

Table 1: Chromatogram Feature Definitions

\mathcal{M}_i	The neutral mass of the i th chromatogram
\mathcal{I}_i	The total intensity array assigned to the i th chromatogram
$\mathcal{I}_{i,j}$	The sum of all peak intensities for peaks observed in the j th scan for the i th chromatogram
$\mathcal{I}_{i,j,k}$	The intensity assigned to the k th peak at the j th scan for the i th chromatogram
\mathbf{c}_i	The set of charge states observed for the i th chromatogram
$\mathcal{I}_{i,c=j}$	The total intensity assigned to the i th chromatogram with charge state j
$\mathbf{t}_{i,j}$	The time of the j th scan of the i th chromatogram
$\mathbf{env}_{i,j,k}$	The normalized experimental isotopic envelope composing the k th peak of the j th scan of the i th chromatogram, whose members sum to 1
\mathbf{a}_i	The set of adduction states observed for the i th chromatogram
$\mathcal{I}_{i,a=j}$	The total intensity assigned to the i th chromatogram with adduct j
\hat{g}_i	The glycan composition assigned to the i th chromatogram, or \emptyset if there was no matched glycan composition

1 Chromatographic Feature Evaluation

For each candidate feature, we computed several statistics to estimate how distinguishable the observed signal was from random noise. We use the following quantities from each LC-MS feature:

1.1 Chromatographic Peak Shape

An LC-MS elution profile should be composed of one or more peak-like components, each following a bi-Gaussian peak shape model (Yu and Peng (2010)) or in less ideal chromatographic circumstances, a skewed Gaussian peak shape model. We fit these models using non-linear least squares (NLS). As measures of goodness of fit are not generally available for NLS, we use the following criterion:

$$\begin{aligned}
 \hat{y}_i &= NLS(\mathcal{I}_i, \mathbf{t}_i) \\
 e_{i,NLS} &= \mathcal{I}_i - \hat{y}_i \\
 \bar{y}_i &= \mathbf{t}_i \left((\mathbf{t}_i^t \mathbf{t}_i)^{-1} \mathbf{t}_i \mathcal{I}_i \right) \\
 e_{i,null} &= \mathcal{I}_i - \bar{y}_i \\
 \mathcal{L}_i &= 1 - \frac{\sum e_{i,NLS}^2}{\sum e_{i,null}^2}
 \end{aligned} \tag{1}$$

where line score describes how much the peak shape fit improves on a straight line fit null model.

We apply two competitive peak fitting strategies to address distorted, overlapping, or multimodal elution profiles. The first works iteratively by finding a best-matching peak shape using non-linear least squares, subtracting the fitted signal and checks if there is another peak with at least half as tall as the removed peak, if so repeating the process until no peak can be found, saving each peak model so constructed. The second approach starts by locating local minima between putative peaks, and partitioning the chromatogram into sub-groups which would be fit independently. This method generates a candidate list of minima, and selects the case which has the greatest difference between the minimum and its pair of maxima to split the feature at. The strategy which produces the maximum \mathcal{L}_i is chosen.

1.2 Composition Dependent Charge State Distribution

As the number of monosaccharides composing a glycan increases, the number of possible sites for charge localization increases. Under normal conditions, we would expect to observe the same molecule in multiple charge states (Maxwell *et al.* (2012)).

