

1 **The genetic architecture of human infectious diseases and pathogen-**
2 **induced cellular phenotypes**

3
4 **Authors:** Andrew T. Hale^{1,2}, Dan Zhou², Lisa Bastarache³, Liuyang Wang^{4,5}, Sandra S. Zinkel⁶,
5 Steven J. Schiff⁷, Dennis C. Ko⁴, and Eric R. Gamazon^{2,8,9,10*}

6
7 ¹Vanderbilt University School of Medicine, Medical Scientist Training Program, Nashville, TN.

8 ²Vanderbilt Genetics Institute & Division of Genetic Medicine, Vanderbilt University Medical
9 Center, Nashville, TN.

10 ³Department of Bioinformatics, Vanderbilt University School of Medicine, Nashville, TN.

11 ⁴Department of Molecular Genetics and Microbiology, Duke University School of Medicine,
12 Durham, NC.

13 ⁵Division of Infectious Diseases, Department of Medicine, Duke University School of Medicine,
14 Durham, NC.

15 ⁶Division of Hematology and Oncology, Vanderbilt University Medical Center, Nashville, TN

16 ⁷Center for Neural Engineering and Infectious Disease Dynamics, Departments of
17 Neurosurgery, Engineering Science and Mechanics, and Physics, Penn State University.
18 University Park, PA.

19 ⁸Clare Hall, University of Cambridge, Cambridge, UK.

20 ⁹MRC Epidemiology Unit, University of Cambridge, Cambridge, UK

21 ¹⁰Lead Contact.

22

23 * Correspondence and requests for materials should be addressed to E.R.G.

24 (eric.gamazon@vanderbilt.edu)

25

26 **SUMMARY**

27

28 Infectious diseases (ID) represent a significant proportion of morbidity and mortality across the
29 world. Host genetic variation is likely to contribute to ID risk and downstream clinical outcomes,
30 but there is a need for a genetics-anchored framework to decipher molecular mechanisms of
31 disease risk, infer causal effect on potential complications, and identify instruments for drug
32 target discovery. Here we perform transcriptome-wide association studies (TWAS) of 35 clinical
33 ID traits in a cohort of 23,294 individuals, identifying 70 gene-level associations with 26 ID traits.
34 Replication in two large-scale biobanks provides additional support for the identified
35 associations. A phenome-scale scan of the 70 gene-level associations across hematologic,
36 respiratory, cardiovascular, and neurologic traits proposes a molecular basis for known
37 complications of the ID traits. Using Mendelian Randomization, we then provide causal support
38 for the effect of the ID traits on adverse outcomes. The rich resource of genetic information
39 linked to serologic tests and pathogen cultures from bronchoalveolar lavage, sputum,
40 sinus/nasopharyngeal, tracheal, and blood samples (up to 7,699 positive pathogen cultures
41 across 92 unique genera) that we leverage provides a platform to interrogate the genetic basis
42 of compartment-specific infection and colonization. To accelerate insights into cellular
43 mechanisms, we develop a TWAS repository of gene-level associations in a broad collection of
44 human tissues with 79 pathogen-exposure induced cellular phenotypes as a discovery and
45 replication platform. Cellular phenotypes of infection by 8 pathogens included pathogen
46 invasion, intercellular spread, cytokine production, and pyroptosis. These rich datasets will
47 facilitate mechanistic insights into the role of host genetic variation on ID risk and
48 pathophysiology, with important implications for our molecular understanding of potentially
49 severe phenotypic outcomes.

50

51

52 **Keywords:** PrediXcan; Human Genetics; Infectious Disease; Transcriptomics; TWAS; GWAS;

53 Electronic Health Records; Hi-HOST; BioVU; UK Biobank; FinnGen; Functional Genomics;

54 GTEx; Phenome Scan: PheWAS; Mendelian Randomization; Clinical Microbiology

55

56 **HIGHLIGHTS**

57

58 • Atlas of genome-wide association studies (GWAS) and transcriptome-wide association
59 studies (TWAS) results for 35 clinical infectious disease (ID) phenotypes, with genome-
60 wide and transcriptome-wide significant results for 13 and 26 clinical ID traits,
61 respectively

62

63 • Phenome-scale scan of ID-associated genes across 197 hematologic, respiratory,
64 cardiovascular, and neurologic traits, facilitating identification of genes associated with
65 known complications of the ID traits

66

67 • Mendelian Randomization analysis, leveraging naturally occurring DNA sequence
68 variation to perform “randomized controlled trials” to test the causal effect of ID traits on
69 potential outcomes and complications

70

71 • A genomic resource of TWAS associations for 79 pathogen-induced cellular traits from
72 High-throughput Human *in vitro* Susceptibility Testing (Hi-HOST) across 44 tissues as a
73 discovery and replication platform to enable *in silico* cellular microbiology and functional
74 genomic experiments

75

76

77

78 INTRODUCTION

79 Genome-wide association studies (GWAS) and large-scale DNA biobanks with
80 phenome-scale information are making it possible to identify the genetic basis of a wide range
81 of complex traits in humans (Bycroft et al., 2018; Roden et al., 2008). A parallel development is
82 the increasing availability of GWAS summary statistics, facilitating genetic analyses of entire
83 disease classes and promising considerably improved resolution of genetic effects on human
84 disease (Cotsapas et al., 2011; Gamazon et al., 2019). Recent analysis involving 558 well-
85 powered GWAS results found that trait-associated loci cover ~50% of the genome, enriched in
86 both coding and regulatory regions, and of these, ~90% are implicated in multiple traits
87 (Watanabe et al., 2019). However, the breadth of clinical and biological information in these
88 datasets will require new methodologies and additional high-dimensional data to advance our
89 understanding of the genetic architecture of complex traits and relevant molecular mechanisms
90 (Bulik-Sullivan et al., 2015; Gamazon et al., 2018; Shi et al., 2016). Approaches to
91 understanding the functional consequences of implicated loci and genes are needed to
92 determine causal pathways and potential mechanisms for pharmacological intervention.

93 The genetic basis of infectious disease (ID) risk and severity has been relatively
94 understudied, and its implications for etiological understanding of human disease and drug
95 target discovery may be investigated using phenome-scale information increasingly available in
96 these biobanks. ID risk and pathogenesis is likely to be multifactorial, resulting from a complex
97 interplay of host genetic variation, environmental exposure, and pathogen-specific molecular
98 mechanisms. With few exceptions, the extent to which susceptibility to ID is correlated with host
99 genetic variation remains poorly understood (de Bakker and Telenti, 2010). However, for at
100 least some ID traits, including poliomyelitis, hepatitis, and *Helicobacter pylori* (Burgner et al.,
101 2006; Herndon and Jennings, 1951; Hohler et al., 2002; Malaty et al., 1994), disease risk is
102 heritable, based on twin studies. Although monogenic mechanisms of ID risk have been

103 demonstrated (Casanova, 2015a, b), the contribution of variants across the entire allele
104 frequency spectrum to interindividual variability in ID risk remains largely unexplored.

105 Here we conduct genome-wide association studies (GWAS) and transcriptome-wide
106 association studies (TWAS) of 35 ID traits. To implement the latter, we apply PrediXcan
107 (Gamazon et al., 2018; Gamazon et al., 2015), which exploits the genetic component of gene
108 expression to probe the molecular basis of disease risk. We combine information across a
109 broad collection of tissues to determine gene-level associations using a multi-tissue approach,
110 which displays markedly improved statistical power over a single-tissue approach (Barbeira et
111 al., 2019; Gamazon et al., 2018; Gamazon et al., 2015). Notably, we identify 70 gene-level
112 associations for 26 of 35 ID traits, i.e., heretofore referred to as ID-associated genes, and
113 conduct replication using the corresponding traits in the UK Biobank and FinnGen consortia
114 data (Bycroft et al., 2018; Locke et al., 2019). The rich resource of genetic information linked to
115 clinical microbiology information that we leverage provides a platform to interrogate the genetic
116 basis of compartment-specific infection and colonization. To gain insights into the phenotypic
117 consequences of ID-associated genes, including adverse outcomes and complications, we
118 perform a phenome-scale scan across hematologic, respiratory, cardiovascular, and neurologic
119 traits. To extend these findings, we use a Mendelian Randomization framework (Lawlor et al.,
120 2008) to conduct causal inference on the effect of a clinical ID trait on an adverse clinical
121 outcome. To elucidate the cellular mechanisms through which host genetic variation influences
122 disease risk, we generate an atlas of gene-level associations with 79 pathogen-induced cellular
123 phenotypes determined by High-throughput Human *in vitro* Susceptibility Testing (Hi-HOST)
124 (Wang et al., 2018) as a discovery and replication platform. The rich genomic resource we
125 generate and the methodology we develop promise to accelerate discoveries on the molecular
126 mechanisms of infection, improve our understanding of adverse outcomes and complications,
127 and enable prioritization of new therapeutic targets.

128

129 RESULTS

130 A schematic diagram illustrating our study design and the reference resource we provide
131 can be found in Figure 1. Here we analyzed 35 clinical ID traits, 79 pathogen-exposure-induced
132 cellular traits, and 197 (cardiovascular, hematologic, neurologic, and respiratory) traits. We
133 performed GWAS and TWAS (Gamazon et al., 2015; Gusev et al., 2016) to investigate the
134 genetic basis of the ID traits and their potential adverse outcomes and complications. We
135 conducted causal inference within a Mendelian Randomization framework (Davey Smith and
136 Hemani, 2014), exploiting genetic instruments for naturally “randomized controlled trials” to
137 evaluate the causality of an observed association between a modifiable exposure or risk factor
138 and a clinical phenotype. We generate a rich resource for understanding the genetic and
139 molecular basis of infection and potential adverse effects and complications.

140

141 *GWAS and TWAS of 35 infectious disease clinical phenotypes implicate broad range of*
142 *molecular mechanisms*

143 We sought to characterize the genetic determinants of 35 ID traits, including many which
144 have never been investigated using a genome-wide approach. First, we performed GWAS of
145 each of these phenotypes using a cohort of 23,294 patients of European ancestry with
146 extensive EHR information from the BioVU (Roden et al., 2008). We identified genome-wide
147 significant associations ($p < 5 \times 10^{-8}$) for 13 ID traits (Figure 2A and Supplementary Table 1). The
148 SNP rs17139584 on chromosome 7 was our most significant association ($p = 1.21 \times 10^{-36}$)
149 across all traits, with bacterial pneumonia. A LocusZoom plot shows several additional genome-
150 wide significant variants in the locus (Figure 2B), in low linkage disequilibrium ($r^2 < 0.20$) with
151 the sentinel variant rs17139584, including variants in the *MET* gene and in *CFTR*. The *MET*
152 gene acts as a receptor to *Listeria monocytogenes* internalin InIB, mediating entry into host
153 cells; interestingly, listeriosis, a bacterial infection caused by this pathogen, can lead to
154 pneumonia (García-Montero et al., 1995). Given the observed associations in the cystic fibrosis

155 gene *CFTR* (~650 Kb downstream of *MET*), we also asked whether the rs17139584 association
156 was driven by cystic fibrosis. Notably, the SNP remained nominally significant, though its
157 significance was substantially reduced, after adjusting for cystic fibrosis status ($p = 0.007$; see
158 Methods) or excluding the cystic fibrosis cases ($p = 0.02$). The LD profile of the genome-wide
159 significant results in this locus (Figure 2B) is consistent with the involvement of multiple gene
160 mechanisms (e.g., *MET* and *CFTR*) underlying bacterial pneumonia risk. The rs17139584
161 association replicated ($p = 5.3 \times 10^{-3}$) in the UK Biobank (Bycroft et al., 2018) (Supplementary
162 Table 2). Eighty percent to ninety percent of patients with cystic fibrosis suffer from respiratory
163 failure due to chronic bacterial infection (with *Pseudomonas aeruginosa*) (Lyczak et al., 2002).
164 Thus, future studies on the role of this locus in lung infection associated with cystic fibrosis may
165 provide germline predictors of this complication; alternatively, the locus may confer susceptibility
166 to lung inflammation, regardless of cystic fibrosis status. Collectively, our analysis shows strong
167 support for allelic heterogeneity, with likely multiple independent variants in the locus
168 contributing to interindividual variability in bacterial pneumonia susceptibility.

169 Additional examples of genome-wide significant associations with other ID traits were
170 identified. For example, rs192146294 on chromosome 1 was significantly associated ($p =$
171 1.23×10^{-9}) with *Staphylococcus* infection. In addition, 10 variants on chromosome 8 were
172 significantly associated ($p < 1.17 \times 10^{-8}$) with Mycoses infection.

173 Next, to improve statistical power, we performed multi-tissue PrediXcan (Barbeira et al.,
174 2019; Gamazon et al., 2018; Gamazon et al., 2015). We constructed an atlas of TWAS
175 associations with these ID traits in separate European and African American ancestry cohorts
176 (Supplementary Data File 1) as a resource to facilitate mechanistic studies. Notably, 70 genes
177 reached experiment-wide or individual ID-trait significance for 26 of the 35 clinical ID traits
178 (Figure 3A and Table 1). Sepsis, the clinical ID trait with the largest sample size in our data
179 (Figure 3B; Phecode 994; number of cases 2,921; number of controls 22,874), was significantly
180 associated ($p = 8.16 \times 10^{-7}$) with *IKZF5* after Bonferroni correction for the number of genes

181 tested. The significant genes (Table 1) were independent of the sentinel variants from the
182 GWAS (Supplementary Table 1), indicating that the gene-based test was identifying additional
183 signals.

