

1 Evolution of the Umbilical Cord Blood Proteome Across Gestational Development

2

3 Leena B. Mithal¹, Nicola Lancki², Ted Ling-Hu^{3,4}, Young Ah Goo⁵, Sebastian Otero¹, Nathaniel J.
4 Rhodes^{6,7,8}, Byoung-Kyu Cho⁵, William A. Grobman⁹, Judd F. Hultquist^{3,4}, Denise Scholtens², Karen
5 G. Mestan¹⁰, Patrick C. Seed¹

6 ¹ Department of Pediatrics, Division of Infectious Diseases, Ann & Robert H. Lurie Children's
7 Hospital of Chicago, Northwestern University Feinberg School of Medicine, Chicago, IL, USA.

8 ² Department of Preventive Medicine, Division of Biostatistics, Northwestern University Feinberg
9 School of Medicine, Chicago, IL, USA.

10 ³Department of Medicine, Division of Infectious Diseases, Northwestern University Feinberg School
11 of Medicine, Chicago, IL, USA.

12 ⁴Center for Pathogen Genomics and Microbial Evolution, Havey Institute for Global Health,
13 Northwestern University Feinberg School of Medicine, Chicago, IL, USA.

14 ⁵ Mass Spectrometry Technology Access Center at McDonnell Genome Institute (MTAC@MGI),
15 Washington University in Saint Louis School of Medicine, MO, USA.

16 ⁶Department of Pharmacy Practice, Midwestern University, College of Pharmacy, Downers Grove,
17 IL, USA.

18 ⁷ Pharmacometrics Center of Excellence, Midwestern University, Downers Grove, IL, USA.

19 ⁸ Department of Pharmacy, Northwestern Memorial Hospital, Chicago, IL, USA.

20 ⁹ Department of Obstetrics and Gynecology, Ohio State University, Columbus, OH, USA.

21 ¹⁰ Department of Pediatrics, Division of Neonatology, University of California San Diego, CA, USA.

22

23 * Correspondence:

24 Leena B. Mithal, MD MSCI
25 lmithal@luriechildrens.org

26 **Keywords:** cord blood, proteomics, neonatal immunology, prematurity, immune development,
27 biomarker development

28

29 Abstract:

30 Neonatal health is dependent on early risk stratification, diagnosis, and timely management of
31 potentially devastating conditions, particularly in the setting of prematurity. Many of these conditions
32 are poorly predicted in real-time by clinical data and current diagnostics. Umbilical cord blood may
33 represent a novel source of molecular signatures that provides a window into the state of the fetus at
34 birth. In this study, we comprehensively characterized the cord blood proteome of infants born

35 between 24 to 42 weeks using untargeted mass spectrometry and functional enrichment analysis. We
36 determined that the cord blood proteome at birth varies significantly across gestational development.
37 Proteins that function in structural development and growth (e.g., extracellular matrix organization,
38 lipid particle remodeling, and blood vessel development) are more abundant earlier in gestation. In
39 later gestations, proteins with increased abundance are in immune response and inflammatory
40 pathways, including complements and calcium-binding proteins. Furthermore, these data contribute
41 to the knowledge of the physiologic state of neonates across gestational age, which is crucial to
42 understand as we strive to best support postnatal development in preterm infants, determine
43 mechanisms of pathology causing adverse health outcomes, and develop cord blood biomarkers to
44 help tailor our diagnosis and therapeutics for critical neonatal conditions.

45 **Manuscript text:**

46 **Introduction**

47 Neonatal health is dependent on early risk stratification, diagnosis, and timely management of many
48 potentially devastating conditions. Preterm infants are at increased risk of prematurity-related
49 complications, including: early-onset sepsis, chronic lung disease, intraventricular hemorrhage,
50 necrotizing enterocolitis, and neurodevelopmental impairment.¹⁻³ Many of these conditions are
51 poorly predicted in real-time by clinical data, including currently available diagnostic testing. Thus,
52 biomarkers have been sought to aid early and targeted treatment and prognosis for these conditions.

53 Umbilical cord blood may represent a novel source of molecular signatures that provides a window
54 into the state of the fetus at birth. Umbilical cord blood inflammatory markers have been studied as
55 diagnostic indicators of early-onset sepsis⁴⁻⁶. Specific cord blood cytokines have been identified as
56 predictors or correlates of retinopathy of prematurity⁷, atopic disease⁸, infantile hemangioma⁹,
57 placental histopathology¹⁰, and more⁴. However, few of these cord blood biomarkers have been
58 translated into diagnostic tools in clinical practice.

59 “Omics” methodologies have been previously used to profile amniotic fluid and infant blood to
60 predict pre-eclampsia, preterm birth, and late-onset sepsis^{11,12}. Mass spectrometry (MS)-based
61 proteomics approaches have emerged as a particularly powerful technology for the comprehensive
62 profiling of proteins comprising the plasma microenvironment¹³. For example, longitudinal profiling
63 of postnatal proteomic changes has provided insights into the development of the immune system
64 over the first weeks to months of life¹⁴. Untargeted proteomic analyses furthermore provide an
65 unbiased approach to biomarkers discovery by removing the need to identify proteins of interest *a*
66 *priori*.¹⁵

67 Proteomic profiling of neonatal cord blood provides a molecular snapshot at variable timepoints
68 throughout neonatal development that could be used to reveal the underlying cellular processes
69 occurring at birth, readiness for postnatal life, and for the identification of biomarkers specific to
70 different disease states and prematurity-related complications.

71 While proteomic profiling of cord blood has demonstrated immunologic differences between preterm
72 and term infants¹⁶, prior research has lacked inclusion of preterm infants across the continuum of
73 gestational age and consideration of key perinatal characteristics such as the route of delivery,
74 preeclampsia, intraamniotic infection, and neonatal sepsis that are likely to affect protein abundance.
75 In this study, we have comprehensively characterized the cord blood proteome from infants born
76 between 25 to 42 weeks using MS to provide a benchmark of normative cord blood proteomic profile
77 and examine proteome differences across the developmental range of gestational ages.