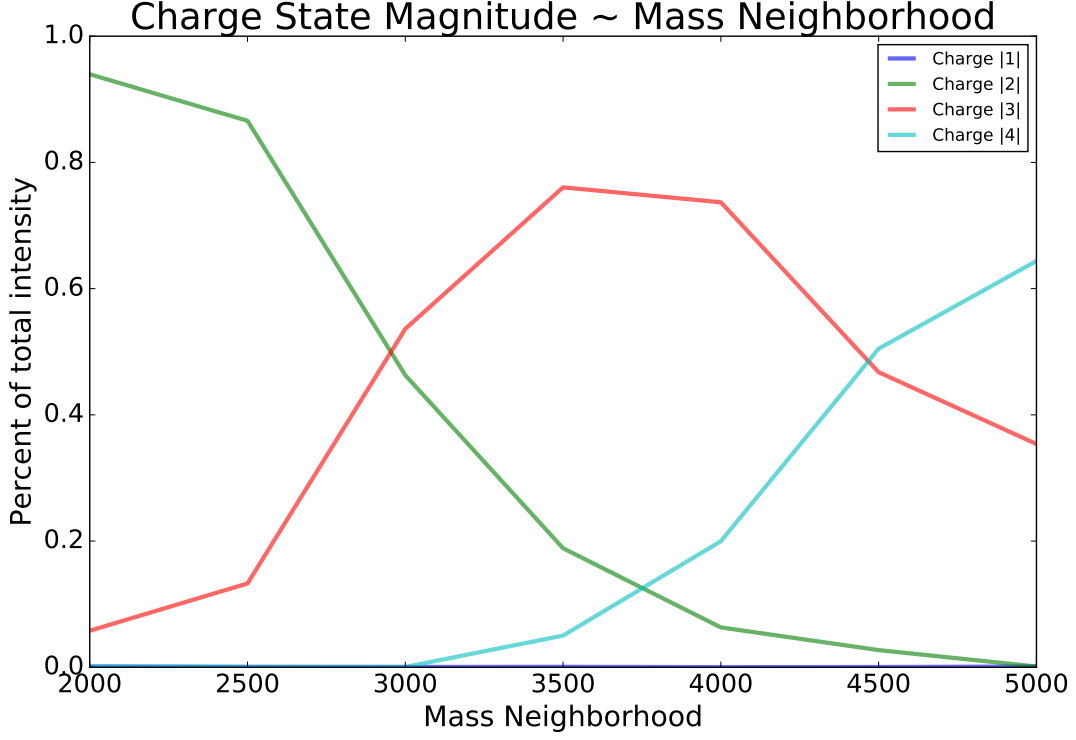


Figure 1: The trend of charge state relative abundance for acidic glycans

Which charge states are expected would depend upon the size of the molecule and it's constituent units' electronegativity. In it's native state, **NeuAc**'s acidic group causes glycans with one or more **NeuAc** to have a propensity for higher negative charge states (Varki and Schauer (2009)). To capture this relationship, we modeled the probability of observing a glycan composition for sialylated and unsialylated compositions separately. For permethylated glycans, charge is carried by protons or metallic cation adducts like sodium, the relationship between acidic monosaccharides and charge state propensities is weaker.

$$\begin{aligned}
m_i &= (\lfloor (\mathcal{M}_i/w)/10 \rfloor + 1) * 10 \\
\mathcal{H}_{i,j} &= \frac{\mathcal{I}_{i,c=j}}{\mathcal{I}_i} \\
P(c, m) &= |m| \sum_{m_i \in m} \mathcal{H}_{i,j} \\
\mathcal{C}_i &= \sum_{c_{i,j} \in \mathbf{c}_i} P(c_{i,j}, m_i)
\end{aligned} \tag{2}$$

where w is the width of the mass bin divided by 10 and $P(c, m)$ is defined as part of the model estimation procedure.

1.3 Adduction Rate

For the samples *AGP-permethylated-2ul-inj-55-SLens* and *Perm-BS-070111-04-Human-Serum* we also include an Adduction Frequency model score \mathcal{A}_i , following the same pattern as the charge state distribution, with the same extension of justification from Maxwell *et al.* (2012). We use one mass scaling model for all glycan compositions as ammonium adduction is not expected to be composition dependent.

$$\begin{aligned}
\mathcal{H}_{i,j} &= \frac{\mathcal{I}_{i,a=j}}{\mathcal{I}_i} \\
P(a, m) &= |m| \sum_{m_i \in m} \mathcal{H}_{i,j} \\
\mathcal{A}_i &= \sum_{a_{i,j} \in \mathbf{a}_i} P(a_{i,j}, m_i)
\end{aligned} \tag{3}$$

We fit the adduction rate model on *AGP-permethyated-2ul-inj-55-SLens* in order to make our comparison to third-party data less biased given limited sample data.

1.4 Isotopic Pattern Consistency

Our ahead-of-time deconvolution procedure uses an average isotopic model and does not capture the consistency of the isotopic pattern that was fit with the isotopic pattern of the glycan composition that matched that peak. The criterion

$$\mathcal{J}_i = 1 - 2\mathcal{I}_i^{-t} \mathbf{I}_i \sum_j^J \sum_k^K \mathcal{I}_{i,j,k} \mathbf{env}_{i,j,k}^t (\ln \mathbf{env}_{i,j,k} - \ln \mathbf{tid}_i) \quad (4)$$

where \mathbf{tid} is the theoretical isotopic pattern derived from either \hat{g}_i or an average interpolated for \mathcal{M}_i if $\hat{g}_i = \emptyset$. This computes a per-peak intensity weighted mean G-test comparing the goodness of fit between the experimental envelope and the theoretical isotopic pattern.