184 Our analysis identified previously implicated genes for the specific ID traits but also
185 proposes novel genes and mechanisms. ID-associated genes include *NDUFA4* for intestinal
186 infection, a component of the cytochrome oxidase and regulator of the electron transport chain
187 (Balsa et al., 2012); *AKIRIN2* for candidiasis, an evolutionarily conserved regulator of
188 inflammatory genes in mammalian innate immune cells (Tartey et al., 2015; Tartey et al., 2014);
189 *ZNF577* for viral hepatitis C, a gene previously shown to be significantly hypermethylated in
190 hepatitis C related hepatocellular carcinoma (Revill et al., 2013); and epithelial cell adhesion
191 molecule (*EPCAM*) for tuberculosis, a known marker for differentiating malignant tuberculous
192 pleurisy (Sun et al., 2014), among many others. These examples of ID-associated genes
193 highlight the enormous range of molecular mechanisms that may contribute to susceptibility and
194 complication phenotypes.

195
196 *Replication of gene-level associations with infectious diseases in the UK Biobank and FinnGen*

197 To bolster our genetic findings and show that our results were not driven by biobank-
198 specific confounding, we performed replication analysis for a subset of ID traits available in the
199 independent UK Biobank and FinnGen consortia datasets (see Methods). Individual gene-level
200 replication results are provided in Supplementary Table 3. Notably, the genes associated with
201 intestinal infection ($p < 0.05$, Phencode 008) in BioVU – the ID trait with the largest sample size in
202 BioVU *and* with a replication dataset in the independent FinnGen biobank – showed a
203 significantly greater level of enrichment for gene-level associations with the same trait in
204 FinnGen compared to the remaining set of genes (Figure 3C). Thus, higher significance (i.e.,
205 lower p-value) was observed in FinnGen for the intestinal infection associated genes identified
206 in BioVU, which included the top association *NDUFA4* (discovery $p = 1.83 \times 10^{-9}$, replication $p =$

207 0.044). These results illustrate the value of exploiting large-scale biobank resources for genetic
208 studies of ID traits- despite well-known caveats (Ko and Urban, 2013; Power et al., 2017).

209

210 *Tissue expression profile of infectious disease associated genes suggests tissue-dependent*
211 *mechanisms*

212 The ID-associated genes tend to be less tissue-specific (i.e., more ubiquitously
213 expressed) than the remaining genes (Figure S1A, Mann Whitney U test on the τ statistic, $p =$
214 7.5×10^{-4}), possibly reflecting the multi-tissue PrediXcan approach we implemented, which
215 prioritizes genes with multi-tissue support to improve statistical power, but also the genes'
216 pleiotropic potential. We hypothesized that tissue expression profiling of ID-associated genes
217 can provide additional insights into disease etiologies and mechanisms. For example, the
218 intestinal infection associated gene *NDUFA4* is expressed in a broad set of tissues, including
219 the alimentary canal, but displays relatively low expression in whole blood (Figure S1B). In
220 addition, *TOR4A*, the most significant association with bacterial pneumonia (Table 1), is most
221 abundantly expressed in lung, consistent with the tissue of pathology, but also in spleen (Figure
222 S1C), whose rupture is a lethal complication of the disease (Domingo et al., 1996; Gerstein et
223 al., 1967). These examples illustrate the diversity of tissue-dependent mechanisms that may
224 contribute in complex and dynamic ways to interindividual variability in ID susceptibility and
225 progression. We therefore provide a resource of single-tissue gene-level associations with the
226 ID traits to facilitate molecular or clinical follow-up studies.

227

228 *Genetic overlap reveals host gene expression programs and common pathways as targets for*
229 *pathogenicity*

230 We hypothesized that ID-associated genes implicate shared functions and pathways,
231 which may reflect common targeted host transcriptional programs. Among the 70 gene-level
232 associations with the 35 clinical ID traits, 40 proteins are post-translationally modified by

233 phosphorylation (Supplementary Table 4), a significant enrichment (Benjamini-Hochberg
234 adjusted $p < 0.10$ on DAVID annotations (Huang da et al., 2009)) relative to the rest of the
235 genome, indicating that phosphoproteomic profiling can shed substantial light on activated host
236 factors and perturbed signal transduction pathways during infection (Soderholm et al., 2016;
237 Stahl et al., 2013). In addition, 16 proteins are acetylated, consistent with emerging evidence
238 supporting this mechanism in the host antiviral response (Murray et al., 2018) (Supplementary
239 Table 5). These data identify specific molecular mechanisms across ID traits with critical
240 regulatory roles (e.g., protein modifications) in host response among the ID-associated genes.

241 We tested the hypothesis that distinct infectious agents exploit common pathways to find
242 a compatible intracellular niche in the host, potentially implicating shared genetic risk factors.
243 Notably, 64 of the 70 ID-associated genes (Table 1) were nominally associated ($p < 0.05$) with
244 multiple ID traits (Supplementary Table 6). These genes warrant further functional study as
245 broadly exploited mechanisms targeted by pathogens or as broadly critical to pathogen-elicited
246 immune response. Gene Set Enrichment Analysis (GSEA) of these genes implicated a number
247 of significant ($FDR < 0.05$) gene sets (Figure 4A), including those involved in actin-based
248 processes and cytoskeletal protein binding, processes previously demonstrated to mediate host
249 response to pathogen infection (Taylor et al., 2011). Since diverse bacterial and viral pathogens
250 target host regulators that control the cytoskeleton (which plays a key role in the biology of
251 infection) or modify actin in order to increase virulence, intracellular motility, or intercellular
252 spread (Aktories and Barbieri, 2005; Yu et al., 2011; Zahm et al., 2013), these results
253 reassuringly lend support to the involvement of the genes in infectious pathogenesis.

254 Notably, we identified an enrichment ($FDR = 9.68 \times 10^{-3}$) for a highly conserved motif
255 ("TCCCRNNRTGC"), within 4 kb of transcription start site (TSS) of multi-ID associated genes
256 (Figure 4A-B), that does not match any known transcription factor binding site (Xie et al., 2005)
257 and may be pivotal for host-pathogen interaction for the diversity of infectious agents included in
258 our study. In addition, we found that several of the multi-ID associated genes (with the

259 sequence motif near the TSS) have been observed in host-pathogen protein complexes (by
260 both coimmunoprecipitation and affinity chromatography approaches) for the specific pathogens
261 responsible for the ID traits (Ammari et al., 2016). See Supplementary Data File 2 for complete
262 list of host-pathogen interactions for these genes/proteins. One example is *CDK5*, a gene
263 significantly associated with Gram-positive septicemia (Table 1) and nominally associated with
264 multiple ID traits, including herpes simplex. CDK5 is activated by p35, whose cleaved form p25
265 results in subcellular relocation of CDK5. The CDK5-p25 complex regulates inflammation
266 (Na et al., 2015) (whose large-scale disruption is characteristic of septicemia) and induces
267 cytoskeletal disruption in neurons (Patrick et al., 1999) (where the herpes virus is responsible
268 for lifelong latent infection). The A and B chains of the CDK5-p25 complex (Figure 4C for
269 structure diagram (Tarricone et al., 2001)) are required for cytoskeletal protein binding (CDK5),
270 whereas the D and E chains (p25) are involved in actin regulation and kinase function, all
271 molecular processes implicated in our pathway analysis. Intriguingly, blocking CDK5 can have a
272 substantial impact on the outcome of inflammatory diseases including sepsis (Pfänder et al.,
273 2019), enhancing the anti-inflammatory potential of immunosuppressive treatments, and has
274 been shown to attenuate herpes virus replication (Man et al., 2019), suggesting that modulation
275 of this complex is important for viral pathogenesis.

276 CDK5 is also altered by several other viruses, identified using unbiased mass
277 spectrometry analysis (Davis et al., 2015) (Figure 4D), indicating a broadly exploited mechanism
278 (across pathogens) that is consistent with the gene's multi-ID genetic associations in our TWAS
279 data (Figure 4D). The CDK5-interaction proteins include: 1) M2_134A1 (matrix protein 2,
280 influenza A virus), a component of the proton-selective ion channel required for viral genome
281 release during cellular entry and is targeted by the anti-viral drug amantadine (Hay et al., 1985);
282 2) VE7_HP16, a component of human papillomavirus (HPV) required for cellular
283 transformation and trans-activation through disassembly of E2F1 transcription factor from RB1
284 leading to impaired production of type I interferons (Barnard et al., 2000; Chellappan et al.,

285 1992; Phelps et al., 1988); 3) VE7_HP31, which has been shown to engage histone
286 deacetylases 1 and 2 to promote HPV31 genome maintenance (Longworth and Laimins, 2004);
287 4) VCYCL_HHV8P (cyclin homolog within the human herpesvirus 8 genome), which has been
288 shown to control cell cycle through CDK6 and induce apoptosis through Bcl2 (Duro et al., 1999;
289 Ojala et al., 1999; Ojala et al., 2000); and 5) F5HC81_HHV8, predicted to act as a viral cyclin
290 homolog. Overall, these data underscore the evolutionary strategies that pathogens have
291 evolved to promote infection, including the hijacking of the host transcriptional machinery and
292 the biochemical alterations of the host proteome.

293

294 *Serology and culture data reveal insights into clinical infection and pathogen colonization*

295 We exploited extensive clinical microbiological laboratory analysis of blood (Figure 5A),
296 bronchoalveolar lavage, sputum, sinus/nasopharyngeal, and tracheal cultures for bacterial and
297 fungal pathogen genus identification (Supplemental Figure 2A-F), as well as respiratory viral
298 genus identification (Supplemental Figure 5G) (see Methods) to evaluate phenotype resolution
299 and algorithm. For example, we found that *Staphylococcus* infection (Phecode = 041.1)
300 performed well in classifying *Staphylococcus aureus* infection based on blood culture data. The
301 area under the Receiver Operating Characteristic (ROC) curve was 0.938 (Figure 5B) with
302 standard error of 0.008 generated from bootstrapping (see Methods). The area under the curve
303 (AUC) quantifies the probability that the Phecode classifier ranks a randomly chosen positive
304 instance of *Staphylococcus aureus* infection in blood higher than a randomly chosen negative
305 one. In comparison, the first principal component (PC) in our European ancestry samples
306 showed AUC of 0.514 (Figure 5B) while sex and age performed even more poorly (AUC \approx
307 0.50). We then tested a logistic model with the Phecode classifier, age, sex, and the first 5 PCs
308 in the model. The Phecode classifier was significantly associated ($p < 2.2 \times 10^{-16}$) after
309 conditioning on the remaining covariates. The fitted value from the joint model consisting of the

310 remaining covariates showed AUC of 0.568 (Figure 5B). Collectively, culture data for improved
311 resolution of clinical infection and pathogen colonization provide validation of our approach.

312

313 *Phenome scan of clinical ID-associated genes identifies adverse outcomes and complications*

314 Electronic Health Records (EHR) linked to genetic data may reveal insights into
315 associated clinical sequelae (Bastarache et al., 2018; Denny et al., 2013; Unlu et al., 2020). To
316 assess the phenomic impact of ID-associated genes (Table 1), we performed a phenome-scale
317 scan across 197 hematologic, respiratory, cardiovascular, and neurologic traits available in
318 BioVU (Figure 6A and Supplementary Data File 3). Correcting for total number of genes and
319 phenotypes tested, we identified four gene-phenotype pairs reaching experiment-wide
320 significance: 1) *WFDC12*, our most significant ($p = 4.23 \times 10^{-6}$) association with meningitis and a
321 known anti-bacterial gene (Hagiwara et al., 2003), is also associated with cerebral edema and
322 compression of brain ($p = 1.35 \times 10^{-6}$), a feared clinical complication of meningitis (Niemöller and
323 Täuber, 1989); 2) *TM7SF3*, the most significant gene with Gram-negative sepsis ($p = 1.37 \times 10^{-6}$),
324 is also associated with acidosis ($p = 1.95 \times 10^{-6}$), a known metabolic derangement associated
325 with severe sepsis (Suetrong and Walley, 2016), and a gene known to play a role in cell stress
326 and the unfolded protein response (Isaac et al., 2017); 3) *TXLNB*, the most significant gene
327 associated with viral warts and human papillomavirus infection ($p = 4.35 \times 10^{-6}$), is also
328 associated with abnormal involuntary movements, $p = 1.39 \times 10^{-6}$; and 4) *RAD18*, the most
329 significant gene associated with Streptococcus infection ($p = 2.01 \times 10^{-6}$), is also associated with
330 anemia in neoplastic disease ($p = 3.10 \times 10^{-6}$). Thus, coupling genetic analysis to EHR data with
331 their characteristic breadth of clinical traits offers the possibility of determining the phenotypic
332 consequences of ID-associated genes, including known (in the case of *WFDC12* and *TM7SF3*)
333 potentially adverse health outcomes and complications.

334

335 *Mendelian Randomization provides causal support for the effect of infectious disease trait on*
336 *identified adverse phenotypic outcomes/complications*

337 Since our gene-level associations with clinical ID diagnoses implicated known adverse
338 complications, we sought to explicitly evaluate the causal relation between the ID traits and the
339 adverse outcomes/complications. We utilized the Mendelian Randomization paradigm (Lawlor
340 et al., 2008) (Figure 6B), which exploits genetic instruments to make causal inferences in
341 observational data, in effect, performing randomized controlled trials to evaluate the causal
342 effect of “exposure” (i.e., ID trait) on “outcome” (e.g., the complication). Specifically, we
343 conducted multiple-instrumental-variable causal inference using GWAS (Davey Smith and
344 Hemani, 2014) and PrediXcan summary results. First, we used independent SNPs ($r^2 = 0.01$)
345 that pass a certain threshold for significance with the ID trait ($p < 1.0 \times 10^{-5}$) as genetic
346 instruments. To control for horizontal pleiotropy and account for the presence of invalid genetic
347 instruments, we utilized MR-Egger regression and weighted-median MR (see Methods)
348 (Bowden et al., 2015; Bowden et al., 2016).