78

79 **Methods**

80 Study cohort and specimen collection

81 We utilized archived cord blood plasma from an ongoing prospective study of infants born at
82 Northwestern Prentice Women’s Hospital between 2008-2019. Parents were consented prior to or
83 after birth; cord blood was centrifuged at 3000 rpm for 10 minutes and was separated into aliquots
84 stored at -80 degrees Celsius until use. Samples in this investigation were selected from the

85 biorepository based on gestational age and the absence of presumed or proven early onset neonatal
86 sepsis (i.e., the infant received no antibiotic treatment course for sepsis within the first 72 hours of
87 life and had no positive microbiologic sterile site cultures). A total of 150 infants were frequency
88 matched within each gestational age (GA) category (epochs: 25-28 weeks, 29-32 weeks, 33-36
89 weeks, 37-42 weeks) with approximately equal numbers by sex, route of delivery (vaginal delivery
90 vs. caesarean delivery with or without labor), and reason for preterm birth (maternal indication such
91 as preeclampsia vs. fetal/pregnancy indication such as spontaneous preterm labor or preterm
92 premature rupture of membranes). Clinical data including birth weight and intraamniotic infection
93 were collected from the electronic medical record. This study was approved by the Institutional
94 Review Boards of Northwestern University (STU00201858) and Lurie Children's Hospital (IRB
95 2018-2145). Parental informed consent was obtained for use of clinical data and infant cord blood
96 samples. All research activities were performed in accordance with the Declaration of Helsinki.

97 Mass spectrometry sample preparation and analysis

98 Samples were thawed on ice and 20 μ l of plasma was utilized for study. Protein concentrations were
99 determined using the Bicinchoninic Acid (BCA) method; untargeted mass spectrometry-based
100 proteomic analysis was applied to 600 μ g of extracted protein from each plasma sample. Samples
101 were first depleted of fourteen known highly abundant proteins (Albumin, IgA, IgD, IgE, IgG, IgG
102 (Light chains), IgM, Alpha-1-acid glycoprotein, Alpha-1-antitrypsin, Alpha-2-macroglobulin,
103 Apolipoprotein A1, Fibrinogen, Haptoglobin, and Transferrin) using the Top 14 Abundant Protein
104 Depletion Spin Columns (Thermo Scientific, Rockford, IL, USA). Remaining proteins were purified
105 by acetone/TCA precipitation, reduced, alkylated, and digested with trypsin. Digested peptides were
106 desalted on C18 columns (Thermo Scientific, Rockford, IL, USA) and eluted in 80% acetonitrile in
107 0.1% formic acid. Peptides were reconstituted with 0.1% formic acid in water and injected onto the
108 in-house C18 trap column (3 cm length, 150 μ m inner diameter, 3 μ m particle size) coupled with an
109 analytical C18 column (10.5 cm length, 75 μ m inner diameter, 2 μ m particle size, PicoChip).
110 Samples were separated using a linear gradient from 5% ACN/0.1% formic acid to 40% ACN/0.1%
111 formic acid over 120 minutes using an UltiMate 3000 Rapid Separation nanoLC coupled to a
112 Orbitrap Elite Mass Spectrometer (Thermo Fisher Scientific Inc, San Jose, CA). The full scans were
113 acquired from 400-2000m/z at 60,000 resolving power and automatic gain control (AGC) set to
114 1×10^6 . The top fifteen most abundant precursor ions in each full scan were selected for
115 fragmentation. Precursors were selected with an isolation width of 1 Da and fragmented by collision-
116 induced dissociation (CID) at 35% normalized collision energy. Previously selected ions were
117 dynamically excluded from re-selection for 58 seconds.

118 Samples were analyzed in duplicate, in a specified run order, across four batches. Samples were
119 randomly assigned to batches using a stratified sampling approach to achieve balance on gestational
120 age and other clinical characteristics (sex, type of delivery). A representative "pooled control,"
121 including samples representing the full spectrum of the cohort, was used as an "internal standard"
122 and run multiple times in each batch. Within each batch, the run order for samples and controls was
123 determined by simple random sampling. MS raw files were analyzed with MaxQuant software
124 (version 1.6.0.16).¹⁷ MS/MS-based peptide identification was carried out against the SwissProt
125 human database with the Andromeda search engine in MaxQuant¹⁸ using a target-decoy approach to
126 identify peptides and proteins at an FDR <1%. For LFQ, the MaxLFQ algorithm was used as part of
127 the MaxQuant environment.¹⁹ The following modifications were set as search parameters: trypsin
128 digestion cleavage after K or R (except when followed by P), 2 allowed missed cleavage sites,
129 carbamidomethylated cysteine (static modification), and oxidized methionine, protein N-term
130 acetylation (variable modification). Search results were validated with peptide and protein FDR, both

131 at 0.01. Transformed (\log_2) LFQ values were used for all statistical analyses. The mass spectrometry
132 proteomics data have been deposited to the ProteomeXchange Consortium via the PRIDE partner
133 repository with the dataset identifier PXD051974.

134 Proteomics data normalization

135 Boxplots representing the median and first and third quartiles were used to visualize the distribution
136 of protein concentration among all proteins in pooled controls and identify the presence of any batch
137 effects. To correct for batch effects demonstrated, batch normalization was conducted as follows.
138 Proteins that were detected in only one batch were excluded. Using the pooled control samples, the
139 average difference in \log_2 LFQ value relative to the first batch was estimated using a linear
140 regression model. Briefly, a beta coefficient for each batch was estimated using linear regression with
141 batch one serving as the referent. The average protein difference from the first batch was then
142 subtracted from the \log_2 LFQ value of each batch to determine the normalized \log_2 LFQ value.
143 Visual inspection of post-normalization protein levels by batch was used to determine the adequacy
144 of the normalization procedure. The batch normalized protein abundance for each sample was
145 averaged across each technical replicate for subsequent inter-patient analyses. If a protein was only
146 detected in one of the replicates, the value of the batch normalized detected protein from the single
147 replicate was used for analyses.

148 Differential protein abundance determination

149 Separate linear regression models were used to examine the association between protein abundance
150 and GA (unadjusted and adjusted for sex, labor, route of delivery, and preeclampsia). The response
151 variable for each model was the batch normalized value \log_2 transformed protein level for the given
152 protein. Proteins were included in adjusted models if found in more than one sex, delivery category,
153 and preeclampsia category. The primary explanatory variable of interest was GA. Scatter plots were
154 examined to determine whether GA demonstrated a linear or nonlinear association (e.g., using splines
155 or quadratic terms) with protein level. The relationships between GA and protein abundance, in
156 general, appeared linear across proteins, so a linear term for GA was included in models. There were
157 seven sets of twins among the 150 controls. One twin from a pair was randomly selected to be
158 included in the models (n=143). To control Type 1 error rate, P-values were adjusted for multiple
159 testing using the Benjamini-Hochberg False Discovery Rate (FDR) method, and associations with
160 FDR-adjusted P values <0.05 were considered statistically significant²⁰.

161 Functional enrichment and visualization

162 The relative expression abundance of all proteins that changes significantly over gestational age was
163 visualized in a heat map. The batch normalized protein values were z-score normalized by
164 subtracting the relative protein abundance within a given specimen by the mean abundance across all
165 specimens in which the protein was detected and then dividing by standard deviation. Proteins that
166 were undetected in more than 50% of specimens were excluded from visualization. Z-score
167 normalized values were visualized in a heatmap using the *clustermap* function in the *seaborn* (v
168 0.11.1) package within the python (v 3.8.8) environment with specimens ordered left to right by GA
169 and proteins clustered by z-score profile from top to bottom. The *clustermap* function uses
170 hierarchical clustering with average linkage and Euclidean distance.

