

1 **Small Sample-Big Data: Integrative Indexed Systems Biology**  
2 **Reveals Dramatic Molecular Ontogeny over the First Week of**  
3 **Human Life**

5 **ABSTRACT**

6 Systems biology provides a powerful approach to unravel complex biological processes  
7 yet it has not been applied systematically to samples from newborns, a group highly  
8 vulnerable to a wide range of diseases. Published methods rely on blood volumes that  
9 are not feasible to obtain from newborns. We optimized methods to extract  
10 transcriptomic, proteomic, metabolomic, cytokine/chemokine, and single cell immune  
11 phenotyping data from <1ml of blood, a volume readily obtained from newborns.  
12 Furthermore, indexing to baseline and applying innovative integrative computational  
13 methods that address the challenge of few data points with many features enabled  
14 identification of robust findings within a readily achievable sample size. This approach  
15 uncovered dramatic changes along a stable developmental trajectory over the first week  
16 of life. The ability to extract information from ‘big data’ and draw key insights from such  
17 small sample volumes will enable and accelerate characterization of the molecular  
18 ontogeny driving this crucial developmental period.

19

20 **INTRODUCTION**

21 The first week of life is characterized by heightened susceptibility to infections and a  
22 major determinant of overall health for the entire human lifespan<sup>1-3</sup>. Knowledge of the  
23 molecular drivers involved in these processes in newborns (defined as those  $\leq 28$  days  
24 of life) is fragmentary<sup>4</sup>. In this context, systems biology methods employing high-  
25 dimensional molecular and cellular measurements along with unbiased analytic  
26 approaches enable the identification of molecular networks and signatures relevant to  
27 host defense. Although systems biology approaches have already increased  
28 understanding of baseline and altered molecular states in adults<sup>5</sup>, such approaches  
29 have yet to be applied systematically to characterize molecular ontogeny in early life<sup>6</sup>.  
30 This likely relates to the limited amount of biosample that can be obtained<sup>6, 7, 8</sup> and the  
31 many rapid physiological changes around birth<sup>1, 2</sup>, resulting in variance of biological  
32 measurements necessitating a large sample size, which increases complexity and cost<sup>9</sup>.  
33

34 To overcome these hurdles, we developed a robust experimental and analytical  
35 approach feasible with <1ml of newborn blood. Our data represent the most  
36 comprehensive systems biology study yet performed. We found that despite substantial  
37 between-subject variation, normalizing ('indexing') all samples from a given newborn  
38 enabled identification of robust within-subject changes over the first week of life across  
39 the entire cohort. Furthermore, data integration using multiple independent strategies  
40 confirmed signatures across methodologically- and biologically-distinct datasets. Such  
41 convergence of signals assessing different layers of information validated key findings  
42 internally. The results highlight that contrary to the relatively biological steady-state  
43 observed in healthy adults<sup>7, 10</sup>, the first week of human life is highly dynamic; however,  
44 despite these dramatic changes, ontogeny appears to follow a common path. This  
45 method now enables future studies to focus on the impact of perturbations such as  
46 immunization and infection on this developmental trajectory.

47 **RESULTS**

48 **Blood Processing SOP**

49 A primary objective of this project was to develop a robust standard operating procedure  
50 (SOP) to allow the extraction and analysis of data using systems biology ('big data')  
51 approaches from small blood sample volumes that can be obtained for research study  
52 purposes from newborns and infants (Figure 1; see also 'Protocol'). Our experimental  
53 SOP utilized important sample-sparing modifications whereby we obtained samples for  
54 immune phenotyping, transcriptomic, proteomic, and metabolomic analysis from <1ml of  
55 blood (see also Supplementary Text)<sup>11, 12</sup>.

56 **Immune Phenotyping across 1<sup>st</sup> week of life**

57 Determining the cellular composition of blood samples has been recognized as crucial in  
58 systems biology studies, not only because the relative and absolute cell numbers predict  
59 e.g. vaccine responses with high accuracy<sup>13</sup>, but also because knowledge about their  
60 relative proportion helps deconvolute OMICS data<sup>9, 14</sup>. Analysis of our pre-defined  
61 targeted cell populations revealed substantial between-subject variability  
62 (Supplementary Figure 1, Supplementary Text, and Supplementary MIFlowCyt  
63 document). However, consistent within-subject changes over the first week of life

64 amongst the entire cohort of 30 newborns emerged when samples were indexed to the  
65 average value obtained from all time points for a given subject. Principle component  
66 analysis (PCA) identified the following discriminating features over the first week of life:  
67 Basophils, plasmacytoid dendritic cells (DC), natural killer cells, and neutrophils  
68 decreased, while myeloid DCs increased from Day of Life 0 (DOL0). We also detected  
69 dramatic but consistent changes in soluble immune markers including plasma cytokines  
70 and chemokines over the first week of life (Supplementary Figure 2). Based on the  
71 component coefficients from the relevant PCA, we found that plasma concentrations of  
72 ADAM11, CXCL10, IL-17A, IFNy, and Flt3L increased, while IL-10, CCL5, granulocyte  
73 colony stimulating factor 2 (GCSF), IL-6, and TGFr decreased with age.  
74

75 The ability to reveal a robust developmental trajectory by indexing each subject was  
76 surprising, given the biological heterogeneity across subjects at any given time point.  
77 We thus applied a similar analytic approach to the other OMICS analyses as well.

### 78 **Transcriptomic analysis across 1<sup>st</sup> week of life**

79 Analysis of gene expression by mRNA quantification is common to most systems  
80 biology studies <sup>7, 15</sup>. We found that we needed ≥ 500 µl of adult blood to consistently  
81 obtain sufficient high-quality RNA, but as little as 100 µl of newborn blood sufficed  
82 (Supplementary Figure 3A and 3B). This likely reflects that newborn blood has a  
83 relatively high content of WBC as well as nucleated RBC, which contain abundant globin  
84 mRNA (Supplementary Figure 3C)<sup>16</sup>. The higher yield and quality of total RNA extracted  
85 from newborn vs. adult whole blood was confirmed across different RNA extraction  
86 platforms (RNALater and PAXGene; both of which yielded similar results (data not  
87 shown)). We chose to conduct the rest of our study using RNALater (Supplementary  
88 Text).

89 Similar to the immune phenotyping data, we detected substantial between-subject  
90 variability in our RNA-Seq data (Supplementary Figure 3D-G). This was resolved by  
91 indexing each subject to baseline, which allowed dramatic yet consistent developmental  
92 signals across the entire cohort to emerge that relate to age (i.e. ontogeny). In  
93 comparing DOL1 vs. DOL0, there were few (12) differentially expressed (DE) genes,  
94 however, dramatic developmental changes started to appear when comparing later days  
95 of life to DOL0 (Figure 2). Specifically, for DOL3 vs. DOL0 we detected 1125 DE genes,  
96 while on DOL7 vs. DOL0, there were 1864 DE genes. All DE genes, pathway  
97 enrichment and statistics are listed in Supplementary Table 1 and Supplementary Text.  
98 In particular, genes with decreased expression as a function of age are involved with  
99 cellular responses to stress, detoxification of reactive oxygen, as well as heme  
100 biosynthesis and iron uptake. Conversely genes involved in interferon signaling,  
101 negative regulation of RIG-I and complement activation were up-regulated over the first  
102 week of life.  
103

### 104 **Proteomic analysis across 1<sup>st</sup> week of life**

105 We analyzed changes in the plasma proteome across the first week of life utilizing 5 µl  
106 of plasma in a 96-well plate format, allowing large numbers of samples to be analyzed in  
107 parallel. A total of 684 different proteins were identified across samples (false discovery  
108 rate (FDR) <1%). Of these, 199 proteins met our criteria (which included the detection of

109 at least two unique peptides per protein) for further detailed quantification. Substantial  
110 between-subject variability was again noted in the plasma proteomic analysis  
111 (Supplementary Figure 4A), but indexing each subject to baseline enhanced the  
112 detection of signatures that again indicated a common developmental trajectory over the  
113 first week of life, with differences in plasma protein composition compared to DOL0  
114 increasing with increasing age (Supplementary Figure 4B). Differentially abundant  
115 plasma proteins and their respective pathways are listed in Supplementary Table 2 and  
116 Supplementary Text. At DOL3 vs. DOL0, three pathways were up-regulated that center  
117 around the complement cascade. At DOL7 vs. DOL0, five additional pathways were up-  
118 regulated including scavenging heme from plasma and signaling to RAS.

### 119 **Metabolomic analysis across 1<sup>st</sup> week of life**

120 We used a mass spectrometry-based global discovery metabolomics approach to  
121 analyze ~50 µl of plasma, which enabled detection of ~700 metabolites (Supplementary  
122 Figure 5A). While initial analysis revealed substantial between-subject variation,  
123 transforming the analysis from a between-subject to a within-subject comparison over  
124 the first week of life revealed a steady but dramatic developmental trajectory  
125 (Supplementary Figure 5B and C). Few differences in plasma metabolites were  
126 discernible comparing DOL1 vs. DOL0, but increasing differences were noted when  
127 contrasting DOL3 and DOL7 vs. DOL0; interestingly the differences detected focused on  
128 metabolic pathways related to sphingolipid biosynthesis, carbohydrate and androgen  
129 metabolism likely reflecting rapid cell proliferation in newborns (Supplementary Table 3  
130 and Supplementary Text).

### 131 **Data integration**

132 Each methodologically- and biologically-distinct data type we examined revealed  
133 substantial changes over the first week of life. We next sought to determine if the  
134 observed changes were related to one another across data types, reflecting consistent  
135 age-dependent changes in functional pathways. In order to minimize limitations inherent  
136 in any single analytical approach, and to detect the most robust signatures, we  
137 addressed data integration using three independent strategies. To this end, we  
138 employed a function-based strategy based on biologically known Molecular Interactions  
139 Networks using NetworkAnalyst<sup>17</sup>; the data-driven (unbiased) multivariate matrix  
140 factorization approach DIABLO (Data Integration Analysis for Biomarker discovery using  
141 Latent cOmponents)<sup>18-20</sup>; and the previously described multi-scale, multifactorial  
142 response network (MMRN) that estimates correlations across data types<sup>7,21</sup>.

### 143 Molecular Interaction Networks using NetworkAnalyst

144 NetworkAnalyst enables the creation of networks based on a framework of known PPI  
145 captured in publically curated databases (specifically InnateDB/IMeX)<sup>22</sup>. Minimum  
146 connected networks were constructed from seed nodes (i.e. from genes or proteins that  
147 changed with age in our data set), as well as first-order interactors that served to  
148 connect the seed nodes with each other. To include metabolomics data, metabolic  
149 enzymes (synthetic and degradative) that would determine the levels of differentially  
150 detected metabolites were used as seed nodes in the network construction. Overall  
151 metabolomics, proteomics and transcriptomics data fit well into a single functional  
152 network (Figure 3), indicating that these techniques reported on different facets of the

153 same biological processes. This PPI-based integration strategy recapitulated many key  
154 findings that had been identified for each of the individual data types, confirming our  
155 expectation that many but not all findings would be validated by different OMICS  
156 methods (Supplementary Table 4). For example, integrating transcriptomic with  
157 proteomic data confirmed the rapid increase in type 1 interferon (IFN)-related functions  
158 over the first week of life. Our integration however also revealed new biological insights  
159 not revealed by any single-data domain analysis. For example, only following the  
160 integration of transcriptomic and metabolomic data did we identify down-regulation of  
161 very long-chain fatty acyl-CoA synthesis activities across the first week of life.

## 162 DIABLO

163 DIABLO is a multivariate approach to address two of the major concerns faced when  
164 integrating multilevel datasets: the dimensionality of the data, particularly with few  
165 samples, each with many observations, and the heterogeneous nature of our data  
166 measured on different scales and technological platforms<sup>18-20</sup>. DIABLO constructs  
167 components (linear combinations of the original features) that are maximally correlated  
168 across any number of input data types and a specified outcome variable (in this case,  
169 DOL), while simultaneously performing feature selection via L1 penalization<sup>23</sup>  
170 (Supplementary Text). We created matrices from our five data types as input to DIABLO  
171 to identify major ontogeny-related features (Figure 4 and Supplementary Table 5). The  
172 resulting model discriminated well between DOLs, with component one of the model  
173 separating birth (DOL0) from all other time points, while component two separated  
174 DOL1, 3, and 7 from each other (Figure 4C).

175

176 We next investigated the relationship between features selected by DIABLO across data  
177 types and visualized the selected features in an integrative network (Supplementary  
178 Text). We compared this integrative network (Figure 4A vs. 4B) to one derived from  
179 features identified using a non-integrative sparse discriminant analysis approach  
180 (Supplementary Text). The integrative network was more densely connected (global  
181 clustering coefficient = 0.91 vs. 0.68) and composed of few, more tightly connected  
182 modules (network modularity = 0.26 vs. 0.09), indicating that DIABLO selected features  
183 that were discriminant and well correlated across data types, while the non-integrative  
184 approach favored features that were discriminant but not well correlated across data  
185 types. The two components of the DIABLO model were composed of distinct sets of  
186 features (Figure 4D, blue bars), representing distinct biology (Figure 4E, blue bars;  
187 Supplementary Text).

## 188 MMRN

189 Multi-scale, multifactorial response networks (MMRNs) are a recently published  
190 framework for data integration<sup>7</sup>. Using MMRN, we found that associations between data  
191 types were strongest at DOL1 and decreased across the first week of life, as partial  
192 least squares regression scores were significantly higher at DOL1 compared to all other  
193 time points (Student's t-test, p-value << 0.01). Stable clusters were more strongly  
194 associated with DOL when compared to transient ones (Figure 5C and Supplementary  
195 Text), with 15/21 clusters significantly associated with DOL being part of stable networks.  
196 This confirmed our already noted robust trajectory of development.

197

198 Most significant clusters were transcriptomic (16/21), but we also identified metabolomic  
199 (1/21) and flow cytometry-derived (4/21) clusters associated with DOL (Supplementary  
200 Table 6). The stable clusters most significantly associated with DOL were composed of  
201 various B-cell sub-populations, blood transcriptomic modules (BTMs) related to dendritic  
202 cells and monocytes, CCR1, CCR7, cell signaling, heme biosynthesis, TLR and  
203 inflammatory signaling, as well as metabolic pathways such as purine metabolism  
204 (Figure 5D)<sup>7</sup>.

