

1 **Elucidating Human Milk Oligosaccharide biosynthetic genes through network-based multi-**
2 **omics integration**

3 Benjamin P. Kellman^{1,3,5,*}, Anne Richelle^{1,*}, Jeong-Yeh Yang⁶, Digantkumar Chapla⁶, Austin W. T.
4 Chiang^{1,2}, Julia Najera¹, Bokan Bao^{1,3,5}, Natalia Koga¹, Mahmoud A. Mohammad⁴, Anders Bech
5 Bruntse¹, Morey W. Haymond⁴, Kelley W. Moremen⁶, Lars Bode^{1,7}, Nathan E. Lewis^{1,2,5,§}

6

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

8 2. The Novo Nordisk Foundation Center for Biosustainability at the University of California, San Diego, La
9 Jolla, CA 92093, USA

10 3. Bioinformatics and Systems Biology Graduate Program, University of California, San Diego, La Jolla, CA
11 92093, USA

12 4. Department of Pediatrics, Children's Nutrition Research Center, US Department of Agriculture/Agricultural
13 Research Service, Baylor College of Medicine, Houston, Texas 77030, USA

14 5. Department of Bioengineering, University of California, San Diego, La Jolla, CA 92093, USA

15 6. Complex Carbohydrate Research Center, University of Georgia, Athens, GA, USA.

16 7. Larsson-Rosenquist Foundation Mother-Milk-Infant Center of Research Excellence (MOMI CORE),
17 University of California, San Diego, La Jolla, CA 92093, USA

18

19 * These authors contributed equally

20 §Correspondence to: Nathan E. Lewis, nlewisres@ucsd.edu

21

22 **Keywords (maximum 6):** Human milk oligosaccharide; mathematical modeling;
23 glycosyltransferases; biosynthesis; systems biology

24 ABSTRACT

25 Human Milk Oligosaccharides (HMOs) are abundant carbohydrates fundamental to infant health and
26 development. Although these oligosaccharides were discovered more than half a century ago, their
27 biosynthesis in the mammary gland remains largely uncharacterized. Here, we used a systems
28 biology framework that integrated glycan and RNA expression data to construct an HMO biosynthetic
29 network and predict glycosyltransferases involved. To accomplish this, we constructed models
30 describing the most likely pathways for the synthesis of the oligosaccharides accounting for >95% of
31 the HMO content in human milk. Through our models, we propose candidate genes for elongation,
32 branching, fucosylation, and sialylation of HMOs. We further explored selected enzyme activities
33 through kinetic assay and their co-regulation through transcription factor analysis. These results
34 provide the molecular basis of HMO biosynthesis necessary to guide progress in HMO research and
35 application with the ultimate goal of understanding and improving infant health and development.

36 SIGNIFICANCE STATEMENT

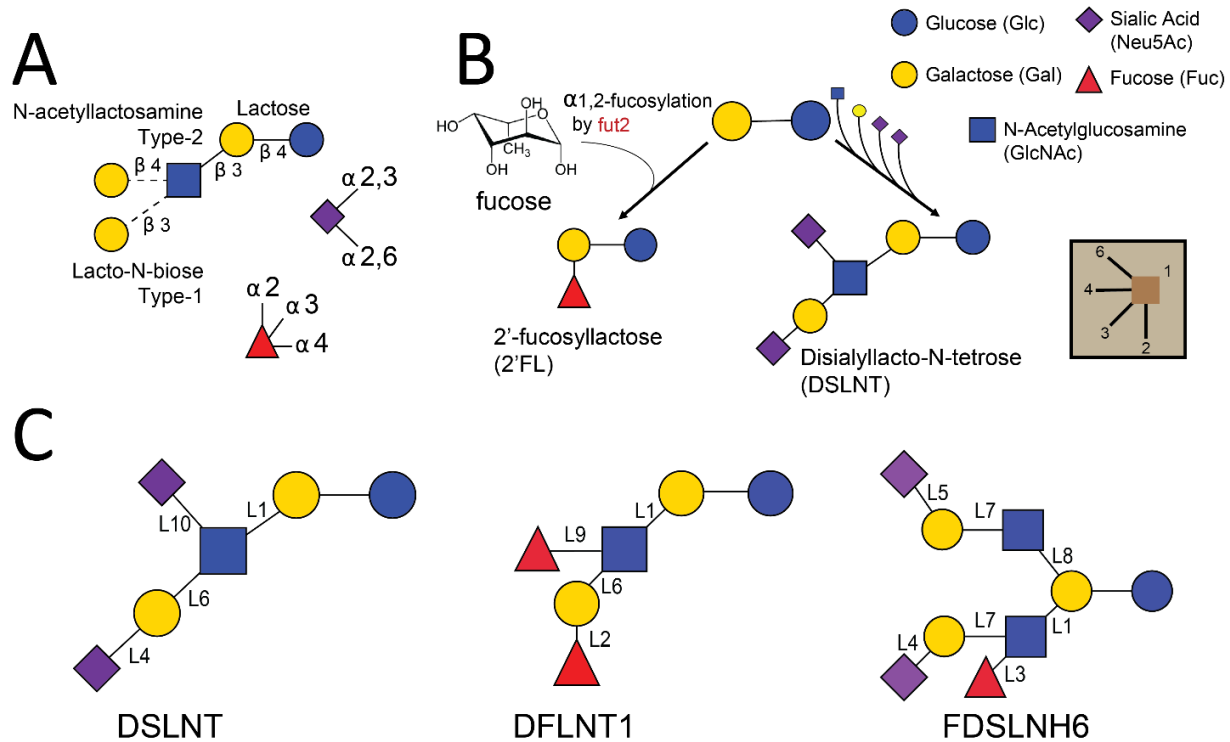
37 With the HMO biosynthesis network resolved, we can begin to connect genotypes with milk types
38 and thereby connect clinical infant, child and even adult outcomes to specific HMOs and HMO
39 modifications. Knowledge of these pathways can simplify the work of synthetic reproduction of these
40 HMOs providing a roadmap for improving infant, child, and overall human health with the specific
41 application of a newly limitless source of nutraceuticals for infants and people of all ages.

42

43 1 INTRODUCTION

44 Human milk is the “gold standard” of nutrition during early life ¹⁻³. Beyond lactose, lipids, and
45 proteins, human milk contains 11-17% (dry weight) oligosaccharides (Human Milk Oligosaccharides,
46 HMOs)^{4,5}. HMOs are milk bioactives known to improve infant immediate and long-term health and
47 development^{2,6}. HMOs are metabolic substrates for specific beneficial bacteria (e.g., *Lactobacillus* spp.
48 and *Bifidobacter* spp.), and shape the infant’s gut microbiome ^{2,7}. HMOs also impact the infant’s
49 immune system, protect the infant from intestinal and immunological disorders (e.g., necrotizing
50 enterocolitis, HIV, etc.), and may aid in proper brain development and cognition ^{2,6,8,9}. In addition,
51 recent discoveries show that some HMOs can be beneficial to humans of all ages, e.g. the HMO 2’-
52 fucosyllactose (2’FL) protecting against alcohol-induced liver disease¹⁰.

53 The biological functions of HMOs are determined by their structures ⁶. HMOs are unconjugated
54 glycans consisting of 3–20 total monosaccharides draw from 3-5 unique monosaccharides: galactose
55 (Gal, A), glucose (Glc, G), N-acetylglucosamine (GlcNAc, GN), fucose (Fuc, F) and the sialic acid N-
56 acetyl-neuraminic acid (NeuAc, NN) (**Figure 1A**). All HMOs extend from a common lactose (Gal β 1-
57 4Glc) core. The core lactose can be extended at the nonreducing end, with a β -1,3-GlcNAc to form a
58 trisaccharide. That intermediate trisaccharide is quickly extended on its non-reducing terminus with
59 a β -1,3-linked galactose to form a type-I tetrasaccharide (LNT) or a β -1,4-linked galactose to form a
60 type-II tetrasaccharide (LNnT). Additional branching of the trisaccharide or tetrasaccharide typically
61 occurs at the lactose core by addition of a β -1,6-linked GlcNAc to the Gal residue. Lactose or the
62 elongated oligosaccharides can be further fucosylated in an α -1,2-linkage to the terminal Gal residue,
63 or α 1,3/4-fucosylated on internal Glc or GlcNAc residues, and α -2,3-sialylated on the terminal Gal
64 residue or α -2,6-sialylated on external Gal or internal GlcNAc residues^{6,8}(**Figure 1B**).



65 **Figure 1 - HMO blueprint and synthesis (A)** HMOs are built from a combination of the five
66 monosaccharides D-glucose (Glc, blue circle), D-galactose (Gal, yellow circle), N-acetyl-glucosamine
67 (GlcNAc, blue square), L-fucose (Fuc, red triangle), and sialic acid (N-acetyl-neuraminic acid (NeuAc),
68 purple diamond). Lactose (Gal- β -1,4-Glc) forms the reducing end and can be elongated with several
69 Lacto-N-biose or N-acetyllactosamine repeat units (Gal- β -1,3/4-GlcNAc). Lactose or the
70 polylactosamine backbone can be fucosylated with α -1,2-, α -1,3-, or α -1,4- linkages or sialylated in α -
71 2,3- or α -2,6- linkages². **(B)** Small HMOs can be fucosylated to make 2'FL while larger HMOs can be
72 synthesized by the extension of the core lactose with N-acetyllactosamine (type-I) or lacto-N-biose
73 (type-II) and subsequent decoration of the extended core with sialic acid to make more complex
74 HMOs, such as DSLNT. **(C)** Three HMOs in this study: DSLNT, isomer 1 of DFLNT, isomer 6 of FDSLNH;
75 isomer structures represent predictions from this study (see Methods, **Figure S 12**). Each
76 monosaccharide-linking glycosidic bond is labeled (L1, L2,...L10) according to the linkage reactions
77 listed in **Table 1**.

78 Despite decades of study, many details of HMO biosynthesis remain unclear. While the many possible
79 monosaccharide addition events above are known, the order of the biosynthetic steps and many of
80 the enzymes involved are not known (**Table 1**). For example, the lactose core is extended by
81 alternating actions of β -1,3-N-acetylglucosaminyltransferases (b3GnT) and β -1,4-
82 galactosaminyltransferases (b4GalT) while β -galactoside sialyltransferases (SGalT) and α -1,2-
83 fucosyltransferases (including the FUT2 'secretor' locus) are responsible for some sialylation and
84 fucosylation of a terminal galactose, respectively¹¹. However, each enzymatic activity in HMO
85 extension and branching can potentially be catalyzed by multiple isozymes in the respective gene
86 family. Direct evidence of the specific isozymes performing each reaction *in vivo* is extremely limited.

87

88 **Table 1 - Glycosylation reactions examined.** We studied here several candidate glycosyltransferases
 89 expressed in our samples to identify candidates for 10 elementary reactions (see Methods, **Table S 1**).
 90 Acceptor, product and constraint are represented in LiCoRR¹²: monosaccharides include Gal (A), Fuc (F),
 91 Glc (G), GlcNAc (GN), Neu5Ac (NN). Additionally, “)” and “(” indicate initiation and termination of a
 92 branch respectively, “[X/Y]” indicates either monosaccharide, and “~” indicates a negation. An asterisk
 93 “*” indicates an imperfect match between the EC number and reaction. Background colors correspond to
 94 the monosaccharide added: GlcNAc (blue), Fuc (red), Neu5Ac (purple), and Gal (yellow).

Linkage	Reaction	EC Identifier	Acceptor {Constraint}	Product	Candidates
L1:b3GnT	b-1,3 N-acetylglucosamine	2.4.1.149	(A	(GNb3A	B3GNT2-6,8-9
L2:a2FucT	a-1,2 fucosyltransferase	(2.4.1.69,344)	(A	(Fa2A	FUT1-2
L3:a3FucT	a-1,3 fucosyltransferase	(2.4.1.152)	G/GN {~Ab3GN}	Fa3G/GN	FUT3-7,9-11
L4:ST3GalT	(b-Gal) a-2,3 sialyltransferase	(2.4.99.4)	(A	(NNA3A	ST3GAL1-6
L5:ST6GalT	(b-Gal) a-2,6 sialyltransferase	2.4.99.1	(A	(NNA6A	ST6GAL1-2
L6:b3GalT	b-1,3 galactotransferase	2.4.1.86	(GN	(Ab3GN	B3GALT1-2,4-5
L7:b4GalT	b-1,4 galactotransferase	2.4.1.90	(GN	(Ab4GN	B4GALT1-6
L8:b6GnT	b-1,6 N-acetylglucosamine	(2.4.1.150)	GNb3Ab4G	GNb3(GNb6)Ab4G	GCNT1-4,7
L9:a4FucT	a-1,4 fucosyltransferase	2.4.1.65	Ab3GNb3A {~GNb4Ab3GNb3A}	Ab3(Fa4)GNb3A	FUT3,5
L10:ST6GnT	(b-1,3-GlcNAc) a-2,6 sialyltransferase	(2.4.99.3,7)	Ab3GNb3A	Ab3(NNA6)GNb3A	ST6GALNAC1-6

95

96 Here we leverage the heterogeneity in HMO composition and gene expression across human subjects
 97 to refine our knowledge of the HMO biosynthetic network. Milk samples were collected from 11
 98 lactating women across two independent cohorts between the 1st and 42nd day post-partum. (see
 99 methods). Gene expression profiling of mammary epithelial cells was obtained from mRNA present
 100 in the milk fat globule membrane interspace. Absolute and relative concentrations of the 16 most
 101 abundant HMOs was measured. Starting from a scaffold of all possible reactions¹³⁻¹⁸, we used
 102 constraint-based modeling^{19,20} to reduce the network to a set of relevant reactions and most plausible

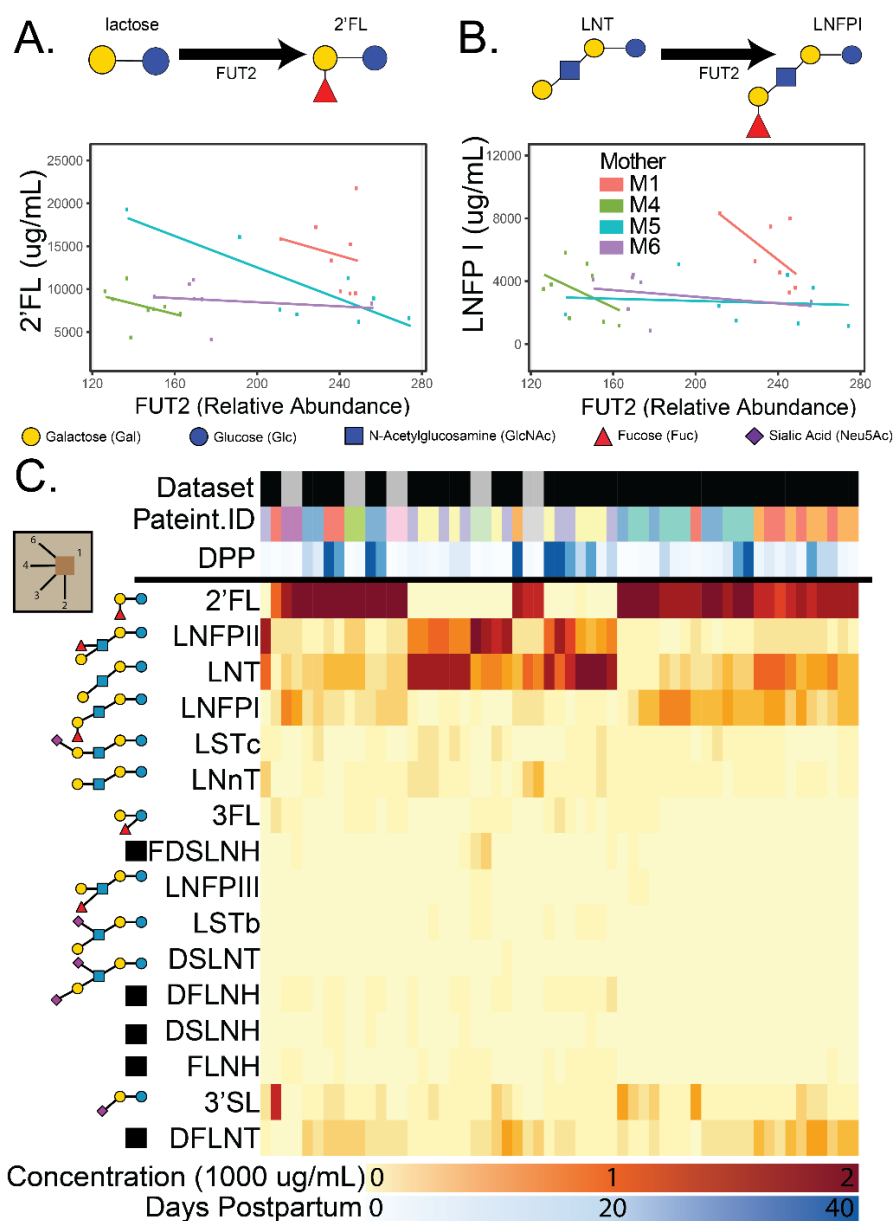
103 HMO structures when not known²¹ to form the basis for a mechanistic model^{13,14,22}. This resulted in a
104 ranked ensemble of candidate biosynthetic pathway topologies. We then ranked 44 million candidate
105 biosynthesis networks to identify the most likely network topologies and candidate enzymes for each
106 reaction by integrating sample-matched transcriptomic and glycoprofiling data from the 11 subjects.
107 For this we simulated all reaction fluxes and tested the consistency between changes in flux and gene
108 expression to determine the most probable gene isoforms responsible for each linkage type. We
109 followed with direct observations through fluorescence activity assays to confirm our predictions.
110 Finally, we performed transcription factor analysis to delineate regulators of the system. The
111 resulting knowledge of the biosynthetic network can guide efforts to unravel the genetic basis of
112 variations in HMO composition across subjects, populations, and disorders using systems biology
113 modeling techniques.

114 2 RESULTS

115 2.1 ABUNDANCES OF HMOs AND THEIR KNOWN ENZYMES DO NOT CORRELATE

116 While α -1,2-fucosylation of glycans in humans can be accomplished by both FUT1 and FUT2, only
117 FUT2 is expressed in mammary gland epithelial cells (**Table S 1**). FUT2, the “secretor” gene, is
118 essential to ABH antigens^{23–25} as well as HMO ^{2,26,27} expression. We confirmed that non-functional
119 FUT2 in “non-secretor” subjects guarantees the near-absence of α -1,2-fucosylated HMOs like 2’FL
120 and LNFP1 (Fig2C). But, examining only subjects with functional FUT2 (Secretors), we found FUT2
121 expression levels and the concentration (nmol/ml) of HMOs containing α -1,2-fucosylation do not
122 correlate in sample-matched microarray and glycomic measurements by HPLC (**Figure 2**).
123 Generalized Estimating Equations (GEE) showed no significant positive association (2’FL Wald p =
124 0.056; LNFP1 Wald p = 0.34). FUT1 could catalyze this reaction but its expression was not detected
125 in these samples. We hypothesized that to successfully connect gene expression to HMO synthesis,
126 one must account for all biosynthetic steps and not solely rely on direct correlations.