171 Functional enrichment analysis of the proteins found to be significantly increased or decreased in
172 abundance was performed using MetaScape v3.5

173 (<https://metascape.org/gp/index.html#/main/step1>)²¹. UniProt IDs were used as unique identifiers;
174 two isoforms of APOB (P04114), PLG (P00747), and FGA (P02671) were consolidated and
175 immunoglobulins were excluded (P0DOX5, P0DOX7, P0DOY3, and P01859) for final analysis of
176 64 proteins. All proteins detected in the overall proteomic dataset (n = 465) were set as the
177 background gene set before enrichment. Protein-protein interaction networks were visualized using
178 STRING v11.5 (<https://string-db.org/>)²². Network visualization was limited to physical subnetworks
179 based on experiment and database active interaction sources with a 0.15 minimum interaction score
180 required. Nodes were colored by an increase or decrease in abundance with edge width reflective of
181 protein interaction confidence score. Proteins contributing to significantly enriched pathways were
182 annotated with colored boxes.

183

184 **Results**

185 Patient demographics

186 The distribution of GA and associated clinical/demographic details for the 150 infants included in
187 this study are displayed in **Table 1**. The mean GA across all infants was 33.2 weeks (standard
188 deviation 4.5, range 25.9-41.4). 17 infants (11%) were 25-28 weeks, 43 (29%) were 29-32 weeks, 50
189 (33%) were 33-36 weeks, and 40 (27%) were 37 weeks and greater. 77 (51%) of the infants were
190 female. 34 (23%) infants were born to women with preeclampsia. 44 (29%) infants were from 22
191 individuals with multiple gestations, all of whom were born at less than 37 weeks GA. Only one
192 infant was born to an individual who had clinical chorioamnionitis.

193 Differential protein abundance across GA

194 The total BCA, representative of protein abundance, is positively correlated with GA (**Figure 1**). Of
195 the 465 unique proteins identified in control plasma samples, 391 were included in the adjusted
196 regression models (adjusted for sex, preeclampsia, and delivery route). Proteins were excluded from
197 adjusted multivariable regression models if they were only found in one group of the covariate
198 categories (for example found in only male infants). Gestational age was associated with protein
199 abundance in 70 proteins with FDR-adjusted P-value <0.05 (**Supplemental Table 1**). To visualize
200 each protein's change over GA, the normalized protein abundance in each specimen was plotted
201 relative to GA of the infant (**Figure 2**). The slope ('beta' value) of the linear fitted model for each
202 protein is provided in **Supplemental Table 1**. Representative plots in **Figure 2** demonstrate
203 examples of proteins with positive (e.g., plasminogen; **Figure 2A**) and negative (e.g., alpha-
204 fetoprotein; **Figure 2B**) correlation between protein abundance and GA. These changes are
205 summarized in a volcano plot (**Figure 3**) depicting the log₁₀ of the FDR-adjusted P-values and the
206 associated betas from linear regression models. Proteins such as alpha-fetoprotein, collagen alpha-
207 1(V) chain, and basement membrane-specific heparan sulfate proteoglycan core protein are highly
208 abundant earlier in gestational development while many immunologically active proteins are more
209 abundant later, including IgG-1 chain C region, complement C1q subunit C, and protein S100-A9, a
210 calcium- and zinc-binding protein which plays a prominent role in the regulation of inflammatory
211 response (all aforementioned proteins with p<0.0001).

212 Visualization and pathway analysis

213 To better visualize the differences in protein levels across GA, we plotted a heatmap of the
214 normalized protein levels for each significantly changing protein identified above in each specimen

215 ordered along the x-axis by GA (**Figure 4A**). Proteins without detectable levels in more than 50% of
216 specimens (n = 15) were excluded from visualization and hierarchical clustering was used to group
217 proteins by similarity in abundance trends over GA. This highlights several distinct groups of
218 proteins where levels change over time. For example, COL5A1, CD14, HSPG2, QSOX1, FCGBP
219 seem to be abundant in early GA, but decrease as GA increases. This trend is also apparent in the
220 cluster located at the bottom half of the heatmap that includes CD109, COL1A1, APOC3, APOE,
221 TGFBI, AFP, AGT, APOB, LUM, SERPINA1, B2M, FGA, THBS4, F13A1 and SERPINA5.
222 However, several proteins also follow the opposite trend with lower abundance early and higher
223 abundance late, including HBA1, HBB, HPX, IGFALS, CP, AFM, SERPINF2, SERPIND1, A2M,
224 ATRN, PGLYRP2, IGHG1, C7, ITIH1, PLG, F2, SERPINC1, C1QC and C1QA.

225 Functional enrichment analysis was performed for those proteins found to be significantly decreased
226 in relative abundance (n = 29, **Figure 4B**) or increased in abundance (n = 34, **Figure 4C**) with
227 increasing GA (7 identifiers did not map back to unique proteins and were excluded from analysis,
228 see **Methods**). Proteins that decreased in abundance were enriched for eight pathways: NABA core
229 matrisome, extracellular matrix organization, lipid particle remodeling, smooth muscle proliferation,
230 blood vessel development, glycosaminoglycan metabolism, insulin-like growth factor regulation, and
231 amyloid fiber formation. These enriched protein sets include several components of known protein
232 complexes. For example, the proteoglycan LUM and the collagen proteins COL5A1 and COL1A1
233 form high confidence protein interactions and are all associated with extracellular matrix
234 organization pathways. Likewise, CETP, APOC3, APOB, and APOE are known to interact and play
235 critical roles in lipid particle remodeling. Several factors implicated in insulin-like growth factor
236 regulation were also decreased, including SERPINF2, SERPIND1, FGA, and SERPINA5.

237 More proteins were found to increase in abundance over GA than decrease. However, these were
238 associated with a narrower set of pathways, specifically: protein nitrosylation, metal ion homeostasis,
239 humoral immune response, NABA core matrisome, and positive regulation of cell death. One well-
240 known transition that occurs in the serum throughout development is the swapping of hemoglobin
241 subunits from γ -globin in neonates to β -globin and δ -globin gene expression in pediatric and adult
242 patients²³. Consistent with this transition, we see increased abundance of β -globin (HBB) and δ -
243 globin (HBD), as well as α -globin (HBA2). We additionally see the increased abundance of several
244 proteins associated with the immune response, including several complement proteins (C1QA,
245 C1QC, C7, CFP, and C8G) and several S100 calcium binding proteins (S100A9, S100A12, and
246 S100A8). Notably, several immunoglobulins were also increased over the course of GA (specifically
247 kappa light chain, lambda light chain, and gamma heavy chains), though these are not visualized
248 here.