349 Here, we performed Mendelian Randomization on the ID trait and a complication trait
350 identified through the unbiased phenome scan. This analysis yields causal support for the effect
351 of 1) Gram-negative sepsis on acidosis (Figure 6C, weighted-median estimator $p = 2.0 \times 10^{-7}$);
352 and 2) meningitis on cerebral edema and compression of brain (Figure 5C, weighted-median
353 estimator $p = 2.7 \times 10^{-3}$). Our resource establishes a framework to elucidate the genetic
354 component of an ID trait and its impact on the human disease phenome, enabling causal
355 inference on the effect of an ID trait on potential complications.

356

357

358 *TWAS of 79 pathogen-exposure induced cellular traits highlights cellular mechanisms and*
359 *enables validation of ID gene-level associations*

360 Elucidating how the genes influence infection-related cellular trait variation may provide
361 a mechanistic link to ID susceptibility. We thus performed TWAS of 79 pathogen-induced
362 cellular traits – including infectivity and replication, cytokine levels, and host cell death, among
363 others (Wang et al., 2018) (Supplementary Data File 4). Across all cellular traits, we identified
364 38 gene-level associations reaching trait-level significance ($p < 2.87 \times 10^{-6}$, correcting for number
365 of statistical tests; Figure 7A). In addition, we identified significantly more replicated SNP
366 associations than expected by chance (binomial test, $p < 0.05$) across all ID traits
367 (Supplementary Data File 4) and ID traits which map directly to cellular phenotypes
368 (Supplementary Data File 5).

369 Integration of EHR data into Hi-HOST (Wang et al., 2018) may enable replication of
370 gene-level associations with a clinical ID trait. Indeed, we observed a marked enrichment for
371 genes associated with direct *Staphylococcus* toxin exposure cellular response in Hi-HOST
372 among the human Gram-positive septicemia associated genes from BioVU (see Supplementary
373 Data File 4 for genes with FDR < 0.05) (Figure 7B). In addition, integration of EHR data into Hi-
374 HOST may improve the signal-to-noise ratio in Hi-HOST TWAS data. Indeed, the top 300 genes
375 nominally associated ($p < 0.016$) with *Staphylococcus* infection (Phecode 041.1) in BioVU
376 departed from null expectation for their associations with *Staphylococcus* toxin exposure in Hi-
377 HOST compared to the full set of genes, which did not (Figure 7C), as perhaps expected due to
378 the modest sample size. Collectively, these results demonstrate that integrating the EHR-
379 derived TWAS results into TWAS of the cellular trait can greatly improve identification of
380 potentially relevant pathogenic mechanisms.

381

382 *Phenome scan of TWAS findings from Hi-HOST*

383 To identify potential adverse effects of direct pathogen exposure, we performed a
384 phenome-scan across the 197 cardiovascular, hematologic, neurologic, and respiratory traits as
385 described above. Our top gene-phenotype pairs include: 1) *FAM171B*, our most significant

386 association with interleukin 13 (IL-13) levels is also associated with alveolar and parietoalveolar
387 pneumonopathy ($p = 4.04 \times 10^{-5}$), a phenotype known to be modulated by IL-13 dependent
388 signaling (Zheng et al., 2008); 2) *OSBPL10*, the most significant gene associated with cell death
389 caused by *Salmonella enterica serovar Typhimurium*, is also associated with intracerebral
390 hemorrhage ($p = 4.99 \times 10^{-5}$), a known complication of *S. Typhimurium* endocarditis (Gómez-
391 Moreno et al., 2000). These data highlight the utility of joint genetic analysis of pathogen-
392 exposure-induced phenotypes and clinical ID traits to gain insights into the molecular and
393 cellular basis of complications and adverse outcomes. However, more definitive conclusions will
394 require larger sample sizes and functional studies.

395

396 **DISCUSSION**

397 ID susceptibility is a complex interplay between host genetic variation and pathogen-
398 exposure induced mechanisms. While GWAS has begun to identify population-specific loci
399 conferring ID risk (Tian et al., 2017), the underlying function of identified variants, predominantly
400 in non-coding regulatory regions, remains poorly understood. Molecular characterization of
401 infectious processes has been, in general, agnostic to the genetic architecture of clinical
402 infection. Although pathogen exposure is requisite to display clinical ID traits, the role of host
403 genetic variation remains largely unexplored.

404 Our study provides a reference atlas of genetic variants and genetically-determined
405 expression traits associated with 35 clinical ID traits from BioVU. We identified 70 gene-level
406 associations, with replication for a subset of ID traits in the UK Biobank and FinnGen. A
407 phenome scan across 197 hematologic, respiratory, cardiovascular, and neurologic traits
408 proposes a molecular basis for the link between certain ID traits and outcomes. Using
409 Mendelian Randomization, we determined the ID traits which, as exposure, show significant
410 causal effect on outcomes. Finally, we developed a TWAS catalog of 79 pathogen-exposure

411 induced cellular traits (Hi-HOST) in a broad collection of tissues, which provides a platform to
412 interrogate mediating cellular and molecular mechanisms.

413 Genetic predisposition to ID onset and progression is likely to be complex (Casanova,
414 2015a). Monogenic mechanisms conferring ID risk have been proposed, but these mechanisms
415 are unlikely to explain the broad contribution of host genetic influence on ID risk (Casanova,
416 2015b). Thus, a function-centric methodology is necessary to disentangle potentially causal
417 pathways. Our approach builds on PrediXcan, which estimates the genetically-determined
418 component of gene expression (Gamazon et al., 2015). The genetic component of gene
419 expression can then be tested for association with the trait, enabling insights into potential
420 pathogenic mechanisms (Gamazon et al., 2019) and novel therapeutic strategies (So et al.,
421 2017).

422 Our study identified genes with diverse functions, including roles in mitochondrial
423 bioenergetics (Balsa et al., 2012; El-Bacha and Da Poian, 2013), regulation of cell death (Labbé
424 and Saleh, 2008), and of course links to host immune response (Brouwer et al., 2019; Liang et
425 al., 2019; Pan et al., 2017; Saitoh et al., 2009; Sharfe et al., 1997; Tsuboi and Meerloo, 2007;
426 Walenna et al., 2018; Willis et al., 2009; Yu et al., 2017; Zhang et al., 2015). These diverse
427 functions may therefore contribute to pleiotropic effects on clinical outcomes and complications.

428 In addition, we identified genes implicated in Mendelian diseases, for which susceptibility
429 to infection is a predominant feature, including *WIPF1* (OMIM #614493; recurrent infections and
430 reduced natural killer cell activity (Lanzi et al., 2012)), *IL2RA* (OMIM #606367; recurrent
431 bacterial infections, recurrent viral infections, and recurrent fungal infections (Sharfe et al.,
432 1997)), and *TBK1* (OMIM #617900; herpes simplex encephalitis (HSE), acute infection, and
433 episodic HSE (Herman et al., 2012)). These examples show that the identified genes may also
434 confer predisposition, with near-complete penetrance, to an infectious disease related trait
435 displaying true Mendelian segregation.

436 Enrichment analysis of 64 of the 70 ID-associated genes with nominal support for
437 associations with other clinical ID traits identified modulation of the actin cytoskeleton as a
438 potential shared mechanism of host susceptibility to infection (Figure 4). While manipulation of
439 the actin cytoskeleton by pathogens is hardly a new concept, our study identified specific host
440 genetic variation in actin regulatory genes that is potentially causative of clinical ID
441 manifestations. In addition to pathogen interaction with the cytoskeletal transport machinery,
442 efficient exploitation of host gene expression program is crucial for successful invasion and
443 colonization, and here we mapped several pathogenicity-relevant targets. Notably, we observed
444 a significant enrichment for a highly conserved sequence motif, within 4 kb of a multi-ID-
445 associated gene's TSS, that is not a known transcription factor binding site. The motif's
446 presence near multi-ID associated genes suggests a broad regulatory role in host-pathogen
447 interaction, involving the diversity of pathogens examined here, towards successful
448 reprogramming of host gene expression. Furthermore, we identified a significant enrichment for
449 phosphorylated host proteins, suggesting the value of global phosphoproteomic profiling, which
450 has recently been used to prioritize pharmacological targets for the novel SARS-CoV-2 virus
451 (Bouhaddou et al., 2020). These data provide several potential avenues by which host
452 susceptibility can be breached by a pathogen's requirement to maintain a niche through
453 manipulation of host cellular machinery.

454 To obtain additional support for our gene-level associations, we leveraged two genomic
455 resources with rich phenotypic information (UK Biobank (Bycroft et al., 2018) and FinnGen
456 (Locke et al., 2019)). These data will prove increasingly useful to characterizing the genetic
457 basis of the ID-associated adverse outcomes and complications. Despite the caveats for the
458 use of EHR in genetic analyses of ID traits (Ko and Urban, 2013; Power et al., 2017), the
459 growing availability of such independent datasets will facilitate identification of robust genetic
460 associations. Perhaps more importantly, the breadth of clinical phenotypes in these EHR

461 datasets should enable identification of associated adverse outcomes and complications for the
462 ID-associated genes.

463 The primary challenges in conducting GWAS of ID traits include phenotype definition
464 and case-control misclassification. Obstacles to accurate phenotype definition include the
465 requirement of specialized laboratory testing to identify specific pathogens and administration of
466 prophylactic therapeutics complicating identification of potentially causative pathogens.
467 Seropositivity may result from the complex genetic properties of the pathogen and the particular
468 mechanisms governing host-pathogen interaction. However, seropositivity may not indicate
469 clinical manifestations of the disease. On the other hand, seronegativity may imply lack of
470 exposure to the pathogen, the absence of infection even in the presence of exposure, or host
471 resistance to infection. Anchoring the analysis to host genetic information (as in our use of
472 genetically-determined expression) and replication of discovered associations may address
473 some aspects of this challenge. Here we exploit an extensive resource of culture data (for
474 identification of pathogens from clinical specimens) linked to whole-genome genetic information
475 to provide additional support to our gene-level associations. Future studies may implement
476 more complex GWAS models, including incorporating the pathogen genome.

477 Mendelian Randomization provides a framework to perform causal inference on the
478 effect of the exposure on the outcome (Davey Smith and Hemani, 2014; Lawlor et al., 2008).
479 We leveraged a summary statistics based approach to test the causal effect of an ID trait on
480 potential adverse outcomes, using genetic instruments. Mendelian Randomization requires
481 three assumptions: 1) the genetic instrument is associated with exposure (i.e., ID trait); 2) the
482 genetic instrument is associated with the outcome (i.e., adverse outcome or complication) only
483 through the exposure of interest; and 3) the genetic instrument is affecting the outcome
484 independent of other factors (i.e., confounders). Violations of these assumptions can have
485 critical implications for the interpretation of the results. Thus, several approaches have been
486 developed that are robust to these violations. In the case of ID traits, a methodology that

487 distinguishes causality from comorbidity is critical. While many phenotypes are highly comorbid
488 and suspected to have a causal relationship (e.g., smoking and depression/anxiety), Mendelian
489 Randomization does not necessarily support the causal hypothesis (Taylor et al., 2014).
490 Furthermore, since RCTs cannot be ethically conducted for ID traits and adverse outcomes, the
491 methodology offers an approach for elucidating the role of an infection phenotype or pathogen
492 exposure in disease causation using an observational study design. Here, we found strong
493 causal support for the effect of certain clinical ID traits on potential adverse complications
494 identified through a phenome scan of the ID-associated genes: 1) meningitis - cerebral edema
495 and compression of brain; and 2) Gram-negative sepsis - acidosis. These data indicate that
496 genetic risk factors for select adverse outcomes and complications exert their phenotypic effect
497 through the relevant ID traits.

498 To enable investigations into mediating cellular and molecular traits for the ID-associated
499 genes, we provide a functional genomics resource built on a high-throughput *in vitro* pathogen
500 infection screen (Hi-HOST) (Wang et al., 2018). Integration of EHR data into Hi-HOST facilitates
501 replication of gene-level associations with clinical ID traits and greatly improves the signal-to-
502 noise ratio. This discovery and replication platform, encompassing human phenomics and
503 cellular microbiology, provides a high-throughput approach to linking host cellular processes to
504 clinical ID traits and adverse outcomes.

505 Although additional mechanistic studies are warranted, our study lays the foundation for
506 anchoring targeted molecular studies in human genetic variation. Elucidation of host
507 mechanisms exploited by pathogens requires multi-disciplinary approaches. Here, we show the
508 broader role of host genetic variation, implicating diverse disease mechanisms. Our study
509 generates a rich resource and a genetics-anchored methodology to facilitate investigations of
510 ID-associated clinical outcomes and complications, with important implications for the
511 development of preventive strategies and more effective therapeutics. Causal inference on the
512 clinical ID traits and potential complications promises to expand our understanding of the

513 molecular basis for the link and, crucially, enable prediction and prevention of serious adverse

514 events.