## 205 META-INTEGRATION

206 To assess convergence of signals across methodologically-distinct integration methods,  
207 we simplified the output of each integration approach to a list of features associated with  
208 DOLs. We then carried out gene set enrichment to determine which biological processes  
209 were identified by each method (Supplementary Tables 7 and 8). The selected features  
210 for Molecular Interactions Networks (NetworkAnalyst) were the component nodes of  
211 minimum connected networks when differentially abundant features (transcripts,  
212 proteins and metabolites) for DOL3 vs. DOL0 and DOL7 vs. DOL0 were used as seed  
213 nodes (2317 features). For DIABLO, the selected features were the component nodes of  
214 minimum connected molecular interaction networks when the features of the two model  
215 components were used as seed nodes (184 features). For MMRN, selected features  
216 were those that composed the small stable network shown in Figure 5D (564 features).  
217 We observed limited overlap at the individual feature level (Figure 6A). However,  
218 assessing gene set enrichment against the Reactome annotation system (674  
219 pathways), 472, 251 and 122 pathways were significantly over-represented in the  
220 NetworkAnalyst, DIABLO and MMRN feature lists, respectively (all with Benjamini-  
221 Hochberg corrected FDR  $\leq 0.05$ , Supplementary Table 7). Importantly, 162 and 26  
222 pathways were identified by NetworkAnalyst and either DIABLO or the MMRN,  
223 respectively, and 72 pathways were identified by all 3 approaches, demonstrating  
224 convergence of signals across different analytical platforms (Figure 6B). This degree of  
225 overlap was unlikely to occur by chance alone as determined by bootstrapping (p-value  
226  $< 0.001$ ; Supplementary Text). The specific pathways identified by meta-integration  
227 indicated as driving molecular ontogeny over the first week of life were interleukin 1  
228 signaling, Toll-like receptor signaling, the NOTCH signaling axis, the DHX RNA helicase  
229 pathway connecting to type 1 interferon gene expression, as well as regulators of the  
230 complement cascade (Supplementary Table 8). Given the striking convergence  
231 observed, we next assessed relevant functional interactions between the in silico  
232 identified pathways. To this end, we used NetworkAnalyst and found that 80% of the  
233 selected common features fit into minimum connected networks of experimentally-  
234 validated interactions (Figure 6C). This suggests that our meta-integration approach  
235 identified functional biological interactions relevant to neonatal ontogeny.

## 236 DISCUSSION

237 Here we present a holistic suite of complementary methods that address key hurdles of  
238 applying systems biology to newborns in the following ways: *i*) overcome limitations in  
239 sample volume via an efficient, field-compatible, sample-sparing SOP to process  
240 peripheral whole blood; *ii*) determine the developmental trajectories from each subject's  
241 own baseline (indexing) and *iii*) reduce the dimensionality of the dataset by integrating  
242 methodologically- and biologically-distinct data types (multi-OMICS integration), which

243 analytically validated signatures of important pathways. This approach revealed for the  
244 first time a dramatic molecular ontogeny over the first week of human life.

245  
246 Multi-OMICS integration poses an exceptional problem with the increase in  
247 dimensionality ( $p$  features) relative to the typically achievable sample size ( $n$  samples).  
248 We here addressed this  $p \gg n$  problem by reducing the dimensionality of our data  
249 through selection of only a subset of variables (MMRN, DIABLO)<sup>24</sup>, and to impose  
250 functionally-relevant, known molecular interaction information (NetworkAnalyst)<sup>21</sup>.  
251 Furthermore, as the integrative methods we applied focused on extracting aggregated  
252 information from each data set *independently* of each other, confidence in the results, if  
253 they independently converged on the same set of key molecular features, was high. This  
254 approach confirmed that integrating across multiple technological platforms representing  
255 different biological data types can dramatically increase robust biological insight. For  
256 example, this approach identified already known changes in early life, including those  
257 related to the composition of hemoglobin mRNA, increases in complement protein C9<sup>25</sup>,  
258 and reductions in plasma steroids<sup>26</sup>, supporting the validity of our findings. However,  
259 our approach also identified pathways never before identified as relevant to ontogeny.  
260 For example, many pathways consistently identified by all three of our analytical  
261 strategies had previously been shown to be relevant for host defense, but they have  
262 never before been identified as central to early human development. These pathways  
263 include interleukin 1 signaling, Toll-like receptor signaling, NOTCH signaling, the DHX  
264 RNA helicase pathway, as well as regulators of the complement cascade<sup>1, 27-30</sup>. Data  
265 integration also led to the identification of novel and surprising but biologically plausible  
266 findings regarding ontogeny. For example, prostaglandin-endoperoxide synthase 2  
267 (PTGS2 or PGHS2 or COX-2) appeared as centrally important in all of our integrative  
268 networks, but not in any of the single OMICS data types analyzed individually. PTGS2 is  
269 expressed abundantly in hematopoietic progenitors<sup>31</sup> and is clinically relevant during  
270 premature labor and in necrotizing enterocolitis<sup>32, 33</sup>. Taken together, our findings being  
271 to outline a developmental trajectory that may serve as a reference akin to the stable  
272 steady state in adults<sup>7, 10, 13</sup>. Deviations from this developmental trajectory could  
273 potentially identify subjects at risk prior to the onset of disease or guide the re-  
274 establishment of homeostasis once derailed.  
275

276 In summary, our method presents an integrative analysis across the broadest range of  
277 multivariate data sets published to date. Our integration over time (indexing data  
278 longitudinally for each newborn) and biological space (multi-OMICS integration) allowed  
279 exploration of the dynamic molecular and cellular developmental characteristics of early  
280 life. In parallel, we identified signals of potential physiological importance using a modest  
281 sample size readily obtainable in most studies. Coupled with our field-tested SOP for  
282 sample-sparing pre-analytical processing, we thus provide the methods to overcome key  
283 challenges in applying systems biology to neonates. Our data revealed a compelling  
284 ‘biological narrative’, where out of the apparently noisy age-dependent change across  
285 the first week of life, a consistent developmental trajectory emerged. These approaches  
286 and observations will serve as a crucial backdrop for future studies that characterize the  
287 impact of a broad array of factors, including genetics, epigenetics, maternal influences,  
288 microbiota, diet, and disease, as well as biomedical interventions such as vaccination.

289 **METHODS (ONLINE)**

290 **Peripheral blood processing**

291 30 healthy, term newborns were enrolled at the Medical Research Council (MRC) Unit  
292 The Gambia in accordance with a local Ethics Committee-approved protocol (SCC  
293 1436). Following informed consent, mothers were screened for HIV-I and -II and  
294 Hepatitis B with positivity for either virus representing an exclusion criterion. Inclusion  
295 criteria were a healthy appearing infant as determined by physical examination, born by  
296 vaginal delivery at gestational age of  $\geq$ 36 weeks, 5 minute Apgar scores  $\geq$ 8, and a birth  
297 weight of  $\geq$  2.5 kilograms. Peripheral blood samples were obtained from all infants on  
298 the day of birth (DOL0) and then again either at DOL1, DOL3 or DOL7, in order to  
299 reduce venipuncture to a maximum of twice in the first week of life. Peripheral venous  
300 blood was drawn from infants via sterile venipuncture directly into heparinized collection  
301 tubes (Becton Dickinson (BD) Biosciences; San Jose, CA, USA). Aliquots (200  $\mu$ l) were  
302 immediately placed in RNA-later (Ambion ThermoFisher, Waltham, MA, USA) with the  
303 remaining blood kept in the collection tubes at room temperature until further processing  
304 within 4 hours. All samples were processed at the MRC Unit The Gambia as described  
305 below and subsequently shipped to collaborating laboratories on dry ice, under  
306 temperature controlled and monitored conditions (World Courier, New Hyde Park, NY,  
307 USA).

308 **Indexing**

309 In this study, we profiled the blood of each subject twice over their first week of life, at  
310 DOL0 (baseline) and additionally at either DOL 1, 3, or 7, and sought to identify  
311 variables that differed between the baseline and latter time points across all subjects.  
312 That is, we were interested in the within-subject variation associated with DOL observed  
313 consistently across the entire cohort. For univariate analyses, we considered (indexed)  
314 paired differences, either implicitly (e.g. paired t-test), or, in the case of the flow  
315 cytometry and luminex cytokine data, explicitly by transforming the data beforehand  
316 using a multilevel approach to separate the between- and within-subject variation as  
317 described<sup>34</sup>. The same approach was used for all multivariate analyses.

318 **Immune Phenotyping**

319 *Flow Cytometry (FCM).* Whole blood was centrifuged on site at 500 xg for 10 minutes at  
320 room temp and plasma harvested and stored at -80°C for later analysis of plasma  
321 cytokines, proteins and metabolites. The amount of plasma removed from the whole  
322 blood after centrifugation was subsequently replaced with RPMI. For assessment of  
323 cellular composition by FCM, EDTA (0.2 mM final concentration) was added to the  
324 whole blood/RPMI mixture to ensure adherent cells were not lost. In parallel, cells were  
325 stained with fixable viability dye (FVD) at 4°C for 15 min prior to red blood cell lysis  
326 followed by storage at -80°C in Smart Tube reagents (Smart Tube Inc.; San Carlos, CA,  
327 USA). At the immunophenotyping laboratory samples were thawed, washed in staining  
328 buffer (PBSAN; 0.5% BSA, 0.1% sodium azide) and stained on ice in PBSAN with a  
329 cocktail of anchor markers to determine frequency of cell populations contained in  
330 peripheral blood (for list of cell types and anchor markers, see Supplemental Figure 1;  
331 for list of clone/ fluorochrome combination see MiFlowCyt section in Supplemental  
332 Information). Flow cytometric analysis employed a custom-built LSRII (for machine

333 settings and compensation settings see the MIFlowCyt section in Supplemental  
334 Information)<sup>35</sup>. Our gating strategy is shown in the Supplemental Figure 1. FCM data  
335 was analyzed in an automated fashion using R/BioConductor packages (Supplemental  
336 Figure 1). Specifically, *flowCore* supported the analysis in single files according to the  
337 Flow Cytometry Standard (FCS), providing the infrastructure to support sub-setting of  
338 data, data transformations and gating<sup>36, 37</sup>. Cell population identification was then  
339 conducted using *flowDensity*, a supervised gating tool, that was customized to provide  
340 threshold calculations designed for each cell subset based on expert knowledge of  
341 hierarchical gating order and one-dimensional density estimation<sup>38</sup>. Lastly,  
342 *flowType/RchyOptimyx* identified cell populations that correlated with outcome, in this  
343 case DOL at the time of blood draw<sup>39</sup>. *flowType* uses partitioning of cells, either  
344 manually or by clustering, into positive or negative for each marker to enumerate all cell  
345 types in a sample. *RchyOptimyx* measures the importance of these cell types by  
346 correlating their abundance to external outcomes, such as DOL, and distils the identified  
347 phenotypes to their simplest possible form.

348  
349 *Luminex*. Plasma (25 µl) was used to measure cytokine concentrations using a custom-  
350 designed multi-analyte Cytokine Human Magnetic Panel bead array, (Invitrogen/Life  
351 Technologies, Carlsbad; CA) consisting of CCL2, CCL3, CCL5, CXCL8, CXCL10, GM-  
352 CSF, IFN- $\alpha$ 2, IL-10, IL-12p40, IL-12p70, IL-1 $\beta$ , IL-6, and TNF $\alpha$ . Results were obtained  
353 with a Flexmap 3D system with Luminex xPONENT software version 4.2 (both from  
354 Luminex Corp.; Austin, TX, USA). Cytokine concentrations were determined using  
355 Milliplex Analyst software (version 3.5.5.0, Millipore).

356  
357 *Immunophenotypic Analysis*. For flow cytometric as well as Luminex raw values were  
358 normalized with a 1 + Log2-transformation. WithinVariation matrices were computed for  
359 each data matrix using the WithinVariation function in R package *mixOmics* version  
360 6.1.2. The Wilcoxon rank-sum test (*wilcox.test* in base R) was used to determine  
361 differentially regulated features within each data type, using the WithinVariation values  
362 for each feature. P-values were adjusted for each data type separately using the  
363 Benjamini-Hochberg method (*p.adjust* function, base R). Features were considered  
364 statistically different from DOL0 if their adjusted p-values were below 0.1. All analyses  
365 were performed in R version 3.3.2 (2016-10-31).

### 366 **RNA-Seq**

367 Total RNA was extracted from each sample using the RiboPure RNA purification kit  
368 (Ambion ThermoFisher; Waltham, MA, USA) following the manufacturer's protocol.  
369 Quantification and quality assessment of total RNA was performed using an Agilent  
370 2100 Bioanalyzer (Santa Clara, CA, USA). Poly-adenylated RNA was captured using  
371 the NEBNext Poly(A) mRNA Magnetic Isolation Module (catalogue no.: E7409L, NEB;  
372 Ipswich, MA, USA). Strand-specific cDNA libraries were generated from poly-adenylated  
373 RNA using the KAPA Stranded RNA-Seq Library Preparation Kit (cat. no.: 07277253001,  
374 Roche; Basel, Switzerland). All cDNA libraries were prepared at the same time and  
375 sequenced on the HiSeq 2500 (Illumina; San Diego, CA, USA), using one Rapid v2 and  
376 two lanes of High Output single-read run of 100 bp-long sequence reads (+  
377 adapter/index sequences). Sequence quality was assessed using FastQC v0.11.5 and  
378 MultiQC v0.8.dev0<sup>40</sup>. The FASTQ sequence reads were aligned to the hg19 human

379 genome (Ensembl GRCh38.86) using STAR v2.5 and mapped to Ensembl GRCh38  
380 transcripts<sup>41</sup>. Read-counts were generated using htseq-count (HTSeq 0.6.1p1)<sup>42</sup>. All  
381 data processing and subsequent differential gene expression analyses were performed  
382 using R version 3.3.0 and DESeq2 version 1.14.1<sup>43</sup>. Genes with very low counts (with  
383 less than 10 counts in eight or more samples, or the smallest number of biological  
384 replicates within each treatment group) and globin transcripts were pre-filtered and  
385 removed *in silico*. Differentially expressed genes were identified using paired analysis  
386 with the Wald statistics test and filtering for any genes that showed 2-fold change and  
387 adjusted p-value < 0.05 (cut-off at 5% FDR) as the threshold. Functional discovery of  
388 pathway enrichment and network analyses was performed using Sigora 2.0.1 and  
389 NetworkAnalyst, respectively<sup>17, 44</sup>.

