

GlyCompare

Bokan, Ben

Take home message

- Study the motif based on glycoprofile from the synthesis network perspective.
- The biosynthesis of glycan is determined by glycosyltransferases that assemble monosaccharides moieties into linear and branched glycan chains step by step . Since there is no template, the synthesis is totally stochastic process. A glycoprofile has multiple glycan and each glycan are the product of a complex network. It is very hard to simulate the synthesis of each glycan individually by using stochastic model. However, the glycoprofile provides a macroscopic description of this stochastic system. If we break down all glycan into substructures, we can study the progress of this complex network comprehensively.

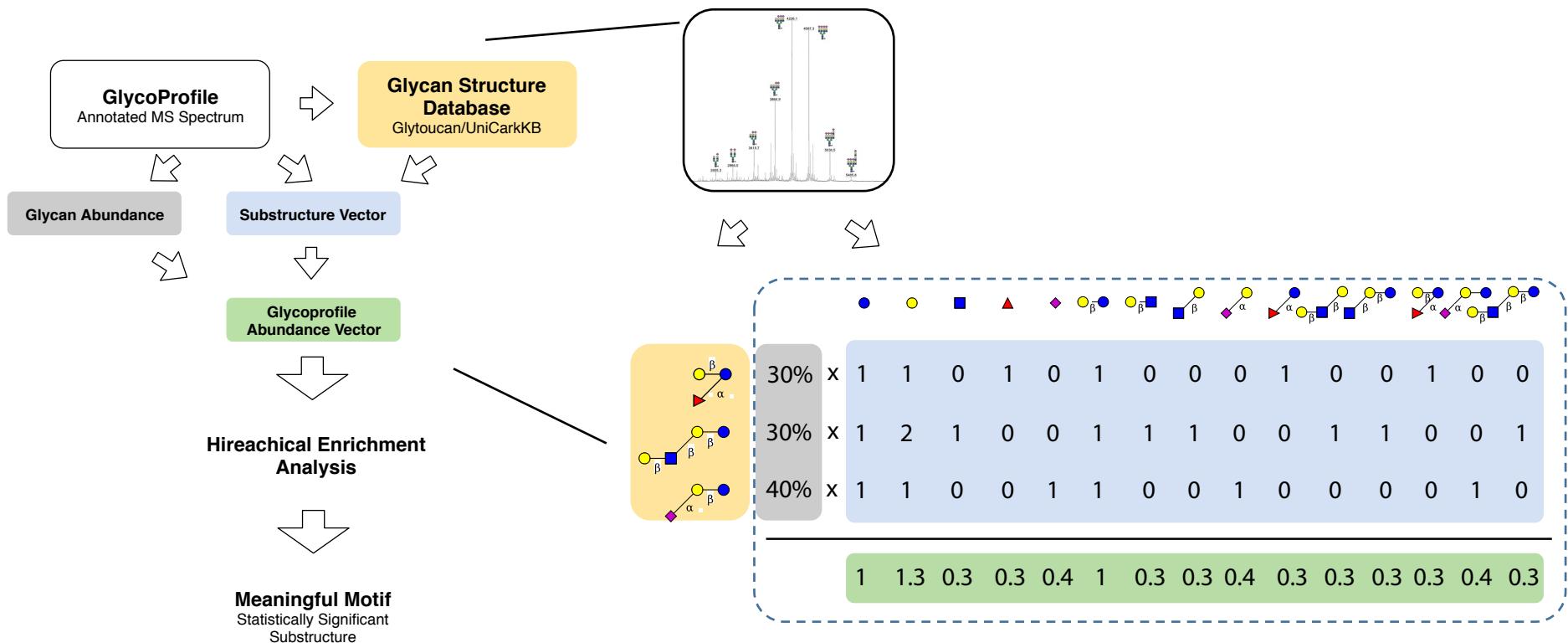
Compare the glycoprofile

- Comparing the glycoprofile is important. It tells us how glycan synthesis process is regulated or influenced by the genome modification. Difference, substructure(motif), synthesis perturbation.
- Reading glycoprofile is hard. Each glycoprofile has multiple glycan structures and each glycan structure has overlapping chemical structure.
- Comparing glycoprofile is even harder. Because each glycoprofile has different glycans, you cannot directly compare them.

GlyCompare

- Comparing large set of glycoprofiles
- Tracking the substructure(motif) abundance change
- Correlating the phenotypes with the motif's abundance.
- Illustrating the perturbation of the synthesis netowork.

Figure 1. Workflow and idea of data preprocessing.

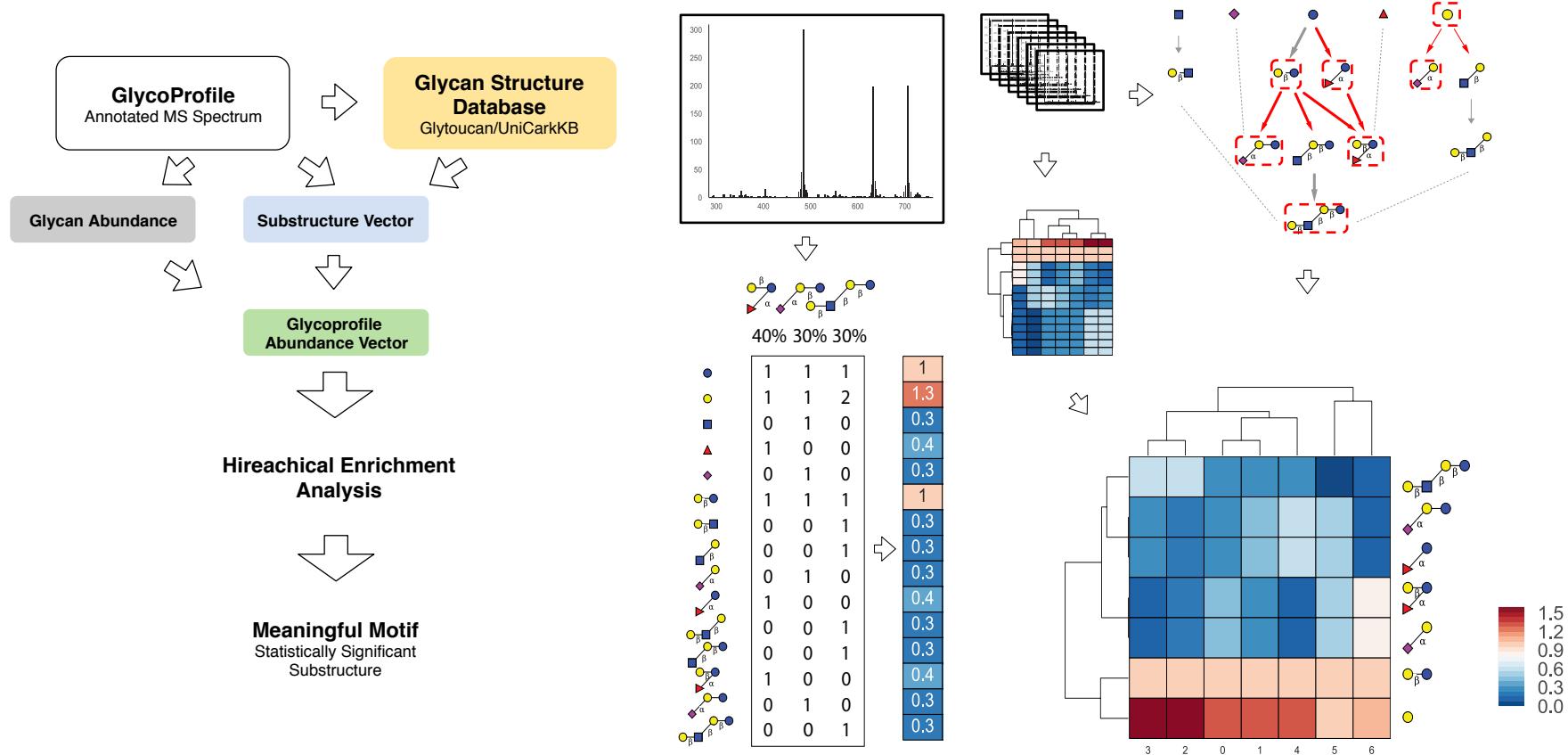


The workflow of the analysis pipeline. The data from Glycoprofile and Glycan structure are vectored to numeric vector and then go to the hierarchical analysis.

White box: pseudo glycoprofile
 Orange box : glycan structure
 Green box : glycoprofile abundance vector
 Grey box : glycan abundance

This plot shows that how the glycoprofile data is transformed to glycan **Substructure Vector** and **Glycoprofile Abundance Vector**

Figure1 Working pipeline



2 Result

We demonstrate our analysis on 16 EPO glycoprofiles curated from NBT paper, the CHO cell lines includes the single or joint knockout from B4galT1/2/3/4/ family, st3gal4/6 family and mgat1/2/4A/4B/5 family. For each glycoprofile, the glycan structure and the relative abundance are obtained.

Result 2.1 Clustering of 16 glycoprofiles

Figure 2A

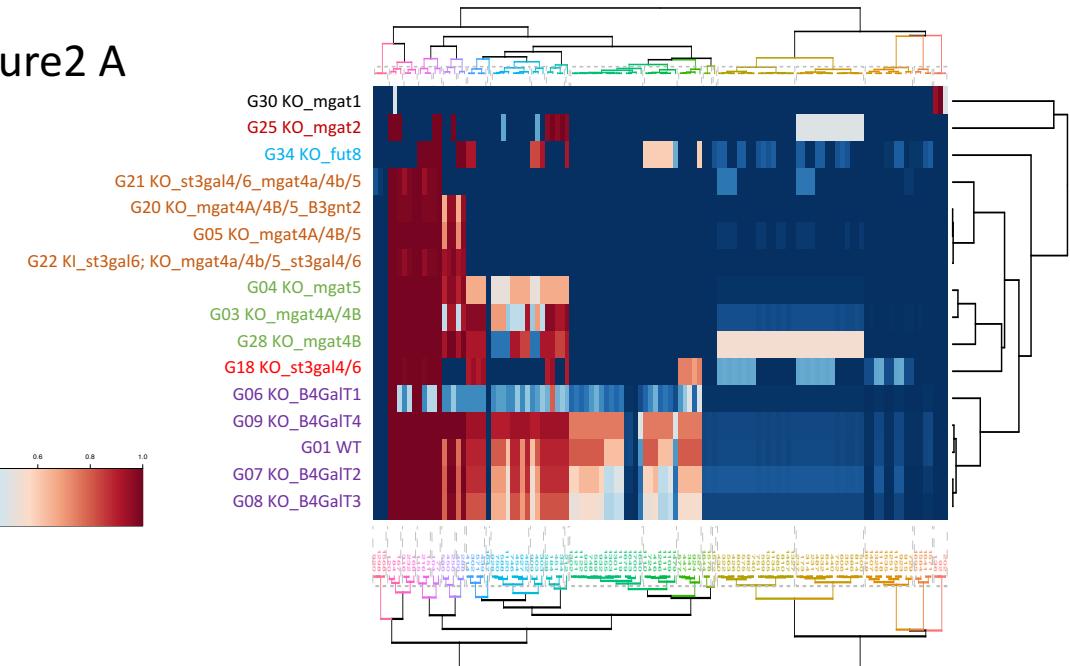
A glycoprofile vector with 722 glycan substructure are generated and then reduced to 118 glycan substructures. The cluster map (Figure 2 A) shows there are 7 different types of glycoprofile. Each cluster has distinguished glycan structure patterns and the distance between each cluster demonstrate the structure difference between each groups of knockout. The cluster contains WT also has B4GalT1, B4GalT4, B4GalT2, B4GalT3 knockout. The mgat cluster mainly has mgat4B, mgat4A/4B, mgat5 knockout. The st3gal4/6 knockout forms one cluster. The cluster with joint knockout has KO_st3gal4/6_mgat4a/4b/5, KO_mgat4A/4B/5_B3gnt2, KO_mgat4A/4B/5, KI_st3gal6; KO_mgat4a/4b/5_st3gal4/6. The glycoprofile with Fuc 8, mgat2, mgat 1 forms individual cluster.

