

# Multiomics integration of 22 immune-mediated monogenic diseases reveals an emergent axis of human immune health

John Tsang (✉ [john.tsang@nih.gov](mailto:john.tsang@nih.gov))

National Institutes of Health <https://orcid.org/0000-0003-3186-3047>

Rachel Sparks

Multiscale Systems Biology Section, Laboratory of Immune System Biology, NIAID, NIH

<https://orcid.org/0000-0001-8216-5084>

Nick Rachmaninoff

NIH

---

## Article

**Keywords:**

**Posted Date:** March 20th, 2023

**DOI:** <https://doi.org/10.21203/rs.3.rs-2070975/v1>

**License:**   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

**Additional Declarations:** There is **NO** Competing Interest.

---

# Multimomics integration of 22 immune-mediated monogenic diseases reveals an emergent axis of human immune health

Rachel Sparks<sup>1,18</sup>, Nicholas Rachmaninoff<sup>1,2,18</sup>, Dylan C. Hirsch<sup>1,18</sup>, Neha Bansal<sup>1</sup>, William W. Lau<sup>1,3</sup>, Andrew J. Martins<sup>1</sup>, Jinguo Chen<sup>4</sup>, Candace C. Liu<sup>1</sup>, Foo Cheung<sup>4</sup>, Laura E. Failla<sup>1</sup>, Angelique Biancotto<sup>4</sup>, Giovanna Fantoni<sup>4</sup>, Brian A. Sellers<sup>4</sup>, Daniel G. Chawla<sup>5</sup>, Katherine N. Howe<sup>6</sup>, Darius Mostaghimi<sup>1</sup>, Rohit Farmer<sup>4</sup>, Yuri Kotliarov<sup>4</sup>, Katherine R. Calvo<sup>7</sup>, Cindy Palmer<sup>6</sup>, Janine Daub<sup>6</sup>, Ladan Foruraghi<sup>6</sup>, Samantha Kreuzburg<sup>6</sup>, Jennifer Treat<sup>6</sup>, Amanda K. Urban<sup>8</sup>, Anne Jones<sup>9</sup>, Tina Romeo<sup>9</sup>, Natalie T. Deutch<sup>9</sup>, Natalia Sampaio Moura<sup>9</sup>, Barbara Weinstein<sup>10</sup>, Susan Moir<sup>15</sup>, Luigi Ferrucci<sup>11</sup>, Karyl S. Barron<sup>12</sup>, Ivona Aksentijevich<sup>9</sup>, Steven H. Kleinstein<sup>5,13,14</sup>, Danielle M. Townsley<sup>10</sup>, Neal S. Young<sup>10</sup>, Pamela A. Frischmeyer-Guerrero<sup>16</sup>, Gulbu Uzel<sup>6</sup>, Gineth Paola Pinto-Patarroyo<sup>9</sup>, Cornelia D. Cudrici<sup>17</sup>, Patrycja Hoffmann<sup>9</sup>, Deborah L. Stone<sup>9</sup>, Amanda K. Ombrello<sup>9</sup>, Alexandra F. Freeman<sup>6</sup>, Christa S. Zerbe<sup>6</sup>, Daniel L. Kastner<sup>9</sup>, Steven M. Holland<sup>6</sup>, John S. Tsang<sup>1,4,19,#</sup>

<sup>1</sup>Multiscale Systems Biology Section, Laboratory of Immune System Biology, NIAID, NIH, Bethesda, MD 20892, USA

<sup>2</sup>Graduate Program in Biological Sciences, University of Maryland, College Park, MD 20742, USA

<sup>3</sup>Office of Intramural Research, CIT, NIH, Bethesda, MD 20892, USA

<sup>4</sup>NIH Center for Human Immunology, NIAID, NIH, Bethesda, MD 20892, USA

<sup>5</sup>Interdepartmental Program in Computational Biology and Bioinformatics, Yale University, New Haven, CT 06511, USA

<sup>6</sup>Laboratory of Clinical Immunology and Microbiology, NIAID, NIH, Bethesda, MD 20892, USA

<sup>7</sup>Hematology Section, Department of Laboratory Medicine, NIH Clinical Center, Bethesda, MD 20892, USA

<sup>8</sup>Clinical Research Directorate, Frederick National Laboratory for Cancer Research, National Cancer Institute, NIH, Frederick, MD 21701, USA.

<sup>9</sup>Inflammatory Diseases Section, National Human Genome Research Institute, NIH, Bethesda, MD 20892, USA

<sup>10</sup>Hematology Branch, National Heart, Lung, and Blood Institute, NIH, Bethesda, MD 20892, USA

<sup>11</sup>Translational Gerontology Branch, National Institute on Aging, Baltimore, MD 21224, USA

<sup>12</sup>Division of Intramural Research, NIAID, NIH, Bethesda, MD 20892, USA

<sup>13</sup>Department of Immunobiology, Yale University School of Medicine, New Haven, CT 06510, USA

<sup>14</sup>Department of Pathology, Yale University School of Medicine, New Haven, CT 06510, USA

<sup>15</sup>Laboratory of Immunoregulation, NIAID, NIH, Bethesda, MD 20892, USA

<sup>16</sup>Laboratory of Allergic Diseases, NIAID, NIH, Bethesda, MD 20892, USA

<sup>17</sup>National Institute of Arthritis and Musculoskeletal and Skin Diseases, NIH, Bethesda MD 20892, USA

<sup>18</sup>These authors contributed equally

<sup>19</sup>Lead contact

# Correspondence: [john.tsang@nih.gov](mailto:john.tsang@nih.gov)

44 **Summary**

45

46 Monogenic diseases are often studied in isolation due to their rarity. Here we utilize multiomics  
47 to assess 22 monogenic immune-mediated conditions with age- and sex-matched healthy  
48 controls. Despite clearly detectable disease-specific and “pan-disease” signatures, individuals  
49 possess stable personal immune states over time. Temporally stable differences among  
50 subjects tend to dominate over differences attributable to disease conditions or medication  
51 use. Unsupervised principal variation analysis of personal immune states and machine learning  
52 classification distinguishing between healthy controls and patients converge to a metric of  
53 immune health (IHM). The IHM discriminates healthy from multiple polygenic autoimmune and  
54 inflammatory disease states in independent cohorts, marks healthy aging, and is a pre-  
55 vaccination predictor of antibody responses to influenza vaccination in the elderly. We  
56 identified easy-to-measure circulating protein biomarker surrogates of the IHM that capture  
57 immune health variations beyond age. Our work provides a conceptual framework and  
58 biomarkers for defining and measuring human immune health.

## 59 Introduction

60

61 Immune system dysregulation is central to diverse pathologies, including cancer, chronic  
62 inflammation, cardiovascular, and neurological diseases<sup>1</sup>. Immune-mediated disease results  
63 from a complex interplay of environmental, exposure history, and genetic factors. In contrast to  
64 polygenic diseases such as rheumatoid arthritis (RA) and systemic lupus erythematosus (SLE),  
65 monogenic diseases offer unique opportunities to highlight important mechanisms by which  
66 individual genes and associated pathways contribute to immune function in humans. For  
67 example, the study of patients with immunodeficiencies has illuminated the critical roles of the  
68 JAK-STAT network in orchestrating microbial defense and inflammatory processes at the  
69 organismal level in humans<sup>2,3</sup>; similarly, monogenic periodic fever syndromes have deepened  
70 our molecular understanding of inflammasomes and their roles in innate immunity and  
71 autoinflammatory diseases<sup>4</sup>.

72

73 Aside from comparison of genetic associations and gene expression quantitative trait loci in  
74 polygenic diseases<sup>5-8</sup>, immune-mediated diseases, in particular those of monogenic origin, have  
75 often been studied in isolation. Molecular and cellular attributes and biomarkers shared across  
76 diseases remained poorly defined, knowledge of which could help advance our understanding  
77 of both common and disease-specific pathophysiology and immune dysregulation, potentially  
78 pointing to multi-disease therapeutic targets. Importantly, the contribution of genetics to  
79 human immune variations can be highly variable and tends to wane by adulthood<sup>9</sup>; even  
80 monogenic disease patients with primary causal defects in the same gene can exhibit extensive  
81 clinical heterogeneity<sup>10</sup> with poorly understood molecular and cellular drivers. Thus, dissecting  
82 the inter- and intra-patient variations in diverse immune parameters both within and across  
83 diseases is critical to understanding disease- and patient-specific dysregulation beyond the  
84 causal gene and proximal pathways. Analyzing diverse monogenic diseases may also  
85 simultaneously reveal features of a normal, healthy immune system, which remains ill-defined  
86 because parameters quantifying immunological health remain elusive<sup>11</sup>. In principle, immune  
87 health metrics should not be defined based on features of the immune systems among healthy  
88 individuals alone, but also incorporate common features of immune pathologies as "negative"  
89 indicators of health. Simultaneous assessment of immune states in monogenic disease patients  
90 and matching healthy subjects may thus reveal quantifiable parameters of human immune  
91 health.

92

93 Here we have integrated multiomics profiling and clinical information to comparatively analyze  
94 22 monogenic immune-mediated disease cohorts together with age- and sex-matched healthy  
95 controls. Using this new dataset, we identified both disease-specific and shared ("pan-disease")  
96 signatures, and importantly, found that both patients and healthy subjects possessed  
97 temporally stable personal immune states independent of disease condition or medication  
98 use<sup>12-14</sup>. Integration of transcriptomic, serum protein, and peripheral blood cell frequency data  
99 revealed a quantitative metric of immune health through both bottom-up, unsupervised  
100 principal variation analysis of personal immune states and supervised machine learning  
101 analyses that discriminated between healthy individuals and sick patients. This metric also  
102 marks healthy aging and is associated with the antibody responses to influenza vaccination in

103 the elderly. We also uncovered easy-to-measure serum protein surrogates of this metric that  
104 capture immune health variations among healthy individuals beyond age. Beyond our specific  
105 findings, this rich dataset can serve as a resource for the research community to probe these  
106 specific monogenic disorders more deeply, for example, by generating new hypotheses. Our  
107 work paves the way for a more quantitative understanding of human immune health and  
108 provides a unique dataset for further exploration.

109  
110

## 111 **Results**

112

### 113 **A multiomics compendium of 22 monogenic immune-mediated diseases reveals temporally** 114 **stable individual differences tend to be the dominant source of variation**

115

116 We employed multiomics analyses of circulating immune cells involving whole blood  
117 transcriptomics, measurements of more than 1300 circulating proteins from serum (using the  
118 Somalogic platform), as well as immune cell frequencies and hematological parameters from a  
119 complete blood count (CBC) and clinical flow cytometry [TBNK: CD4+ and CD8+ T-cells, B-cells,  
120 natural killer (NK) cells] to comparatively analyze samples collected from 364 visits of 228  
121 patients (some patients had multiple samples collected at different visits/timepoints)—  
122 spanning 22 monogenic immune-mediated diseases—and 42 age- and sex-matched healthy  
123 subjects (Fig. 1a-c, Extended Data Fig. 1a-c, Table 1, Extended Data Table 1). Once data were  
124 generated, we set aside a set of subjects including patients from the majority of disease groups  
125 and matched healthy controls (see Table 1) to enable potential future independent validation  
126 or follow-up analyses (see Methods). This monogenic disease compendium includes primary  
127 immunodeficiencies, autoinflammatory disorders, and defects in hematopoiesis, each with  
128 known causal genetic mutations affecting major molecular and cellular networks and functions  
129 of the innate [e.g., NOD-, LRR- and pyrin domain-containing protein 3 (NLRP3)] and adaptive  
130 [e.g., signal transducer and activator of transcription 1 (STAT1)] immune systems. Disease  
131 manifestations cover a spectrum of features including frequent and severe infections,  
132 autoimmunity, allergy, and recurrent fever with inflammation (autoinflammation). Thus, this  
133 multi-disease cohort offers unique opportunities for examining the shared and distinct features  
134 of these natural genetic perturbations in humans at the molecular and cellular levels. To the  
135 best of our knowledge, this constitutes the first and largest multiomics/multimodal  
136 comparative map of diverse monogenic, immune-mediated diseases in humans.

137

138 To reduce data dimensionality and assess the correlation among parameters, weighted gene  
139 correlation network analysis (WGNCA)<sup>15</sup> was applied to the serum protein and transcriptomic  
140 data to derive co-expression modules separately for each data modality. This resulted in 12  
141 blood transcriptomic modules (TMs; Fig. 1d, Extended Data Table 2) and 10 protein modules  
142 (PMs; Fig. 1e, Extended Data Table 3). Most of the TMs were enriched for signatures of major  
143 immune cell types (e.g., B-cells in TM7; Extended Data Table 4, Extended Data Fig. 1d) or  
144 intracellular processes (Extended Data Table 4). A subset of the proteins also formed modules  
145 based on co-expression (Fig. 1e; Extended Data Table 3, which contains the full list of the 1300  
146 proteins), including a PM enriched for platelet and lymphocyte activation (PM6; Extended Data

147 Table 5), as well as other PMs enriched for tissue-specific proteins as annotated in the Human  
148 Protein Atlas<sup>16</sup>, such as bone marrow proteins in PM3 (OR = 23.70, adj.  $p = 1.7 \times 10^{-6}$ ) and spleen  
149 proteins in PM2 (OR = 11.18, adj.  $p = 4.6 \times 10^{-5}$ ) (Extended Data Table 6). In contrast to the highly  
150 modular nature of blood transcriptomic measurements (Fig. 1d), a large fraction (48%) of the  
151 proteins fell into the “gray” module, which contains “singleton” proteins that did not exhibit  
152 sufficient correlation with other proteins to be incorporated in a module (Fig. 1e). Interestingly,  
153 the gray module proteins were enriched for those expressed in the liver (OR = 4.67, adj.  $p =$   
154  $9.68 \times 10^{-8}$ ), small intestine (OR = 3.71, adj.  $p = 0.011$ ), and adipose tissue (OR = 4.00, adj.  $p =$   
155  $0.045$ ) (Extended Data Table 6). These observations are consistent with the notion that whole  
156 blood transcriptomic data mainly capture variation in circulating immune cell frequencies and  
157 cellular states that give rise to correlated, modular gene expression structures, while circulating  
158 protein levels reflect more diverse sources of variation, including those from circulating blood  
159 cells but also from tissues and potentially their status such as inflammation. The blood  
160 transcriptomic and serum protein measurements thus provide orthogonal, complementary  
161 information and together enable comprehensive assessment of phenotypically diverse  
162 individuals.

163  
164 Multiple sources contribute to variations in the level of a parameter (e.g., cell frequency or  
165 WGCNA module score), including those associated with disease and medications as well as  
166 inter-subject and temporal differences within individuals. Leveraging data from 63, 62, and 64  
167 subjects for the cell frequencies, whole blood transcriptomics, and serum proteins, respectively,  
168 from whom we had collected more than one sample over time (spanning 5 days to roughly 1  
169 year from 19 disease groups and healthy subjects, 25% quantile = 86 days, median = 130 days,  
170 75% quantile = 181 days), we fit a variance partition model<sup>17</sup> to estimate the relative  
171 contributions from the following sources: differences associated with disease, differences  
172 among patients with the same disease, medication/treatment effects, and intra-patient  
173 variations over time (Fig. 1f,g). A large fraction of the parameters, including blood transcripts  
174 and especially circulating proteins, was temporally stable within individual patients, i.e., the  
175 systematic differences between patients were larger than those in the same patient over time  
176 as indicated by the larger variance explained by the patient covariate (Fig. 1f,g; Extended Data  
177 Fig. 1e-g). Major medication categories, including steroids and immunosuppressants, could only  
178 account for a small fraction of the variance in most parameters (Extended Data Fig. 1h),  
179 suggesting that immune states of individuals were not broadly affected by these medications.  
180 Also unexpectedly, but consistent with the substantial temporally stable inter-subject  
181 variations, the differences between patients with the same disease (inter-subject variance  
182 explained by the patient) were often larger than the disease effects (i.e., group level average  
183 differences between disease and healthy: variance explained by the disease/condition label) for  
184 most of the serum protein and transcriptomic parameters (Fig. 1g). Jackknife analysis indicated  
185 that the variance explained by subject for all features is robust to sampling noise, particularly  
186 for the features with the highest variation explained by subject (Extended Data Fig. 2).  
187 Consistently, patients did not cluster by disease labels based on CBC/TBNC data alone, with  
188 healthy subjects intermixed with disease groups (Extended Data Fig. 3a,b), indicating that CBC  
189 and basic immune cell frequency data alone are insufficient to delineate health and disease.  
190 Together, these data suggest that factors such as the environment and exposure history play an

191 important role in shaping the immune state of an individual, even in adult patients with highly  
192 penetrant monogenic conditions.

193

### 194 **Pan-disease and disease-specific signatures**

195

196 We next derived and compared disease signatures, although our aim was to generate new  
197 hypotheses rather than “deep diving” mechanistically into any specific monogenic disease. We  
198 used linear models to derive signatures of individual disease conditions in comparison to  
199 matching healthy subjects accounting for age, sex, and major medication groups (Fig. 2a;  
200 Extended Data Fig. 3c). Despite the diversity of conditions, we detected signatures shared  
201 across diseases. These shared signatures had consistent directions of change across multiple  
202 diseases, including increases in red cell distribution width (RDW; a measure of the variation of  
203 erythrocyte volume<sup>18</sup>), TM2 (enriched for heme biosynthesis), and PM2, as well as decreases in  
204 TM6 (enriched for NK cells and CD8+ T-cells), NK cell frequencies, and PM6 (enriched for  
205 platelet related factors) (Fig. 2a,b; Extended Data Tables 7-9). RDW is known to be associated  
206 with all-cause mortality and several common diseases, including cardiovascular disease and  
207 cancer<sup>19</sup>, but it has not been assessed simultaneously across multiple pathologies including the  
208 monogenic diseases analyzed here. Proteins in PM2 spanned several inflammatory pathways  
209 (Extended Data Table 3), including interleukin-23 (IL-23), tumor necrosis factor  $\alpha$  soluble  
210 receptors 1 and 2, interferon (IFN)-related or -induced proteins [e.g., IP-10/CXCL10, I-  
211 TAC/CXCL11, monokine induced by gamma (MIG)/CXCL9], and the shed receptor sCD163 that  
212 might reflect macrophage activation in tissues<sup>20</sup>. Together, these signals may reflect both  
213 systemic and tissue inflammation shared across diseases.

214

215 As an example of how our comparative analysis may be explored to reveal disease-specific  
216 insights, we identified signatures more specific to individual or subgroups of diseases. For  
217 example, the PM2 score was highly elevated in deficiency of adenosine deaminase 2 (DADA2)  
218 patients and several PIDs such as STAT1 gain-of-function (STAT1 GOF) and X-linked chronic  
219 granulomatous disease (X-CGD), relative to healthy subjects (Fig. 2a,b; Extended Data Table 8).  
220 IL-23, a member of PM2, was elevated in DADA2 (Fig. 2c,d; Extended Data Table 10), even  
221 though it is not a known marker of this disease. IL-23 was positively correlated with IFN- $\gamma$  in  
222 DADA2 patients (Fig. 2e), consistent with the fact that IL-23 can induce IFN- $\gamma$  production in  
223 several cell types such as  $\gamma\delta$  and CD8+ T-cells<sup>21</sup>. Although we verified that this increase in IL-23  
224 was not driven purely by changes in cell frequencies by fitting an additional model controlling  
225 for major cell subset frequencies (Extended Data Table 12, see Methods), DADA2 patients with  
226 high IL-23 tended to have decreased platelets, neutrophils, and total B-cells (Fig. 2e). These  
227 phenotypes are consistent with bone marrow biopsies from some of these DADA2 patients that  
228 showed decreased cellularity and B-cell precursors. Interestingly, like DADA2, some GATA2  
229 deficiency (GATA2) patients also had lower peripheral blood cell counts but *decreased* levels of  
230 circulating IL-23 (Fig. 2c), suggesting that the connection between circulating IL-23 level and  
231 bone marrow status in DADA2 patients is distinct from that in other diseases with bone marrow  
232 failure or low peripheral cell count phenotypes.

233

234 Elevated type I IFN (IFN-I) blood transcriptional signatures have been found in monogenic and  
235 polygenic inflammatory diseases such as Aicardi-Goutières syndrome and SLE, respectively<sup>22,23</sup>.  
236 Here DADA2, STAT1 GOF, X-CGD, and p47<sup>phox</sup>CGD (p47-CGD) had clear IFN-I signatures as  
237 reflected by elevation in TM1 (FDR < 0.2; Fig. 2a, Extended Data Table 9). This is to be expected  
238 for the STAT1 GOF patients given their elevated STAT1-dependent signaling<sup>24</sup>. However, the  
239 CGDs, not typically known as interferonopathies<sup>22</sup>, had the most elevated TM1 scores  
240 compared to healthy (Extended Data Table 9), which were also significantly higher than STAT1  
241 GOF (X-CGD vs STAT1 GOF: logFC = 0.83,  $p = 0.001$ ; p47-CGD vs STAT1 GOF: logFC = 0.82,  $p =$   
242 0.002). Relative serum concentrations of the IFN-inducible protein I-TAC/CXCL11, as well as  
243 STAT1 itself, were higher in X-CGD and STAT1 GOF patients relative to healthy subjects  
244 (Extended Data Table 10), with circulating STAT1 protein concentrations significantly higher in  
245 X-CGD compared to STAT1 GOF (logFC = 0.83,  $p = 0.006$ ). Consistently, IFN-inducible transcripts  
246 in TM1 tended to be elevated in both the CGDs and STAT1 GOF patients compared to healthy,  
247 but again the elevations appeared stronger in the CGDs than the STAT1 GOF (Fig. 2f, Extended  
248 Data Table 11). We additionally verified that this increase in TM2 score was not driven purely by  
249 changes in cell frequencies by fitting an additional model controlling for major cell frequencies  
250 (Extended Data Table 12, see Methods). Together, these results suggest that IFN-I signatures  
251 and related pathways may be a good source of biomarkers and therapeutic targets for CGD.  
252

253 In addition to examining differences in relation to healthy subjects, we also compared each  
254 disease against all other diseases excluding the healthy subjects. Surprisingly, this other-  
255 disease-as-background map was qualitatively similar to the healthy-as-background map  
256 (Extended Data Fig. 3d). For example, the autoinflammatory diseases tumor necrosis factor  
257 receptor-associated periodic syndrome (TRAPS), familial cold autoinflammatory syndrome  
258 (FCAS; NLRP3-associated autoinflammatory disease-mild) and familial Mediterranean fever  
259 (FMF) as a group differed from the healthy subjects and other diseases by similar signatures,  
260 including lymphocyte and B-cell counts that trended higher than other diseases, which to the  
261 best of our knowledge has not been described for this group of diseases. These disease-specific  
262 signatures suggest that predictive models could also be built to help identify possible diagnoses  
263 for patients. Indeed, Random Forest (RF) classifiers built for the major disease groups (Extended  
264 Data Fig. 3e,f) revealed that STAT3 dominant-negative (STAT3 DN) disease patients (also known  
265 as autosomal dominant hyper-IgE syndrome or Job's Syndrome) could easily be differentiated  
266 from other patients in the cohort based on cross-validation analysis (0.98 AUC, STAT3 DN  $n =$   
267 21, Other  $n = 127$ ), as could the p47-CGD/X-CGD patients (0.99 AUC, CGD  $n = 37$ , Other  $n =$   
268 111). In contrast, predictive performance was poorer for STAT1 GOF (0.64 AUC, STAT1 GOF  $n =$   
269 15, Other  $n = 133$ ) and FMF (0.56 AUC, FMF  $n = 10$ , Other  $n = 138$ ), which may reflect disease  
270 and patient heterogeneity, some of which might not be well captured by the parameters  
271 measured, or because FMF patients may have been sampled largely at clinically quiescent time  
272 points<sup>25</sup>. Together, our data provide a rich resource for the biomedical community and highlight  
273 shared and disease-specific cellular, transcriptional, and serum protein signatures of diverse  
274 monogenic immune-mediated diseases. The shared signatures in particular point to commonly  
275 dysregulated pathways and processes in the immune system independent of disease-specific  
276 pathologies.  
277



278 **Integration of transcriptomic and serum protein personal immune profiles revealed an**  
279 **emergent axis of immune health**

280  
281 Our disease signature analyses suggest that both overlapping and unique information is  
282 provided by blood transcriptomic and circulating serum protein data. To assess whether the  
283 shared information between them can provide more integrated measures to examine individual  
284 patient-to-patient heterogeneity without knowledge of disease labels (Fig. 1b), we used JIVE<sup>26</sup>  
285 to infer latent components shared among the temporally stable transcriptomic and serum  
286 protein parameters (Fig. 3a, see Methods). JIVE decomposes the data into components,  
287 including the shared information between both data types reported as “joint principal  
288 components” (jPCs) and information captured uniquely by each data type (individual principal  
289 components; iPCs).

290  
291 JIVE revealed that approximately 20% of the variation (or information) in each data type was  
292 shared (Fig. 3b) with jPCs 1, 2 and 3 capturing 56%, 28% and 16% of the joint variation,  
293 respectively. The unique information in each data type could be further decomposed into 25  
294 and 18 iPCs for the transcriptomic and serum protein data, respectively (Extended Data Fig.  
295 4a,b; Extended Data Table 13). The top two transcriptomic data-specific iPCs reflected diverse  
296 processes and cell types, such as enrichments of neutrophil degranulation, monocytes, and IFN-  
297 I signatures. The top two protein-specific iPCs similarly exhibited enrichments for several  
298 functions, including extracellular matrix proteins, neurological processes and certain signaling  
299 pathway components (Extended Data Tables 14 and 15). These JIVE results suggest that not  
300 only can blood transcriptomic and serum protein data mutually reinforce each other based on  
301 the shared information present in jPCs (see below), each on its own can provide potentially  
302 non-redundant information and should thus be collected and analyzed together in human  
303 immune profiling studies.

304  
305 We next focused on the shared jPC components because they captured information from both  
306 data modalities and thus provide robust information regarding personal immune states and  
307 patient-to-patient heterogeneity. jPC1 appeared to quantify the extent of *attenuation* in  
308 inflammation-related processes as evident by: 1) jPC1 was negatively correlated with the  
309 neutrophil-to-lymphocyte ratio, which is a known marker of systemic inflammation and  
310 elevated in acute infections and cancer<sup>27,28</sup>, and positively correlated with B- and T-cell  
311 frequencies (Fig. 3c; Extended Data Table 16); and 2) jPC1 was negatively associated with innate  
312 immunity, inflammation, and IFN related processes (Fig. 3c, Extended Data Fig. 4c, Extended  
313 Data Table 15). jPC2 was negatively associated with the counts of multiple cell lineages,  
314 including WBC, platelet, neutrophils, monocytes, lymphocytes, and hemoglobin (Fig. 3c,  
315 Extended Data Table 16), suggesting that it captured hematopoietic output capacity. Indeed, it  
316 was also negatively associated with a combined score derived from the above immune cell  
317 populations (Extended Data Fig. 4d). This negative association was especially apparent within  
318 the DADA2, GATA2, and activated PI3K delta syndrome 1 (p110 $\delta$ ; APDS1) patient groups  
319 (Extended Data Fig. 4d), consistent with the loss of one or more cell lineages being a shared  
320 characteristic of these diseases<sup>29–32</sup>. Interestingly, for GATA2, patients with the highest jPC2

321 scores were also more likely to have dysplastic marrow (Extended Data Fig. 4d), a known  
322 complication of the disease<sup>30</sup>.

323  
324 We next placed individual patients onto the two-dimensional jPC1 vs. jPC2 space to visually  
325 examine inter-patient and inter-disease heterogeneity (Fig. 3d). Most disease groups and  
326 healthy subjects displayed narrower or comparable within-group variations along jPC2 than  
327 jPC1, but a few (DADA2, APDS1, CTLA4 haploinsufficiency) appeared to have higher jPC2  
328 differences among patients (Extended Data Fig. 4e), which, at least for DADA2 and APDS1, is  
329 expected given that jPC2 reflects hematopoietic output and bone marrow pathologies are  
330 known to be variable in both groups of patients<sup>33,34</sup>. Consistent with the notion that jPC1 might  
331 reflect systemic inflammatory burden (or immune “health”) and the expectation that patients  
332 would have elevated inflammation and potentially poorer immune health, jPC1 score is  
333 significantly higher in healthy subjects than patients (Fig. 3e), and this was robust to adjusting  
334 for major cell frequencies (Extended Data Table 12). Intriguingly, however, healthy subjects  
335 alone spanned a wide range along jPC1, similar to or even exceeding that of patients within  
336 individual disease groups, suggesting that jPC1 might provide quantitative information on  
337 systemic inflammation among even clinically healthy individuals.