### 390 **Plasma Proteomics**

391 Plasma samples were prepared for proteome analysis using the in-house (Boston  
392 Children's Hospital; Steen laboratory)-developed MStern blotting sample processing and  
393 trypsinization protocol<sup>45</sup>, which was adapted for plasma samples<sup>46</sup>. To this end, 5 µL  
394 plasma was first diluted in 100 µL sample buffer (8 M urea in TRIS-HCl, pH 8.5). Protein  
395 disulfide bonds were then reduced with dithiothreitol (10 mM final concentration) for 30  
396 min, and alkylated with iodoacetamide (50 mM final concentration) for 30 min in sample  
397 buffer. Three µL (approximately 10 µg) of this protein solution was then transferred to a  
398 96 well plate with a polyvinylidene fluoride (PVDF) membrane at the bottom. Protein  
399 digestion was performed with sequencing-grade modified trypsin (V5111, Promega;  
400 Madison, WI, USA) at a nominal protease to protein ratio of 1:25 w/w. After incubation  
401 for two hours at 37°C, the peptides were eluted from the PVDF membrane, and  
402 concentrated to dryness in a vacuum centrifuge. To monitor retention time stability and  
403 system performance, iRT peptides (Biognosys, Schlieren, Switzerland) were spiked into  
404 all samples. Samples were analyzed using a nanoLC system (Eksigent; Dublin, CA)  
405 equipped with a LCchip system (cHiPLC nanoflex, Eksigent) coupled online to a Q  
406 Exactive mass spectrometer (Thermo Scientific; Bremen, Germany). From each sample,  
407 0.2 µg peptide material was separated using a linear gradient from 93% solvent A (0.1%  
408 formic acid in water), 7% solvent B (0.1% formic acid in acetonitrile) which was  
409 increased to 32% solvent B over 60 min. The mass spectrometer was operated in data-  
410 dependent mode, selecting up to 12 of the most intense precursors for fragmentation  
411 from each precursor scan. Label-free protein quantitation analysis employed MaxQuant  
412 1.5.3.30<sup>47</sup>. Raw-data were downloaded and used to build a matching library and  
413 searched against the UniProt Human Reference Proteome as described<sup>48</sup>. Standard  
414 search settings were employed with the following modifications: Max missed cleavage 3;  
415 variable modifications Deamidation (NQ) and Oxidation (M)<sup>49</sup>. A revert decoy search  
416 strategy was employed to filter all proteins and peptides to < 1% FDR<sup>50</sup>. The list of  
417 proteins was further processed in Perseus 1.5.5.3, log2-transformed, and proteins with  
418 less than two peptides (razor) were filtered out as described<sup>51</sup>. The samples were  
419 grouped according to DOL, and proteins which were not quantifiable in at least five of  
420 the samples in any day were filtered out. Remaining missing values were imputed using  
421 numbers drawn from a normal distribution with the standard parameters in Perseus  
422 (width 0.3, downshift 1.8) to simulate signals from low abundant proteins<sup>52</sup>. The R-script  
423 ComBat was used to correct for batch effects for samples run on different LC-MS  
424 columns<sup>53</sup>. Proteins with a statistically significant change of abundance between

425 different DOL were identified by paired two-samples t-test. To correct for multiple  
426 hypothesis testing, permutation-based false positive control was applied using standard  
427 parameters in Perseus (FDR=0.05, s0=0.1)<sup>54</sup>. Significant proteins were further analyzed  
428 in Cytoscape, SIGORA<sup>44</sup> with the Reactome<sup>55</sup> gene annotation system and DAVID<sup>56</sup>.

#### 429 **Plasma Metabolomics**

430 Plasma samples were analyzed using the non-targeted metabolomics platform of  
431 Metabolon Inc. (Durham, NC, USA). Samples were extracted and prepared for analysis  
432 using Metabolon's solvent extraction method<sup>57</sup>. Each plasma sample was stored at -  
433 80°C and accessioned into the Metabolon Laboratory Information Management System  
434 (LIMS). Recovery standards were added at the first step in the extraction process to  
435 allow quality control for the entire process. Protein removal employed methanol  
436 precipitation (Glen Mills GenoGrinder 2000) followed by centrifugation. The  
437 supernatants were divided into five fractions: two for analysis by two separate reverse  
438 phase (RP)/UPLC-MS/MS methods with positive ion mode electrospray ionization (ESI),  
439 one for analysis by RP/UPLC-MS/MS with negative ion mode ESI, one for analysis by  
440 HILIC/UPLC-MS/MS with negative ion mode ESI, and one sample was reserved for  
441 backup using Waters ACQUITY ultra-performance liquid chromatography (UPLC) and a  
442 Thermo Scientific Q-Exactive high resolution/accurate mass spectrometer interfaced  
443 with a heated electrospray ionization (HESI-II) source and Orbitrap mass analyzer  
444 operated at 35,000 mass resolution. Compounds were identified by comparison to  
445 Metabolon library entries of purified standards or recurrent unknown entities.  
446 Furthermore, biochemical identifications were based on three criteria: retention index  
447 (RI) within a narrow RI window of the proposed identification, accurate mass match to  
448 the library ± 10 ppm, and MS/MS forward and reverse scores between the experimental  
449 data and authentic standards found in the LIMS. MS/MS scores were based on a  
450 comparison of the ions present in the experimental spectrum to the ions present in the  
451 library spectrum. Exact molecular mass data from redundant m/z peaks corresponding  
452 to the formation of different parent and product ions were first used to help confirm the  
453 metabolite molecular mass. Metabolon's MassFragmentTM application manager  
454 (Waters MassLynx v4.1, Waters corp.; Milford, USA) was used to facilitate the MS/MS  
455 fragment ion analysis process using peak-matching algorithms and quantified using an  
456 area-under-the-curve in-house algorithm. Finally, the identities of the specific  
457 metabolites were confirmed by comparing their mass spectra and chromatographic  
458 retention times with commercially available reference standards.  
459 Metabolomic data were then pre-processed to original scale values by Bradford protein  
460 assay of cellular protein levels and to volume extracted values obtained on a per sample  
461 basis. Data was normalized to the interquartile range, in which the metabolites with an  
462 IQR of zero are uninformative and can be deleted from further analyses. Missing values  
463 below the threshold level were imputed with the minimum value obtained for the  
464 respective biochemical from within-all-groups. Each biochemical value was then  
465 transformed using generalized logarithmic transformation for distributional assumptions  
466 of chemometrics. Statistical analyses were performed using R statistical packages and  
467 MetaboAnalyst 3.0<sup>58</sup>. Paired t-tests with p ≤0.05 were employed to identify altered  
468 metabolites between paired samples. Regulation of individual metabolites was  
469 compared across groups and time using repeated measures ANOVA and generalized  
470 linear mixed model (GLMM) for longitudinal models. Multiple comparisons were

471 accounted for by estimating the false discovery rate (FDR) using q-values<sup>59</sup>.  
472 Significantly altered pathways were determined by pathway set enrichment analysis  
473 within Metabolon Pathway Analysis (MPA) software which was determined by the  
474 following equation: # of significant metabolites in pathway ( $k$ )/total # of detected  
475 metabolites in pathway ( $m$ )/total # of significant metabolites ( $n$ )/total # of detected  
476 metabolites ( $N$ ) or  $(k/m)/(n/N)$ . Pathways with a higher number of experimentally  
477 regulated compounds relative to the overall study in a follow-up day (e.g., DOL7 vs.  
478 DOL0) received a pathway impact score > 1 suggesting that these pathways may be of  
479 interest to the metabolic perturbations observed. Pathway analysis and visualization  
480 employed the KEGG Pathway<sup>60</sup>, The Human Metabolome Database Version 3.6  
481 (HMDB)<sup>61</sup> and Metabolnc Pathway Analysis using Cytoscape plugin<sup>62</sup>. In addition,  
482 metabolite IDs were converted to HMDB<sup>61</sup>, KEGG<sup>60</sup>, and PubChem IDs via the  
483 Metaboanalyst ID mapping tool<sup>58</sup>. Metabolites were mapped to their directly interacting  
484 enzymes via BioCyc (executed in R)<sup>63</sup>. Metabolites were also mapped to their directly  
485 interacting enzymes via KEGG (executed in R via functional web scraping). The union of  
486 genes from the previous two steps was then used to proceed with metabolome data  
487 integration. Two pathway mapping databases were used, as neither appeared to contain  
488 a complete list of all metabolites and reactions. The list of proteins derived from the  
489 metabolite data processing (above) and lists of differentially expressed genes or  
490 proteins were then combined. Pathway enrichment analysis employed SIGORA<sup>44</sup> using  
491 the Reactome gene annotation system<sup>55</sup>.

## 492 Data integration

493 Data obtained via the immune phenotyping (cellular composition and plasma cytokines),  
494 transcriptomic, proteomic and metabolomic methods was integrated to identify  
495 correlations of signatures across these methodologically- and biologically-distinct  
496 datasets, since convergence of signatures across such diverse biological domains  
497 provides an independent assessment and approaches functional validation<sup>7</sup>. In addition,  
498 we aimed to determine whether we could derive any novel biological information via the  
499 integration of multiple OMICS data types that was not revealed in a single dataset alone.  
500 To cross-validate results we employed 3 data integration platforms, each applying a  
501 different analytical strategy; these independent but complementary data-driven vs.  
502 knowledge/network-driven strategies, pursued in parallel, decreased the chance of false  
503 discoveries. Specifically, data integration strategies included: *i*) a novel method for  
504 integrating multiple OMICS data types into known protein-protein interaction networks  
505 (Molecular Interaction Network) using NetworkAnalyst that provided context based on  
506 annotated molecular interactions<sup>17</sup>; *ii*) a new approach to identify the underlying key  
507 drivers of ontogeny using specialized multivariate methods capable of identifying  
508 relevant features from high-dimensional datasets, namely sparse generalized canonical  
509 correlation discriminant analysis via the '*Data Integration Analysis for Biomarker*  
510 *discovery using a Latent component method for Omics*' (DIABLO) framework, which is  
511 part of the mixOmics R package<sup>18-20</sup>; and *iii*) the recently published multi-scale,  
512 multifactorial response network (MMRN), querying the informatically derived correlations  
513 within a network for statistically significant features<sup>7</sup>. Finally, we also 'integrated the  
514 integration' by combining the protein-protein-interaction driven approach  
515 (NetworkAnalyst) with the feature-mining power of DIABLO ('meta-integration').

516

517 *Molecular Interaction Network.* NetworkAnalyst integrates data using protein-protein  
518 interactions as a biological framework. Metabolomic data required pre-processing in  
519 order to integrate into a network with transcriptomic and proteomic data, as metabolites  
520 must be associated with proteins that are involved in their creation and/or degradation.  
521 To identify such proteins, the following steps were taken: Metabolite IDs were converted  
522 to HMDB, KEGG, and PubChem IDs via the Metaboanalyst ID mapping tool<sup>58</sup>.  
523 Metabolites were mapped to their directly interacting enzymes via BioCyc (executed in  
524 R)<sup>63</sup>. Metabolites were also mapped to their directly interacting enzymes via KEGG  
525 (executed in R via functional web scraping). The union of genes from the previous two  
526 steps was then used to proceed with metabolome data integration. Two pathway-  
527 mapping databases were used, as neither alone appeared to contain a complete list of  
528 all metabolites and reactions. The list of proteins derived from the metabolite data  
529 processing (above) and lists of differentially expressed genes or proteins were then  
530 combined using NetworkAnalyst to produce a zero-order or minimum-connected  
531 network, depending on the size of the dataset<sup>17</sup>. Networks consist of seed nodes  
532 (proteins/genes that were used as input to build the network), novel nodes  
533 (proteins/genes that are 1<sup>st</sup> order interactors of the seed nodes and are used as  
534 “scaffolding” to build the network) and edges (links that join the nodes together and are  
535 indicative of a molecular interaction between nodes). Our focus was on the identification  
536 of novel nodes as emergent information that can be derived from a biological network.  
537 Node lists were then downloaded to identify median degree of connectivity for nodes of  
538 each data type and to identify novel nodes stemming from data integration.  
539 Transcriptomes were integrated via their respective encoded gene products.  
540

541 A minimum connected network for transcriptomic, proteomic and metabolomic data was  
542 constructed using NetworkAnalyst based on differentially present genes/proteins on  
543 each DOL vs DOL0. Novel nodes were identified as those, which were not present in the  
544 set of differentially present genes/proteins originally entered (i.e. not part of the seed  
545 nodes). Note that for DOL3 vs DOL0 and DOL7 vs DOL0, the transcriptomic datasets  
546 were simplified by first producing a zero order network, and only using these nodes for  
547 minimum connected network construction. Next, pairwise lists of each data type were  
548 constructed and minimum connected networks were constructed. Novel nodes in each  
549 network were identified as nodes that were not in either set of differentially present  
550 genes/proteins or either minimum connected network constructed from a single data  
551 type in the pairwise set. A minimum connected network for all data types was then  
552 constructed and novel nodes were identified as nodes that were not in any differentially  
553 present list of genes/proteins/metabolites, any novel nodes in a minimum network  
554 constructed from any single data type, or any novel nodes in a minimum network  
555 constructed from any pairwise data. Pathway enrichment analysis was run using  
556 SIGORA software<sup>44</sup> with the Reactome ontology system on the node lists from each  
557 minimum connected network built during the novel node identification process outlined  
558 above. For each data type alone, (transcriptomics, proteomics and metabolomics), the  
559 number of unique pathways identified were counted as novel pathways. For each  
560 pairwise combination of data types, previously identified pathways that had been  
561 identified in either data type alone were subtracted. For the three-way combination, all  
562 pairwise pathways and unique data type pathways were subtracted.