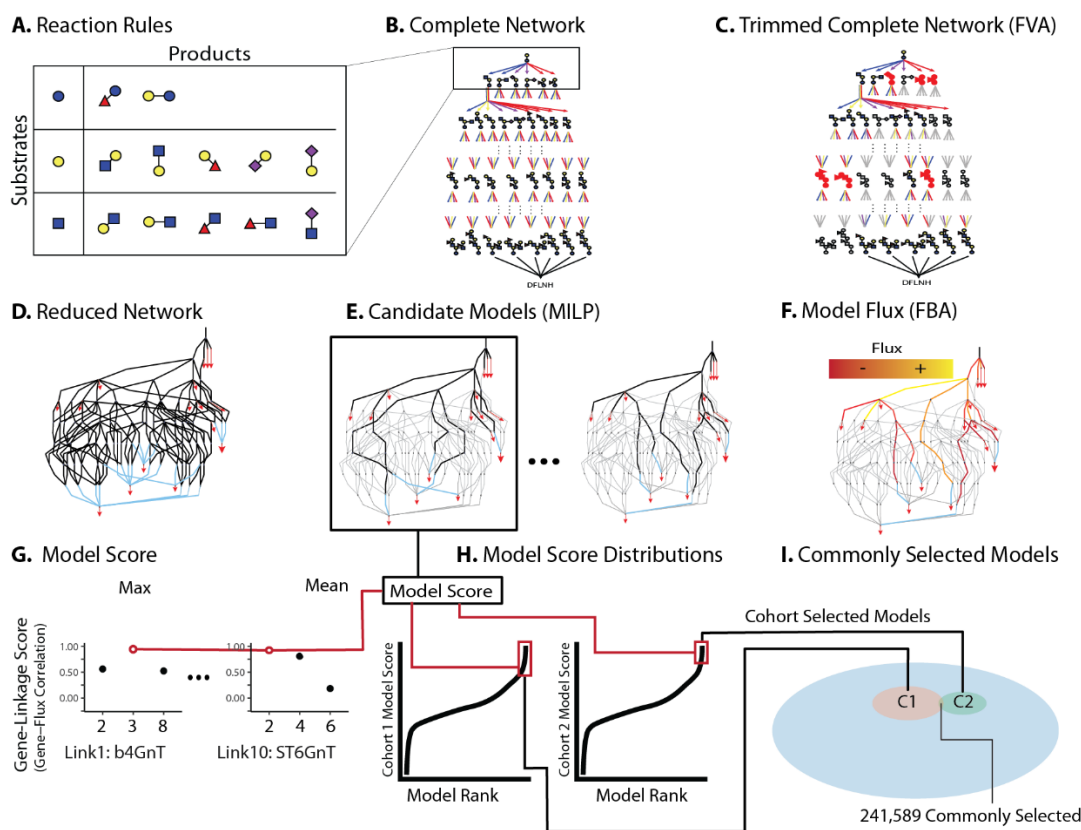
127



128

129 **Figure 2 - FUT2 expression should increase 2'FL and LNFPI which require the enzyme but there**
 130 **is no significant positive association.** Direct comparison of FUT2 gene expression and
 131 concentrations (nmol/mL) of α -1,2-fucose containing HMOs, 2'FL **(A)** and LNFPI **(B)**, in sample-
 132 matched microarray and HPLC reveal no significant association in secretor women from cohort 1
 133 sampled between day 1 and 42 post-partum. Trendlines and points are colored by subject. Linear
 134 trends were used to illustrate the intuition of the GEE approach used to estimate these associations
 135 across subjects. Non-secretor mothers were excluded due to non-functional FUT2. **(C)** A heatmap of
 136 all HMO concentrations across cohort 1 and cohort 2 (top-bar black and grey respectively). Known
 137 HMO structures are shown to the left of each row while uncharacterized structures are indicated with
 138 a black box. For proposed isomers of uncharacterized structures, see **Figure S 12**.

139



140

141 **Figure 3 - Overview of Computational Methods for Model Assembly (A-F) and Assessment (G-I).**

142 (A) To build the candidate models of HMO biosynthesis, reaction rules were defined to specify all possible
 143 monosaccharide additions. (B) The Complete Network includes all oligosaccharides and reactions
 144 resulting from the iterative addition of monosaccharides to a root lactose. (C) Using Flux Variability
 145 Analysis, the Complete Network was trimmed, removing reactions that cannot reach experimentally-
 146 measured HMOs, to produce a (D) Reduced Network **Figure S 9**; red triangles are observed HMOs blue
 147 lines are “sink reactions” joining alternative isomers (**Figure S 12**). (E) From the Reduced Network, Mixed
 148 Integer Linear Programming (MILP) was used to extract Candidate Models, each representing a
 149 subnetwork capable of uniquely synthesizing the observed oligosaccharide profile using a minimal
 150 number of reactions; black clines are reactions retained in a candidate model. (F) Flux Balance Analysis
 151 was used to estimate flux through each reaction necessary to simulate the measured oligosaccharide
 152 concentrations. (G) Model scores were computed as the average maximum correlation between linkage-
 153 specific candidate genes and normalized flux through that linkage (**Figure S 11, S1.1.4**). (H) Model scores
 154 were parameterized on cohort 1 (left) and cohort 2 (right) data (see Methods). High-performing models,
 155 95th percentile of scores, are highlighted in red. (I) Of the >40 million models considered (blue), 2.66
 156 and 2.32 million models were high-performing when parameterized on data from cohort 1 or cohort 2,
 157 respectively. Nearly 250,000 models consistently explained the relationship between predicted flux and
 158 expression data from both cohort 1 and cohort 2. These commonly selected models were analyzed for
 159 common structural features.

160 **2.2 HIGH-PERFORMING CANDIDATE BIOSYNTHETIC MODELS ARE SUPPORTED BY GENE EXPRESSION**
 161 **AND PREDICTED MODEL FLUX ACROSS SUBJECTS**

162 We built and examined models for HMO biosynthesis in human mammary gland epithelial cells. From
 163 the basic reaction set (**Figure 3A**), we generated the complete reaction network (**Figure 3B**) containing

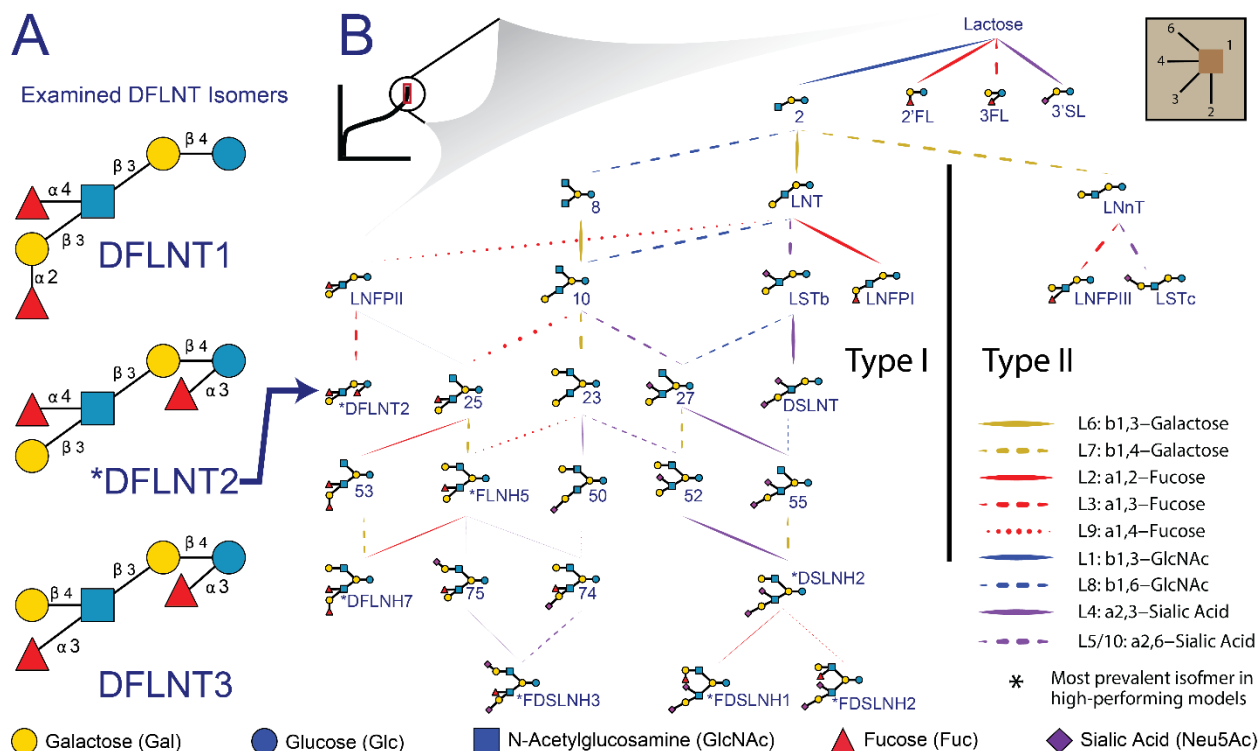
164 all possible reactions and HMOs with up to nine monosaccharides. The Complete Network was
165 trimmed to obtain a Reduced Network (**Figure 3D**; **Figure S 9**, **Table S 2**) by removing reactions
166 unnecessary for producing the observed oligosaccharides. Candidate models (**Figure 3E**) were built,
167 capable of uniquely recapitulating the glycoprofiling data from milk using two independent cohorts-
168 cohort 1 with 8 samples from 6 mothers between 6 hours and 42 days postpartum^{28,29} and cohort 2
169 with 2 samples per mother on the 1st and second day after birth³⁰. Mixed integer linear programming
170 was used to identify subnetworks with the minimal number of reactions from the Reduced Network.
171 We identified 44,984,988 candidate models that can synthesize the measured oligosaccharides. Each
172 candidate model contains 43-54 reactions (19.5-24.4% of the reactions in the Reduced Network
173 (**Table S 3**)). These models covered all the feasible combinations of HMO synthesis by the 10 known
174 glycosyltransferase families (**Figure 1D**) that could describe the synthesis of the HMOs in this study.

175 To identify the most likely biosynthetic pathways for HMOs, we computed a model score for each
176 candidate model using the glycoprofiling and transcriptomic data from the two independent cohorts,
177 after excluding low-expression gene candidates. Genes were excluded when expression was
178 undetected in over 75% of microarray samples and the independent RNA-seq³¹ measured low
179 expression relative to the GTEx³²: TPM<2 and 75th percentile Lemay < GTEx Median TPM. Specificity
180 and expression filtration reduced the candidate genes from 54 to 24 (see supplemental results, **Table**
181 **S 1**, **Figure S 7**); three linkages (L2, L5 and L9) were resolved by filtration alone indicating that FUT2,
182 ST6GAL1 and FUT3 respectively perform these reactions.

183 Following low-expression filtering, we compared flux-expression correlation. Leveraging sample-
184 matched transcriptomics and glycomics datasets, we computed model scores indicating the capacity
185 of each candidate gene to support corresponding reaction flux. The model score was computed by
186 first identifying for each reaction, the candidate gene that shows the best Spearman correlation
187 between gene expression and normalized flux; flux was normalized as a fraction of the input flux to
188 limit the influence of upstream reactions (**Figure S 11**, S1.1.4). The highest gene-linkage scores, for
189 each reaction, for each model were averaged to obtain a model score (**Figure 3G**, see Methods). The
190 model scores indicate consistency between gene expression and model-predicted flux. The high-
191 performing models ($z(\text{model score}) > 1.646$) were selected for further examination (**Figure 3H**, see
192 Methods). Though quantile-quantile plots indicated the model score distributions were pseudo-
193 gaussian, variation in skew resulted in slightly different numbers of high-performing models for the
194 two different subject cohorts. Specifically, we found 2,658,052 high-performing models from cohort
195 1 and 2,322,262 high-performing models using cohort 2 (**Figure 3I**, **Table S 4**). We found 241,589 high-
196 performing models common to cohort 1 and cohort 2. The model scores of commonly high-
197 performing models are significantly correlated (Spearman $R_s = 0.2$, $p < 2.2e-16$) and a hypergeometric
198 enrichment of cohort 1 and cohort 2 selected models shows the overlap is significant relative to the
199 background of 44 million models ($p < 2.2e-16$). We analyzed these 241,589 commonly high-
200 performing models and determined which candidate genes were common in high-performing
201 models.

202 To determine the most important reactions (**Figure 4**) in the reduced network, we asked which
203 reactions were most significantly and frequently represented among the top 241,589 high-
204 performing models. We then filtered to retain only the top 5% of most important paths from lactose

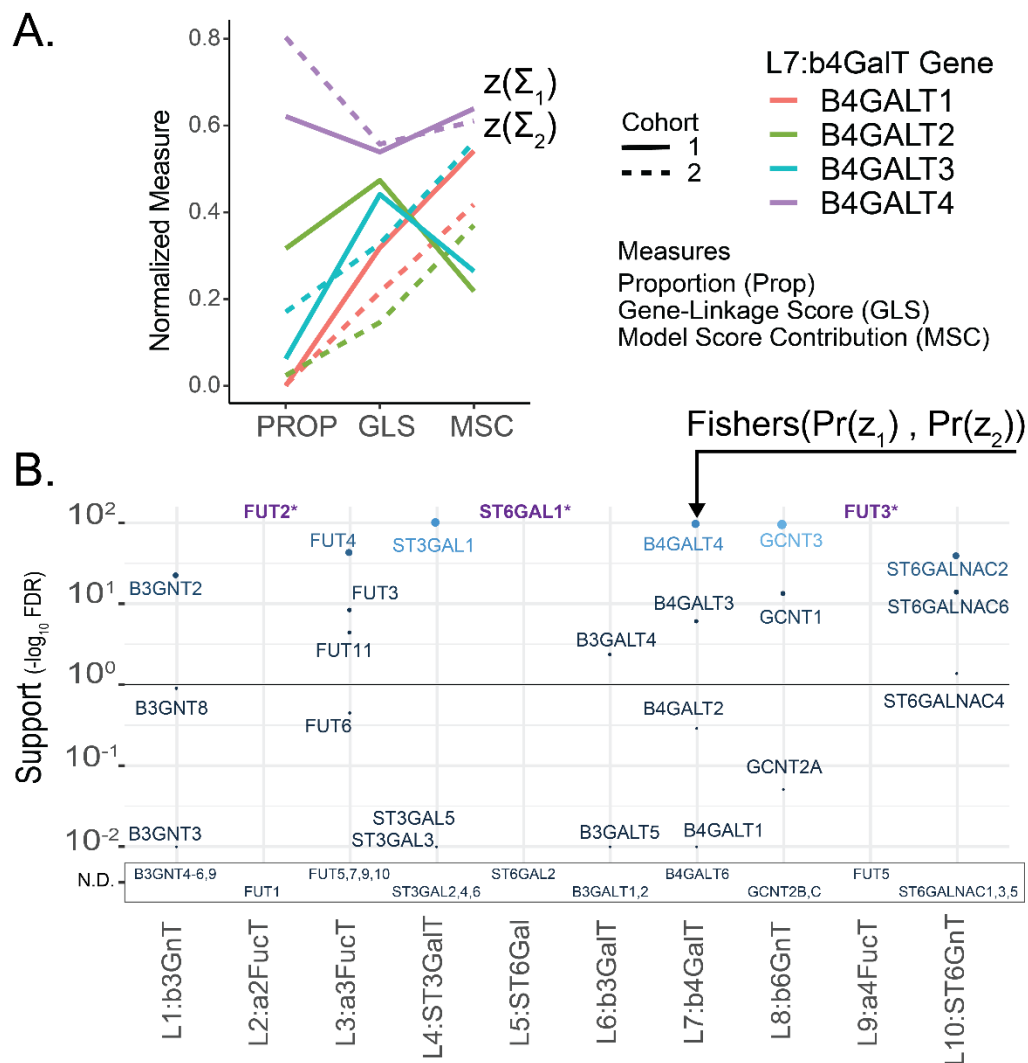
205 to each observed HMO (see Methods). The most important reactions form the Summary Network
206 (**Figure 4**). Here, HMO biosynthesis naturally segregates into type-I backbone structures, with β -1,3-
207 galactose addition to the GlcNAc-extended core lactose, and type-II structures, with β -1,4-galactose
208 addition to the GlcNAc-extended core lactose. As expected, LNFPI, LNFPII, LSTb and DSLNT segregate
209 to the type-I pathway while LNFPIII and LSTc are found in the type-II pathway (see Methods for HMO
210 definitions). The Summary Network suggests resolutions to large structurally ambiguous HMOs
211 (FLNH5, DFLNT2, DFLNH7, and DSLNH2) by highlighting their popularity in high-performing
212 models. The Summary Network also shows three reactions of high comparable strength projecting
213 from GlcNAc- β 1,3-lactose to LNT, LNnT and a bi-GlcNAc-ylated lactose (HMO8, **Figure 4**, **Table S 2**)
214 suggesting LNT may be bypassed through an early β -1,3-GlcNAc branching event; a previously
215 postulated alternative path³³. We checked for consistency with previous work³⁴ and found that (1)
216 the single fucose on the reducing-end Glc residue is always α -1,3 linked, (2) for monofucosylated
217 structures, the non-reducing terminal β -1,3-galactose is α -1,2-fucosylated, (3) all galactose on the β -
218 1,6-GlcNAc is always β -1,4 linked while all galactose on the β -1,3-GlcNAc are either β -1,3/4 linked.
219 With the exception FDSLNH1, (4) no fucose is found at the reducing end of a branch and (5) all α -1,2-
220 fucose appear on a β -1,3-galactose and not β -1,4-galactose in monofucosylated structures with more
221 than four monosaccharides; suggesting that FDSLNH1 is an unlikely isomer. The summary network
222 also suggests that most HMOs have type-I LacNAc backbones. To address the potential over-
223 representation of type-I HMOs in our models, we examined the distribution of type-I and type-II in
224 tetra- and pentasaccharides with known structures. Across samples, the median abundance of type-
225 II HMOs, LNnT, LNFPIII and LSTc were 3.33%, 0.041%, and 2.68% of total nmol/mL while type-I
226 HMOs of the same size, LNT, LNFPI, LNFPII, and LSTc, was 15.3%, 9.39%, 7.45% and 0.45%
227 respectively. This confirms the greater abundance of type-I HMOs compared with the type-II
228 structures in the glycomic profiles (**Figure 2C**). This Summary Network thus provides orientation in
229 this underspecified space.
230



231 **Figure 4 Summary Network of the most important reactions in the reduced network.** Observed,
 232 intermediate and candidate HMOs most important to commonly high-performing networks were
 233 selected from the Reduced Network (**Figure 3D**; Supplemental Methods S1.1.2). (A) Several
 234 ambiguous isomers (**Figure S 12**) were preferred (**Figure S 1**) in the commonly high-performing
 235 models. (B) A summary network was constructed from reaction importance; an aggregation of the
 236 proportion of high-performing models that include a reaction, and the enrichment of a reaction in
 237 the high-performing model set (see Methods). Line weight indicates the relative importance of each
 238 reaction. Line color corresponds to the monosaccharide added at each step and line type corresponds
 239 to the linkage type. The Summary Network naturally segregates into type-I and type-II backbone
 240 structures. For measured HMO definitions (e.g. FDSLNH and DSLNT) see Methods, for intermediate
 241 HMO definitions (e.g. 8, 10, or 25) see **Table S 2**, for uncertain structures (e.g. DFLNH7, FLNH5) see
 242 **Figure S 12**.