1.5 Observation Spacing Score

The less time between observations of a glycan composition the less likely the chromatogram is to contain peaks missing or caused by isotopic pattern interference or missing information.

$$\mathcal{T}_i = 1 - 2\mathcal{I}_i^{-t} \mathbf{I}_i \sum_{j=1}^J \mathcal{I}_{i,j} (\mathbf{t}_{i,j} - \mathbf{t}_{i,j-1}) \quad (5)$$

1.6 Summarization Score

Each scoring feature $\in [\mathcal{L}_i, \mathcal{C}_i, \mathcal{J}_i, \mathcal{T}_i]$ is penalized by $\epsilon = 1e-6$ bounded in the range $[0, 1)$, with values below 0 set to ϵ .

$$s_i = \sum_{f_{i,j} \in \text{features}_i} \ln \frac{f_{i,j}}{1 - f_{i,j}} \quad (6)$$

producing a value between $(-\infty, \infty)$. $s_i < 8$ reflects multiple poor feature scores and is unexpected to be real, while $s_i > 15$ is consistent with model expectations.

2 MS^n Signature Ion Criterion

When MS^n scans are present, it may be useful to consider only those MS^1 features which are associated with MS^n scans that contain glycan-like signature ions. We include an algorithm for classifying an MS^n scan as being "glycan-like":

$$I = \max(\text{intensity}(p)) \quad (7)$$

$$t = I * 0.01 \quad (8)$$

$$p_{\text{oxonium}} = \{p_i \leftarrow |ppmerror(\text{mass}(p_j), \text{mass}(f_g))| < e, f_g \in \text{oxonium}(g), f_g \neq \text{Fucose}, \text{intensity}(p_i) > t\} \quad (9)$$

$$p_{\text{edges}} = \{(p_i, p_j) \leftarrow |ppmerror(\text{mass}(p_j) - \text{mass}(p_i), \text{mass}(f_g))| < e, \text{oxonium}(f_g) \in g, \text{intensity}(p_i) > t, \text{intensity}(p_j) > t\} \quad (10)$$

$$s_{\text{oxonium}} = \frac{1}{|p_{\text{oxonium}}|} \sum_{p_i \in p_{\text{oxonium}}} \left(\frac{\text{intensity}(p_i)}{I} \right) * \min(\log_4 |p_{\text{oxonium}}|, 1) \quad (11)$$

$$s_{\text{edges}} = \frac{1}{|p_{\text{edges}}|} \sum_{p_i, p_j \in p_{\text{edges}}} \left(\frac{\text{intensity}(p_i) + \text{intensity}(p_j)}{I} \right) * \min(\log_4 |p_{\text{edges}}|, 1) \quad (12)$$

$$s_g = \max(s_{\text{oxonium}}, s_{\text{edges}}) \quad (13)$$

$$(14)$$

Where p is the set of peaks in the scan, g is the glycan composition, e the required parts-per-million mass accuracy. $\text{oxonium}()$ is a function that given a glycan composition g , produces fragments f_g of g composed of between one and three monosaccharides, commonly observed as oxonium ions alone, or as the mass difference between two peaks formed from consecutive fragmentation of a glycosidic bond. This method is not intended to identify a glycan structure, just detect patterns in the signal peaks of the MS^n scan that could indicate the fragmentation of a glycan.

3 A more complete derivation of $\hat{\phi}$

To obtain the optimal ϕ , we take the partial derivative of ℓ w.r.t ϕ_m

$$0 = \frac{\partial \ell}{\partial \phi_m} \left((\mathbf{s} - \phi_o)^t (\mathbf{s} - \phi_o) + \lambda [\phi_o - \tau_o, \phi_m - \tau_m] \begin{bmatrix} \mathbf{L}_{oo} & \mathbf{L}_{om} \\ \mathbf{L}_{mo} & \mathbf{L}_{mm} \end{bmatrix} \begin{bmatrix} \phi_o - \tau_o \\ \phi_m - \tau_m \end{bmatrix} \right) \quad (15)$$