The clustering with annotated glycan is bad because there are only a few glycans.

Figure 2B/Figure2C

- The clustering result is very consistent with NCM paper (Figure2 B). NCM clusters the glycoprofile with the native mass spectrum data. The disagreement comes with the rearrange of the KO mgat2 and KO fu8. Our method take care of the structure difference across isomers. We can distinguish very different structure from very similar molecular mass. Thus, we can find the structure profile of KO mgat2 and KO fut8 are very different from the rest of glycoprofile. Our result is better than NCM 's result.
- On the same time, the glycans abundant table are also used for clustering as a comparison. (Figure2 C) It is worse than the clustering with the glyco-substructure. The glycan abundant table used for clustering is very sparse. The cluster is mainly based on the existence of the glycans and the distances between profiles are not well linearly characterized. The cluster is not consistent with the NCM and our result.

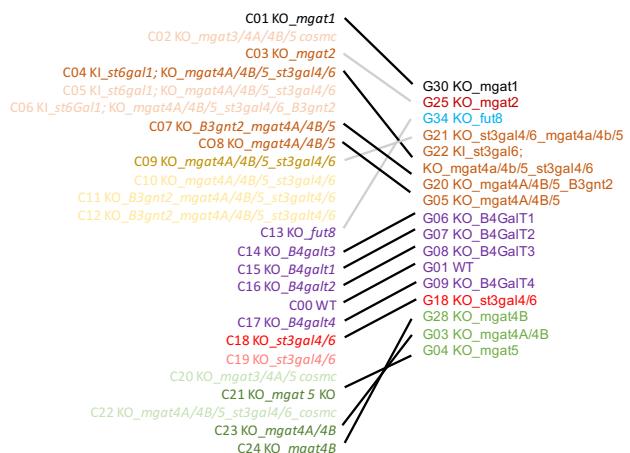
Figure2 A



Left side is the glycoprofiles in the NCM paer, the right side is the shared glycoprofile from our pool. The major difference comes from the position of the KOfu8.

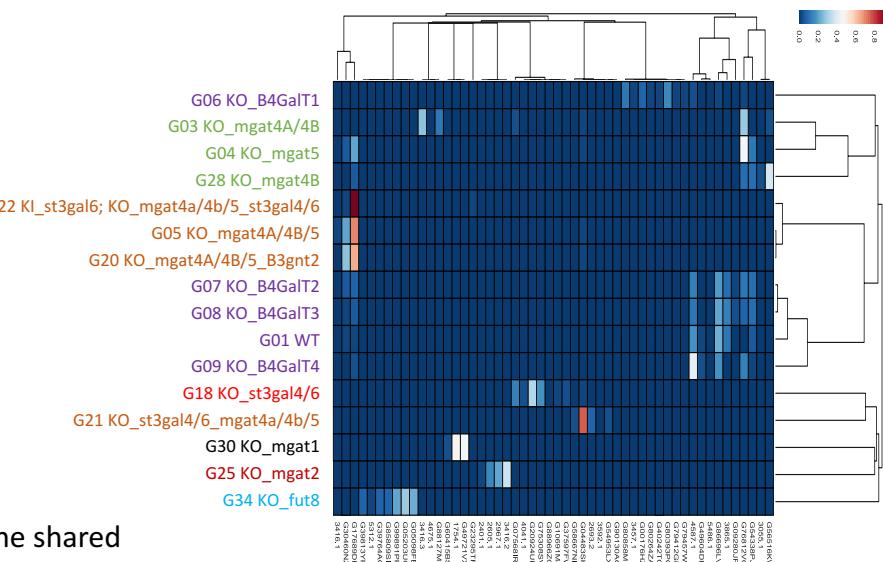
77 glycan structure

Figure2 B



Left side is the glycoprofiles in the NCM paer, the right side is the shared glycoprofile from our pool. The major difference comes from the position of the KOfu8.

Figure2 C



Result 1.2 Characterization of substructure abundance change

- By setting the threshold for cutting the motif cluster. 24 clusters of substructure are picked and the representative substructure for each cluster are generated(see supplementary). The mean abundance of a cluster is used to represent the abundance of the representative substructure. This is the first time to have a comprehensive, automatic quantification of the abundance variation of substructure across multiple glycoprofile.

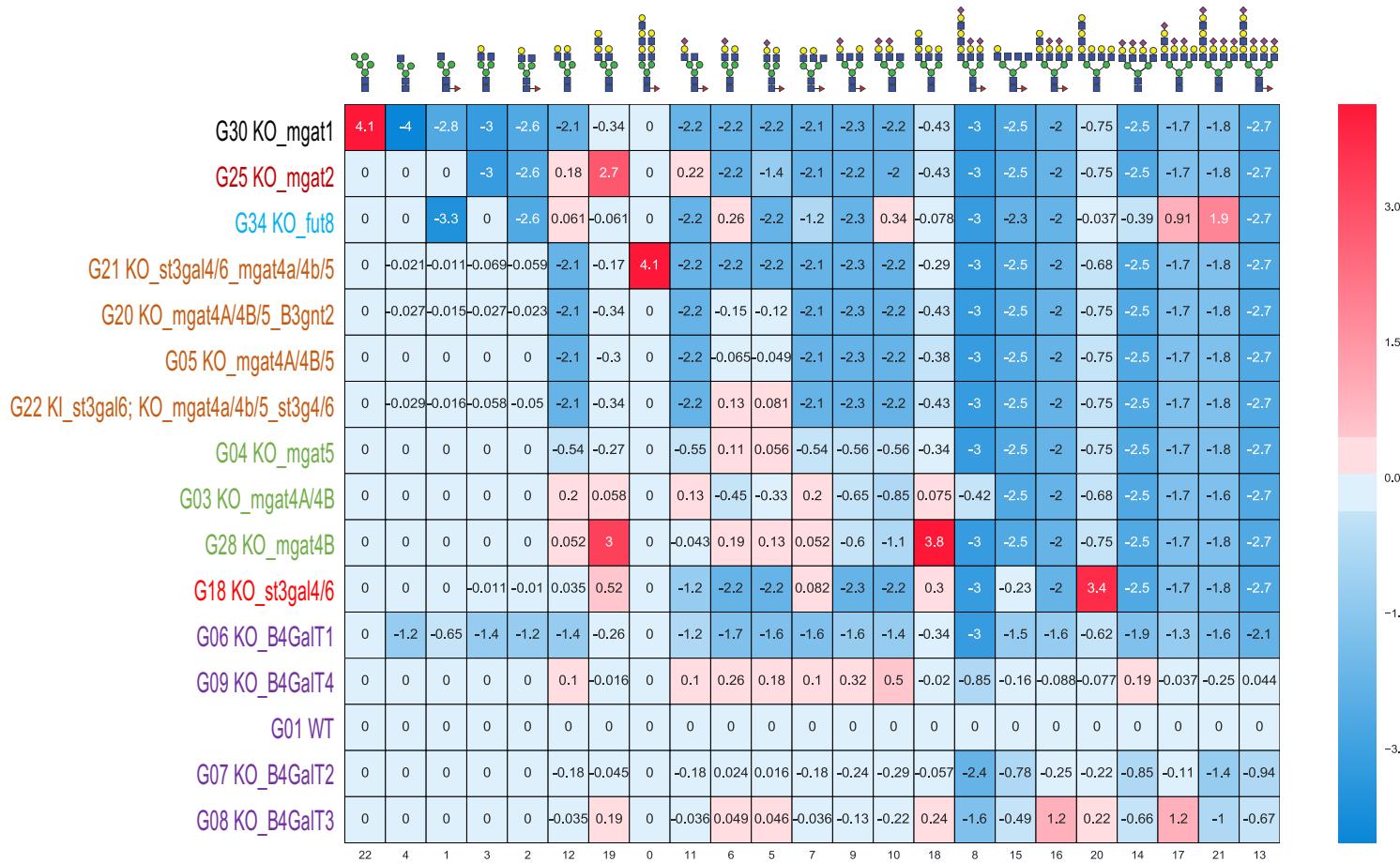


Figure 2 D: We are able to quantify the relative abundance of the representative substructure across 16 glycotypes. Each column is rescaled with z score. And adjusted with the WT abundance (the abundance of the WT is 0). The substructures from left to right are reordered from single-antennary to tetra-antennary. It shows the decreasing of the complex glycan structure when the knockout becomes more complicate, from bottom to the top. The substructure in deep red implies the significantly increase and several of them have not been noticed before.

Figure 2 EFG

- Figure2 E shows in the KO_fu8 profile, the relative abundance of structures without fucose have significant decrease. But the tetra-antennary polyLacNac elongated N-glycan with no fucose increases as well as the tri-antennary one.
- Figure2 F shows in the KO_stgal4/6 profile, the relative abundance of structures with sialylation have significant decrease. But the tetra-antennary and tri-antennary polyLacNac elongated N-glycan with no sialylation increase.
- Figure2 G shows the knockout in mgat 4B, mgat 4A/4B and mgat 5. Most of the tetra-antennary structures decrease. The Mgat 4B and mgat 5 both have significantly decrease in tri-antennary polyLacNac elongated N-glycan. While mgat 4B has significantly increase in tri-antennary LacNac elongated N-glycan.

Figure2 E

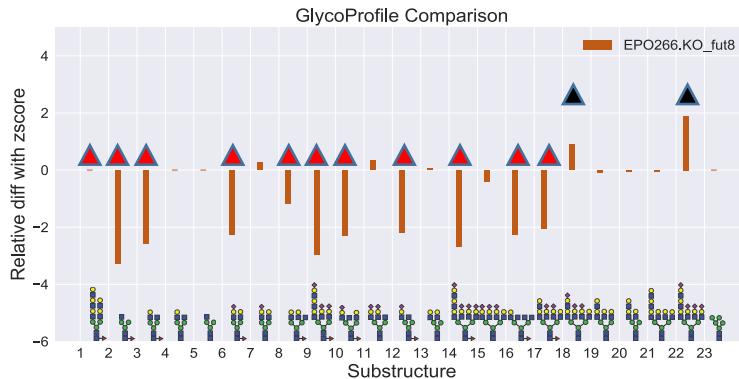


Figure2 F

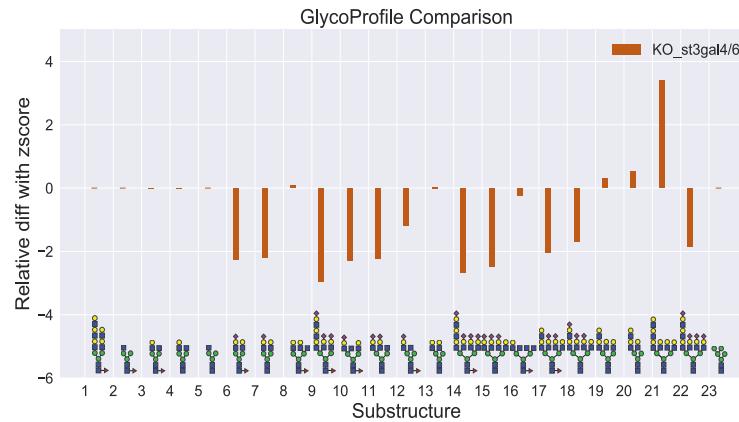
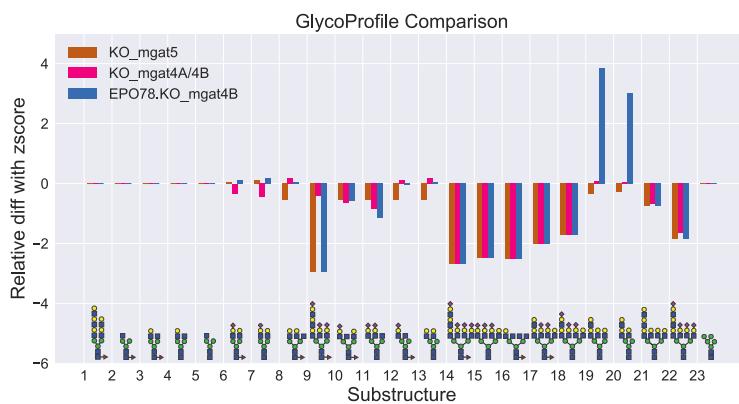