243

244



245

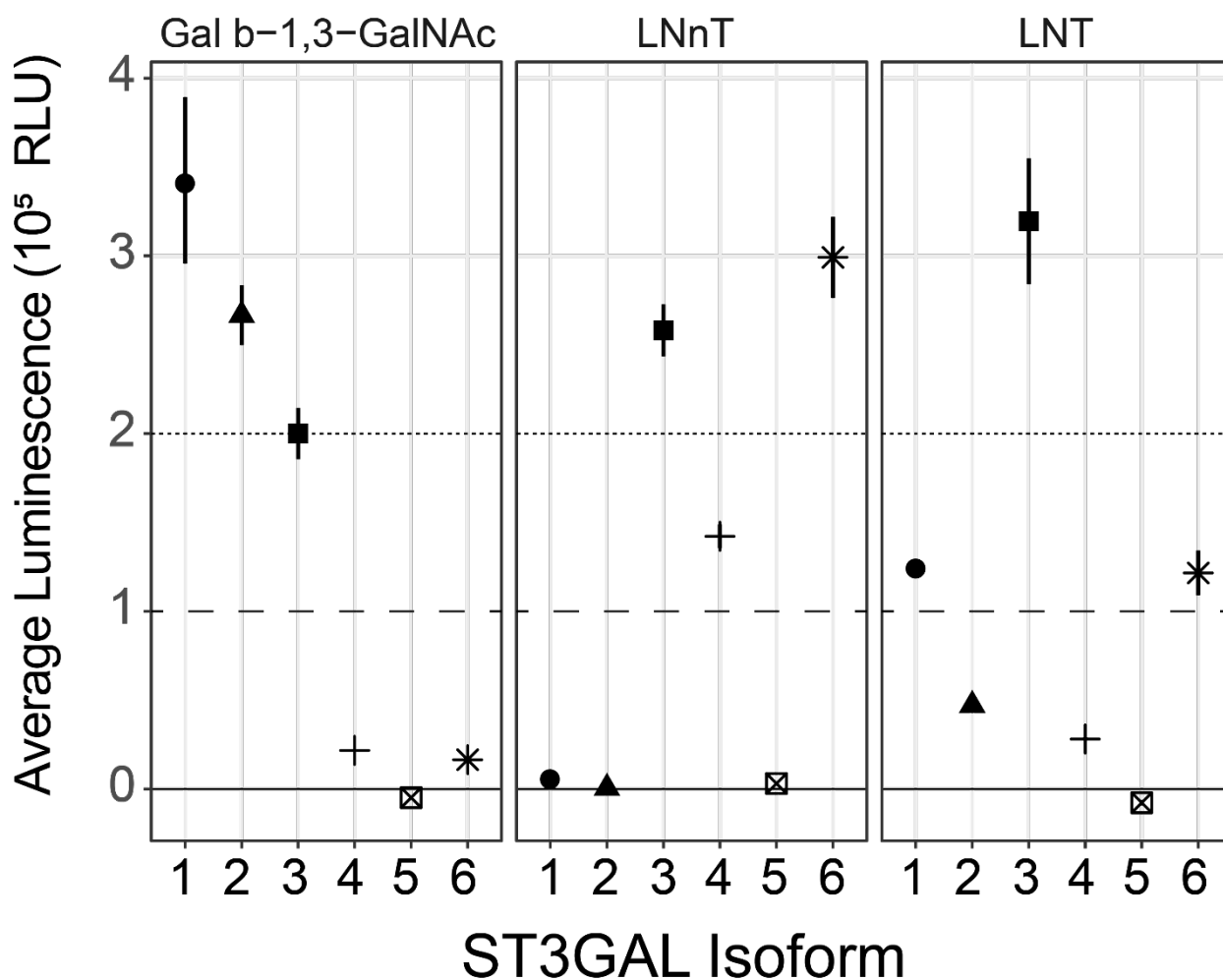
246 **Figure 5 – Gene expression correlation with model flux predicts enzymes involved in HMO**
 247 **biosynthesis.** (A) To determine the gene expression that best explains flux through each reaction in
 248 each glycomics-transcriptomics matched sample, we examined the proportion of high-performing
 249 models were each gene was most flux-correlated (PROP, **Figure S 4**), we also examined the gene-linkage
 250 score (GLS) and Model Score Contribution (MSC). For this visual, each measure was max-min normalized
 251 between 0 and 1. Genes were selected based on high performance on all three measures across cohorts
 252 (line type). (B) We summarize the three performance scores from panel A across cohorts into a single
 253 support score (see Methods). Briefly, “Support” is p-value for the sum of PROP, GLS and MSC z-scores
 254 (relative to a permuted background), Fisher-pooled across cohorts then False Discovery Rate (FDR)
 255 corrected across genes (see Methods). Unmeasured genes appear below the plot in the Not Determined
 256 (N.D.) box. Genes selected by default (purple, “*”) as the only measured gene candidate (**Table 1**)

257

258 **2.3 GLYCOSYLTRANSFERASES ARE RESOLVED BY RANKING REACTION CONSISTENCY ACROSS SEVERAL** 259 **METRICS**

260 We further analyzed the high-performing models to identify the glycosyltransferases responsible for
261 each step in HMO biosynthesis (**Table 1**). As previously described, not all members of a gene family
262 were examined in this analysis. Some genes were excluded due to their well characterized irrelevance
263 (e.g. FUT8) and others, like FUT1, were excluded due to low expression in lactating breast epithelium
264 (see **Table S 1**, methods and supplemental results for the detailed inclusion criteria). To determine
265 the genes preferred for each reaction, we used three metrics to quantify the association between
266 candidate gene expression and predicted flux. These were (1) *proportion* (*PROP* - the relative
267 proportion of models best explained by a candidate gene, **Figure S 4**), (2) *gene linkage score* (*GLS* - the
268 average Spearman correlation between gene expression and flux), and (3) *model score contribution*
269 (*MSC* - an estimate of the gene-influence indicated by the Pearson correlation between model score
270 and gene linkage score) (**Figure 5A**, **Figure S 5**). For each candidate gene, we generated a reaction
271 support score (**Figure 5B**, see Methods); the pooled significance of the maxima of PROP, GLS and MSC
272 across both cohorts.

273 Three reactions, L2 (FUT2), L5 (ST3GAL1) and L9 (FUT3), were matched to genes by default as they
274 were the only gene candidates remaining following gene expression filtering (**Table S 1**,
275 Supplementary Results). At least one gene showed significant support ($q < 0.1$) for each remaining
276 reaction. GCNT3 shows highly significant support ($q < 0.001$) and nearly 100% of models selected this
277 isoform over GCNT2C or GCNT1 (**Figure S 4**). B4GALT4 is the most significantly supporting gene for
278 the L7: b4GalT reaction (**Figure 5B**). In both cohort 1 and 2, B4GALT4 outperforms all other isoforms
279 in all three metrics. B4GALT4 expression best explains flux in 62% and 80% (PROP) of high-
280 performing models using cohort 1 and 2 data respectively (**Figure S 4**). B4GALT4 also has the highest
281 MSC and GLS ($z > 5.6$) of any isoforms. Interestingly, while B4GALT1 is highly expressed and
282 fundamental to lactose synthesis in the presence of α -lactalbumin and lactation in general^{35,36}, it
283 showed negligible support for the L7 reaction (**Figure 5B**). Considering the reaction support score, all
284 linkages show at least one gene for each reaction that significantly explains behavior across cohorts
285 (**Figure 5B**).



286

287 **Figure 6 - Results of the CMP-Glo™ Glycosyltransferase Assay to test GT candidates on relevant**
 288 **HMO acceptors.** Average luminescence below 10,000 is considered weak activity, and activity above
 289 200,000 is considered very high activity. Reported luminescence values were background corrected and
 290 95% confidence intervals are shown. For complete details see **Table S 6**. Shapes correspond to ST3GALT
 291 isoforms

292

293 2.4 KINETIC ASSAYS CONFIRM SELECTED GENES AND EXPAND OUR SCOPE

294 Towards validating and expanding our gene-reaction predictions, glycosyltransferase enzyme
 295 activity assays were performed using the NTP-Glo™ Glycosyltransferase assay format from Promega.
 296 We used linkage L1:b3GnT and L10:ST6GnT to validate our selections and examined every plausible
 297 isoform of the ST3GAL for its ability to perform the linkage L4:ST3GalT reaction. Five acceptors were
 298 used: (1) lactose to examine activity on the initial HMO acceptor, (2) LNT and (3) LNnT to establish
 299 which enzymes would act on larger type-I and type-II tetrasaccharides, (4) Gal β 1,3-GalNAc to
 300 determine specificity for non-HMO O-type glycans, and (5) a GlcNAc- β 1,3-Gal- β 1,4-GlcNAc- β 1,3-Gal-
 301 β 1,4-Glc pentasaccharide structure to test the formation of a non-reducing terminal type-I (Gal-b1,3-

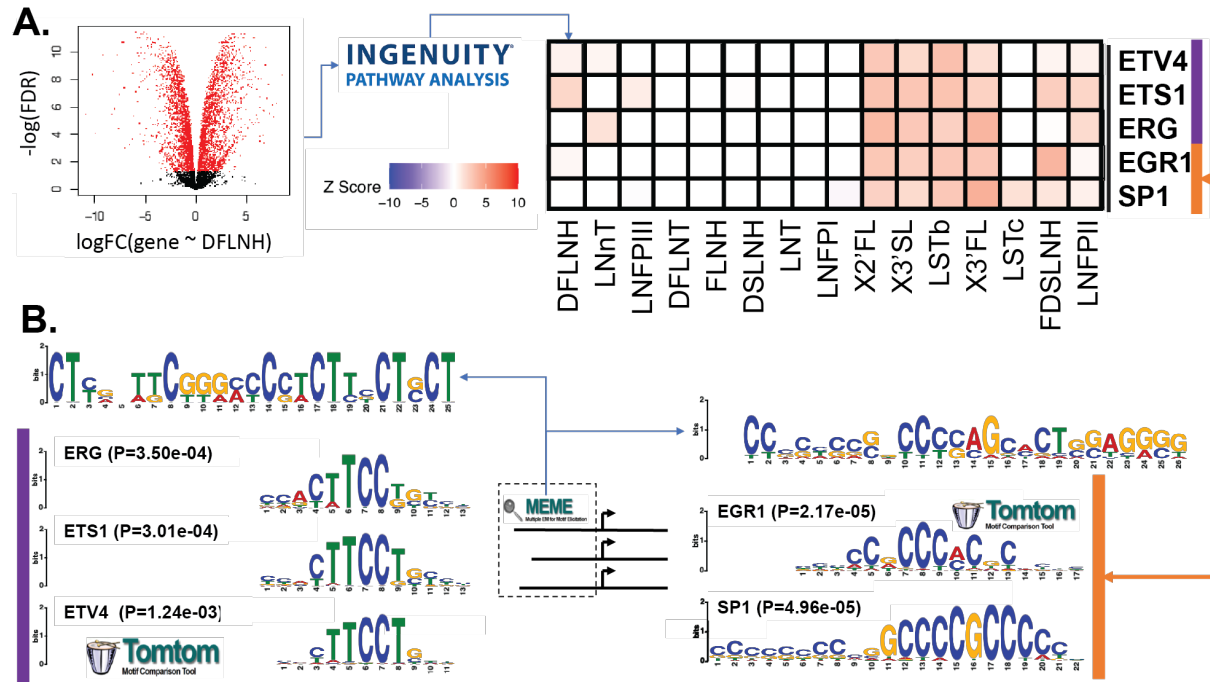
302) cap on a longer acceptor. We explored the activities of various gene products to perform specific
303 glycosyltransferase reactions crucial to HMO biosynthesis (**Figure 6, Table S 5**).

304 In the cross-cohort aggregate analysis (**Figure 5B**), B3GNT2 is selected as a reasonable candidate to
305 catalyze flux through the L1:b3GnT reaction. The B3GNT2 support score is nearly 100 times more
306 significant than B3GNT8, the next most associated gene. Consistent with the predictions that b3GnT
307 should convert lactose into the precursor to LNT and LNnT, the UDP-Glo™ assay showed B3GNT2 had
308 high activity toward lactose as an acceptor. We further found that B3GNT2 could add a β 1,3-GlcNAc
309 to LNnT as is necessary for poly-lacNAc HMOs. The cross-cohort aggregate analysis (**Figure 5B**)
310 selected ST6GALNAC2 to perform L10, the α 2,6 addition of sialic acid to the internal β 1,3-GlcNAc;
311 necessary for the biosynthesis of LSTb from LNT and possibly DSLNT from LSTa. However, the CMP-
312 GLO™ assay highlighted a negligible activity of ST6GALNAC2 toward LNT even at very high enzyme
313 input indicating that this enzyme does not convert LNT to LSTb. We did not test if it can convert LSTa
314 to DSLNT. In contrast, ST6GALNAC5 was effectively able to use LNT as an acceptor, although we did
315 not confirm the formation of the LSTb structure. ST6GALNAC5 could not be considered in the support
316 score calculation because it was only measured in cohort 2; expression was greater than zero in 1 of
317 12 samples.

318 Finally, we tested the affinities of plausible ST3GAL isoforms to sialylate LNT, LNnT or β 1,3-GlcNAc
319 (**Table S 5**). The multi-cohort analysis (**Figure 5B**) implicates ST3GAL1 as the best candidate for this
320 reaction. The CMP-Glo™ assay indicated that ST3GAL1 has limited activity toward LNT but high
321 activity toward Gal β 1,3-GlcNAc suggesting ST3GAL1, *in vitro*, is more involved in non-HMO O-type
322 glycan biosynthesis. ST3GAL2 showed a similar but less substantial pattern. ST3GAL3 showed the
323 strongest activity for sialylation both LNT and LNnT suggesting it could synthesize LSTa from LNT.
324 ST3GAL6 shares a similar but lesser activity for LNT and LNnT.

325 We analyzed the original expression profiles to determine which genes were sufficiently expressed
326 to actuate this activity. STGAL1, 3 and 5 were strongly expressed in nearly 100% of samples across
327 both cohorts; ST3GAL2 and 4 show zero expression in 75% of samples in at least one cohort (**Figure**
328 **S 6**). ST3GAL3 was highly expressed and effective at catalyzing the L4 reaction for LNT and LNnT
329 while ST3GAL1 was highly expressed and weakly catalyzed sialylation of LNT making ST3GAL3 the
330 most likely candidate for L4 reaction on LNT and LNnT.

331



332

333 **Figure 7 - de novo promoter-enriched TF motifs and IPA predicted TFs using differential**
 334 **expression analyses with respect to 16 HMOs. (A) MEME identified TF motifs and 5 known TFs (ETV4,**
 335 **ETS1, EGR1, SP1, and ERG) associated with them (see Table S 6). MEME-discovered TFs were cross-**
 336 **referenced with known TF binding sites using TOMTOM. Logos for the matched known and discovered**
 337 **motifs are shown in the top and bottom of each subpanel; the p-value is a logo matching significance**
 338 **calculated by TOMTOM. (B) Subset of a biclustering of activation z-score computed by IPA indicating**
 339 **the likelihood that a TF activates ($z > 0$) or inhibits ($z < 0$) an HMO concentration signature (gene**
 340 **expression associated with changes in HMO concentration). The full biclustering can be found in the**
 341 **supplement (Figure S 16)**

342 **Table 2 - TF motif (MEME) and IPA upstream regulator integrated results**

LINKAGE	REACTION SUPPORT SCORE SELECTED CANDIDATE	MEME TF MOTIF (P-VALUE)*1	JASPAR TF (P-VALUE)*2	IPA PREDICTED TF*3
L1:B3GNT	B3GNT2	TF Motif-II	SP1 (4.96e-05)	Y
		(1.39e-12)	EGR1 (2.17e-05)	Y
L2:A2FUCT	FUT2	TF Motif-III	IKZF1 (7.62e-04)	Y
L3:A3FUCT	FUT11	TF Motif-II	SP1 (4.96e-05)	Y
		(1.00e-7)	EGR1 (2.17e-05)	Y
L4:ST3GAL T	ST3GAL1	TF Motif-I	ETV4 (1.24e-03)	Y
		(2.76e-11)	ETS1 (3.01e-04)	Y
			ERG (3.50e-04)	Y
L7:B4GALT	B4GALT4	TF Motif-II	SP1 (4.96e-05)	Y
		(7.67e-11)	EGR1 (2.17e-05)	Y
L10:ST6GNT	ST6GALNAC2	TF Motif-II	SP1 (4.96e-05)	Y
		(1.08e-7)	EGR1 (2.17e-05)	Y

343 *1. The P-value (see **Figure S 15**) is the significance of the selected GT to the MEME identified TF motif.

344 *2. The P-value (see **Table S 7**) is the significance of known TF associated with the MEME identified TF
345 motif.

346 *3. The IPA upstream regulator analyses were conducted on the three different sets of DEGs: 16 HMOs,
347 19 glycan motifs, and 4 differential motifs (see **Methods** for details). Based on the Z-score predicted by
348 IPA using the gene expression data, we selected the significant TFs with IPA predicted activation score
349 $|Z \text{ value}| \geq 3$ in this study. Note, 'Y' denotes the known TF is presented in the indicated dataset (HMO
350 (**Figure 7**, **Figure S 16**), Motif (**Figure S 17**), or differential motif (**Figure S 18**)) of the IPA predicted TF
351 and 'N' means the TF doesn't present in the dataset of IPA predicted TF.

352 2.5 SELECTED GLYCOSYLTRANSFERASES SHARE TRANSCRIPTIONAL REGULATORS ACROSS 353 INDEPENDENT PREDICTIONS

354 To explore the transcriptional regulation during lactation, we used two orthogonal approaches for
355 transcription factor (TF) discovery. We used Ingenuity Pathway Analysis (IPA) to predict upstream
356 regulatory factors based on differential expression associated with each HMO. IPA analyzed all genes
357 differentially expressed with HMO abundance, not only HMO glycoconjugates; these differential
358 expression patterns formed HMO specific gene expression signatures. Additionally, we used MEME
359 for *de novo* motif discovery in the promoter regions of HMO glycoconjugates and TOMTOM to map those
360 discovered motifs to known TFs. We validated these predictions by examining transcriptional
361 regulators selected by both MEME and IPA (**Figure S 13**, see **Methods**).

362 IPA discovered 57 TFs significantly ($|z| \geq 3$; $p < 0.001$) associated with the 16 HMO-specific gene
363 expression signatures. We performed differential expression on HMO substructure abundance and
364 substructure abundance ratios¹⁷; IPA found 66 and 49 TFs significantly ($|z| \geq 3$; $p < 0.001$) associated
365 with HMO substructure and substructure ratio specific gene expression signatures. Using MEME, we

366 identified three putative TF regulatory sites (TF motifs I, II and III) for 6 selected glycosyltransferases
367 responsible for the HMO biosynthesis (**Table 2** and **Figure S 15**). TOMTOM calculated that these
368 putative binding sites were significantly associated with six known TFs (IKZF1, SP1, EGR1, ETS1,
369 ETV4 and ERG) that were also predicted by IPA as regulators of gene signatures associated with HMO
370 concentration (**Figure 7**, **Figure S 16**) or HMO glycan substructures abundance (**Figure S 17**). SP1,
371 EGR1, ETS1, ETV4 and ERG are all predicted to positively influence expression associated with the
372 biosynthetically related HMOs: 3'SL, 3FL, LSTb and DSLNT; 3'SL and 3FL share a common substrate
373 (lactose) while LSTb is a likely precursor to DSLNT. The motif-level analysis showed opposing
374 regulation between IKZF1: upregulating gene expression signatures associated with the 3'SL and
375 LSTb substructure abundance¹⁷ (X34 and X62 respectively, see **Figure S 19**) and downregulating
376 gene expression associated with GlcNAC-lactose, LNT and LNFPI substructure abundance (X18, X40
377 and X65 respectively, see **Figure S 19**), while EGR1, ERG and ETS1 have the opposite predicted
378 impact (**Figure S 17**). The motif-level predictions are consistent with the HMO-level predictions of
379 upregulation on 3'SL and LSTb while adding an additional point of contrast. While EGR1, ERG and
380 ETS1 are predicted to increase production of sialylated HMOs, they may have the opposite impact on
381 LNFPI. Thus, we detect signatures of multiple transcription factors that could coordinate the
382 regulation of the genes we identified to contribute to HMO biosynthesis (see supplemental
383 discussion).