515

516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553

REFERENCES

- (2015). Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science (New York, NY)* *348*, 648-660.
- Aktorjies, K., and Barbieri, J.T. (2005). Bacterial cytotoxins: targeting eukaryotic switches. *Nat Rev Microbiol* *3*, 397-410.
- Ammari, M.G., Gresham, C.R., McCarthy, F.M., and Nanduri, B. (2016). HPIDB 2.0: a curated database for host-pathogen interactions. *Database : the journal of biological databases and curation* *2016*.
- Auton, A., Brooks, L.D., Durbin, R.M., Garrison, E.P., Kang, H.M., Korbel, J.O., Marchini, J.L., McCarthy, S., McVean, G.A., and Abecasis, G.R. (2015). A global reference for human genetic variation. *Nature* *526*, 68-74.
- Balsa, E., Marco, R., Perales-Clemente, E., Szklarczyk, R., Calvo, E., Landázuri, M.O., and Enríquez, J.A. (2012). NDUFA4 is a subunit of complex IV of the mammalian electron transport chain. *Cell Metab* *16*, 378-386.
- Barbeira, A.N., Dickinson, S.P., Bonazzola, R., Zheng, J., Wheeler, H.E., Torres, J.M., Torstenson, E.S., Shah, K.P., Garcia, T., Edwards, T.L., *et al.* (2018). Exploring the phenotypic consequences of tissue specific gene expression variation inferred from GWAS summary statistics. *Nat Commun* *9*, 1825.
- Barbeira, A.N., Pividori, M., Zheng, J., Wheeler, H.E., Nicolae, D.L., and Im, H.K. (2019). Integrating predicted transcriptome from multiple tissues improves association detection. *PLoS Genet* *15*, e1007889.
- Barnard, P., Payne, E., and McMillan, N.A. (2000). The human papillomavirus E7 protein is able to inhibit the antiviral and anti-growth functions of interferon-alpha. *Virology* *277*, 411-419.
- Bastarache, L., Hughey, J.J., Hebring, S., Marlo, J., Zhao, W., Ho, W.T., Van Driest, S.L., McGregor, T.L., Mosley, J.D., Wells, Q.S., *et al.* (2018). Phenotype risk scores identify patients with unrecognized Mendelian disease patterns. *Science (New York, NY)* *359*, 1233-1239.
- Battle, A., Brown, C.D., Engelhardt, B.E., and Montgomery, S.B. (2017). Genetic effects on gene expression across human tissues. *Nature* *550*, 204-213.
- Bouhaddou, M., Memon, D., Meyer, B., White, K.M., Rezelj, V.V., Marrero, M.C., Polacco, B.J., Melnyk, J.E., Ulferts, S., Kaake, R.M., *et al.* (2020). The Global Phosphorylation Landscape of SARS-CoV-2 Infection. *Cell*.
- Bowden, J., Davey Smith, G., and Burgess, S. (2015). Mendelian randomization with invalid instruments: effect estimation and bias detection through Egger regression. *Int J Epidemiol* *44*, 512-525.
- Bowden, J., Davey Smith, G., Haycock, P.C., and Burgess, S. (2016). Consistent Estimation in Mendelian Randomization with Some Invalid Instruments Using a Weighted Median Estimator. *Genetic epidemiology* *40*, 304-314.

- 554 Brouwer, W.P., Chan, H.L., Lampertico, P., Hou, J., Tangkijvanich, P., Reesink, H.W., Zhang,
555 W., Mangia, A., Tanwandee, T., Montalto, G., *et al.* (2019). Genome Wide Association Study
556 Identifies Genetic Variants Associated With Early And Sustained Response To (Peg)Interferon
557 In Chronic Hepatitis B Patients: The GIANT-B Study. *Clinical infectious diseases : an official*
558 *publication of the Infectious Diseases Society of America.*
- 559 Bulik-Sullivan, B., Finucane, H.K., Anttila, V., Gusev, A., Day, F.R., Loh, P.R., Duncan, L.,
560 Perry, J.R., Patterson, N., Robinson, E.B., *et al.* (2015). An atlas of genetic correlations across
561 human diseases and traits. *Nat Genet* 47, 1236-1241.
- 562 Burgner, D., Jamieson, S.E., and Blackwell, J.M. (2006). Genetic susceptibility to infectious
563 diseases: big is beautiful, but will bigger be even better? *Lancet Infect Dis* 6, 653-663.
- 564 Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L.T., Sharp, K., Motyer, A., Vukcevic, D.,
565 Delaneau, O., O'Connell, J., *et al.* (2018). The UK Biobank resource with deep phenotyping and
566 genomic data. *Nature* 562, 203-209.
- 567 Casanova, J.L. (2015a). Human genetic basis of interindividual variability in the course of
568 infection. *Proc Natl Acad Sci U S A* 112, E7118-7127.
- 569 Casanova, J.L. (2015b). Severe infectious diseases of childhood as monogenic inborn errors of
570 immunity. *Proc Natl Acad Sci U S A* 112, E7128-7137.
- 571 Chellappan, S., Kraus, V.B., Kroger, B., Munger, K., Howley, P.M., Phelps, W.C., and Nevins,
572 J.R. (1992). Adenovirus E1A, simian virus 40 tumor antigen, and human papillomavirus E7
573 protein share the capacity to disrupt the interaction between transcription factor E2F and the
574 retinoblastoma gene product. *Proc Natl Acad Sci U S A* 89, 4549-4553.
- 575 Consortium, G. (2013). The Genotype-Tissue Expression (GTEx) project. *Nat Genet* 45, 580-
576 585.
- 577 Cotsapas, C., Voight, B.F., Rossin, E., Lage, K., Neale, B.M., Wallace, C., Abecasis, G.R.,
578 Barrett, J.C., Behrens, T., Cho, J., *et al.* (2011). Pervasive sharing of genetic effects in
579 autoimmune disease. *PLoS Genet* 7, e1002254.
- 580 Datsenko, K.A., and Wanner, B.L. (2000). One-step inactivation of chromosomal genes in
581 *Escherichia coli* K-12 using PCR products. *Proc Natl Acad Sci U S A* 97, 6640-6645.
- 582 Davey Smith, G., and Hemani, G. (2014). Mendelian randomization: genetic anchors for causal
583 inference in epidemiological studies. *Human molecular genetics* 23, R89-98.
- 584 Davis, Z.H., Verschueren, E., Jang, G.M., Kleffman, K., Johnson, J.R., Park, J., Von Dollen, J.,
585 Maher, M.C., Johnson, T., Newton, W., *et al.* (2015). Global mapping of herpesvirus-host
586 protein complexes reveals a transcription strategy for late genes. *Molecular cell* 57, 349-360.
- 587 de Bakker, P.I., and Telenti, A. (2010). Infectious diseases not immune to genome-wide
588 association. *Nat Genet* 42, 731-732.
- 589 Denny, J.C., Bastarache, L., Ritchie, M.D., Carroll, R.J., Zink, R., Mosley, J.D., Field, J.R.,
590 Pulley, J.M., Ramirez, A.H., Bowton, E., *et al.* (2013). Systematic comparison of phenome-wide

- 591 association study of electronic medical record data and genome-wide association study data.
592 *Nature biotechnology* 31, 1102-1110.
- 593 Denny, J.C., Ritchie, M.D., Basford, M.A., Pulley, J.M., Bastarache, L., Brown-Gentry, K., Wang,
594 D., Masys, D.R., Roden, D.M., and Crawford, D.C. (2010). PheWAS: demonstrating the
595 feasibility of a phenome-wide scan to discover gene-disease associations. *Bioinformatics* 26,
596 1205-1210.
- 597 Derks, E.M., Zwinderman, A.H., and Gamazon, E.R. (2017). The Relation Between Inflation in
598 Type-I and Type-II Error Rate and Population Divergence in Genome-Wide Association Analysis
599 of Multi-Ethnic Populations. *Behavior genetics* 47, 360-368.
- 600 Domingo, P., Rodriguez, P., Lopez-Contreras, J., Rebasa, P., Mota, S., and Matias-Guiu, X.
601 (1996). Spontaneous rupture of the spleen associated with pneumonia. *European journal of*
602 *clinical microbiology & infectious diseases : official publication of the European Society of*
603 *Clinical Microbiology* 15, 733-736.
- 604 Duro, D., Schulze, A., Vogt, B., Bartek, J., Mitnacht, S., and Jansen, D.r.P. (1999). Activation of
605 cyclin A gene expression by the cyclin encoded by human herpesvirus-8. *J Gen Virol* 80 (Pt 3),
606 549-555.
- 607 Efron, B. (1979). Bootstrap Methods: Another Look at the Jackknife. *Ann Statist* 7, 1-26.
- 608 El-Bacha, T., and Da Poian, A.T. (2013). Virus-induced changes in mitochondrial bioenergetics
609 as potential targets for therapy. *The international journal of biochemistry & cell biology* 45, 41-
610 46.
- 611 Gamazon, E.R., Segre, A.V., van de Bunt, M., Wen, X., Xi, H.S., Hormozdiari, F., Ongen, H.,
612 Konkashbaev, A., Derks, E.M., Aguet, F., *et al.* (2018). Using an atlas of gene regulation across
613 44 human tissues to inform complex disease- and trait-associated variation. *Nat Genet* 50, 956-
614 967.
- 615 Gamazon, E.R., Wheeler, H.E., Shah, K.P., Mozaffari, S.V., Aquino-Michaels, K., Carroll, R.J.,
616 Eyer, A.E., Denny, J.C., Nicolae, D.L., Cox, N.J., *et al.* (2015). A gene-based association
617 method for mapping traits using reference transcriptome data. *Nat Genet* 47, 1091-1098.
- 618 Gamazon, E.R., Zwinderman, A.H., Cox, N.J., Denys, D., and Derks, E.M. (2019). Multi-tissue
619 transcriptome analyses identify genetic mechanisms underlying neuropsychiatric traits. *Nat*
620 *Genet* 51, 933-940.
- 621 García-Montero, M., Rodríguez-García, J.L., Calvo, P., González, J.M., Fernández-Garrido, M.,
622 Loza, E., and Serrano, M. (1995). Pneumonia caused by *Listeria monocytogenes*. *Respiration*
623 62, 107-109.
- 624 Gerstein, A.R., Riegel, N., and Dennis, M. (1967). Ruptured Spleen Simulating Pneumonia.
625 *JAMA* 199, 589-589.
- 626 Gómez-Moreno, J., Moar, C., Román, F., Pérez-Maestu, R., and López de Letona, J.M. (2000).
627 *Salmonella* endocarditis presenting as cerebral hemorrhage. *Eur J Intern Med* 11, 96-97.

- 628 Gusev, A., Ko, A., Shi, H., Bhatia, G., Chung, W., Penninx, B.W.J.H., Jansen, R., de Geus,
629 E.J.C., Boomsma, D.I., Wright, F.A., *et al.* (2016). Integrative approaches for large-scale
630 transcriptome-wide association studies. *Nature genetics* 48, 245-252.
- 631 Hagiwara, K., Kikuchi, T., Endo, Y., Huqun, Usui, K., Takahashi, M., Shibata, N., Kusakabe, T.,
632 Xin, H., Hoshi, S., *et al.* (2003). Mouse SWAM1 and SWAM2 are antibacterial proteins
633 composed of a single whey acidic protein motif. *J Immunol* 170, 1973-1979.
- 634 Hay, A.J., Wolstenholme, A.J., Skehel, J.J., and Smith, M.H. (1985). The molecular basis of the
635 specific anti-influenza action of amantadine. *Embo j* 4, 3021-3024.
- 636 Herman, M., Ciancanelli, M., Ou, Y.H., Lorenzo, L., Klaudel-Dreszler, M., Pauwels, E., Sancho-
637 Shimizu, V., Pérez de Diego, R., Abhyankar, A., Israelsson, E., *et al.* (2012). Heterozygous
638 TBK1 mutations impair TLR3 immunity and underlie herpes simplex encephalitis of childhood.
639 *The Journal of experimental medicine* 209, 1567-1582.
- 640 Herndon, C.N., and Jennings, R.G. (1951). A twin-family study of susceptibility to poliomyelitis.
641 *Am J Hum Genet* 3, 17-46.
- 642 Hohler, T., Reuss, E., Evers, N., Dietrich, E., Rittner, C., Freitag, C.M., Vollmar, J., Schneider,
643 P.M., and Fimmers, R. (2002). Differential genetic determination of immune responsiveness to
644 hepatitis B surface antigen and to hepatitis A virus: a vaccination study in twins. *Lancet* 360,
645 991-995.
- 646 Huang da, W., Sherman, B.T., and Lempicki, R.A. (2009). Systematic and integrative analysis of
647 large gene lists using DAVID bioinformatics resources. *Nat Protoc* 4, 44-57.
- 648 Isaac, R., Goldstein, I., Furth, N., Zilber, N., Streim, S., Boura-Halfon, S., Elhanany, E., Rotter,
649 V., Oren, M., and Zick, Y. (2017). TM7SF3, a novel p53-regulated homeostatic factor,
650 attenuates cellular stress and the subsequent induction of the unfolded protein response. *Cell*
651 *Death Differ* 24, 132-143.
- 652 Jordan, D.M., Verbanck, M., and Do, R. (2019). HOPS: a quantitative score reveals pervasive
653 horizontal pleiotropy in human genetic variation is driven by extreme polygenicity of human traits
654 and diseases. *Genome Biol* 20, 222.
- 655 Ko, D.C., Gamazon, E.R., Shukla, K.P., Pfuetzner, R.A., Whittington, D., Holden, T.D.,
656 Brittnacher, M.J., Fong, C., Radey, M., Ogohara, C., *et al.* (2012). Functional genetic screen of
657 human diversity reveals that a methionine salvage enzyme regulates inflammatory cell death.
658 *Proc Natl Acad Sci U S A* 109, E2343-2352.
- 659 Ko, D.C., Shukla, K.P., Fong, C., Wasnick, M., Brittnacher, M.J., Wurfel, M.M., Holden, T.D.,
660 O'Keefe, G.E., Van Yserloo, B., Akey, J.M., *et al.* (2009). A genome-wide in vitro bacterial-
661 infection screen reveals human variation in the host response associated with inflammatory
662 disease. *Am J Hum Genet* 85, 214-227.
- 663 Ko, D.C., and Urban, T.J. (2013). Understanding human variation in infectious disease
664 susceptibility through clinical and cellular GWAS. *PLoS Pathog* 9, e1003424.
- 665 Labbé, K., and Saleh, M. (2008). Cell death in the host response to infection. *Cell Death Differ*
666 15, 1339-1349.