338  
339 To test whether jPC1 emerged solely because of differences between sick patients and healthy  
340 subjects, we removed healthy subjects from our cohort and repeated the JIVE analysis.  
341 Strikingly, the resultant jPCs were highly correlated with those previously computed with HCs  
342 included (Fig. 3f;  $r = 0.98, 0.97, 0.92$ , respectively, for jPCs 1, 2, and 3). In fact, even if only  
343 healthy subjects were used to derive the jPCs, the resultant jPC1 was still significantly  
344 correlated with the original jPC1 derived from patients and HCs together (Fig. 3f). These results  
345 together suggest that the major emergent axis of immune variation within healthy subjects  
346 alone (i.e., derived in a totally unsupervised manner) is surprisingly similar to that obtained  
347 from sick patients with diverse monogenic immune-mediated diseases. These observations  
348 provide further support that this axis captures important information about immune health in  
349 diverse individuals.

350  
351 In addition to the healthy subjects, most disease groups such as STAT3 DN, GATA2, and STAT1  
352 GOF, spanned a wide range along jPC1 (Fig. 3d). The extensive overlap of healthy subjects and  
353 STAT3 DN patients is notable given that these patients could be easily distinguished from  
354 healthy subjects based on a few parameters as described in the disease classification analysis  
355 above (Extended Data Fig. 3e,f), suggesting that jPC1 captures immune health related  
356 phenotypes distinct from disease-specific deviations from healthy. On the “less healthy”, lower  
357 end of the jPC1 spectrum were CGDs; they also had extensive heterogeneity along jPC1, which  
358 is consistent with their wide spectrum of clinical presentations, including frequent infections,  
359 colitis, and pulmonary disease<sup>35</sup>, although further assessment would be needed to ascertain  
360 potential correlations between jPC1 and clinical phenotypes in larger patient cohorts. Patients  
361 with p47-CGD also trended higher than X-CGDs ( $p = 0.09$ , Wilcoxon test), consistent with the  
362 tendency for less severe disease in p47-CGD compared to X-CGD patients<sup>36</sup>. Together, our  
363 unbiased integration of blood transcriptomic and circulating protein data revealed an emergent

364 axis of immune health that delineates both inter-disease and inter-subject heterogeneity in  
365 patient and healthy populations.

366

### 367 **A quantitative metric of human immune health**

368

369 The emergence of pan-disease signatures (Fig. 2a) and an immune health axis, jPC1, (Fig. 3d)  
370 prompted us to assess whether supervised machine learning could help refine our immune  
371 health metric and the associated correlates of health and disease. We tested several RF  
372 healthy-versus-all-disease classifiers using temporally stable parameters as inputs, each using a  
373 different combination of data modalities (Fig. 4a) and assessed its performance with leave one  
374 out cross-validation (LOOCV). The classifier using all data modalities [including the use of  
375 singleton, grey module proteins (Fig. 1e)] had the best performance (Fig. 4b, Extended Data Fig.  
376 5a). It showed similar prediction performance in the independent (thus never-been-seen) set of  
377 patients and healthy subjects we set aside immediately after data generation but before any  
378 analysis began (these subjects were not included in the initial LOOCV evaluation or any of the  
379 analyses described in this manuscript except here in this independent robustness check;  
380 Extended Data Fig. 5b). This classifier revealed top parameters that contributed to the  
381 prediction [as measured by permutation tests of the global variable importance (GVI) –  
382 Extended Data Table 17]. These include RDW and parameters capturing systemic inflammation  
383 (sialoadhesin, C-reactive protein, PM2) and myeloid cell/macrophage signals (MIP-1 $\alpha$ , LD78 $\beta$ ),  
384 as well as the frequency of circulating NK cells (Fig. 4c, Extended Data Fig. 5c,d). These together  
385 revealed common deviations of disease from normal and are broadly concordant with the  
386 qualitative pan-disease signatures above (Fig. 2a).

387

388 In essence, our RF classifier had learned from a diverse set of monogenic diseases (i.e., as  
389 “negative” examples of health) against healthy subjects (“positive” examples) what a healthy  
390 immune system should (or should not) look like. Thus, we next used our classifier to assign each  
391 sample an “immune health metric” (IHM) score that reflects the probability that the sample  
392 belongs to the healthy group (see Methods, Extended Data Table 18). Despite jPC1 being  
393 derived in an unsupervised manner (i.e., without labeling the subjects with their  
394 disease/condition or healthy status), the IHM was highly correlated with jPC1 in patients with  
395 disease alone or in the healthy subjects only (Fig. 4d), but less so with the other jPCs (Extended  
396 Data Fig. 5e). As seen with jPC1 (Fig. 3d,e), the healthy subjects displayed a broad range of IHM  
397 scores (ranging from the very healthy to presumably the less healthy), but their median IHM  
398 score was significantly higher than that of most disease groups (Fig. 4e,f). Furthermore,  
399 consistent with the intuitive notion that immune health declines with age given that older  
400 individuals have elevated risk of immune-mediated diseases and tend to respond more poorly  
401 to infections and vaccinations compared to the young<sup>37</sup>, the IHM score and jPC1 were both  
402 negatively correlated with age in healthy individuals (Fig. 4g). Since certain cell frequencies are  
403 known to decline with age<sup>37</sup>, we verified that the IHM was correlated with age in healthy  
404 individuals even after controlling for cell-frequencies (Extended Data Table 12). Additionally,  
405 the IHM classifier could not have directly learned age-associated signals by training on patients  
406 versus healthy subjects because these two groups had indistinguishable age distributions in our  
407 cohort (KS test,  $D = 0.17$ ,  $p = 0.41$ , Extended Data Fig. 1a,b). This negative age association also

408 suggests that older healthy subjects resembled sick patients according to the IHM and age is a  
409 major contributor to IHM variability in the clinically healthy population. Thus, supervised  
410 (resulted in the IHM) and unsupervised (resulted in jPC1) analyses converged to a concordant  
411 metric of immune health.

412

413 **IHM is associated with common immune-mediated disease, vaccine responses in the elderly,**  
414 **and serum protein changes in healthy aging**

415

416 To assess the generalizability of the IHM beyond the monogenic diseases we studied, we sought  
417 to validate and further characterize the biological relevance of the IHM using independent  
418 datasets (Fig. 5a). First, we assessed the IHM in common autoimmune/inflammatory diseases  
419 distinct from the rare monogenic ones we examined above by using blood transcriptomic data  
420 from a published meta-analysis of 21 independent human datasets of type 1 diabetes,  
421 sarcoidosis, RA, and multiple sclerosis (Extended Data Table 19)<sup>38-40</sup>. We estimated the  
422 coherent deviation (meta-effect size) between disease and healthy subjects across the four  
423 diseases for every transcript and the transcriptional signature scores of the IHM, jPC1, and the  
424 top predictive markers from the IHM (the IHM and jPC1 signatures comprise blood transcripts  
425 correlated with the IHM or jPC1 – herein referred to as the “IHM and jPC1 blood transcriptional  
426 signatures”; Extended Data Table 20; see Methods). We found that these transcriptional  
427 signature scores were both significantly different between the four common diseases and  
428 healthy controls in the expected directions (Fig. 5b; Extended Data Fig. 6a,b; Extended Data  
429 Tables 21 and 22). Thus, the IHM can delineate health vs. disease in a different set of diseases  
430 common in the human population.

431

432 We next evaluated whether pre-vaccination immune health as reflected by the IHM might be  
433 predictive of responses to vaccination, a well-defined immune perturbation, and a potential “*in*  
434 *vivo*” readout of the consequences of having different levels of the IHM (Fig. 5a). We focused  
435 on the elderly population only because the extensive immune variability among the elderly is  
436 less well understood and baseline predictors of responses have been elusive in this population  
437 despite the fact that older individuals are known to have attenuated vaccination responses  
438 compared to the young<sup>41</sup>. Using meta-analysis of publicly available pre-vaccination blood  
439 transcriptomic data from four cohorts of older adults (61-96 years)<sup>42</sup>, we found that the IHM is  
440 indeed positively associated with antibody responses to influenza vaccination [summary effect  
441 size = 0.45 (weighted Hedge’s *g* between high and low responders across data sets), *p* = 0.046;  
442 Fig. 5c, Extended Data Fig. 6c]. Thus, the IHM could delineate baseline immune variation  
443 associated with vaccination outcomes among the elderly.

444

445 We next further assessed IHM-age associations in a published independent proteomic study  
446 (the “Baltimore Aging Study”) of 240 healthy subjects evenly distributed between the ages of  
447 20 and 90<sup>43</sup> (Fig. 5a). We derived circulating protein surrogates of the IHM (Extended Data  
448 Table 23) and found that the IHM protein surrogate score was indeed negatively correlated  
449 with age in this cohort (Fig. 5d). Interestingly, there was only a small overlap between the IHM  
450 circulating protein surrogates and those identified as associated with age in the original  
451 Baltimore study (Extended Data Fig. 6d), perhaps because the IHM is more reflective of aging-

452 related immune health and inflammation<sup>37</sup> while those identified in the original study captured  
453 aging signals from more biologically diverse sources. Furthermore, the IHM was not correlated  
454 with the level of circulating interleukin-6 (IL-6), a widely-studied cytokine linked to aging-  
455 related inflammation<sup>44</sup>, in healthy individuals from either the Baltimore Aging Study (Extended  
456 Data Fig. 6e) or our cohort (Extended Data Fig. 6f). However, IL-6 was correlated with the IHM  
457 when assessed in patients in our cohort (i.e., excluding healthy subjects; Extended Data Fig. 6f),  
458 partly because it was substantially elevated in some X-CGD and STAT1 GOF patients who had  
459 low IHM scores (data not shown). Thus, aspects of IL-6 related inflammation may be captured  
460 by the IHM in sick patients. In contrast, we did find that CXCL9/MIG, a marker known to be  
461 downstream of IFN- $\gamma$  signaling and associated with aging-related inflammation<sup>45</sup>, is correlated  
462 with the IHM in both healthy subjects and patients alone (Extended Data Fig. 6g). However, the  
463 IHM remained negatively correlated with age independent of CXCL9/MIG (Extended Data Table  
464 24) and its negative association with age did not change even when PM2, the protein module in  
465 the IHM that contained CXCL9/MIG (Fig. 2c), was removed during the derivation of the IHM  
466 (Extended Data Fig. 6h). Together, our results validate the utility and biological relevance of the  
467 IHM in distinct settings using independent datasets: a signature shared among common  
468 autoimmune and inflammatory diseases, a baseline correlate of vaccination responses in the  
469 elderly, and a biomarker of healthy aging.

470

#### 471 **The cellular origin of the IHM transcriptional signature**

472

473 To better understand the cellular origins of the IHM/jPC1 blood transcriptional signature, we  
474 utilized gene expression data of sorted peripheral immune cells from an independent study of  
475 10 immune-mediated diseases (including RA and SLE) and healthy controls<sup>5</sup>. We computed the  
476 signature scores for the IHM and jPC1 within each cell type and tested whether these  
477 signatures were elevated in healthy controls compared to patients with immune-mediated  
478 diseases in the cohort (Fig. 6a; Extended Data Table 25). We found higher IHM and jPC1  
479 signature scores in healthy individuals across nearly all the evaluated cell types (Fig. 6b,c),  
480 suggesting that the IHM and jPC1 reflect conserved transcriptional differences across a broad  
481 range of peripheral immune cells present in individuals with both polygenic and idiopathic  
482 immunological disease. These findings also further support the notion that the IHM/jPC1 and  
483 their constituent parameters are robust biomarkers of immune health beyond rare monogenic  
484 immune diseases.

485

486 Since the IHM was associated with healthy aging (Fig. 4g, 5d), we also used only the healthy  
487 subjects from the gene expression data of sorted immune cells<sup>5</sup> to assess what type of cells  
488 might have contributed to the age association. Compared to the disease-versus-healthy  
489 observations above, the IHM and jPC1 signature scores were negatively correlated with age in a  
490 subset of the cell types, most prominently in low density granulocytes (LDGs), a subset of naïve  
491 regulatory T-cells (Fr. 1 nTregs in Ota *et al*<sup>5</sup>), and certain T-cell subsets such as CD8+ effector  
492 memory T-cells expressing CD45RA (TEMRA) (Fig. 6d,e). These results suggest that while  
493 common blood transcriptional changes associated with immunological diseases are conserved  
494 broadly across multiple peripheral immune cell types (Fig. 6b; Extended data table 26), healthy  
495 aging-related decline in the IHM could be attributed to a more specific subset of these cell

496 types. However, this observed difference could be partly driven by differences in statistical  
497 power given the larger effect and sample sizes in the disease-versus-healthy comparison. Taken  
498 together, the IHM blood transcriptional signature captures shared signals from multiple  
499 peripheral immune cell types and subsets.

500

### 501 **IHM captures immune variation in healthy individuals beyond age**

502

503 Given the broad cell-type origin of the IHM, some of its serum protein surrogates/correlates  
504 (Extended Data Table 27) may represent cell extrinsic factors that could induce similar  
505 transcriptional profiles across different cell types – circulating serum proteins also represent  
506 easy-to-assay biomarker targets for routine clinical monitoring. Among the circulating protein  
507 correlates of the IHM, we noticed that some proteins were highly correlated with the IHM in  
508 both healthy subjects only and in patients (Extended Data Fig. 7a, Extended Data Table 27).  
509 These proteins include the IFN-induced IP-10/CXCL10 and beta-2 microglobulin, suggesting that  
510 interferons and related factors may be among the underlying cell-extrinsic inducers.

511

512 Given that age is a key contributor to IHM (and jPC1) variation, particularly in healthy subjects,  
513 and yet unexplained variation remains beyond age (Fig. 4g, 5d), we next assessed the extent by  
514 which the associations between serum proteins and the IHM depended on age (Fig. 6f).

515 Surprisingly, they were largely independent of age (Fig. 6g). For example, certain proteins were  
516 highly correlated with the IHM, including IP-10/CXCL10 and other negative indicator of immune  
517 health (lower left-hand corner in Fig. 6g), regardless of age in healthy individuals (Fig. 6g,  
518 Extended Data Table 27) or in sick patients alone (Extended Data Fig. 7b, Extended Data Table  
519 27). Interestingly, the positive correlates of the IHM (i.e., positive indicators of immune health –  
520 upper right-hand corner in Fig. 6g) were also independent of age. These include neurotrophin-3  
521 (Fig. 6h) and GDF11/GDF8 (GDF11 is also known as BMP-11), both of which have critical  
522 developmental and potentially “rejuvenation” functions such as neurodevelopment, patterning,  
523 and angiogenesis<sup>46–49</sup>. Together, these observations suggest that factors beyond those linked to  
524 aging are shaping immune health (as reflected by the IHM) in clinically healthy individuals and  
525 the IHM variation among healthy subjects alone reflects both age-dependent and age-  
526 independent biology. Thus, learning from diverse rare diseases as “negative” examples of  
527 health also revealed a quantitative metric that captures meaningful variations in clinically  
528 healthy individuals.

529

530

### 531 **Discussion**

532

533 Monogenic diseases are often studied in isolation due to their rarity, and thus the data and  
534 insight obtained from one condition cannot be easily compared to those of others. Here a  
535 unified approach was taken to simultaneously compare multiple rare immune-mediated  
536 conditions with natural genetic perturbations disrupting key pathways. To our surprise, despite  
537 penetrant genetic defects and clearly detectable common and disease-specific signatures, we  
538 observed that temporally stable, between-subject variation in cellular, transcriptomic, and  
539 circulating protein parameters dominates relative to the variation attributable to disease

540 condition, medication, age, and sex. This observation is consistent with the clinical  
541 heterogeneity often observed even within single monogenic disorders<sup>10</sup>, suggesting that  
542 environmental, exposure history, and other genetic factors [e.g., genetic modifiers of primary  
543 causal mutations<sup>50</sup>] together play important roles in setting and maintaining personal immune  
544 states. Indeed, various immune parameters have been found to be temporally stable over  
545 months in healthy individuals; some of these inter-subject differences were associated with  
546 responses to perturbations such as vaccination and autoimmune disease flares<sup>12-14</sup>. Here we  
547 have extended these concepts and observations to diverse monogenic patients with high-  
548 penetrance deleterious mutations affecting immune functions.

549 In general, there were both shared and modality-specific information provided by the  
550 transcriptomic and circulating protein data, suggesting that both should be measured to  
551 capture personal biological states when possible. Importantly, our results using the protein and  
552 transcriptional signatures were largely independent of circulating immune cell frequency, which  
553 is a major driver of blood transcriptomic profiles. Some of the circulating protein modules we  
554 uncovered may also reflect tissue status, as was postulated previously in a large proteomic  
555 study of older individuals<sup>51</sup>. Our findings raise the possibility that a targeted set of parameters  
556 comprising select blood immune cell frequencies, proteins, and transcripts could be developed  
557 from a multi-disease cohort like ours with the goal of optimizing both information overlap (to  
558 increase robustness) and uniqueness (to capture diverse, informative biological states) to track  
559 the health and disease status of individuals in the general population.

560 Our dataset serves as a valuable resource for hypothesis generation and exploratory analyses  
561 by the research community. As an example, we revealed that IFN-stimulated gene transcripts  
562 were elevated in the blood of CGD patients and often at higher levels than in STAT1 GOF  
563 patients. This was unexpected given that STAT1 GOF patients are known to have increased  
564 STAT1 signaling and transcription of IFN-stimulated genes due to their gain-of-function  
565 mutations in the STAT1 gene<sup>24</sup>. This observation suggests that JAK inhibitors, which have been  
566 successfully used to treat some inflammatory complications of STAT1 GOF patients<sup>52</sup>, may also  
567 be a therapeutic option for inflammatory complications of CGD. While IFN signatures have been  
568 reported in some inflammatory conditions<sup>53,54</sup>, their presence and relative magnitude have not  
569 been comparatively analyzed across multiple monogenic disorders. These observations and  
570 hypotheses highlight the power of the comparative approach taken to study monogenic  
571 diseases in this study.

572 Our bottom-up analysis of subject-level immune states revealed an axis (jPC1) of natural  
573 subject-to-subject variation captured by both blood transcriptomic and circulating protein data.  
574 Surprisingly, this was not driven by differences among diseases or between healthy and sick  
575 patients because a similar, correlated principal axis emerged from the data of sick patients or  
576 healthy subjects alone. This axis was also highly concordant with the IHM derived through a  
577 supervised machine learning analysis for differentiating healthy from sick patients in our  
578 cohort. Thus, the unsupervised and supervised analyses independently converged on a measure  
579 of immune health potentially applicable to diverse populations. Supporting this notion, the  
580 applicability of the IHM was validated in three independent and biologically distinct datasets.

581 First, we showed that the IHM signature was lower (associated with poor immune health) in  
582 patients from a meta-analysis of several polygenic autoimmune and inflammatory diseases.  
583 Second, it was associated, when evaluated pre-vaccination, with the antibody response to  
584 seasonal influenza vaccination in older individuals, pointing to a potential baseline determinant  
585 of vaccine responsiveness in this population. This is notable because the baseline immune  
586 statuses of the elderly are often highly heterogeneous and shaped by myriad complex factors  
587 (e.g., medications and comorbidities)<sup>41,55</sup>. Finally, it was negatively correlated with age in  
588 healthy subjects in our cohort and in a large independent cohort of healthy adults age ~20-90,  
589 consistent with the expectation that immune health declines with age. The IHM is based on a  
590 relatively small number of parameters and can be evaluated using circulating proteins from  
591 serum alone, and thus can potentially serve as an inexpensive tool for monitoring immune  
592 states and functions in diverse populations.

593 Given the applicability of the IHM in a range of biological scenarios, it is perhaps not surprising  
594 that IHM transcriptional scores appeared lower in nearly every peripheral immune cell type  
595 from patients with various polygenic or idiopathic immunological diseases. This coherent  
596 signature could be, at least partly, driven by cell-extrinsic factors, such as some of cytokines  
597 (interferons) and tissue growth/homeostatic factors (e.g., Neurotrophin-3) revealed by the IHM  
598 circulating protein correlate analysis. This result obtained using another independent dataset  
599 further validates the notion that the IHM likely has applicability beyond the monogenic  
600 conditions explored in this study. Interestingly, these coherent IHM signals across cell types  
601 were seen in only a subset of cell types when assessing the cell type specific correlation  
602 between the IHM transcriptional score and age in healthy subjects, including LDGs and some  
603 regulatory and effector memory T-cell subsets. LDGs (which includes low density neutrophils)  
604 and these T-cell subsets have been implicated in a spectrum of immunological and  
605 inflammatory conditions, including autoimmunity, cancer, and cardiovascular disease<sup>56-59</sup>. The  
606 age-related signals that we detected in Tregs and neutrophils confirm previous reports that  
607 aging contributes to their pathologic potential<sup>56,60</sup>.

608 Markers of systemic inflammation (e.g., CRP and serum amyloid A), RDW, and NK cell  
609 frequencies were some of the key constituents of the IHM. RDW and inflammatory markers  
610 were negative indicators of immune health. Increased RDW has been associated with human  
611 aging and several pathologies, including heart disease and cancer<sup>19</sup>, as well as mortality and  
612 morbidity risks (e.g., in Coronavirus Disease 2019<sup>61</sup>). While the mechanisms behind these  
613 associations are not entirely clear, increased RDW is known to reflect dysregulation of  
614 erythropoiesis and potential reductions in the rate of RBC turnover<sup>18,62</sup>. Conversely, higher NK  
615 cell numbers were associated with higher IHM scores. Aging, which is associated with the IHM  
616 in our study, is known to be associated with decreased NK cell production in the bone marrow.  
617 While it is unclear whether decreased bone marrow output or reduced expansion capacity of  
618 specific NK cell subsets played a role in the lower NK cell numbers we observed across multiple  
619 diseases, the association of both RDW and NK cell frequency with the IHM suggests that  
620 disruption of hematologic homeostasis may be involved.  
621



622 Inflammaging (chronic, sterile inflammation that increases with age) has been linked to age-  
623 related adverse outcomes such as cardiovascular disease. However, the inflammatory  
624 mechanisms or molecules responsible have not been well characterized<sup>37,44,63</sup>. Inflammaging  
625 has been linked to increased IL-6 in the literature, although there has been conflicting data<sup>63</sup>; IL-  
626 6 was neither correlated with the IHM in our study nor a key feature of an inflammatory aging  
627 (iAge) “clock” recently developed from ~1000 healthy individuals<sup>45</sup>. That study identified  
628 CXCL9/MIG as an informative feature of age-related inflammation. In our data, CXCL9 is a  
629 member of the protein module PM2, a key component of the IHM. PM2 also includes other  
630 inflammatory cytokines (e.g., IL-23) and IFN-related or -induced proteins (e.g., IP-10/CXCL10, I-  
631 TAC/CXCL11). As expected, the IHM was negatively correlated with CXCL9/MIG, but it remained  
632 correlated with age even when CXCL9/MIG and PM2 were removed, consistent with our  
633 findings that the protein IP-10/CXCL10 was negatively correlated with the IHM independent of  
634 age in healthy individuals only. More broadly, the IHM (and jPC1) was surprisingly variable even  
635 among apparently healthy subjects; the correlation between circulating proteins (including  
636 both negative and positive indicators of immune health) and the IHM in healthy subjects is also  
637 independent of age, suggesting that the IHM captures aspects of immune health not linked to  
638 age and inflammaging. Thus, the IHM, as measured by easy-to-assay serum protein parameters  
639 for example, could be applicable to the healthy population.

640  
641 It has been recognized that despite ample clinical tools for assessing general physiologic and  
642 organ system function and health (e.g., cardiovascular function), aside from the CBC, such tools  
643 are largely missing for the immune system<sup>11,64</sup>. This is partly because the function and  
644 pathology of the immune system are wide ranging and thus unified definitions and metrics of  
645 general immunological health have been elusive<sup>11,65,66</sup>. Here we have developed a framework  
646 for defining and quantifying immune health by searching for personal, temporally stable  
647 immune parameters enriched in health (i.e., in healthy subjects) but depleted in patients across  
648 diverse pathologies due to perturbations of normal immune functions. The resulting measure  
649 was surprisingly generalizable to different patient populations and healthy individuals. Further  
650 refinement and development of such approaches, e.g., by increasing the diversity and number  
651 of studied subjects including the incorporation of additional pathologies, utilizing  
652 measurements from tissues, and modeling potential modifiers such as sex and genetic factors,  
653 hold promise for the development of clinically useful immune health monitoring tools to  
654 advance personalized and preventative medicine<sup>67,68</sup>.

## 655 656 **Limitations of the Study**

657  
658 As expected, some of the observed immune variations across individuals in our cohort are  
659 reflected by information shared across correlated data modalities (e.g., circulating proteins,  
660 whole blood transcripts, and cell frequencies); however, all major results presented were  
661 robust to variations in circulating immune cell frequencies and still significant when controlling  
662 explicitly for cell-frequencies. Our analysis of temporal stability by estimating between-subject  
663 variations was limited by a relatively small number of patients with repeat samples. Despite this  
664 we observed consistent temporally stable, between-subject variations among data modalities,  
665 including cellular, transcriptomic, and circulating protein parameters, that dominate relative to

666 those attributable to disease condition, medication, age, and sex; these results are also robust  
667 to resampling noise as suggested by Jackknifing analysis. Although achieving mechanistic  
668 insights into any specific monogenic disease was not our goal, we demonstrated how this  
669 multimodal data could be used to yield new observations and hypotheses concerning disease  
670 etiology and therapeutic targets. For example, through our comparative study of interferon-  
671 related transcriptional signatures among several diseases, we were able to suggest JAK  
672 inhibitors as a possible therapeutic to further explore for CGD. Lastly, some of the major signals  
673 related to the IHM may partially reflect age-related decline of immune health and increase in  
674 inflammation in healthy individuals<sup>69</sup>. However, even when we examined the jPCs, which  
675 represent principal components of variation shared by the transcriptomic and serum protein  
676 data, there was considerable variation unexplained by age. Furthermore, similar positive and  
677 negative circulating protein correlates of the IHM emerged regardless of whether age was  
678 included as a co-variate. Thus, our work provides a broadly useful dataset and a conceptual  
679 framework and markers for defining and measuring human immune health.

## 680 **Acknowledgements**

681 We thank the patients and their families who participated in this study, as well as the NIH  
682 phlebotomy staff for their help and contribution to this project. We thank Philip Johnson and  
683 Ronald Germain for critical reading of the manuscript, and Cassie Seamon for assistance with  
684 healthy subject recruitment. Illustrations in Fig. 1a, 3a, 4a, 5a, 6a, and 6f were created using  
685 BioRender.com. This research was supported by: 1) the Intramural Research Programs of the  
686 NIAID, NHLBI and NHGRI, 2) the Intramural Research Programs of the NIH supporting the NIH  
687 Center for Human Immunology, and 3) federal funds from the National Cancer Institute, NIH,  
688 under Contract No. 75N91019D00024, Task Order No. 75N91019F00130. The content of this  
689 publication does not necessarily reflect the views or policies of the Department of Health and  
690 Human Services, nor does mention of trade names, commercial products, or organizations  
691 imply endorsement by the U.S. Government.

692

## 693 **Data and code availability**

694 The analysis ready data will be available under controlled access in dbGaP upon publication.  
695 NIH review of the clinical study protocols under which these samples were collected  
696 determined that dbGaP is the appropriate repository under which the data should be  
697 deposited. A dbGaP PHS number and BioProject number will be provided when the manuscript  
698 is accepted for publication at a peer reviewed journal. Software code for reproducing our  
699 analyses will be available at: <https://github.com/niaid/monogenic-immune-health>.

700

701

## 702 **Declaration of Interests**

703 The authors declare no competing interests.

704

## 705 **References**

- 706 1. Zhong, J. & Shi, G. Editorial: Regulation of Inflammation in Chronic Disease. *Front. Immunol.*  
707 **10**, (2019).
- 708 2. Casanova, J.-L., Holland, S. M. & Notarangelo, L. D. Inborn Errors of Human JAKs and STATs.  
709 *Immunity* **36**, 515–528 (2012).
- 710 3. Leonard, W. J., Lin, J.-X. & O’Shea, J. J. The  $\gamma$ c Family of Cytokines: Basic Biology to  
711 Therapeutic Ramifications. *Immunity* **50**, 832–850 (2019).
- 712 4. Manthiram, K., Zhou, Q., Aksentijevich, I. & Kastner, D. L. The monogenic autoinflammatory  
713 diseases define new pathways in human innate immunity and inflammation. *Nat. Immunol.*  
714 **18**, 832–842 (2017).
- 715 5. Ota, M. *et al.* Dynamic landscape of immune cell-specific gene regulation in immune-  
716 mediated diseases. *Cell* **184**, 3006-3021.e17 (2021).