563  
564 *DIABLO*. Prior to integration with DIABLO, data were transformed. Specifically, cell  
565 proportions from the flow cytometry were normalized to total cell counts; the resulting  
566 relative cell proportions were then transformed using centered log-ratios<sup>64</sup>. Normalized  
567 transcriptomic, proteomic, cytokine, and metabolomic data were log-transformed. We  
568 further decomposed the within-subject from the between-subject variance in the data  
569 sets to account for repeated measures<sup>65</sup>. This is analogous to normalizing all samples  
570 to their DOL0 time point. Finally, a broad, unsupervised, variance-based filter was  
571 applied to the transcriptomic data, retaining the 50% most relevant features. To integrate  
572 across data types, we applied sparse generalized canonical correlation discriminant  
573 analysis via the DIABLO framework, part of the mixOmics package<sup>18, 19, 66</sup>. DIABLO  
574 constructs components across any number of input matrices, maximizing their  
575 covariance with each other and a given response variable (in this case, DOL), while  
576 simultaneously performing feature selection by applying an L1 penalty (LASSO) on the  
577 model loadings<sup>23</sup>. Importantly, DIABLO identifies key drivers associated with the  
578 response variable of interest across all input data matrices jointly. Cross-validation (20 x  
579 5-fold) was used to determine the optimal model hyperparameters (number of  
580 components, number of features per component), as well as to provide an estimate of  
581 the ability of the model to generalize to new data. Selected model features, i.e. key  
582 correlates of change across DOL, were subjected to pathway over-representation  
583 analysis against the Reactome pathways database (obtained via MSigDB<sup>67</sup>) and BTM<sup>21</sup>  
584 annotated gene set libraries, using the hypergeometric test. Obtained p-values were  
585 adjusted for multiple comparisons using the Benjamini-Hochberg procedure<sup>68</sup>, either  
586 separately per component, or after enrichment to include their first degree neighbours  
587 sub-graph in the protein:protein interaction data<sup>17</sup>.  
588

589 *Multi-scale, Multifactorial response Network*. Examples of systematic integration across  
590 many large OMICS datasets remain rare. Recently, Li and others<sup>7</sup> used multi-scale,  
591 multifactorial response networks to study the immune response to varicella zoster  
592 vaccine across four OMICS datasets. We wanted to compare insights generated by our  
593 chosen approaches to theirs and, to that end, created a multi-scale interaction network  
594 as they described, with the following important modifications. First, instead summarizing  
595 the higher dimensional data blocks (transcriptomics, metabolomics) to either gene  
596 modules or metabolic pathways (Blood Transcriptomic Modules (BTM])<sup>69</sup>; Metabolon Inc.  
597 annotation, respectively) using a simple average, we used eigen-gene summarization<sup>70</sup>,  
598 a weighted average based on the 1st principal component of the data, in order to  
599 maximize the variation explained. Second, clusters were identified at each DOL using  
600 Euclidean distance and Ward's method<sup>71</sup>, as recommended by Li (personal  
601 communication). The optimal number of clusters was determined using the elbow  
602 criterion<sup>72</sup>. Stable clusters were identified by comparing cluster membership using the  
603 Szymkiewicz-Simpson coefficient<sup>73</sup>. Finally, cluster association with DOL was assessed  
604 using the Correlation Adjusted MEan RAnk (Camera) gene set test<sup>74</sup>.  
605

606 *Meta-Integration*. Minimum connected networks for both components 1 and 2 of the  
607 DIABLO model were constructed using NetworkAnalyst, using pairwise lists of each data  
608 type for both components to construct a minimum connected network.  
609

610 *Identifying convergence across data integration strategies.* The various data integration  
611 approaches resulted in outputs that were not directly comparable. We simplified their  
612 outputs to lists of features of interests to enable direct comparisons. NetworkAnalyst: list  
613 of features that made up the minimal connected network. DIABLO: features selected by  
614 the model with hyperparameterization that minimized cross-validation error rate. MMRN:  
615 features that made up the inter-connected node-networks that were stably identified  
616 across DOLs and significantly associated with DOL. We then compared the identified  
617 feature sets and carried out gene set enrichment analysis against the Reactome  
618 collection.

619 FIGURE LEGENDS.  
620

621 **Figure 1. Blood Processing Overview.** Newborn peripheral venous blood was drawn  
622 directly into heparinized collection tubes. Aliquots (200  $\mu$ L) were removed for  
623 transcriptomic analysis. Plasma was then harvested after a spin, and cryopreserved for  
624 cytokine, proteomic and metabolomic analysis. The remaining cellular fraction was  
625 diluted to replace the plasma removed, and 100  $\mu$ L aliquots from this mixture were  
626 processed for single cell immunophenotyping by flow cytometry. With a starting volume  
627 of 1 ml, this SOP still left the cellular fraction contained in 400  $\mu$ L of blood that could be  
628 used for other analysis.  
629

630 **Figure 2. Transcriptomic analyses identified dramatic increases in the number of**  
631 **differentially expressed genes during the first week of life.** In **A**, up- and down-  
632 regulated differentially expressed genes are plotted by DOL (vs DOL0) and numbers of  
633 genes are listed above each point except for down-regulated genes at DOL1 vs DOL0,  
634 where the number is zero. **B** shows zero order interaction networks for genes that  
635 differentially expressed at DOL3 vs DOL0 and DOL7 vs DOL0. Within networks, up-  
636 regulated nodes are displayed in red and down-regulated nodes in green. **C** shows  
637 selected pathways enriched among up- and down-regulated genes along with  
638 Bonferroni-corrected  $p$ -values for DOL3 and DOL7.  
639

640 **Figure 3. Integration of Multiple Data Types via NetworkAnalyst provided novel**  
641 **insight.** In **A**, the number of *novel nodes* derived from molecular networks integration  
642 from each combination of data types on each DOL vs DOL0 are shown. **B** shows  
643 minimum connected networks containing all three individual data types, where nodes  
644 derived from the transcriptome are shown in blue, nodes from the metabolome in red,  
645 and nodes from the proteome in green. Novel nodes are shown in orange. **C** shows  
646 enrichment of *novel pathways* on DOL1, 3 and 7 vs DOL0 from each data type and their  
647 respective combinatorial analysis. **D** lists selected novel pathways along with their  
648 Bonferroni corrected  $p$ -values and the data type combinations from which they were  
649 derived.  
650

651 **Figure 4: DIABLO uncovered biologically relevant features by integrating**  
652 **information across data types.** Schematic representation of two contrasting  
653 integration approaches using multi-variate techniques: **A** shows that DIABLO selects  
654 features jointly across data types, resulting in the identification of features with strong  
655 associations across data types. Conversely, as shown in **B**, ensembles of multi-variate  
656 models, constructed independently of each other, result in a selection of features that  
657 are poorly associated across data types. This is visualized in correlation heatmaps of  
658 the selected features (middle) and corresponding networks (right), with dense sub-  
659 graphs, or network modules, encircled. In particular, the network modules identified in **A**  
660 (**right**) include a number of features selected from all data types. This is not the case in  
661 **B (right).** The minimal set of features selected by DIABLO across data types as shown  
662 in **C** could discriminate between DOL and distinct sets of these features separated birth  
663 from all other DOLs (DIABLO component 1) and DOL 1, 3, and 7 (DIABLO component  
664 2). Features identified by DIABLO (blue bars) were largely distinct from those identified

665 by more traditional single-OMICS multi-variate approaches (red bars; overlaps in grey);  
666 shown in **D** using an UpSet plot. Moreover, features identified by DIABLO were more  
667 strongly enriched for known biological (functional) pathways; shown in **E** using an UpSet  
668 plot (blue vs. red bars). Horizontal bars are mapped to the number of elements in each  
669 set of features being compared. Vertical bars correspond to the number of elements in  
670 the intersections when carrying out various set comparisons.

671  
**Figure 5. Multi-scale, multifactorial response networks identified a small set of**  
672 **biological processes consistently associated with each other across all DOLs.**  
673 Networks were constructed separately at each DOL and stable clusters identified. The  
674 partial least squares regression scores, a measure of the strength of association  
675 between data types, were highest at DOL 1 and decreased across the first week of life  
676 are shown in **A**; the transcriptomic and metabolomic data were most strongly associated  
677 as indicated in **B**, followed by the proteomics and cytokine data, metabolomic and  
678 proteomic, and flow cytometry and transcriptomic. Association of stable and transient  
679 clusters with DOL was assessed using the Correlation-Adjusted Mean RAnk (CAMERA)  
680 method. **C** shows the comparison of the distribution of the observed Benjamini-  
681 Hochberg false discovery rates to expected quantiles is shown in a Quantile-quantile  
682 (QQ) plot. Stable clusters were more strongly associated with DOL compared to  
683 transient ones. Finally, a small set of stable clusters were consistently associated at all 3  
684 DOLs, forming a stable sub-network; the features that make up this stable sub-network  
685 are visualized in a heatmap in **D**.

686  
**Figure 6. Meta-integration confirmed findings across analytical methods.**  
687 **A** is an UpSet plot that depicts the concordance in features selected by the three  
688 different integration methods. Features identified by at least two methods are shown in  
689 blue, while features identified by all three methods are shown in red. The concordance  
690 between significantly enriched pathways (Benjamini-Hochberg FDR  $\leq 0.05$ ) identified by  
691 the three methods are shown in the upset plot in **B**, with gene sets identified by at least  
692 2 of the 3 (blue) or all 3 methods (red) highlighted. Minimum connected networks  
693 containing luminex (purple nodes), proteome (green nodes) and transcriptome (blue  
694 nodes) features from DIABLO components 1 and 2 are shown in **C**.  
695  
696