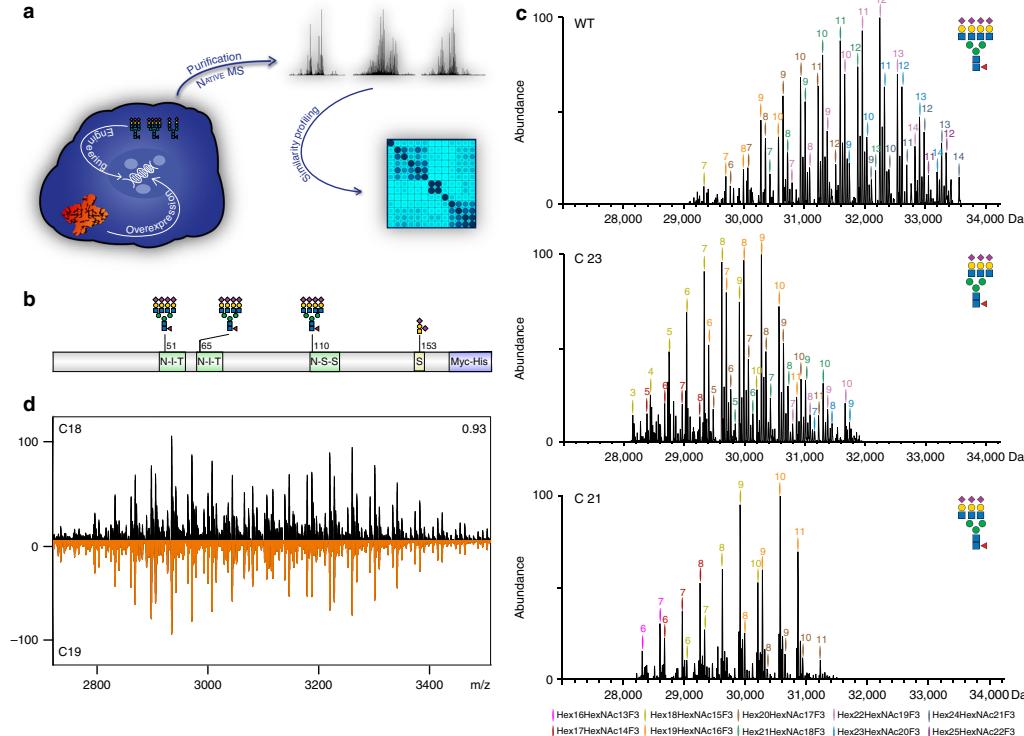
Figure2 G



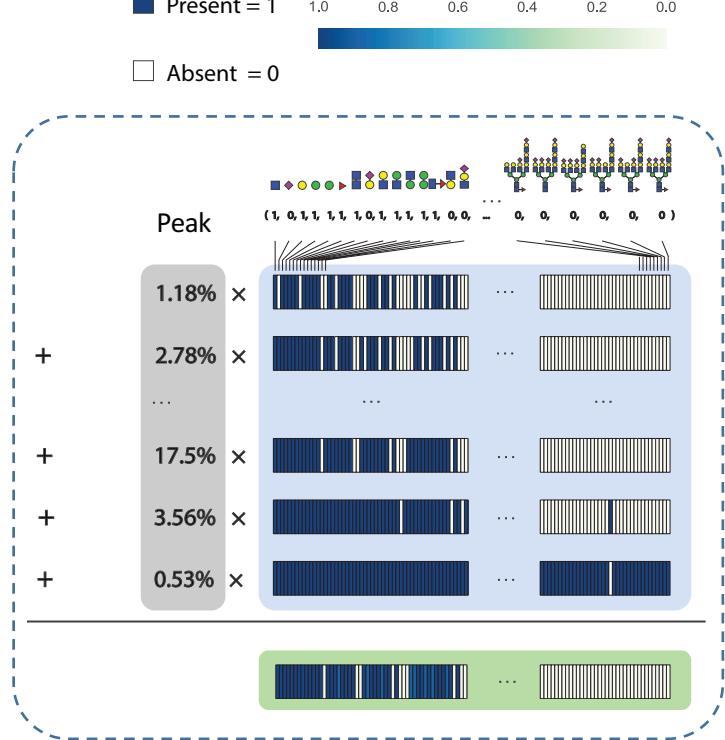
Data 1 Part 3 Inferring the glycoprofile abundance vector from native mass spectrum (ongoing)

NATURE COMMUNICATIONS | DOI: 10.1038/s41467-018-05536-3

ARTICLE

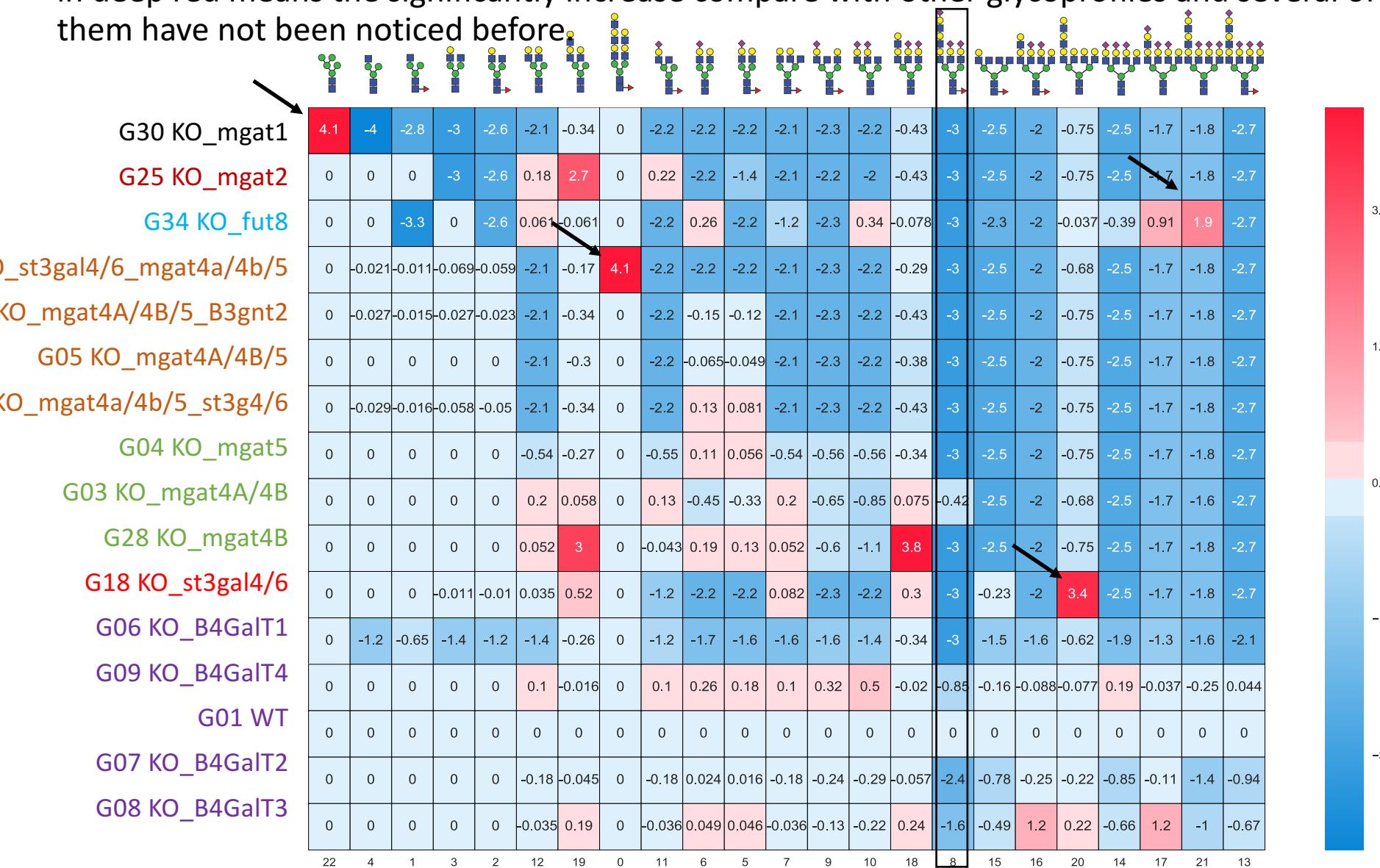


The plot on the left side comes from the **Nature Communication** paper, which shows the native mass spectrum profile for WT, KO mgat4A, KO mgat4B profile. The plot on the right shows it can be transferred to motif vector.