- 667 Lanzi, G., Moratto, D., Vairo, D., Masneri, S., Delmonte, O., Paganini, T., Parolini, S., Tabellini,
668 G., Mazza, C., Savoldi, G., *et al.* (2012). A novel primary human immunodeficiency due to
669 deficiency in the WASP-interacting protein WIP. *The Journal of experimental medicine* 209, 29-
670 34.
- 671 Lawlor, D.A., Harbord, R.M., Sterne, J.A., Timpson, N., and Davey Smith, G. (2008). Mendelian
672 randomization: using genes as instruments for making causal inferences in epidemiology. *Stat*
673 *Med* 27, 1133-1163.
- 674 Liang, X., Gupta, K., Quintero, J.R., Cernadas, M., Kobzik, L., Christou, H., Pier, G.B., Owen,
675 C.A., and Cataltepe, S. (2019). Macrophage FABP4 is required for neutrophil recruitment and
676 bacterial clearance in *Pseudomonas aeruginosa* pneumonia. *Faseb j* 33, 3562-3574.
- 677 Locke, A.E., Steinberg, K.M., Chiang, C.W.K., Service, S.K., Havulinna, A.S., Stell, L., Pirinen,
678 M., Abel, H.J., Chiang, C.C., Fulton, R.S., *et al.* (2019). Exome sequencing of Finnish isolates
679 enhances rare-variant association power. *Nature* 572, 323-328.
- 680 Longworth, M.S., and Laimins, L.A. (2004). The binding of histone deacetylases and the
681 integrity of zinc finger-like motifs of the E7 protein are essential for the life cycle of human
682 papillomavirus type 31. *J Virol* 78, 3533-3541.
- 683 Lyczak, J.B., Cannon, C.L., and Pier, G.B. (2002). Lung infections associated with cystic
684 fibrosis. *Clin Microbiol Rev* 15, 194-222.
- 685 Malaty, H.M., Engstrand, L., Pedersen, N.L., and Graham, D.Y. (1994). *Helicobacter pylori*
686 infection: genetic and environmental influences. A study of twins. *Annals of internal medicine*
687 120, 982-986.
- 688 Man, A., Slevin, M., Petcu, E., and Fraefel, C. (2019). The Cyclin-Dependent Kinase 5 Inhibitor
689 Peptide Inhibits Herpes Simplex Virus Type 1 Replication. *Scientific reports* 9, 1260.
- 690 Murray, L.A., Sheng, X., and Cristea, I.M. (2018). Orchestration of protein acetylation as a
691 toggle for cellular defense and virus replication. *Nat Commun* 9, 4967.
- 692 Na, Y.R., Jung, D., Gu, G.J., Jang, A.R., Suh, Y.-H., and Seok, S.H. (2015). The early synthesis
693 of p35 and activation of CDK5 in LPS-stimulated macrophages suppresses interleukin-10
694 production. *Science signaling* 8, ra121-ra121.
- 695 Niemöller, U.M., and Täuber, M.G. (1989). Brain edema and increased intracranial pressure in
696 the pathophysiology of bacterial meningitis. *European journal of clinical microbiology &*
697 *infectious diseases* : official publication of the European Society of Clinical Microbiology 8, 109-
698 117.
- 699 Odds, F.C., Brown, A.J., and Gow, N.A. (2004). *Candida albicans* genome sequence: a platform
700 for genomics in the absence of genetics. *Genome Biol* 5, 230.
- 701 Ojala, P.M., Tiainen, M., Salven, P., Veikkola, T., Castaños-Vélez, E., Sarid, R., Biberfeld, P.,
702 and Mäkelä, T.P. (1999). Kaposi's sarcoma-associated herpesvirus-encoded v-cyclin triggers
703 apoptosis in cells with high levels of cyclin-dependent kinase 6. *Cancer Res* 59, 4984-4989.

- 704 Ojala, P.M., Yamamoto, K., Castaños-Vélez, E., Biberfeld, P., Korsmeyer, S.J., and Mäkelä,
705 T.P. (2000). The apoptotic v-cyclin-CDK6 complex phosphorylates and inactivates Bcl-2. *Nature*
706 *cell biology* 2, 819-825.
- 707 Pan, Y., Tian, T., Park, C.O., Lofftus, S.Y., Mei, S., Liu, X., Luo, C., O'Malley, J.T., Gehad, A.,
708 Teague, J.E., *et al.* (2017). Survival of tissue-resident memory T cells requires exogenous lipid
709 uptake and metabolism. *Nature* 543, 252-256.
- 710 Patrick, G.N., Zukerberg, L., Nikolic, M., de la Monte, S., Dikkes, P., and Tsai, L.H. (1999).
711 Conversion of p35 to p25 deregulates Cdk5 activity and promotes neurodegeneration. *Nature*
712 402, 615-622.
- 713 Pfänder, P., Fidan, M., Burret, U., Lipinski, L., and Vettorazzi, S. (2019). Cdk5 Deletion
714 Enhances the Anti-inflammatory Potential of GC-Mediated GR Activation During Inflammation.
715 *Frontiers in Immunology* 10.
- 716 Phelps, W.C., Yee, C.L., Münger, K., and Howley, P.M. (1988). The human papillomavirus type
717 16 E7 gene encodes transactivation and transformation functions similar to those of adenovirus
718 E1A. *Cell* 53, 539-547.
- 719 Power, R.A., Parkhill, J., and de Oliveira, T. (2017). Microbial genome-wide association studies:
720 lessons from human GWAS. *Nature Reviews Genetics* 18, 41-50.
- 721 Price, A.L., Patterson, N.J., Plenge, R.M., Weinblatt, M.E., Shadick, N.A., and Reich, D. (2006).
722 Principal components analysis corrects for stratification in genome-wide association studies. *Nat*
723 *Genet* 38, 904-909.
- 724 Pujol, C., and Bliska, J.B. (2003). The ability to replicate in macrophages is conserved between
725 *Yersinia pestis* and *Yersinia pseudotuberculosis*. *Infect Immun* 71, 5892-5899.
- 726 Revill, K., Wang, T., Lachenmayer, A., Kojima, K., Harrington, A., Li, J., Hoshida, Y., Llovet,
727 J.M., and Powers, S. (2013). Genome-wide methylation analysis and epigenetic unmasking
728 identify tumor suppressor genes in hepatocellular carcinoma. *Gastroenterology* 145, 1424-
729 1435.e1421-1425.
- 730 Roden, D.M., Pulley, J.M., Basford, M.A., Bernard, G.R., Clayton, E.W., Balsler, J.R., and
731 Masys, D.R. (2008). Development of a large-scale de-identified DNA biobank to enable
732 personalized medicine. *Clinical pharmacology and therapeutics* 84, 362-369.
- 733 Saitoh, T., Fujita, N., Hayashi, T., Takahara, K., Satoh, T., Lee, H., Matsunaga, K., Kageyama,
734 S., Omori, H., Noda, T., *et al.* (2009). Atg9a controls dsDNA-driven dynamic translocation of
735 STING and the innate immune response. *Proc Natl Acad Sci U S A* 106, 20842-20846.
- 736 Saka, H.A., Thompson, J.W., Chen, Y.S., Kumar, Y., Dubois, L.G., Moseley, M.A., and Valdivia,
737 R.H. (2011). Quantitative proteomics reveals metabolic and pathogenic properties of *Chlamydia*
738 *trachomatis* developmental forms. *Molecular microbiology* 82, 1185-1203.
- 739 Sharfe, N., Dadi, H.K., Shahar, M., and Roifman, C.M. (1997). Human immune disorder arising
740 from mutation of the alpha chain of the interleukin-2 receptor. *Proc Natl Acad Sci U S A* 94,
741 3168-3171.

- 742 Shi, H., Kichaev, G., and Pasaniuc, B. (2016). Contrasting the Genetic Architecture of 30
743 Complex Traits from Summary Association Data. *Am J Hum Genet* 99, 139-153.
- 744 So, H.C., Chau, C.K., Chiu, W.T., Ho, K.S., Lo, C.P., Yim, S.H., and Sham, P.C. (2017).
745 Analysis of genome-wide association data highlights candidates for drug repositioning in
746 psychiatry. *Nat Neurosci* 20, 1342-1349.
- 747 Soderholm, S., Kainov, D.E., Ohman, T., Denisova, O.V., Schepens, B., Kuleskiy, E., Imanishi,
748 S.Y., Corthals, G., Hintsanen, P., Aittokallio, T., *et al.* (2016). Phosphoproteomics to
749 Characterize Host Response During Influenza A Virus Infection of Human Macrophages.
750 *Molecular & cellular proteomics : MCP* 15, 3203-3219.
- 751 Stahl, J.A., Chavan, S.S., Sifford, J.M., MacLeod, V., Voth, D.E., Edmondson, R.D., and
752 Forrest, J.C. (2013). Phosphoproteomic analyses reveal signaling pathways that facilitate lytic
753 gammaherpesvirus replication. *PLoS Pathog* 9, e1003583.
- 754 Subramanian, A., Tamayo, P., Mootha, V.K., Mukherjee, S., Ebert, B.L., Gillette, M.A.,
755 Paulovich, A., Pomeroy, S.L., Golub, T.R., Lander, E.S., *et al.* (2005). Gene set enrichment
756 analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc*
757 *Natl Acad Sci U S A* 102, 15545-15550.
- 758 Suetrong, B., and Walley, K.R. (2016). Lactic Acidosis in Sepsis: It's Not All Anaerobic:
759 Implications for Diagnosis and Management. *Chest* 149, 252-261.
- 760 Sun, W., Li, J., Jiang, H.G., Ge, L.P., and Wang, Y. (2014). Diagnostic value of MUC1 and
761 EpCAM mRNA as tumor markers in differentiating benign from malignant pleural effusion. *QJM :*
762 *monthly journal of the Association of Physicians* 107, 1001-1007.
- 763 Tarricone, C., Dhavan, R., Peng, J., Areces, L.B., Tsai, L.H., and Musacchio, A. (2001).
764 Structure and regulation of the CDK5-p25(ncK5a) complex. *Molecular cell* 8, 657-669.
- 765 Tartey, S., Matsushita, K., Imamura, T., Wakabayashi, A., Ori, D., Mino, T., and Takeuchi, O.
766 (2015). Essential Function for the Nuclear Protein Akirin2 in B Cell Activation and Humoral
767 Immune Responses. *J Immunol* 195, 519-527.
- 768 Tartey, S., Matsushita, K., Vandenbon, A., Ori, D., Imamura, T., Mino, T., Standley, D.M.,
769 Hoffmann, J.A., Reichhart, J.M., Akira, S., *et al.* (2014). Akirin2 is critical for inducing
770 inflammatory genes by bridging I κ B-zeta and the SWI/SNF complex. *Embo j* 33, 2332-
771 2348.
- 772 Taylor, A.E., Fluharty, M.E., Bjørngaard, J.H., Gabrielsen, M.E., Skorpen, F., Marioni, R.E.,
773 Campbell, A., Engmann, J., Mirza, S.S., Loukola, A., *et al.* (2014). Investigating the possible
774 causal association of smoking with depression and anxiety using Mendelian randomisation
775 meta-analysis: the CARTA consortium. *BMJ Open* 4, e006141.
- 776 Taylor, M.P., Koyuncu, O.O., and Enquist, L.W. (2011). Subversion of the actin cytoskeleton
777 during viral infection. *Nat Rev Microbiol* 9, 427-439.
- 778 Tian, C., Hromatka, B.S., Kiefer, A.K., Eriksson, N., Noble, S.M., Tung, J.Y., and Hinds, D.A.
779 (2017). Genome-wide association and HLA region fine-mapping studies identify susceptibility
780 loci for multiple common infections. *Nat Commun* 8, 599.

- 781 Tsuboi, S., and Meerloo, J. (2007). Wiskott-Aldrich syndrome protein is a key regulator of the
782 phagocytic cup formation in macrophages. *The Journal of biological chemistry* 282, 34194-
783 34203.
- 784 Unlu, G., Qi, X., Gamazon, E.R., Melville, D.B., Patel, N., Rushing, A.R., Hashem, M., Al-Faifi,
785 A., Chen, R., Li, B., *et al.* (2020). Phenome-based approach identifies RIC1-linked Mendelian
786 syndrome through zebrafish models, biobank associations and clinical studies. *Nat Med* 26, 98-
787 109.
- 788 Walenna, N.F., Kurihara, Y., Chou, B., Ishii, K., Soejima, T., Itoh, R., Shimizu, A., Ichinohe, T.,
789 and Hiromatsu, K. (2018). Chlamydia pneumoniae exploits adipocyte lipid chaperone FABP4 to
790 facilitate fat mobilization and intracellular growth in murine adipocytes. *Biochem Biophys Res*
791 *Commun* 495, 353-359.
- 792 Wang, L., Pittman, K.J., Barker, J.R., Salinas, R.E., Stanaway, I.B., Williams, G.D., Carroll, R.J.,
793 Balmat, T., Ingham, A., Gopalakrishnan, A.M., *et al.* (2018). An Atlas of Genetic Variation
794 Linking Pathogen-Induced Cellular Traits to Human Disease. *Cell Host Microbe* 24, 308-
795 323.e306.
- 796 Watanabe, K., Stringer, S., Frei, O., Umicevic Mirkov, M., de Leeuw, C., Polderman, T.J.C., van
797 der Sluis, S., Andreassen, O.A., Neale, B.M., and Posthuma, D. (2019). A global overview of
798 pleiotropy and genetic architecture in complex traits. *Nat Genet.*
- 799 Wei, W.-Q., Bastarache, L.A., Carroll, R.J., Marlo, J.E., Osterman, T.J., Gamazon, E.R., Cox,
800 N.J., Roden, D.M., and Denny, J.C. (2017). Evaluating phecodes, clinical classification software,
801 and ICD-9-CM codes for phenome-wide association studies in the electronic health record. *PloS*
802 *one* 12, e0175508-e0175508.
- 803 Willis, K.L., Patel, S., Xiang, Y., and Shisler, J.L. (2009). The effect of the vaccinia K1 protein on
804 the PKR-eIF2alpha pathway in RK13 and HeLa cells. *Virology* 394, 73-81.
- 805 Xie, X., Lu, J., Kulbokas, E.J., Golub, T.R., Mootha, V., Lindblad-Toh, K., Lander, E.S., and
806 Kellis, M. (2005). Systematic discovery of regulatory motifs in human promoters and 3' UTRs by
807 comparison of several mammals. *Nature* 434, 338-345.
- 808 Yavorska, O.O., and Burgess, S. (2017). MendelianRandomization: an R package for
809 performing Mendelian randomization analyses using summarized data. *Int J Epidemiol* 46,
810 1734-1739.
- 811 Yu, B., Cheng, H.C., Brautigam, C.A., Tomchick, D.R., and Rosen, M.K. (2011). Mechanism of
812 actin filament nucleation by the bacterial effector VopL. *Nat Struct Mol Biol* 18, 1068-1074.
- 813 Yu, Z., Song, H., Jia, M., Zhang, J., Wang, W., Li, Q., Zhang, L., and Zhao, W. (2017). USP1-
814 UAF1 deubiquitinase complex stabilizes TBK1 and enhances antiviral responses. *The Journal*
815 *of experimental medicine* 214, 3553-3563.
- 816 Zahm, J.A., Padrick, S.B., Chen, Z., Pak, C.W., Yunus, A.A., Henry, L., Tomchick, D.R., Chen,
817 Z., and Rosen, M.K. (2013). The bacterial effector VopL organizes actin into filament-like
818 structures. *Cell* 155, 423-434.