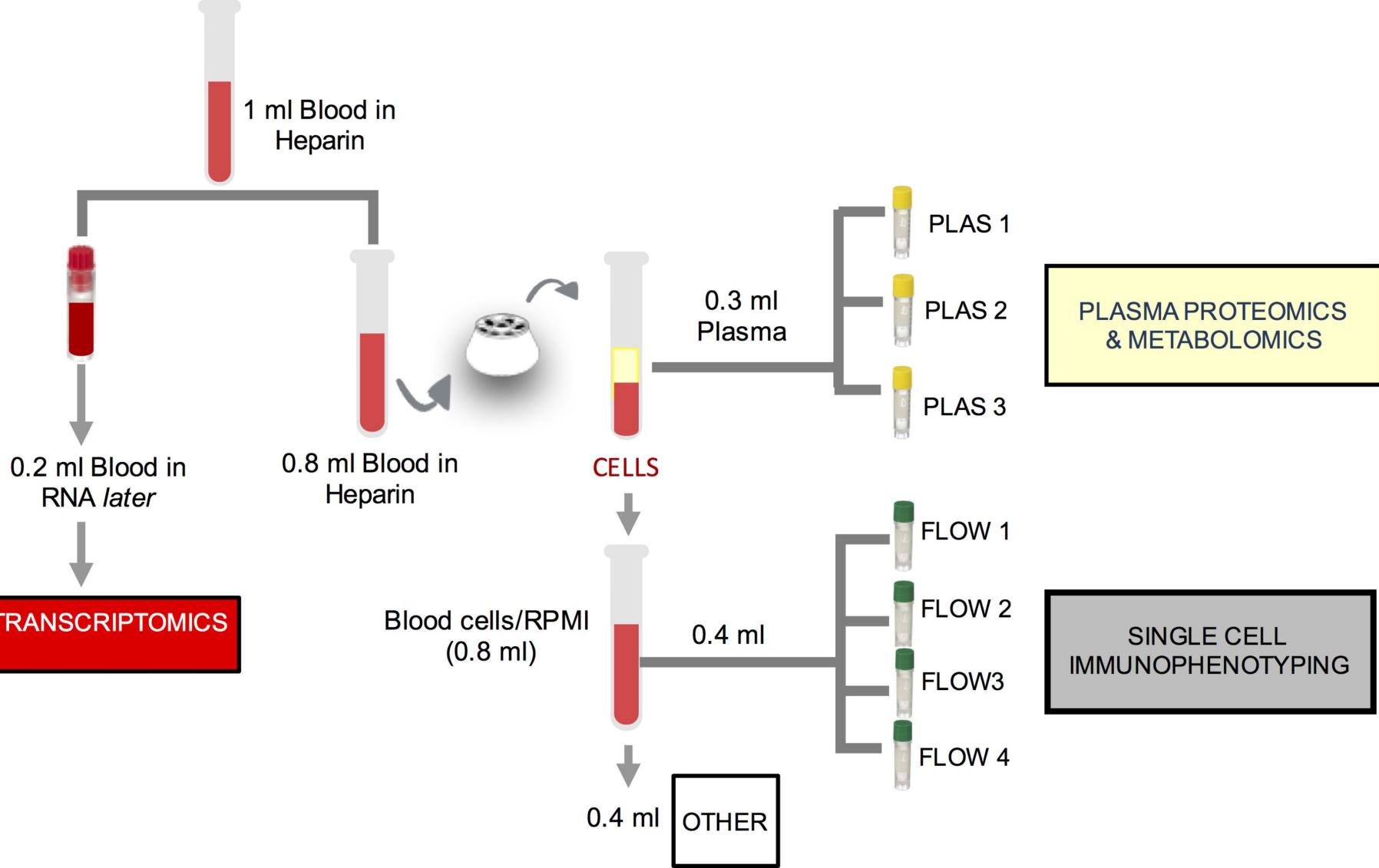
698 **REFERENCES**

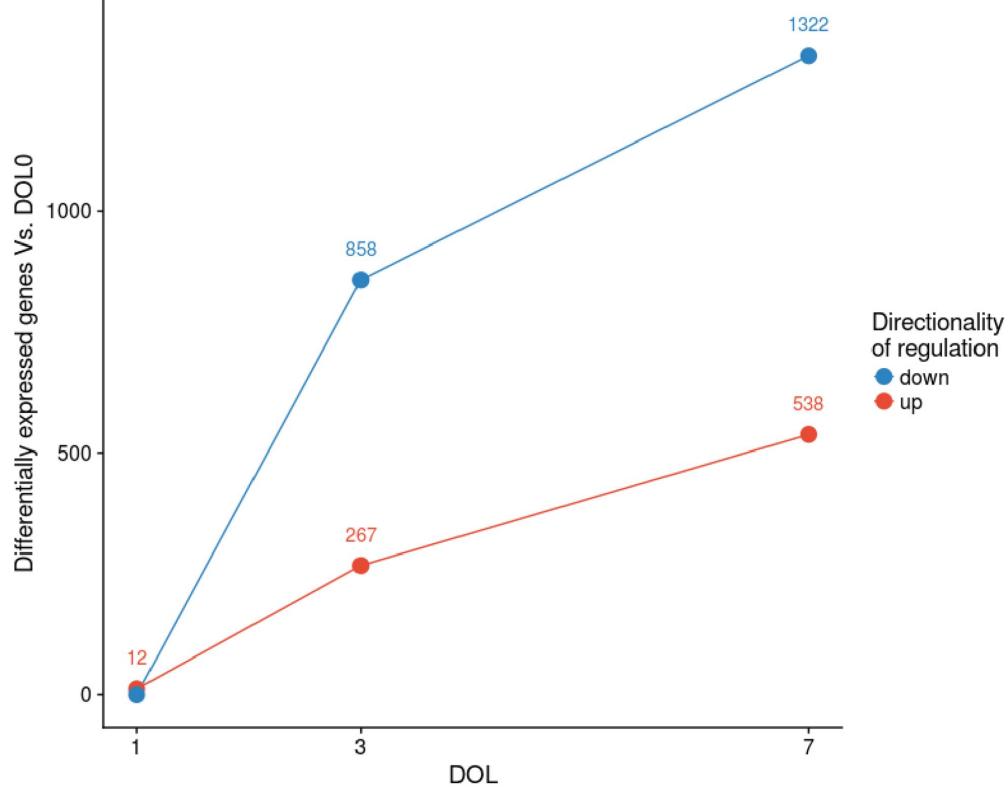
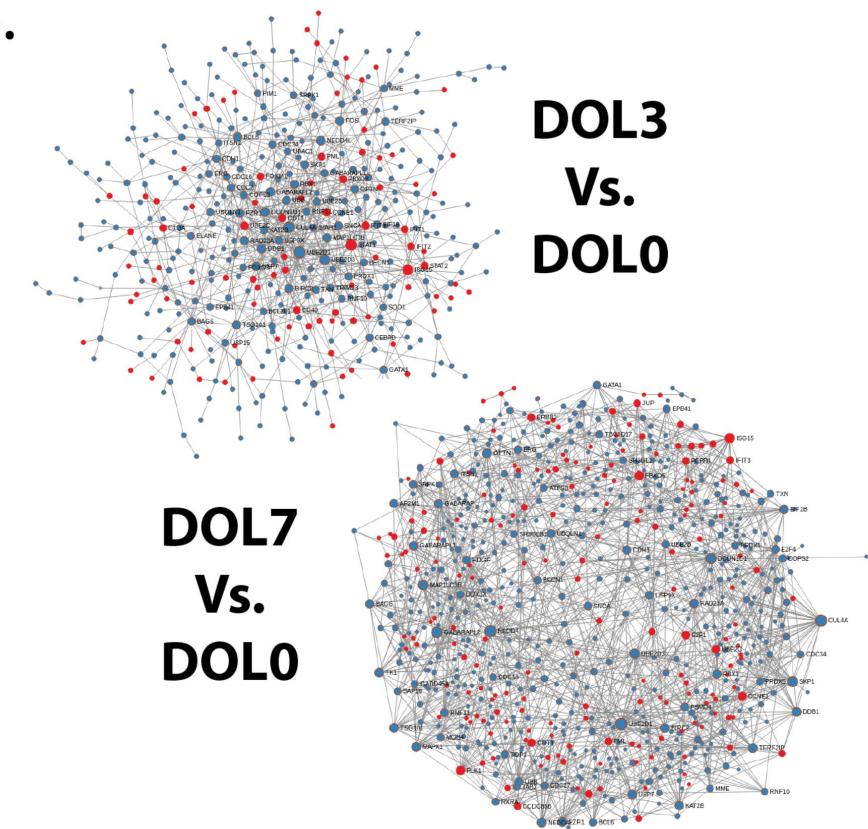
- 699 1. Kollmann, T.R., Kampmann, B., Mazmanian, S.K., Marchant, A. & Levy, O. Protecting  
700 the Newborn and Young Infant from Infectious Diseases: Lessons from Immune  
701 Ontogeny. *Immunity* **46**, 350-363 (2017).
- 702 2. Zhang, X., Zhivaki, D. & Lo-Man, R. Unique aspects of the perinatal immune system.  
703 *Nature reviews* (2017).
- 704 3. Balbus, J.M. et al. Early-life prevention of non-communicable diseases. *Lancet* **381**, 3-  
705 4 (2013).
- 706 4. Amenyogbe, N., Kollmann, T.R. & Ben-Othman, R. Early-Life Host-Microbiome  
707 Interphase: The Key Frontier for Immune Development. *Frontiers in pediatrics* **5**, 111  
708 (2017).
- 709 5. Chaussabel, D. & Pulendran, B. A vision and a prescription for big data-enabled  
710 medicine. *Nat Immunol* **16**, 435-439 (2015).
- 711 6. Amenyogbe, N., Levy, O. & Kollmann, T.R. Systems vaccinology: a promise for the  
712 young and the poor. *Philos Trans R Soc Lond B Biol Sci* **370** (2015).
- 713 7. Li, S. et al. Metabolic Phenotypes of Response to Vaccination in Humans. *Cell* **169**,  
714 862-877.e817 (2017).
- 715 8. Howie, S.R. Blood sample volumes in child health research: review of safe limits.  
716 *Bulletin of the World Health Organization* **89**, 46-53 (2011).
- 717 9. Tsang, J.S. Utilizing population variation, vaccination, and systems biology to study  
718 human immunology. *Trends Immunol.* **36**, 479-493. doi:  
719 410.1016/j.it.2015.1006.1005. Epub 2015 Jul 1014. (2015).
- 720 10. Carr, E.J. et al. The cellular composition of the human immune system is shaped by  
721 age and cohabitation. *Nat Immunol* **17**, 461-468 (2016).
- 722 11. Smolen, K.K. et al. Single-Cell Analysis of Innate Cytokine Responses to Pattern  
723 Recognition Receptor Stimulation in Children across Four Continents. *J Immunol*  
724 (2014).
- 725 12. Smolen, K.K. et al. Pattern recognition receptor-mediated cytokine response in  
726 infants across 4 continents. *The Journal of allergy and clinical immunology* **133**, 818-  
727 826 e814 (2014).
- 728 13. Tsang, J.S. et al. Global analyses of human immune variation reveal baseline  
729 predictors of postvaccination responses. *Cell* **157**, 499-513 (2014).
- 730 14. Shannon, C.P. et al. Two-stage, in silico deconvolution of the lymphocyte  
731 compartment of the peripheral whole blood transcriptome in the context of acute  
732 kidney allograft rejection. *PloS one* **9**, e95224 (2014).
- 733 15. Smith, C.L. et al. Identification of a human neonatal immune-metabolic network  
734 associated with bacterial infection. *Nature communications* **5**, 4649 (2014).
- 735 16. Henry, E. & Christensen, R.D. Reference Intervals in Neonatal Hematology. *Clinics in  
736 perinatology* **42**, 483-497 (2015).
- 737 17. Xia, J., Gill, E.E. & Hancock, R.E. NetworkAnalyst for statistical, visual and network-  
738 based meta-analysis of gene expression data. *Nature protocols* **10**, 823-844 (2015).
- 739 18. Le Cao, K.A., Boitard, S. & Besse, P. Sparse PLS discriminant analysis: biologically  
740 relevant feature selection and graphical displays for multiclass problems. *BMC  
741 bioinformatics* **12**, 253 (2011).
- 742 19. Singh, A. et al. DIABLO - an integrative, multi-omics, multivariate method for multi-  
743 group classification. *bioRxiv* (2016).

- 744 20. Rohart, F., Gautier, B., Singh, A. & Le Cao, K.-A. mixOmics: an R package for 'omics  
745 feature selection and multiple data integration. *bioRxiv* (2017).
- 746 21. Chaussabel, D. & Baldwin, N. Democratizing systems immunology with modular  
747 transcriptional repertoire analyses. *Nature reviews* **14**, 271-280 (2014).
- 748 22. Breuer, K. et al. InnateDB: systems biology of innate immunity and beyond--recent  
749 updates and continuing curation. *Nucleic acids research* **41**, D1228-1233 (2013).
- 750 23. Tibshirani, R. Regression Shrinkage and Selection via the Lasso. (1996).
- 751 24. Fan, J. & Lv, J. A Selective Overview of Variable Selection in High Dimensional Feature  
752 Space. *Statistica Sinica* **20**, 101-148 (2010).
- 753 25. Pettengill, M.A., van Haren, S.D. & Levy, O. Soluble mediators regulating immunity in  
754 early life. *Frontiers in immunology* **5**, 457 (2014).
- 755 26. Ismail, A.A., Walker, P.L., Macfaul, R. & Gindal, B. Diagnostic value of serum  
756 testosterone measurement in infancy: two case reports. *Annals of clinical  
757 biochemistry* **26 ( Pt 3)**, 259-261 (1989).
- 758 27. Radtke, F., Wilson, A., Mancini, S.J. & MacDonald, H.R. Notch regulation of lymphocyte  
759 development and function. *Nat Immunol* **5**, 247-253 (2004).
- 760 28. Abdelhaleem, M. RNA helicases: regulators of differentiation. *Clinical biochemistry*  
761 **38**, 499-503 (2005).
- 762 29. Loo, Y.M. & Gale, M., Jr. Immune signaling by RIG-I-like receptors. *Immunity* **34**, 680-  
763 692 (2011).
- 764 30. Schmidt, C.Q., Lambris, J.D. & Ricklin, D. Protection of host cells by complement  
765 regulators. *Immunological reviews* **274**, 152-171 (2016).
- 766 31. Romero-Moya, D. et al. Cord blood-derived CD34+ hematopoietic cells with low  
767 mitochondrial mass are enriched in hematopoietic repopulating stem cell function.  
768 *Haematologica* **98**, 1022-1029 (2013).
- 769 32. Lugo, B., Ford, H.R. & Grishin, A. Molecular signaling in necrotizing enterocolitis:  
770 regulation of intestinal COX-2 expression. *Journal of pediatric surgery* **42**, 1165-1171  
771 (2007).
- 772 33. Reinebrant, H.E. et al. Cyclo-oxygenase (COX) inhibitors for treating preterm labour.  
773 *Cochrane Database Syst Rev*, Cd001992 (2015).
- 774 34. Westerhuis, J.A., van Velzen, E.J., Hoefsloot, H.C. & Smilde, A.K. Multivariate paired  
775 data analysis: multilevel PLSDA versus OPLSDA. *Metabolomics : Official journal of the  
776 Metabolomic Society* **6**, 119-128 (2010).
- 777 35. Lee, J.A. et al. MIFlowCyt: the minimum information about a Flow Cytometry  
778 Experiment. *Cytometry* **73**, 926-930 (2008).
- 779 36. Hahne, F. et al. flowCore: a Bioconductor package for high throughput flow  
780 cytometry. *BMC Bioinformatics*. **10:106.**, 10.1186/1471-2105-1110-1106. (2009).
- 781 37. Spidlen, J. et al. Data File Standard for Flow Cytometry, version FCS 3.1. *Cytometry A*  
782 **77**, 97-100 (2010).
- 783 38. Malek, M. et al. flowDensity: reproducing manual gating of flow cytometry data by  
784 automated density-based cell population identification. *Bioinformatics*. **31**, 606-607.  
785 doi: 610.1093/bioinformatics/btu1677. Epub 2014 Oct 1016. (2015).
- 786 39. O'Neill, K., Jalali, A., Aghaeepour, N., Hoos, H. & Brinkman, R.R. Enhanced  
787 flowType/RchyOptimyx: a BioConductor pipeline for discovery in high-dimensional  
788 cytometry data. *Bioinformatics* **30**, 1329-1330 (2014).

- 789 40. Ewels, P., Magnusson, M., Lundin, S. & Kaller, M. MultiQC: summarize analysis results  
790 for multiple tools and samples in a single report. *Bioinformatics* **32**, 3047-3048  
791 (2016).
- 792 41. Dobin, A. et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21  
793 (2013).
- 794 42. Anders, S., Pyl, P.T. & Huber, W. HTSeq--a Python framework to work with high-  
795 throughput sequencing data. *Bioinformatics* **31**, 166-169 (2015).
- 796 43. Love, M.I., Huber, W. & Anders, S. Moderated estimation of fold change and  
797 dispersion for RNA-seq data with DESeq2. *Genome biology* **15**, 550 (2014).
- 798 44. Foroushani, A.B., Brinkman, F.S. & Lynn, D.J. Pathway-GPS and SIGORA: identifying  
799 relevant pathways based on the over-representation of their gene-pair signatures.  
800 *PeerJ* **1**, e229 (2013).
- 801 45. Berger, S.T. et al. MStern Blotting-High Throughput Polyvinylidene Fluoride (PVDF)  
802 Membrane-Based Proteomic Sample Preparation for 96-Well Plates. *Molecular &*  
803 *cellular proteomics : MCP* **14**, 2814-2823 (2015).
- 804 46. Bennike, T.B. & Steen, H. High-Throughput Parallel Proteomic Sample Preparation  
805 Using 96-Well Polyvinylidene Fluoride (PVDF) Membranes and C18 Purification  
806 Plates. *Methods Mol Biol* **1619**, 395-402 (2017).
- 807 47. Cox, J. et al. Accurate proteome-wide label-free quantification by delayed  
808 normalization and maximal peptide ratio extraction, termed MaxLFQ. *Molecular &*  
809 *cellular proteomics : MCP* **13**, 2513-2526 (2014).
- 810 48. Keshishian, H. et al. Multiplexed, Quantitative Workflow for Sensitive Biomarker  
811 Discovery in Plasma Yields Novel Candidates for Early Myocardial Injury. *Molecular*  
812 & *cellular proteomics : MCP* **14**, 2375-2393 (2015).
- 813 49. Bennike, T. et al. A normative study of the synovial fluid proteome from healthy  
814 porcine knee joints. *Journal of proteome research* **13**, 4377-4387 (2014).
- 815 50. Reiner, A., Yekutieli, D. & Benjamini, Y. Identifying differentially expressed genes  
816 using false discovery rate controlling procedures. *Bioinformatics* **19**, 368-375 (2003).
- 817 51. Bennike, T.B. et al. Proteome Analysis of Rheumatoid Arthritis Gut Mucosa. *Journal of*  
818 *proteome research* **16**, 346-354 (2017).
- 819 52. Deeb, S.J., D'Souza, R.C., Cox, J., Schmidt-Suprian, M. & Mann, M. Super-SILAC allows  
820 classification of diffuse large B-cell lymphoma subtypes by their protein expression  
821 profiles. *Molecular & cellular proteomics : MCP* **11**, 77-89 (2012).
- 822 53. Johnson, W.E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray expression  
823 data using empirical Bayes methods. *Biostatistics (Oxford, England)* **8**, 118-127  
824 (2007).
- 825 54. Camargo, A., Azuaje, F., Wang, H. & Zheng, H. Permutation - based statistical tests for  
826 multiple hypotheses. *Source code for biology and medicine* **3**, 15 (2008).
- 827 55. Croft, D. et al. The Reactome pathway knowledgebase. *Nucleic acids research* **42**,  
828 D472-477 (2014).
- 829 56. Shannon, P. et al. Cytoscape: a software environment for integrated models of  
830 biomolecular interaction networks. *Genome research* **13**, 2498-2504 (2003).
- 831 57. Evans, A.M., DeHaven, C.D., Barrett, T., Mitchell, M. & Milgram, E. Integrated,  
832 nontargeted ultrahigh performance liquid chromatography/electrospray ionization  
833 tandem mass spectrometry platform for the identification and relative quantification