384

385 **3 DISCUSSION**

386 By integrating sample-matched quantitative oligosaccharide measurements and gene expression
387 data using computational models of HMO biosynthesis, we resolved genes responsible for 10
388 elementary reactions in human mammary gland epithelial cells. The modeling-based strategy was
389 essential since simple correlations failed to capture the simplest HMO-gene associations, given the
390 complex interactions of glycosyltransferases in the HMO biosynthetic pathway. Because the pathway
391 characterization is still incomplete, we built >44 million candidate models that uniquely recapitulate
392 glycoprofiling data in two independent cohorts. Candidate model flux, i.e. activity of each reaction,
393 was predicted for each model and compared to sample-matched gene expression data. We used the
394 consistency between gene expression and predicted flux across cohorts in high-performing models
395 to select genes for each fundamental reaction. Analysis of these models suggested glycosyltransferase
396 genes, thus providing a clearer picture of the enzymes and regulators of HMO biosynthesis in
397 mammary epithelial cells. The clarification of the pathways and enzymes involved in HMO
398 biosynthesis will be an invaluable resource to help (1) discover the maternal genetic basis of health-
399 impacting^{1,2,5,6,37-46} HMO composition heterogeneity^{7,26,47,48} and (2) drive chemoenzymatic synthesis
400⁴⁹⁻⁵³ and metabolic engineering for manufacturing HMOs for food ingredients, supplements and
401 potential therapeutics⁵⁴⁻⁵⁹ (see supplemental discussion).

402 Of the three fucosylation reactions, two were determined using expression data alone while the third
403 required additional insight from the flux-expression comparison or, support score. Consistent with
404 studies in blood²³⁻²⁵ and milk^{26,47,60} types, we selected FUT2 as the gene supporting the α 1,2-
405 fucosylation (L2:a2FucT) linkage reaction. FUT1 was ruled out due to non-expression (**Table S 1**,

406 supplemental results). In the second fucosylation reaction, FUT3, FUT4 and FUT11 all show
407 significant support for α 1,3-fucosylation (L3:a3FucT) linkage formation. FUT11 is more commonly
408 considered an N-glycan-specific transferase⁶¹ and therefore a less likely candidate. Both FUT3 and
409 FUT4 prefer to fucosylate the inner GlcNAc of a type-I polylectosamine⁶². FUT3 prefers neutral type-
410 I polylectosamine while FUT4 also fucosylates the sialylated form^{63,64}; the charge preferences are
411 inverted for type-II polylectosamine acceptors⁶⁵. Prudden et. al.⁵² used FUT9 to perform this reaction,
412 consistent with its ability to transfer α 1,3 fucose to the distal GlcNAc of a neutral polylectosamine⁶¹⁻
413 ⁶³. The four HMO structures with α 1,3-Fucose in the Summary Network (**Figure 4**) include 3FL
414 (neutral inner fucosylation), LNFPIII (neutral distal fucosylation), DFLNT2 (neutral inner
415 fucosylation), and FDSLNH2 (sialylated and neutral distal fucosylation). FUT9 showed negligible
416 expression in RNA-Seq (3rd Quartile TPM=0.37, **Table S 1**), yet it is highly expressed (TPM>10) brain
417 and stomach³². Therefore, it is likely that the distal fucosylation is conducted by another enzyme *in*
418 *vivo* while the inner fucosylation is likely performed by either FUT3 or FUT4. FUT3 was also chosen
419 for the α 1,4-fucosylatoin (L9:a4FucT) by default due to the non-expression of FUT5, confirmed by
420 RNA-Seq (**Table S 1**, supplemental results). FUT3 adds an α 1,4-fucose to the GlcNAc of a neutral type-
421 I chain to form the Lewis-A or Lewis-B group and adds an α 1,3-fucose to the GlcNAc of a type-II
422 chain^{63,64}. Usage of FUT3 would provide a parsimonious explanation for the fucosylation of both type-
423 I and type-II HMOs like LNFPII (Fuc- α 1,4-LNT (type-I)) and LNFPIII (Fuc- α 1,3-LNnT (type-II)).

424 One of two sialyltransferases was clearly resolved with expression data alone, the other required
425 additional examination. ST6GAL1 was chosen by default to support the α 2,6-sialylation (L5:ST6GalT)
426 reaction due to the non-expression of ST6GAL2 (**Table S 1**). ST6GAL1 sialylates galactose in HMOs⁵².
427 For the second sialylation reaction, our flux-expression comparison selected ST6GALNAC2 and
428 ST6GALNAC6 as the significant supporters of α 2,6 sialylation (L10:ST6GnT). Through a kinetic assay,
429 we confirmed that ST6GALNAC2 (previously shown to accept core-1 O-glycans^{66,67}) fails to sialylate
430 LNT. Though our kinetic assay shows that ST6GALNAC5 (known to sialylate GM1b⁶⁸) can sialylate
431 LNT, it was not expressed in this context (**Table S 1**, supplemental results). ST6GALNAC3 expression
432 was not observed in microarrays but could not be ruled out due to RNA-Seq expression (**Table S 1**,
433 supplemental results); it sialylates the GalNAc of NeuAc- α 2,3-Gal- β 1,3-GalNAc- α 1-O-Ser/Thr and
434 NeuAc- α 2,3-Gal- β 1,3-GalNAc- β 1,4-Gal- β 1,4-Glc- β 1-Cer when the inner galactose is not sialylated
435 (e.g. GD1a or GT1b)⁶⁹⁻⁷² but has not been shown to transfer to a GlcNAc. The last ganglioside-
436 accepting family gene, ST6GALNAC6, has broader activity accepting several gangliosides (GM1b,
437 GD1a, and GT1b)⁶⁹ and sialylating the GlcNAc of LNT-ceramide⁷³. Considering the broader activity,
438 clear expression and computational selection, ST6GALNAC6 is the most likely candidate, though
439 ST6GALNAC3 should not be ruled out. In the third reaction, ST3GAL1 shows significant support for
440 α 2,3-sialylation (L4:ST3GalT) reactions while ST3GAL3 shows negligible consistency in the flux-
441 expression comparison. Yet, *in vitro*, ST3GAL3 was most effective at sialylating both LNT and LNnT
442 in kinetic assays while ST3GAL1 weakly sialylated LNT. ST3GAL4, which prefers type-II acceptors⁷⁴⁻
443 ⁷⁶, was used previously to perform this reaction *in vitro*⁵², but it was not expressed on the microarrays
444 nor RNA-Seq. ST3GAL3 can accept type-I, type-II and type-III acceptors including LNT and prefer
445 type-I acceptors^{74,75,77} while ST3GAL1 accepts type-I, type-III and core-1 acceptors but not type-
446 II^{74,75,78}. The kinetic assays and previous literature show ST3GAL3 is more capable than ST3GAL1 at
447 catalyzing this reaction, while ST3GAL1 expression was found to be the only plausible candidate

448 based on estimated flux through this reaction. If ST3GAL1 were responsible for this reaction, its
449 inability to sialylate type-II HMO could partially explain the lack of sialylation and larger structures
450 in the type-II HMO branch. Both ST3GAL1 and ST3GAL3 remain plausible candidate genes, and
451 further *in vivo* studies are needed. Both galactosylation reactions required further examination of
452 flux-expression relationships. We found B3GALT4 to significantly support the type-I β -1,3-galactose
453 addition (L6:b3GalT). B3GALT4 can transfer a galactose to GalNAc in the synthesis of GM1 from
454 GM2⁷⁹. Unlike B3GALT5, there is no evidence that B3GALT4 can transfer galactose to a GlcNAc⁸⁰.
455 B3GALT5, has been shown to transfer a β -1,3-galactose to GlcNAc to form LNT *in vitro*⁸¹. B3GALT5
456 expression measured for cohort 1 microarray was much lower than expression in cohort 2 and the
457 independent RNA-Seq³¹ suggesting that the probes in the first microarray may have failed (**Table S 1**,
458 supplemental results). While both B3GALT4 and B3GALT5 seem plausible, given the historical
459 failures of B3GALT4 to perform this reaction and our likely failure to measure and evaluate B3GALT5,
460 B3GALT5 may be the stronger candidate for this reaction. In the second galactosylation reaction, the
461 flux-expression comparison found B4GAL4 and B3GALT3 most significantly supports the type-II
462 definitive β -1,4-galactose addition (L7:b4GalT). These gene-products can synthesize LNnT-
463 ceramide⁸². Additionally, in the presence of α -lactalbumin (highly expressed during lactation),
464 B4GALT4 shows an increased affinity for GlcNAc acceptors suggesting during lactation it is more
465 likely to perform the L7 reaction^{82,83}. B4GALT1 and B4GALT2 synthesize lactose in the presence of α -
466 lactalbumin during lactation^{35,36}, but B4GALT1 expression was not correlated with L7 flux and
467 B4GALT2 was not expressed (**Table S 1**). We note that associations between B4GALT1 expression L7
468 flux may be masked due to its consistent high. Therefore, flux-expression correlation should not be
469 used to exclude B4GALT1 as a candidate for the L7 reaction. Doing so, B4GALT4, B4GALT3 and
470 possibly B3GALT1 remain the most plausible candidates.

471 Finally, both GlcNAc additions required flux-expression examinations. B3GNT2 showed significant
472 support in the flux-expression comparison. In our kinetic assays, B3GNT2 demonstrated high activity
473 towards lactose as an acceptor. Previously, B3GNT2 has performed the β -1,3-GlcNAc addition
474 (L1:b3GnT) on multiple glycan types including several HMOs: lactose, LNnT, polylectosamine-
475 LNnT⁸⁴. The agreement of literature, kinetic assays and flux-expression analysis indicate B3GNT2 is
476 an appropriate choice for this reaction. In the second GlcNAc reaction, GCNT3 and GCNT1 most
477 significantly support the branching β -1,6-GlcNAc addition (L8:b6GnT). While GCNT2B can effectively
478 transfer the branching GlcNAc to the inner galactose of LNnT^{52,85}, it was not expressed in the cohort
479 microarrays or independent RNA-Seq. GCNT1 transfers a branching GlcNAc to the GalNAc of a core-
480 1 O-glycan^{86,87} while GCNT3 acts on core-1 and the galactose of the LNT-like core-3 structure^{88,89}.
481 GCNT3 is also specifically expressed in mucus-producing tissues^{88,89} like lactating mammary gland
482 epithelium. Interestingly, GCNT3 acts on galactose of the GlcNAc- β 1,3-Gal- β 1,4-Glc trisaccharide
483 (predistally) while GCNT2 acts on the central galactose of the LNnT or LNT tetrasaccharide
484 (centrally)⁸⁵. Therefore, reliance on GCNT3 for the branching reaction would explain the
485 noncanonical branched tetrasaccharide (HMO8, **Figure 4**) suggesting a third major branch from
486 GlcNAc- β 1,6-lactose, distinct from LNT and LNnT. Predistal addition of the branched GlcNAc may
487 also explain the lack of branched type-II structures since B4GALT4 cannot act on branched core-4
488 structures⁹⁰. HMO biosynthesis with GCNT3 and B4GALT4 could explain the type-I bias seen in
489 the Summary Network (**Figure 4**).

490 Our results show consistency with experimental validation here and the published literature. Further
491 direct empirical studies will be invaluable to confirm each gene-reaction association and the
492 complete biosynthesis network. Such studies would include further clinical cohort studies and the
493 development of mammary organoid models capable of producing HMOs. Such experimental systems
494 can clarify the impact of mammary-tissue specific genes, cofactors, and HMO chaperones like α -
495 lactalbumin^{82,83} on glycosyltransferase activity. Therefore, further development of authentic *in vitro*
496 cell and organoid models will be invaluable to finalizing our model of HMO biosynthesis.

497 **4 CONCLUSION**

498 By using systems biology approaches, different omics data can be integrated, as shown here to
499 predict gene-reaction relations even in highly uncertain and underdetermined networks. Of the ten
500 fundamental reactions we aimed to resolve and reduce (**Table 1**), we succeeded in narrowing the
501 candidate substantially for each one. The newly reduced space of HMO biosynthetic pathways and
502 knowledge of the enzymes and their regulation will enable mechanistic insights into the relationship
503 of maternal genotype and infant development. Finally, once essential HMOs are identified, the
504 knowledge presented here on the HMO biosynthetic network can provide insights for large-scale
505 synthesis of HMOs as ingredients, supplements, or potential therapeutics to further help improve the
506 health of infants, mothers, and people of all ages.

507 **5 AUTHOR CONTRIBUTION**

508 BK, AR, LB and NEL designed and performed the study and wrote the manuscript. ABB performed
509 preliminary analysis. AR performed modeling analyses. BK analyzed and interpreted the models
510 analyses. MAM and MWH provided samples. DC, JYY, JN, KM and LB performed expression,
511 purification of glycosyltransferases and kinetic assays. AR, BK, and NK performed literature surveys
512 to determine appropriate candidate genes for each reaction. BK and BB performed motif-level
513 analysis. BK and AWTC performed transcription factor analysis.

514 **6 ACKNOWLEDGMENTS**

515 Special thanks to Frederique Lisacek and Andrew McDonald for their input on navigating this
516 interdisciplinary topic. Additional thanks to Philip Spahn, Hooman Hefzi, Krystyna Kolodziej and Caressa
517 Robinson for help editing this manuscript. This work was supported by a Lilly Innovation Fellows Award
518 (A.R.), the Novo Nordisk Foundation provided to the Center for Biosustainability at the Technical
519 University of Denmark (NNF10CC1016517: N.E.L.), NIGMS (R35 GM119850: N.E.L., P41GM103390,
520 P01GM107012, and R01GM130915 to K.W.M.), NICHD (R21 HD080682: L.B.) and USDA (USDA/ARS 6250-
521 6001; M.W.H). This work is a publication of the U.S. Department of Agriculture/Agricultural Research
522 Service, Children's Nutrition Research Center, Department of Pediatrics, Baylor College of Medicine,
523 Houston, Texas. The contents of this publication do not necessarily reflect the views or policies of the U.S.
524 Department of Agriculture, nor does mention of trade names, commercial products, or organizations
525 imply endorsement from the U.S. government.

526 7 MATERIALS AND METHODS

527 7.1 MILK SAMPLE COLLECTION

528 Samples were collected following Institutional Review Board approval (Baylor College of Medicine,
529 Houston, TX). Lactating women 18-35 years of age with uncomplicated singleton pregnancy, vaginal
530 delivery at term (>37 weeks), Body Mass Index <26 kg/m² without diabetes, impaired glucose
531 tolerance, anemia, or renal or hepatic dysfunction were given informed consent before sample
532 collection. Description of the protocols used to collect milk samples and the diversity of subjects
533 present in both datasets. Cohort 1 consists of 8 samples for each of the 6 subjects (48 samples total)
534 including milk from 4 secretor mothers and 2 non-secretor mothers spanning from 6 hrs to 42 days
535 postpartum. Sample collection was previously described^{28,29}. Cohort 2 consists of 2 samples over each
536 of the 5 (10 samples total) including samples from 4 secretor mothers and 1 non-secretor mother
537 spanning 1 to 2 days postpartum. Sample collection was previously described³⁰.

538 7.2 ILLUMINA MRNA MICROARRAYS & GLYCOPROFILING

539 All expression and glycoprofiling measurements were sample-matched. Therefore, comparisons
540 across data-types occurred within each individual sample described in the previous section. Not all
541 samples in these studies have both microarray and glycoprofile measurements, only the samples
542 described in the previous section have matched glycomics and transcriptomics data.

543 mRNA was isolated from TRIzol-treated milk fat in each sample. Expression in cohort 1 was
544 measured using HumanHT-12 v4 Expression Beadchip microarrays (Illumina, Inc.) with ~44k
545 probes. Extraction of mRNA and measurement of expression in milk samples was performed as
546 previously described^{28,29}. Gene expression data for cohort 1 were retrieved from the Gene Expression
547 Omnibus at accession: GSE36936. Cohort 2 gene expression data were measured using a Human Ref-
548 8 BeadChip array (Illumina, Inc) with ~22k probes. Extraction of mRNA and related methods were
549 previously described³⁰. Expression data for cohort 1 can be accessed at accession: GSE12669. Both
550 microarrays were background corrected. The cohort 1 microarray was normalized using cubic spline
551 normalization and the cohort 2 microarray was normalized using the robust spline normalization.

552 As previously described^{41,91}, HMO composition and abundance data were collected using high-
553 performance liquid chromatography (HPLC) with 2-aminobenzamide (CID: 6942) derivatization and
554 a raffinose (CID:439242) standard. 16 HMOs were measured using retention time and commercial
555 standards including 2-fucosyllactose (2'FL), 3-fucosyllactose (3FL), 3-sialyllactose (3'SL), lacto-N-
556 tetraose (LNT), lacto-N-neotetraose (LNnT), lacto-N-fucopentaose (LNFP1, LNFP2 and LNFP3),
557 sialyl-LNT (LSTb and LSTc), difucosyl-LNT (DFLNT), disialyllacto-N-tetraose (DSLNT), fucosyl-lacto-
558 N-hexaose (FLNH), difucosyl-lacto-N-hexaose (DFLNH), fucosyl-disialyl-lacto-N-hexaose (FDSLNH)
559 and disialyl-lacto-N-hexaose (DSLNH). Technicians were blinded to sample metadata. HMO
560 composition and abundance measurement for cohort 1 were fully described in¹⁷. Measurements for
561 cohort 2 are previously unpublished and used the same methodology.

562 7.3 SOFTWARE

563 Modeling of HMO biosynthesis was performed in Matlab 2016b using the CobraToolbox⁹². All
564 analysis of biosynthetic models, interpretation and statistics were performed in R v3.5 and v3.6. In
565 R, we used *bigmemory*, *bigalgebra* and *biganalytics* to handle the millions of models and associated
566 statistics⁹³. We used *metap* for pooling p-values⁹⁴.

567 7.4 GENERATION AND SCORING OF GLYCOSYLATION NETWORK MODELS

568 Here we attempt to determine the genes responsible for making HMOs through the construction and
569 interrogation of models of their biosynthesis. Similar to the other biosynthetically constrained
570 glycomic models like the milk metaglycome²¹, Cartoonist⁹⁵ and several N-glycome simulations^{13,96–98},
571 we began with a set of elementary reactions. Enumerating all feasible permutations of the elementary
572 reaction (**Figure 3A**; S1.1.1), we delineated every possible reaction series from lactose to each of the
573 16 most abundant HMOs. Of the measured HMOs, 11 have fully determined molecular structures,
574 while the remaining five have multiple candidate structures (**Figure 1C**, **Figure S 12**)^{6,8,34,99–101}. The set
575 of all possible reactions leading to characterized and ambiguous structures formed the Complete
576 Network (**Figure 3B**; Supplemental Methods S1.1.1). Though non-lysosomal glycosidase^{102–104}
577 reactions are not explicitly specified, they are implicitly encoded in the flux. To reduce the Complete
578 Network to a more manageable size, we identified and removed all reactions that do not lead to
579 observed oligosaccharides using Flux Variability Analysis (FVA; Supplemental Methods S1.2.4;^{105–107}).
580 This trimming (**Figure 3C**; Supplemental Methods S1.1.2) defines the Reduced Network (**Figure 3D**;
581 Supplemental Methods S1.1.2). The Reduced Network describes many candidate models that can
582 uniquely simulate the HMO abundance collected through High-Performance Liquid Chromatography
583 (HPLC). A mixed integer linear programming (Supplemental Methods S1.2.5;^{108,109}) approach was
584 employed to extract candidate models from the Reduced Network capable of uniquely recapitulating
585 the HPLC data with minimal reactions (**Figure 3E**; Supplemental Methods S1.1.3). The reactions of
586 each candidate model were parameterized to determine the necessary flow of material (flux) through
587 each reaction to reproduce the measured oligosaccharide profiles (**Figure 3F**; Supplemental Methods
588 S1.1.3; S). The models were ranked by the consistency between the predicted flux and the expression
589 of genes believed to be associated with each reaction (**Figure 3G**; Supplemental Methods S1.1.4). This
590 consistency is evaluated by the Spearman correlation of changes in flux and gene expression across
591 subjects (**Figure 3H**; Supplemental Methods S1.1.4.1).