- 819 Zhang, X., Bogunovic, D., Payelle-Brogard, B., Francois-Newton, V., Speer, S.D., Yuan, C.,
820 Volpi, S., Li, Z., Sanal, O., Mansouri, D., *et al.* (2015). Human intracellular ISG15 prevents
821 interferon-alpha/beta over-amplification and auto-inflammation. *Nature* 517, 89-93.
- 822 Zheng, T., Liu, W., Oh, S.Y., Zhu, Z., Hu, B., Homer, R.J., Cohn, L., Grusby, M.J., and Elias,
823 J.A. (2008). IL-13 receptor alpha2 selectively inhibits IL-13-induced responses in the murine
824 lung. *J Immunol* 180, 522-529.
825

826 **AUTHOR CONTRIBUTIONS**

827 Conceptualization, A.T.H. and E.R.G.; Methodology, A.T.H., D.Z., L.B., S.J.S., D.C.K., and
828 E.R.G.; Investigation A.T.H., D.Z., L.B., L.W., S.S.Z., S.J.S., D.C.K., and E.R.G. Writing –
829 Original Draft, A.T.H. and E.R.G.; Writing – Review and Editing, A.T.H., D.Z., L.B., L.W., S.S.Z.,
830 S.J.S., D.C.K., and E.R.G, Funding Acquisition- A.T.H. and E.R.G., Supervision, E.R.G.

831

832 **ACKNOWLEDGEMENTS**

833 A.T.H. is supported by the National Institutes of Health (F30HL143826) and Vanderbilt
834 University Medical Scientist Training Program (T32GM007347). E.R.G. is supported by the
835 National Human Genome Research Institute of the National Institutes of Health under Award
836 Number R35HG010718. E.R.G and S.S.Z. are funded by the National Heart, Lung, & Blood
837 Institute of the National Institutes of Health under Award Number R01HL133559. The content is
838 solely the responsibility of the authors and does not necessarily represent the official views of
839 the National Institutes of Health. E.R.G. has also significantly benefitted from a Fellowship at
840 Clare Hall, University of Cambridge (UK) and is grateful to the President and Fellows of the
841 college for a stimulating intellectual home. Genomic data are also supported by individual
842 investigator-led projects including U01-HG004798, R01-NS032830, RC2-GM092618, P50-
843 GM115305, U01-HG006378, U19-HL065962, and R01-HD074711. Additional funding sources
844 for BioVU are listed at <https://victr.vanderbilt.edu/pub/biovu/>. L.B. is supported by R01-
845 LM010685. S.J.S. is supported R01-EB1019804, R01-AI145057, R01-EB014641, R01-
846 HD085853, and DP1-HD086071. D.C.K. is supported by R01-AI118903, R21-AI144586, and
847 R21-AI146520. D.C.K. and L.W. are supported by R21-AI133305.




848

849 **DECLARATION OF INTERESTS**

850 E.R.G. receives an honorarium from the journal *Circulation Research* of the American Heart
851 Association, as a member of the Editorial Board. He performed consulting on pharmacogenetic
852 analysis with the City of Hope / Beckman Research Institute.
853

854 **FIGURE LEGENDS**

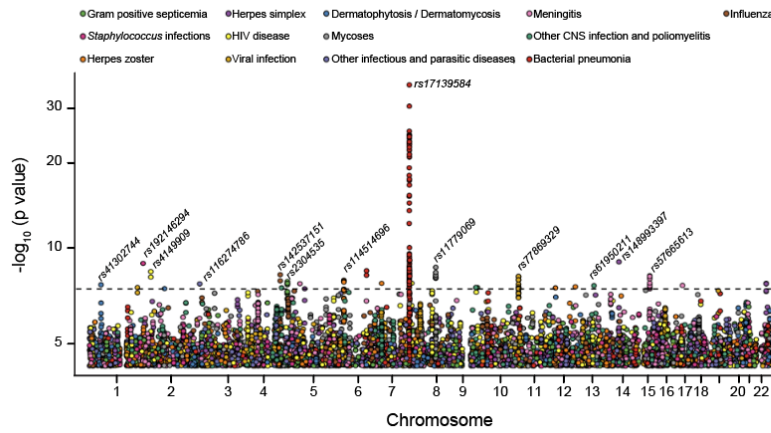
855 **Figure 1.** Overview of ID atlas resource. List of ID traits tested with corresponding Phecode
 856 (phewascatalog.org) in parentheses.
 857
 858

Infectious diseases	Study population	Genetic analyses
<p>Bacterial pneumonia (480.1) Candidiasis (112) Dermatophytosis and dermatomycosis (110) Encephalitis (323) Escherichia coli infection (041.4) Graft vs. host disease (081.1) Gram-negative septicemia (038.1) Gram-positive septicemia (038.2) Helicobacter pylori infection 041.8 Helminthiasis (134) Hepatitis A infection (070.1) Hepatitis B infection (070.2) Hepatitis C infection (070.3) Herpes simplex infection (054) Herpes Zoster (053) Human Immunodeficiency Virus infection (071) Human papillomavirus infection (078) Infection with drug-resistant organism (041.9) Infectious mononucleosis (079.2) Influenza infection (481) Intestinal infection (008) Meningitis (320) Mycoses (117) Other central nervous system infections and poliomyelitis (324) Other infectious and parasitic diseases (136) Protozoan infection (131) Sepsis (994) Sexually transmitted disease excluding HIV and hepatitis (090) Spirochetal infection (130) Staphylococcus infection (041.1) Streptococcus infection (041.2) Tuberculosis (010) Varicella infection (079.1) Viral infection (079) Viral pneumonia (480.2)</p>	  	<p>GWAS for 35 ID traits</p> <p>TWAS for 35 ID traits</p> <p>Development of tissue-specific TWAS catalog of ID traits</p> <p>Gene-level replication in the UK Biobank & FinnGen</p> <p>Pathogen genus determined by serology and culture linked to host genetic information</p> <p>Phenome scan across cardiac, hematologic, respiratory, and neurologic traits</p> <p>Development of tissue-specific TWAS catalog of pathogen-induced cellular phenotypes</p> <p>Mendelian Randomization to infer causal relationships</p>

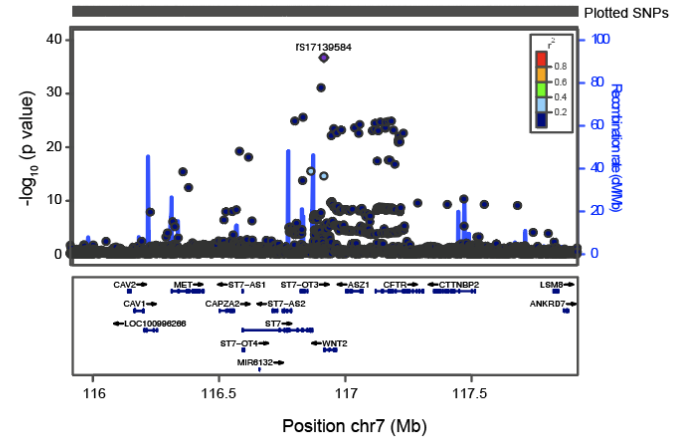
859
 860

861 **Figure 2.** Genome-wide association study (GWAS) of ID traits. (A) Threshold for inclusion of
 862 SNP associations was set at 1.0×10^{-4} . Genome-wide significance for an ID trait was set at $p =$
 863 5.0×10^{-8} , as indicated by the horizontal dotted line. The subset of 13 ID traits (among the full set
 864 tested) with variants that meet the traditional genome-wide significance threshold are included.
 865 The top variant association for each of the 13 traits is labeled. The most significant variant
 866 association is with bacterial pneumonia ($p < 1.0 \times 10^{-30}$). (B) LocusZoom plot at the sentinel
 867 variant, rs17139584, associated with bacterial pneumonia. Several variants in low LD ($r^2 < 0.20$)
 868 with the sentinel variant, including variants in the cystic fibrosis gene *CFTR* and in the *MET*
 869 gene >650 Kb upstream, are genome-wide significant for bacterial pneumonia. The sentinel
 870 variant remains statistically significant ($p = 0.007$) after adjusting for a diagnosis of cystic
 871 fibrosis.
 872

A

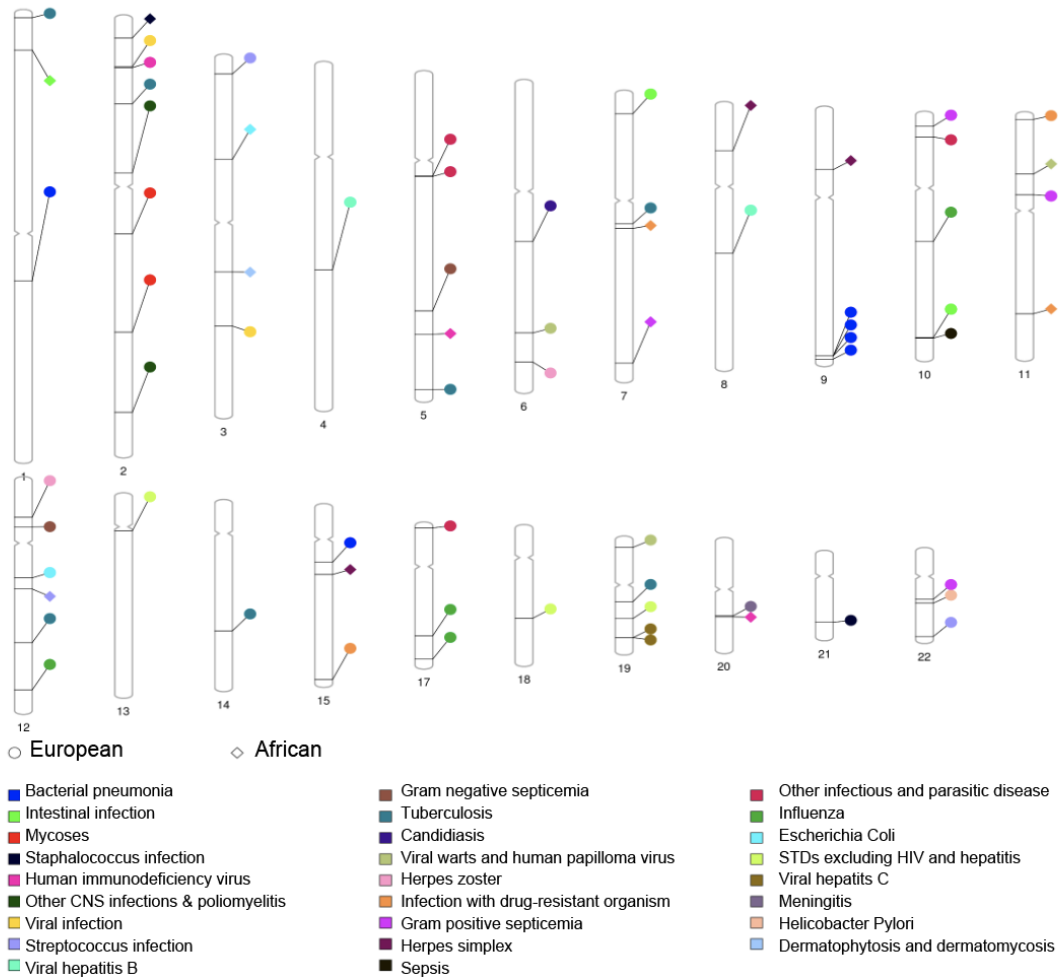


B

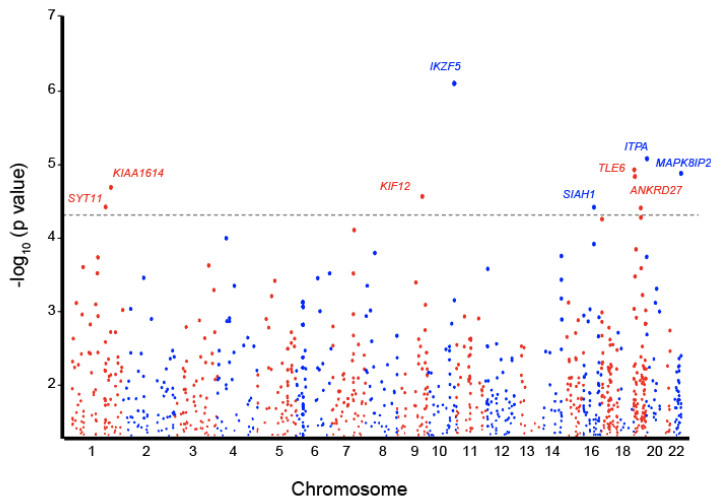


874 **Figure 3.** Transcriptome-wide association studies (TWAS) of 35 ID traits reveal novel ID-
875 associated genes. The genetic component of gene expression for autosomal genes was
876 individually tested for association with each of 35 ID traits (see Methods). (A) Experiment-wide
877 or ID-specific significant genes are displayed on the ideogram using their chromosomal
878 locations and color-coded using the associated ID traits. Most associations represent unique
879 genes within the implicated loci, suggesting the genes are not tagging another causal gene. A
880 locus on chromosome 9, by contrast, shows multiple associations with the same ID trait, which
881 may indicate correlation of the expression traits with a single causal gene in the locus. (B)
882 Manhattan plot shows the PrediXcan associations with sepsis (Phecode 994; number of cases
883 2,921; number of controls 22,874). Dashed line represents $p < 5 \times 10^{-5}$. The gene *IKZF5* was
884 significant ($p = 8.16 \times 10^{-7}$) after Bonferroni correction for the number of genes tested. (C) Q-Q
885 plot of FinnGen replication p-values for genes associated with intestinal infection ($p < 0.05$) in
886 BioVU (red) compared to the remaining set of genes (black). The ID-associated genes tended to
887 be more significant in the independent dataset than the remaining genes, as evidenced by the
888 leftward shift in the Q-Q plot.
889