249

250 **Discussion**

251 Our data demonstrate that the abundance of several cord blood proteins varies significantly across
252 GA. Proteins that function in structural development and growth, including extracellular matrix
253 organization, lipid particle remodeling, blood vessel development, and insulin-like growth factor
254 regulation, are more abundant earlier in gestation. Later in gestation, proteins involved in immune
255 response pathways, including complements, and calcium-binding proteins involved in inflammation
256 are higher in abundance. These data highlight the differences in immunologic state across GA and
257 provide insights into the higher risk of invasive infections among preterm infants. Furthermore, these
258 data contribute to the knowledge of the physiologic state of neonates across GA, which is crucial to

259 understand as we: 1) strive to emulate the *in utero* environment to best support the developmental
260 process of those born preterm, 2) understand mechanisms of pathology that cause adverse health
261 outcomes for preterm infants, and 3) develop cord blood markers for neonatal disease conditions that
262 can predict and help tailor medical management.

263 In a 2021 review of proteomic studies that attempted to identify biomarkers for prematurity-related
264 diseases, Letunica *et al.* determined that only 13% of studies investigated cord blood even though
265 cord blood is a readily available specimen at birth.¹¹ Suski *et al.* investigated the cord blood proteome
266 of preterm infants in three GA groups (≤ 26 weeks, 27-28 weeks, and 29-30 weeks) and compared
267 them to the proteomes of a full term control group. They reported differences in inflammatory,
268 immunomodulation, coagulation, and complement systems in preterm versus term infants.¹⁶
269 Specifically, they found that preterm infants had decreased levels of anti-inflammatory proteins (e.g.,
270 orsomucoid isoforms) and B-cell mediated immunity markers, and increased abundance of
271 inflammatory proteins such as leucine-rich alpha-2-glycoprotein (LRG1) and complement activation
272 cascades, a finding that complements our results of lower proteins related to humoral immunity at
273 earlier GA.

274 However, the authors also suggest an increase in inflammatory mediators in preterm infants, whereas
275 our results showed increased inflammatory and immune response proteins and complement
276 components later in gestational age. Our pathway analysis revealed many of the proteins that function
277 in inflammatory signaling and immune response are lower in preterm infants with no infection. A
278 likely explanation for this notable difference is that our research excluded preterm infants with early
279 onset sepsis, and thus represents the state of the cord blood proteome in the absence of infection.
280 Given that treatment of early-onset sepsis is common in very preterm infants, the analysis of cord
281 blood inflammatory proteins may be skewed if one has not accounted for infection.

282 Other types of immune phenotyping have been reported in cord blood across GA. Olin *et al.* noted
283 differences in both cord blood proteins and decreased neutrophil proportions in preterm compared to
284 term infants. They reported an increase in inflammatory cord blood proteins, attributed to the role of
285 inflammation and infection in preterm birth, but not reflecting gestational norms without infection.¹⁴
286 Anderson *et al.* utilized flow cytometry and cytokine assays of cord blood to compare preterm infants
287 (30-34 weeks GA) to full-term infants.²⁴ They found that preterm infants had lower frequencies of
288 monocytes, NK cells, CD8+ T-cells and gamma-delta T-cells than their term full term counterparts.
289 There were increased intermediate monocytes, CD4 T cells, Tregs, and transitional B-cells in preterm
290 infants indicating immaturity of the innate immune system and a skewed cellular landscape related to
291 increased susceptibility of preterm infants to bacterial and viral infections. They also noted lower
292 levels of pro-inflammatory cytokines and chemokines in preterm infants, further confirming preterm
293 infants impaired ability to fight off infection. Finally, Peterson *et al.* applied single-cell
294 immunoprofiling of cord blood for 45 infants (20 preterm) after excluding infants exposed to clinical
295 chorioamnionitis or with active infection.²⁵ The study also controlled for potential other clinical
296 confounders, including steroid administration. They found a strong relationship between GA and the
297 neonatal immune profile at birth. Specifically, increasing GA was associated with a progressive
298 increase in the ligand-specific responsiveness to immune system stimulation. This finding aligns with
299 our finding of increased cell-signaling, calcium binding, and immune response proteins with later
300 gestational age. Our work supports the conclusion that decreased antigen- and cytokine-specific
301 immune responses may contribute to preterm infant susceptibility to infection.

302 Furthermore, differences in proteins across GA may provide insight into underlying pathophysiology
303 and risk of pathology. Functional analysis identified several pathways associated with increased

304 abundance of proteins that are implicated in vascular development, lipid metabolism, smooth muscle
305 proliferation, insulin-like growth factor regulation, and the matrisome. For example, afamin, an anti-
306 inflammatory protein previously hypothesized to be a hallmark of detrimental oxidative stress and
307 related to retinopathy of prematurity, is less abundant in the cord blood of preterm infants.¹⁶ The
308 process of in utero development represents a complex and dynamic system between the pregnant
309 person and fetus. Through this study of neonates born from 25-42 weeks GA, we aim to help
310 establish the baseline state of the developmental continuum.

311 The strengths of this study include: 1) the analysis of cord blood proteomics on a large sample size
312 across the GA spectrum; 2) precise clinical categorization and consideration of covariates that may
313 impact the cord blood proteome including exclusion of infants with early onset infection and
314 adjustment for labor, preeclampsia, and sex; and 3) careful methodologic and data normalization,
315 both in design and analysis of discovery mass spectrometry proteomics (distribution and
316 normalization across batches, pooled control, addressing missingness). Additionally, functional
317 pathway analysis strengthens our ability to parse key pathways of relevance and provide validation
318 through the demonstration of known GA-related differences in hemoglobin and immunoglobulin
319 proteins²⁶. Limitations include that mass spectrometry proteomics does not provide absolute
320 quantitation of protein, but rather spectral counts and relative abundance. The detectable protein
321 abundance reflects the level after potential clearance, degradation, or transport/localization of
322 expressed proteins to compartments. For this reason, we highlight relative abundance and levels of
323 proteins rather than using terms akin to protein expression (i.e., “up/down-regulation”). Thus,
324 specific biomarker development warrants quantitative validation methods. Additionally, the cord
325 blood specimen used in this analysis was intended to be obtained at the time of birth from the
326 umbilical vein. However, it is possible that there is some mixing of umbilical arterial and venous
327 blood. Prior literature raises questions about mediating cord blood markers by placental clearance
328 and whether cord blood proteins may reflect maternal serum. In multiple studies, paired analysis of
329 maternal and fetal cord blood biomarkers has shown weak or no correlation.^{27,28}

330 In conclusion, our study utilizing untargeted proteomics has demonstrated that the cord blood
331 proteome varies significantly with GA at birth. There are meaningful differences in several pathways,
332 including crucial aspects of inflammation and immune response. Future research can apply this
333 knowledge of the baseline state to find methods to develop more precise, GA-specific cord blood
334 diagnostic markers of short and perhaps long-term²⁹ health and disease.