$$\begin{aligned} &= \lambda (\phi_o - \tau_o)^t \mathbf{L}_{om} + \lambda \mathbf{L}_{mo} (\phi_o - \tau_o) + \lambda (\phi_m - \tau_m)^t (\mathbf{L}_{mm}^t + \mathbf{L}_{mm}) \\ &= 2\lambda \mathbf{L}_{mo} (\phi_o - \tau_o) + 2\lambda \mathbf{L}_{mm} (\phi_m - \tau_m) \\ -\mathbf{L}_{mm} (\phi_m - \tau_m) &= \mathbf{L}_{mo} (\phi_o - \tau_o) \\ (\phi_m - \tau_m) &= -\mathbf{L}_{mm}^{-1} \mathbf{L}_{mo} (\phi_o - \tau_o) \\ \hat{\phi}_m &= -\mathbf{L}_{mm}^{-1} \mathbf{L}_{mo} (\phi_o - \tau_o) + \tau_m \end{aligned} \quad (16)$$

and w.r.t. ϕ_o

$$0 = \frac{\partial \ell}{\partial \phi_o} \left((\mathbf{s} - \phi_o)^t (\mathbf{s} - \phi_o) + \lambda [\phi_o - \tau_o, \phi_m - \tau_m] \begin{bmatrix} \mathbf{L}_{oo} & \mathbf{L}_{om} \\ \mathbf{L}_{mo} & \mathbf{L}_{mm} \end{bmatrix} \begin{bmatrix} \phi_o - \tau_o \\ \phi_m - \tau_m \end{bmatrix} \right) \quad (17)$$

$$\begin{aligned} &= -2\mathbf{s} + 2\phi_o + \lambda (\mathbf{L}_{oo} + \mathbf{L}_{oo}^t) (\phi_o - \tau_o) + \lambda \mathbf{L}_{om} (\phi_m - \tau_m) + \lambda \mathbf{L}_{mo}^t (\phi_m - \tau_m) \\ &= -2\mathbf{s} + 2\phi_o + 2\lambda \mathbf{L}_{oo} (\phi_o - \tau_o) + 2\lambda \mathbf{L}_{om} (\phi_m - \tau_m) \\ \mathbf{s} &= \phi_o + \lambda (\mathbf{L}_{oo} (\phi_o - \tau_o) + \mathbf{L}_{om} (\phi_m - \tau_m)) \\ &= \phi_o + \lambda (\mathbf{L}_{oo} (\phi_o - \tau_o) + \mathbf{L}_{om} (-\mathbf{L}_{mm}^{-1} \mathbf{L}_{mo} (\phi_o - \tau_o) + \tau_m - \tau_m)) \\ &= \phi_o + \lambda (\mathbf{L}_{oo} (\phi_o - \tau_o) - \mathbf{L}_{om} \mathbf{L}_{mm}^{-1} \mathbf{L}_{mo} (\phi_o - \tau_o)) \\ \mathbf{s} - \tau_o &= \phi_o - \tau_o + \lambda (\mathbf{L}_{oo} (\phi_o - \tau_o) - \mathbf{L}_{om} \mathbf{L}_{mm}^{-1} \mathbf{L}_{mo} (\phi_o - \tau_o)) \\ &= \mathbf{I} (\phi_o - \tau_o) + \lambda (\mathbf{L}_{oo} (\phi_o - \tau_o) - \mathbf{L}_{om} \mathbf{L}_{mm}^{-1} \mathbf{L}_{mo} (\phi_o - \tau_o)) \\ &= [\mathbf{I} + \lambda (\mathbf{L}_{oo} - \mathbf{L}_{om} \mathbf{L}_{mm}^{-1} \mathbf{L}_{mo})] (\phi_o - \tau_o) \\ (\phi_o - \tau_o) &= [\mathbf{I} + \lambda (\mathbf{L}_{oo} - \mathbf{L}_{om} \mathbf{L}_{mm}^{-1} \mathbf{L}_{mo})]^{-1} (\mathbf{s} - \tau_o) \\ \hat{\phi}_o &= [\mathbf{I} + \lambda (\mathbf{L}_{oo} - \mathbf{L}_{om} \mathbf{L}_{mm}^{-1} \mathbf{L}_{mo})]^{-1} (\mathbf{s} - \tau_o) + \tau_o \end{aligned} \quad (18)$$

4 Estimation of Laplacian Regularization Parameters

We model the relationship between \mathbf{s} , ϕ_o , and τ as a set of gaussian distribution.