- 717 6. Parkes, M., Cortes, A., van Heel, D. A. & Brown, M. A. Genetic insights into common  
718 pathways and complex relationships among immune-mediated diseases. *Nat. Rev. Genet.*  
719 **14**, 661–673 (2013).
- 720 7. Pickrell, J. K. *et al.* Detection and interpretation of shared genetic influences on 42 human  
721 traits. *Nat. Genet.* **48**, 709–717 (2016).
- 722 8. Sahni, N. *et al.* Widespread Macromolecular Interaction Perturbations in Human Genetic  
723 Disorders. *Cell* **161**, 647–660 (2015).
- 724 9. Brodin, P. *et al.* Variation in the Human Immune System Is Largely Driven by Non-Heritable  
725 Influences. *Cell* **160**, 37–47 (2015).
- 726 10. Tangye, S. G. *et al.* Human Inborn Errors of Immunity: 2019 Update on the Classification  
727 from the International Union of Immunological Societies Expert Committee. *J. Clin.*  
728 *Immunol.* **40**, 24–64 (2020).
- 729 11. Davis, M. M. A Prescription for Human Immunology. *Immunity* **29**, 835–838 (2008).
- 730 12. Brodin, P. & Davis, M. M. Human immune system variation. *Nat. Rev. Immunol.* **17**, 21–29  
731 (2017).
- 732 13. Kotliarov, Y. *et al.* Broad immune activation underlies shared set point signatures for  
733 vaccine responsiveness in healthy individuals and disease activity in patients with lupus.  
734 *Nat. Med.* **26**, 618–629 (2020).
- 735 14. Tsang, J. S. *et al.* Global Analyses of Human Immune Variation Reveal Baseline Predictors of  
736 Postvaccination Responses. *Cell* **157**, 499–513 (2014).
- 737 15. Langfelder, P. & Horvath, S. WGCNA: an R package for weighted correlation network  
738 analysis. *BMC Bioinformatics* **9**, 559 (2008).
- 739 16. Uhlén, M. *et al.* Tissue-based map of the human proteome. *Science* **347**, (2015).
- 740 17. Hoffman, G. E. & Schadt, E. E. variancePartition: interpreting drivers of variation in complex  
741 gene expression studies. *BMC Bioinformatics* **17**, 483 (2016).
- 742 18. Salvagno, G. L., Sanchis-Gomar, F., Picanza, A. & Lippi, G. Red blood cell distribution width:  
743 A simple parameter with multiple clinical applications. *Crit. Rev. Clin. Lab. Sci.* **52**, 86–105  
744 (2015).
- 745 19. Pan, J., Borné, Y. & Engström, G. The relationship between red cell distribution width and  
746 all-cause and cause-specific mortality in a general population. *Sci. Rep.* **9**, 16208 (2019).
- 747 20. Møller, H. J. Soluble CD163. *Scand. J. Clin. Lab. Invest.* **72**, 1–13 (2012).
- 748 21. Martínez-Barricarte, R. *et al.* Human IFN- $\gamma$  immunity to mycobacteria is governed by both  
749 IL-12 and IL-23. *Sci. Immunol.* **3**, (2018).
- 750 22. Lee-Kirsch, M. A. The Type I Interferonopathies. *Annu. Rev. Med.* **68**, 297–315 (2017).
- 751 23. Muskardin, T. L. W. & Niewold, T. B. Type I interferon in rheumatic diseases. *Nat. Rev.*  
752 *Rheumatol.* **14**, 214–228 (2018).
- 753 24. Okada, S. *et al.* Human STAT1 Gain-of-Function Heterozygous Mutations: Chronic  
754 Mucocutaneous Candidiasis and Type I Interferonopathy. *J. Clin. Immunol.* **40**, 1065–1081  
755 (2020).
- 756 25. Onen, F. Familial Mediterranean fever. *Rheumatol. Int.* **26**, 489–496 (2006).
- 757 26. Lock, E. F., Hoadley, K. A., Marron, J. S. & Nobel, A. B. JOINT AND INDIVIDUAL VARIATION  
758 EXPLAINED (JIVE) FOR INTEGRATED ANALYSIS OF MULTIPLE DATA TYPES. *Ann. Appl. Stat.* **7**,  
759 523–542 (2013).

- 760 27. Templeton, A. J. *et al.* Prognostic Role of Neutrophil-to-Lymphocyte Ratio in Solid Tumors: A  
761 Systematic Review and Meta-Analysis. *JNCI J. Natl. Cancer Inst.* **106**, (2014).
- 762 28. Russell, C. D. *et al.* The utility of peripheral blood leucocyte ratios as biomarkers in  
763 infectious diseases: A systematic review and meta-analysis. *J. Infect.* **78**, 339–348 (2019).
- 764 29. Lee, P. Y. Vasculopathy, Immunodeficiency, and Bone Marrow Failure: The Intriguing  
765 Syndrome Caused by Deficiency of Adenosine Deaminase 2. *Front. Pediatr.* **6**, 282 (2018).
- 766 30. McReynolds, L. J., Calvo, K. R. & Holland, S. M. Germline GATA2 Mutation and Bone Marrow  
767 Failure. *Hematol. Oncol. Clin. North Am.* **32**, 713–728 (2018).
- 768 31. Coulter, T. I. *et al.* Clinical spectrum and features of activated phosphoinositide 3-kinase  $\delta$   
769 syndrome: A large patient cohort study. *J. Allergy Clin. Immunol.* **139**, 597-606.e4 (2017).
- 770 32. Kallen, M. E., Dulau-Florea, A., Wang, W. & Calvo, K. R. Acquired and germline  
771 predisposition to bone marrow failure: Diagnostic features and clinical implications. *Semin.*  
772 *Hematol.* **56**, 69–82 (2019).
- 773 33. Aksentijevich, I., Sampaio Moura, N. & Barron, K. Adenosine Deaminase 2 Deficiency. in  
774 *GeneReviews*<sup>®</sup> (eds. Adam, M. P. *et al.*) (University of Washington, Seattle, 2019).
- 775 34. Dulau Florea, A. E. *et al.* Abnormal B-Cell Maturation in the Bone Marrow of Patients with  
776 Germline Mutations in PIK3CD. *J. Allergy Clin. Immunol.* **139**, 1032-1035.e6 (2017).
- 777 35. Arnold, D. E. & Heimall, J. R. A Review of Chronic Granulomatous Disease. *Adv. Ther.* **34**,  
778 2543–2557 (2017).
- 779 36. Kuhns, D. B. *et al.* Residual NADPH Oxidase and Survival in Chronic Granulomatous Disease.  
780 *N. Engl. J. Med.* **363**, 2600–2610 (2010).
- 781 37. Nikolich-Žugich, J. The twilight of immunity: emerging concepts in aging of the immune  
782 system. *Nat. Immunol.* **19**, 10–19 (2018).
- 783 38. Lau, W. W., Sparks, R., OMiCC Jamboree Working Group & Tsang, J. S. Meta-analysis of  
784 crowdsourced data compendia suggests pan-disease transcriptional signatures of  
785 autoimmunity. *F1000Research* **5**, 2884 (2016).
- 786 39. Shah, N. *et al.* A crowdsourcing approach for reusing and meta-analyzing gene expression  
787 data. *Nat. Biotechnol.* **34**, 803–806 (2016).
- 788 40. Sparks, R., Lau, W. W. & Tsang, J. S. Expanding the Immunology Toolbox: Embracing Public-  
789 Data Reuse and Crowdsourcing. *Immunity* **45**, 1191–1204 (2016).
- 790 41. HIPC-CHI Signatures Project Team & HIPC-I Consortium. Multicohort analysis reveals  
791 baseline transcriptional predictors of influenza vaccination responses. *Sci. Immunol.* **2**,  
792 eaal4656 (2017).
- 793 42. Avey, S. *et al.* Seasonal Variability and Shared Molecular Signatures of Inactivated Influenza  
794 Vaccination in Young and Older Adults. *J. Immunol.* **204**, 1661–1673 (2020).
- 795 43. Tanaka, T. *et al.* Plasma proteomic signature of age in healthy humans. *Aging Cell* **17**,  
796 e12799 (2018).
- 797 44. Ferrucci, L. & Fabbri, E. Inflammageing: chronic inflammation in ageing, cardiovascular  
798 disease, and frailty. *Nat. Rev. Cardiol.* **15**, 505–522 (2018).
- 799 45. Sayed, N. *et al.* An inflammatory aging clock (iAge) based on deep learning tracks  
800 multimorbidity, immunosenescence, frailty and cardiovascular aging. *Nat. Aging* **1**, 598–615  
801 (2021).
- 802 46. Chao, M. V., Rajagopal, R. & Lee, F. S. Neurotrophin signalling in health and disease. *Clin.*  
803 *Sci.* **110**, 167–173 (2006).

- 804 47. Omar, N. A., Kumar, J. & Teoh, S. L. Neurotrophin-3 and neurotrophin-4: The unsung heroes  
805 that lies behind the meninges. *Neuropeptides* **92**, 102226 (2022).
- 806 48. Rochette, L. & Malka, G. Neuroprotective Potential of GDF11: Myth or Reality? *Int. J. Mol.*  
807 *Sci.* **20**, 3563 (2019).
- 808 49. Schafer, M. J. & LeBrasseur, N. K. The influence of GDF11 on brain fate and function.  
809 *GeroScience* **41**, 1–11 (2019).
- 810 50. Rahit, K. M. T. H. & Tarailo-Graovac, M. Genetic Modifiers and Rare Mendelian Disease.  
811 *Genes* **11**, 239 (2020).
- 812 51. Emilsson, V. *et al.* Co-regulatory networks of human serum proteins link genetics to disease.  
813 *Science* **361**, 769–773 (2018).
- 814 52. Weinacht, K. G. *et al.* Ruxolitinib reverses dysregulated T helper cell responses and controls  
815 autoimmunity caused by a novel signal transducer and activator of transcription 1 (STAT1)  
816 gain-of-function mutation. *J. Allergy Clin. Immunol.* **139**, 1629-1640.e2 (2017).
- 817 53. Kaleviste, E. *et al.* Interferon signature in patients with STAT1 gain-of-function mutation is  
818 epigenetically determined. *Eur. J. Immunol.* **49**, 790–800 (2019).
- 819 54. Rodero, M. P. & Crow, Y. J. Type I interferon–mediated monogenic autoinflammation: The  
820 type I interferonopathies, a conceptual overview. *J. Exp. Med.* **213**, 2527–2538 (2016).
- 821 55. Pereira, B., Xu, X.-N. & Akbar, A. N. Targeting Inflammation and Immunosenescence to  
822 Improve Vaccine Responses in the Elderly. *Front. Immunol.* **11**, 2670 (2020).
- 823 56. Carrasco, E. *et al.* The role of T cells in age-related diseases. *Nat. Rev. Immunol.* **22**, 97–111  
824 (2022).
- 825 57. Lucca, L. E. & Dominguez-Villar, M. Modulation of regulatory T cell function and stability by  
826 co-inhibitory receptors. *Nat. Rev. Immunol.* **20**, 680–693 (2020).
- 827 58. Wang, X., Qiu, L., Li, Z., Wang, X.-Y. & Yi, H. Understanding the Multifaceted Role of  
828 Neutrophils in Cancer and Autoimmune Diseases. *Front. Immunol.* **9**, 2456 (2018).
- 829 59. Liu, Y. & Kaplan, M. J. Cardiovascular disease in systemic lupus erythematosus: an update.  
830 *Curr. Opin. Rheumatol.* **30**, 441–448 (2018).
- 831 60. Tseng, C. W. & Liu, G. Y. Expanding roles of neutrophils in aging hosts. *Curr. Opin. Immunol.*  
832 **29**, 43–48 (2014).
- 833 61. Foy, B. H. *et al.* Association of Red Blood Cell Distribution Width With Mortality Risk in  
834 Hospitalized Adults With SARS-CoV-2 Infection. *JAMA Netw. Open* **3**, e2022058 (2020).
- 835 62. Patel, H. H., Patel, H. R. & Higgins, J. M. Modulation of red blood cell population dynamics is  
836 a fundamental homeostatic response to disease. *Am. J. Hematol.* **90**, 422–428 (2015).
- 837 63. Furman, D. *et al.* Chronic inflammation in the etiology of disease across the life span. *Nat.*  
838 *Med.* **25**, 1822–1832 (2019).
- 839 64. Shen-Orr, S. S. Challenges and Promise for the Development of Human Immune Monitoring.  
840 *Rambam Maimonides Med. J.* **3**, e0023 (2012).
- 841 65. Ayres, J. S. The Biology of Physiological Health. *Cell* **181**, 250–269 (2020).
- 842 66. López-Otín, C. & Kroemer, G. Hallmarks of Health. *Cell* **184**, 33–63 (2021).
- 843 67. Collins, F. S. & Varmus, H. A New Initiative on Precision Medicine. *N. Engl. J. Med.* **372**, 793–  
844 795 (2015).
- 845 68. Hood, L. & Friend, S. H. Predictive, personalized, preventive, participatory (P4) cancer  
846 medicine. *Nat. Rev. Clin. Oncol.* **8**, 184–187 (2011).

847 69. Bektas, A., Schurman, S. H., Sen, R. & Ferrucci, L. Human T cell immunosenescence and  
848 inflammation in aging. *J. Leukoc. Biol.* **102**, 977–988 (2017).  
849

**Table 1. Patient Characteristics**

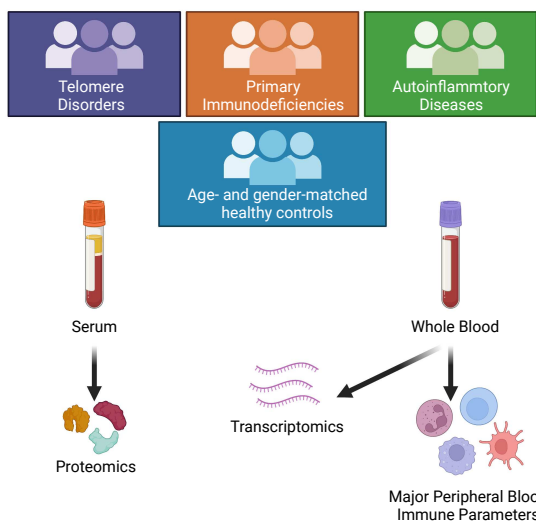
Some patients had multiple samples collected over time at different visits, thus the number of samples can exceed the number of patients indicated.

Condition	Subject Count		Sample Count			Age at Sample Drawn	Sex	Race					
	Primary	Set Aside	Serum Proteomics	CBC + TBNK immune cell phenotyping	Whole Blood Transcriptomics	median [min-max] (Years)	Male	Asian	Black/African American	Hawaiian/Pacifier Islander	Multiple Race	White	Unknown
<b>p47-CGD</b>	18	4	31	33	32	36.8 [14.9-58.3]	12 (54.5%)	-	4 (18.2%)	-	-	17 (77.3%)	1 (4.5%)
<b>X-CGD</b>	23	6	41	51	49	31.3 [7.6-52]	28 (96.6%)	1 (3.4%)	4 (13.8%)	-	1 (3.4%)	22 (75.9%)	1 (3.4%)
<b>CARD14 DN</b>	2	0	2	2	1	13.25 [12.4-14.1]	1 (50%)	-	2 (100%)	-	-	-	-
<b>CTLA4</b>	4	1	7	8	10	31.6 [18.3-57.9]	4 (80%)	-	-	-	-	5 (100%)	-
<b>DADA2</b>	8	2	13	13	13	15.2 [7.4-26.3]	7 (70%)	1 (10%)	-	-	-	8 (80%)	1 (10%)
<b>FCAS</b>	6	1	7	7	6	21.2 [2.7-55.8]	3 (42.9%)	-	-	-	-	4 (57.1%)	3 (42.9%)
<b>FMF</b>	10	2	12	12	13	53.6 [14.2-77.6]	7 (58.3%)	-	-	-	-	12 (100%)	-
<b>GATA2</b>	14	4	19	21	17	41.9 [16.4-81.8]	4 (22.2%)	-	-	-	1 (5.6%)	15 (83.3%)	2 (11.1%)
<b>HIDS</b>	4	1	6	6	7	19.4 [10.4-20.4]	2 (40%)	-	-	-	-	5 (100%)	-
<b>IL-12R</b>	2	1	3	4	4	21.4 [6.5-43.5]	1 (33.3%)	-	-	-	1 (33.3%)	2 (66.7%)	-
<b>LAD1</b>	2	0	3	4	5	30.5 [30.3-38.4]	2 (100%)	-	-	-	-	2 (100%)	-
<b>Muckle-Wells</b>	3	1	5	5	5	36.5 [7.9-43.8]	2 (50%)	-	-	-	1 (25%)	3 (75%)	-
<b>NEMO</b>	2	1	6	6	7	29.9 [8.9-39.2]	3 (100%)	-	-	-	-	3 (100%)	-
<b>NEMO carrier</b>	2	0	2	2	2	24.1 [15.3-32.9]	0 (0%)	-	-	-	-	2 (100%)	-
<b>PAPA Syndrome</b>	6	2	14	14	11	29.3 [17.5-60.1]	5 (62.5%)	-	-	1 (12.5%)	1 (12.5%)	6 (75%)	-
<b>PGM3</b>	6	1	9	11	10	15.5 [3.9-38.7]	6 (85.7%)	-	-	-	-	7 (100%)	-
<b>PI3K</b>	9	2	13	17	15	14.75 [9.4-25.9]	3 (27.3%)	1 (9.1%)	2 (18.2%)	-	-	8 (72.7%)	-
<b>STAT1 GOF</b>	15	4	31	34	32	29 [16.7-71.1]	5 (26.3%)	-	1 (5.3%)	-	-	18 (94.7%)	-
<b>STAT3 DN</b>	32	8	39	50	44	25.7 [6.2-59.9]	21 (52.5%)	1 (2.5%)	5 (12.5%)	-	-	30 (75%)	4 (10%)
<b>TERC</b>	2	0	2	2	2	36.65 [29.3-44]	1 (50%)	-	-	-	-	2 (100%)	-
<b>TERT</b>	3	1	5	5	3	53.3 [28.5-59.3]	3 (75%)	-	-	-	-	4 (100%)	-
<b>TRAPS</b>	10	3	14	14	13	30.7 [12-67.9]	6 (46.2%)	-	-	-	-	12 (92.3%)	1 (7.7%)
<b>Healthy</b>	34	8	42	43	44	33.2 [6.1-67.8]	20 (47.6%)	3 (7.1%)	8 (19%)	-	2 (4.8%)	28 (66.7%)	1 (2.4%)

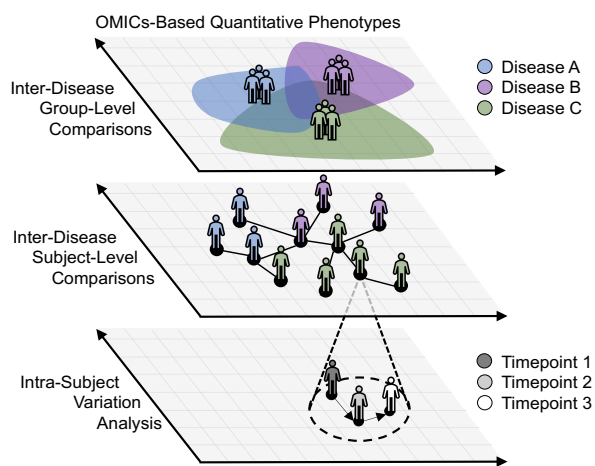


**Figure 1**

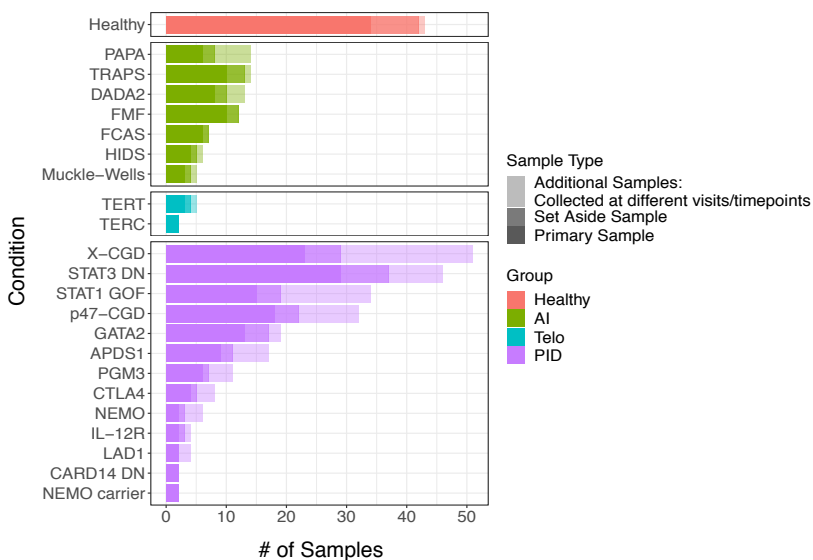
**a**



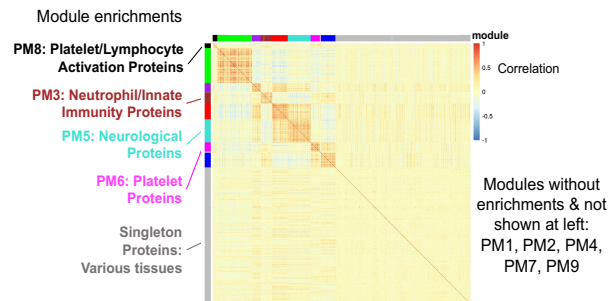
**b**



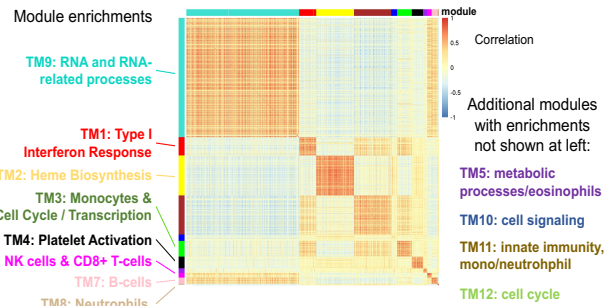
**c**



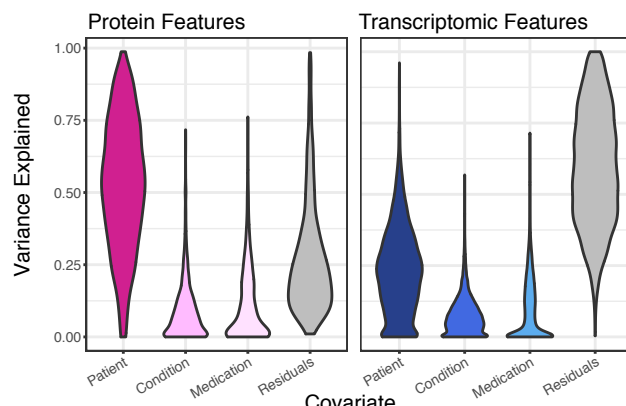
**d**



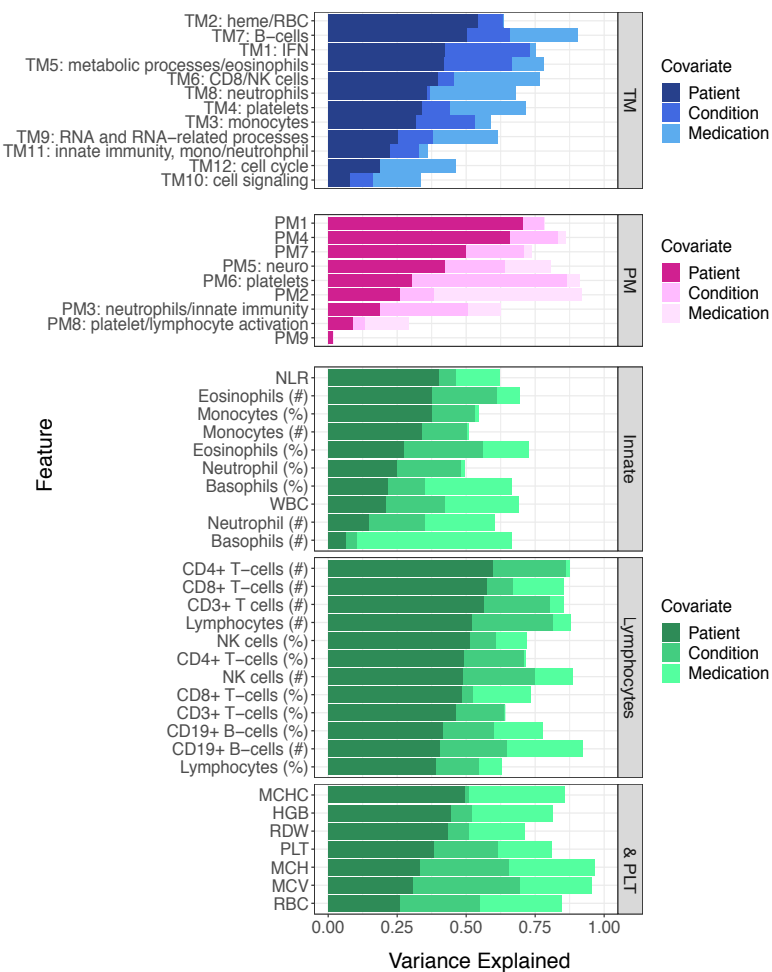
**e**



**f**



**g**



850 **Figure 1. Study and data overview.**

851 **a**, Patient groups and data collected. Individual disease groups are shown in (c).

852 **b**, Conceptual overview of the study and analysis approaches. Both disease group centric (top-  
853 down, disease label based) and individual subject based (bottom-up, unbiasedly starting from  
854 subject-subject similarities) analyses are pursued.

855 **c**, Breakdown of cohort by disease and sample type. Data are broken down into the number of  
856 “primary” samples (equal to the number of subjects analyzed in this study), subjects reserved  
857 (“set aside”) up front immediately after data generation and before any data analyses for  
858 potential independent follow-up analyses (see Methods), and samples from the primary  
859 subjects (“repeat”) but collected at additional timepoints. AI = autoimmune diseases. Telo  
860 = telomere disorders. PID = primary immunodeficiencies.

861 **d**, Gene-gene correlation heatmap of whole blood transcriptomic data. Modules of correlated  
862 genes [or “transcriptional modules” (TMs);  $k = 12$ ] are annotated by color at the top and left.  
863 Modules were created using all transcriptional features; however, only the temporally stable  
864 genes are shown in the heatmap (see (f) and (g) below). Only modules with significant  
865 enrichments are labeled/annotated.

866 **e**, Similar to (d) but for serum protein data. Modules of correlated proteins (PMs;  $k = 10$ ) are  
867 annotated by color at the top and left. The serum protein data contains a large, weakly  
868 correlated set of proteins (grey module). Modules were created using all features; however,  
869 only the temporally stable proteins are shown in the heatmap [see (f) and (g) below]. Only  
870 modules with significant enrichments are labeled/annotated.

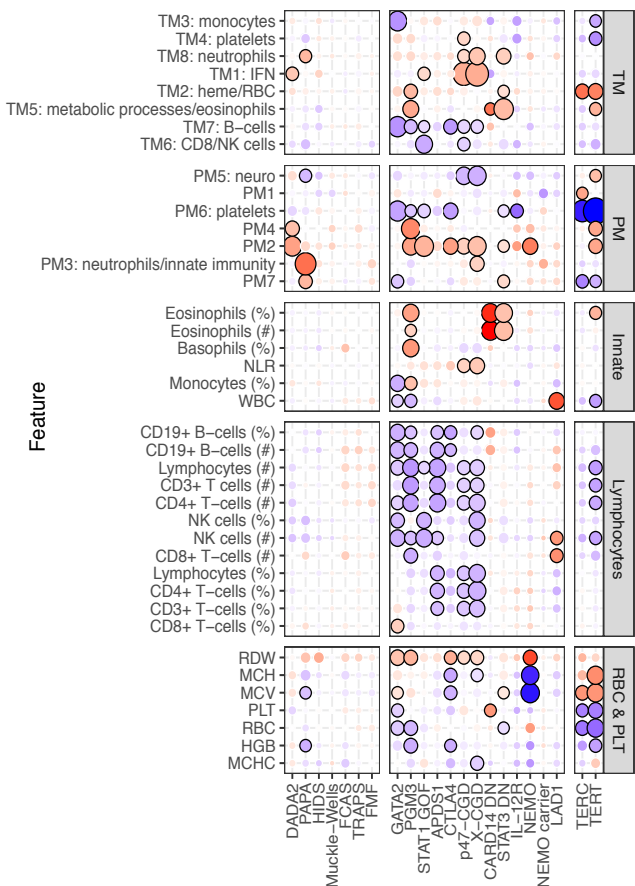
871 **f**, Violin plots showing the distribution, across all measured proteins (1,305) and transcripts  
872 (15,729), of the percent of variance assigned to each variable in the variance partition analysis.  
873 The transcriptomic data had 276 samples with 62 subjects with repeated sampling. The serum  
874 protein data consisted of 271 samples with 64 subjects with repeated sampling.

875 **g**, Barplots of the percent of variance assigned to each variable in the variance partition  
876 analysis, run across each transcriptomic module (blue), serum protein module (magenta), and  
877 CBC parameter (green). This analysis used subjects with repeat samples collected at different  
878 timepoints. The CBC/TBNK data consisted of 271 samples with 63 subjects with repeated  
879 sampling. TM = whole blood transcriptomic modules. PM = serum protein modules. IFN =  
880 interferon. NLR = neutrophil-to-lymphocyte ratio. WBC = white blood cell count. MCHC = mean  
881 corpuscular hemoglobin concentration. HGB = hemoglobin. RDW = red cell distribution width.  
882 PLT = platelet count. MCH = mean corpuscular hemoglobin. MCV = mean corpuscular volume.  
883 RBC = red blood cell count. NK = natural killer.