335 References

- 336 1 Siffel, C., Hirst, A. K., Sarda, S. P., Kuzniewicz, M. W. & Li, D.-K. The clinical burden of
337 extremely preterm birth in a large medical records database in the United States: Mortality
338 and survival associated with selected complications. *Early human development* **171**, 105613
339 (2022). <https://doi.org/10.1016/j.earlhumdev.2022.105613>
- 340 2 Manuck, T. A. *et al.* Preterm neonatal morbidity and mortality by gestational age:
341 a contemporary cohort. *American journal of obstetrics and gynecology* **215**, 103.e101-
342 103.e114 (2016). <https://doi.org/10.1016/j.ajog.2016.01.004>
- 343 3 Ohuma, E. O. *et al.* National, regional, and global estimates of preterm birth in 2020, with
344 trends from 2010: a systematic analysis. *Lancet* **402**, 1261-1271 (2023).
345 [https://doi.org/10.1016/s0140-6736\(23\)00878-4](https://doi.org/10.1016/s0140-6736(23)00878-4)
- 346 4 Dongen, O. R. E. *et al.* Umbilical Cord Procalcitonin to Detect Early-Onset Sepsis in
347 Newborns: A Promising Biomarker. *Frontiers in pediatrics* **9**, 779663 (2021).
348 <https://doi.org/10.3389/fped.2021.779663>
- 349 5 Mithal, L. B., Palac, H. L., Yogev, R., Ernst, L. M. & Mestan, K. K. Cord Blood Acute Phase
350 Reactants Predict Early Onset Neonatal Sepsis in Preterm Infants. *PloS one* **12**, e0168677
351 (2017). <https://doi.org/10.1371/journal.pone.0168677>
- 352 6 Su, H. *et al.* Inflammatory markers in cord blood or maternal serum for early detection of
353 neonatal sepsis—a systemic review and meta-analysis. *Journal of perinatology : official*
354 *journal of the California Perinatal Association* **34**, 268-274 (2014).
355 <https://doi.org/10.1038/jp.2013.186>
- 356 7 Park, Y. J. *et al.* Immune and Inflammatory Proteins in Cord Blood as Predictive Biomarkers
357 of Retinopathy of Prematurity in Preterm Infants. *Investigative ophthalmology & visual*
358 *science* **60**, 3813-3820 (2019). <https://doi.org/10.1167/iovs.19-27258>
- 359 8 Soti, A. L. *et al.* Can biomarkers in umbilical cord blood predict atopic disease at school age?
360 *Pediatric research* **89**, 389-392 (2021). <https://doi.org/10.1038/s41390-019-0686-z>
- 361 9 Jiang, C. H. *et al.* Metabolic Profiling Revealed Prediction Biomarkers for Infantile
362 Hemangioma in Umbilical Cord Blood Sera: A Prospective Study. *J Proteome Res* **21**, 822-
363 832 (2022). <https://doi.org/10.1021/acs.jproteome.1c00430>
- 364 10 Mestan, K. *et al.* Cord blood biomarkers of the fetal inflammatory response. *The journal of*
365 *maternal-fetal & neonatal medicine : the official journal of the European Association of*
366 *Perinatal Medicine, the Federation of Asia and Oceania Perinatal Societies, the International*
367 *Society of Perinatal Obstet* **22**, 379-387 (2009). <https://doi.org/10.1080/14767050802609759>
- 368 11 Letunica, N. *et al.* The use of proteomics for blood biomarker research in premature infants: a
369 scoping review. *Clin Proteomics* **18**, 13 (2021). <https://doi.org/10.1186/s12014-021-09316-y>
- 370 12 Law, K. P., Han, T. L., Tong, C. & Baker, P. N. Mass spectrometry-based proteomics for pre-
371 eclampsia and preterm birth. *International journal of molecular sciences* **16**, 10952-10985
372 (2015). <https://doi.org/10.3390/ijms160510952>
- 373 13 Shuken, S. R. An Introduction to Mass Spectrometry-Based Proteomics. *Journal of Proteome*
374 *Research* **22**, 2151-2171 (2023). <https://doi.org/10.1021/acs.jproteome.2c00838>
- 375 14 Olin, A. *et al.* Stereotypic Immune System Development in Newborn Children. *Cell* **174**,
376 1277-1292 e1214 (2018). <https://doi.org/10.1016/j.cell.2018.06.045>

- 377 15 Geyer, P. E., Holdt, L. M., Teupser, D. & Mann, M. Revisiting biomarker discovery by
378 plasma proteomics. *Mol Syst Biol* **13**, 942 (2017). <https://doi.org/10.15252/msb.20156297>
- 379 16 Suski, M. *et al.* Plasma proteome changes in cord blood samples from preterm infants. *J*
380 *Perinatol* **38**, 1182-1189 (2018). <https://doi.org/10.1038/s41372-018-0150-7>
- 381 17 Cox, J. & Mann, M. MaxQuant enables high peptide identification rates, individualized
382 p.p.b.-range mass accuracies and proteome-wide protein quantification. *Nature Biotechnology*
383 **26**, 1367-1372 (2008). <https://doi.org/10.1038/nbt.1511>
- 384 18 Cox, J. *et al.* Andromeda: A Peptide Search Engine Integrated into the MaxQuant
385 Environment. *Journal of Proteome Research* **10**, 1794-1805 (2011).
386 <https://doi.org/10.1021/pr101065j>
- 387 19 Cox, J. *et al.* Accurate proteome-wide label-free quantification by delayed normalization and
388 maximal peptide ratio extraction, termed MaxLFQ. *Mol Cell Proteomics* **13**, 2513-2526
389 (2014). <https://doi.org/10.1074/mcp.M113.031591>
- 390 20 Benjamini, Y. H., Y. Controlling the false Discovery Rate: A Practical and Powerful
391 Approach to Multiple Testing. *Journal of the Royal Statistical Society* **57**, 289-300 (1995).
- 392 21 Zhou, Y. *et al.* Metascape provides a biologist-oriented resource for the analysis of systems-
393 level datasets. *Nat Commun* **10**, 1523 (2019). <https://doi.org/10.1038/s41467-019-09234-6>
- 394 22 Szklarczyk, D. *et al.* The STRING database in 2023: protein-protein association networks and
395 functional enrichment analyses for any sequenced genome of interest. *Nucleic Acids Res* **51**,
396 D638-d646 (2023). <https://doi.org/10.1093/nar/gkac1000>
- 397 23 Wang, X. & Thein, S. L. Switching from fetal to adult hemoglobin. *Nat Genet* **50**, 478-480
398 (2018). <https://doi.org/10.1038/s41588-018-0094-z>
- 399 24 Anderson, J. *et al.* Immune Profiling of Cord Blood From Preterm and Term Infants Reveals
400 Distinct Differences in Pro-Inflammatory Responses. *Frontiers in Immunology* **12** (2021).
401 <https://doi.org/10.3389/fimmu.2021.777927>
- 402 25 Peterson, L. S. *et al.* Single-Cell Analysis of the Neonatal Immune System Across the
403 Gestational Age Continuum. *Frontiers in Immunology* **12** (2021).
404 <https://doi.org/10.3389/fimmu.2021.714090>
- 405 26 Herkner, K. R. *et al.* Pediatric and perinatal reference intervals for immunoglobulin light
406 chains kappa and lambda. *Clin Chem* **38**, 548-550 (1992).
- 407 27 Blohm, M. E. *et al.* Cardiovascular biomarkers in paired maternal and umbilical cord blood
408 samples at term and near term delivery. *Early human development* **94**, 7-12 (2016).
409 <https://doi.org/https://doi.org/10.1016/j.earlhumdev.2016.01.001>
- 410 28 Sivan, E. *et al.* Adiponectin in Human Cord Blood: Relation to Fetal Birth Weight and
411 Gender. *The Journal of Clinical Endocrinology & Metabolism* **88**, 5656-5660 (2003).
412 <https://doi.org/10.1210/jc.2003-031174>
- 413 29 Hansmeier, N., Chao, T. C., Goldman, L. R., Witter, F. R. & Halden, R. U. Prioritization of
414 biomarker targets in human umbilical cord blood: identification of proteins in infant blood
415 serving as validated biomarkers in adults. *Environ Health Perspect* **120**, 764-769 (2012).
416 <https://doi.org/10.1289/ehp.1104190>