592 7.5 CANDIDATE MODEL RANKING, MODEL SELECTION AND SELECTION VALIDATION

593 Model scores, indicating the consistency between flux and gene expression (S1.1.4.1), were used to
594 rank candidate models (S1.1.4.2). The distribution of model scores computed from each dataset were
595 approximately normal, as evidenced by their linear Q-Q plots. This permitted the construction of a
596 background normal distribution of model scores (**Figure S 20**). We then selected high-performing
597 models, those with z-score normalized model scores greater than 1.646 (i.e., greater than the top 5%
598 of scores from a normal distribution) for further study. Model selection was performed on scores
599 computed independently for cohort 1 and cohort 2. Commonly high-performing models were those
600 that perform well in both cohort 1 and cohort 2. Hypergeometric enrichment was used to confirm

601 that the top cohort 1 and cohort 2 models significantly overlapped. (see Supplemental Methods
602 S1.1.4.2)

603 7.6 SUMMARY NETWORK EXTRACTION FROM THE REDUCED NETWORK

604 The Summary Network relates a heuristic selection of the most important reactions in the HMO
605 biosynthesis network as measured by proportion of inclusion in the commonly high-performing
606 models and enrichment in the commonly high-performing models relative to the background. Paths
607 drawn from observed HMOs to the root lactose were scored for their aggregate importance. The top
608 5% of paths leading to each observed HMO were retained to form the Summary Network
609 (Supplemental Methods see S1.1.4.3).

610 7.7 AMBIGUOUS GENE SELECTION

611 We aimed to match 10 elementary glycosyltransferase reactions to the supporting genes (**Table 1**).
612 Candidate genes were filtered from the relevant gene families to exclude gene products well known
613 to perform unrelated reactions (**Table 1**). Candidate genes were first evaluated for expression in
614 breast epithelium samples including microarrays in this study, independent RNA-Seq (GSE45669)³¹
615 and comparison to global expression distributions in GTEx³²; genes unmeasured by microarray in at least
616 75% of microarray samples (3rd Quartile, Q3) within each cohort were excluded unless they were
617 non-negligibly expressed in the independent RNA-Seq ($TPM_{Lemay} > 2$ or $TPM_{Lemay} > \text{Median}(TPM_{GTEx})$)
618 (see supplemental results, **Table S 1**, **Figure S 7**).

619 We used the model score definition, which quantifies how well the genes explain a model, i.e., if the
620 expression of the genes are best correlated to the normalized flux of the reaction (**Figure S 11**, S1.1.4)
621 they are proposed to support. We examined each gene contribution to the overall model score in
622 three ways to determine a consensus support score for each gene-reaction association (see S1.1.5.2).

623 The first metric we examined was the proportion (PROP) of commonly high-performing models best
624 explained by an isoform relative to the proportion of background models that select that same
625 isoform. The second metric was the average gene-linkage score (GLS) in high-performing models, i.e.,
626 the Spearman correlation between the normalized flux (**Figure S 11**, S1.1.4) and gene expression of
627 corresponding candidate genes. The gene-linkage score is a continuous measure of the consistency
628 between each gene with the flux it was proposed to support. Because it considers every gene, not just
629 the most flux-consistent gene, it is helpful for judging performance when the most flux-consistent
630 gene is more ambiguous. The third metric was the model-score contribution (MSC). MSC quantifies
631 the Pearson correlation between the gene-linkage score, the gene expression consistency with the
632 normalized flux, and the overall model score (i.e., the average correlation of all most-flux-consistent
633 genes). The model score indicates the frequency with which a gene is the most flux-consistent gene
634 normalized by its contribution relative to the other most flux-consistent genes in that model.

635 An aggregate reaction support score was constructed to describe performance within each individual
636 score (PROP, GLS, and MSC) and consistency across cohorts. To measure significance, the gene-
637 linkage score matrix (i.e., Spearman correlation between each candidate gene and the corresponding
638 normalized flux for each model) was shuffled ($n=27$) and all analyses rerun on each shuffle to
639 generate a permuted background distribution for PROP, GLS and MSC; shuffling of the GLS matrix
640 was done using a perfect minimal hash to remap all entries back to the GLS matrix in a random

641 order¹¹⁰. Performance within each independent cohort was described as the sum of z-scores for each
642 of three measures; z-score was calculated relative to the mean and standard deviations of these
643 scores in the permutation results. Consistency across cohorts was determined by pooling p-values
644 using the Fisher's log-sum method^{94,111}. The score presented in **Figure 5B** is the $-\log_{10}(\text{FDR}(\text{cohort-}$
645 $\text{pooled-p}))$.

646 **7.8** *IN VITRO* GLYCOSYLTRANSFERASE ACTIVITY ASSAYS

647 Recombinant forms of the respective glycosyltransferases were expressed and purified as previously
648 described¹¹². Enzyme activity was determined using the UDP-Glo™ or UMP/CMP-Glo™
649 Glycosyltransferase Assay (Promega) that determined UDP/CMP concentration formed as a by-
650 product of the glycosyltransferase reaction. Assays were performed according to the manufacturer's
651 instructions using reactions (10 μL) that consisted of a universal buffer containing 100 mM each of
652 MES, MOPS, and TRIS, pH 7.0, donor (1 mM UDP-GlcNAc (Promega) for B3GNT2; 1 mM UDP-Gal
653 (Promega) for B3GALT2; 0.2 mM CMP-SA (Nacalai USA Inc.) for ST3GAL1-6, ST6GALNAC2, and
654 ST6GALNAC5), 1 mM acceptor (lactose (Sigma) and lacto-N-neotetraose (LNnT) (Carbosynth) for
655 B3GNT2; lacto-N-tetraose (LNT, Bode lab) and pentasaccharide (GlcNAc-b1,3-Gal-b1,4-GlcNAc-b1,3-
656 Gal- b1,4-Glc, Boons lab, University of Georgia) for B3GALT2; LNnT, LNT, and Gal- β 1,3-GalNAc
657 (Carbosynth) for ST3GAL1-6; LNT for ST6GALNAC2 and ST6GALNAC5. The B3GNT2 and B3GALT2
658 assays also contained 1 mg/ml BSA and 5 mM MnCl_2 . Assays were carried out for 1 h (B3GNT2,
659 B3GALT2, ST6GALNAC2, and ST6GALNAC5) or 30 min (ST3GAL1-6) at 37 °C. Reactions (5 μL) were
660 stopped by mixing with an equal volume of Detection Reagent (5 μL) in white polystyrene, low-
661 volume, 384-well assay plates (Corning) and incubated for 60 min at room temperature. After
662 incubation, luminescence measurements were performed using a GloMax Multi Detection System
663 plate reader (Promega). The average luminescence was subtracted from the average luminescence
664 of respective blank to correct for background. Background and reaction measurements were
665 performed in triplicate.

666 **7.9** DIFFERENTIAL EXPRESSION (DE) ANALYSIS

667 The differential expression analysis was conducted on three different datasets: 1) 16 different HMOs
668 (2'FL, 3'SL, 3FL, FLNH, LNT, LNnT, LSTb, LNFP-III, LNFP-II, LNFP-I, DFLNT, LSTc, DSLNT, FDSLNH,
669 DSLNH, DFLNH), 2) 19 glycan motifs (X18, X32, X34, X35, X37, X40, X62, X63, X64, X65, X66, X94,
670 X106, X113, X120, X127, X141, X142, X143, see **Figure S 19**), and 3) 4 differential motifs for the
671 difference ("conversion rate") between related motifs (X65-X40, X106-X62, X63-X37, X62-X40, see
672 **Figure S 19**). Substructure abundance for glycan motifs and conversion ratios were computed using
673 Glycompare v1¹⁷. The gene expression data were downloaded from the Gene Expression Omnibus¹¹³
674 (GSE36936). Specifically, for each HMO, motif or differential motif, we used concentration (e.g., HMO-
675 3FL) as the predictor for gene expression in the differential expression analysis (e.g., "gene
676 expression \sim [3FL]"). The differential expression analysis was performed by fitting linear models
677 using empirical Bayes method as implemented in the *limma* v3.40.6 in R v3.6.1 package¹¹⁴ and p-
678 values were adjusted for multiple testing using Benjamini-Hochberg (BH) method¹¹⁵. In this way, we
679 determined gene-expression signatures indicative of each HMO and motif abundance.

680 7.10 INGENUITY PATHWAYS ANALYSIS (IPA) UPSTREAM REGULATOR

681 Differential expression signatures indicative of differential abundance in 16 HMOs, 19 motifs and 4
682 differential motifs were analyzed to predict upstream regulators using Ingenuity Pathway Analysis
683 (IPA, QIAGEN Inc.). Gene expression signatures indicative of HMO and motif abundance were defined
684 as genes differentially expressed with abundance in the previous *limma* analysis (FDR $q < 0.05$ and
685 $|\text{Fold Change}| > 1.5$).

686 7.11 DE NOVO TF BINDING SITE MOTIFS DISCOVERY AND KNOWN TF BINDING SITE IDENTIFICATION

687 We downloaded promoter sequences (file: “*upstream1000.fa.gz*”; version: GRCH38) from UCSC
688 Genome Browser public database (<https://genome.ucsc.edu/>) for the O-glycosyltransferase genes
689 used in this study (Table S 1). These promoter sequences included 1,000 bases upstream of annotated
690 transcription starts of RefSeq genes with annotated 5' UTRs. To conduct *de novo* TF binding site
691 motifs discovery, we first applied the motif discovery program MEME¹¹⁶ to identify candidate TF
692 binding site motifs on the downloaded promoter sequences with default parameters. The 10 TF
693 binding site motifs found by MEME were analyzed further for matches to known TF binding sites for
694 mammalian transcription factors in the motif databases, JASPAR Vertebrates¹¹⁷, via motif
695 comparison tool, TOMTOM¹¹⁸. The resulting discovered TF binding site motifs and their significantly
696 associated known TF binding sites (Table S 6, Table S 7) for mammalian transcription factors were
697 used further to compare with the IPA predicted upstream regulators.

698

699 8 REFERENCES

- 700 1. Edmond, K. M. *et al.* Delayed Breastfeeding Initiation Increases Risk of Neonatal Mortality.
701 *Pediatrics* **117**, e380–e386 (2006).
- 702 2. Bode, L. Human milk oligosaccharides: every baby needs a sugar mama. *Glycobiology* **22**, 1147–
703 1162 (2012).
- 704 3. Jantscher-Krenn, E. & Bode, L. Human milk oligosaccharides and their potential benefits for the
705 breast-fed neonate. *Minerva Pediatr.* **64**, 83–99 (2012).
- 706 4. Coppa, G. V. *et al.* Changes in Carbohydrate Composition in Human Milk Over 4 Months of
707 Lactation. *Pediatrics* **91**, (1993).
- 708 5. Picciano, M. F. Nutrient Composition of Human Milk. *Pediatr. Clin. North Am.* **48**, 53–67 (2001).
- 709 6. Bode, L. The functional biology of human milk oligosaccharides. *Early Hum. Dev.* **91**, 619–622

- 710 (2015).
- 711 7. Azad, M. B. *et al.* Human Milk Oligosaccharide Concentrations Are Associated with Multiple
712 Fixed and Modifiable Maternal Characteristics, Environmental Factors, and Feeding Practices. *J.*
713 *Nutr.* **148**, 1733–1742 (2018).
- 714 8. Kobata, A. Structures and application of oligosaccharides in human milk. *Proc. Jpn. Acad. Ser. B*
715 *Phys. Biol. Sci.* **86**, 731–747 (2010).
- 716 9. Etzold, S. & Bode, L. Glycan-dependent viral infection in infants and the role of human milk
717 oligosaccharides. *Curr. Opin. Virol.* **7**, 101–107 (2014).
- 718 10. Zhou, R. *et al.* Deficiency of intestinal α 1-2-fucosylation exacerbates ethanol-induced liver
719 disease in mice. *Alcohol. Clin. Exp. Res.* (2020) doi:10.1111/acer.14405.
- 720 11. Kellman, B. P. *et al.* A consensus-based and readable extension of Linear Code for Reaction
721 Rules (LiCoRR). *bioRxiv* 2020.05.31.126623 (2020) doi:10.1101/2020.05.31.126623.
- 722 12. Kobata, A. Possible application of milk oligosaccharides for drug development. *Chang Gung*
723 *Med. J.* **26**, 621–636 (2003).
- 724 13. Spahn, P. N. *et al.* A Markov chain model for N-linked protein glycosylation – towards a low-
725 parameter tool for model-driven glycoengineering. *Metabolic Engineering* vol. 33 52–66
726 (2016).
- 727 14. Liang, C. *et al.* A Markov model of glycosylation elucidates isozyme specificity and
728 glycosyltransferase interactions for glycoengineering. *Current Research in Biotechnology* **in**
729 **press**, (2020).
- 730 15. Liu, G. & Neelamegham, S. A computational framework for the automated construction of
731 glycosylation reaction networks. *PLoS One* **9**, e100939 (2014).
- 732 16. McDonald, A. G., Tipton, K. F. & Davey, G. P. A Knowledge-Based System for Display and
733 Prediction of O-Glycosylation Network Behaviour in Response to Enzyme Knockouts. *PLoS*
734 *Comput. Biol.* **12**, e1004844 (2016).

- 735 17. Bao, B. *et al.* Correcting for sparsity and non-independence in glycomic data through a systems
736 biology framework. *bioRxiv* 693507 (2019) doi:10.1101/693507.
- 737 18. Akune, Y. *et al.* Comprehensive analysis of the N-glycan biosynthetic pathway using
738 bioinformatics to generate UniCorn: A theoretical N-glycan structure database. *Carbohydr. Res.*
739 **431**, 56–63 (2016).
- 740 19. Lewis, N. E., Nagarajan, H. & Palsson, B. O. Constraining the metabolic genotype-phenotype
741 relationship using a phylogeny of in silico methods. *Nat. Rev. Microbiol.* **10**, 291–305 (2012).
- 742 20. Burgard, A. P., Vaidyaraman, S. & Maranas, C. D. Minimal reaction sets for Escherichia coli
743 metabolism under different growth requirements and uptake environments. *Biotechnol. Prog.*
744 **17**, 791–797 (2001).
- 745 21. Agravat, S. B., Song, X., Rojsajakul, T., Cummings, R. D. & Smith, D. F. Computational approaches
746 to define a human milk metaglycome. *Bioinformatics* **32**, 1471–1478 (2016).
- 747 22. Spahn, P. N., Hansen, A. H., Kol, S., Voldborg, B. G. & Lewis, N. E. Predictive glycoengineering of
748 biosimilars using a Markov chain glycosylation model. *Biotechnol. J.* **12**, (2017).
- 749 23. Nishihara, S. *et al.* Molecular genetic analysis of the human Lewis histo-blood group system. *J.*
750 *Biol. Chem.* **269**, 29271–29278 (1994).
- 751 24. Kudo, T. *et al.* Molecular genetic analysis of the human Lewis histo-blood group system. II.
752 Secretor gene inactivation by a novel single missense mutation A385T in Japanese nonsecretor
753 individuals. *J. Biol. Chem.* **271**, 9830–9837 (1996).
- 754 25. Koda, Y., Soejima, M., Liu, Y. & Kimura, H. Molecular basis for secretor type alpha(1,2)-
755 fucosyltransferase gene deficiency in a Japanese population: a fusion gene generated by
756 unequal crossover responsible for the enzyme deficiency. *Am. J. Hum. Genet.* **59**, 343–350
757 (1996).
- 758 26. Thurl, S., Henker, J., Siegel, M., Tovar, K. & Sawatzki, G. Detection of four human milk groups
759 with respect to Lewis blood group dependent oligosaccharides. *Glycoconj. J.* **14**, 795–799

- 760 (1997).
- 761 27. Stahl, B. *et al.* Detection of four human milk groups with respect to Lewis-blood-group-
762 dependent oligosaccharides by serologic and chromatographic analysis. *Adv. Exp. Med. Biol.*
763 **501**, 299–306 (2001).
- 764 28. Mohammad, M. A., Hadsell, D. L. & Haymond, M. W. Gene regulation of UDP-galactose synthesis
765 and transport: potential rate-limiting processes in initiation of milk production in humans. *Am.*
766 *J. Physiol. Endocrinol. Metab.* **303**, E365-76 (2012).
- 767 29. Mohammad, M. A. & Haymond, M. W. Regulation of lipid synthesis genes and milk fat
768 production in human mammary epithelial cells during secretory activation. *Am. J. Physiol.*
769 *Endocrinol. Metab.* **305**, E700-16 (2013).
- 770 30. Maningat, P. D. *et al.* Gene expression in the human mammary epithelium during lactation: the
771 milk fat globule transcriptome. *Physiol. Genomics* **37**, 12–22 (2009).
- 772 31. Lemay, D. G. *et al.* RNA sequencing of the human milk fat layer transcriptome reveals distinct
773 gene expression profiles at three stages of lactation. *PLoS One* **8**, e67531 (2013).
- 774 32. Carithers, L. J. *et al.* A novel approach to high-quality postmortem tissue procurement: the
775 GTEx project. *Biopreserv. Biobank.* **13**, 311–319 (2015).
- 776 33. Blank, D., Dotz, V., Geyer, R. & Kunz, C. Human milk oligosaccharides and Lewis blood group:
777 individual high-throughput sample profiling to enhance conclusions from functional studies.
778 *Adv. Nutr.* **3**, 440S–9S (2012).
- 779 34. Wu, S., Tao, N., German, J. B., Grimm, R. & Lebrilla, C. B. Development of an annotated library of
780 neutral human milk oligosaccharides. *J. Proteome Res.* **9**, 4138–4151 (2010).
- 781 35. Brodbeck, U. & Ebner, K. E. Resolution of a soluble lactose synthetase into two protein
782 components and solubilization of microsomal lactose synthetase. *J. Biol. Chem.* **241**, 762–764
783 (1966).
- 784 36. Nakhasi, H. L. & Quasba, P. K. Quantitation of milk proteins and their mRNAs in rat mammary