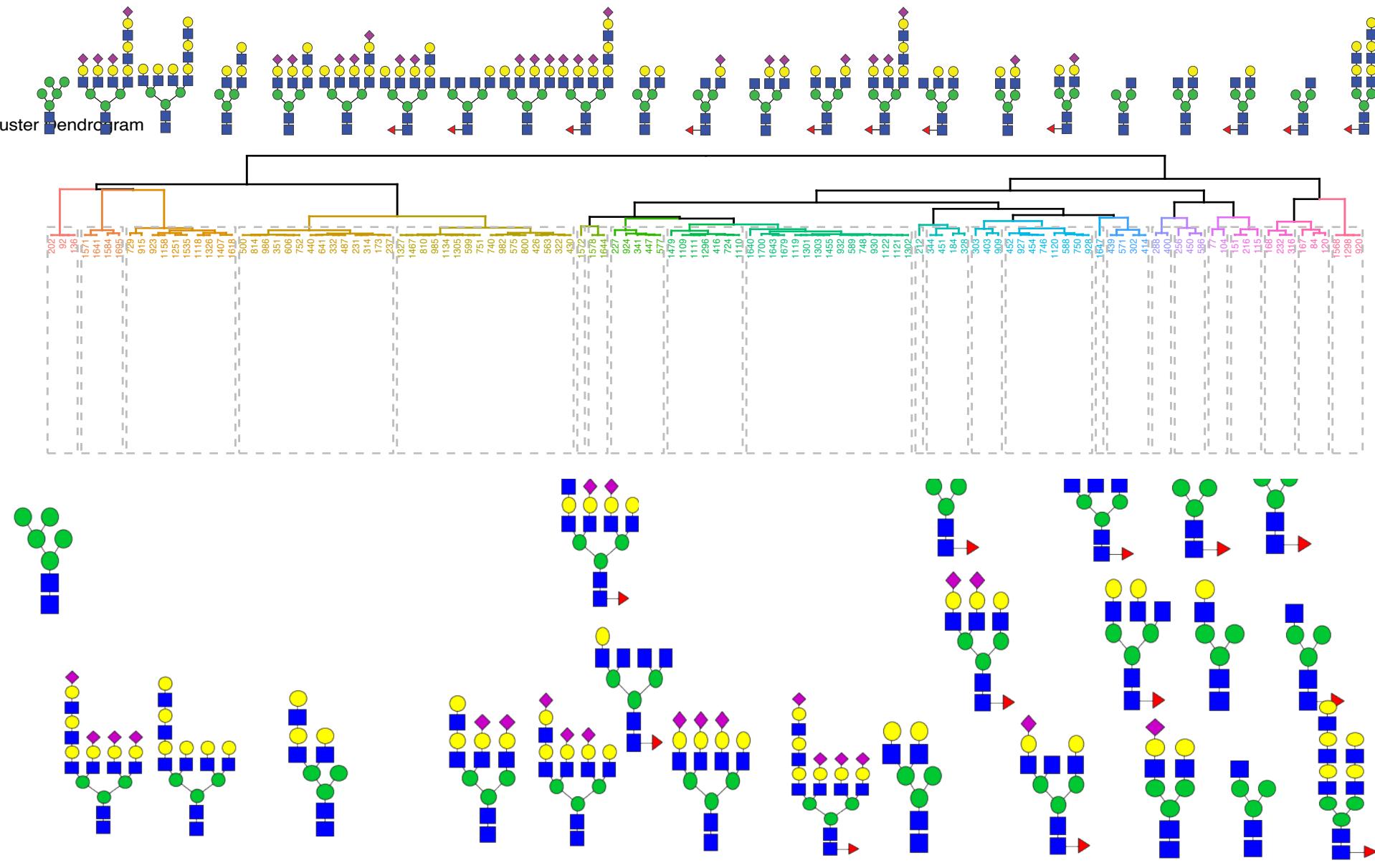


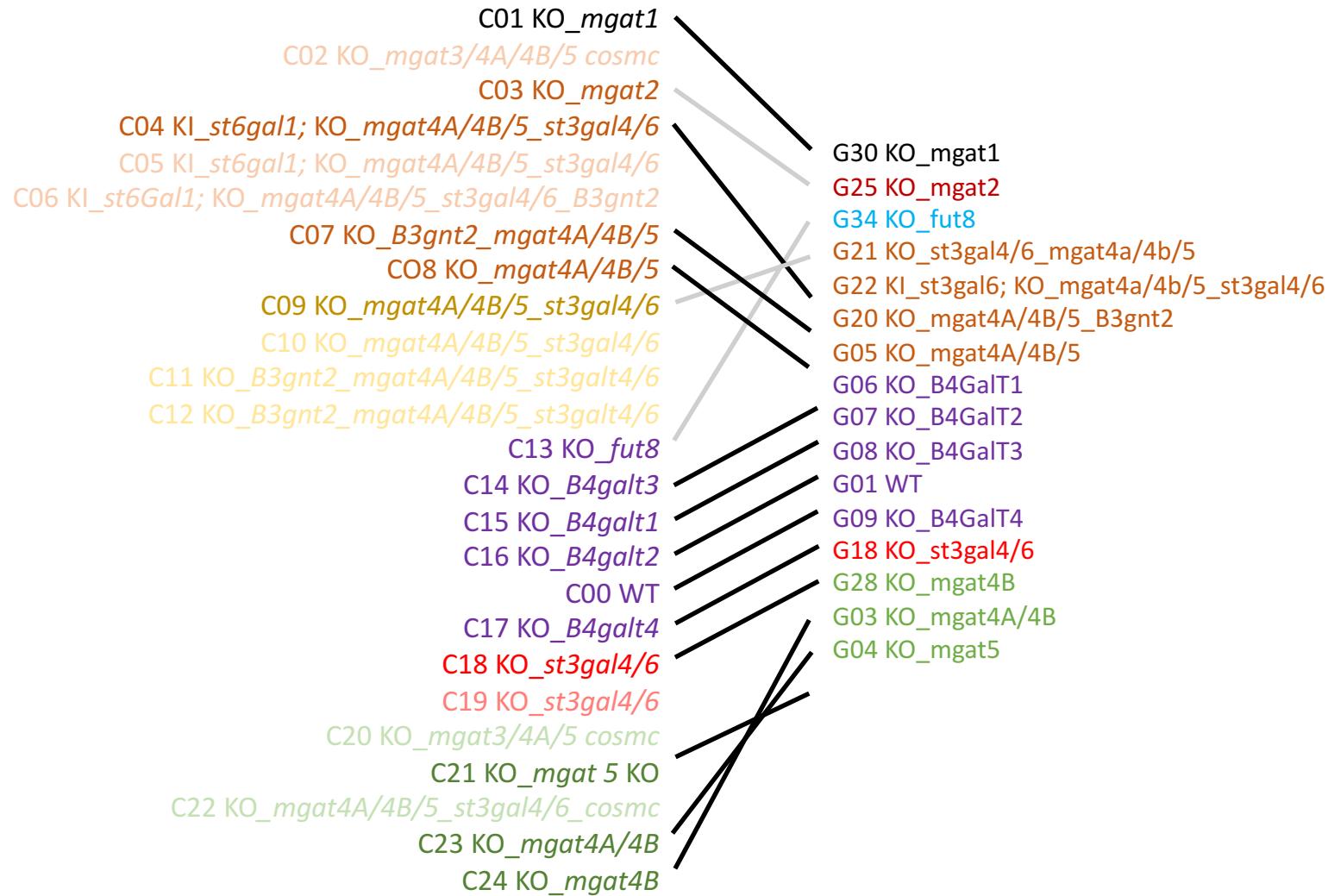
Data 1 part2

Note: each column is rescaled with z_score. And adjusted with the WT abundance. The substructures from left to right are becoming more complicate. It shows the decreasing of the complex glycan structure when the knockout becomes more complicate. Several substructure in deep red means the significantly increase compare with other glycoprofiles and several of them have not been noticed before.

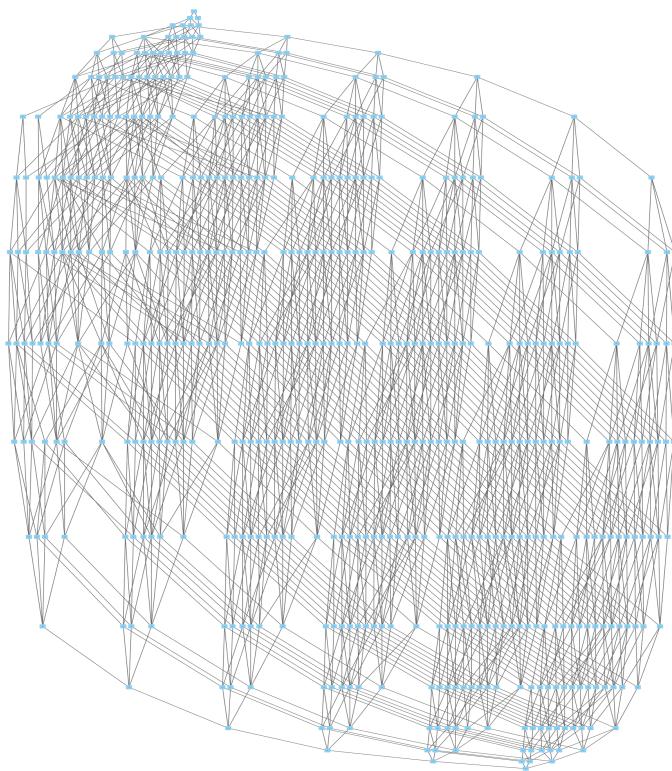
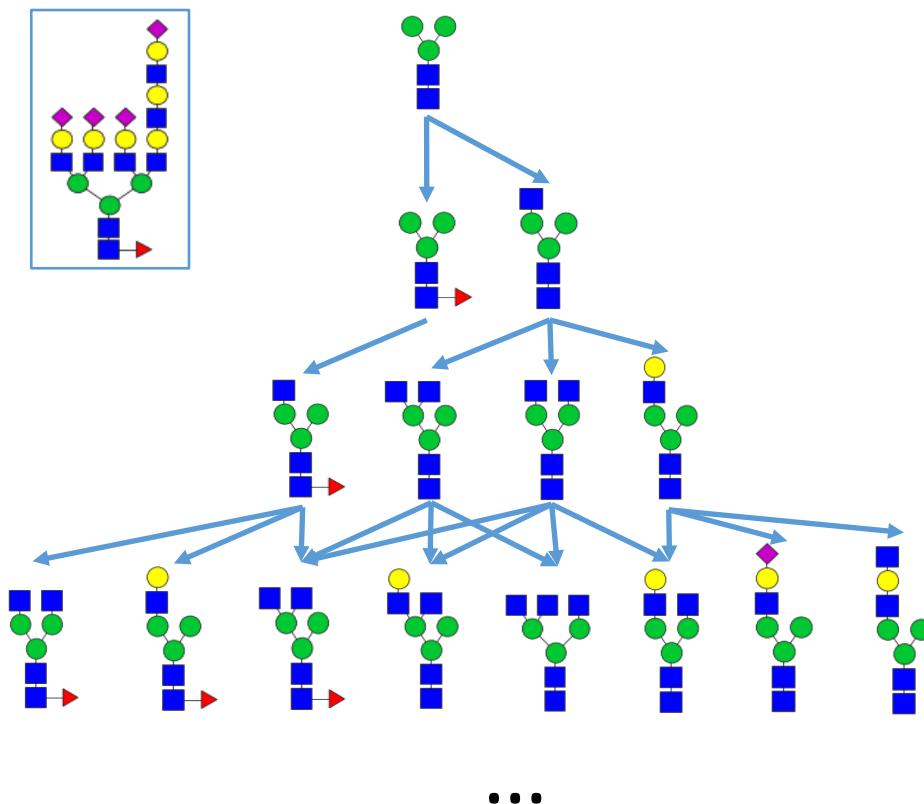


(Supplementary Figure)

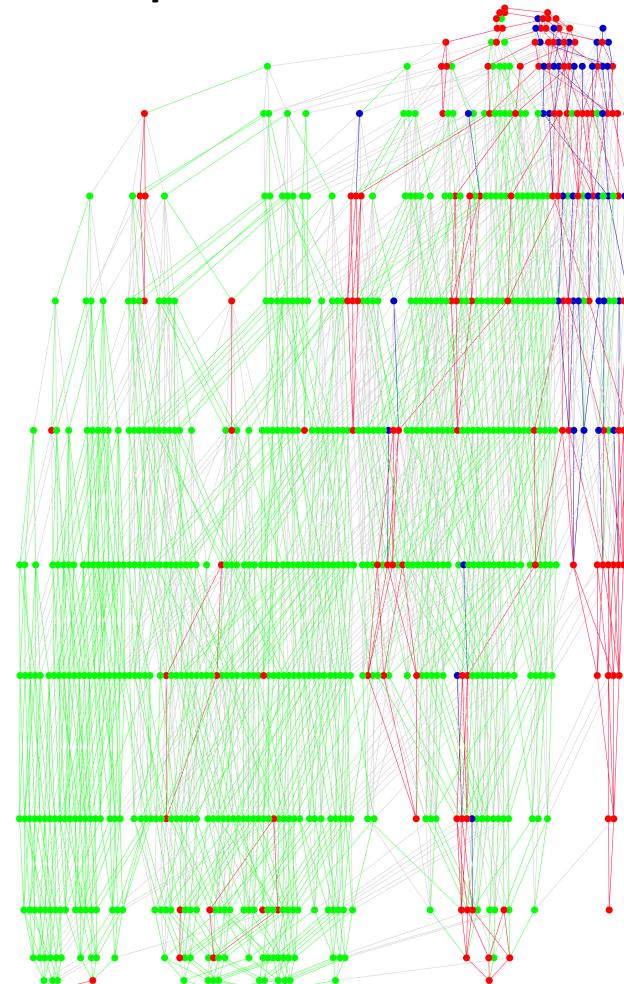
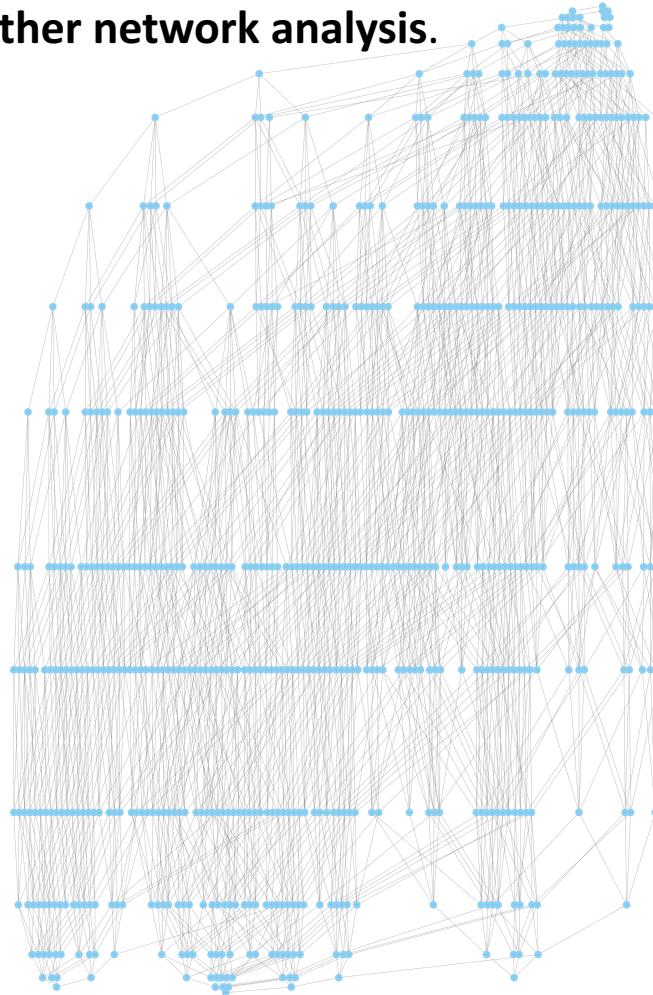




Note: a glycan can be break down to substructure and then form a synthesis network. The left side shows the top several substructures from the right side network.