- 834 of the small-molecule complement of biological systems. *Analytical chemistry* **81**,  
835 6656-6667 (2009).
- 836 58. Xia, J., Sinelnikov, I.V., Han, B. & Wishart, D.S. MetaboAnalyst 3.0--making  
837 metabolomics more meaningful. *Nucleic acids research* **43**, W251-257 (2015).
- 838 59. Storey, J.D. & Tibshirani, R. Statistical significance for genomewide studies. *Proc Natl  
839 Acad Sci USA* **100**, 9440-9445 (2003).
- 840 60. Ogata, H., Goto, S., Fujibuchi, W. & Kanehisa, M. Computation with the KEGG pathway  
841 database. *Bio Systems* **47**, 119-128 (1998).
- 842 61. Wishart, D.S. et al. HMDB 3.0--The Human Metabolome Database in 2013. *Nucleic  
843 acids research* **41**, D801-807 (2013).
- 844 62. Smoot, M.E., Ono, K., Ruscheinski, J., Wang, P.L. & Ideker, T. Cytoscape 2.8: new  
845 features for data integration and network visualization. *Bioinformatics* **27**, 431-432  
846 (2011).
- 847 63. Caspi, R. et al. The MetaCyc database of metabolic pathways and enzymes and the  
848 BioCyc collection of Pathway/Genome Databases. *Nucleic acids research* **42**, D459-  
849 471 (2014).
- 850 64. Aitchison, J. The Statistical Analysis of Compositional Data. (1982).
- 851 65. Liquet, B., Le Cao, K.A., Hocini, H. & Thiebaut, R. A novel approach for biomarker  
852 selection and the integration of repeated measures experiments from two assays.  
*BMC bioinformatics* **13**, 325 (2012).
- 853 66. Singh, A. et al. Identifying Molecular Mechanisms of the Late-Phase Asthmatic  
854 Response by Integrating Cellular, Gene, and Metabolite Levels in Blood. *Annals of the  
855 American Thoracic Society* **13 Suppl 1**, S98 (2016).
- 856 67. Subramanian, A. et al. Gene set enrichment analysis: a knowledge-based approach  
857 for interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A* **102**,  
858 15545-15550 (2005).
- 859 68. Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and  
860 Powerful Approach to Multiple Testing. (1995).
- 861 69. Li, S., Rouphael, N., Duraisingham, S., Romero-Steiner, S. & Presnell, S. Molecular  
862 signatures of antibody responses derived from a systems biology study of five  
863 human vaccines. **15**, 195-204 (2014).
- 864 70. Langfelder, P. & Horvath, S. Eigengene networks for studying the relationships  
865 between co-expression modules. *BMC systems biology* **1**, 54 (2007).
- 866 71. Ward, J.H. Hierarchical Grouping to Optimize an Objective Function. *J. Am. Stat. Assoc.*  
867 **58**, 236-244 (1963).
- 868 72. Thorndike, R.L. Who belongs in the family? *Psychometrika* **18**, 267-276 (1953).
- 869 73. Shannon, C.P. et al. SABRE: a method for assessing the stability of gene modules in  
870 complex tissues and subject populations. *BMC bioinformatics* **17**, 460 (2016).
- 871 74. Wu, D. & Smyth, G.K. Camera: a competitive gene set test accounting for inter-gene  
872 correlation. *Nucleic acids research* **40**, e133 (2012).
- 873
- 874



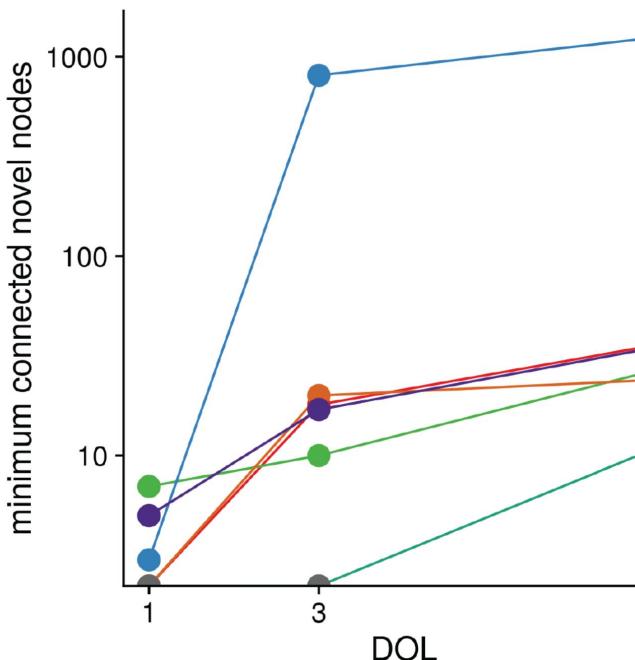
**A.****B.****C.****DOL3 Vs. DOLO Enriched Pathways**

Up-regulated		Down-regulated	
Pathway Name	Bonferroni	Pathway Name	Bonferroni
Interferon alpha/beta signaling	1.46E-86	Heme biosynthesis	6.53E-35
Negative regulators of RIG-I/MDA5 signaling	2.35E-31	Cellular responses to stress	8.31E-27
Interferon gamma signaling	8.18E-23	Macroautophagy	2.36E-22
ISG15 antiviral mechanism	1.19E-03	Iron uptake and transport	5.70E-11
		Downregulation of SMAD2/3:SMAD4 transcriptional activity	5.76E-07

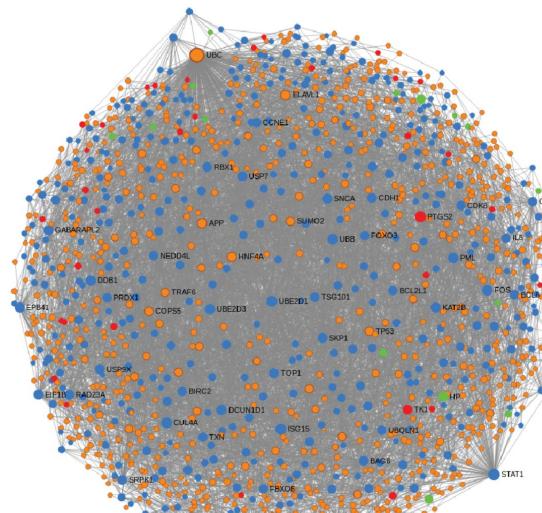
**DOL7 Vs. DOLO Enriched Pathways**

Up-regulated		Down-regulated	
Pathway Name	Bonferroni	Pathway Name	Bonferroni
Interferon alpha/beta signaling	1.91E-72	Detoxification of Reactive Oxygen Species	4.01E-31
Classical antibody-mediated complement activation	3.76E-19	Golgi Associated Vesicle Biogenesis	4.83E-23
Signaling by FGFR1	7.13E-04	MyD88-independent TLR3/TLR4 cascade	1.86E-09
Chemokine receptors bind chemokines	1.74E-03	Antigen processing: Ubiquitination & Proteasome degradation	6.51E-09
SHC1 events in EGFR signaling	1.47E-02	Interleukin-6 signaling	2.19E-05

A.

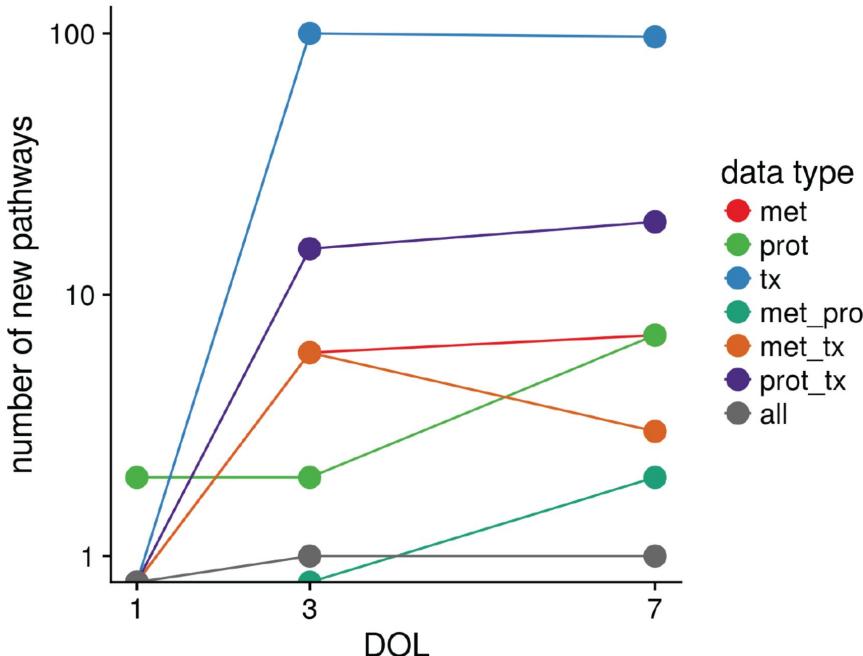


B.



**DOL3  
Vs.  
DOL0**

C.

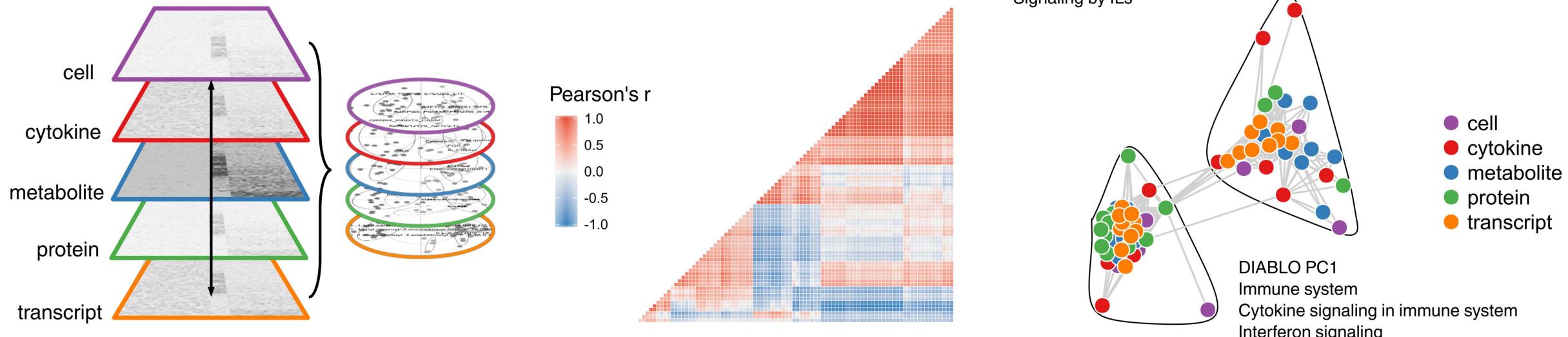
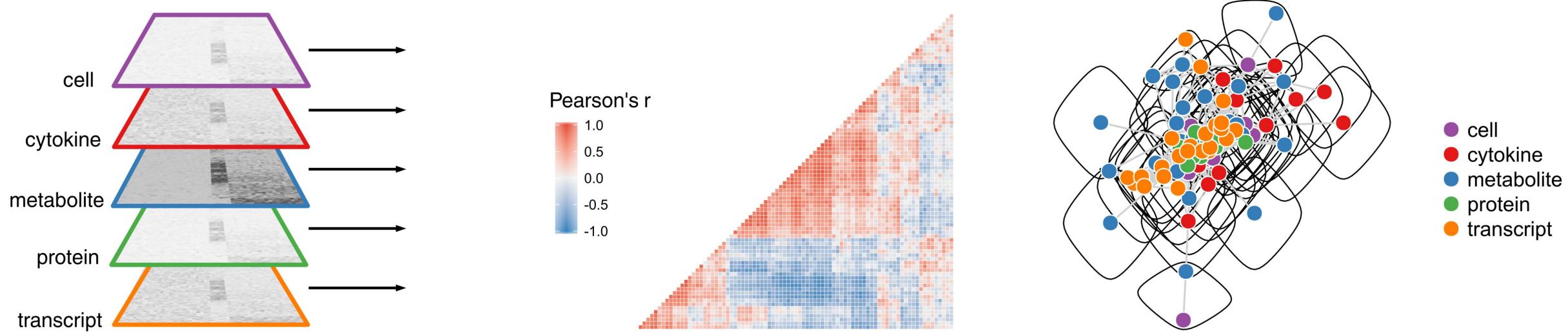
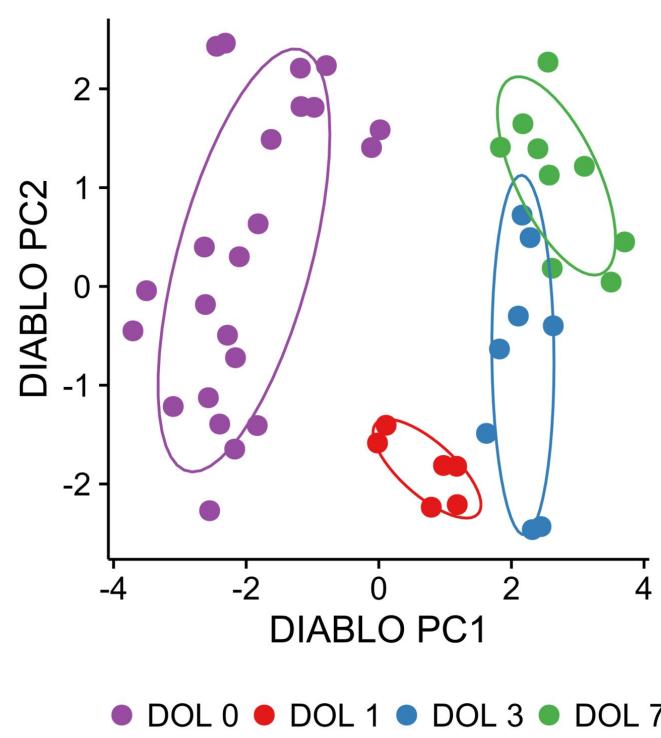
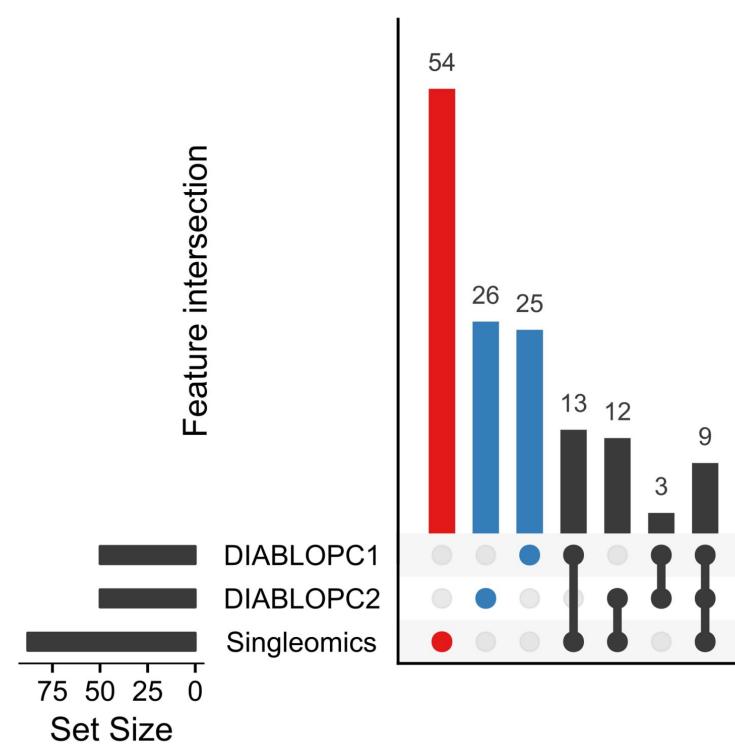
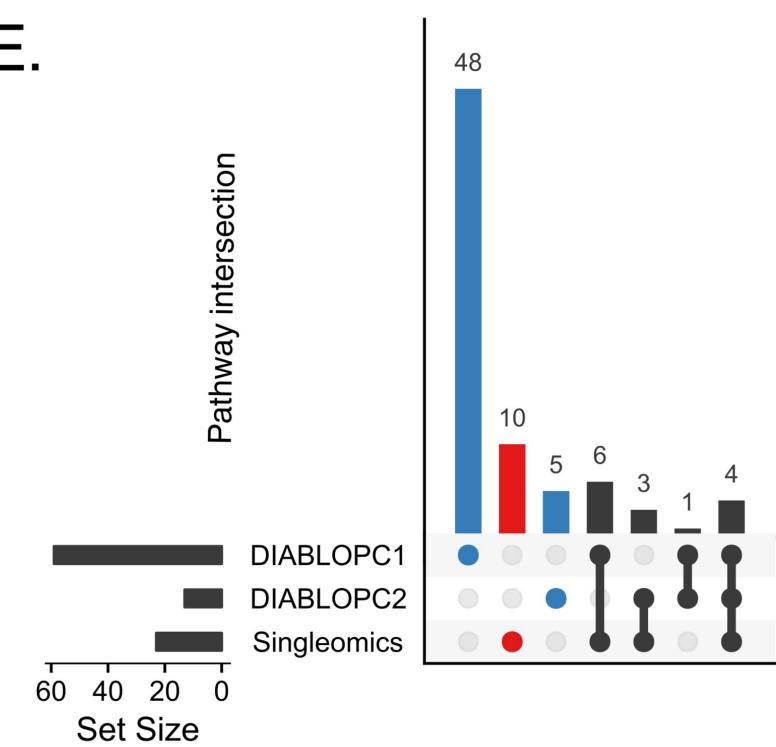