- 785 gland at various stages of gestation and lactation. *J. Biol. Chem.* **254**, 6016–6025 (1979).
- 786 37. Morrow, A. L. *et al.* Fucosyltransferase 2 non-secretor and low secretor status predicts severe
787 outcomes in premature infants. *J. Pediatr.* **158**, 745–751 (2011).
- 788 38. Autran, C. A. *et al.* Human milk oligosaccharide composition predicts risk of necrotising
789 enterocolitis in preterm infants. *Gut* **67**, 1064–1070 (2018).
- 790 39. Morrow, A. L. *et al.* Human Milk Oligosaccharide Blood Group Epitopes and Innate Immune
791 Protection against Campylobacter and Calicivirus Diarrhea in Breastfed Infants. in 443–446
792 (Springer, Boston, MA, 2004).
- 793 40. Yu, Z.-T., Nanda Nanthakumar, N. & Newburg, D. S. The Human Milk Oligosaccharide 2'-
794 Fucosyllactose Quenches Campylobacter jejuni-Induced Inflammation in Human Epithelial
795 Cells HEP-2 and HT-29 and in Mouse Intestinal Mucosa. *The Journal of Nutrition* vol. 146 1980–
796 1990 (2016).
- 797 41. Alderete, T. L. *et al.* Associations between human milk oligosaccharides and infant body
798 composition in the first 6 mo of life. *Am. J. Clin. Nutr.* **102**, 1381–1388 (2015).
- 799 42. Uwaezuoke, S. N., Eneh, C. I. & Ndu, I. K. Relationship Between Exclusive Breastfeeding and
800 Lower Risk of Childhood Obesity: A Narrative Review of Published Evidence. *Clin. Med. Insights*
801 *Pediatr.* **11**, 1179556517690196 (2017).
- 802 43. Uwaezuoke, S. *et al.* Maternal diet during exclusive breastfeeding can predict food preference
803 in preschoolers: A cross-sectional study of mother- child dyads in Enugu, south-east Nigeria.
804 *Int. J. Child Health Nutr.* **6**, 70–79 (2017).
- 805 44. Moro, G. *et al.* Dosage-related bifidogenic effects of galacto- and fructooligosaccharides in
806 formula-fed term infants. *J. Pediatr. Gastroenterol. Nutr.* **34**, 291–295 (2002).
- 807 45. Costalos, C., Kapiki, A., Apostolou, M. & Papathoma, E. The effect of a prebiotic supplemented
808 formula on growth and stool microbiology of term infants. *Early Hum. Dev.* **84**, 45–49 (2008).
- 809 46. Vos, A. P. *et al.* A specific prebiotic oligosaccharide mixture stimulates delayed-type

- 810 hypersensitivity in a murine influenza vaccination model. *Int. Immunopharmacol.* **6**, 1277–
811 1286 (2006).
- 812 47. Viverge, D., Grimmonprez, L., Cassanas, G., Bardet, L. & Solere, M. Discriminant carbohydrate
813 components of human milk according to donor secretor types. *J. Pediatr. Gastroenterol. Nutr.*
814 **11**, 365–370 (1990).
- 815 48. McGuire, M. K. *et al.* What's normal? Oligosaccharide concentrations and profiles in milk
816 produced by healthy women vary geographically. *Am. J. Clin. Nutr.* **105**, 1086–1100 (2017).
- 817 49. Furuike, T., Yamada, K., Ohta, T., Monde, K. & Nishimura, S.-I. An efficient synthesis of a
818 biantennary sialooligosaccharide analog using a 1,6-anhydro- β -lactose derivative as a key
819 synthetic block. *Tetrahedron* **59**, 5105–5113 (2003).
- 820 50. Fair, R. J., Hahm, H. S. & Seeberger, P. H. Combination of automated solid-phase and enzymatic
821 oligosaccharide synthesis provides access to α (2,3)-sialylated glycans. *Chem. Commun.* **51**,
822 6183–6185 (2015).
- 823 51. Yao, W., Yan, J., Chen, X., Wang, F. & Cao, H. Chemoenzymatic synthesis of lacto-N-
824 tetrasaccharide and sialyl lacto-N-tetrasaccharides. *Carbohydr. Res.* **401**, 5–10 (2015).
- 825 52. Prudden, A. R. *et al.* Synthesis of asymmetrical multiantennary human milk oligosaccharides.
826 *Proc. Natl. Acad. Sci. U. S. A.* **114**, 6954–6959 (2017).
- 827 53. Prudden, A. R., Chinoy, Z. S., Wolfert, M. A. & Boons, G.-J. A multifunctional anomeric linker for
828 the chemoenzymatic synthesis of complex oligosaccharides. *Chem. Commun.* **50**, 7132–7135
829 (2014).
- 830 54. Bode, L. *et al.* Overcoming the limited availability of human milk oligosaccharides: challenges
831 and opportunities for research and application. *Nutr. Rev.* **74**, 635–644 (2016).
- 832 55. Guan, N. & Chen, R. Recent Technology Development for the Biosynthesis of Human Milk
833 Oligosaccharide. *Recent patents on biotechnology* (2018).
- 834 56. Lee, W.-H. *et al.* Whole cell biosynthesis of a functional oligosaccharide, 2\textasciicute-

- 835 fucosyllactose, using engineered Escherichia coli. *Microb. Cell Fact.* **11**, 48 (2012).
- 836 57. Chin, Y.-W., Kim, J.-Y., Lee, W.-H. & Seo, J.-H. Enhanced production of 2'-fucosyllactose in
837 engineered Escherichia coli BL21star(DE3) by modulation of lactose metabolism and
838 fucosyltransferase. *J. Biotechnol.* **210**, 107–115 (2015).
- 839 58. Baumgärtner, F., Seitz, L., Sprenger, G. A. & Albermann, C. Construction of Escherichia coli
840 strains with chromosomally integrated expression cassettes for the synthesis of
841 2'-fucosyllactose. *Microb. Cell Fact.* **12**, 40 (2013).
- 842 59. Baumgärtner, F., Conrad, J., Sprenger, G. A. & Albermann, C. Synthesis of the Human Milk
843 Oligosaccharide Lacto-N-Tetraose in Metabolically Engineered, Plasmid-Free E. coli.
844 *Chembiochem* **15**, 1896–1900 (2014).
- 845 60. Kumazaki, T. & Yoshida, A. Biochemical evidence that secretor gene, Se, is a structural gene
846 encoding a specific fucosyltransferase. *Proc. Natl. Acad. Sci. U. S. A.* **81**, 4193–4197 (1984).
- 847 61. Mollicone, R. *et al.* Activity, Splice Variants, Conserved Peptide Motifs, and Phylogeny of Two
848 New α 1,3-Fucosyltransferase Families (FUT10 and FUT11). *J. Biol. Chem.* **284**, 4723–4738
849 (2009).
- 850 62. Kaneko, M. *et al.* Assignment¹ of the human α 1,3-fucosyltransferase IX gene (FUT9) to
851 chromosome band 6q16 by in situ hybridization. *Cytogenetic and Genome Research* vol. 86
852 329–330 (1999).
- 853 63. Nishihara, S. *et al.* α 1, 3-Fucosyltransferase 9 (FUT9; Fuc-TIX) preferentially fucosylates the
854 distal GlcNAc residue of polylactosamine chain while the other four α 1, 3FUT members
855 preferentially fucosylate the inner GlcNAc residue. *FEBS Lett.* **462**, 289–294 (1999).
- 856 64. Niemelä, R. *et al.* Complementary acceptor and site specificities of Fuc-TIV and Fuc-TVII allow
857 effective biosynthesis of sialyl-TriLex and related polylactosamines present on glycoprotein
858 counterreceptors of selectins. *J. Biol. Chem.* **273**, 4021–4026 (1998).
- 859 65. Mondal, N. *et al.* Distinct human α (1,3)-fucosyltransferases drive Lewis-X/sialyl Lewis-X

- 860 assembly in human cells Downloaded from. (2018) doi:10.1074/jbc.RA117.000775.
- 861 66. Kurosawa, N., Inoue, M., Yoshida, Y. & Tsuji, S. Molecular Cloning and Genomic Analysis of
862 Mouse Gal β 1,3GalNAc-specific GalNAc α 2,6-Sialyltransferase. *Journal of Biological Chemistry*
863 vol. 271 15109–15116 (1996).
- 864 67. Kurosawa, N., Kojima, N., Inoue, M., Hamamoto, T. & Tsuji, S. Cloning and expression of Gal beta
865 1,3GalNAc-specific GalNAc alpha 2,6-sialyltransferase. *J. Biol. Chem.* **269**, 19048–19053 (1994).
- 866 68. Okajima, T. *et al.* Molecular Cloning of Brain-specific GD1 α Synthase (ST6GalNAc V) Containing
867 CAG/Glutamine Repeats. *J. Biol. Chem.* **274**, 30557–30562 (1999).
- 868 69. Okajima, T. *et al.* Expression cloning of human globoside synthase cDNAs. Identification of beta
869 3Gal-T3 as UDP-N-acetylgalactosamine:globotriaosylceramide beta 1,3-N-
870 acetylgalactosaminyltransferase. *J. Biol. Chem.* **275**, 40498–40503 (2000).
- 871 70. Sjoberg, E. R., Kitagawa, H., Glushka, J., van Halbeek, H. & Paulson, J. C. Molecular Cloning of a
872 Developmentally Regulated N-Acetylgalactosamine 2,6-Sialyltransferase Specific for Sialylated
873 Glycoconjugates. *J. Biol. Chem.* **271**, 7450–7459 (1996).
- 874 71. Tsuchida, A. *et al.* Molecular cloning and expression of human ST6GalNAc III: restricted tissue
875 distribution and substrate specificity. *J. Biochem.* **138**, 237–243 (2005).
- 876 72. Lee, Y.-C. *et al.* Molecular Cloning and Functional Expression of Two Members of Mouse
877 NeuAc α 2,3Gal β 1,3GalNAc GalNAc α 2,6-Sialyltransferase Family, ST6GalNAc III and IV. *Journal*
878 *of Biological Chemistry* vol. 274 11958–11967 (1999).
- 879 73. Tsuchida, A. *et al.* Synthesis of Disialyl Lewis a (Lea) Structure in Colon Cancer Cell Lines by a
880 Sialyltransferase, ST6GalNAc VI, Responsible for the Synthesis of α -Series Gangliosides. *J. Biol.*
881 *Chem.* **278**, 22787–22794 (2003).
- 882 74. Kitagawa, H. & Paulson, J. C. Cloning of a novel alpha 2,3-sialyltransferase that sialylates
883 glycoprotein and glycolipid carbohydrate groups. *J. Biol. Chem.* **269**, 1394–1401 (1994).
- 884 75. Kono, M. *et al.* Mouse beta-galactoside alpha 2,3-sialyltransferases: comparison of in vitro

- 885 substrate specificities and tissue specific expression. *Glycobiology* **7**, 469–479 (1997).
- 886 76. Blixt, O. *et al.* Glycan microarrays for screening sialyltransferase specificities. *Glycoconj. J.* **25**,
887 59–68 (2008).
- 888 77. Weinstein, J., de Souza-e-Silva, U. & Paulson, J. C. Sialylation of glycoprotein oligosaccharides N-
889 linked to asparagine. Enzymatic characterization of a Gal beta 1 to 3(4)GlcNAc alpha 2 to 3
890 sialyltransferase and a Gal beta 1 to 4GlcNAc alpha 2 to 6 sialyltransferase from rat liver. *J. Biol.*
891 *Chem.* **257**, 13845–13853 (1982).
- 892 78. Gillespie, W., Kelm, S. & Paulson, J. C. Cloning and expression of the Gal beta 1, 3GalNAc alpha
893 2,3-sialyltransferase. *J. Biol. Chem.* **267**, 21004–21010 (1992).
- 894 79. Miyazaki, H. *et al.* Expression Cloning of Rat cDNA Encoding UDP-galactose:GD2 β 1,3-
895 galactosyltransferase That Determines the Expression of GD1b/GM1/GA1. *J. Biol. Chem.* **272**,
896 24794–24799 (1997).
- 897 80. Amado, M. *et al.* A family of human beta3-galactosyltransferases. Characterization of four
898 members of a UDP-galactose:beta-N-acetyl-glucosamine/beta-nacetyl-galactosamine beta-1,3-
899 galactosyltransferase family. *J. Biol. Chem.* **273**, 12770–12778 (1998).
- 900 81. Isshiki, S. *et al.* Cloning, Expression, and Characterization of a Novel UDP-galactose: β -N-
901 Acetylglucosamine β 1,3-Galactosyltransferase (β 3Gal-T5) Responsible for Synthesis of Type 1
902 Chain in Colorectal and Pancreatic Epithelia and Tumor Cells Derived Therefrom. *Journal of*
903 *Biological Chemistry* vol. 274 12499–12507 (1999).
- 904 82. Schwientek, T. *et al.* Cloning of a Novel Member of the UDP-Galactose: β -N-Acetylglucosamine
905 β 1,4-Galactosyltransferase Family, β 4Gal-T4, Involved in Glycosphingolipid Biosynthesis. *J.*
906 *Biol. Chem.* **273**, 29331–29340 (1998).
- 907 83. Sato, T., Aoki, N., Matsuda, T. & Furukawa, K. Differential effect of alpha-lactalbumin on beta-
908 1,4-galactosyltransferase IV activities. *Biochem. Biophys. Res. Commun.* **244**, 637–641 (1998).
- 909 84. Shiraishi, N. *et al.* Identification and Characterization of Three Novel β 1,3-N-

- 910 Acetylglucosaminyltransferases Structurally Related to the β 1,3-Galactosyltransferase Family.
911 *J. Biol. Chem.* **276**, 3498–3507 (2001).
- 912 85. Chen, G. Y., Kurosawa, N. & Muramatsu, T. A novel variant form of murine beta-1, 6-N-
913 acetylglucosaminyltransferase forming branches in poly-N-acetyllactosamines. *Glycobiology*
914 **10**, 1001–1011 (2000).
- 915 86. Bierhuizen, M. F. & Fukuda, M. Expression cloning of a cDNA encoding UDP-GlcNAc: Gal beta 1-
916 3-GalNAc-R (GlcNAc to GalNAc) beta 1-6GlcNAc transferase by gene transfer into CHO cells
917 expressing polyoma large tumor antigen. *Proceedings of the National Academy of Sciences* **89**,
918 9326–9330 (1992).
- 919 87. Schwientek, T. *et al.* Control of O-glycan branch formation. Molecular cloning of human cDNA
920 encoding a novel beta1,6-N-acetylglucosaminyltransferase forming core 2 and core 4. *J. Biol.*
921 *Chem.* **274**, 4504–4512 (1999).
- 922 88. Yeh, J. C., Ong, E. & Fukuda, M. Molecular cloning and expression of a novel beta-1, 6-N-
923 acetylglucosaminyltransferase that forms core 2, core 4, and I branches. *J. Biol. Chem.* **274**,
924 3215–3221 (1999).
- 925 89. Schwientek, T. *et al.* Control of O-Glycan Branch Formation: MOLECULAR CLONING OF HUMAN
926 cDNA ENCODING A NOVEL β 1,6-N-ACETYLGLUCOSAMINYLTRANSFERASE FORMING CORE 2
927 AND CORE 4. *J. Biol. Chem.* **274**, 4504–4512 (1999).
- 928 90. Ujita, M., Misra, A. K., McAuliffe, J., Hindsgaul, O. & Fukuda, M. Poly-N-acetyllactosamine
929 Extension in N-Glycans and Core 2- and Core 4-branched O-Glycans Is Differentially Controlled
930 by i-Extension Enzyme and Different Members of the β 1, 4-Galactosyltransferase Gene Family.
931 *J. Biol. Chem.* **275**, 15868–15875 (2000).
- 932 91. Bode, L. *et al.* Human milk oligosaccharide concentration and risk of postnatal transmission of
933 HIV through breastfeeding. *Am. J. Clin. Nutr.* **96**, 831–839 (2012).
- 934 92. Heirendt, L. *et al.* Creation and analysis of biochemical constraint-based models using the

- 935 COBRA Toolbox v.3.0. *Nat. Protoc.* **14**, 639–702 (2019).
- 936 93. Kane, M. J., Emerson, J. W., Haverty, P. & Others. bigmemory: Manage massive matrices with
937 shared memory and memory-mapped files. *R package version 4*, (2010).
- 938 94. Dewey, M. metap: Meta-analysis of significance values. R package version 0.7. (2016).
- 939 95. Goldberg, D., Sutton-Smith, M., Paulson, J. & Dell, A. Automatic annotation of matrix-assisted
940 laser desorption/ionization N-glycan spectra. *Proteomics* **5**, 865–875 (2005).
- 941 96. Hossler, P., Mulukutla, B. C. & Hu, W.-S. Systems analysis of N-glycan processing in mammalian
942 cells. *PLoS One* **2**, e713 (2007).
- 943 97. Krambeck, F. J. *et al.* A mathematical model to derive N-glycan structures and cellular enzyme
944 activities from mass spectrometric data. *Glycobiology* **19**, 1163–1175 (2009).
- 945 98. McDonald, A. G. *et al.* Galactosyltransferase 4 is a major control point for glycan branching in N-
946 linked glycosylation. *J. Cell Sci.* **127**, 5014–5026 (2014).
- 947 99. Mantovani, V., Galeotti, F., Maccari, F. & Volpi, N. Recent advances on separation and
948 characterization of human milk oligosaccharides. *Electrophoresis* **37**, 1514–1524 (2016).
- 949 100. Ninonuevo, M. R. *et al.* A strategy for annotating the human milk glycome. *J. Agric. Food Chem.*
950 **54**, 7471–7480 (2006).
- 951 101. Wu, S., Grimm, R., German, J. B. & Lebrilla, C. B. Annotation and structural analysis of sialylated
952 human milk oligosaccharides. *J. Proteome Res.* **10**, 856–868 (2011).
- 953 102. Wiederschain, G. Y. & Newburg, D. S. Glycoconjugate stability in human milk: glycosidase
954 activities and sugar release. *J. Nutr. Biochem.* **12**, 559–564 (2001).
- 955 103. Miura, K., Hakamata, W., Tanaka, A., Hirano, T. & Nishio, T. Discovery of human Golgi β -
956 galactosidase with no identified glycosidase using a QMC substrate design platform for exo-
957 glycosidase. *Bioorg. Med. Chem.* **24**, 1369–1375 (2016).
- 958 104. Dudzik, D. *et al.* Activity of N-acetyl- β -D-hexosaminidase (HEX) and its isoenzymes A and B in
959 human milk during the first 3 months of breastfeeding. *Adv. Med. Sci.* **53**, (2008).