Note: multiple glycans' synthesis network can be merged together and multiple glycoprofiles can also be merged together. Thus, the comparison across the glycoprofiles is feasible and the important change of substructure abundance can be highlighted (red and blue nodes). **Only the nodes with blue and red color are picked for further network analysis.**

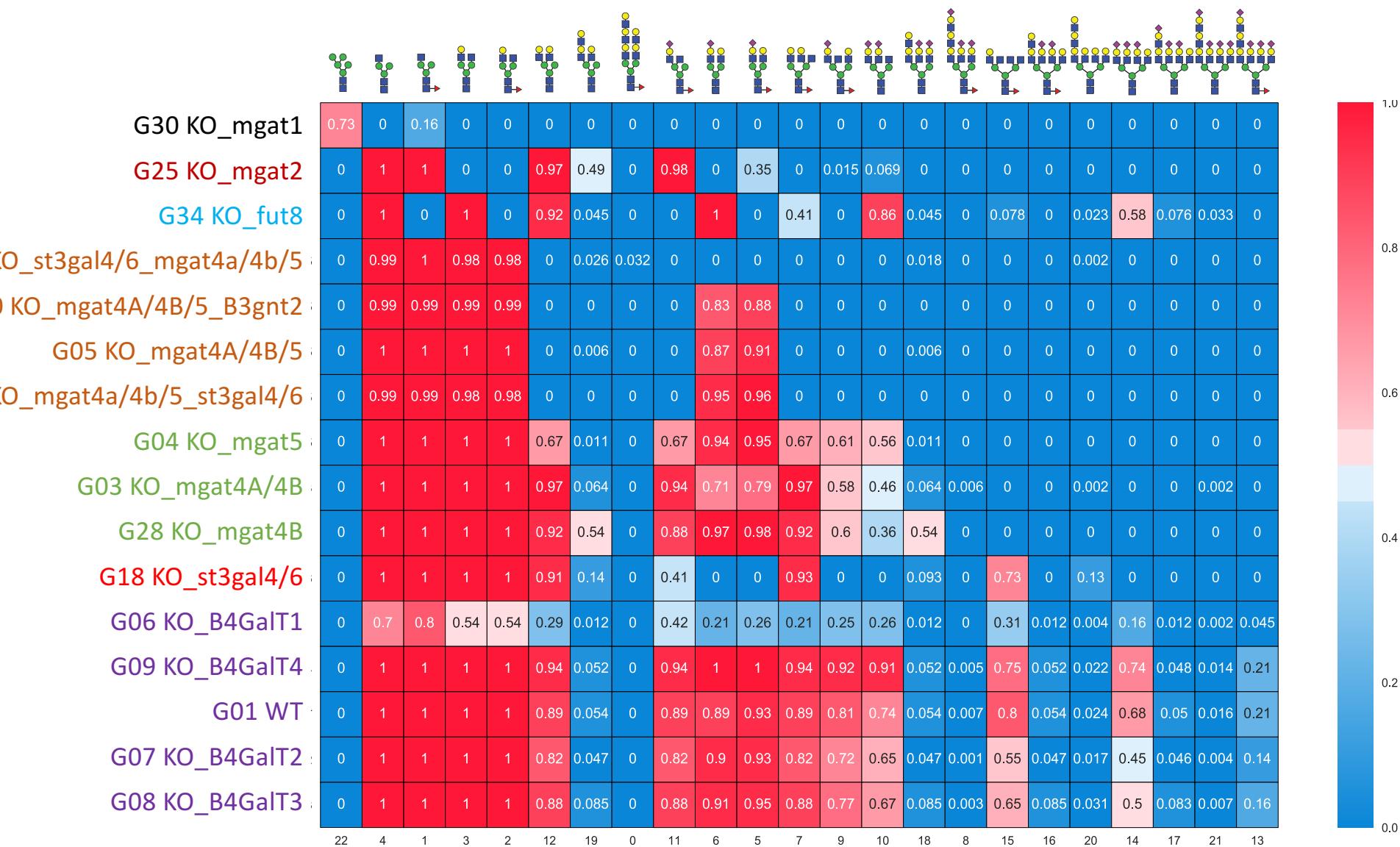


GlyCompare

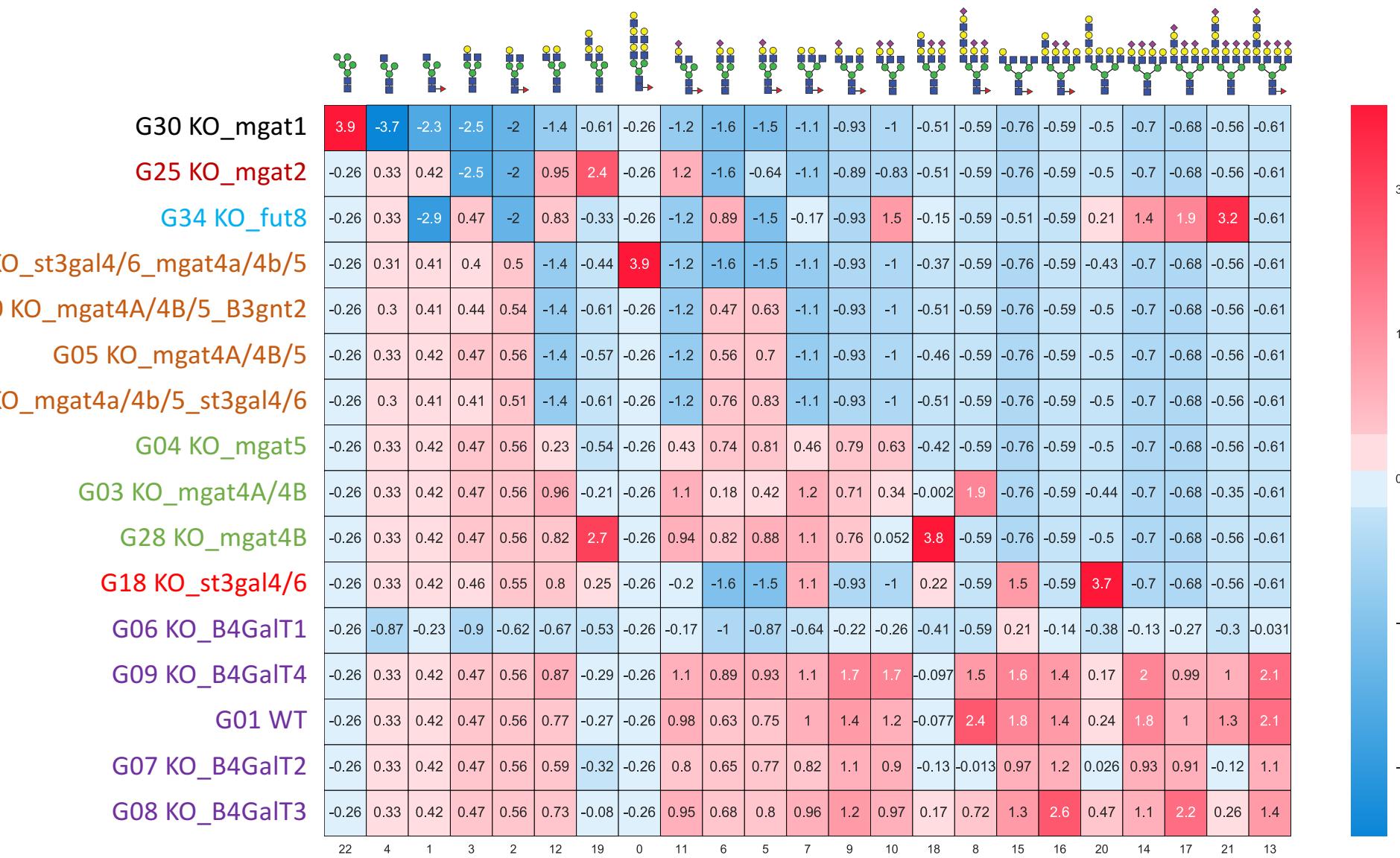
- Comparing large set of glycoprofiles
- Tracking the substructure(motif) abundance change
- Correlating the phenotypes with the motif's abundance.
- Illustrating the perturbation of the synthesis netowork.

Data 1 part2

Note: the mean of each cluster is showed here.

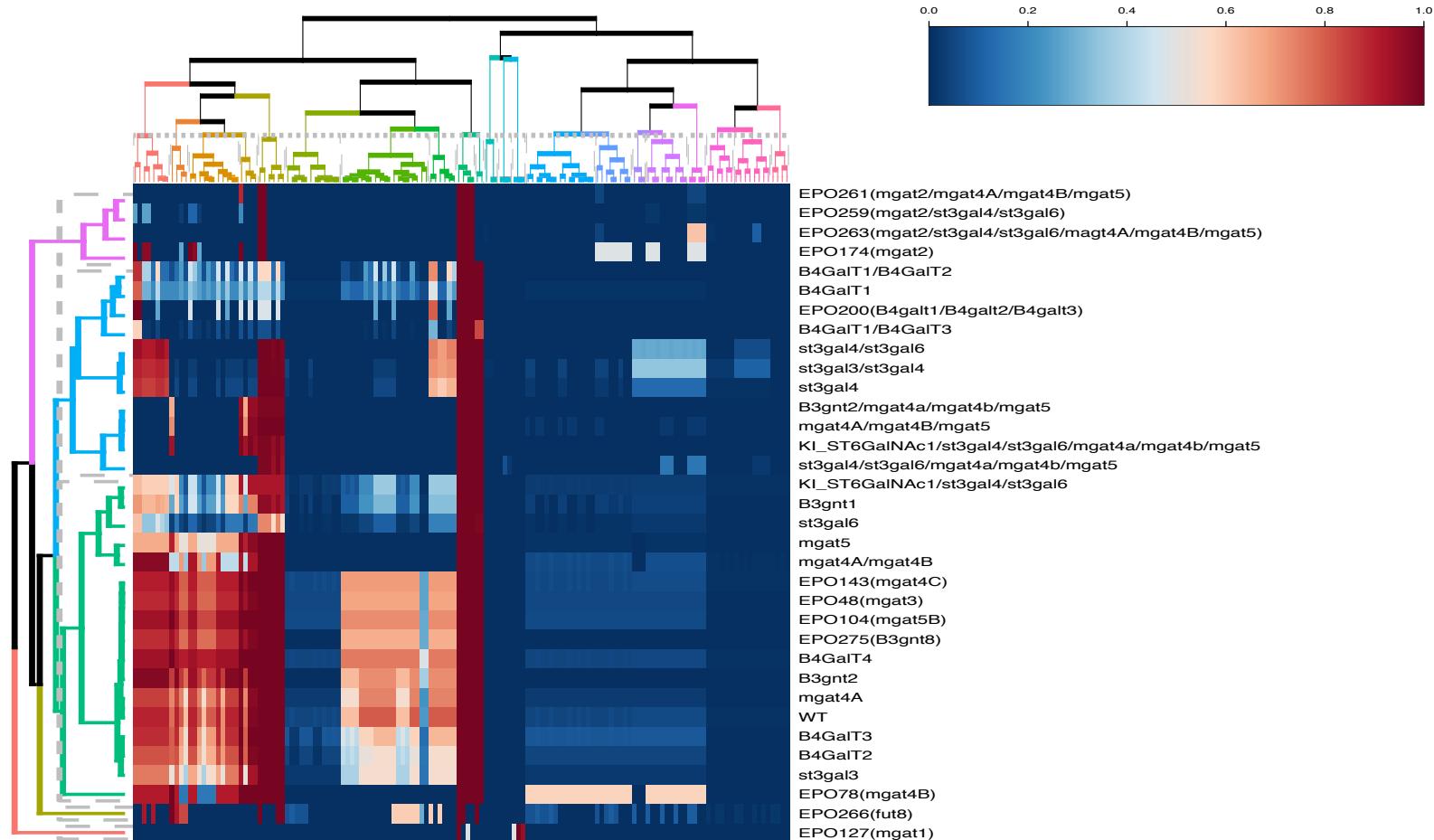


Note: each column is rescaled with z_score.