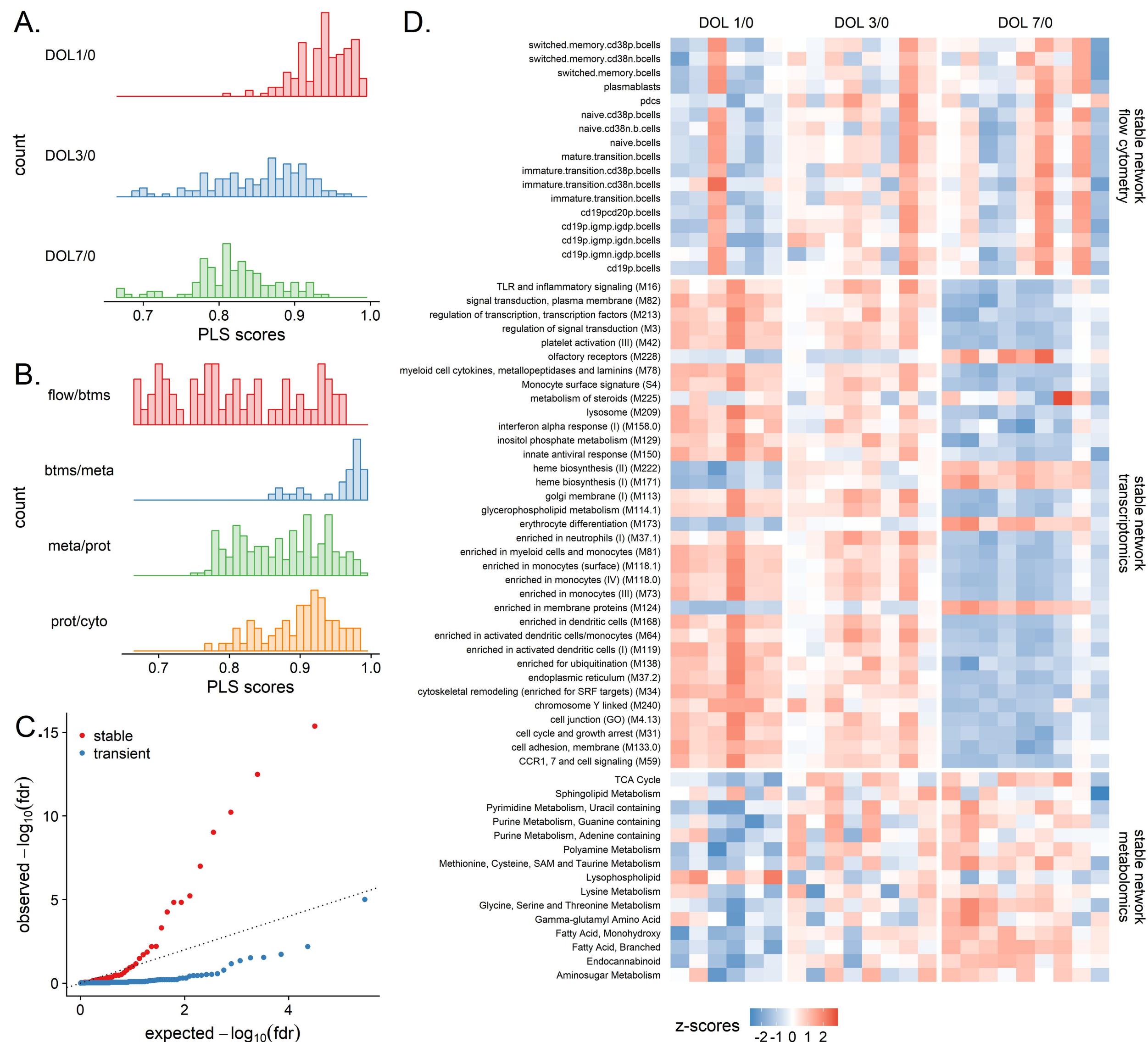
D.

Novel Pathways from DOL3 Vs. DOL0 Data Integration		
Pathway Name	Bonferroni	Data Types
Transcriptional regulation of white adipocyte differentiation	2.03E-13	Prot + Tx
Metabolism of porphyrins	1.14E-07	Met + Tx
RHO GTPases activate PKNs	1.82E-04	Prot + Tx
Chromatin organization	1.96E-03	Prot + Tx
Sialic acid metabolism	1.62E-02	Met + Tx
Cytosolic sensors of pathogen-associated DNA	3.71E-02	All

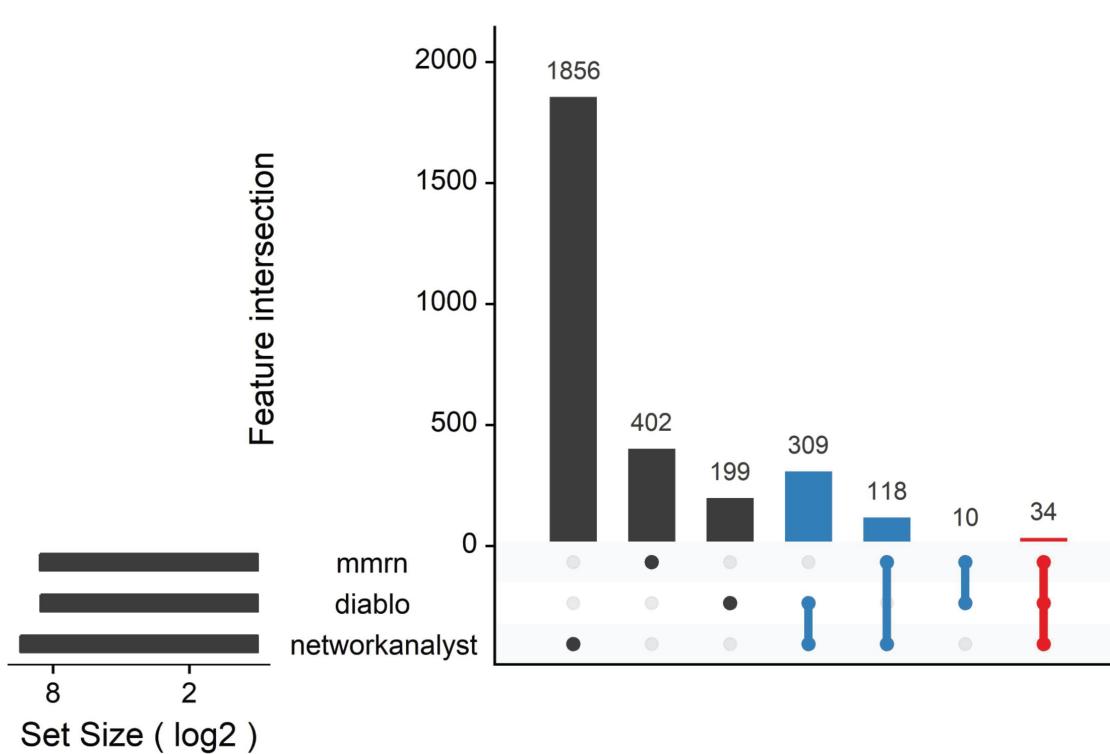
  

Novel Pathways from DOL7 Vs. DOL0 Data Integration		
Pathway Name	Bonferroni	Data Types
Endosomal/Vacuolar pathway	3.65E-11	Met + Prot
Clathrin derived vesicle budding	1.04E-06	Prot + Tx
Signaling by EGFR	1.47E-06	Prot + Tx
GLI3 is processed to GLI3R by the proteasome	3.50E-02	Prot + Tx
Death Receptor Signalling	4.16E-02	Prot + Tx
Platelet activation, signaling and aggregation	4.60E-02	Prot + Tx

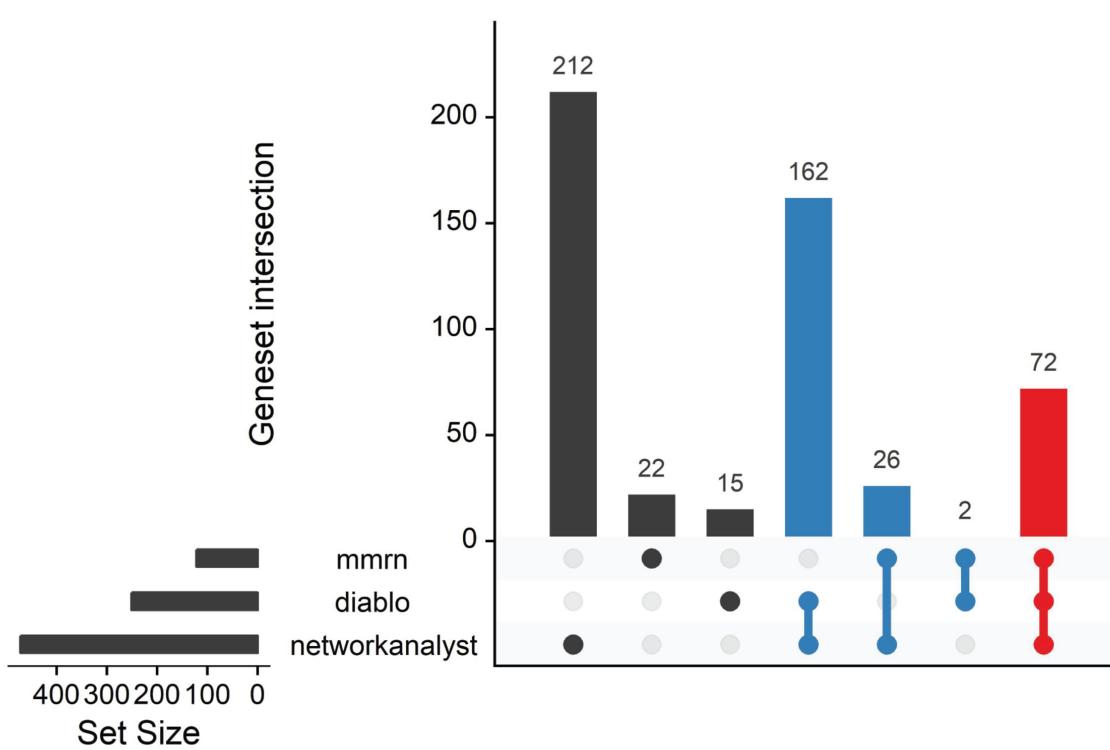
**A.****B.****C.****D.****E.**



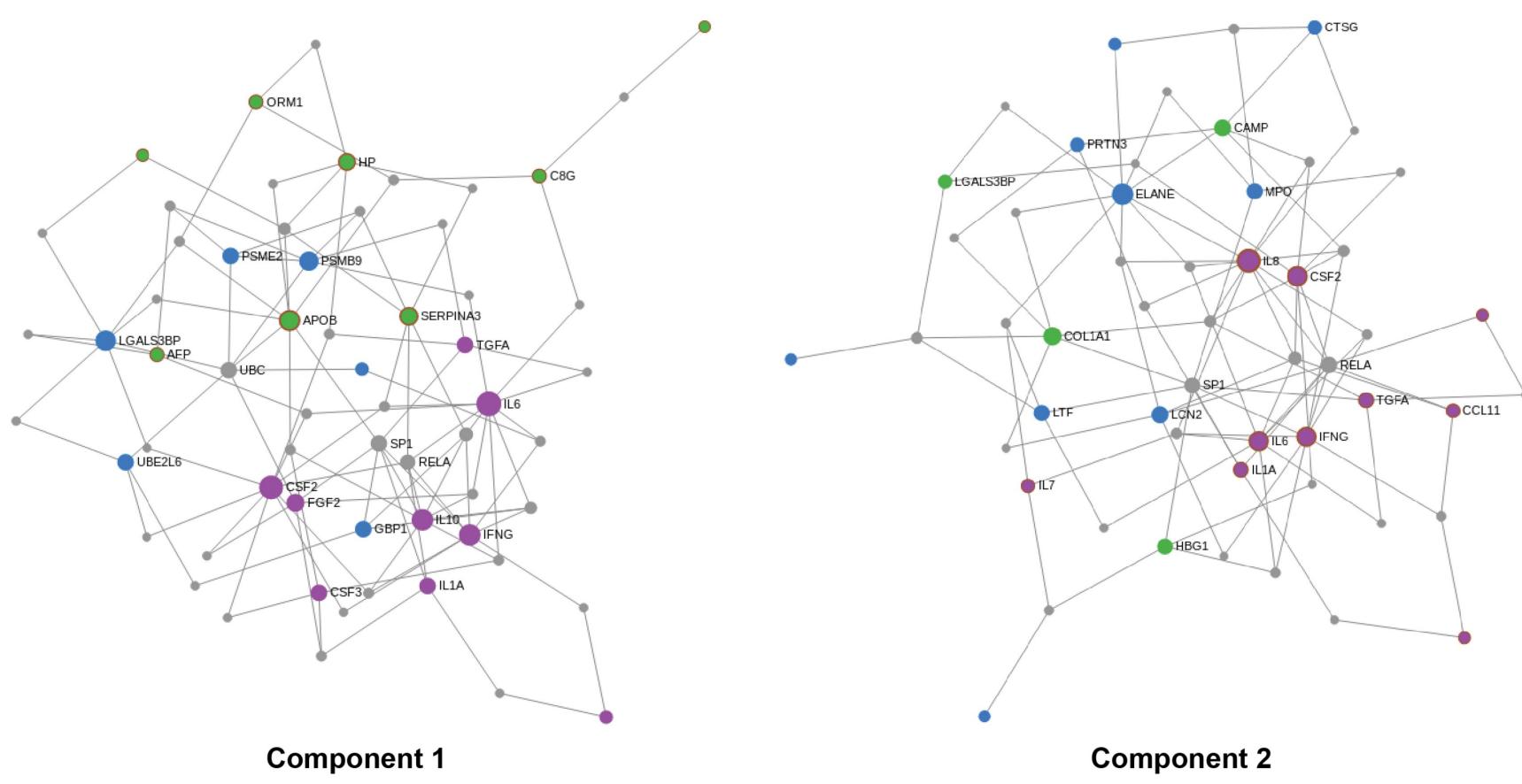
A.