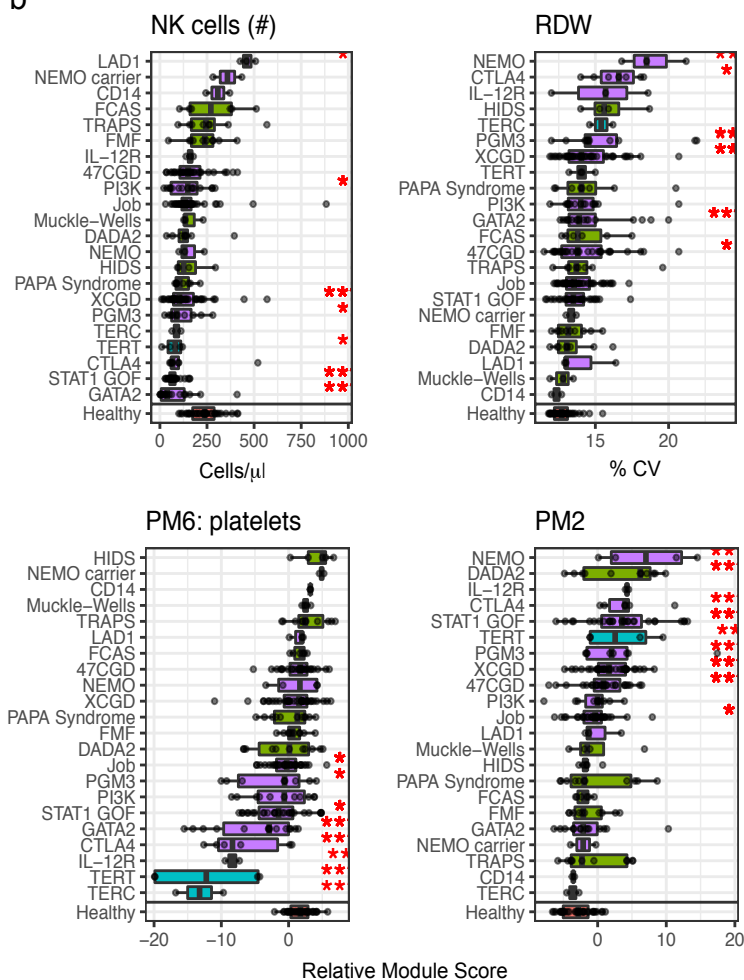
884

**Figure 2**

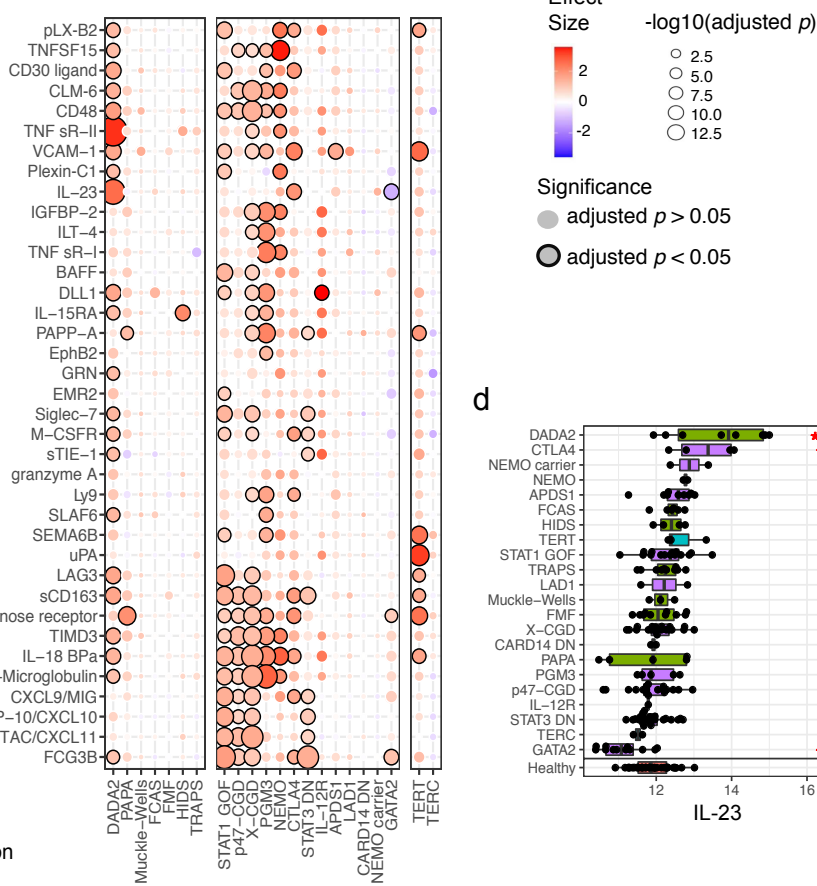
**a**



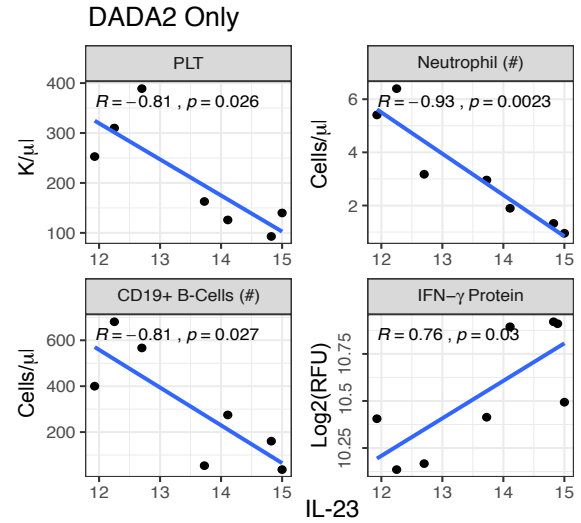
**b**



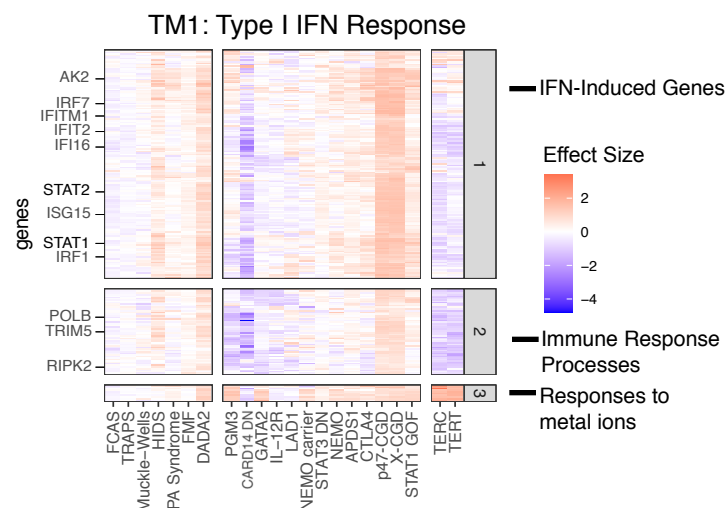
**c**



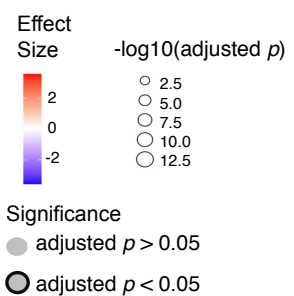
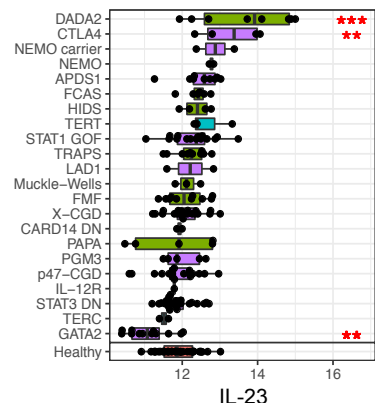
**e**



**f**



**d**



885 **Figure 2. Molecular and cellular signatures of individual monogenic diseases.**

886 **a**, A bubble plot of temporally stable (>50% variance explained by subject) complete blood  
887 count (CBC) and lymphocyte (T, B, NK cell) phenotyping (TBNK) parameters, and serum protein  
888 and transcriptomic module scores (rows) vs. the disease groups (columns). Columns and rows  
889 are ordered by hierarchical clustering (columns/diseases were clustered within major groups,  
890 i.e. primary Immunodeficiencies, autoinflammatory diseases, and telomere disorders). The  
891 bubble color corresponds to the effect size (estimated difference between patients in the  
892 disease group vs. matching healthy subjects via a linear model) for each group while controlling  
893 for age, gender, and whether the patient was acutely ill during sampling. The size of the bubble  
894 reflects the adjusted  $p$  value associated with the fitted t-statistic and the presence of black  
895 outlines around the bubble denotes an adjusted  $p$  value < 0.05. Red boxes highlight specific  
896 parameters discussed in the text. TM = whole blood transcriptomic modules. PM = serum  
897 protein modules. IFN = interferon. NLR = neutrophil-to-lymphocyte ratio. WBC = white blood  
898 cell count. MCHC = mean corpuscular hemoglobin concentration. HGB = hemoglobin. RDW =  
899 red cell distribution width. PLT = platelet count. MCH = mean corpuscular hemoglobin. MCV =  
900 mean corpuscular volume. RBC = red blood cell count. NK = natural killer.

901 **b**, Boxplots of NK cell count, RDW, and module scores of PM2, and PM6 (enriched for platelet-  
902 related factors) across all disease and healthy groups in the study. The healthy subject group is  
903 shown separately at the bottom.  $P$  values computed from linear models used in (a). \*adjusted  $p$   
904 value < 0.05, \*\*adjusted  $p$  value < 0.01, \*\*\*adjusted  $p$  value < 0.001. Box plot center lines  
905 correspond to the median value; lower and upper hinges correspond to the first and third  
906 quartiles (the 25th and 75th percentiles), and lower and upper whiskers extend from the box to  
907 the smallest or largest value correspondingly, but no further than 1.5X inter-quantile range.

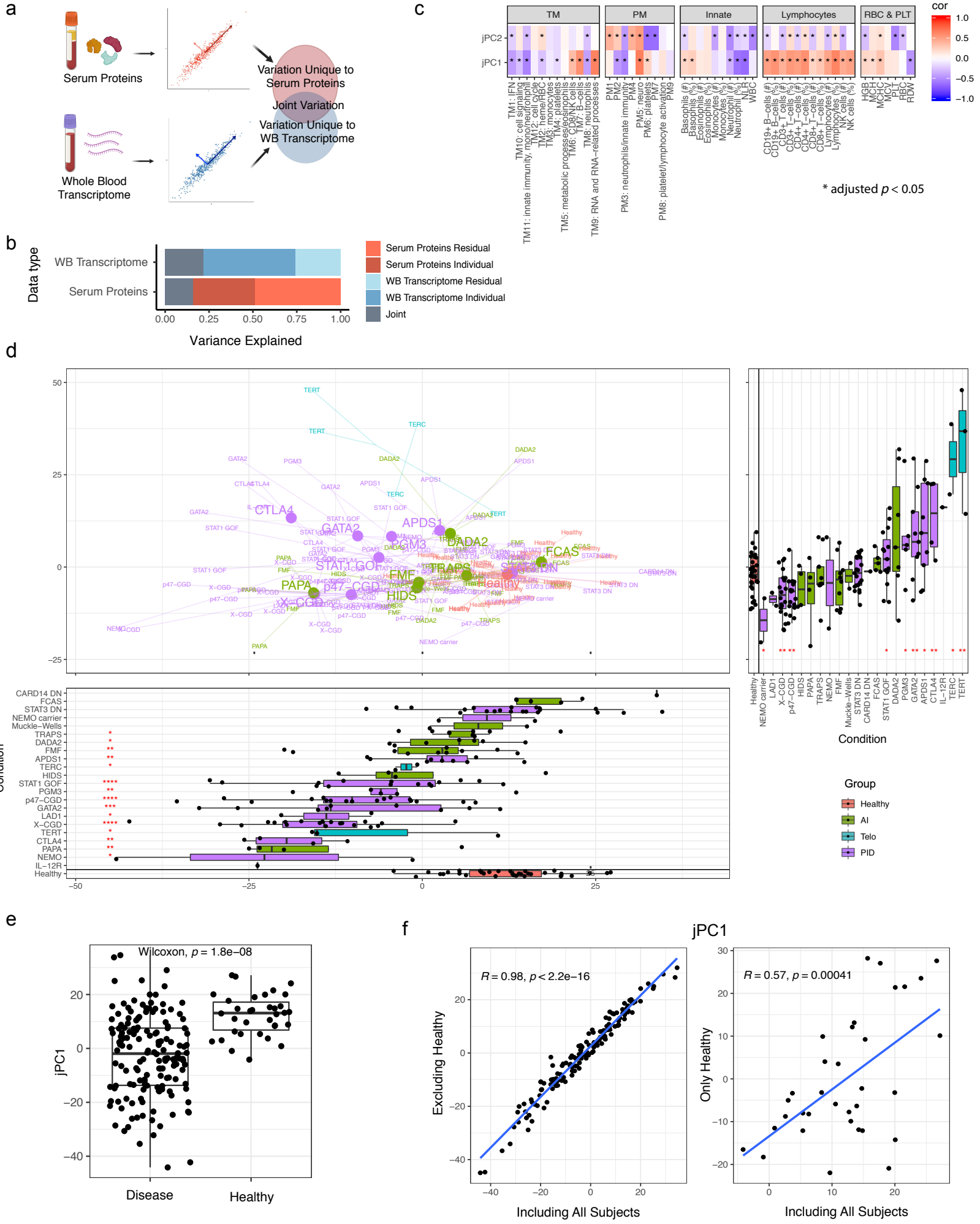
908 **c**, Similar to (a) but limited to the PM2 member proteins (rows). The red box highlights IL-23,  
909 the distribution of which is shown in boxplot in (d).

910 **d**, Similar to (b) but for IL-23 relative serum protein level (as measured by the Somalogic  
911 platform) across all disease conditions and healthy subjects in the study.

912 **e**, Scatterplots showing the correlation between the relative serum protein level of IL-23 (as  
913 measured by the Somalogic platform) and the indicated peripheral blood cell  
914 frequencies/counts and the IFN- $\gamma$  relative serum protein level (lower right plot) for DADA2  
915 patients in the study. Pearson correlation coefficient and associated  $p$  value shown.

916 **f**, Heatmap of effect sizes from linear models of individual transcripts (rows) from TM1  
917 (enriched for interferon-stimulated genes) transcriptomic module. All transcripts in the module  
918 are shown without filtering based on significance. The cell color corresponds to the effect size  
919 (estimated log fold-change relative to healthy subjects) for each disease group (columns) while  
920 controlling for age, sex, and whether the patient was acutely ill during sampling. The genes are  
921 clustered into three groups as indicated on the right. Example gene names are highlighted on  
922 the left. IFN = interferon.

**Figure 3** Find axes of variation that are shared between and unique to each datatype



923 **Figure 3. Bottom-up integration of transcriptomic and serum protein personal immune**  
924 **profiles reveals an emergent axis of immune health.**

925 **a**, Conceptual overview of JIVE analysis integrating whole blood transcriptome and serum  
926 protein data. JIVE was performed using the subject-level data (n=188 subjects who had both  
927 serum protein and whole blood transcriptomic data).

928 **b**, Variation explained by the joint (grey – shared by both data types), individual data type  
929 (darker blue and red for transcriptome and protein data, respectively), and residual latent  
930 factors (lighter blue and red for transcriptome and protein data, respectively) in JIVE analysis.

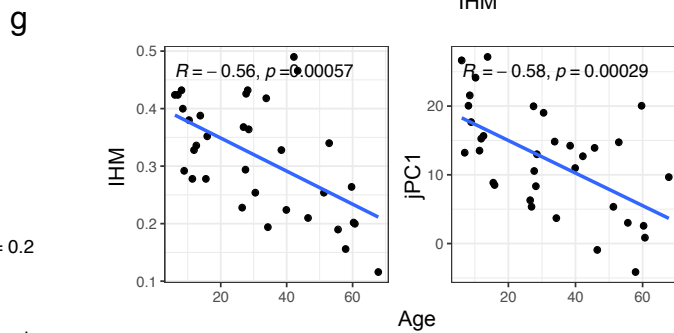
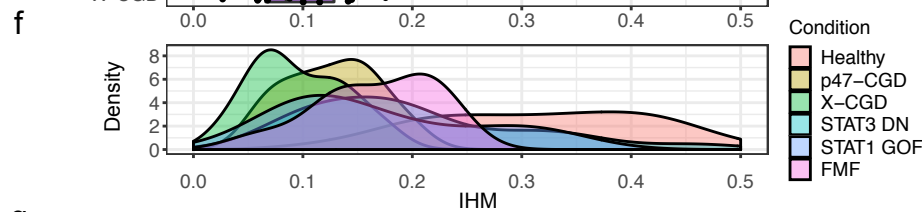
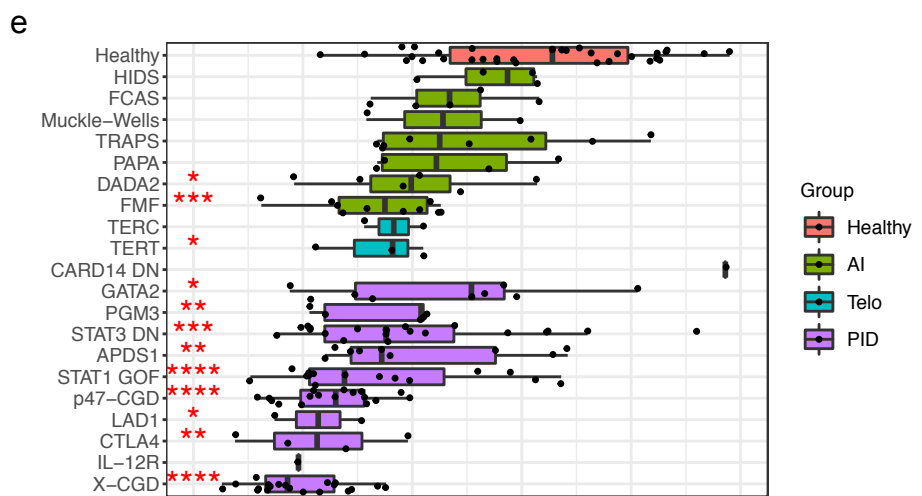
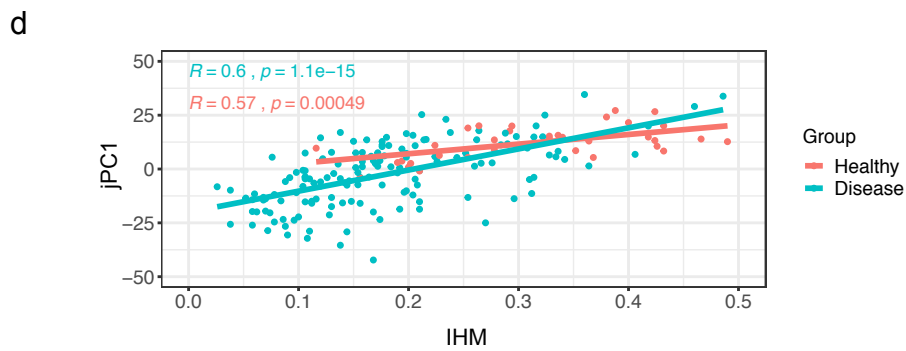
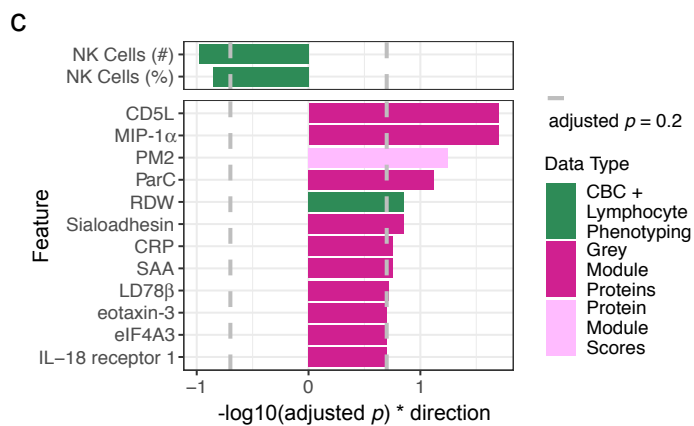
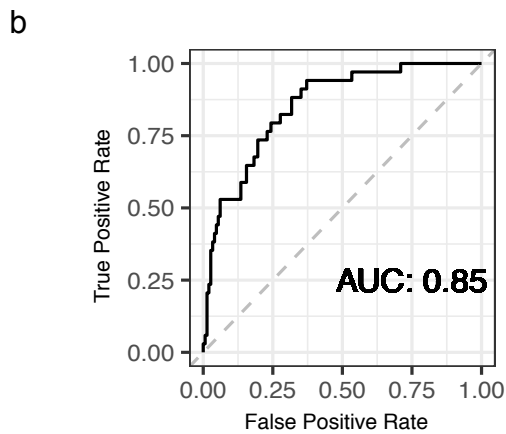
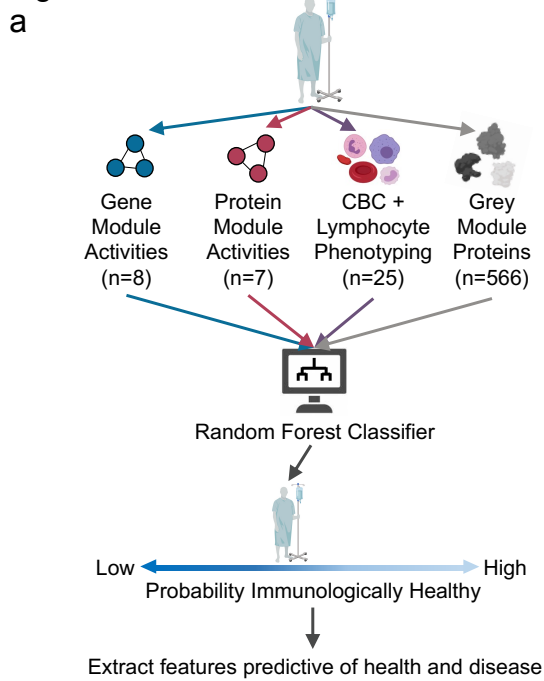
931 **c**, Heatmaps showing Pearson correlation between jPCs (rows) and major peripheral immune  
932 parameters and module scores (columns). Red denotes positive correlation and blue denotes  
933 negative correlation (\*adjusted  $p$  value < 0.05, FDR adjustment performed across all  
934 comparisons together). Correlation was computed using the subject-level data (n = 182 subjects  
935 who had serum protein, whole blood transcriptomic, and CBC/TBNK data). TM = whole blood  
936 transcriptomic modules. PM = serum protein modules. IFN = interferon. NLR = neutrophil-to-  
937 lymphocyte ratio. WBC = white blood cell count. MCHC = mean corpuscular hemoglobin  
938 concentration. HGB = hemoglobin. RDW = red cell distribution width. PLT = platelet count. MCH  
939 = mean corpuscular hemoglobin. MCV = mean corpuscular volume. RBC = red blood cell count.  
940 NK = natural killer.

941 **d**, Projection of patients and healthy subjects onto the jPC1 vs. jPC2 space. N = 154 and 34  
942 disease and healthy subjects, respectively. Text label shows the disease group to which the  
943 patient belongs. Colors denote disease categories involving larger groups of conditions. Large  
944 dots and text denote the centroid (mean jPC1 and jPC2 values) of the indicated disease group.  
945 Only conditions with greater than three subjects have a centroid shown. Boxplots show  
946 projections onto single PC dimensions with patients grouped by disease condition (jPC1 below  
947 the centroid plot; jPC2 to the right of the centroid plot). Each subject's score is represented as a  
948 single point. The healthy subject group is shown in red. (\*  $p$  < 0.05, \*\*  $p$  < 0.01, \*\*\*  $p$  < 0.001,  $p$   
949 values from two-sided Wilcoxon test). Box plot center lines correspond to the median value;  
950 lower and upper hinges correspond to the first and third quartiles (the 25th and 75th  
951 percentiles), and lower and upper whiskers extend from the box to the smallest or largest value  
952 correspondingly, but no further than 1.5X inter-quantile range. The healthy subject group is  
953 shown in red. (\* $p$  < 0.05, \*\* $p$  < 0.01, \*\*\* $p$  < 0.001,  $p$  values from two-sided Wilcoxon test). AI =  
954 autoinflammatory diseases. Telo = telomere disorders. PID = primary immunodeficiencies.

955 **e**, Boxplot of jPC1 scores comparing patients (all disease conditions combined) with healthy  
956 subjects [ $p$  value computed using two-sided Wilcoxon test; same set of subjects in panel (d)].  
957 Box plot center lines correspond to the median value; lower and upper hinges correspond to  
958 the first and third quartiles (the 25th and 75th percentiles), and lower and upper whiskers  
959 extend from the box to the smallest or largest value correspondingly, but no further than 1.5X  
960 inter-quantile range.

961 **f**, Scatterplot of JIVE PCs derived using all subjects vs. JIVE PCs derived using patients only by  
962 removing healthy subjects (left) or only healthy subjects alone (right). Spearman correlation  
963 and associated  $p$  value shown [ $n = 154$  and  $34$  patients and healthy subjects, respectively; same  
964 as in panels **(d)** and **(e)**].

Figure 4





965 **Figure 4. Top-down supervised machine learning classification analysis independently reveals**  
966 **an immune health metric highly concordant with that from unsupervised analysis.**

967 **a**, Conceptual overview of the supervised machine learning analysis of healthy vs. disease  
968 patients using Random Forest classifiers to obtain a probability score of immunological health  
969 [the Immune Health Metric (IHM)]. The number of temporally stable features used from each  
970 data modality is shown. Models were trained using the subject-level data ( $n = 182$  subjects with  
971 serum protein, whole blood transcriptomic, and CBC/TBNK data).

972 **b**, Receiver Operating Characteristic (ROC) curve for distinguishing healthy subjects vs. patients  
973 using the approach shown in **(a)**.

974 **c**, Barplot of the  $-\log_{10}$  adjusted  $p$  values for features passing a 0.2 FDR significance cutoff (grey  
975 dashed line;  $p$  values estimated through permutation testing of Global Variable Importance  
976 from the Random Forest classifiers); these are top features contributed to the classifier used to  
977 derive the IHM. Direction was determined as the sign of the average difference between healthy  
978 subjects and patients from all disease groups.

979 **d**, Scatterplot showing correlation between IHM score and the jPC1 scores across subjects.  
980 Least squares regression lines included for healthy subjects with correlation statistics  
981 shown. 95% confidence interval of the estimated conditional mean is shown.  $N = 148$  and 34  
982 disease patients and healthy subjects, respectively.

983 **e**, Boxplots of IHM scores of individual subjects grouped by condition (disease and healthy  
984 groups). The healthy group (top row) is shown in red; the statistical significance of the  
985 comparison between the condition and the healthy groups is shown for conditions that tested  
986 significant ( $*p < 0.05$ ,  $**p < 0.01$ ,  $***p < 0.001$ ,  $p$  values from two-sided Wilcoxon test). Box  
987 plot center lines correspond to the median value; lower and upper hinges correspond to the  
988 first and third quartiles (the 25th and 75th percentiles), and lower and upper whiskers extend  
989 from the box to the smallest or largest value correspondingly, but no further than 1.5X inter-  
990 quantile range. AI = autoinflammatory diseases. Telo = telomere disorders. PID = primary  
991 immunodeficiencies.

992 **f**, Similar to **(e)**, but here showing smoothed density of IHM scores for each of the groups with  
993 at least 10 subjects.

994 **g**, Scatterplots with trendlines showing the age dependence of the IHM and jPC1 in healthy  
995 individuals only (Spearman correlation and  $p$  values shown;  $n = 34$  healthy subjects with serum  
996 protein, whole blood transcriptomic, and CBC/TBNK data).