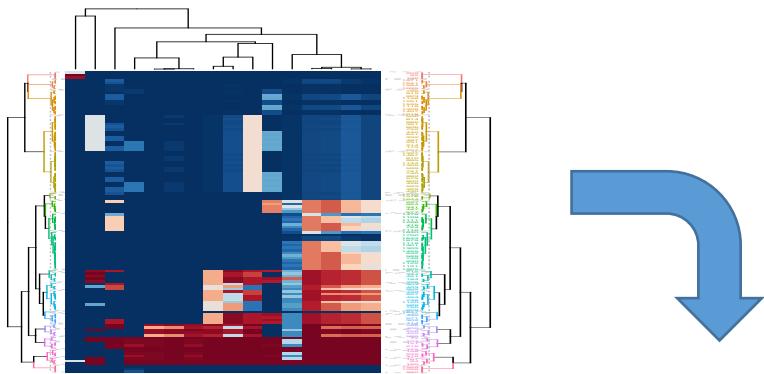


Data 1 part1

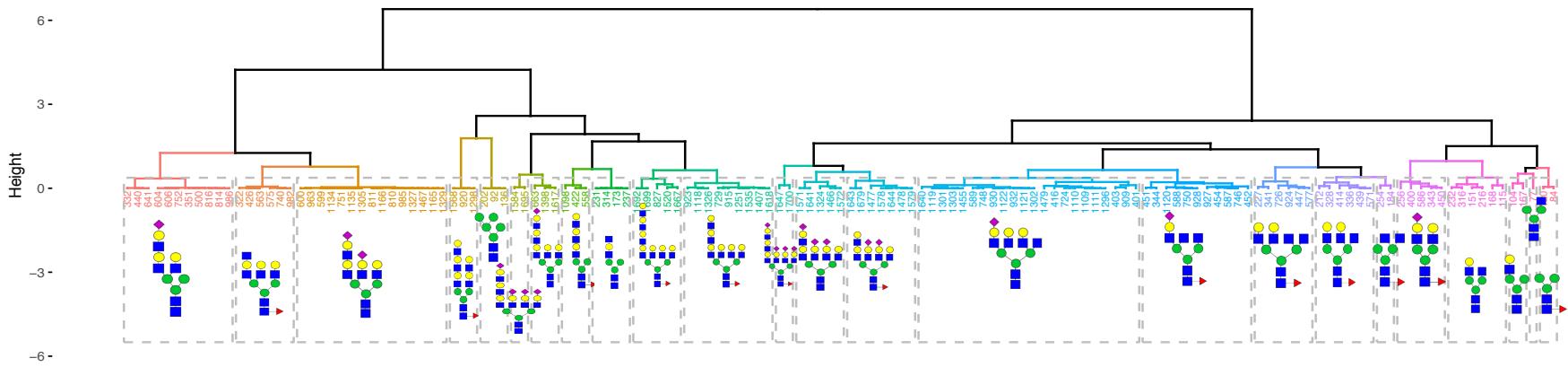
Note: preliminary analysis, clustering across 33 glycoprofiles,
based on the abundance of the substructures and glycoprofiles.



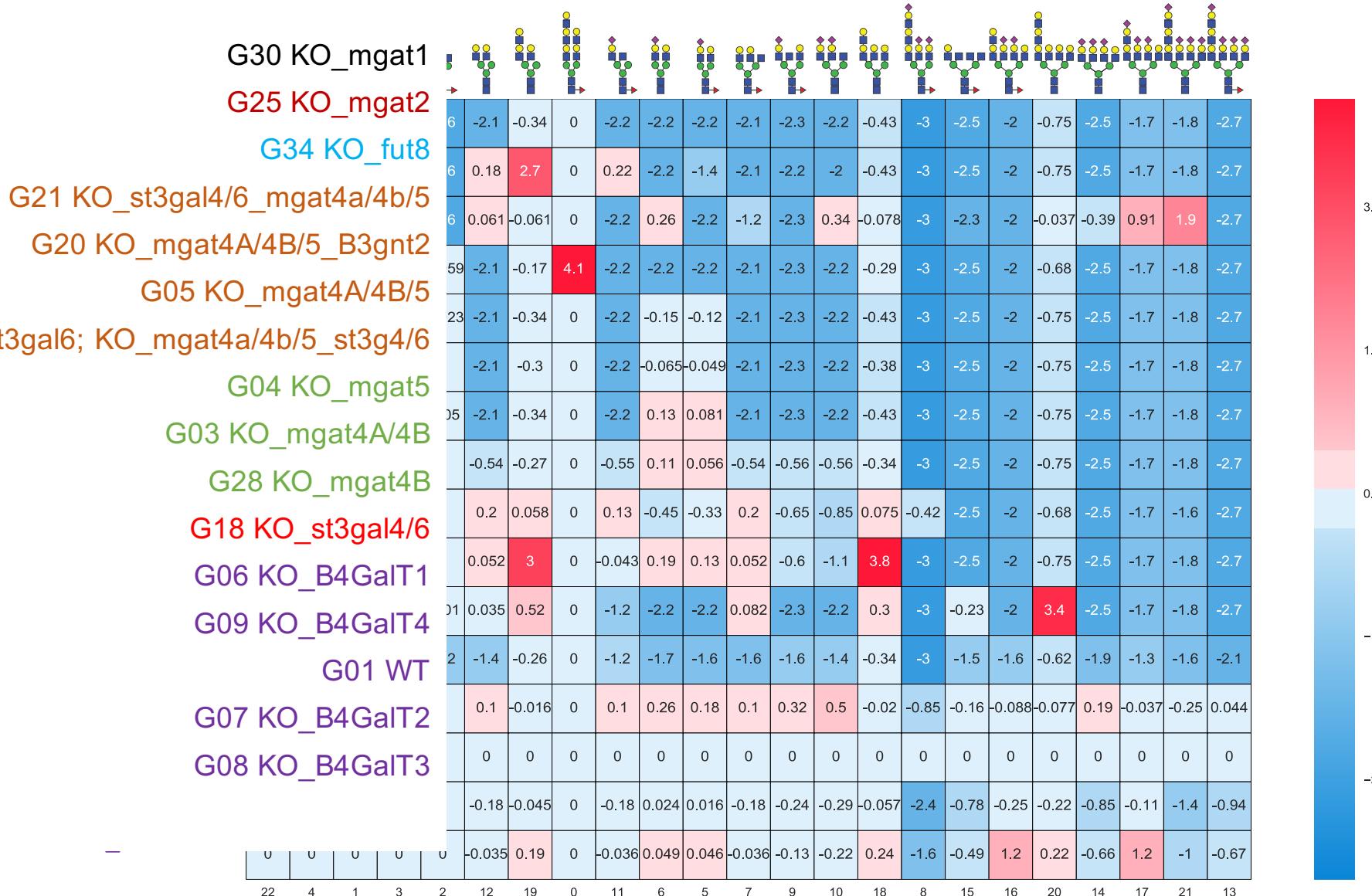
Note: the cluster of the substructure. The representative structures are generated for each cluster.



Cluster Dendrogram

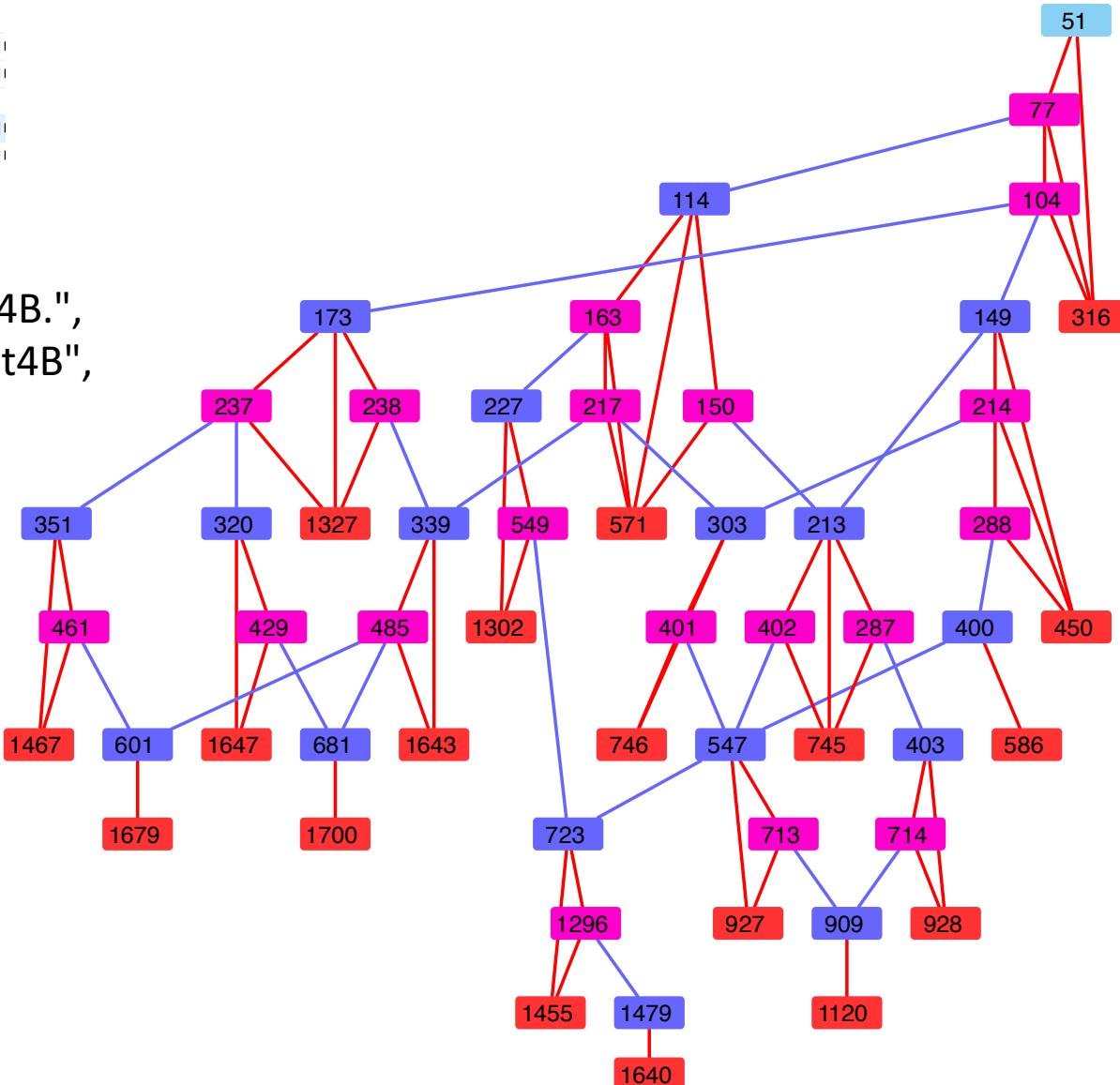


Note: each column is rescaled with z_score. And adjusted with the WT abundance. The substructures from left to right are becoming more complicate. It shows the decreasing of the complex glycan structure when the knockout becomes more complicate. Several substructure in deep red means the significantly increase compare with other glycoprofiles and several of them have not been noticed before.