417

418 **Acknowledgments**

419 The authors would like to acknowledge Erin Cullather, BS; Paul Martin Thomas, PhD; Aaron
420 Hamvas, MD; and Thomas Shanley, MD for their support toward this work. Proteomics services
421 were performed by the Northwestern Proteomics Core Facility, generously supported by NCI CCSG
422 P30 CA060553 awarded to the Robert H Lurie Comprehensive Cancer Center, instrumentation award
423 (S10OD025194) from NIH Office of Director, and the National Resource for Translational and
424 Developmental Proteomics supported by P41 GM108569.

425 **Author Contributions**

426 LBM: conceptualization, methodology, data curation, analysis, funding acquisition, writing - original
427 draft. B-KC and YG: methodology, investigation, data curation, writing - review & editing. DS, NL,
428 TL-H and JFH: methodology, data curation, formal analysis, visualization, writing - review &
429 editing. SO and NJR: methodology, investigation, writing - review & editing. WAG:
430 conceptualization, methodology, supervision, writing - review & editing. KM: conceptualization,
431 methodology, data curation, supervision, writing - review & editing. PCS: conceptualization,
432 methodology, data curation, supervision, writing - review & editing, funding acquisition. All authors
433 contributed to the article and approved the submitted version.

434 **Data Availability Statement**

435 The original contributions presented in the study are included in the article/supplementary material,
436 further inquiries can be directed to the corresponding author/s. The raw data supporting the
437 conclusions of this article will be made available by the authors, without undue reservation and is
438 also available on Proteome Xchange Consortium via the PRIDE partner repository (dataset
439 PXD051974).

440 **Conflicts of Interest**

441 JFH has received research support, paid to Northwestern University, from Gilead Sciences, and is a
442 paid consultant for Merck. The authors declare that the research was conducted in the absence of any
443 commercial or financial relationships that could be construed as a potential conflict of interest.

444 **Funding**

445 This work was supported by funding from the NIH (NIAID K23AI139337 to LBM and NHLBI
446 K23HL093302 to KGM), Gerber Foundation, Friends of Prentice, and Thrasher Research Fund.
447 Additional support was provided by Northwestern University Clinical and Translational Sciences
448 Institute (UL1TR001422), Perinatal Origins of Disease Research Program at Lurie Children's, and
449 the NUCord Biorepository. Salary support for JFH and TL-H was provided by NIH/NIAID grants
450 (R21AI163912, R01AI165236, R01AI150455, R01AI150998, and U19AI135964) and additional
451 institutional support for the Center for Pathogen Genomics and Microbial Evolution. The funding
452 sources had no role in the study design, data collection, analysis, interpretation, or writing of the
453 report.

454

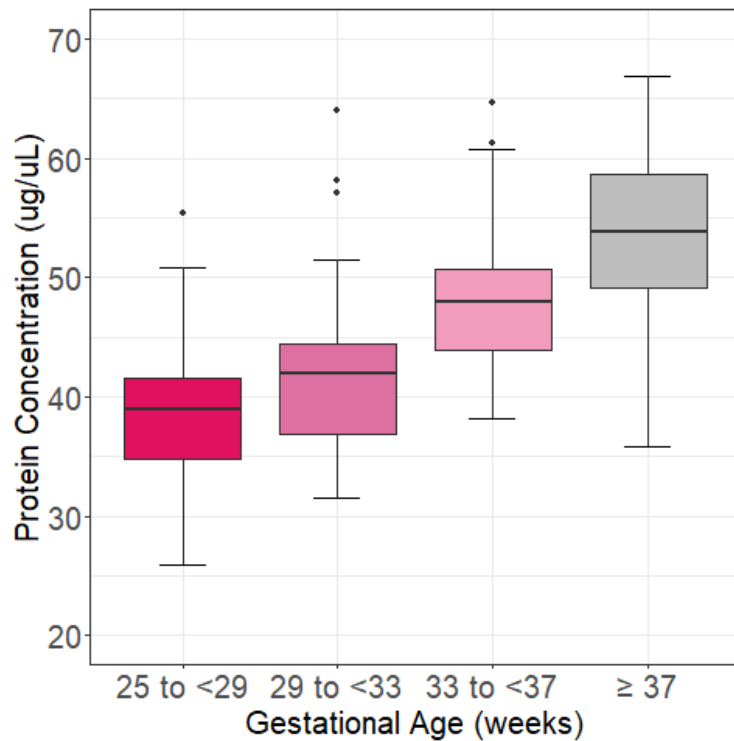
455

456 **Ethics Statement**

457 Northwestern University and Lurie Children’s Institutional Review Boards reviewed and approved
458 the studies involving human participants (Northwestern STU00201858, Lurie IRB 2018-2145). The
459 patients/participants provided their written informed consent to participate in this study.

