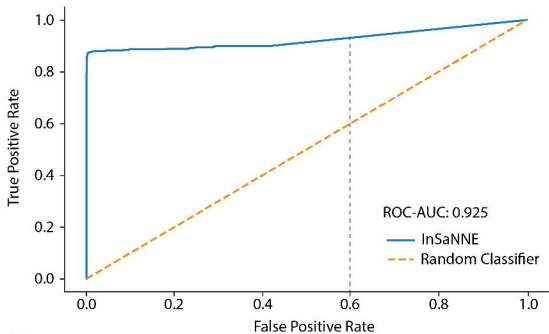
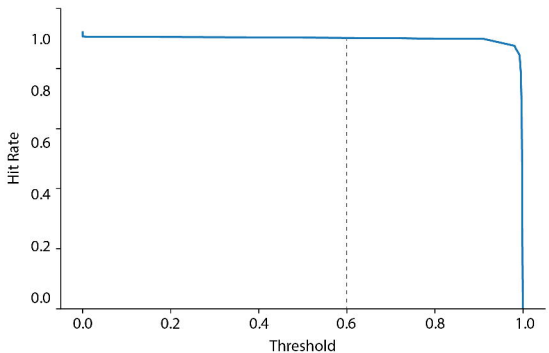
A**B****C**

Average Accuracy

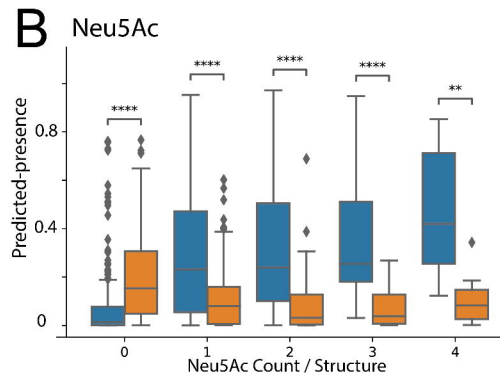
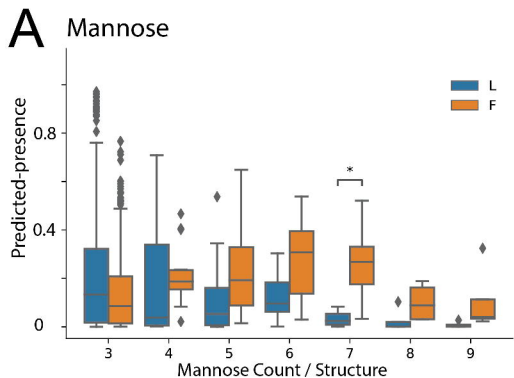
- 0.0
- 0.2
- 0.4
- 0.6
- 0.8

D**E**

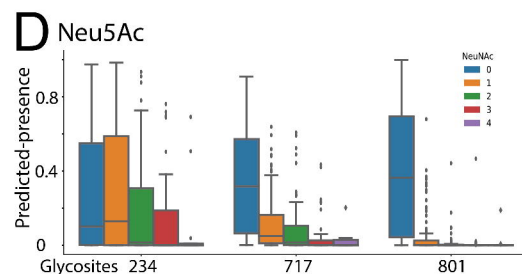
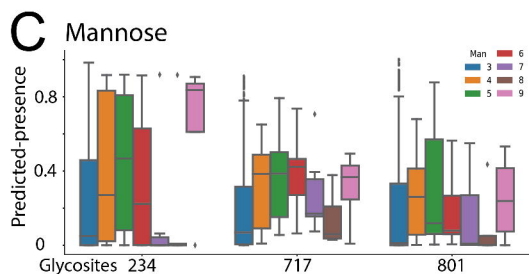


A**B**

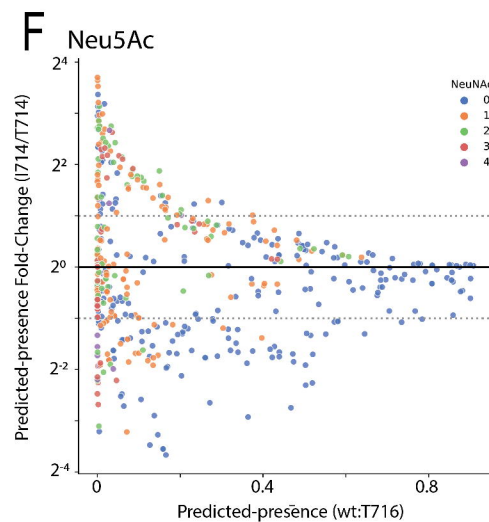
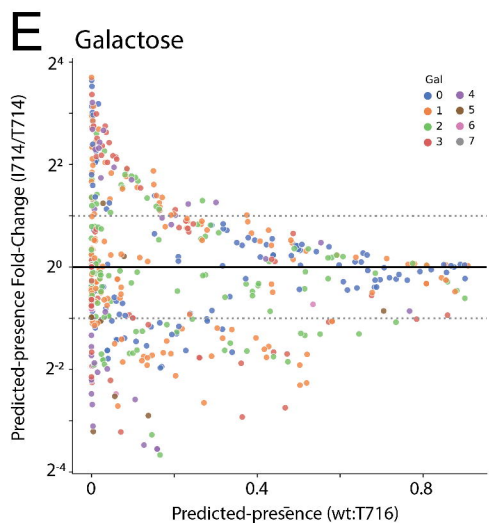
Predicted Complex and Oligomannose Structures: L to F Mutant