immd	■
med_root	■
root	■
ttest	■
yes	■

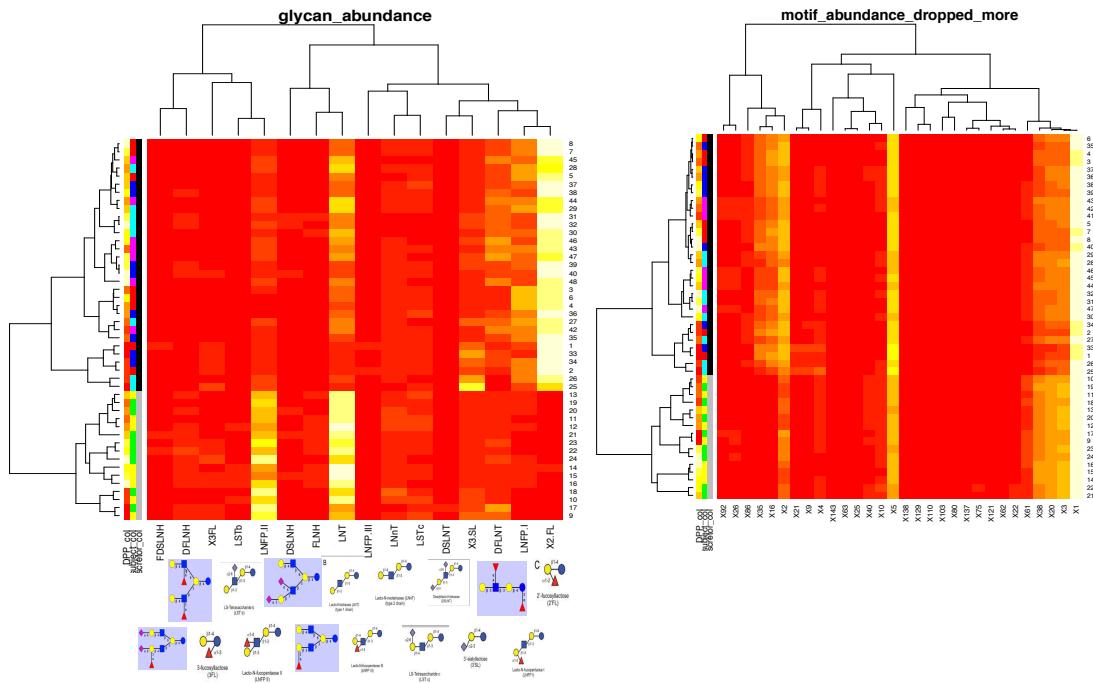
"EPO78.mgat4B.",
 "mgat4A.mgat4B",
 "WT"

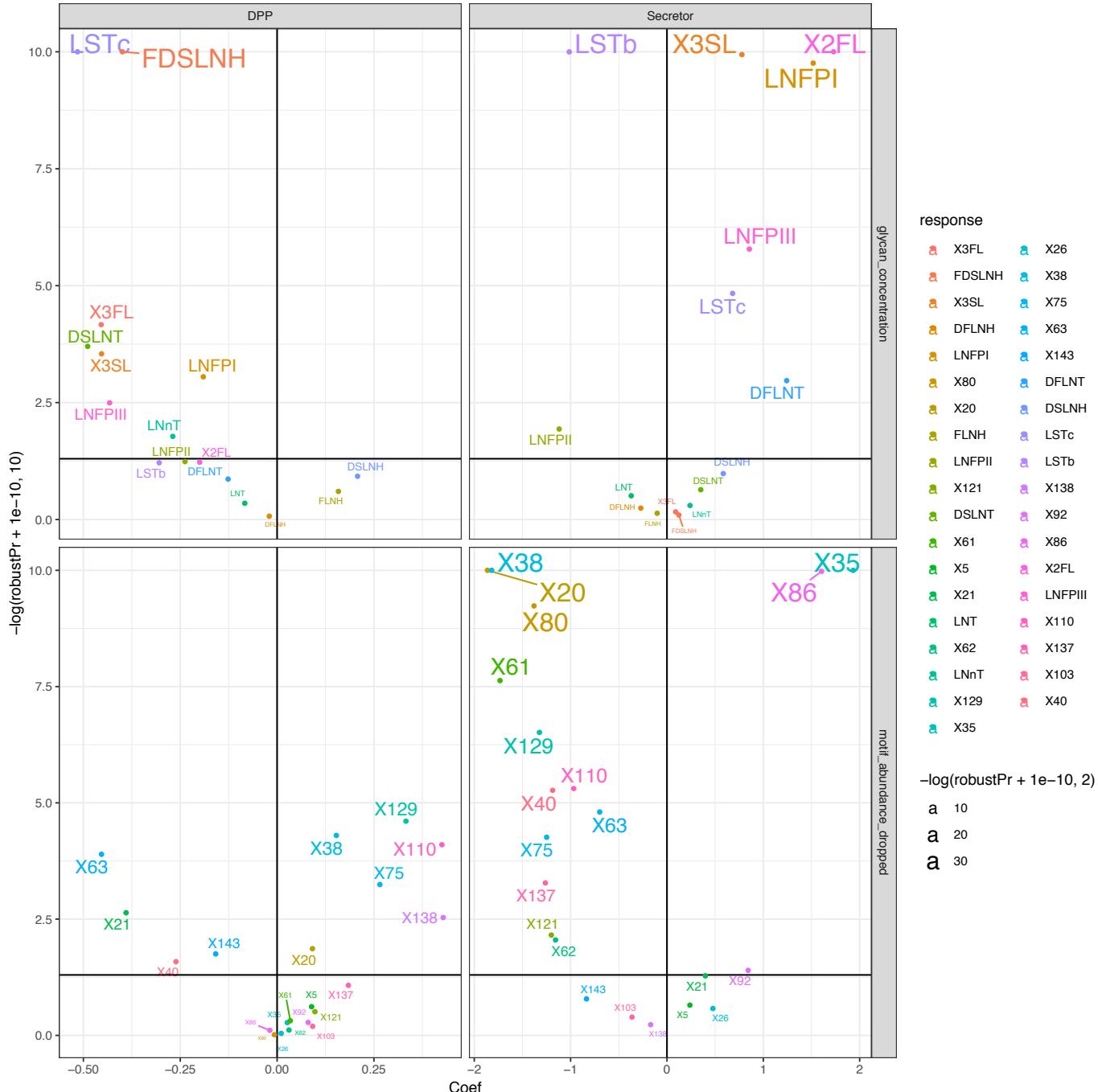


Data 2: HMO data 16 structure, 47 (48-1) glycoprofile

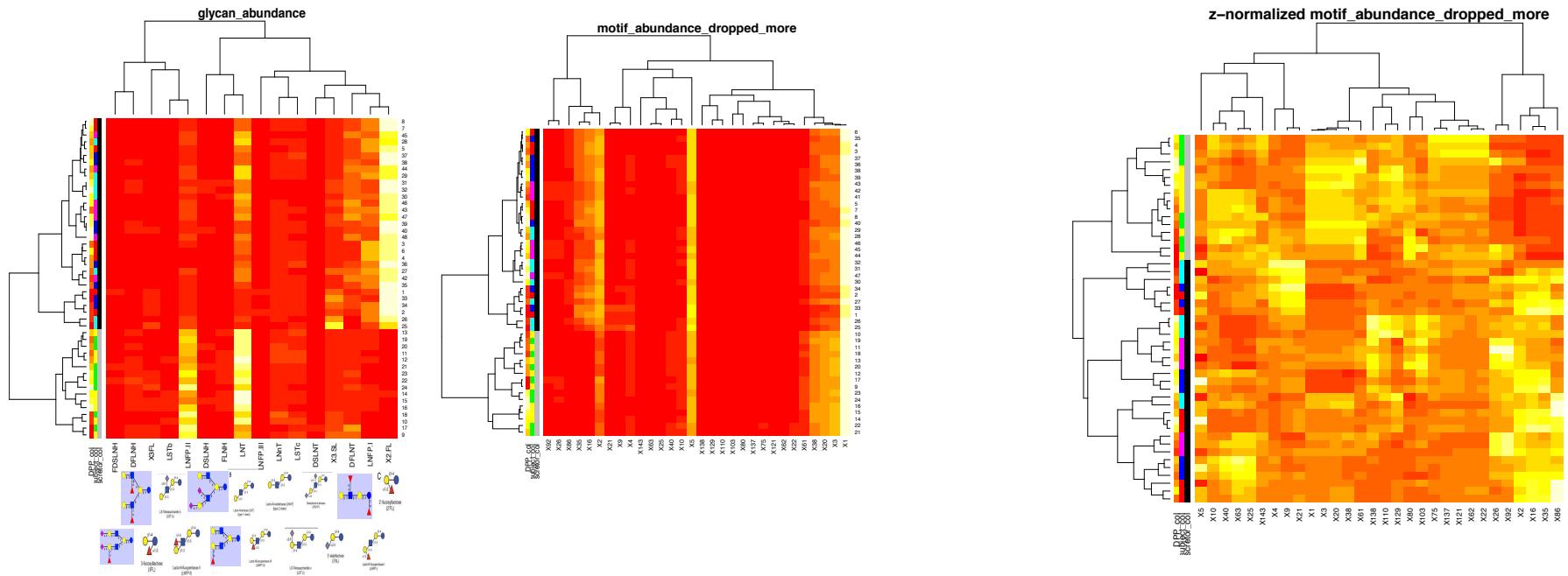
- Results
- Heatmaps
 - Motifs and glycans segregate secretor vs non-secretor
 - Second major organizing trend is time
 - Sialylation decreases over time
 - Secretors enriched & non secretors for: 26, 92, 2, 16, 35, 86
- Regression:
 - X80 is anticorrelated with secretor status (log odds, -.015, p=1e-2), X21 is positively correlated with secretor status.
 - X129 & X110 are positively associated with time while sialylated motifs (X?, X??, and X???) are negatively associated with time.
- Scatterplots
 - X80 abundance is anticorrelated with time in non-secretors and positively associated with log time in secretors.
 - This trend is not retrievable in DSLNH, DSLNH and LSTb, the X80 containing HMOs.
- Secretor/non-secretor spearman correlations
 - X?, X??... are positively correlated in secretors and negatively correlated in non-secretors
 - X?, X??... are positively correlated in non-secretors and negatively correlated in secretors
 - motifs positively correlated in secretors are (fucosylated, fuc & sialylated, galactosylated...?)
 - glycan leve observations
 - LSTb is correlated to LSTc in nonsecretors but not secretors suggesting a secretor-status dependent competition