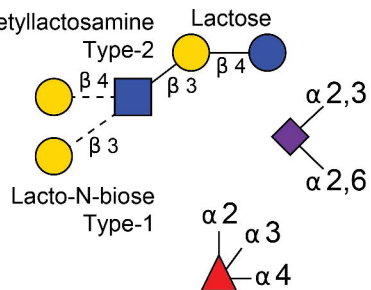
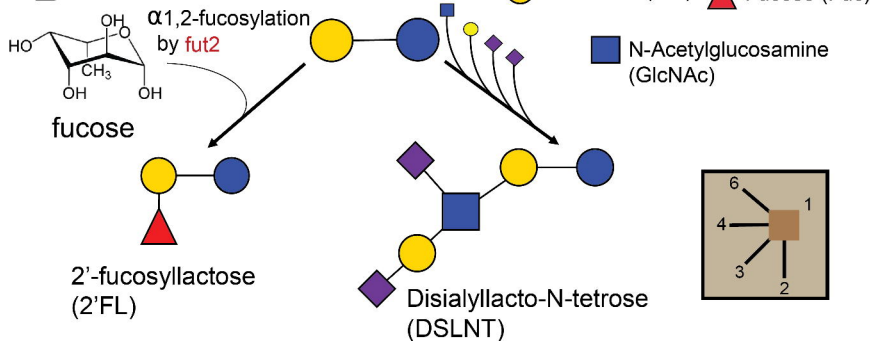
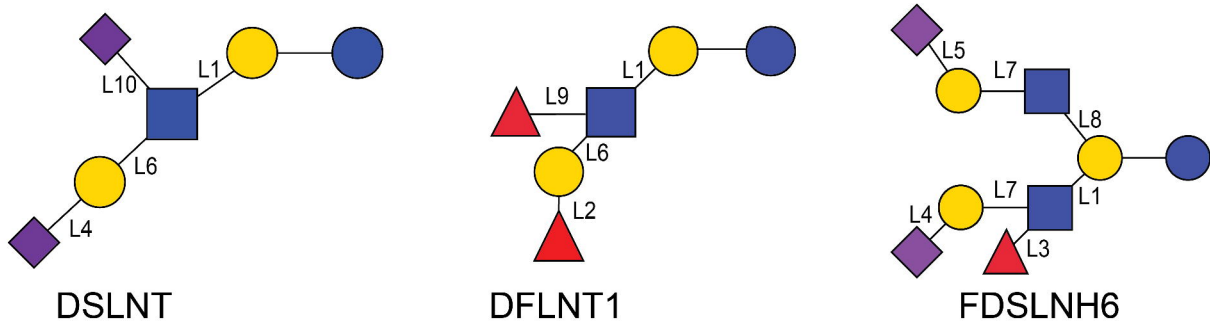
- 960 105. Orth, J. D., Thiele, I. & Palsson, B. Ø. What is flux balance analysis? *Nat. Biotechnol.* **28**, 245–248
961 (2010).
- 962 106. Mahadevan, R. & Schilling, C. H. The effects of alternate optimal solutions in constraint-based
963 genome-scale metabolic models. *Metab. Eng.* **5**, 264–276 (2003).
- 964 107. Gudmundsson, S. & Thiele, I. Computationally efficient flux variability analysis. *BMC*
965 *Bioinformatics* **11**, 489 (2010).
- 966 108. Reed, J. L. & Palsson, B. Ø. Genome-scale in silico models of E. coli have multiple equivalent
967 phenotypic states: assessment of correlated reaction subsets that comprise network states.
968 *Genome Res.* **14**, 1797–1805 (2004).
- 969 109. Lee, S., Phalakornkule, C., Domach, M. M. & Grossmann, I. E. Recursive MILP model for finding
970 all the alternate optima in LP models for metabolic networks. *Comput. Chem. Eng.* **24**, 711–716
971 (2000).
- 972 110. Fredman, M. L., Komlós, J. & Szemerédi, E. Storing a Sparse Table with 0 (1) Worst Case Access
973 Time. *Journal of the ACM (JACM)* vol. 31 538–544 (1984).
- 974 111. E., W. P. & W., P. E. Statistical Methods for Research Workers. By Fisher R. A. [Pp. 239 ix VI
975 Tables. Edinburgh and London: Oliver & Boyd. 1925. Price 15s.]The Fundamentals of Statistics.
976 By Thurstone L. L.. [Pp. 237 xvi. New York: The Macmillan Company. 1925. Price 8s. 6d.].
977 *Journal of the Institute of Actuaries* vol. 56 326–327 (1925).
- 978 112. Moremen, K. W. *et al.* Expression system for structural and functional studies of human
979 glycosylation enzymes. *Nat. Chem. Biol.* **14**, 156–162 (2018).
- 980 113. Edgar, R., Domrachev, M. & Lash, A. E. Gene Expression Omnibus: NCBI gene expression and
981 hybridization array data repository. *Nucleic Acids Res.* **30**, 207–210 (2002).
- 982 114. Smyth, G. K., Thorne, N. P. & Wettenhall, J. LIMMA: Linear Models for Microarray Data Version
983 1.6. 6. *User's Guide* (2004).
- 984 115. Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and Powerful

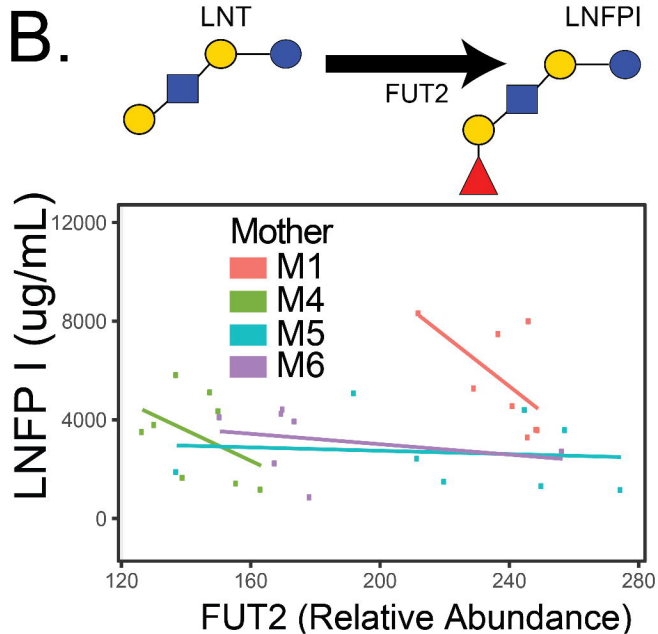
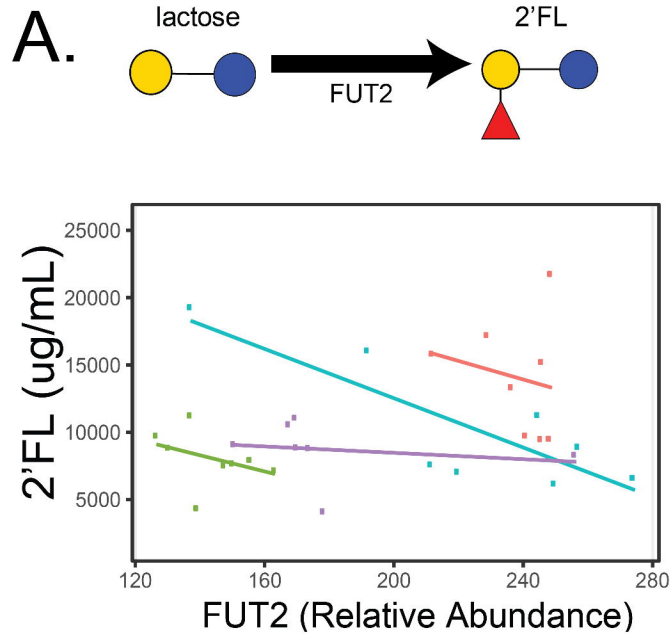
- 985 Approach to Multiple Testing. *J. R. Stat. Soc. Series B Stat. Methodol.* **57**, 289–300 (1995).
- 986 116. Bailey, T. L. & Elkan, C. The value of prior knowledge in discovering motifs with MEME. *Proc.*
987 *Int. Conf. Intell. Syst. Mol. Biol.* **3**, 21–29 (1995).
- 988 117. Fornes, O. *et al.* JASPAR 2020: update of the open-access database of transcription factor
989 binding profiles. *Nucleic Acids Res.* **48**, D87–D92 (2020).
- 990 118. Gupta, S., Stamatoyannopoulos, J. A., Bailey, T. L. & Noble, W. S. Quantifying similarity between
991 motifs. *Genome Biol.* **8**, R24 (2007).
- 992 119. Taniguchi, N., Honke, K. & Fukuda, M. *Handbook of Glycosyltransferases and Related Genes.*
993 (Springer Science & Business Media, 2011).
- 994 120. Narimatsu, H. Construction of a human glycogene library and comprehensive functional
995 analysis. *Glycoconj. J.* **21**, 17–24 (2004).
- 996 121. Schomburg, I. *et al.* The BRENDA enzyme information system-From a database to an expert
997 system. *J. Biotechnol.* **261**, 194–206 (2017).
- 998 122. Magrane, M. & UniProt Consortium. UniProt Knowledgebase: a hub of integrated protein data.
999 *Database* **2011**, bar009 (2011).
- 1000 123. Caspi, R. *et al.* The MetaCyc database of metabolic pathways and enzymes and the BioCyc
1001 collection of Pathway/Genome Databases. *Nucleic Acids Res.* **42**, D459–71 (2014).
- 1002 124. Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. & Morishima, K. KEGG: new perspectives on
1003 genomes, pathways, diseases and drugs. *Nucleic Acids Res.* **45**, D353–D361 (2017).
- 1004 125. Praissman, J. L. *et al.* B4GAT1 is the priming enzyme for the LARGE-dependent functional
1005 glycosylation of α -dystroglycan. *Elife* **3**, (2014).
- 1006 126. Willer, T. *et al.* The glucuronyltransferase B4GAT1 is required for initiation of LARGE-mediated
1007 α -dystroglycan functional glycosylation. *Elife* **3**, (2014).
- 1008 127. Yanagidani, S. *et al.* Purification and cDNA Cloning of GDP-L-Fuc:N-Acetyl- β -D-
1009 Glucosaminide: α 1-6 Fucosyltransferase (α 1-6 FucT) from Human Gastric Cancer MKN45 Cells.

- 1010 *J. Biochem.* **121**, 626–632 (1997).
- 1011 128. Uozumi, N. *et al.* Purification and cDNA cloning of porcine brain GDP-L-Fuc: N-acetyl- β -D-
1012 glucosaminide α 1 \rightarrow 6fucosyltransferase. *J. Biol. Chem.* **271**, 27810–27817 (1996).
- 1013 129. Kataoka, K. & Huh, N.-H. A novel beta1,3-N-acetylglucosaminyltransferase involved in invasion
1014 of cancer cells as assayed in vitro. *Biochem. Biophys. Res. Commun.* **294**, 843–848 (2002).
- 1015 130. Kitayama, K., Hayashida, Y., Nishida, K. & Akama, T. O. Enzymes responsible for synthesis of
1016 corneal keratan sulfate glycosaminoglycans. *J. Biol. Chem.* **282**, 30085–30096 (2007).
- 1017 131. Seko, A. & Yamashita, K. β 1,3-N-Acetylglucosaminyltransferase-7 (β 3Gn-T7) acts efficiently on
1018 keratan sulfate-related glycans. *FEBS Lett.* **556**, 216–220 (2004).
- 1019 132. Bai, X. *et al.* Biosynthesis of the linkage region of Glycosaminoglycans cloning and activity of
1020 galactosyltransferase ii, the sixth member of the β 1, 3-galactosyltransferase family (β 3GalT6).
1021 *J. Biol. Chem.* **276**, 48189–48195 (2001).
- 1022 133. Ju, T., Brewer, K., D’Souza, A., Cummings, R. D. & Canfield, W. M. Cloning and Expression of
1023 Human Core 1 β 1,3-Galactosyltransferase. *J. Biol. Chem.* **277**, 178–186 (2002).
- 1024 134. Almeida, R. *et al.* Cloning and expression of a proteoglycan UDP-galactose: β -Xylose β 1, 4-
1025 galactosyltransferase IA seventh member of the human β 4-galactosyltransferase gene family. *J.*
1026 *Biol. Chem.* **274**, 26165–26171 (1999).
- 1027 135. Okajima, T., Yoshida, K., Kondo, T. & Furukawa, K. Human homolog of *Caenorhabditis elegans*
1028 sqv-3 gene is galactosyltransferase I involved in the biosynthesis of the glycosaminoglycan-
1029 protein linkage region of proteoglycans. *J. Biol. Chem.* **274**, 22915–22918 (1999).
- 1030 136. Inshaw, J. R. J., Cutler, A. J., Burren, O. S., Stefana, M. I. & Todd, J. A. Approaches and advances in
1031 the genetic causes of autoimmune disease and their implications. *Nat. Immunol.* **19**, 674–684
1032 (2018).
- 1033 137. Rouillard, A. D. *et al.* The harmonizome: a collection of processed datasets gathered to serve
1034 and mine knowledge about genes and proteins. *Database* **2016**, (2016).

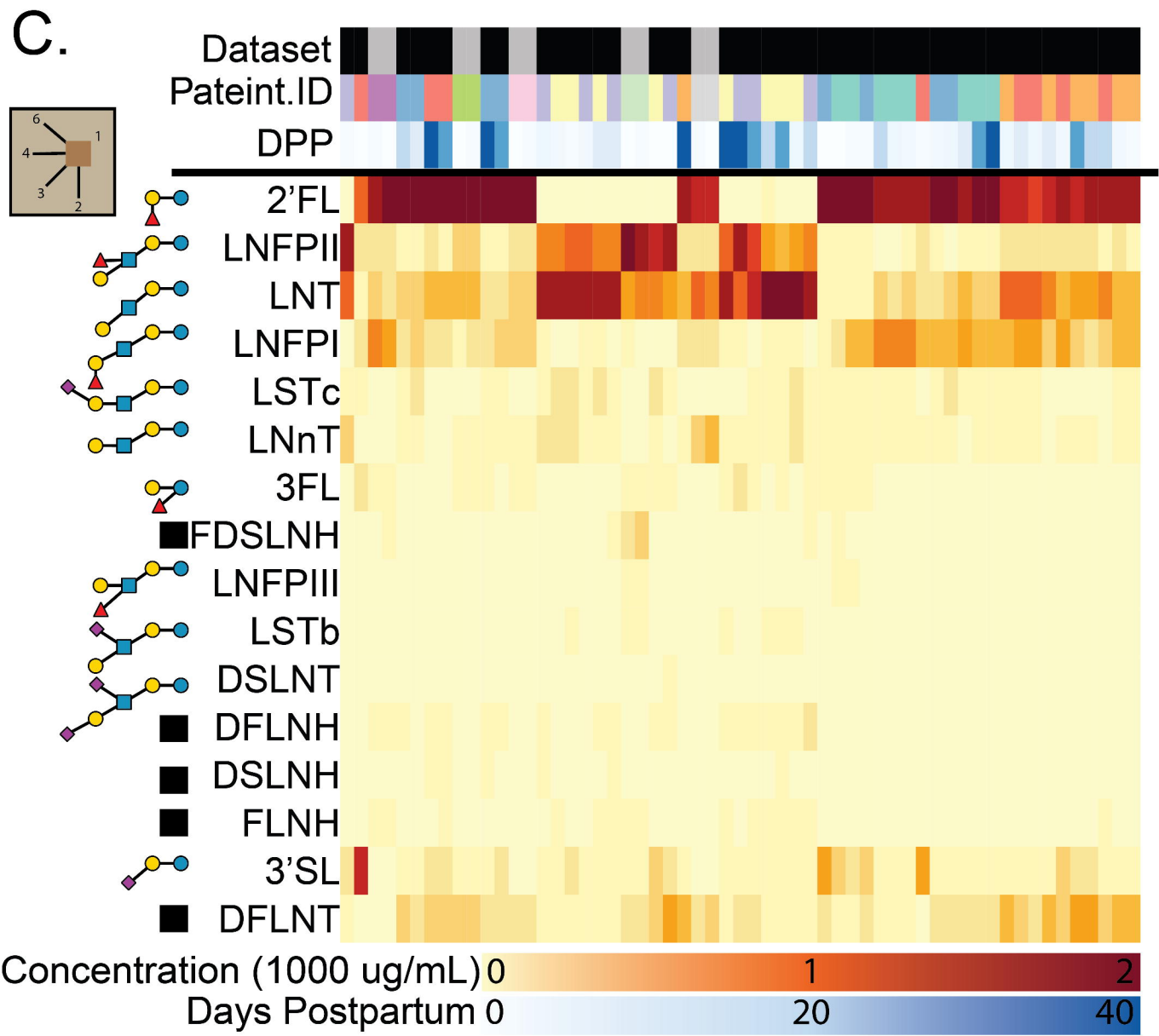
- 1035 138. Hu, Z.-Z., Zhuang, L., Meng, J. & Dufau, M. L. Transcriptional Regulation of the Generic Promoter
1036 III of the Rat Prolactin Receptor Gene by C/EBP β and Sp1. *J. Biol. Chem.* **273**, 26225–26235
1037 (1998).
- 1038 139. Sugiyama, A., Fukushima, N. & Sato, T. Transcriptional Mechanism of the β 4-
1039 Galactosyltransferase 4 Gene in SW480 Human Colon Cancer Cell Line. *Biol. Pharm. Bull.* **40**,
1040 733–737 (2017).
- 1041 140. Sato, T. & Furukawa, K. Transcriptional Regulation of the Human β -1,4-Galactosyltransferase V
1042 Gene in Cancer Cells: ESSENTIAL ROLE OF TRANSCRIPTION FACTOR Sp1. *J. Biol. Chem.* **279**,
1043 39574–39583 (2004).
- 1044 141. Sato, T. & Furukawa, K. Sequential Action of Ets-1 and Sp1 in the Activation of the Human β -
1045 1,4-Galactosyltransferase V Gene Involved in Abnormal Glycosylation Characteristic of Cancer
1046 Cells. *J. Biol. Chem.* **282**, 27702–27712 (2007).
- 1047 142. Zhou, L., Jiang, J. & Gu, J. β 1,4-Galactosyltransferase: Regulation and Signaling in Cancers.
1048 *Glycoscience: Biology and Medicine* 1141–1148 (2015) doi:10.1007/978-4-431-54841-6_74.
- 1049 143. Kurcon, T. *et al.* miRNA proxy approach reveals hidden functions of glycosylation. *Proc. Natl.*
1050 *Acad. Sci. U. S. A.* **112**, 7327–7332 (2015).
- 1051 144. Liu, B. *et al.* MiR-29b/Sp1/FUT4 axis modulates the malignancy of leukemia stem cells by
1052 regulating fucosylation via Wnt/ β -catenin pathway in acute myeloid leukemia. *J. Exp. Clin.*
1053 *Cancer Res.* **38**, 200 (2019).
- 1054 145. Dhordain, P., Dewitte, F., Desbiens, X., Stehelin, D. & Duterque-Coquillaud, M. Mesodermal
1055 expression of the chicken erg gene associated with precartilaginous condensation and cartilage
1056 differentiation. *Mech. Dev.* **50**, 17–28 (1995).
- 1057 146. Taniguchi, A., Itaru, Y. & Matsumoto, K. Genomic structure and transcriptional regulation of
1058 human Gal β 1,3GalNAc α 2,3-sialyltransferase (hST3Gal I) gene. *Glycobiology* **11**, 241–247
1059 (2001).

- 1060 147. Vandenplas, Y. *et al.* Human Milk Oligosaccharides: 2'-Fucosyllactose (2'-FL) and Lacto-N-
- 1061 Neotetraose (LNnT) in Infant Formula. *Nutrients* **10**, (2018).