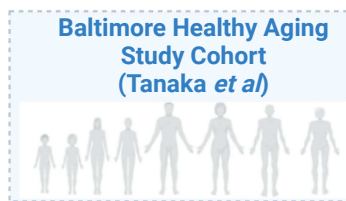
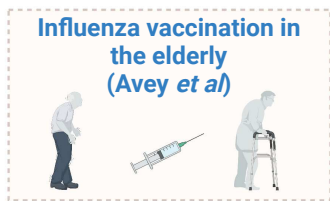
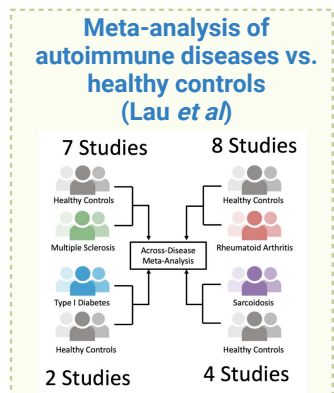
Figure 5

a

Evaluate IHM surrogate signatures in external datasets

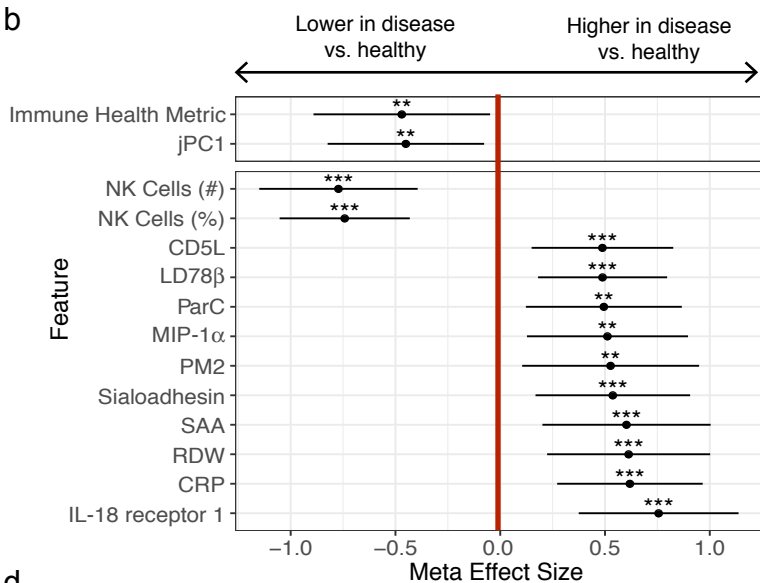
Transcriptomic surrogate of IHM

Proteomic surrogate of IHM

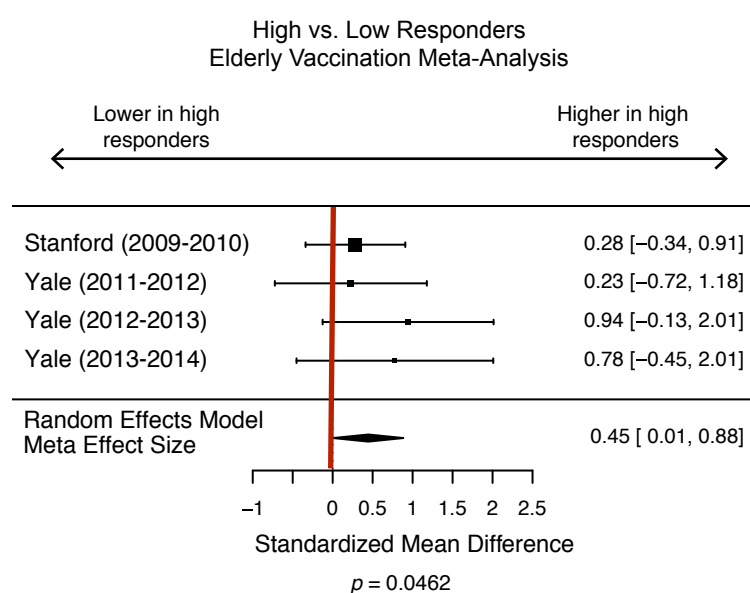


Disease vs. healthy (panel b)

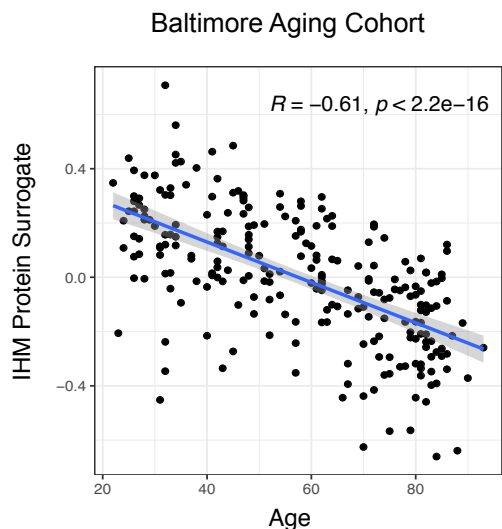
b



c



d



997 **Figure 5. Assessing the IHM in independent datasets**

998 **a**, Graphical depiction of the creation of blood transcriptional and protein surrogate signatures  
999 followed by (from left to right): 1) meta-analysis of four common, non-monogenic  
1000 autoimmune/inflammatory diseases across 21 independent studies, 2) meta-analysis  
1001 comparing high vs. low responders in influenza vaccination in the elderly, and 3) validation of  
1002 the IHM and healthy aging association using an independent cohort.

1003 **b**, Plot of meta effect sizes (average difference between disease and healthy groups) for each  
1004 surrogate gene signature tested using the meta-analysis, including the IHM itself with a  
1005 statistically significant negative effect size (i.e., it is lower in disease than healthy). The point  
1006 shows the estimated effect across all studies used in the meta-analysis and error bars show the  
1007 95% confidence interval ( $1.96 * \text{standard error}$ ) in the meta-analysis.

1008 **c**, Forest plot of effect sizes from the meta-analysis across four independent influenza  
1009 vaccination cohorts of elderly subjects testing whether the IHM transcriptional surrogate  
1010 signature evaluated at baseline before vaccination was associated with antibody titer responses  
1011 to seasonal influenza vaccination in elderly individuals (i.e., whether those with better immune  
1012 health according to the IHM had higher antibody responses.) Effect sizes in each study  
1013 (squares), their 95% confidence interval ( $1.96 * \text{standard error}$ , error bars around square), the  
1014 overall meta effect size (diamond) combining evidence across the four cohorts and the  
1015 standard error of the meta-effect (width of diamond) are shown. Size of square denotes the  
1016 relative number of subjects in that study.

1017 **d**, Scatterplot with trendline showing the negative correlation between chronological age and  
1018 the circulating protein-based IHM surrogate signature scores (see Methods – the circulating  
1019 protein IHM surrogate was developed using data from our cohorts only) in healthy subjects  
1020 from the independent Baltimore Aging Study (Tanaka *et al.*, 2018). N = 240 subjects.

**Figure 6**

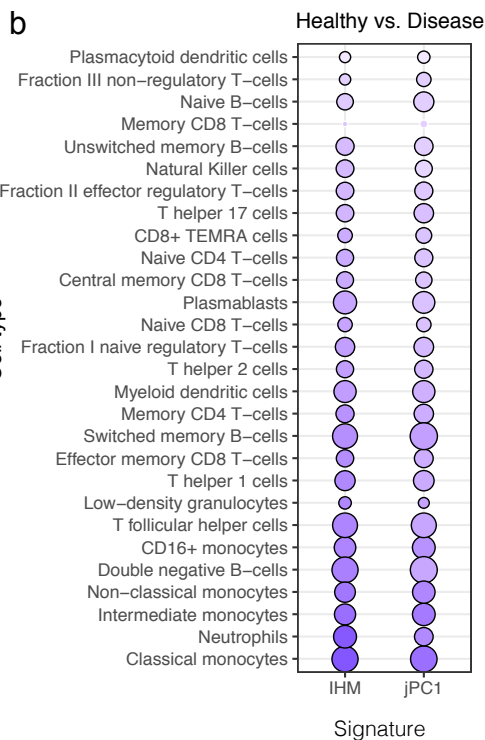
**a**

**Dissect the cellular origins of the IHM**

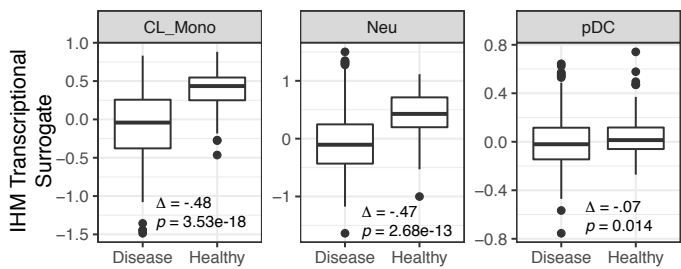
Gene expression from 28 immune cell types in 10 immune-mediated diseases (Ota *et al*)

Compute the IHM transcriptional surrogate signature score within different cell types

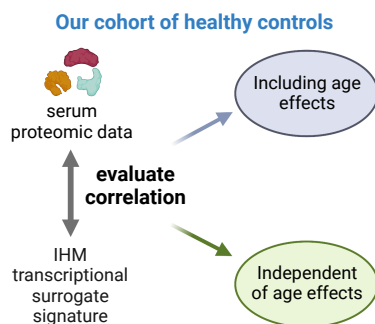
Assess the cell type-specific IHM score  
(1) Healthy vs. disease (panels b and c)  
(2) Correlation with age within healthy (panels d and e)



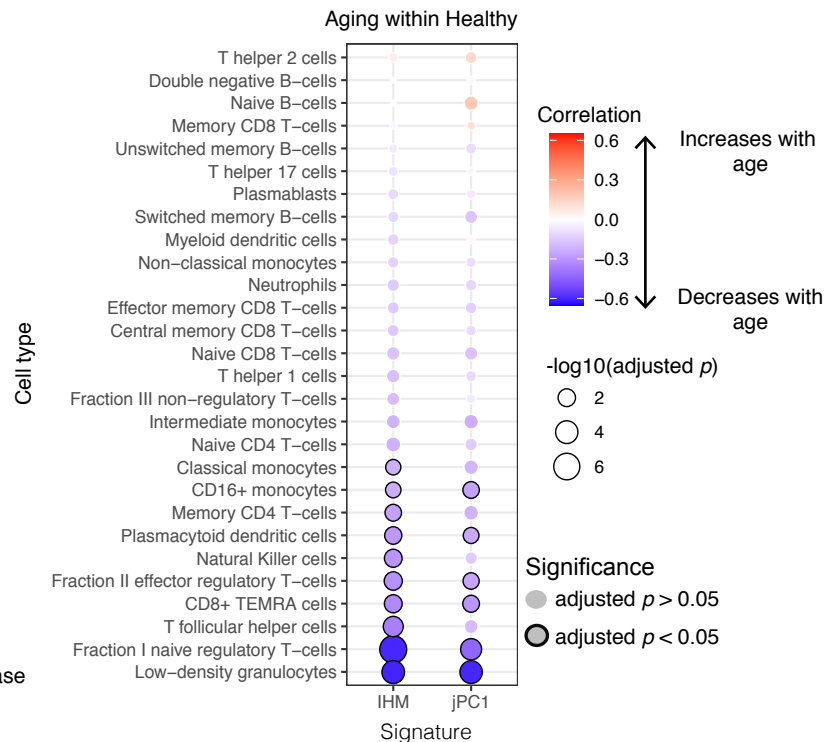
**c**



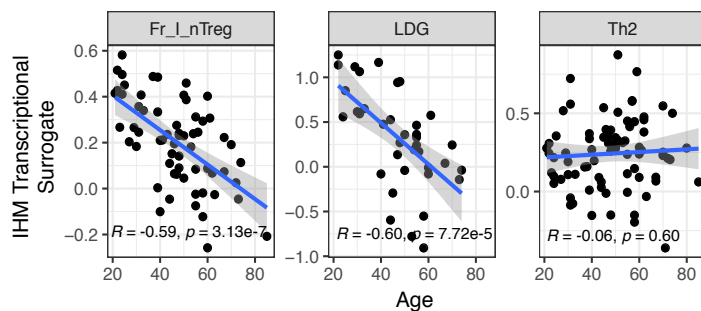
**f**



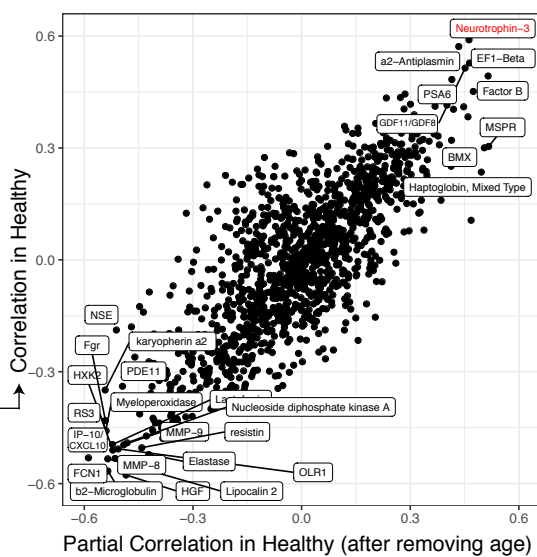
**d**



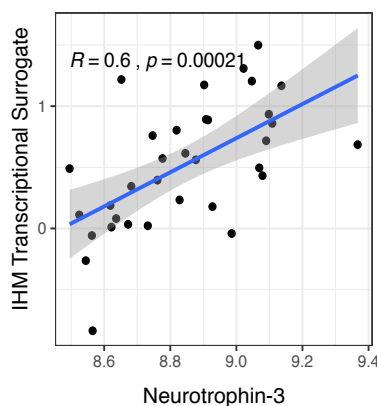
**e**



**g**



**h**



1021 **Figure 6. Cellular origin and circulating protein correlates of the IHM blood transcriptional**  
1022 **surrogate signature**

1023 **a**, Graphical overview of our analysis strategy for assessing 1) the differential expression of the  
1024 IHM's transcriptional surrogates between healthy and autoimmune disease, and 2) association  
1025 with age, in each of 28 cell types from Ota *et al.*

1026 **b**, Bubble plot showing the effect sizes and statistical significance from the comparison of  
1027 autoimmune diseases vs. healthy for the IHM and jPC1 transcriptional signature scores in 28  
1028 cell types from Ota *et al.* Effect sizes are denoted with the color scale shown. Significance is  
1029 denoted by the size of the bubble and the presence of an outline. A negative effect size  
1030 represents a decrease in the signature score in individuals with autoimmune disease relative to  
1031 healthy. CD8+ TEMRA = CD8+ T effector memory CD45RA+ cells.

1032 **c**, Boxplots of IHM transcriptional surrogate signature scores comparing healthy controls vs.  
1033 disease subjects from Ota *et al.* highlighting selected cell types from **(b)** CL\_Mono: classical  
1034 monocytes, Neu: neutrophil, pDC: plasmacytoid dendritic cells. Effect size ( $\Delta$ ) and  $p$  value are  
1035 shown.

1036 **d**, Bubble plot showing Pearson correlation between age and the IHM (and jPC1) transcriptional  
1037 signature scores in healthy individuals only, assessed separately for each one of the 28 cell  
1038 types from Ota *et al.* Correlation strength is denoted by the color scale shown. Significance is  
1039 denoted by the size of the bubble and the presence of an outline. A negative correlation  
1040 represents a decrease in the signature score with older age. A higher signature score is  
1041 associated with higher immune health.

1042 **e**, Scatterplots of IHM transcriptional surrogate signature scores vs. age in healthy controls  
1043 from Ota *et al.* highlighting selected cell types from **(d)** Fr\_I\_nTreg: Fraction I naive regulatory  
1044 T-cells (Ota *et al.*), LDG: low density granulocytes, Th2: T helper cells type 2. Pearson correlation  
1045 and associated  $p$  value are shown.

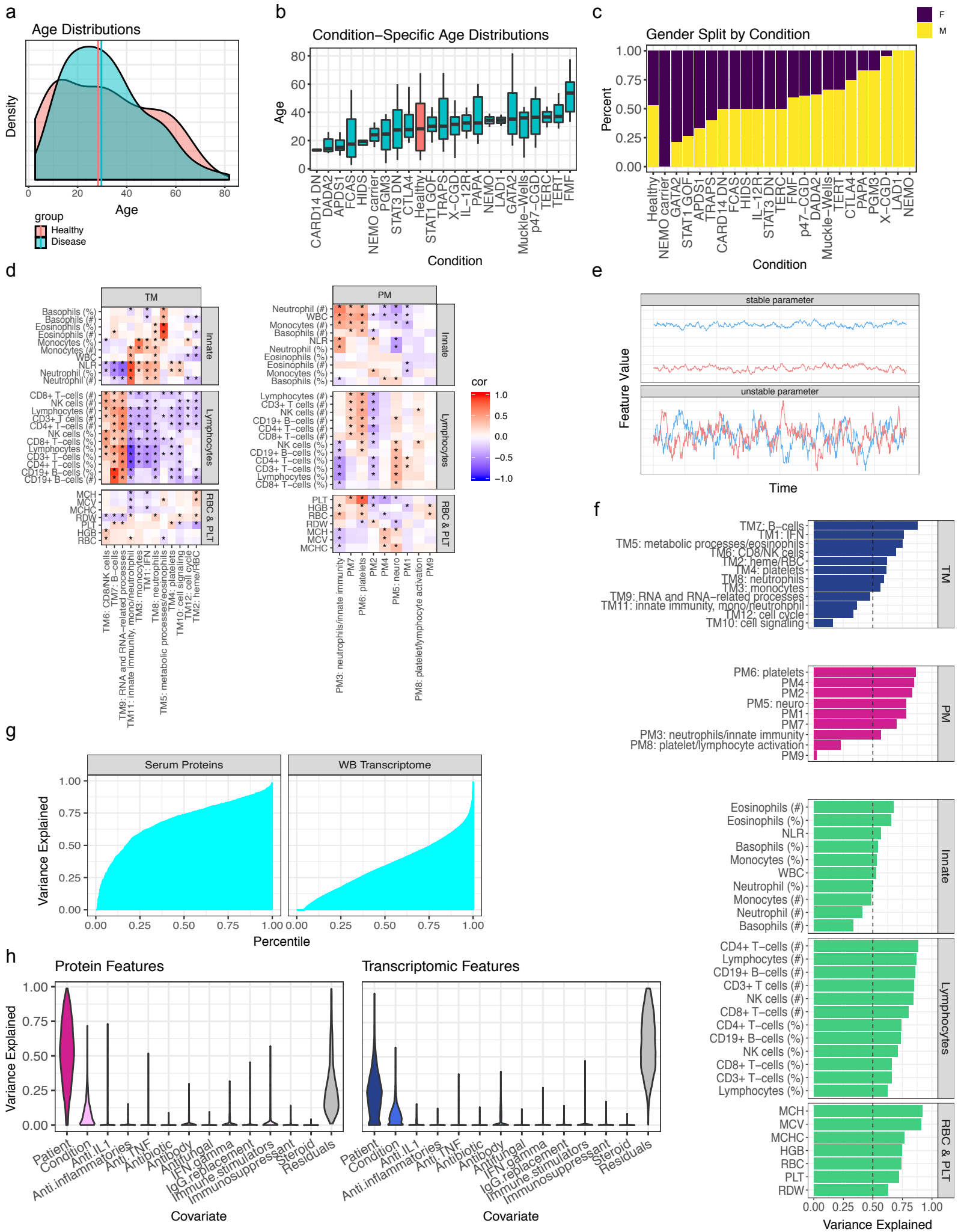
1046 **f**, Graphical overview of the analyses behind the results shown in panel **(g)**. We aim to identify  
1047 circulating proteins that are correlated with the IHM whole blood transcriptional surrogate  
1048 signature in our monogenic patients and assess whether the correlation (and thus the resulting  
1049 protein correlates/surrogates) depends on age (thus without or with age effects removed). The  
1050 age-dependent correlation is simply the correlation between the protein levels and the IHM  
1051 transcriptional surrogate, whereas the age-independent refers to the partial correlation  
1052 between these values after removing the effect of age with a linear regression model.

1053 **g**, Scatterplot showing the Spearman correlation values of serum proteins with the IHM  
1054 transcriptional surrogate signature within healthy individuals only from the monogenic cohort.  
1055 Raw Spearman correlations are shown on the y-axis, and partial correlations after removing the  
1056 effect of age from the protein data and IHM transcriptional signature score are shown on the x-  
1057 axis. The names of the 20 proteins with the highest absolute correlations on the x or y axes are

1058 shown. Neurotrophin-3 is highlighted in red. Correlations were computed with  $n = 34$  healthy  
1059 subjects only.

1060 **h**, Scatterplot of IHM transcriptional surrogate signature score vs. Neurotrophin-3 in healthy  
1061 controls from this study ( $n=34$ ). Spearman correlation and associated  $p$  value shown.

# Extended Data Figure 1



1062 **Extended Data Figure 1. Subject demographics and further characterization of the serum**  
1063 **protein and transcriptomic modules.**

1064 **a**, Density plot of patient and healthy subjects' age distributions (Kolmogorov-Smirnov test  
1065 assessing difference between the two distributions,  $p = 0.41$ ). Extended Data Fig. 1a-c only  
1066 show data for subjects in primary set of subjects; data for set-aside subjects not shown but  
1067 included in Table 1.

1068 **b**, Boxplots of subject ages in each subject group with healthy in red. Box plot center lines  
1069 correspond to the median value; lower and upper hinges correspond to the first and third  
1070 quartiles (the 25th and 75th percentiles), and lower and upper whiskers extend from the box to  
1071 the smallest or largest value correspondingly, but no further than 1.5X inter-quantile range.

1072 **c**, Barplots depicting sex distribution within each group shown as male/female proportions.

1073 **d**, Pearson correlation between the protein (left) or transcriptomic (right) WGCNA modules  
1074 (columns) and cellular [complete blood count (CBC) and lymphocyte (T, B, NK cell) phenotyping  
1075 (TBNK)] parameters (rows). \*adjusted  $p$  value  $< 0.05$ . Computed with 198 subjects with both  
1076 whole blood transcriptome and CBC/TBNK data, and 197 subjects with both serum protein and  
1077 CBC/TBNK data. TM = whole blood transcriptomic modules. PM = serum protein modules. IFN =  
1078 interferon. NLR = neutrophil-to-lymphocyte ratio. WBC = white blood cell count. MCHC = mean  
1079 corpuscular hemoglobin concentration. HGB = hemoglobin. RDW = red cell distribution width.  
1080 PLT = platelet count. MCH = mean corpuscular hemoglobin. MCV = mean corpuscular volume.  
1081 RBC = red blood cell count. NK = natural killer.

1082 **e**, Conceptual illustration of parameter temporal stability, defined by low intra-subject variation  
1083 relative to inter-subject variation.

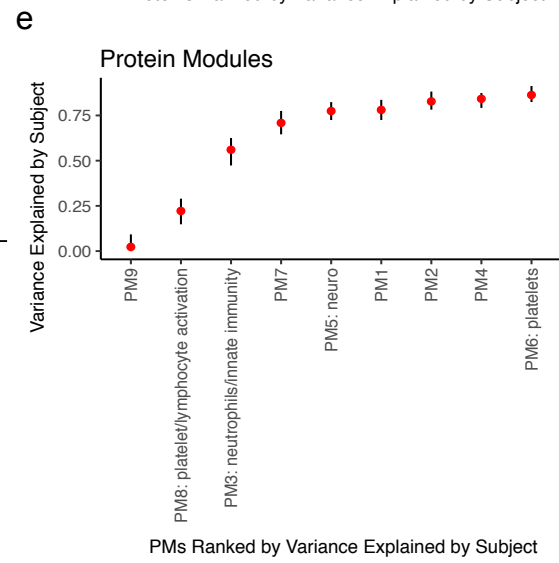
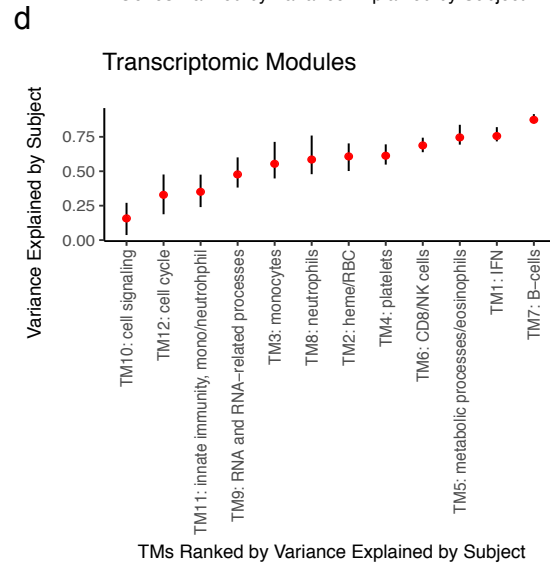
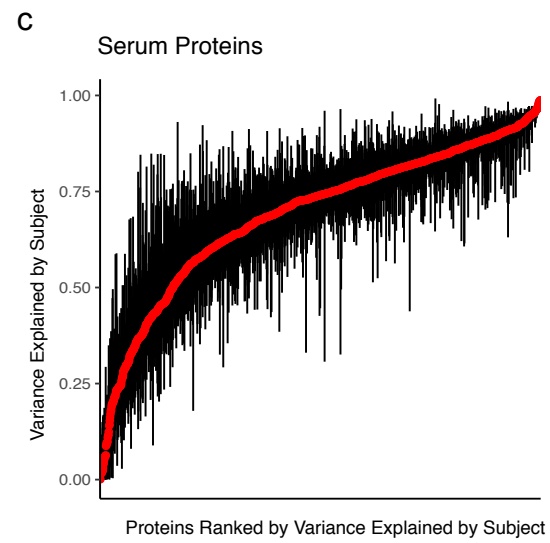
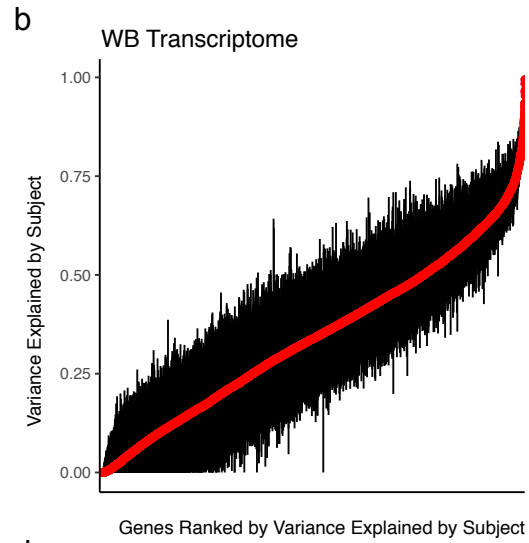
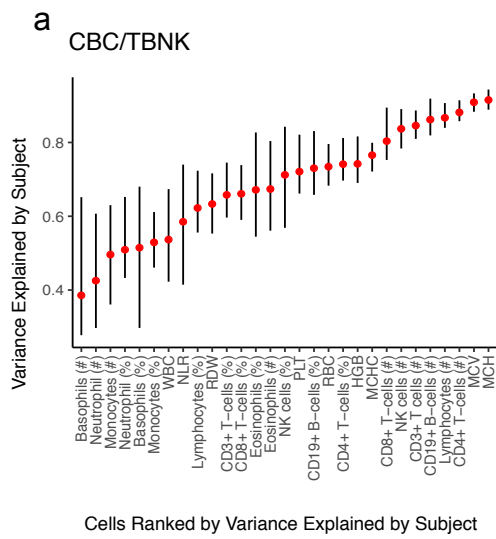
1084 **f**, Barplots of variance assigned to the subject term in the variance partition analysis fit using  
1085 only a subject random intercept (see Methods), run across each CBC parameter, protein  
1086 module, and transcriptomic module. TM = whole blood transcriptomic modules. PM = serum  
1087 protein modules. RBC = red blood cell parameters. PLT = platelets.

1088 **g**, Percent variation explained by the subject term in the variance partition model in the protein  
1089 and transcriptomic features using the variance partition model with only a subject random  
1090 intercept (see Methods) as in (f). Proteins (left) and genes (right) are ordered on the x-axis by  
1091 the percent variation explained by the subject term. WB = whole blood.

1092 **h**, Percent variation explained by the patient and medication covariate (showing effect of each  
1093 medication individually) for each protein (left) and gene (right) measured. Medications were  
1094 included in the model if they were used by many patients and not highly confounded with one  
1095 of the condition groups.



# Extended Data Figure 2



1096 **Extended Data Figure 2. Jackknife resampling shows robustness of variation explained by**  
1097 **subject covariate in mixed effect model**

1098 A jackknife was performed subsampling 80% of subjects with repeat samples and 80% of  
1099 subjects without repeat samples to assess robustness of intra-patient stability estimates for cell  
1100 frequencies (a), gene expression (b), serum protein data (c), gene expression modules (d),  
1101 serum protein modules (e). 100 replicates of subsampling were performed. Points represent  
1102 mean variance explained by subject across all replicates and error bars denote 95% confidence  
1103 intervals (2.5 % and 97.5 % quantiles across jackknife replicates). CBC = complete blood count.  
1104 TBNK = lymphocyte (T, B, NK cell) phenotyping.



1105 **Extended Data Figure 3. Supporting data for the disease-associated molecular and cellular**  
1106 **signatures.**

1107 **a**, Heatmap of complete blood count (CBC) and lymphocyte (T, B, NK cell) phenotyping (TBNK)  
1108 parameters (rows) across patients and healthy subjects (columns); columns and rows are  
1109 ordered by hierarchical clustering. Top annotation row shows the age of the subject, middle  
1110 row shows the large condition groups ( $n > 10$  subjects), and third row shows all condition  
1111 groups regardless of number of subjects.

1112 **b**, Patients and healthy subjects shown in PC1 and PC2 space of CBC and TBNK parameters.  
1113 Each parameter was standardized to unit variance and mean of zero prior to computation of  
1114 the principal components. The text denotes the subject's condition, and the color denotes  
1115 larger condition groups. Large dots and text denote the centroid of that disease group. Only  
1116 conditions with greater than three subjects have a centroid shown. AI = autoinflammatory  
1117 diseases. Telo = telomere disorders. PID = primary immunodeficiencies.

1118 **c**, Table of sample sizes for each data modality-condition group combination. TM: whole blood  
1119 transcriptomic modules; PM: protein modules.