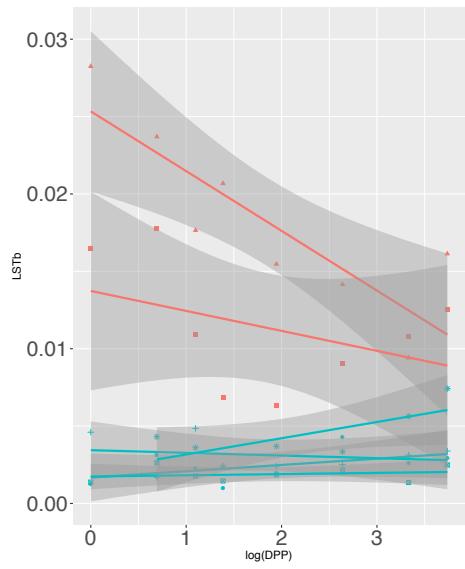
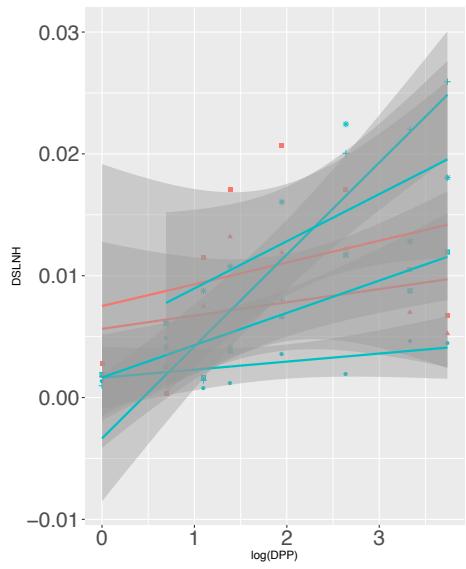
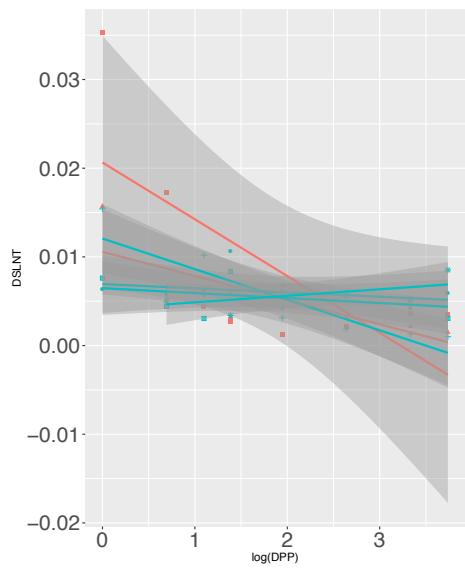
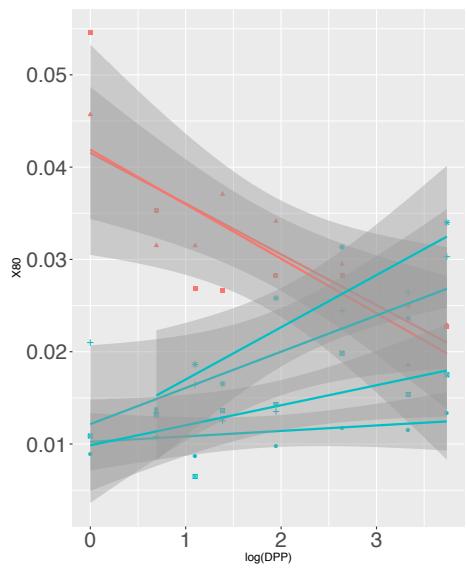
Figure 3 A





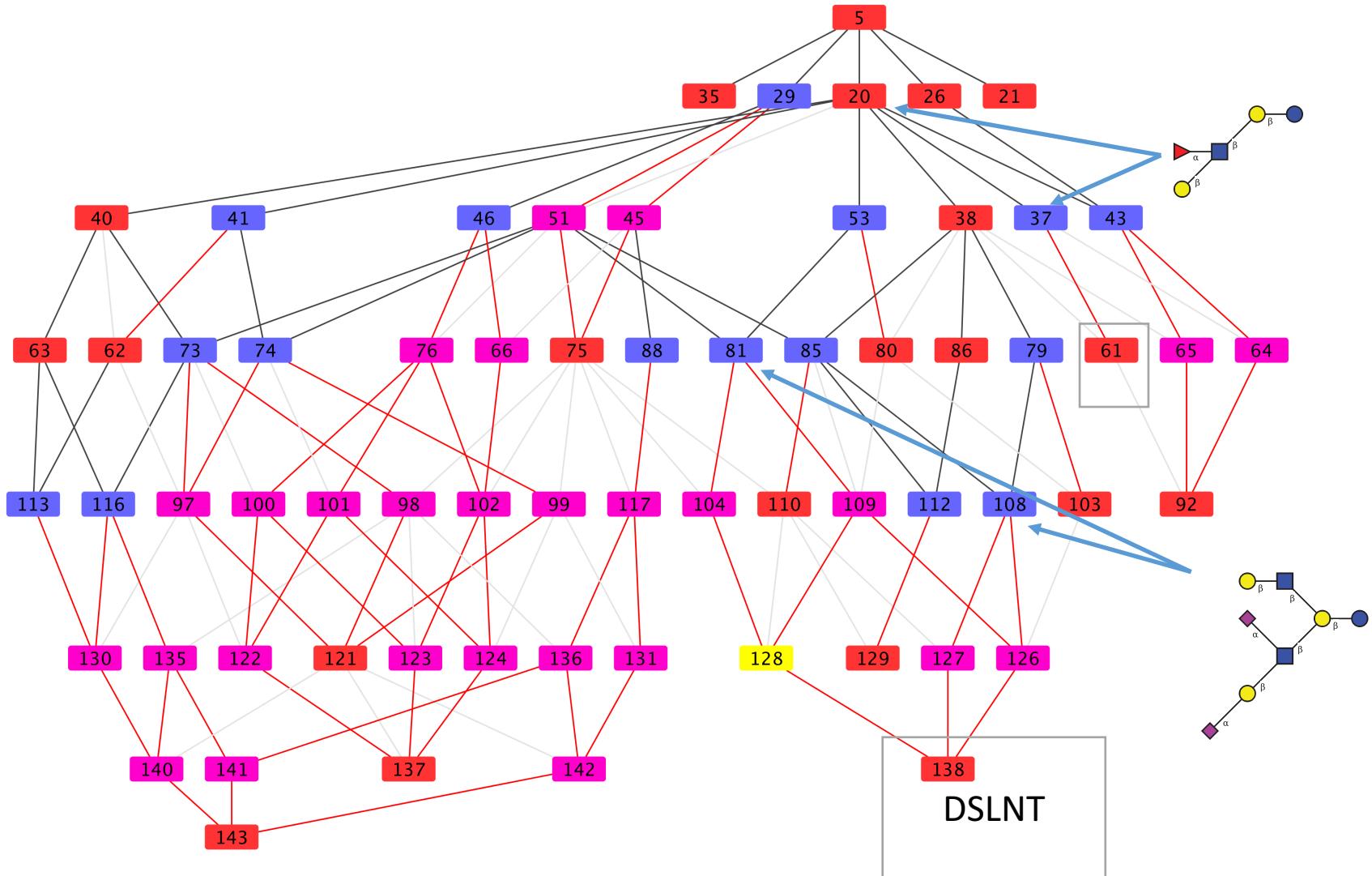
Data 2 part1: Information is more enriched in the motif level. Several motif clusters are highlighted in the plot.





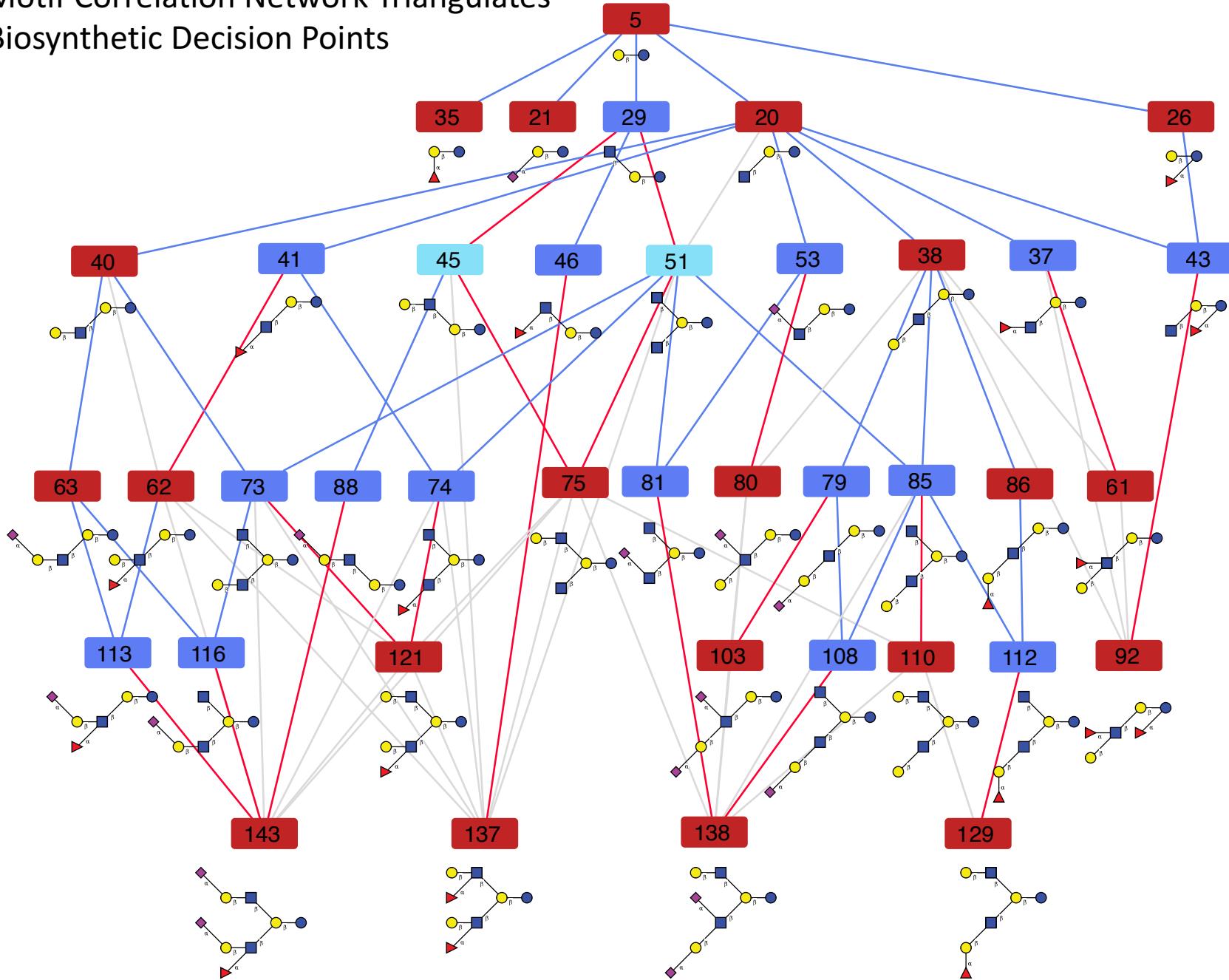
GlyCompare

- Comparing large set of glycoprofiles
- Tracking the substructure(motif) abundance change
- **Illustrating the perturbation of the synthesis network.**
- Correlating the phenotypes with the motif's abundance.



Motif Correlation Network Triangulates
Biosynthetic Decision Points

Motif Correlation Network Triangulates Biosynthetic Decision Points



GlyCompare

- Comparing large set of glycoprofiles
- Tracking the substructure(motif) abundance change
- Illustrating the perturbation of the synthesis network.
- Correlating the phenotypes with the motif's abundance.

Data 2 Part 2: Regression model are used to evaluate the importance of motifs as features to predict the Date after birth and Secretor(or not) phenotypes. Motifs that are enriched with DPP, Secretor phenotype. Sialylation is more abundant in the early date after the birth. And x80 motif which is Lst b, have different abundance between the secretor and non secretor.

Y axis is $-\log(p)$, x axis is the coefficient value.