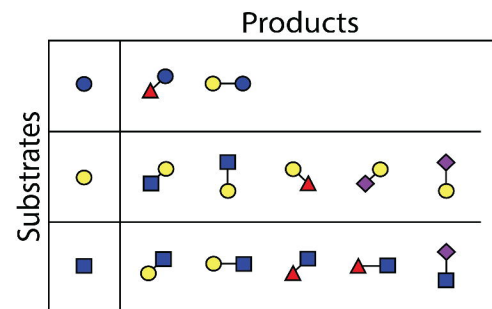
A**B****C**



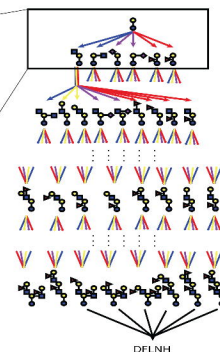
● Galactose (Gal)
 ● Glucose (Glc)
 ■ N-Acetylglucosamine (GlcNAc)
 ▲ Fucose (Fuc)
 ◆ Sialic Acid (Neu5Ac)



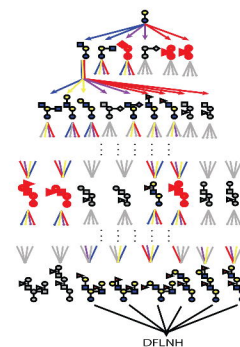
A. Reaction Rules



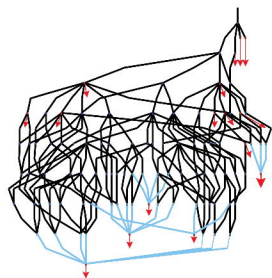
B. Complete Network



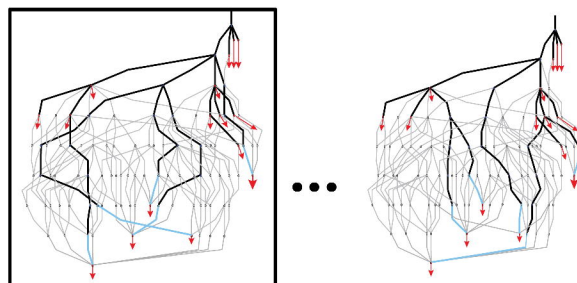
C. Trimmed Complete Network (FVA)



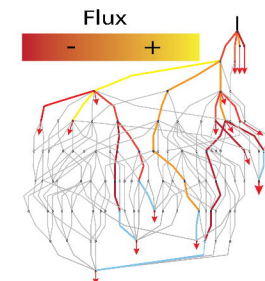
D. Reduced Network



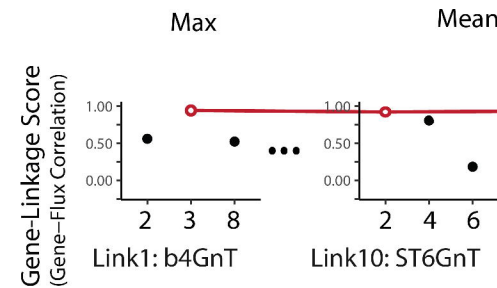
E. Candidate Models (MILP)



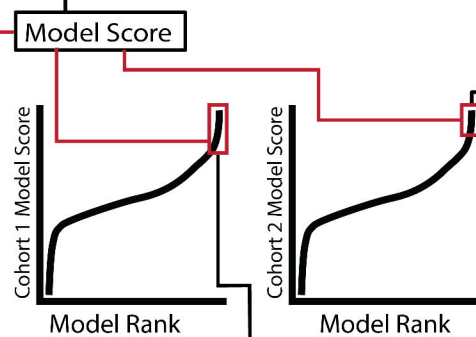
F. Model Flux (FBA)



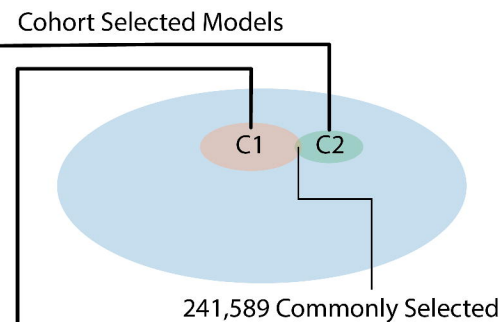
G. Model Score



H. Model Score Distributions

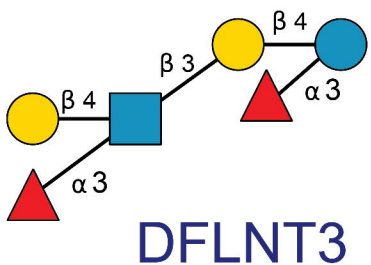
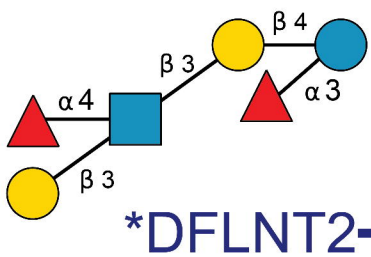
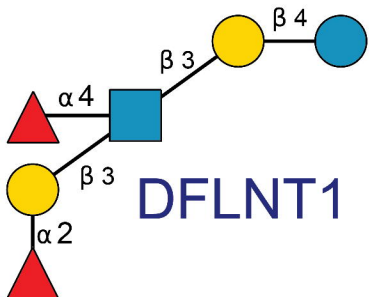


I. Commonly Selected Models



A

Examined DFLNT Isomers



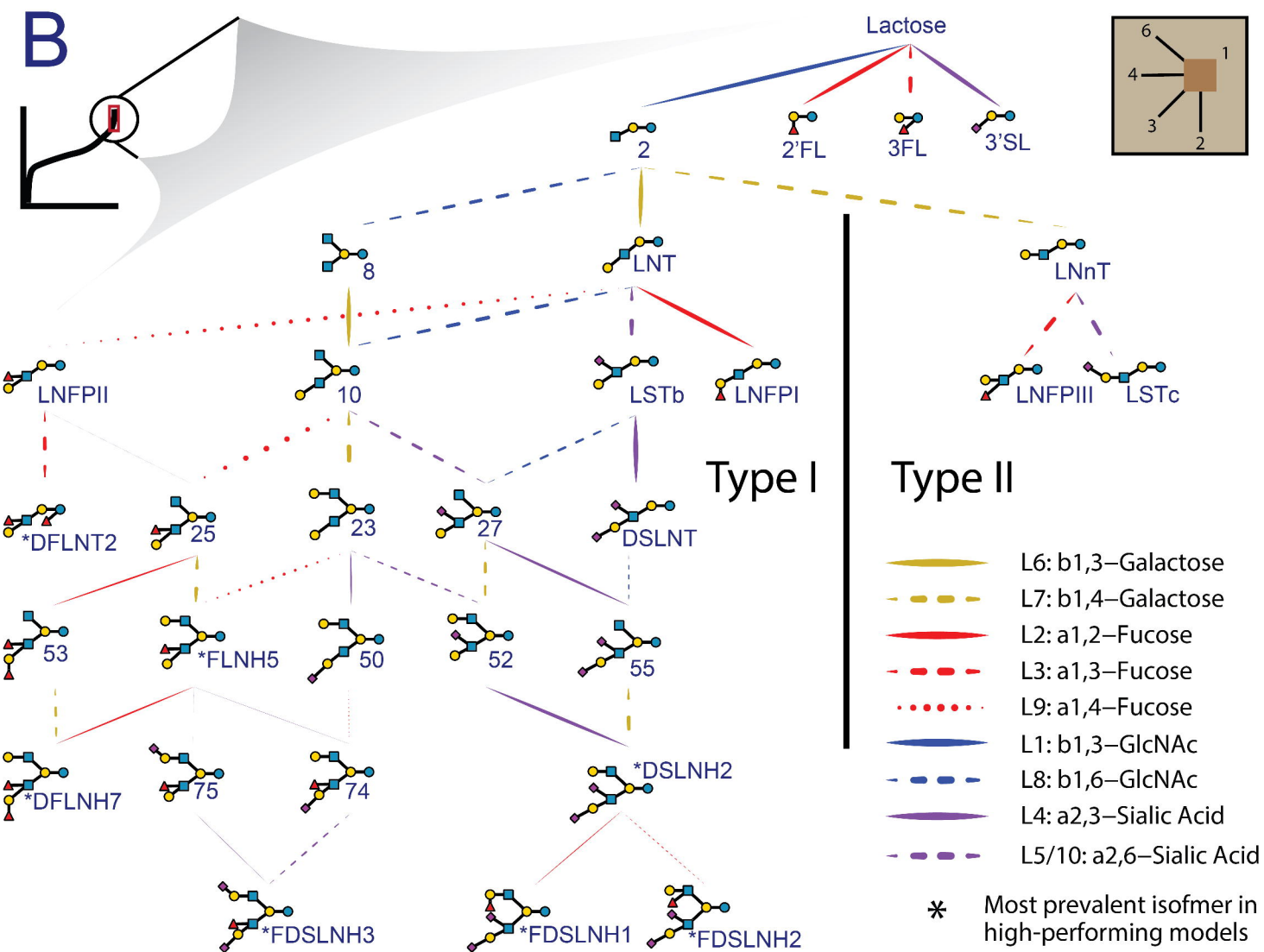
Galactose (Gal)

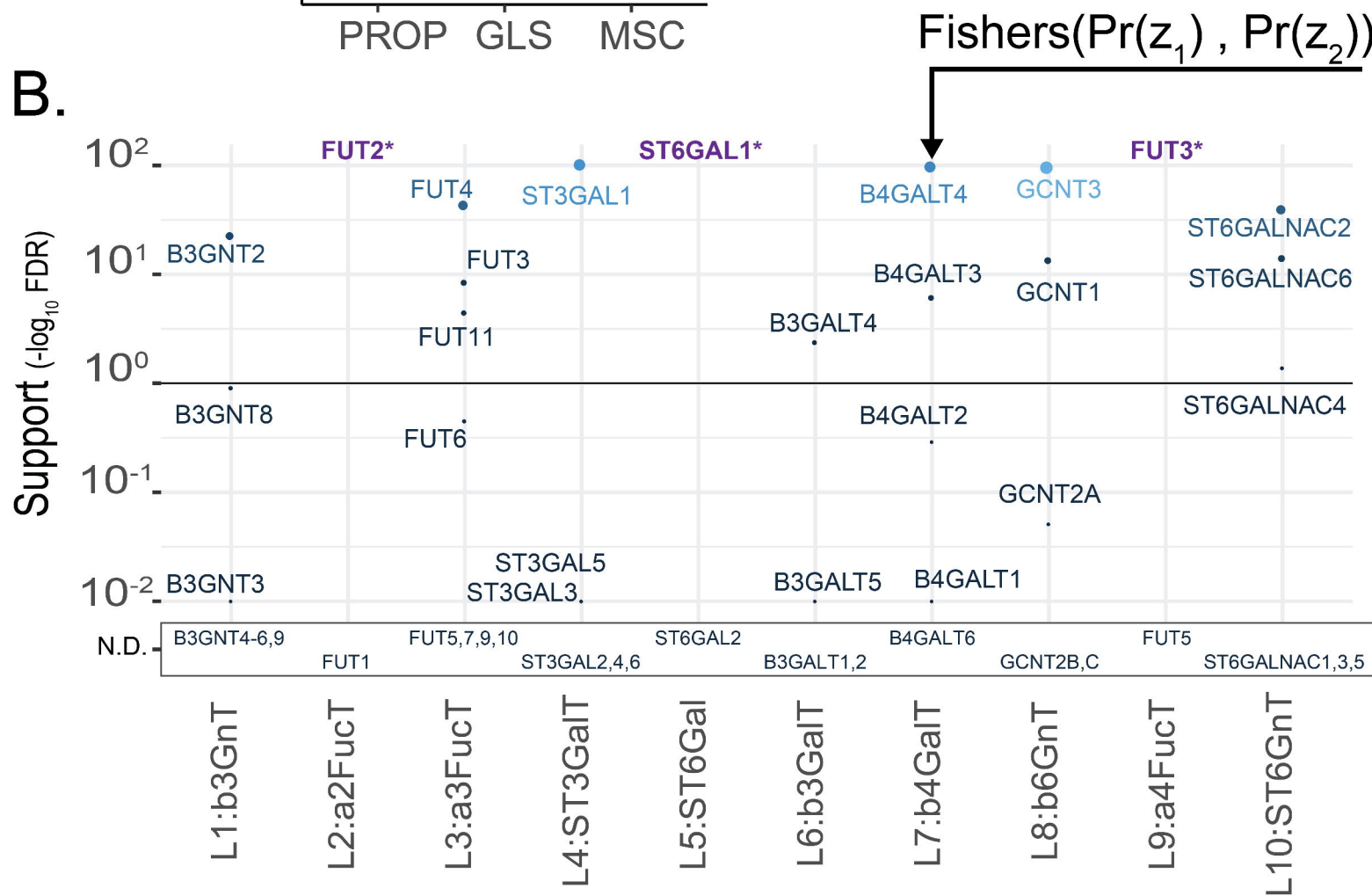
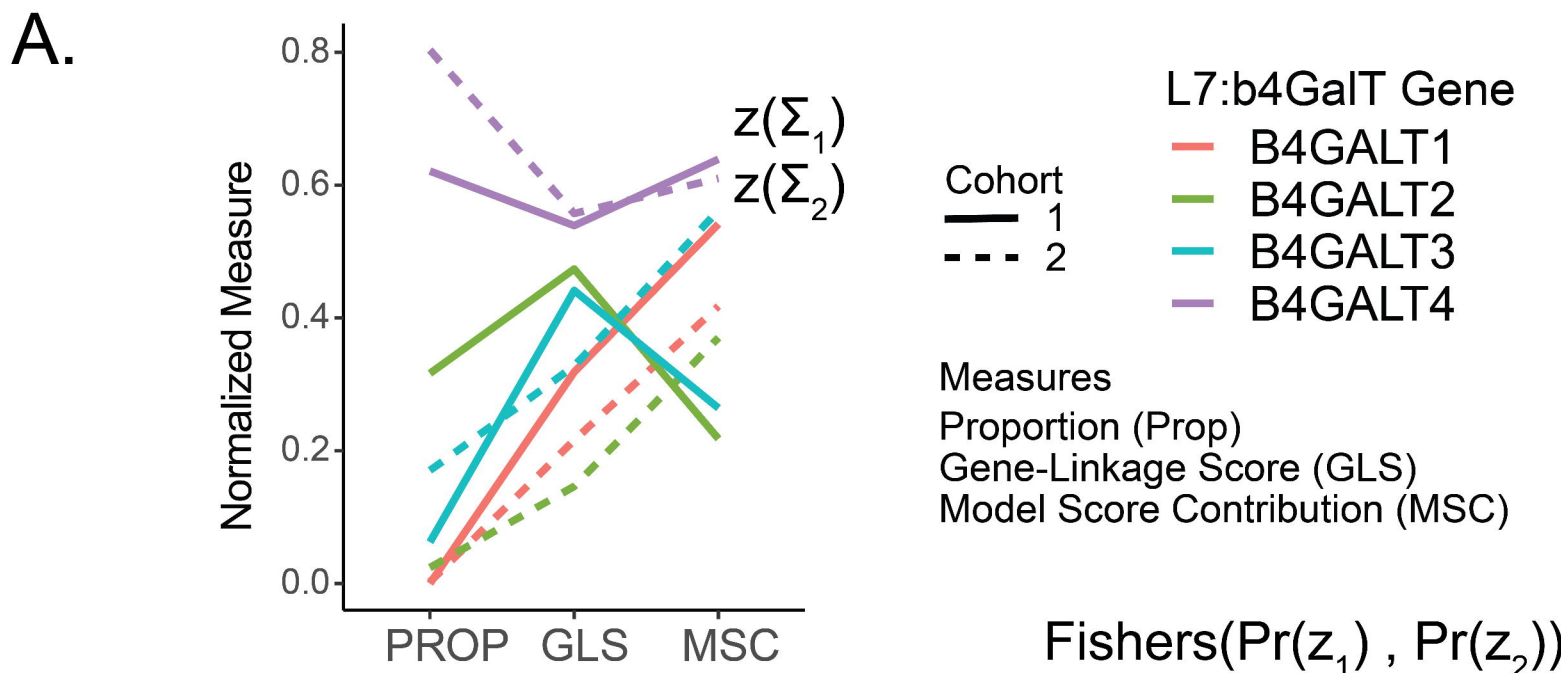
Glucose (Glc)

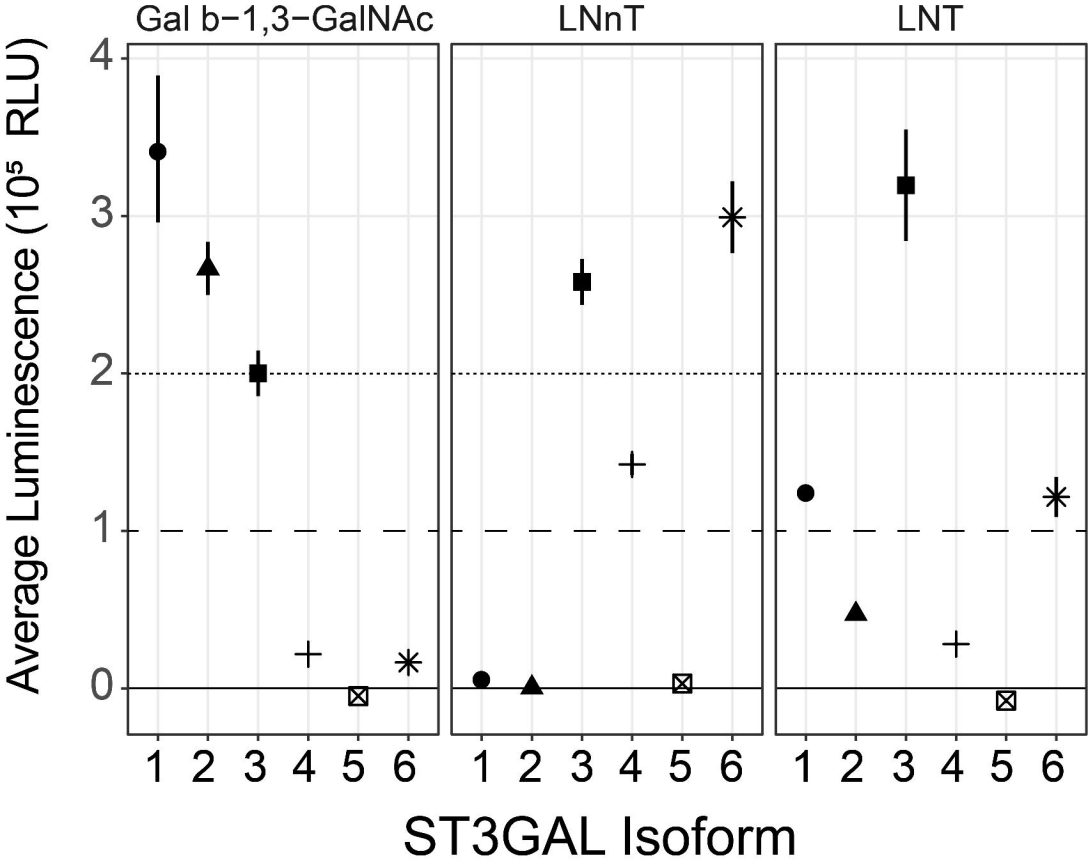
N-Acetylglucosamine (GlcNAc)

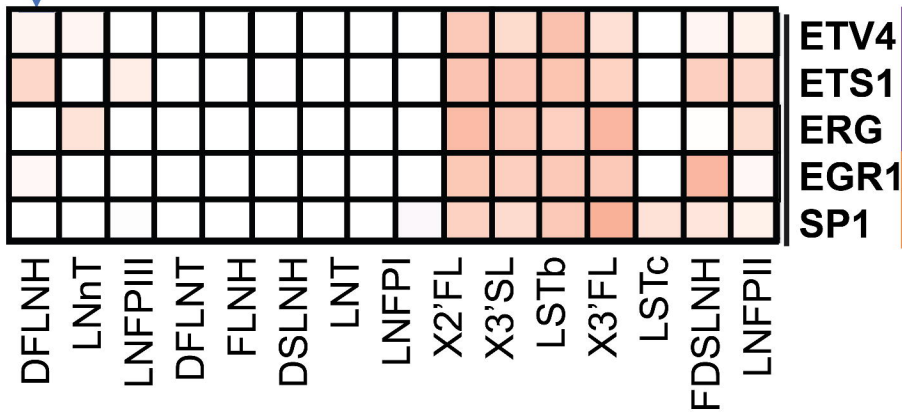
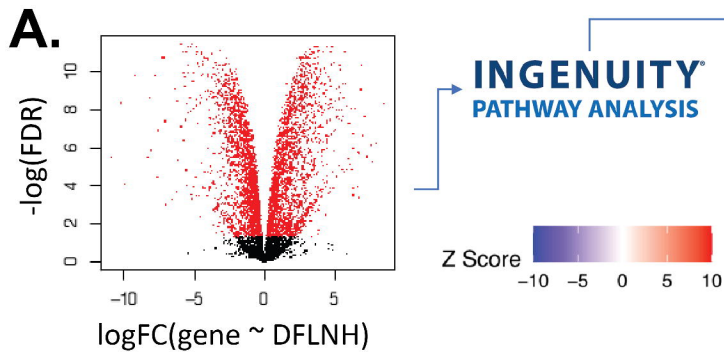
Fucose (Fuc)

Sialic Acid (Neu5Ac)

B







B.