460

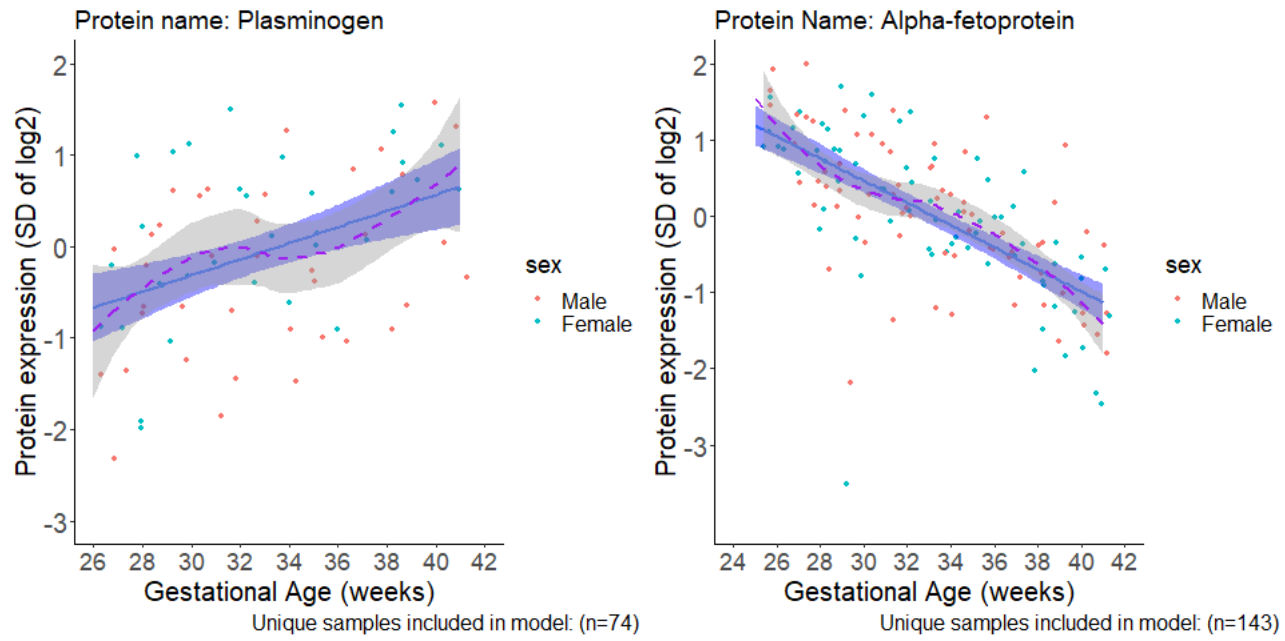
461 **Figures:**



462

463 **Figure 1. Total protein concentration in each plasma sample.** Box and whisker plot of protein
464 concentration (ug/uL) distribution across gestational age (GA) categories. The lower and upper ends
465 of each box correspond to the 25th and 75th percentiles for a given group [shaded area is the
466 interquartile range (IQR)]. The black line in each box is the median. The whiskers represent the
467 largest and smallest observed data points that are no further than +/-1.5 times the IQR, respectively.
468 Points outside of the boundary of the whiskers are outliers. Kruskal-Wallis across GA categories
469 $p < 0.0001$.