1120 **d**, Similar to Fig. 2a but comparing each condition to all other conditions (healthy subjects are  
1121 removed from the analysis).

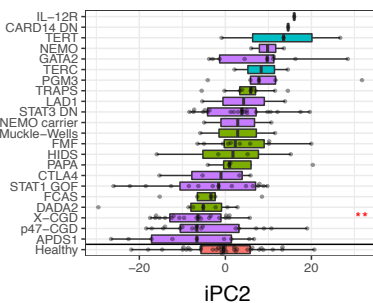
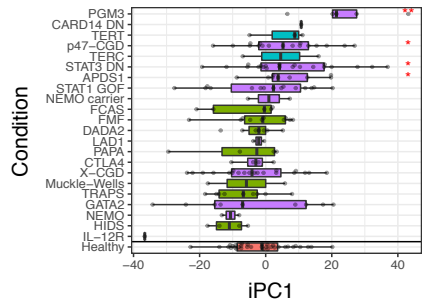
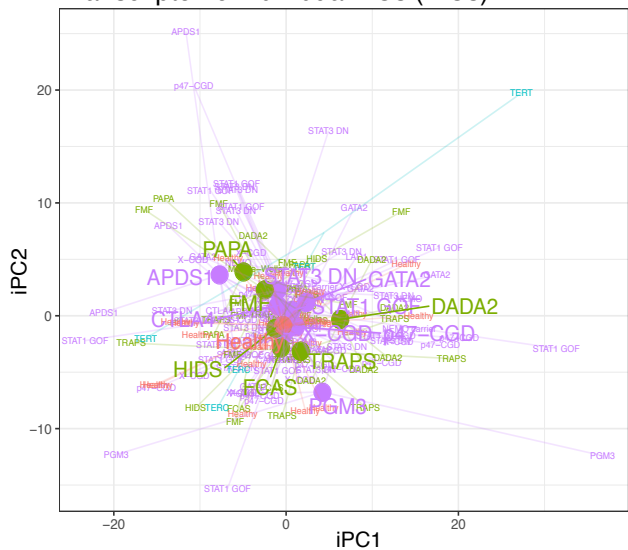
1122 **e**, Barplot of Receiver Operating Characteristic Area Under the Curve (AUC) for conditions-  
1123 versus-all-other-conditions Random Forest classifiers using all features as input. Classifiers were  
1124 trained only for the four condition groups with the most subjects (healthy subjects were  
1125 removed from the analysis); however, subjects from all other disease groups were used as the  
1126 negative samples for each classifier.

1127 **f**, Plot of  $-\log_{10}$  adjusted  $p$  values and global variable importance (GVIs from the Random  
1128 Forest models) of features in the classifiers for the four most represented disease groups. The  
1129 plot is subset to the union of the top five predictive features for each condition.

# Extended Data Figure 4

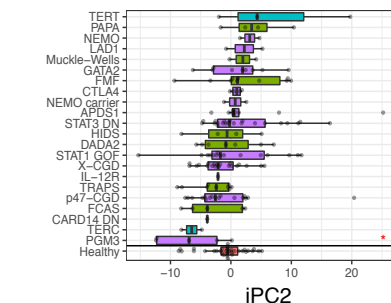
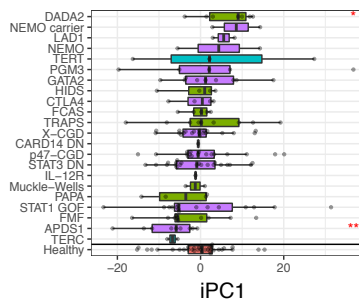
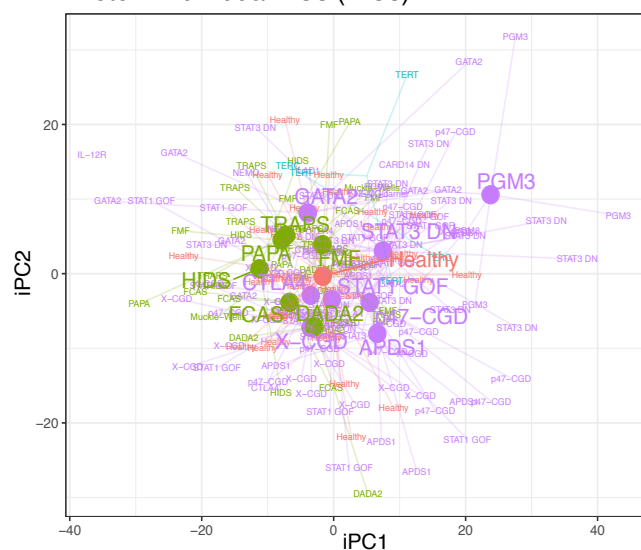
**a**

## Transcriptome Individual PCs (iPCs)



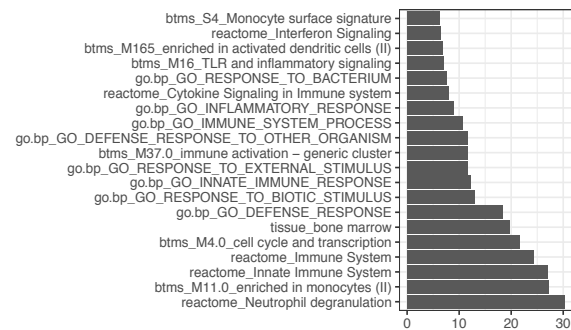
**b**

## Protein Individual PCs (iPCs)

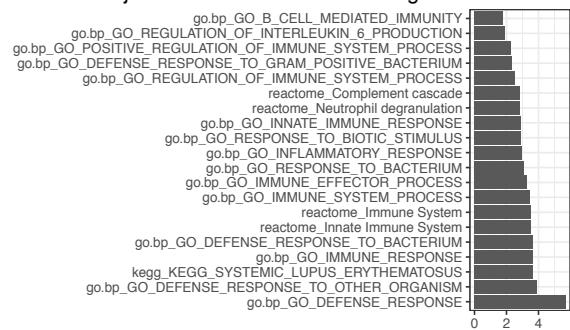


**c**

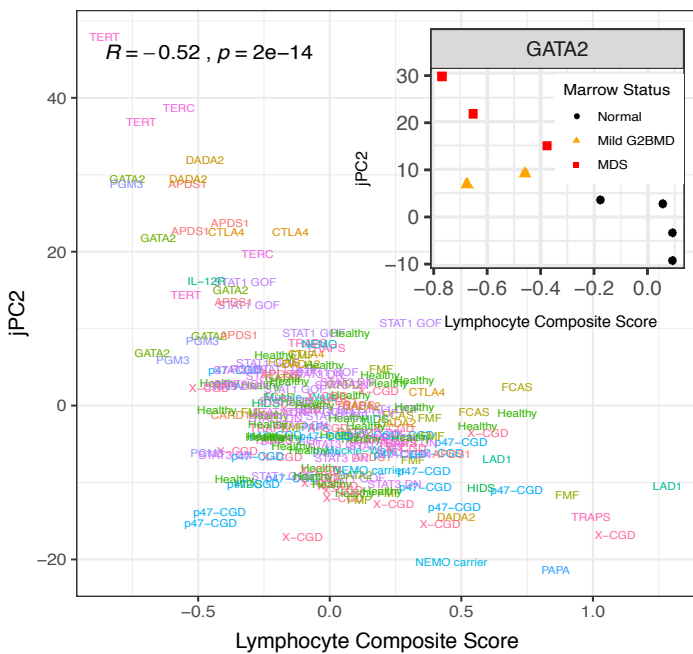
## iPC1 Transcriptomic Enrichments: Negative Correlates



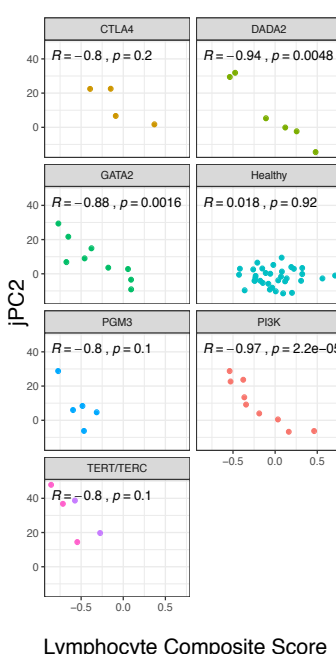
## iPC1 Protein Enrichments: Negative Correlates



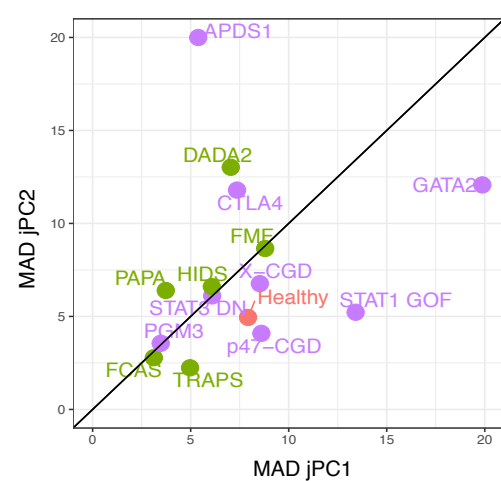
**d**



$-\log_{10}(\text{adjusted } p)$



**e**



1130 **Extended Data Figure 4. Characteristics of the individual and joint PCs from the JIVE analysis.**

1131 **a**, Top panel: patients and healthy subjects shown in transcriptomic individual PC (iPC) 1 vs.  
1132 iPC2 space. Large dots and text denote the centroid of that disease group. Only conditions with  
1133 greater than three subjects have a centroid shown. Bottom panels: boxplots of individual  
1134 transcriptomic iPC1 and iPC2. The rows correspond to the conditions and the color denotes  
1135 larger condition groups. Box plot center lines correspond to the median value; lower and upper  
1136 hinges correspond to the first and third quartiles (the 25th and 75th percentiles), and lower and  
1137 upper whiskers extend from the box to the smallest or largest value correspondingly, but no  
1138 further than 1.5X inter-quantile range. AI = autoinflammatory diseases. Telo = telomere  
1139 disorders. PID = primary immunodeficiencies.

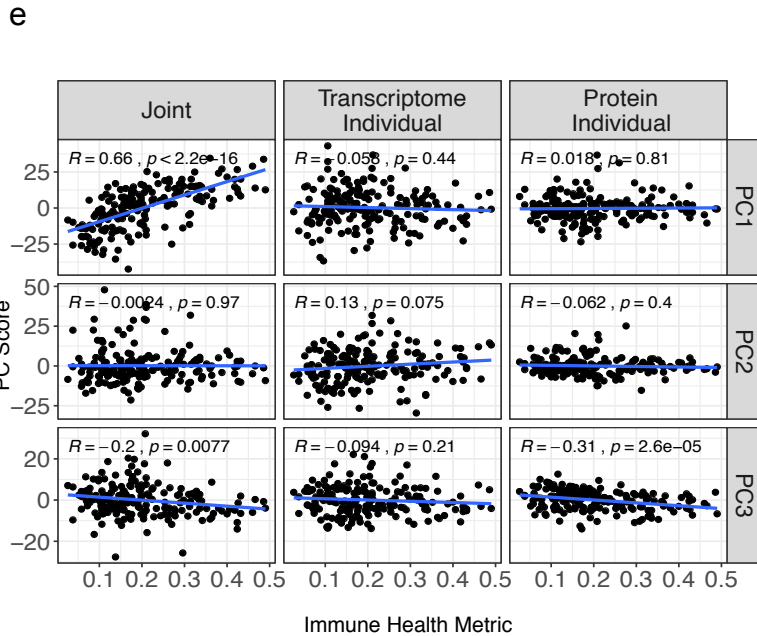
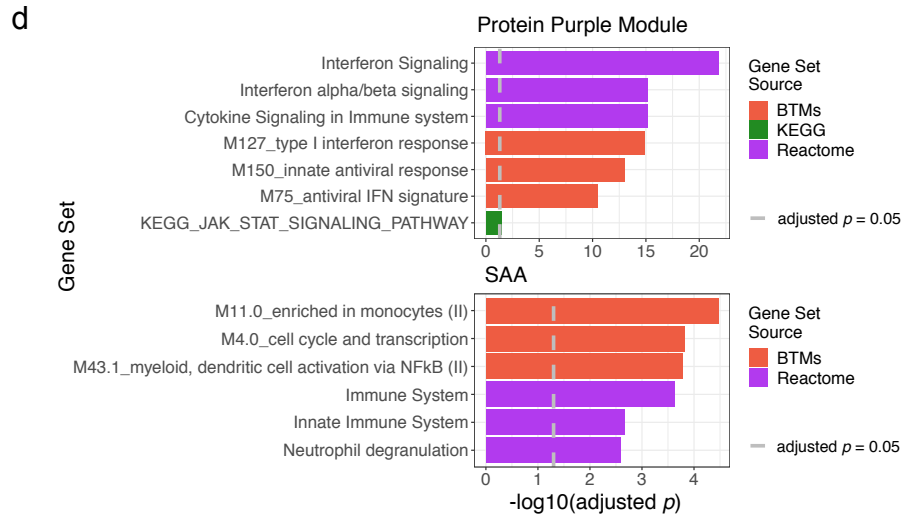
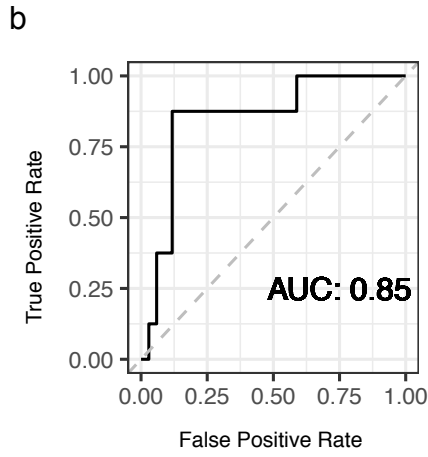
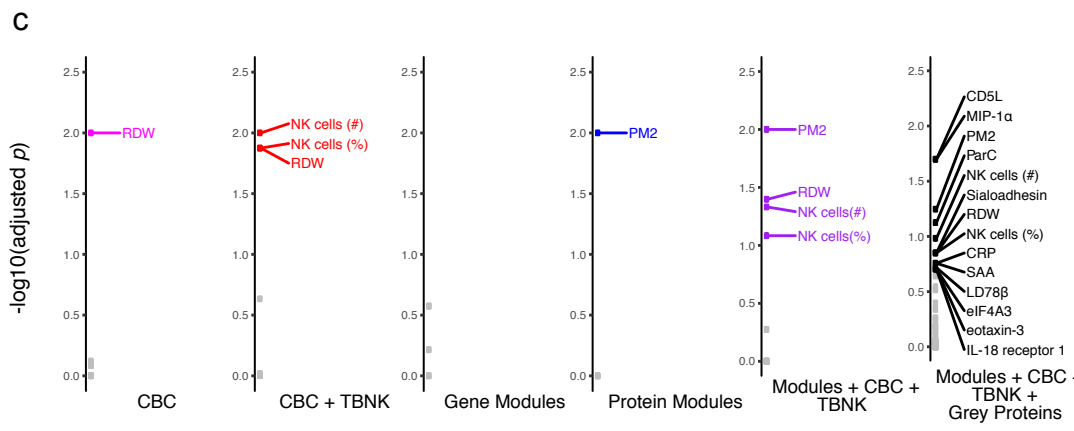
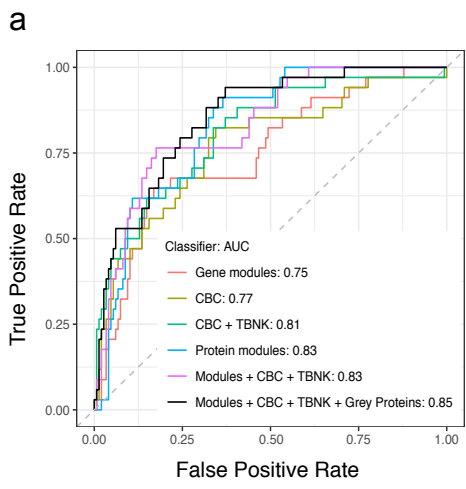
1140 **b**, Similar to (a) but showing the serum protein iPCs.

1141 **c**, Gene set enrichment of transcriptomic (left) and serum protein (right) features negatively  
1142 correlated with jPC1 (enrichment calculated using CameraPR; genes/proteins ranked by the  
1143 Spearman correlation with the JIVE PCs). Gene sets from KEGG pathways, GO biological process  
1144 gene sets, Reactome pathways, and the blood transcriptomic modules and Human Protein Atlas  
1145 tissue gene sets.

1146 **d**, Scatterplot of a hematopoietic composite score (see Methods) vs. jPC2. Left panel displays  
1147 the trend across all patients including healthy subjects and the right set of panels focus on  
1148 individual disease groups whose clinical presentation may include marrow failure or  
1149 lymphopenia. Inset focuses on GATA2 patients, highlighting those with abnormal bone marrow  
1150 biopsies. Spearman correlation and associated  $p$  values are shown. G2BMD = GATA2 deficiency-  
1151 associated bone marrow disorder. MDS = myelodysplastic syndrome.

1152 **e**, Scatterplot of Median Absolute Deviation (MAD) of jPC1 and jPC2 scores for each condition in  
1153 the study. A higher MAD corresponds to greater variation within a disease for that jPC.

# Extended Data Figure 5



1154 **Extended Data Figure 5. Supporting data for the development and characterization of the**  
1155 **Immune Health Metric (IHM).**

1156 **a**, Receiver Operating Characteristic (ROC) curves for Random Forest classifiers from LOOCV  
1157 (leave-one-out-cross-validation) using temporally stable features of individual or the indicated  
1158 combinations of data modalities. CBC = complete blood count. TBNK = lymphocyte (T, B, NK  
1159 cell) phenotyping.

1160 **b**, ROC curve for the Random Forest classifier (the one trained on all data modalities in the  
1161 primary dataset) applied to the set of unseen, independent set-aside patients and healthy  
1162 subjects.

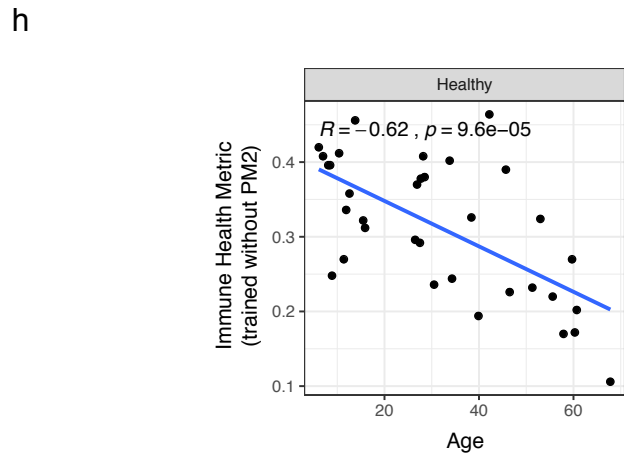
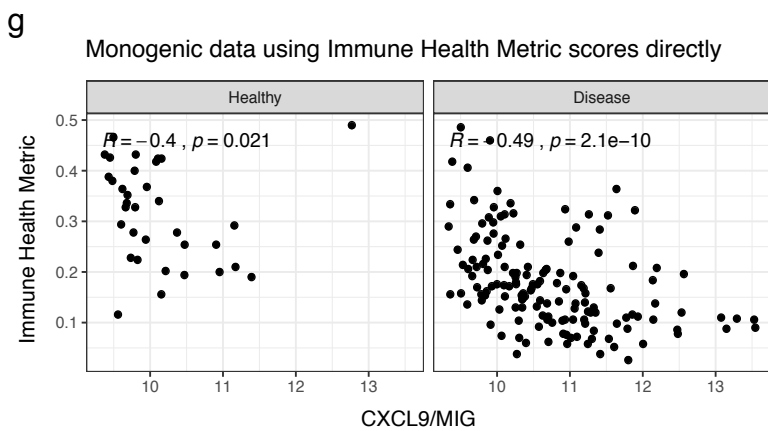
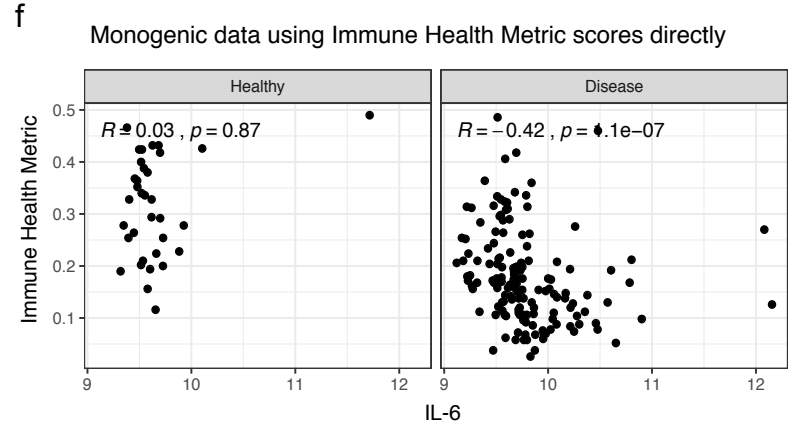
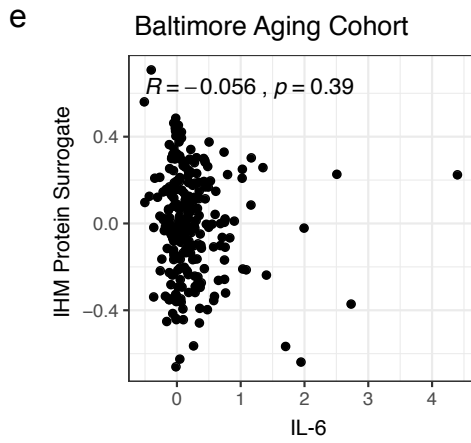
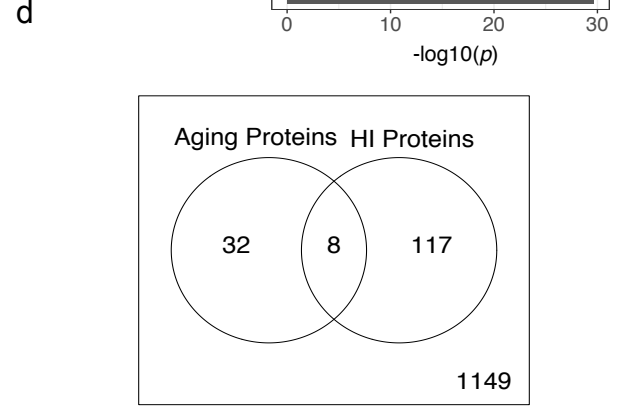
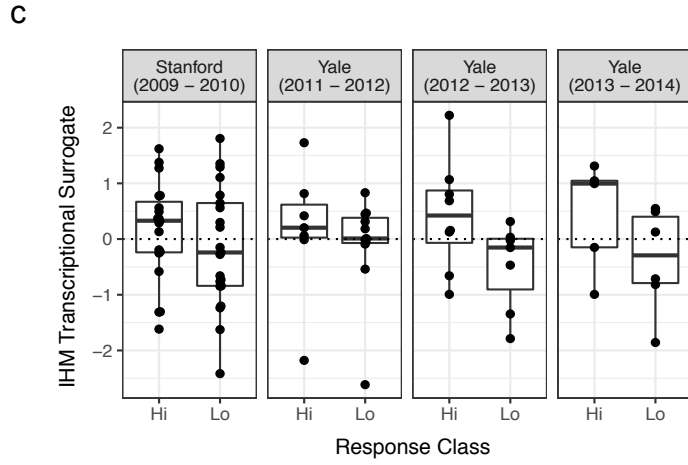
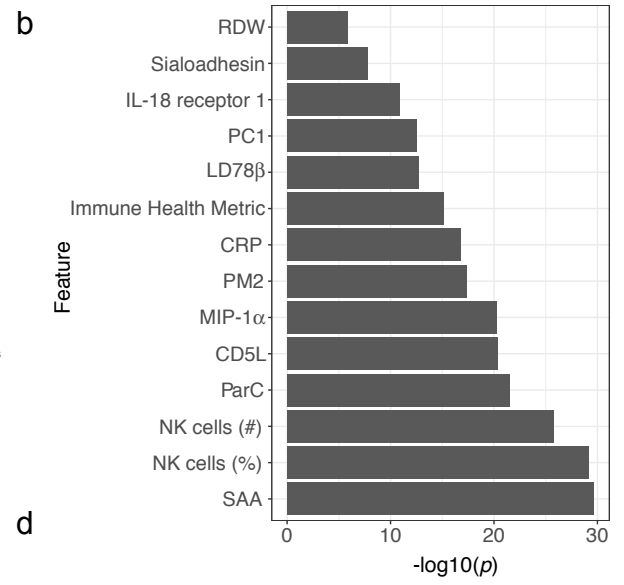
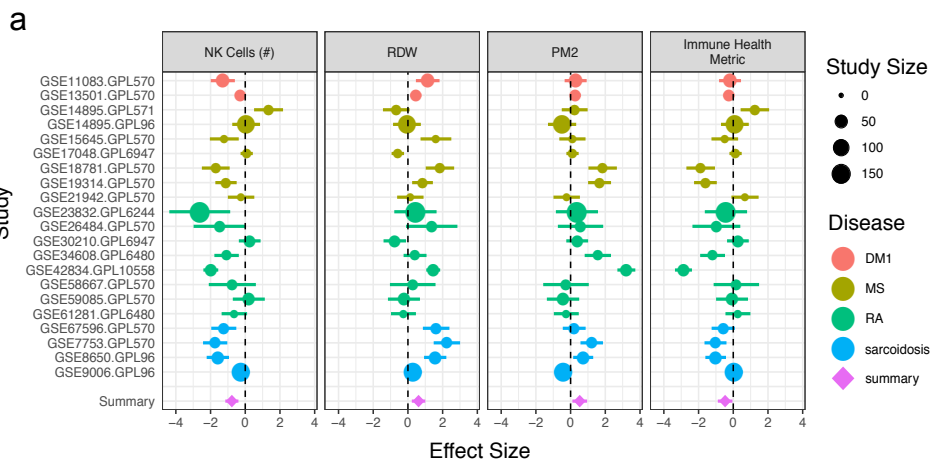
1163 **c**, Negative log<sub>10</sub> adjusted *p* values (FDR) of Global Variable Importance of features in each  
1164 Random Forest classifier. *P* values were determined through permutation (see Methods). Labels  
1165 are shown for parameters passing an FDR cutoff of 0.2 for each classifier. FDR adjustment was  
1166 performed on *p* values for parameters within a classifier. Features used in classifier are shown  
1167 on x-axis. NK = natural killer. RDW = red cell distribution width.

1168 **d**, Enrichment of transcriptional surrogate signatures for the predictive features identified by  
1169 the Random Forest classifier in Fig. 4b; gene sets from KEGG pathways, GO biological  
1170 processes, Reactome pathways, and the blood transcriptomic modules (BTMs) were included  
1171 for the enrichment analysis. SAA = serum amyloid A.

1172 **e**, Scatterplots with regression lines and associated Pearson correlations and *p* values of  
1173 subjects' Immune Health Metric (IHM) scores vs. the first 3 PC scores from the jPCs,  
1174 transcriptomic individual PCs (transcriptomic iPCs), and serum protein individual PCs  
1175 (proteomic iPCs). N = 182 subjects with both jPC and IHM scores. Pearson correlation and  
1176 associated *p* value are shown.



# Extended Data Figure 6



1177 **Extended Data Figure 6. Supporting data for assessing the Immune Health Metric (IHM).**

1178 **a**, Forest plot showing the effect sizes and associated standard errors in each study in the meta-  
1179 analysis for a selection of the transcriptional surrogate signatures capturing the status of the  
1180 indicated parameters (e.g., NK cell number). Summary meta-effect sizes shown at the bottom.  
1181 Size of circles indicates the relative sample numbers of each study. Effect sizes correspond to  
1182 average differences between disease and healthy, thus a positive effect size indicates that the  
1183 parameter was elevated in disease compared to healthy on average. Error bars show the 95%  
1184 confidence interval ( $1.96 * \text{standard error}$ ) in the meta-analysis.

1185 **b**, Barplot of  $-\log_{10} p$  value (two-sided Wilcoxon rank sum test) to assess whether genes in a  
1186 given transcriptional surrogate signature had significantly lower  $p$  values in the meta-analysis  
1187 results compared with genes not in the signature.

1188 **c**, Boxplots showing the transcriptional IHM scores of high and low responders in individual  
1189 studies from elderly vaccine meta-analysis.

1190 **d**, Venn Diagram showing the overlap between proteins in the IHM protein surrogate signature  
1191 and the original aging signature reported in the Baltimore Aging Study (odds ratio and  $p$  value  
1192 from the one-sided Fisher's exact test used to test the significance of the overlap).

1193 **e**, Scatterplot displaying the relationship between the IHM protein surrogate score and serum  
1194 IL-6 relative serum protein concentration (as measured by the Somalogic platform) in the  
1195 Baltimore Aging study (Spearman correlation and associated  $p$  value shown;  $n = 240$ ).