$$(\mathbf{s} | \phi_o, \tau) \sim \mathcal{N}(\phi_o, \Sigma) \quad (19)$$

$$\Sigma = \rho \mathbf{I} \quad (20)$$

$$\left(\begin{bmatrix} \phi_o \\ \phi_m \end{bmatrix} \middle| \tau \right) \sim \mathcal{N}(\mathbf{A}\tau, \lambda^{-1} \mathbf{L}^-) \quad (21)$$

$$(\phi_o | \tau) \sim \mathcal{N}(\mathbf{A}_o \tau, \Sigma_{\phi_o}) \quad (22)$$

$$\Sigma_{\phi_o} = \lambda^{-1} (\mathbf{L}_{oo} - \mathbf{L}_{om} \mathbf{L}_{mm}^{-1} \mathbf{L}_{mo})^{-1} \quad (23)$$

$$\tau \sim \mathcal{N}(0, \sigma^2 \mathbf{I}) \quad (24)$$

Fully expanded, this becomes

$$\begin{bmatrix} \mathbf{s} \\ \phi_o \\ \tau \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \Sigma + \Sigma_{\phi_o} + \sigma^2 \mathbf{A}_o \mathbf{A}_o^t & \Sigma_{\phi_o} + \sigma^2 \mathbf{A}_o \mathbf{A}_o^t & \sigma^2 \mathbf{A}_o \\ \Sigma_{\phi_o} + \sigma^2 \mathbf{A}_o \mathbf{A}_o^t & \Sigma_{\phi_o} + \sigma^2 \mathbf{A}_o \mathbf{A}_o^t & \sigma^2 \mathbf{A}_o \\ \sigma^2 \mathbf{A}_o^t & \sigma^2 \mathbf{A}_o^t & \sigma^2 \mathbf{I} \end{bmatrix} \right) \quad (25)$$

We can form the conditional distribution $\tau | \mathbf{s}$ which has a mean

$$\mu_{\tau | \mathbf{s}} = 0 + (\sigma^2 \mathbf{A}_o^t) (\Sigma + \Sigma_{\phi_o} + \sigma^2 \mathbf{A}_o \mathbf{A}_o^t)^{-1} \mathbf{s} \quad (26)$$

$$= \mathbf{A}_o^t \left(\tilde{\rho} \mathbf{I} + \frac{1}{\lambda} \mathbf{L}_{oo}^- + \mathbf{A}_o \mathbf{A}_o^t \right)^{-1} \mathbf{s} \quad (27)$$

We assume that $\sigma^2 \gg 1$, and treat λ and ρ as relative to σ^2 , as $\tilde{\rho}$ and $\tilde{\lambda}$. This model gives us an estimate for τ given a value for ρ and λ . As ρ has no direct role in the central tendency of ϕ or \mathbf{s} , we choose to fix the value of $\tilde{\rho} = 0.1$, which leaves only $\tilde{\lambda}$. We estimate the optimal $\tilde{\lambda}$ by grid search, minimizing the predicted residual error sum of squares (PRESS) statistic.

$$\arg \min_{\tilde{\lambda}} \frac{\mathbf{s} - \hat{\phi}_{\mathbf{o}}}{\left(1 - \left(\mathbf{I} + \tilde{\lambda}\mathbf{L}\right)^{-1}\right)^2} \quad (28)$$

This formulation depends upon the value of \mathbf{s} and is sensitive to low scoring matches, which can lead to incorrect estimates of τ and PRESS. We therefore perform a grid search over both $\tilde{\lambda}$ and a minimum threshold for \mathbf{s} , γ .

As we increase γ we remodel the graph \mathcal{G} , removing nodes whose score is below γ . For each pair of neighbors of removed node g_m , (g_u, g_v) , if $L_1(g_u, g_v) > L_1(g_u, g_m) + L_1(g_m, g_v)$, we add an edge from g_u to g_v with weight $\frac{1}{L_1(g_u, g_m) + L_1(g_m, g_v)}$, up to a limit of $L_1(g_k, g_m) < 5$. We give the result of this grid search the name \mathbf{r} . At each point, on the grid, we save the value of τ in $r_{\lambda_i, \gamma_j, \tau}$ and the PRESS in $r_{\lambda_i, \gamma_j, PRESS}$. To select the optimal parameters, we traverse the grid along γ , computing τ_{γ_j} :

$$\bar{\lambda}_j = \arg \min_{\lambda_i} r_{\lambda_i, \gamma_j, PRESS} \quad (29)$$

$$\tau_{\gamma_j} = |r_{\bar{\lambda}_j, \gamma_j, \tau}| * \left(\frac{\gamma_j}{b} + \left(1 - \frac{1}{b}\right) \right) \quad (30)$$

where b is a bias factor defining how much weight to give to higher values of γ which correspond to networks made up of higher confidence assignments. We chose $b = 4$. We define $\bar{\tau}_{\gamma} = \max \tau_{\gamma}$ and define the vector $\bar{\gamma} = [\gamma_j \leftarrow \tau_{\gamma_j} \geq \bar{\tau}_{\gamma} * 0.9]$. This favors values of γ where large values of τ are selected, meaning that the neighborhoods are well populated, while also giving an estimate for $\tilde{\lambda}$ that is non-zero. We term the values of γ in $\bar{\gamma}$ the *target thresholds* of \mathbf{s} .