B.



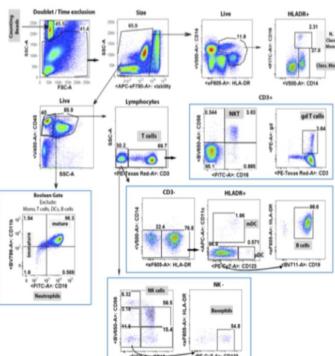
C.



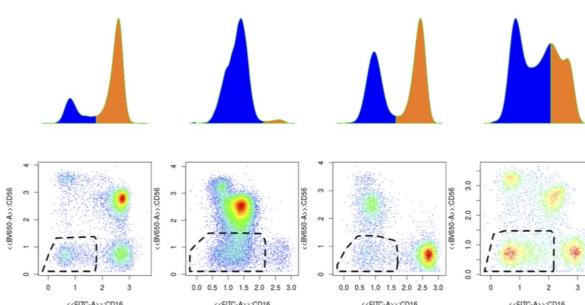
A

Cell population	Anchor Markers
T cells	CD45+/CD3+/CD56-/CD16-
y <sup>b</sup> -T cells	CD45+/CD3+/y <sup>b</sup> +
B cells	CD45+/CD3-/HLADR+/CD11c-/CD123-/CD19+
Classical Monocytes	CD45+/HLADR+/CD14+/CD16-
Non-classical Monocytes	CD45+/HLADR+/CD14+/CD16+
NKT cells	CD45+/CD3+/CD56+/CD16+
Myeloid Dendritic Cells (mDC)	CD45+/CD3-/HLADR+/CD11c+
Plasmacytoid Dendritic Cells (pDC)	CD45+/CD3-/HLADR-/CD123+/CD11c+
NK.CD56hi	CD45+/CD3-/HLADR+/CD16-/CD56hi
NK.CD56dim	CD45+/CD3-/HLADR-/CD16-/CD56dim
Mature Neutrophils	Live CD45+/Boolean gate excluding All other cell types / CD16+/CD11b+
Immature Neutrophils	Live CD45+/Boolean gate excluding All other cell types / CD16-/CD11b <sup>-</sup>
Eosinophils	CD45+/hi/CD16-
Basophils	CD45+/CD3-/CD16-/CD56-/CD20-/CD123+

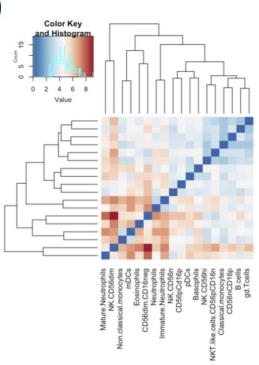
B



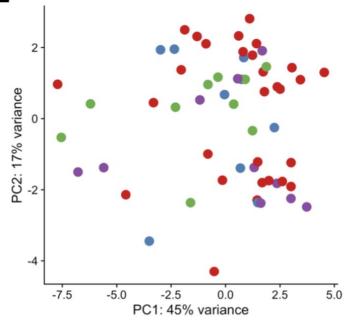
C



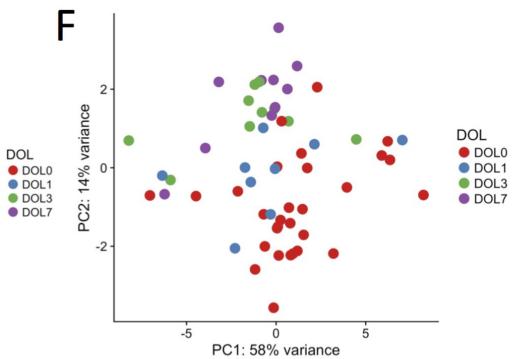
D



E

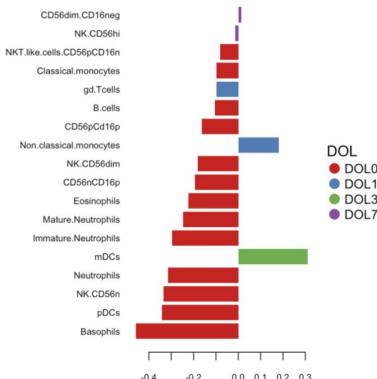


F

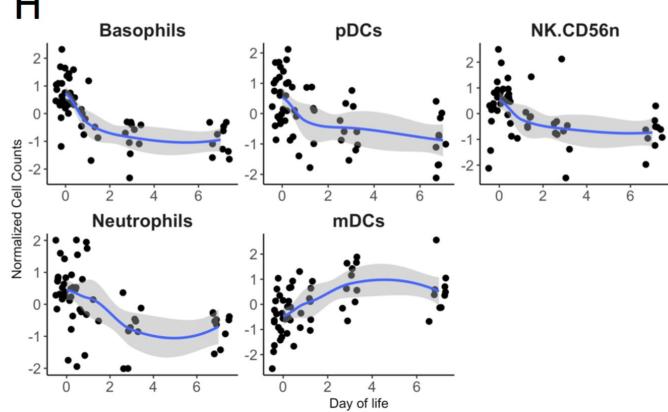


G

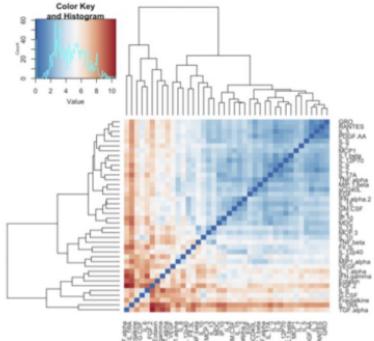
Contribution on comp 1



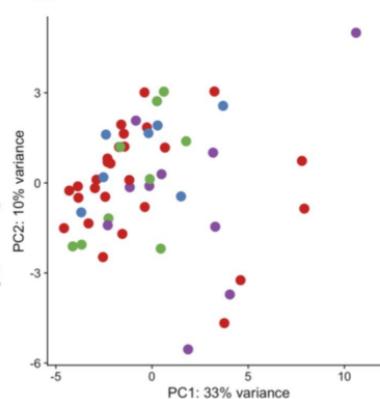
H



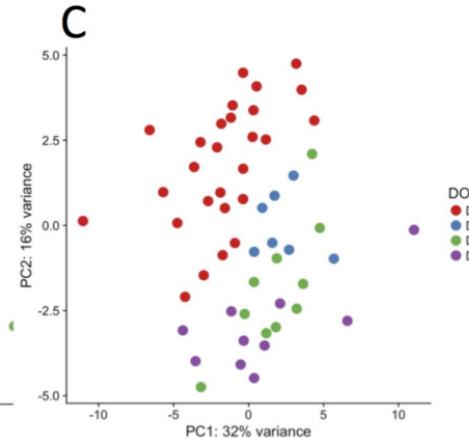
A



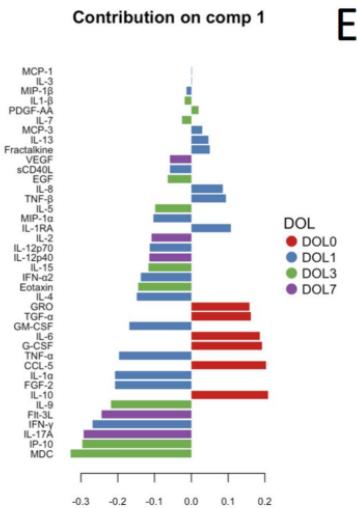
B



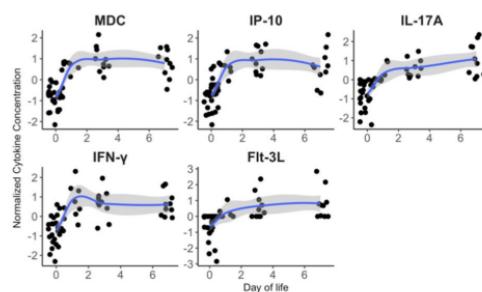
C



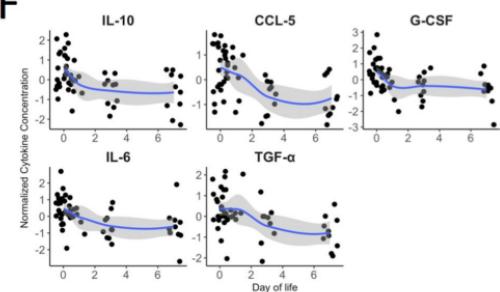
D

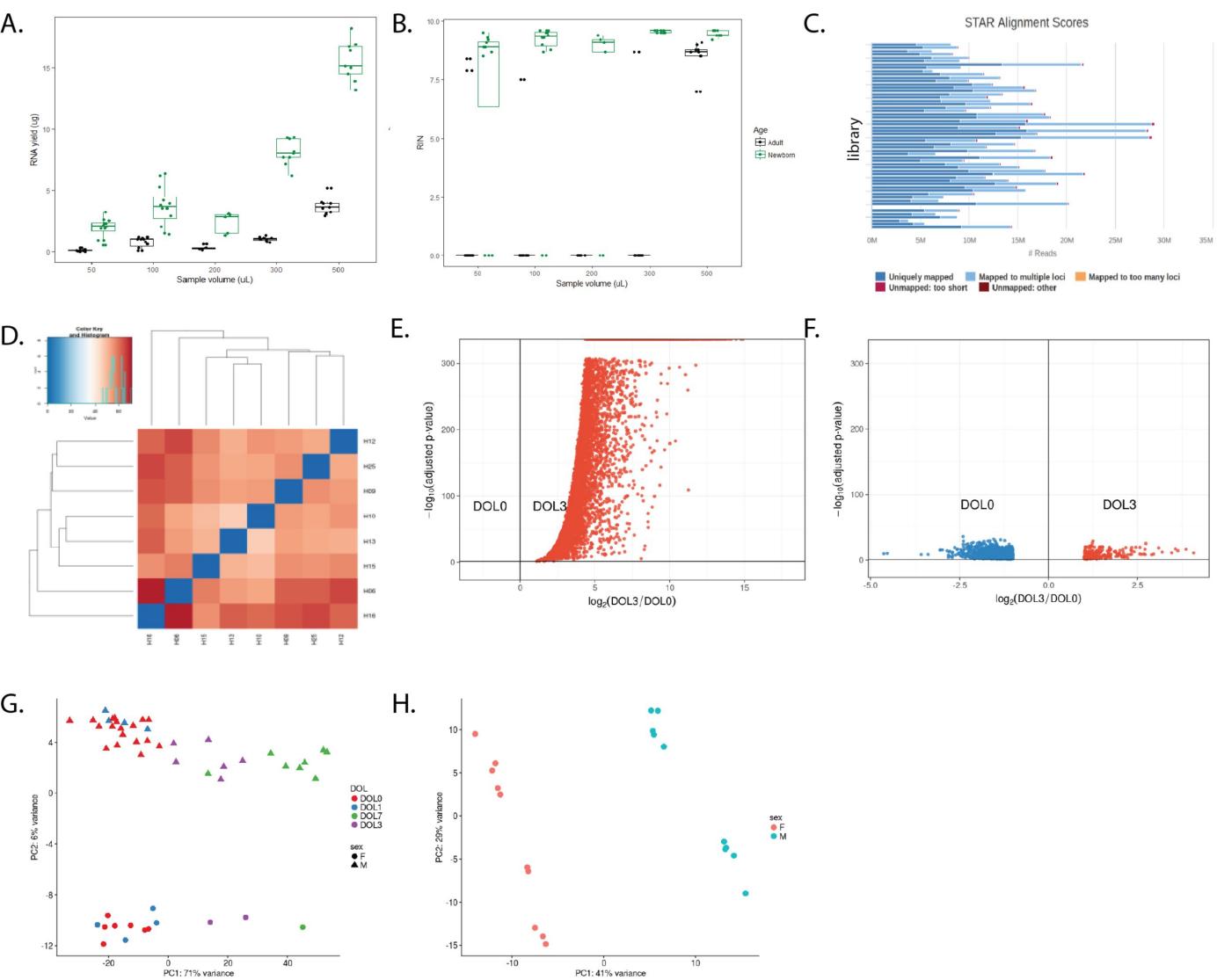


E

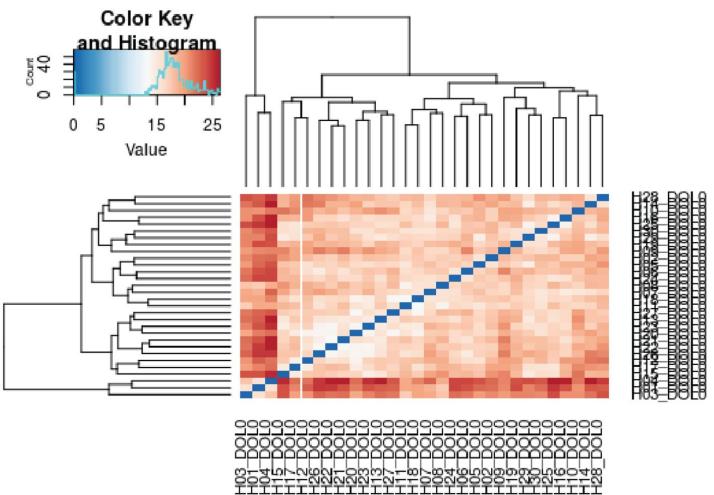


F

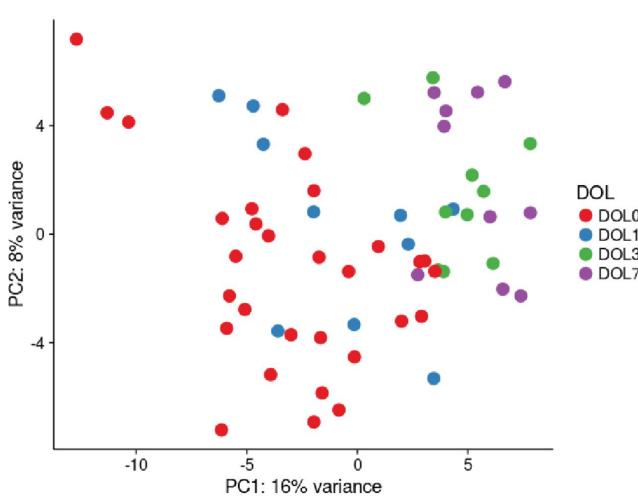




A.



B.



C.