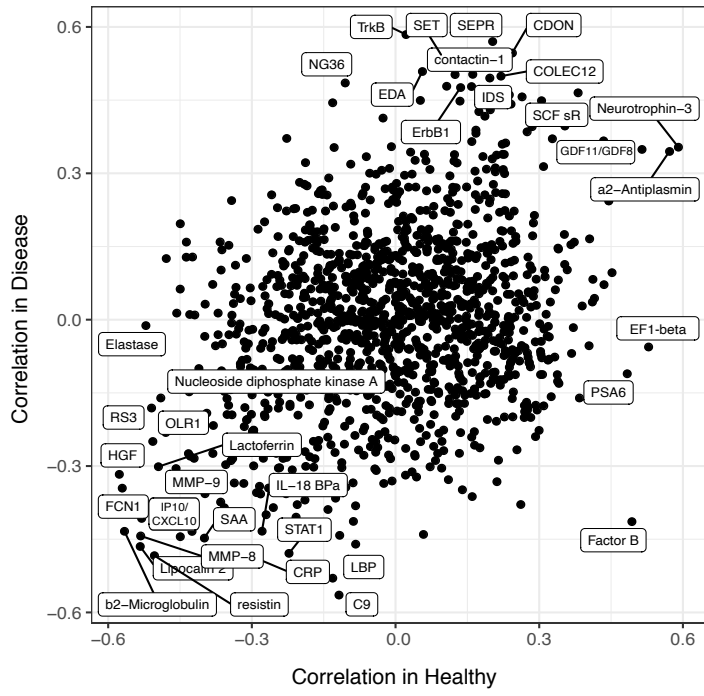
1196 **f**, Scatterplots showing the relative serum level of IL-6 (as measured by the Somalogic platform)  
1197 and the IHM in healthy subjects (left) and patients (right) in this study (Spearman correlation  
1198 and associated  $p$  values shown).  $n = 148$  and  $34$  disease and healthy subjects, respectively.

1199 **g**, Scatterplots showing association between the relative serum level of CXCL9/monokine  
1200 induced by gamma (MIG; as measured by the Somalogic platform) and the IHM in the healthy  
1201 subjects (left) and patients only (right) in our study (with Spearman correlation and  $p$  value  
1202 shown).  $n = 148$  and  $34$  disease and healthy subjects, respectively.

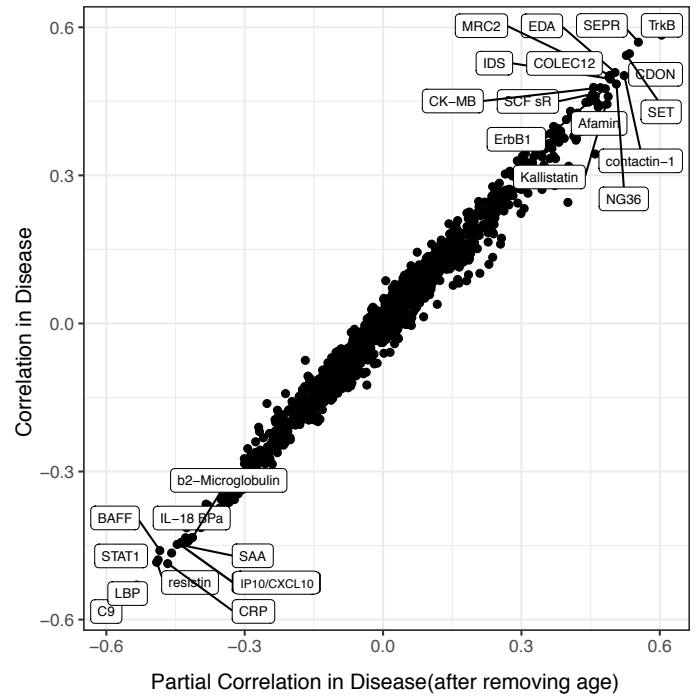
1203 **h**, The IHM was re-derived but without including PM2 (which contains CXCL9/MIG and  
1204 correlated proteins) during training or testing. Scatterplot shows the correlation between age  
1205 and this alternative IHM (without PM2) in the healthy subjects only (with Spearman correlation  
1206 and  $p$  value shown;  $n = 34$ ).

Extended Data Figure 7

a



b



1207 **Extended Data Figure 7. Supporting data for assessing the Immune Health Metric (IHM).**

1208 **a**, Scatterplot showing the Spearman correlation of serum proteins with the IHM transcriptional  
1209 surrogate signature within healthy individuals (x-axis) vs. disease individuals (y-axis) from the  
1210 monogenic cohort. The names of the 20 proteins with the highest absolute correlations on the x  
1211 or y axes are shown. Correlations were computed with  $n = 34$  healthy and  $n = 154$  for disease  
1212 individuals.

1213 **b**, Similar to Fig. 6g but showing the correlation and partial correlation computed in subjects  
1214 with disease only ( $n = 154$ ).

## Extended Data Table 1. Description of monogenic diseases in this study.

### Autoinflammatory Diseases

Disease Acronym	Gene/Protein	Disease Name	OMIM Number	Inheritance; Mutation effect	Phenotype	Pathomechanism of Inflammation	Ref
<b>CAPS</b>	<i>NLRP3</i> / NLRP3	Familial cold autoinflammatory syndrome (FCAS): NLRP3-associated autoinflammatory disease-mild Muckle-Wells syndrome (MWS): NLRP3-associated autoinflammatory disease-moderate	120100, 191900	Autosomal Dominant / De novo; Gain of Function Mutations	Fever, urticaria-like rash, CNS inflammation, bone overgrowth	Constitutively active NLRP3 inflammasome and increased IL-1 $\beta$ production	(Aksentjevich and Schnappauf, 2021; Manthiram et al., 2017; Tangye et al., 2020)
<b>DADA2</b>	<i>ADA2/CECR1</i> / ADA2	Deficiency of Adenosine Deaminase 2	615688	Autosomal Recessive; Loss of Function Mutations	Fever, lacunar strokes, livedo, immunodeficiency, anemia	Decrease in protein expression/activity leads to preferential differentiation of M1 proinflammatory macrophages,	(Aksentjevich and Schnappauf, 2021; Meyts and Aksentjevich, 2018)
<b>FMF</b>	<i>MEFV</i> / Pyrin	Familial Mediterranean Fever	249100	Autosomal Recessive; Gain of Function Mutations	Fever, serositis, rash, SAA amyloidosis	Facilitated activation of pyrin inflammasome leads to increased IL-1 $\beta$ production	(Aksentjevich and Schnappauf, 2021; Manthiram et al., 2017)
<b>HIDS/MKD</b>	<i>MVK</i> / MVK	Hyperimmunoglobulinemia D syndrome / Mevalonate Kinase Deficiency	260920, 610377	Autosomal Recessive; Loss of Function Mutations	Fever, serositis, rash, lymphadenopathy	Decrease in MVK activity enhances IL-1 $\beta$ production through activation of pyrin inflammasome	(Aksentjevich and Schnappauf, 2021; Manthiram et al., 2017)
<b>PAPA</b>	<i>PSTPIP1</i> / PSTPIP1	Pyogenic Arthritis, Pyoderma Gangrenosum and Acne Syndrome	604416	Autosomal Dominant / De novo; Not known	Pyoderma, pyogenic arthritis, severe cystic acne	Increased affinity to pyrin causes enhanced IL-1 $\beta$ production	(Aksentjevich and Schnappauf, 2021; Manthiram et al., 2017; Tangye et al., 2020)
<b>TRAPS</b>	<i>TNFRSF1A</i> / TNFR1	TNFR1-associated Periodic Syndrome	142680	Autosomal Dominant / De novo; Not known	Fever, serositis, rash, myalgia, orbital inflammation, SAA amyloidosis	Misfolding of extracellular domain of the receptor leads to intracellular protein retention and increased endoplasmic reticulum (ER) stress	(Cudrici et al., 2020; Tangye et al., 2020)

Primary Immunodeficiency Diseases (see [Tangye et al, 2020](#) for additional phenotypic and functional details and references)

Disease Acronym	Gene/Protein	Disease Name	OMIM Number	Inheritance; Mutation effect	Phenotype	Ref
<b>STAT1 GOF</b>	<i>STAT1</i> / STAT1	STAT1-gain-of-function	614162	Autosomal Dominant / De novo; Gain of Function Mutations	Chronic mucocutaneous candidiasis, bacterial infections, viral infections, autoimmunity	(Tangye et al., 2020; Toubiana et al., 2016)
<b>GATA2</b>	<i>GATA2</i> / GATA2	GATA2 deficiency / GATA2 haploinsufficiency	614172	Autosomal Dominant / De novo; Loss of Function Mutations	Lymphopenia, monocytopenia, myelodysplastic syndrome/acute myeloid leukemia, viral infections, NTM infection	(Spinner et al., 2014; Tangye et al., 2020)
<b>APDS1</b>	<i>PIK3CD</i> / p110 $\delta$ catalytic subunit of PI3K $\delta$	Activated PI3K delta syndrome 1	615513	Autosomal Dominant / De novo; Gain of Function Mutations	Bacterial infection, lymphoproliferation, herpesvirus infections, autoimmunity	(Coulter et al., 2017; Tangye et al., 2020)
<b>X-CGD</b>	<i>CYBB</i> / p91 <sup>phox</sup>	X-linked chronic granulomatous disease	306400	X-linked recessive; Loss of Function Mutations	Bacterial infection, invasive fungal infection, colitis, inflammatory lung disease, autoimmunity	(Arnold and Heimall, 2017; Henrickson et al., 2018; Tangye et al., 2020)
<b>p47-CGD</b>	<i>NCF1</i> / p47 <sup>phox</sup>	Autosomal recessive chronic granulomatous disease due to p47 <sup>phox</sup> deficiency	233700	Autosomal Recessive; Loss of Function Mutations	Bacterial infection, invasive fungal infection, colitis, inflammatory lung disease, autoimmunity	(Arnold and Heimall, 2017; Henrickson et al., 2018; Tangye et al., 2020)
<b>CTLA4</b>	<i>CTLA4</i> / CTLA4	CTLA4 haploinsufficiency	616100	Autosomal Dominant / De novo; Loss of Function Mutations	Hypogammaglobulinemia, lymphoproliferation, pulmonary infections, autoimmune cytopenias	(Schwab et al., 2018; Tangye et al., 2020)
<b>PGM3</b>	<i>PGM3</i> / PGM3	PGM3 deficiency	615816	Autosomal Recessive; Loss of Function Mutations	Bacterial infections, atopic dermatitis, elevated serum IgE, skeletal abnormalities, developmental delay	(Bergerson and Freeman, 2019; Tangye et al., 2020)
<b>LAD1</b>	<i>ITGB2</i> / integrin subunit $\beta$ 2	Leukocyte Adhesion Deficiency type 1	116920	Autosomal Recessive; Loss of Function Mutations	Periodontitis, skin infections, delayed umbilical cord separation	(Almarza Novoa et al., 2018; Tangye et al., 2020)
<b>IL12R</b>	<i>IL12R<math>\beta</math>1</i> / IL12R $\beta$ 1	IL-12 receptor $\beta$ 1 deficiency	614891	Autosomal Recessive; Loss of Function Mutations	Invasive mycobacterial disease, chronic mucocutaneous candidiasis, <i>Salmonella</i> infection	(Bustamante et al., 2014; Tangye et al., 2020)
<b>CARD14 DN</b>	<i>CARD14</i> / Caspase recruitment domain-containing protein 14	Dominant-negative CARD14 deficiency	607211	Autosomal Dominant / De novo; Dominant Negative Mutations	Severe atopic dermatitis, elevated serum IgE, food allergy, asthma	(Peled et al., 2019)

<b>NEMO</b>	<i>IKBKG</i> / inhibitor of nuclear factor kappa B kinase regulatory subunit gamma	NEMO deficiency	300636	X-linked recessive; Loss of Function Mutations	Ectodermal dysplasia, bacterial, viral, and mycobacterial infections, conical teeth, colitis	(Miot et al., 2017; Tangye et al., 2020)
<b>STAT3 DN</b>	<i>STAT3</i> / STAT3	STAT3-dominant-negative hyper-IgE syndrome / autosomal dominant hyper-IgE syndrome / Job's syndrome	147060	Autosomal Dominant / De novo; Dominant Negative Mutations	Bacterial infections, viral infections, atopic dermatitis, elevated serum IgE, skeletal and vascular abnormalities	(Bergerson and Freeman, 2019; Tangye et al., 2020)

## Telomere disorders

Disease Acronym	Gene/Protein or RNA	Disease Name	OMIM Number	Inheritance; Mutation effect	Phenotype	Ref
<b>TERT</b> <b>TERC</b>	<i>TERT</i> / TERT protein <i>TERC</i> / TERC RNA molecule	Telomere biology disorder, or telomereopathy	614742, 614743	Autosomal Recessive; Loss of Function Mutations	Hypocellular and aplastic anemia, pulmonary fibrosis, liver disease	(Townsend et al., 2014)

## References

- Aksentijevich, I., and Schnappauf, O. (2021). Molecular mechanisms of phenotypic variability in monogenic autoinflammatory diseases. *Nat. Rev. Rheumatol.* *17*, 405–425.
- Almarza Novoa, E., Kasbekar, S., Thrasher, A.J., Kohn, D.B., Sevilla, J., Nguyen, T., Schwartz, J.D., and Bueren, J.A. (2018). Leukocyte adhesion deficiency-I: A comprehensive review of all published cases. *J. Allergy Clin. Immunol. Pract.* *6*, 1418-1420.e10.
- Arnold, D.E., and Heimall, J.R. (2017). A Review of Chronic Granulomatous Disease. *Adv. Ther.* *34*, 2543–2557.
- Bergerson, J.R.E., and Freeman, A.F. (2019). An Update on Syndromes with a Hyper-IgE Phenotype. *Immunol. Allergy Clin. North Am.* *39*, 49–61.
- Bustamante, J., Boisson-Dupuis, S., Abel, L., and Casanova, J.-L. (2014). Mendelian susceptibility to mycobacterial disease: Genetic, immunological, and clinical features of inborn errors of IFN- $\gamma$  immunity. *Semin. Immunol.* *26*, 454–470.
- Coulter, T.I., Chandra, A., Bacon, C.M., Babar, J., Curtis, J., Screaton, N., Goodlad, J.R., Farmer, G., Steele, C.L., Leahy, T.R., et al. (2017). Clinical spectrum and features of activated phosphoinositide 3-kinase  $\delta$  syndrome: A large patient cohort study. *J. Allergy Clin. Immunol.* *139*, 597-606.e4.
- Cudrici, C., Deutch, N., and Aksentijevich, I. (2020). Revisiting TNF Receptor-Associated Periodic Syndrome (TRAPS): Current Perspectives. *Int. J. Mol. Sci.* *21*, 3263.
- Henrickson, S.E., Jongco, A.M., Thomsen, K.F., Garabedian, E.K., and Thomsen, I.P. (2018). Noninfectious Manifestations and Complications of Chronic Granulomatous Disease. *J. Pediatr. Infect. Dis. Soc.* *7*, S18–S24.
- Manthiram, K., Zhou, Q., Aksentijevich, I., and Kastner, D.L. (2017). The monogenic autoinflammatory diseases define new pathways in human innate immunity and inflammation. *Nat. Immunol.* *18*, 832–842.
- Meyts, I., and Aksentijevich, I. (2018). Deficiency of Adenosine Deaminase 2 (DADA2): Updates on the Phenotype, Genetics, Pathogenesis, and Treatment. *J. Clin. Immunol.* *38*, 569–578.



Miot, C., Imai, K., Imai, C., Mancini, A.J., Kucuk, Z.Y., Kawai, T., Nishikomori, R., Ito, E., Pellier, I., Dupuis Girod, S., et al. (2017). Hematopoietic stem cell transplantation in 29 patients hemizygous for hypomorphic IKBKG/NEMO mutations. *Blood* 130, 1456–1467.

Peled, A., Sarig, O., Sun, G., Samuelov, L., Ma, C.A., Zhang, Y., Dimaggio, T., Nelson, C.G., Stone, K.D., Freeman, A.F., et al. (2019). Loss-of-function mutations in caspase recruitment domain-containing protein 14 (CARD14) are associated with a severe variant of atopic dermatitis. *J. Allergy Clin. Immunol.* 143, 173-181.e10.

Schwab, C., Gabrysch, A., Olbrich, P., Patiño, V., Warnatz, K., Wolff, D., Hoshino, A., Kobayashi, M., Imai, K., Takagi, M., et al. (2018). Phenotype, penetrance, and treatment of 133 cytotoxic T-lymphocyte antigen 4–insufficient subjects. *J. Allergy Clin. Immunol.* 142, 1932–1946.

Spinner, M.A., Sanchez, L.A., Hsu, A.P., Shaw, P.A., Zerbe, C.S., Calvo, K.R., Arthur, D.C., Gu, W., Gould, C.M., Brewer, C.C., et al. (2014). GATA2 deficiency: a protean disorder of hematopoiesis, lymphatics, and immunity. *Blood* 123, 809–821.

Tangye, S.G., Al-Herz, W., Bousfiha, A., Chatila, T., Cunningham-Rundles, C., Etzioni, A., Franco, J.L., Holland, S.M., Klein, C., Morio, T., et al. (2020). Human Inborn Errors of Immunity: 2019 Update on the Classification from the International Union of Immunological Societies Expert Committee. *J. Clin. Immunol.* 40, 24–64.

Toubiana, J., Okada, S., Hiller, J., Oleastro, M., Lagos Gomez, M., Aldave Becerra, J.C., Ouachée-Chardin, M., Fouyssac, F., Girisha, K.M., Etzioni, A., et al. (2016). Heterozygous STAT1 gain-of-function mutations underlie an unexpectedly broad clinical phenotype. *Blood* 127, 3154–3164.

Townsley, D.M., Dumitriu, B., and Young, N.S. (2014). Bone marrow failure and the telomeropathies. *Blood* 124, 2775–2783.

## 1215 **Methods**

### 1216 ***Patient population and sample collection***

1217 Samples were collected on patients with monogenic immune disorders enrolled on National  
1218 Institutes of Health (NIH) protocols 00-I-0159 (NCT00006150), 01-I-0202 (NCT00018044), 07-I-  
1219 0033 (NCT00404560), 13-I-0157 (NCT01905826), 93-I-0119 (NCT05104723), 04-H-0012  
1220 (NCT00071045), and 94-HG-0105 (NCT00001373). Samples were collected when patients  
1221 presented to NIH for inpatient or routine outpatient care between September, 2015 and  
1222 November, 2017. Samples from matching healthy subjects were collected from subjects  
1223 enrolled on NIH protocols 91-I-0140 (NCT00001281) and 15-I-0162 (NCT02504853). These  
1224 studies were approved by the NIH Institutional Review Board and complied with all relevant  
1225 ethical regulations. Informed consent was obtained from all participants.

### 1226 ***RNA isolation***

1227 Blood was drawn directly into the Tempus Blood RNA Tube (Thermo Fisher Scientific, Waltham,  
1228 MA) according to manufacturer's protocol. Two Tempus tubes were collected per patient and  
1229 healthy donor. The blood sample from each Tempus tube was aliquoted in to two 4.5mL  
1230 cryovials. These cryovials were directly stored in -80°C freezer for long term.

1231 RNA was isolated from tempus blood samples using the Tempus Spin RNA Isolation kit (Thermo  
1232 Fisher Scientific, Waltham, MA) with following modifications to the manufacturer's protocol:  
1233 For each sample, 4ml of tempus blood sample was added to a 50ml conical tube containing  
1234 1.5ml of 1x PBS. The tubes were vortexed at full speed for 30 seconds, followed by  
1235 centrifugation at 3000 g for 1 hour at 4°C. After centrifugation, the supernatant from the tubes  
1236 was decanted and tubes were placed upside down on clean paper towels for 2 minutes. 400ul  
1237 of RNA Purification buffer was added, vortexed briefly to resuspend the pellet and transferred  
1238 the RNA to a purification filter with a pre-wet purification filter with 100ul wash solution I. The  
1239 tubes were centrifuged at 16,000 g for 30 seconds and liquid waste was discarded. A second  
1240 wash was done with 500ul wash solution I, followed by centrifugation at 16,000g for 30  
1241 seconds. The filter was washed with 500ul of wash solution 2 and centrifuged at 16,000 g for 60  
1242 seconds. DNase treatment was performed by adding 100ul of AbsoluteRNA wash solution  
1243 (Thermo Fisher Scientific, Waltham, MA), followed by 15 mins of incubation at room  
1244 temperature and 5 mins of incubation with wash solution 2. The tubes were spun at 16,000 for  
1245 60 seconds. The liquid waste was discarded, and empty tube was spun at 16,000 g for 30  
1246 seconds to remove any residual liquid and the filter was inserted into a new collection tube.

1247 The Nucleic Acid Purification Elution Solution was pre-warmed at 45°C. 100ul of this pre-  
1248 warmed elution solution was added to the filter and incubated at 37°C for 5 minutes. The tubes  
1249 were spun at 16,000 g for 2 minutes. The eluate was pipetted back to the filter and spun again  
1250 at 16,000 g for 1 minutes such that the eluate was collected in a new collection tube. 90ul of  
1251 the eluate was transferred to a new tube.

1252 RNA QC was performed using Qubit RNA BR assay (Thermo Fisher Scientific, Waltham, MA) and  
1253 Agilent RNA (Agilent Technologies, Santa Clara, CA). The average RIN was 8.26 and average  
1254 yield was 4.69 µg for the RNA samples.

#### 1255 ***Serum isolation***

1256 Serum was collected directly in Serum Separator Tubes and allowed to clot at room  
1257 temperature for a minimum of 30 minutes. Within two hours of blood collection, the tubes  
1258 were spun at 1800 g for 10 minutes at room temperature. The top (serum) layer was removed  
1259 via pipette and stored in individual vials at -80°C.

#### 1260 ***Microarray hybridization***

1261 All blood samples at different time points from the same subject were processed together.  
1262 Before assay, 396 samples were carefully batched into 14 groups according to their age, gender  
1263 and race but run blindly. One in-house reference sample was simultaneously processed with  
1264 the real samples in each batch. RNA was amplified from 300 ng of total RNA using Ambion WT  
1265 Expression Kit (Thermo Scientific, Wilmington, DE). Fragmented single-stranded sense cDNA  
1266 was terminally biotinylated and hybridized to the Affymetrix Human Gene 1.0 ST Arrays with  
1267 the probes for 36,079 RefSeq coding and noncoding transcripts and 466 lncRNA  
1268 transcripts (Affymetrix, Santa Clara, CA). The arrays were then washed and stained on a  
1269 GeneChip Fluidics Station 450 (Affymetrix); scanning was carried out with the GeneChip  
1270 Scanner 3000 and image analyzed with the Affymetrix GeneChip Command Console (AGCC)  
1271 software 4.0.

#### 1272 ***Somalomic SOMAScan Blood proteomic assays***

1273 Proteomic profiles for 1,305 SOMAmers in serum were assessed using the 1.3K SOMAScan  
1274 assay at the Trans-NIH Center for Human Immunology and Autoimmunity, and Inflammation  
1275 (CHI), National Institute of Allergy and Infectious Disease, National Institutes of Health  
1276 (Bethesda, MD, USA). Samples were run according to Somalomic standard operating procedures.  
1277 If operators identified presence of hemolysis in sample, those were marked for presence of  
1278 hemolysis (1 low- 4 high). In addition to Somalomic quality control samples, internal QC of the

1279 runs (cross checked of hemolyzed samples and outliers) was performed using CHI webtools  
1280 (Cheung *et al*). A total of 358 samples were included in this analysis. Two samples with high  
1281 levels of hemolysis (hemolysis score 4) and one sample with odd appearance were removed  
1282 from downstream analysis resulting in 355 total samples. The SOMAscan assay has a total of  
1283 1322 SOMAmer Reagents, and of these 12 are hybridization controls, which were removed  
1284 after hybridization normalization. 5 are nonspecifically-targeted SOMAmers (P05186; ALPL,  
1285 P09871; C1S, Q14126; DSG2, Q93038; TNFRSF25, Q9NQC3; RTN4, P00533; EGFRvIII, leaving  
1286 1305 somamers targeting 1273 unique proteins. The protein panel includes 4 proteins that are  
1287 rat homologues (P05413; FABP3, P48788; TINNI2, P19429; TINNI3, P01160; NPPA) of human  
1288 proteins and 4 viral proteins (HPV type 16, HPV type 18, isolate BEN, isolate LW123).

### 1289 ***Somalogic normalization***

1290 The Somalogic SOMAscan 1.3k assay data was normalized using the procedure outlined in<sup>1</sup>  
1291 followed by additional inter-plate batch correction prior to log transformation. As described in  
1292 <sup>1</sup>, hybridization control normalization (HybNorm) was first performed for each well on a plate,  
1293 and subsequent inter-plate calibration (CalNorm) was used to correct for plate-specific effects  
1294 between plates sharing the same Somalogic control samples. After these steps, median signal  
1295 normalization was performed on each group of samples from Somalogic plates that used the  
1296 same Somalogic control. This median normalization was performed to correct for shifts in the  
1297 median somamer RFUs across samples that may have been due to technical effects rather than  
1298 biological ones.

1299 Additionally, four bridge samples (QC\_CHI), derived from healthy donor blood, were added to  
1300 every run to allow in-house batch calibration normalization. These QC\_CHI samples were mixed  
1301 pools of serum samples of healthy donors from the Center for Human Immunology. In each  
1302 batch, the QC\_CHI controls were used for inter-plate calibration after the initial inter-plate  
1303 calibration with the Somalogic control samples. After this step, all relative protein expression  
1304 values were log<sub>2</sub> transformed.

### 1305 ***Curation of patient medication and medical metadata***

1306 Patient medical records were evaluated at the level of individual patient visits by trained  
1307 medical personnel. Medications used at the time of the visit were documented based on notes  
1308 from that visit; at the time of entry, medications were matched to the closest corresponding  
1309 term in MeSH. Medications were documented to include the route, dose, frequency, potency  
1310 (when applicable), date started and date ended (when available). Medical conditions were  
1311 obtained from chart review and were documented to include past and current medical history.  
1312 The conditions were entered by hand into a SQL database and selected from available terms in

1313 the Human Phenotype Ontology (HPO). Conditions that were unable to be reasonably matched  
1314 to HPO terms were entered with free text. Current medical conditions were denoted as one of  
1315 four options: 1) acute, active; 2) acute, resolved; 3) chronic, flare; 4) chronic, stable; 5) future  
1316 (for planned procedures or therapies).

### 1317 ***Microarray normalization, processing, filtering***

1318 Data were normalized and summarized to the probeset level using the RMA algorithm  
1319 implemented in the oligo R package<sup>2</sup>. Probesets mapping to multiple genes were discarded. To  
1320 select a single probeset for each gene, principal components analysis was performed for every  
1321 group of probesets corresponding to a given gene. The probeset most correlated with the first  
1322 principal component of this group was chosen as the “best” probeset to represent the  
1323 expression of this gene. With the microarray data summarized to the gene level, genes were  
1324 then filtered to remove genes that appeared lowly expressed or showed higher technical  
1325 variation than biological variation. Lowly expressed genes were identified as discussed in<sup>3</sup>;  
1326 briefly, a histogram of the median log<sub>2</sub> expression values were plotted and a lowly expressed  
1327 local maximum was identified. There exists a “plateau” where genes with low median intensity  
1328 are enriched. A manual threshold was selected to remove all genes in this enriched low  
1329 intensity area of the histogram. To determine the relative amounts of biological vs technical  
1330 variation, the variance of a gene in technical control samples (identical runs of same RNA) was  
1331 compared to the variance of the gene across all of the patient/healthy control samples. Those  
1332 genes with higher variance in the technical controls were removed from further analysis.

### 1333 ***Complete Blood Counts and lymphocyte phenotyping***

1334 Subjects had standard complete blood counts (CBCs) performed at the NIH Clinical Center in the  
1335 Department of Laboratory Medicine. Lymphocyte (T cell, B cell, NK cell) flow cytometry  
1336 quantification was performed using the BD FACS Canto11 flow cytometer. The following  
1337 parameters were collected on most patients, but were removed in downstream analysis for the  
1338 given reasons:

- 1339 ● Hematocrit measurements were removed, as they are highly redundant with  
1340 hemoglobin measurements
- 1341 ● Nucleated red blood cell measurements were removed, as they were zero for the  
1342 majority of patients.
- 1343 ● MPV, immature granulocytes (concentration and percent WBCs), CRP, and ESR  
1344 measurements were removed, as they were missing for 14, 62, 53, and 61 samples  
1345 respectively.

1346 Three samples were removed due to inconsistencies found in their data (the sum of the  
1347 absolute counts of cells from the TBNK assay was highly inconsistent with the total lymphocytes  
1348 from the complete blood counts).

1349 Absolute counts of leukocytes (including TBNK) were used for downstream analysis. The  
1350 neutrophil to lymphocyte ratio (NLR), the ratio between the neutrophil absolute counts and  
1351 lymphocyte absolute counts, was included as an additional CBC parameter for classification due  
1352 to its previously described association with multiple medical conditions such as infections and  
1353 cancer<sup>4,5</sup>.