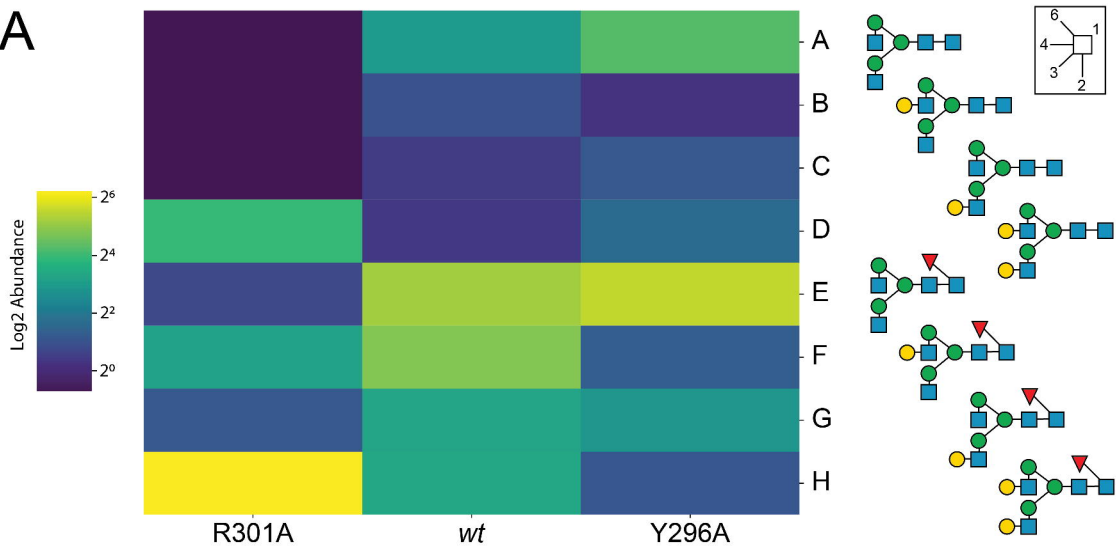


Predicted Glycosylation for SARS-CoV-2 Spike: wild-type/ancestral

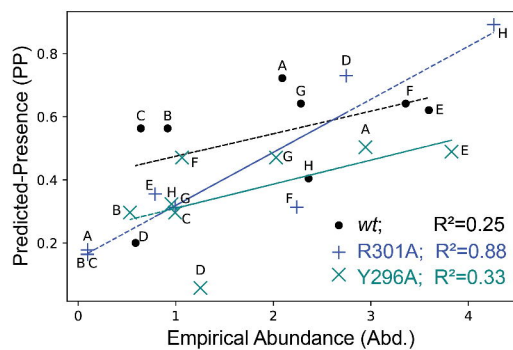


Predicted Differential Glycosylation for SARS-CoV-2 Spike, site N717: B.1.1.7 vs ancestral (T716)

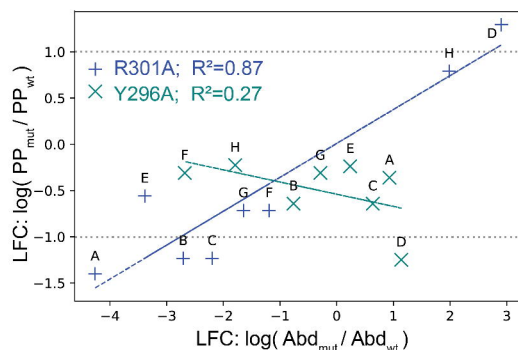


A**B**

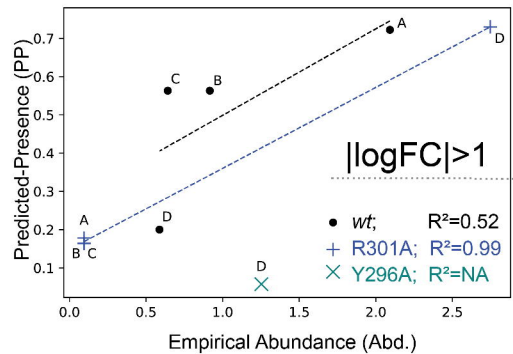
Predicted vs. Empirical Abundance

**C**

Log Fold Change (LFC)

**D**

Predicted vs. Empirical Abundance

**E**Log Fold Change ($\log(\text{mut}/\text{wt})$)