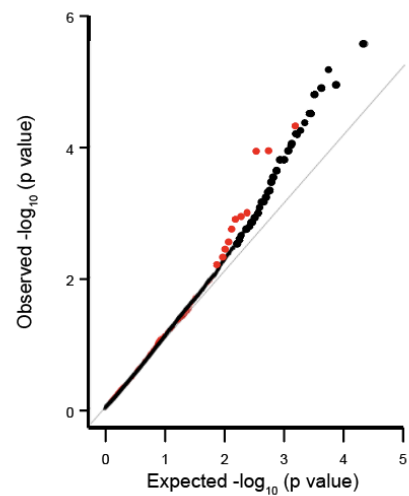
A



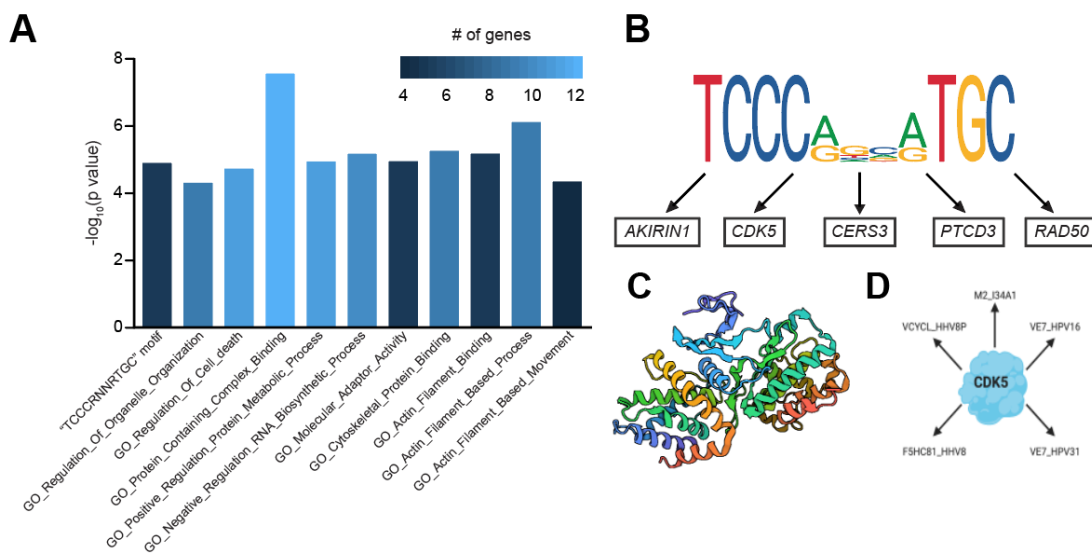
B



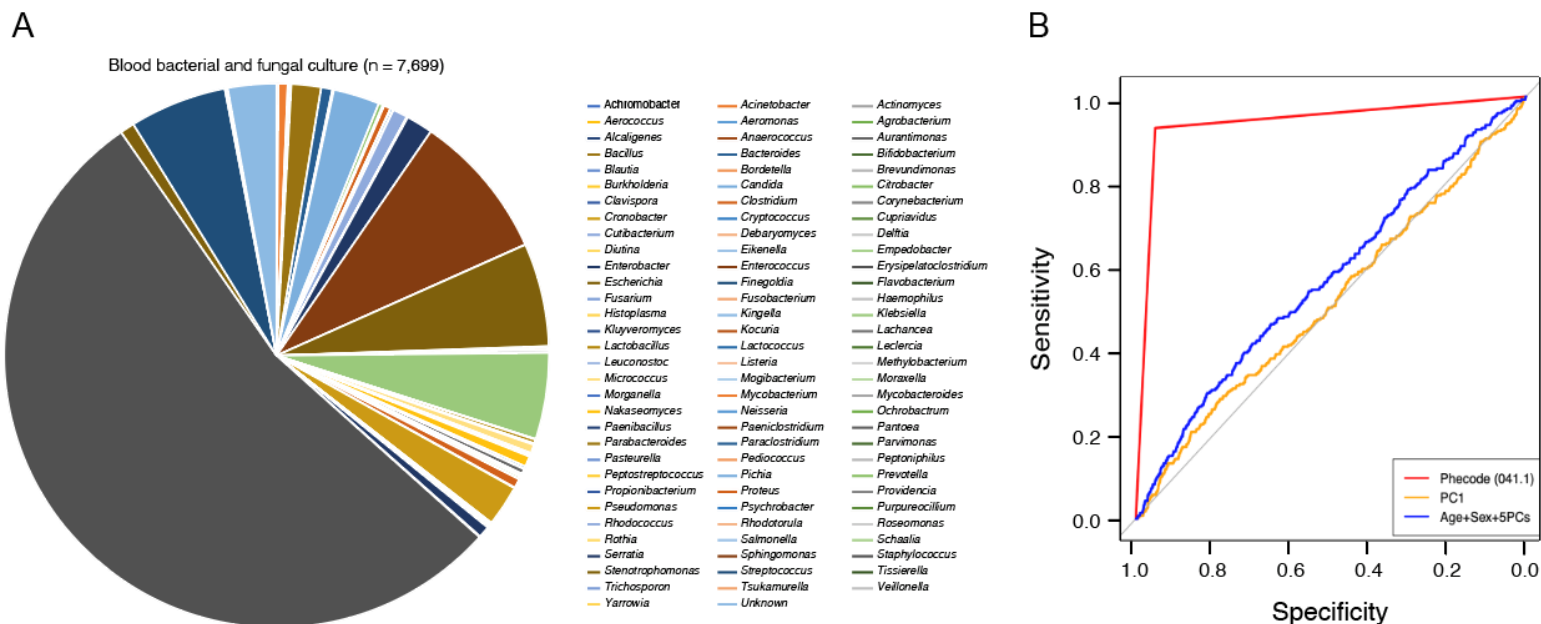
C



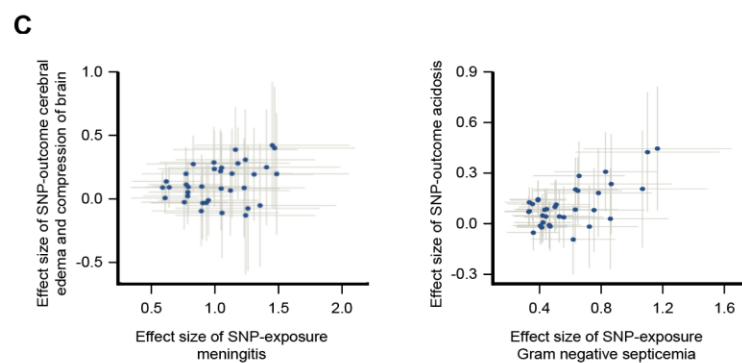
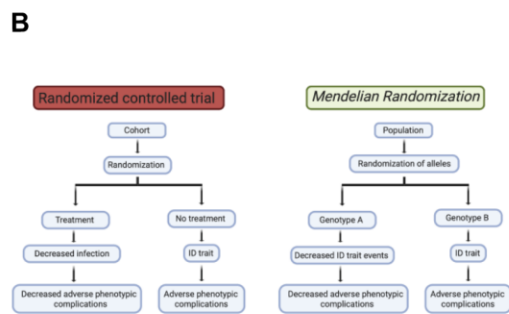
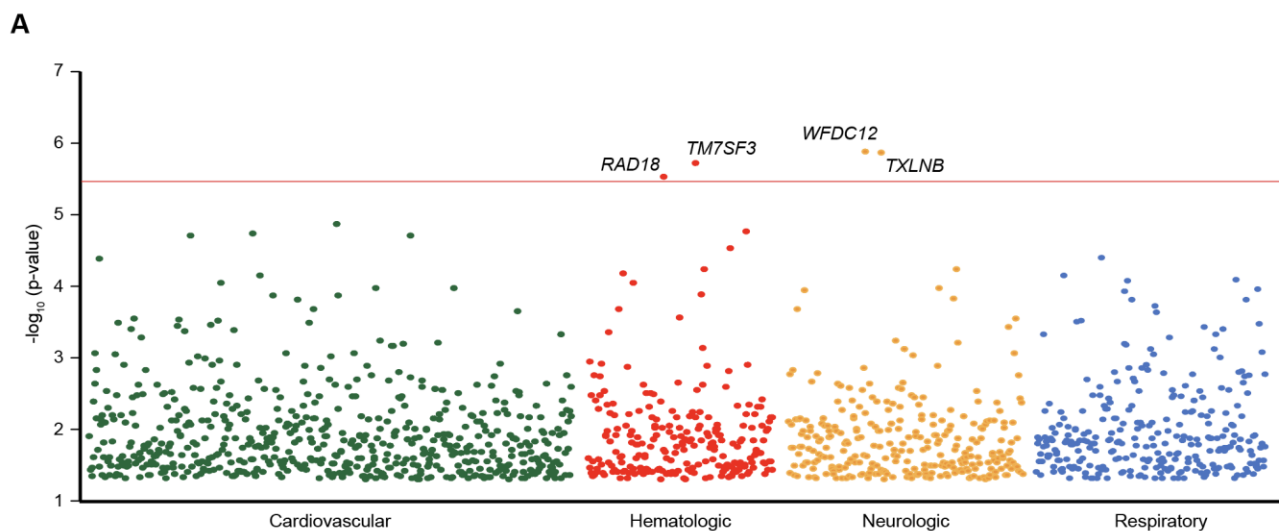
892 **Figure 4.** Enriched pathways across multiple ID traits and pathogen evolutionary strategies to
 893 promote infection. (A) Gene set enrichment analysis of ID associated genes having also
 894 nominal associations with additional ID traits. All gene sets satisfied false discovery rate < 0.05
 895 for pathway enrichment and included known biological processes (e.g. protein complex
 896 formation, cytoskeletal protein binding, cell death, actin motility, etc.) relevant to the biology of
 897 infection. (B) Highly conserved motif “TCCCRNNRTGC”, within 4 kb of TSS of ID-associated
 898 genes, is enriched among the multi-ID associated genes and does not match any known
 899 transcription factor binding site. Genes with this motif near the TSS include *AKIRIN2*, *CDK5*,
 900 *RAD50*, *PTCD3*, and *CERS3*. This suggests a strategy that the pathogens may broadly exploit
 901 to hijack the host transcriptional machinery. (C) CDK5 is an example of a multi-ID associated
 902 gene, significantly associated with Gram-positive septicemia and nominally associated with
 903 other IDs, including herpes simplex virus. CDK5 is activated by its regulatory subunit p35/p25.
 904 The CDK5-p25 complex regulates inflammation (whose large-scale disruption is characteristic
 905 of septicemia) and induces cytoskeletal disruption in neurons (where the herpes virus promotes
 906 lifelong latent infection). Structure of the CDK5-p25 complex (PDB: 1H4L, (Tarricone et al.,
 907 2001)) is shown here. The A and B chains are required for cytoskeletal protein binding (CDK5),
 908 whereas the D and E chains (p25) are involved in actin regulation and kinase function, all
 909 functions implicated in our pathway analysis. (D) Multi-ID associated genes identified by our
 910 study have also been observed in host-pathogen protein complexes (by coimmunoprecipitation,
 911 affinity chromatography, and two-hybrid approaches, among others) for the specific pathogens
 912 responsible for the ID traits. Interactions of pathogen proteins with CDK5 are shown here.
 913 M2_134A1 (UniProt: PO6821) is the matrix protein 2 component of the proton-selective ion
 914 channel required for influenza A viral genome release during cellular entry and is targeted by
 915 the anti-viral drug amantadine (Hay et al., 1985). VE7_HP16 (UniProt: PO3129) is a
 916 component of human papillomavirus (HPV) required for cellular transformation and trans-
 917 activation through disassembly of E2F1 transcription factor from RB1 leading to impaired
 918 production of type I interferons (Barnard et al., 2000; Chellappan et al., 1992; Phelps et al.,
 919 1988). VE7_HP31 (UnitProt: P17387) engages histone deacetylases 1 and 2 to promote
 920 HPV31 genome maintenance (Longworth and Laimins, 2004). VCYCL_HHV8P (UniProt:
 921 Q77Q36) is a cyclin homolog within the human herpesvirus 8 genome that has been shown to
 922 control cell cycle through CDK6 and induce apoptosis through Bcl2 (Duro et al., 1999; Ojala et
 923 al., 1999; Ojala et al., 2000). F5HC81_HHV8 (UniProt: F5HC81) is not well-characterized, but
 924 predicted to act as a viral cyclin homolog. This suggests a second strategy that the pathogens
 925 exploit, i.e., alteration of the host proteome, to promote infection.