#### 1354 ***Assignment of subjects to the main and set-aside cohorts***

1355

1356 From a total of 270 subjects (including 42 healthy controls), two sub-cohorts, namely *main* and  
1357 *set-aside*, were created with the purpose of holding out the *set-aside* group for future  
1358 validation and testing of specific hypotheses. Subjects with multiple visits were assigned to the  
1359 main group to allow for the assessment of temporal, intra-subject stability. The rest of the  
1360 participants were randomly assigned to one of the sub-cohorts to achieve a ratio of  
1361 approximately 80% *main* to 20% *set-aside* for each of the conditions, resulting in 217 and 53  
1362 subjects in the two groups, respectively. All analyses unless explicitly stated utilize only the  
1363 *main* subjects.

#### 1364 ***Averaging of technical replicate samples***

1365 Each measured parameter among technical replicate samples (samples taken from a patient  
1366 during the same visit) were averaged for downstream analysis after normalization (including  
1367 log<sub>2</sub> transformation for the Somalogic and Microarray data). Samples from the same visit were  
1368 considered technical replicates, although a visit could be an inpatient visit spanning several  
1369 days or a one-day outpatient visit (of 364 total visits in the study, 7 visits consisted of blood  
1370 draws over multiple days and 6 consisted of multiple draws on the same day). This was done for  
1371 gene log-intensities, protein log-RFUs, and CBC parameters. We refer to the data after  
1372 averaging across technical replicates as “sample-level” data.

#### 1373 ***Averaging of biological replicate samples***

1374 In situations where we wished to investigate data at the subject level rather than sample level,  
1375 we averaged each parameter over biological replicates in the sample-level data. We refer to the  
1376 result as “subject-level” data. Note that patient ages associated with a subject for a data type  
1377 were assigned to be the average age across all visits for which a sample of that data type was  
1378 collected. The largest time difference between samples from the same subject was 369 days.

1379 ***Gene and protein module creation***

1380 Weighted Gene Correlation Network Analysis (WGCNA)<sup>6</sup> was used to form modules of genes  
1381 and modules of proteins using the subject-level data (see *averaging of biological replicate*  
1382 *samples*). The parameters chosen were the same as the tutorial available at  
1383 <https://horvath.genetics.ucla.edu/html/CoexpressionNetwork/Rpackages/WGCNA/Tutorials/FemaleLiver-02-networkConstr-man.pdf>  
1384 with the following deviations: for the microarray data  
1385 and Somalomic data, a soft-threshold of 12 was manually chosen. Additionally, for the  
1386 Somalomic data the *cutreeDynamic* method parameter was set to ‘tree,’ as this provided  
1387 modules with greater variation explained by the 1st principal component compared to the  
1388 ‘hybrid’ method, as used in the microarray WGCNA analysis.

1389 Prior to module creation with WGCNA, samples were flagged as outliers by cutting an  
1390 agglomerative hierarchical tree formed from distances between samples in the sample-level  
1391 data. Data were scaled to unit variance prior to distance calculation. This was done separately  
1392 for each data type and tree cut heights for the proteomic and transcriptomic data hierarchical  
1393 trees were manually chosen 75 and 250 in each data type respectively. For both data types, the  
1394 minimum branch size required so that samples on the branch were not removed was set to 10.  
1395 The subject-level data was then rederived by averaging as before, but without these outlier  
1396 samples. Although outliers were removed during the module creation process to avoid these  
1397 extreme samples creating undue impact on the modules, these samples were included for  
1398 downstream analyses, as they may have been flagged as outliers due to their extreme  
1399 phenotypes (e.g. marrow failure) rather than technical noise. Thus, module activity scores were  
1400 still computed for these outlier samples, even though they were not used to inform the  
1401 creation of the modules.

1402 ***Gene and protein module activity scores***

1403 Module activity scores (sometimes referred to as module eigengenes) for a gene or protein  
1404 module were calculated for each sample in the following way: First, the subject-level data was  
1405 recomputed (using the same procedure described in ‘*Averaging of biological replicate samples*’)  
1406 from the sample-level data, after removing the outlier samples in the given data type. Next, the  
1407 module’s first principal component axis (PC1) was found through performing PCA on the  
1408 recomputed subject-level data, subsetted to only include features belonging to the module.  
1409 Then, for each sample in the sample-level data (including the outliers not used when deriving  
1410 the modules and principal component axes), the projection of the sample’s feature vector,  
1411 subsetted to only the features in a given module, onto the PC1 for that module was computed.  
1412 This result was assigned to be that sample’s activity score for that module. As the modules  
1413 were derived through signed WGCNA, the features in the modules were designed to be

1414 positively correlated with one another; however, PCA can produce PC's that are positively or  
1415 negatively correlated with the features. If a module's activity scores were negatively correlated  
1416 with more features in the module than were positively correlated, we multiplied that module's  
1417 activity scores (derived via PC1) by -1, such that the scores were positively correlated with most  
1418 of the features in the module. Samples were not assigned a module activity score for the grey  
1419 WGCNA module.

#### 1420 ***Analysis of feature stability***

1421 Variance component models were fit using the variancePartition package<sup>7</sup> to estimate the  
1422 sources of variation from a list of covariates for each feature in the transcriptomic and  
1423 proteomic data, leveraging repeat samples to estimate intra-subject temporal variation in  
1424 parameters. Two variance partition models were fit; The first model (VP\_M1) only includes the  
1425 subject as a random effect, with all other variation being considered "residual." The second  
1426 model (VP\_M2) includes subject, condition, and various binary medication variables as random  
1427 effects. The medication groups included in VP\_M2 were Monoclonal antibodies(not including  
1428 those for TNF and IL1), Anti-fungal, Antibiotic, Anti-TNF, Anti-IL1, Anti-inflammatory, IgG-  
1429 replacement, IFN-gamma, Immune-stimulator, Immunosuppressant, and Steroids. As patients  
1430 often were taking different combinations of medications, which potentially changed between  
1431 repeat samples, the medications were coded as binary variables denoting whether a patient  
1432 was or was not taking a given medication at the time of sampling. The individual variance  
1433 contributions assigned to each of the medications were then summed to a single medication-  
1434 associated variance contribution. Medications were included in the model if they were used by  
1435 many patients and not highly confounded with one of the condition groups.

1436 A feature was deemed to be stable if VP\_M1 estimated that there was more intrer-subject  
1437 variation than intra-subject variation in that feature (i.e. 50% or greater of the variation is  
1438 explained by patient covariate in VP\_M1). This determination was made for all data types  
1439 (transcriptomic measurements, transcriptomic modules activities, proteomic measurements,  
1440 proteomic module activities, and CBC+TBNK parameters). In various downstream analyses, only  
1441 the stable features as determined by this method were used.

1442 To evaluate the robustness of these estimates, VP\_M1 was performed with 100 replicates of  
1443 jackknife resampling in which 80% of subjects with repeats and 80% of subjects without repeats  
1444 were selected. Results were summarized with the mean variance explained by subject across  
1445 jackknife samples and the 95% confidence interval was taken as 2.5% quantile and the 97.%  
1446 quantile across jackknife samples.

#### 1447 ***Disease Signature/Differential expression analyses***



1448 To determine the disease signatures, Limma<sup>8</sup> was used to fit linear models and test differential  
1449 expression for each feature (Somamer, transcript, module or CBC/TBNK parameter). A single  
1450 model was fit for each feature that accounted for Condition, Gender, and Age, and Visit Type  
1451 (whether or not the patient reports feeling sick on a given visit): feature ~ condition + age +  
1452 gender + visit\_type.

1453 T-statistics and p-values were computed for the following contrasts of the coefficients:

- 1454 ● Disease vs. Healthy signatures
  - 1455 ○ Healthy was coded as the reference level and a t-statistics were computed for
  - 1456 the coefficient for each condition
- 1457 ● Each disease vs. all other diseases
  - 1458 ○ A contrast matrix was made such that each disease was compared to all other
  - 1459 diseases (the weights for each 'other' disease group were set to be equal).
- 1460 ● Comparison-specific contrasts were created to compare single diseases to others or
- 1461 groups of diseases to other groups.

1462 For tests involving the gene expression or proteomic modules, standard t-statistics (those  
1463 computed without empirical bayes moderation) were used to compute p-values due to the  
1464 lower number of features. For the individual proteomic or transcriptomic features, the  
1465 empirical Bayes moderated t-statistics<sup>8</sup> were computed and used to compute p values.  
1466 Multiple hypothesis correction was performed using the Benjamini-Hochberg<sup>9</sup> method to  
1467 compute FDR-adjusted p values.

#### 1468 ***Clustering genes within TM1: Interferon***

1469 The genes with TM1: Interferon were subclustered by computing the Euclidean distance matrix  
1470 between all genes based on the T-statistics from the differential expression analysis comparing  
1471 all conditions to Healthy Controls. The genes were clustered using Ward's method (method =  
1472 "Ward.D2") with the hclust function in R. The hierarchical clustering tree was then cut to  
1473 produce three clusters with the cutree function with k = 3.

#### 1474 ***JIVE analysis***

1475 The whole blood microarray and serum proteomics data (Somalogic) were filtered to select only  
1476 stable features (see *determining feature stability*). Data were averaged to the subject level (see  
1477 *averaging of biological replicates*). The JIVE algorithm<sup>10</sup> was then used to partition the data into  
1478 joint (sharing axes variation between the transcriptomic and proteomic data) and individual  
1479 (unique to a data-type) components. Input data were first z-score normalized for each feature

1480 and then each input matrix was scaled by the frobenius norm of that data type so as to not give  
1481 greater weight to data with more features (i.e. the transcriptomic data). The JIVE algorithm  
1482 produces 3 matrices for each data-type, representing joint (shared between data types),  
1483 individual (unique to that data-type) and residual (potentially noise) variation. JIVE PC scores  
1484 were computed for each subject using the prcomp function from R, using the resulting joint,  
1485 and individual matrices as inputs. To compute the joint PC scores (jPC's), the transcriptomic and  
1486 proteomic joint matrices were concatenated to a single joint matrix prior to calculation of the  
1487 PC scores.

#### 1488 ***JIVE variance explained calculations***

1489 To calculate the amount of variation explained by each of the joint and individual components  
1490 from the JIVE analysis, we computed the frobenius norm of the input data (proteomic or  
1491 transcriptomic) to determine the total amount of variation present in a given data matrix. This  
1492 same computation was then applied to the resulting joint and individual matrices. Dividing the  
1493 variation in the joint and individual matrices by the amount of total variation gives the variance  
1494 explained by each of these respectively. Lastly to determine residual variation, the joint and  
1495 individual variation were subtracted from the total variation.

#### 1496 ***JIVE PC geneset enrichment***

1497 To determine the gene set enrichments for the JIVE PC's, the whole blood microarray and  
1498 serum proteome data were separately correlated with each JIVE PC. Genesets were then tested  
1499 for enrichment of correlation to each PC in each data type separately, using the two-sided t-test  
1500 with correlation described in Wu, Di, and Gordon K. Smyth. Nucleic acids research 40.17 (2012),  
1501 using the cameraPR function from limma<sup>8</sup> with use.ranks = FALSE.

#### 1502 ***Leukocyte composite score***

1503 A leukocyte composite score was computed for each patient by first averaging repeated  
1504 observations from a given patient. A Z-score was then computed for the lymphocyte, neutrophil  
1505 and monocyte counts relative to the healthy mean and standard deviation, for that parameter.  
1506 The three Z-scores were then averaged across the cell-types to give the final composite score.

#### 1507 ***Creation of Immune Health Metric***

1508 The Immune Health Metric presented represents the likelihood that a given subject is a healthy  
1509 control according to the leave one out cross validation (LOO CV) prediction probabilities of our  
1510 random forest model.

1511 Prior to training the models, we subsetted the subjects to those that had measurements from  
1512 all of the following data: proteomic, transcriptomic, and CBC/TBNK (and passed respective  
1513 quality checks). Biological replicate samples from the same patients were averaged, so that  
1514 each subject had one associated value for each measured feature. Features included for  
1515 classification were subsetted to those for which the VP\_M1 variance partition model assigned  
1516 at least 50% of the variation to the patient covariate (i.e. the stable features).

1517 Three unimodal classifier schemas were designed: a proteomic module classifier, a  
1518 transcriptomic module classifier, and a CBC parameter classifier, using the stable features from  
1519 each respective data type.

1520 Two multimodal classifiers were also created: the first included all features from the three  
1521 unimodal classifiers. The second included all features from the first, but also included the log-  
1522 RFUs of all singleton proteins (the proteins in the grey Somalogic module). Each classifier  
1523 described above was then evaluated using leave-one-out cross validation, and an ROC curve  
1524 was generated from the LOO CV probabilities of being a healthy subject (the positive class).

1525 Predicting healthy subjects vs. disease using all subjects, we computed the LOO CV prediction  
1526 probabilities that an individual was a healthy control, that we termed the Immune Health  
1527 Metric.

#### 1528 ***Classification accuracy using set aside patients***

1529 The second multimodal classifier incorporating module activity scores, immune cell frequencies,  
1530 and grey module protein RFUs was trained using all subjects in the *main* set of subjects. The  
1531 disease vs. healthy status of set aside subjects was then predicted and an ROC curve was  
1532 generated from the predicted probabilities of being a healthy subject (the positive class).

#### 1533 ***Statistical testing of classification feature global variable importance***

1534 For each classifier, the global variable importance (GVI) of all features were collected after  
1535 training the classifier on all subjects used in the creation of the Immune Health Metric.

1536 To find the significance of the global variable importance (GVI) for each feature, permutation  
1537 testing was performed to determine how often the GVI, as estimated by classifiers training on  
1538 permuted class labels, was higher than the classifier trained on the true labels. A total of  
1539 10,000,000 permutations were performed.

#### 1540 ***Condition-specific classifiers***

1541 One-versus-all-condition binary classifiers were created for the largest groups of patient  
1542 conditions: CGDs (XCGDs and 47CGDs were combined), Job, STAT1 GOF, and FMF. Each one-  
1543 versus-all classifiers for each group were created analogously to the multimodal classifier  
1544 including all modules, CBC +TBNK, and grey module proteins created to differentiate healthy  
1545 subjects from monogenic patients. Feature GVIs were identified and tested analogously as well.  
1546 Note that for the disease-versus-all classifiers, healthy controls were excluded from the LOO CV  
1547 model training, prediction, and calculation of feature GVI.

1548 ***Transcriptional surrogate signatures for autoimmunity meta-analysis validation***

1549 Transcriptional signatures for features from the three following categories were created:

- 1550 ● Immune Health Metric
- 1551 ● jPC1
- 1552 ● Features: all features from multimodal classifier that passed GVI testing with an FDR-  
1553 adjusted p value of less than 0.20

1554 Signatures in the indexes and features categories both were formed by taking the 150 genes  
1555 from the stable microarray features with highest correlation to the feature (based upon  
1556 correlation with all subjects in our training cohort, including healthy controls). Selected genes  
1557 were then subsetted to those with a Spearman correlation to the feature of interest of more  
1558 than 0.35 in magnitude. Genes in the signature were then divided into two groups: those  
1559 positively correlated with the index/feature of interest, and those negatively correlated.  
1560 Module signatures were all simply composed of the genes that the module was comprised of  
1561 (stable and unstable). All these genes were placed in the positive correlates group of the  
1562 signature, as we used a signed WGCNA performed to derive the modules.

1563 To assign each subject in the validation study a signature score, we subsetted the genes in the  
1564 surrogate signatures to those also measured in the validation studies and we then averaged the  
1565 z-scores of each gene/protein (scaled across subjects) for each gene in the signature. Note that  
1566 z-scores of proteins in the ‘negative correlates’ group were flipped in sign prior to averaging.

1567 ***Proteomic Immune Health Metric surrogate signature for aging validation using data from***  
1568 ***Tanaka 2018***

1569 The proteomic IHM surrogate was derived and computed analogously to the transcriptional  
1570 surrogate signatures as described above, with one small modification: to ensure that the  
1571 signature was not reliant on proteins that had substantive relative differential abundance in  
1572 serum compared to plasma (the data in which we planned to test these signatures), we

1573 removed any Somamers that fell into different dilution groups between plasma assays and  
1574 serum assays.

### 1575 ***Autoimmune disease cohort meta-analysis***

1576 Comparison group pairs (CGPs) for the OMiCC Jamboree<sup>11</sup> were used to test our transcriptional  
1577 surrogate signatures in other data sets. Briefly, CGPs from the same study and platform were  
1578 combined to ensure that samples were not being replicated across studies. Samples from the  
1579 same patient in a study were removed manually. Several CGPs used in the OMiCC jamboree  
1580 were removed for the following reasons:

- 1581 • CGPs/studies of systemic lupus erythematosus (SLE) appearing in Lau *et al*<sup>11</sup>  
1582 were removed as many genes in the signatures to be tested were not present in the  
1583 platforms used.
- 1584 • CGP 'GSE9006-Diabetes\_Mellitus,\_Type\_1-PBMC\_newlydiagnosed\_paired with 1  
1585 month follow up::GSE9006-Healthy-PBMC\_unpaired' was not included because  
1586 samples in this CGP were follow up samples from another CGP, GSE9006-  
1587 Diabetes\_Mellitus,\_Type\_1-PBMC\_newly diagnosed\_unpaired::GSE9006-Healthy-  
1588 PBMC\_unpaired
- 1589 • CGPs 'Jam\_human\_RA\_GSE26554-JIA-PBMC::Jam\_human\_RA\_GSE26554-  
1590 Control-PBMC', 'Jam\_Human\_RA\_JIA-PBMC::Jam\_Human\_RA\_Controls-PBMC',  
1591 'Jam\_human\_RA\_GSE26554-OligoarticularJIA-PBMC::Jam\_human\_RA\_GSE26554-  
1592 Control-PBMC', and 'Jam\_Human\_RA\_JIA-PBMC::Jam\_Human\_RA\_Controls-PBMC',  
1593 were removed because the all had many overlapping samples with another CGP  
1594 already included in our study, Jam\_Human\_RA\_JIA-  
1595 PBMC::Jam\_Human\_RA\_Controls-PBMC.
- 1596 • CGP 'Jam\_human\_RA\_GSE61281-Psoriatic\_arthritis-  
1597 Whole\_blood::Cutaneouspsoriasis without arthritis\_GSE61281-  
1598 Cutaneous\_psoriasis\_without\_arthritis-Whole\_blood' was removed because the  
1599 control patients had psoriasis.

1600

1601 Additionally, some samples were removed within certain studies

- 1602 • GSE21942
  - 1603 ○ GSM545843, GSM545845 were removed as these were technical replicates
  - 1604 ○ of other samples in the study
- 1605 • GSE30210
  - 1606 ○ We removed additional biological replicates from patients that were sampled
  - 1607 ○ longitudinally and we selected the last sample for each patient

- 1608           • GSE8650  
1609           ○ We removed additional biological and technical replicates from the same  
1610           individual. The last sample was selected for each patient.  
1611           ○ Samples GSM214490 and GSM214492 were removed as they were believed  
1612           to have unreliable diagnoses according to the original publication  
1613           • GSE15645  
1614           ○ We removed patients who were experiencing clinical remission of symptoms  
1615           • GSE42834  
1616           ○ We removed patients with non-active sarcoid  
1617           A complete listing of the studies and all case/control samples in the meta-analysis can be found  
1618           in Supplementary Table 19

1619           Each study was quantile normalized within the study. The standard pipeline from the  
1620           metaIntegrator package<sup>12</sup> was then used to compute meta effect sizes of each of the surrogate  
1621           signature scores. Meta-analysis was also performed for all genes that overlapped with those in  
1622           our the monogenic microarray data and Wilcoxon tests were also used to determine whether  
1623           genes belonging to each transcriptomic surrogate signature tended to have higher meta-effect  
1624           sizes than genes that did not belong to the signature.

### 1625           ***Overlap of Baltimore Aging signature and Proteomic Immune Health Metric***

1626           We considered the proteins passing an FDR-adjusted significance threshold of 0.05 from  
1627           Supplementary table 3 of Tanaka *et al*<sup>13</sup> as the previous aging signature. These proteins were  
1628           compared to the proteins from the Immune Health metric proteomic surrogate with a one-  
1629           sided Fisher's exact test, with the alternative hypothesis being that the overlap was greater  
1630           than that expected by chance.

### 1631           ***Gene set enrichment analyses***

1632           Gene modules from the transcriptomics data were tested using hypergeometric tests for the  
1633           following collections of gene sets: The Li blood transcriptomic modules<sup>14</sup>, Kyoto Encyclopedia of  
1634           Genes and Genomes<sup>15</sup>, Reactome<sup>16</sup>, and Gene Ontology Biological Processes<sup>17</sup>. For each  
1635           module, FDR multiple hypothesis corrections were performed on all gene sets (pooled across  
1636           collections).

1637           Proteomic modules were tested for gene set enrichments analogously after converting each  
1638           protein targets of Somamers to their respective gene according to the SomaScan assay.  
1639           Proteins that mapped to multiple genes were removed from the analysis. Additionally, some  
1640           genes corresponded to multiple proteins. In this case, when testing a gene module, genes that

1641 mapped to both proteins in and outside of the module were removed from the module and the  
1642 background proteins.

1643 An analogous analysis was performed for the proteomic modules using gene sets from the  
1644 Human Protein Atlas<sup>18</sup>. Gene sets were made for various tissues by looking for proteins  
1645 enriched for that tissue based on the HPA. The following categories were considered for  
1646 enrichments: “enriched”, “enhanced”, and “tissue enriched.”

#### 1647 **Correlation of serum proteins with IHM surrogate transcriptional signature**

1648 The correlation, without removing the effect of age, was computed simply by computing the  
1649 Spearman correlation of every protein with the IHM surrogate signature. We additionally  
1650 computed partial correlations where the effect of age had been removed from both the protein  
1651 data and IHM transcriptional surrogate signature by using the limma removeBatchEffect  
1652 function with age as the single covariate, which fits a linear model (feature ~age) to remove the  
1653 effect of age prior to computing the correlation of each protein with the IHM transcriptional  
1654 signature.

#### 1655 **Testing IHM and jPC1 signatures in Ota *et al*<sup>19</sup> 2021 sorted cell data**

1656 Data were downloaded from [https://ddbj.nig.ac.jp/public/ddbj\\_database/gea/experiment/E-GEAD-000/E-GEAD-397/](https://ddbj.nig.ac.jp/public/ddbj_database/gea/experiment/E-GEAD-000/E-GEAD-397/). For each cell-type, the log cpm values with TMM normalization were  
1657 computed using edgeR. We noted a large batch effect due to the “Phase” of the study and thus  
1658 removed the phase effect at the individual gene level using limma’s removeBatchEffect  
1659 function. After this, genes were z-scored normalized and signature scores were computed as  
1660 described in the section above *Transcriptional surrogate signatures for autoimmunity meta-*  
1661 *analysis validation*. We then tested for differences in signature scores between healthy and  
1662 disease using linear models with limma. The association with age within healthy individuals only  
1663 was assessed using the Pearson correlation as implemented in the cor.test function in R.  
1664

#### 1665 **Vaccination response in elderly meta-analysis**

Cohort	Ages	Source
Stanford (2009-2010)	61-90 years	Furman, 2013 (SDY 212)
Yale (2011-2012)	66-93 years	Avey, 2020 (GSE65442)
Yale (2012-2013)	65-88 years	Avey, 2020 (GSE95584)
Yale (2013-2014)	65-86 years	Avey, 2020 (GSE101709)

1666

1667 Gene expression profiles for Yale vaccination subjects were quantile normalized using the R  
1668 package *limma*. Processed expression data from SDY212 was downloaded from *ImmuneSpace*.  
1669 Each dataset was filtered to baseline, pre-vaccination samples from subjects over the age of 60.  
1670 High and low antibody response labels for each subject were derived from HAI titer  
1671 measurements using the maximum residual after baseline adjustment (maxRBA) end point <sup>20</sup>.  
1672 IHM signature scores were calculated in each subject using the *MetaIntegrator* R package.  
1673 Briefly, the signature score for each subject was calculated from normalized, log<sub>2</sub> transformed  
1674 gene expression data by taking the geometric mean of positive signature genes and subtracting  
1675 the geometric mean of negative signature genes. The standardized mean difference of baseline  
1676 IHM scores between high and low antibody responders was estimated by fitting a random  
1677 effects model using the *metafor* R package.

### 1678 **Checks of robustness to variation in cell frequencies**

1679 Linear models were fit using the *lm* function in R both with and without including cell  
1680 frequencies in the model. Cell frequencies were included as percent of total white blood cells  
1681 and included major cell populations from the CBC/TBNK, specifically neutrophils, monocytes,  
1682 CD4 T-cells, CD8 T-cells, B cells, NK cells, eosinophils, and basophils. The percent mediation,  
1683 which reflects how much of the main effect can be explained by additional covariates, was  
1684 calculated as:  $1 - \text{coefficient\_without\_controlling\_for\_cell\_freq} /$   
1685  $\text{coefficient\_with\_controlling\_for\_cell\_freq}$ .

1686

### 1687 **Methods References**

- 1688 1. Candia, J. *et al.* Assessment of Variability in the SOMAscan Assay. *Sci. Rep.* **7**, 14248  
1689 (2017).
- 1690 2. Carvalho, B. S. & Irizarry, R. A. A framework for oligonucleotide microarray  
1691 preprocessing. *Bioinforma. Oxf. Engl.* **26**, 2363–2367 (2010).
- 1692 3. Klaus, B. & Reisenauer, S. An end to end workflow for differential gene expression using  
1693 Affymetrix microarrays. (2018) doi:10.12688/f1000research.8967.2.
- 1694 4. Templeton, A. J. *et al.* Prognostic Role of Neutrophil-to-Lymphocyte Ratio in Solid  
1695 Tumors: A Systematic Review and Meta-Analysis. *JNCI J. Natl. Cancer Inst.* **106**, (2014).
- 1696 5. Russell, C. D. *et al.* The utility of peripheral blood leucocyte ratios as biomarkers in  
1697 infectious diseases: A systematic review and meta-analysis. *J. Infect.* **78**, 339–348 (2019).
- 1698 6. Langfelder, P. & Horvath, S. WGCNA: an R package for weighted correlation network  
1699 analysis. *BMC Bioinformatics* **9**, 559 (2008).
- 1700 7. Hoffman, G. E. & Schadt, E. E. variancePartition: interpreting drivers of variation in  
1701 complex gene expression studies. *BMC Bioinformatics* **17**, 483 (2016).
- 1702 8. Smyth, G. K. Linear models and empirical bayes methods for assessing differential  
1703 expression in microarray experiments. *Stat. Appl. Genet. Mol. Biol.* **3**, Article3 (2004).



- 1704 9. Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and  
1705 Powerful Approach to Multiple Testing. *J. R. Stat. Soc. Ser. B Methodol.* **57**, 289–300 (1995).
- 1706 10. Lock, E. F., Hoadley, K. A., Marron, J. S. & Nobel, A. B. JOINT AND INDIVIDUAL  
1707 VARIATION EXPLAINED (JIVE) FOR INTEGRATED ANALYSIS OF MULTIPLE DATA TYPES. *Ann. Appl.*  
1708 *Stat.* **7**, 523–542 (2013).
- 1709 11. Lau, W. W., Sparks, R., OMiCC Jamboree Working Group & Tsang, J. S. Meta-analysis of  
1710 crowdsourced data compendia suggests pan-disease transcriptional signatures of  
1711 autoimmunity. *F1000Research* **5**, 2884 (2016).
- 1712 12. Haynes, W. A. *et al.* Empowering Multi-Cohort Gene Expression Analysis to Increase  
1713 Reproducibility. <http://biorxiv.org/lookup/doi/10.1101/071514> (2016) doi:10.1101/071514.
- 1714 13. Tanaka, T. *et al.* Plasma proteomic signature of age in healthy humans. *Aging Cell* **17**,  
1715 e12799 (2018).
- 1716 14. Li, S. *et al.* Molecular signatures of antibody responses derived from a systems biology  
1717 study of five human vaccines. *Nat. Immunol.* **15**, 195–204 (2014).
- 1718 15. Kanehisa, M. & Goto, S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids*  
1719 *Res.* **28**, 27–30 (2000).
- 1720 16. Jassal, B. *et al.* The reactome pathway knowledgebase. *Nucleic Acids Res.* **48**, D498–  
1721 D503 (2020).
- 1722 17. Ashburner, M. *et al.* Gene ontology: tool for the unification of biology. The Gene  
1723 Ontology Consortium. *Nat. Genet.* **25**, 25–29 (2000).
- 1724 18. Uhlén, M. *et al.* Tissue-based map of the human proteome. *Science* **347**, (2015).
- 1725 19. Ota, M. *et al.* Dynamic landscape of immune cell-specific gene regulation in immune-  
1726 mediated diseases. *Cell* **184**, 3006-3021.e17 (2021).
- 1727 20. Avey, S. *et al.* Seasonal Variability and Shared Molecular Signatures of Inactivated  
1728 Influenza Vaccination in Young and Older Adults. *J. Immunol.* **204**, 1661–1673 (2020).
- 1729

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [ExtendedDataTables.xlsx](#)