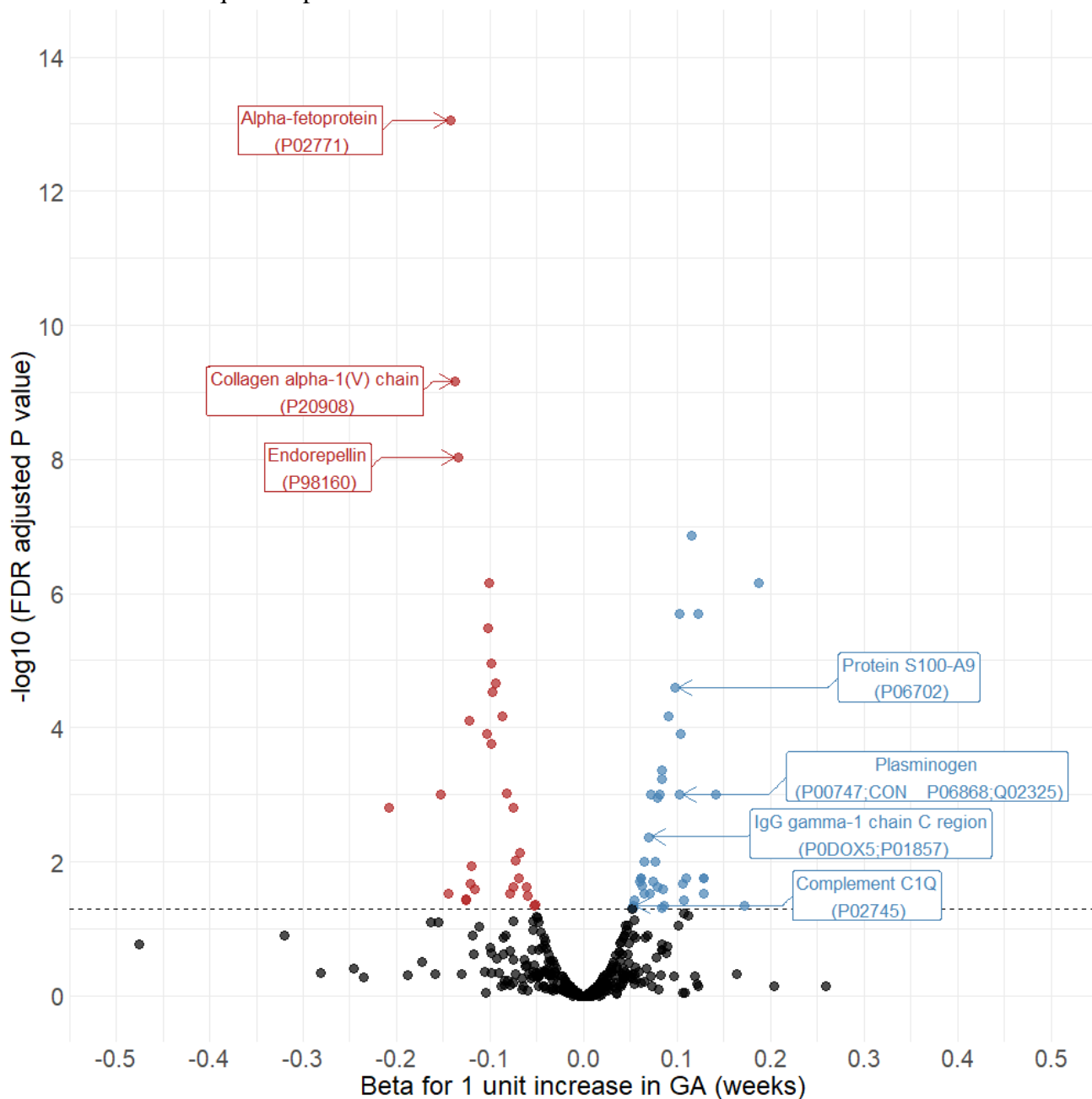
470



471

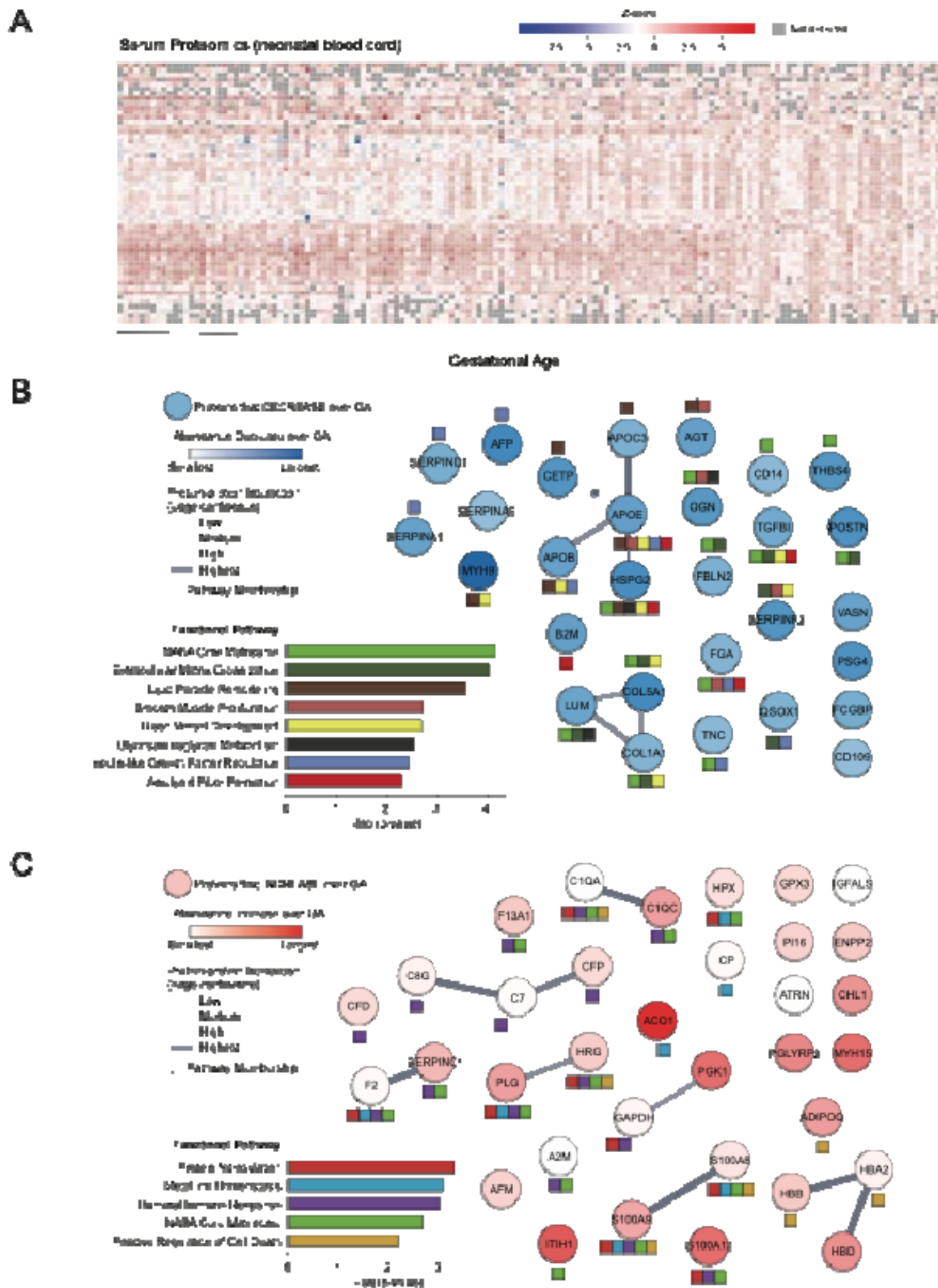
472 **Figure 2. Representative plots of relative protein abundance by gestational age.** The scatter
473 plots show the observed values by sex of newborn. The blue line and shaded blue show the fitted
474 linear model and 95% confidence interval (CI) of the association respectively. The purple dashed line
475 shows a loess smoothed line of the association, and the 95% CI is the shaded gray region (most
476 appeared approximately linear). **A)** Plasminogen model included n=74 samples. **B)** Alpha-fetoprotein

477 included n=143 unique samples in the model.



478

479 **Figure 3. Volcano plot of protein abundance association with gestational age.** Shown are $-\log_{10}$
480 of the FDR adjusted P values and betas from linear regression models for a unit increase in
481 continuous gestational age term on a standard deviation increase in protein abundance for a given
482 protein adjusted for sex, preeclampsia, labor route of delivery. Colors show direction of linear
483 associations (positive [blue] indicates increasing GA associated with increasing protein abundance
484 and negative [red] indicates decreasing GA associated with decreasing protein abundance). Seventy
485 proteins were found to be significantly associated with gestational age in adjusted models. Proteins
486 with FDR adjusted $p < 0.0001$ are labelled. The full list of proteins with FDR adjusted $p < 0.05$ can be
487 found in Supplementary Table 1.



488

489 **Figure 4. Functional enrichment analysis of proteins in neonatal cord blood that change over**
 490 **gestational age. A) Heatmap of the protein Z-scores detected in neonatal cord blood samples**

491 arranged by gestational age from left to right. Proteins are grouped top to bottom by hierarchical
492 clustering. Functional enrichment analysis and protein-protein interaction networks of proteins
493 significantly **B**) decreased or **C**) increased over gestational age are shown below. Network nodes are
494 shaded by abundance change over gestational age with edge width reflecting protein-protein
495 interaction confidence. Significantly enriched pathways are highlighted in colored bar charts to the
496 left; each protein that maps to the identified pathways is indicated by a color-matched box beneath
497 the network node.

498 **Table 1. Demographics and clinical covariates**

| n=150 | Patients n (%) |
|--|---------------------------|
| Gestational age weeks median (IQR) | 33.7 (29.6-37.5) |
| 25-28 | 17 (11%) |
| 29-32 | 43 (29%) |
| 33-36 | 50 (33%) |
| ≥37 | 40 (27%) |
| Infant sex (female) | 77 (51%) |
| Labor and delivery | |
| vaginal with labor | 77 (51%) |
| caesarean with labor | 42 (28%) |
| caesarean without labor | 31 (21%) |
| Preeclampsia | 34 (23%) |
| Multiple gestation | 44 (29%) |
| Clinical chorioamnionitis | 1 (0.7%) |

499