To estimate $\tilde{\lambda}$ and τ from these results, we select the columns of the grid \mathbf{r} at each $\gamma_j \in \bar{\gamma}$ and applied the following procedure:

$$\bar{\tau}_{\gamma} = \max \tau_{\gamma} \quad (31)$$

$$\bar{\gamma} = \{\gamma_j \leftarrow \tau_{\gamma_j} \geq \bar{\tau}_{\gamma} * 0.9\} \quad (32)$$

$$\bar{\lambda} = \{\bar{\lambda}_j \leftarrow \gamma_j \in \bar{\gamma}\} \quad (33)$$

$$\mathbf{s}_{\gamma_j} = \{s_i \leftarrow s_i > \gamma_j\} \quad (34)$$

$$\bar{\tau}_{\mathbf{j}} = \mu_{\tau|\mathbf{s}_{\gamma_j}, \bar{\lambda}_j} \quad (35)$$

$$\hat{\lambda} = \frac{1}{|\bar{\lambda}|} \sum_j \bar{\lambda}_j \quad (36)$$

$$\hat{\tau} = \frac{1}{|\bar{\tau}|} \sum_j \bar{\tau}_{\mathbf{j}} \quad (37)$$

$$\hat{\gamma} = \frac{1}{|\bar{\gamma}|} \sum_j \bar{\gamma}_j \quad (38)$$

where \mathbf{s}_{γ_j} is the set of observed scores which are greater than γ_j , but where the estimation of is carried out with the complete Laplacian \mathbf{L} , not the reduced network used to compute \mathbf{r} . This set of averaged estimates of $\hat{\lambda}$ and $\hat{\tau}$ are then used to estimate $\hat{\phi}_{\mathbf{o}}$ by 18, labeled ?? in the main text.

5 Algorithmic Performance on All Datasets

TODO

6 Differences in Assigned Glycans for *Perm-BS-070111-04-Serum*

Of the compositions assigned by our algorithm that were not mentioned in Yu *et al.* (2013) but were annotated in the original publication of this dataset in Hu and Mechref (2012) include **HexNAc3 Hex4**, **HexNAc3 Hex4 NeuAc1**, and **HexNAc5 Hex3**. Because our database was constructed based on combinatorial rules that did not take into account all biosynthetic constraints, we include infeasible compositions in our search space, such as **HexNAc2 Hex10 Fuc1** and **HexNAc5 Hex3**

Fuc1 NeuAc2. Future work could be done to restrict the database to only biosynthetically feasible glycan compositions. This would also have benefits for the construction of the composition network where only those compositions which have an enzymatic reaction to from one to the other would have an edge connecting them, such that **HexNAc5 Hex6 NeuAc2** would not have an edge to **HexNAc5 Hex7 NeuAc2** as in our current model.

References

- Hu, Y. and Mechref, Y. (2012). Comparing MALDI-MS, RP-LC-MALDI-MS and RP-LC-ESI-MS glycomic profiles of permethylated N-glycans derived from model glycoproteins and human blood serum. *Electrophoresis*, **33**(12), 1768–1777.
- Maxwell, E., Tan, Y., Tan, Y., Hu, H., Benson, G., Aizikov, K., Conley, S., Staples, G. O., Slys, G. W., Smith, R. D., and Zaia, J. (2012). GlycReSoft: a software package for automated recognition of glycans from LC/MS data. *PloS one*, **7**(9), e45474.
- Varki, A. and Schauer, R. (2009). *Sialic Acids*. Cold Spring Harbor Laboratory Press.
- Yu, C.-Y. C.-Y., Mayampurath, A., Hu, Y., Zhou, S., Mechref, Y., and Tang, H. (2013). Automated annotation and quantification of glycans using liquid chromatography-mass spectrometry. *Bioinformatics*, **29**(13), 1706–1707.
- Yu, T. and Peng, H. (2010). Quantification and deconvolution of asymmetric LC-MS peaks using the bi-Gaussian mixture model and statistical model selection. *BMC bioinformatics*, **11**(1), 559.