927 **Figure 5. Pathogen genus identification from clinical blood cultures linked to whole-**
 928 **genome information reveals insights into host colonization and infection.** (A) Bacterial
 929 and fungal pathogens identified from blood (n = 7,699 positive cultures across 94 genera) from
 930 2,417 individuals. (B) Area under the receiver operating characteristic curve (AUC), showing that
 931 the clinical trait *Staphylococcus* infection (Phecode = 041.1) performs well in classifying
 932 *Staphylococcus aureus* infection based on blood culture data from (A), with AUC of 0.938 with
 933 standard error of 0.008. The first PC in the European ancestry samples and a model with age,
 934 sex, and the first 5 PCs, both with substantially lower performance (AUC of 0.514 and 0.568,
 935 respectively), are also shown.
 936
 937

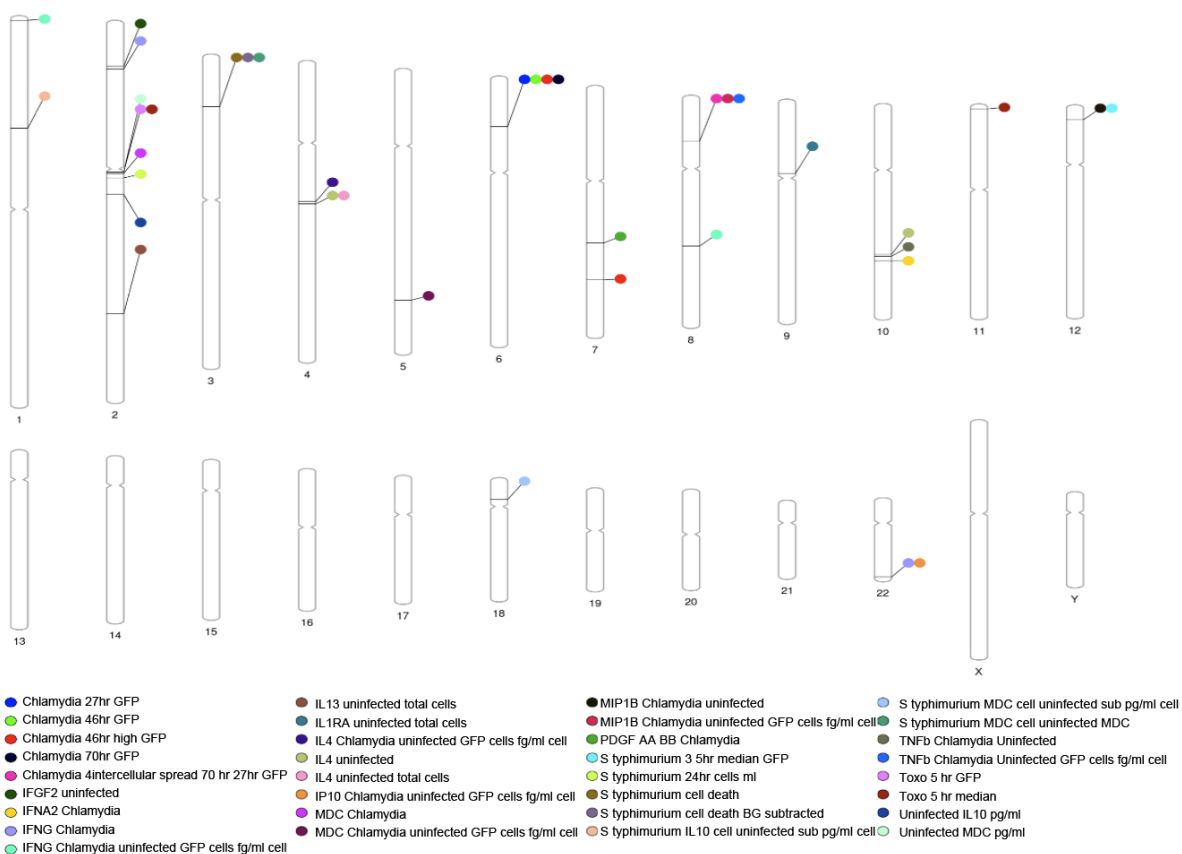


939 **Figure 6.** Phenome-scale scan of 70 ID-associated genes across 197 cardiovascular, 940 hematologic, neurologic, and respiratory phenotypes (cases > 200) in BioVU 941 (phewascatalog.org) identifies genes association with both disease risk and corresponding 942 known complications of the infection. (A) Each dot represents the association of an ID- 943 associated gene with one of the 197 (hematologic, respiratory, cardiovascular, and neurologic) 944 phenotypes. Horizontal red line indicates threshold for statistical significance correcting for 945 number of phenotypes and ID-associated genes tested. We identify four gene-phenotype pairs 946 reaching experiment-wide significance: 1) *WFDC12*, our most significant ($p = 4.23 \times 10^{-6}$) 947 association with meningitis, is also associated with cerebral edema and compression of brain ($p = 1.35 \times 10^{-6}$), a feared clinical complication of meningitis (Niemöller and Täuber, 1989); 2) 948 *TM7SF3*, the most significant gene with Gram-negative sepsis ($p = 1.37 \times 10^{-6}$), is also 949 associated with acidosis ($p = 1.95 \times 10^{-6}$), a known metabolic derangement associated with 950 severe sepsis (Suetrong and Walley, 2016); 3) *TXLNB*, the most significant gene associated 951 with viral warts and human papillomavirus infection ($p = 4.35 \times 10^{-6}$), is also associated with 952 abnormal involuntary movements, $p = 1.39 \times 10^{-6}$; and 4) *RAD18*, the most significant gene 953 associated with Streptococcus infection ($p = 2.01 \times 10^{-6}$), is also associated with anemia in 954 neoplastic disease ($p = 3.10 \times 10^{-6}$). (B) Mendelian randomization framework. P-value threshold 955 used to define an instrumental variable was set at $p < 1.0 \times 10^{-5}$ and variants in linkage 956 equilibrium ($r^2 = 0.01$) were used. (C) Mendelian Randomization provides strong support for 957 causal exposure-outcome relationships for 1) meningitis and compression of brain (left, median- 958 weighted estimator $p = 2.7 \times 10^{-3}$); and 2) gram-negative septicemia and acidosis (right, median- 959 weighted estimator $p = 2.0 \times 10^{-7}$). Grey lines indicate 95% confidence interval. 960

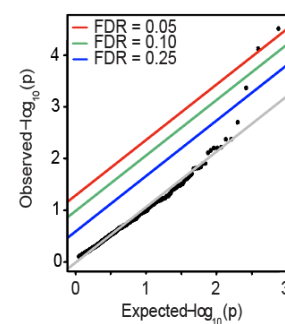


962 **Figure 7.** TWAS of 79 pathogen-exposure induced cellular traits improves identification of
 963 pathogen-induced cellular mechanisms. (A) Genes reaching significance in Hi-HOST after
 964 correction for the total number of genes and cellular phenotypes tested. (B) Integration of EHR
 965 data into Hi-HOST facilitates replication of gene-level associations with a clinical ID trait. Genes
 966 nominally associated ($p < 0.05$) with Gram-positive septicemia (Phecode 038.2) in BioVU show
 967 significant enrichment for *Staphylococcus* toxin exposure, a Hi-HOST phenotype. The Q-Q plot
 968 shows the distribution of TWAS p-values in the Hi-HOST data for the top genes in the BioVU
 969 data. False discovery rate (FDR) thresholds at 0.25 (blue), 0.10 (green), and 0.05 (red) are
 970 shown. (C) Integration of EHR data into Hi-HOST also improves the signal-to-noise ratio in Hi-
 971 HOST. For example, the top 300 genes nominally associated with *Staphylococcus* infection
 972 (Phecode 041.1) in BioVU ($p < 0.016$, red) depart from null expectation for their TWAS
 973 associations with *Staphylococcus* toxin exposure in Hi-HOST compared to the full set of genes
 974 (black).
 975
 976

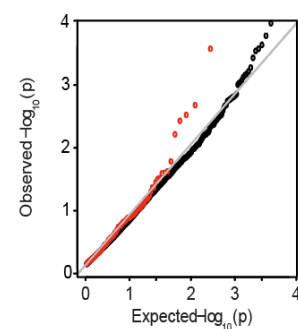
A



B



C



978

979

980 **Table 1.** Significant trait-specific gene-level associations with individual infectious
 981 disease phenotypes (for which number of cases > 100). Experiment-wide findings are
 982 noted in **bold**.

983

Gene	PheCode	Phenotype	Cases	Controls	Ancestry	Odds ratio	P value
IKZF5	994	Sepsis	2,921	22,874	European	0.91	8.16x10 ⁻⁷
AKIRIN2	112	Candidiasis	2,284	21,426	European	0.91	2.83x10 ⁻⁶
PSMG1	041.1	<i>Staphylococcus</i> infection	2,180	19,844	European	0.90	3.13x10 ⁻⁶
AGTR1	079	Viral infection	1,811	20,904	European	1.12	1.49x10 ⁻⁶
SLC35F6	079	Viral infection	1,811	20,904	European	0.89	3.30x10 ⁻⁶
NDUFA4	008	Intestinal infection	1,608	24,187	European	1.16	1.83x10⁻⁹
C10orf120	008	Intestinal infection	1,608	24,187	European	1.13	4.92x10⁻⁸
RAD18	041.2	<i>Streptococcus</i> infection	1,262	19,844	European	1.16	2.01x10 ⁻⁶
MAPK8IP2	041.2	<i>Streptococcus</i> infection	1,262	19,844	European	1.14	3.81x10 ⁻⁶
AVIL	041.4	<i>Escherichia Coli</i>	1,231	19,844	European	1.16	1.58x10 ⁻⁶
STAP2	078	Viral warts and human papilloma virus	1,152	20,904	European	1.15	2.33x10 ⁻⁶
TXLNB	078	Viral warts and human papilloma virus	1,152	20,904	European	0.86	4.35x10 ⁻⁶
SLCO1A2	053	Herpes zoster	989	20,904	European	0.93	1.64x10 ⁻⁷
CLDN20	053	Herpes zoster	989	20,904	European	0.86	4.54x10 ⁻⁶
IGF2	041.9	Infection with drug-resistant organism	893	19,844	European	0.83	4.01x10 ⁻⁷
CERS3	041.9	Infection with drug-resistant organism	893	19,844	European	1.17	4.18x10 ⁻⁶
TOR4A	480.1	Bacterial pneumonia	862	18,054	European	0.84	5.15x10⁻⁸
FAM166A	480.1	Bacterial pneumonia	862	18,054	European	1.19	1.10x10 ⁻⁷
C9orf173	480.1	Bacterial pneumonia	862	18,054	European	1.18	4.48x10 ⁻⁷
PIP5K1A	480.1	Bacterial pneumonia	862	18,054	European	1.16	7.00x10 ⁻⁷
NELFB	480.1	Bacterial pneumonia	862	18,054	European	0.86	1.87x10 ⁻⁶
AVEN	480.1	Bacterial pneumonia	862	18,054	European	0.85	3.31x10 ⁻⁶
TM7SF3	038.1	Gram negative septicemia	820	19,844	European	1.17	1.37x10 ⁻⁶
RAD50	038.1	Gram negative septicemia	820	19,844	European	1.17	4.50x10 ⁻⁶
ZNF577	070.3	Viral hepatitis C	808	20,904	European	0.84	6.21x10 ⁻⁷
ZNF649	070.3	Viral hepatitis C	808	20,904	European	0.85	1.85x10 ⁻⁶
SETD9	136	Other infectious and parasitic diseases	746	24,770	European	0.83	3.04x10⁻⁸
AC022431.1	136	Other infectious and parasitic diseases	746	24,770	European	1.20	7.92x10⁻⁸
MYO1C	136	Other infectious and parasitic diseases	746	24,770	European	1.10	2.97x10 ⁻⁶
NUDT5	136	Other infectious and parasitic diseases	746	24,770	European	0.84	3.52x10 ⁻⁶
MAATS1	110	Dermatophytosis and dermatomycosis	654	3,330	African	0.80	4.82x10 ⁻⁶
PTPN4	117	Mycoses	627	21,426	European	0.79	1.56x10⁻⁷
WIPF1	117	Mycoses	627	21,426	European	1.20	2.72x10 ⁻⁶
ALX4	038.2	Gram positive septicemia	613	19,844	European	1.25	4.21x10⁻⁸
C22orf31	038.2	Gram positive septicemia	613	19,844	European	0.81	2.05x10 ⁻⁶
IL2RA	038.2	Gram positive septicemia	613	19,844	European	1.20	3.88x10 ⁻⁶
VVA5B1	008	Intestinal infection	368	4,060	African	1.30	3.85x10 ⁻⁶
ATP6V1C2	041.1	<i>Staphylococcus</i> infection	358	3,337	African	1.33	1.51x10 ⁻⁶
WDR66	481	Influenza	272	18,054	European	0.71	3.47x10 ⁻⁷
FAM20A	481	Influenza	272	18,054	European	0.77	1.51x10 ⁻⁶
HKDC1	481	Influenza	272	18,054	European	1.34	2.20x10 ⁻⁶
ASPSCR1	481	Influenza	272	18,054	European	1.35	2.76x10 ⁻⁶
FAM208A	041.4	<i>Escherichia Coli</i>	243	3,337	African	1.42	1.15x10 ⁻⁶
TBK1	041.2	<i>Streptococcus</i> infection	229	3,337	African	1.41	4.70x10 ⁻⁶
DNAJC5G	071	Human immunodeficiency virus	196	20,904	European	1.08	7.36x10 ⁻⁷
FABP4	070.2	Viral hepatitis B	166	20,904	European	0.70	4.39x10 ⁻⁶
ANK2	070.2	Viral hepatitis B	166	20,904	European	0.77	4.56x10 ⁻⁶
HIP1	041.9	Infection with drug-resistant organism	165	3,337	African	1.46	6.74x10 ⁻⁷
C11orf53	041.9	Infection with drug-resistant organism	165	3,337	African	0.72	4.98x10 ⁻⁶
EPCAM	010	Tuberculosis	156	19,844	European	1.40	5.04x10⁻⁸
AL589739.1	010	Tuberculosis	156	19,844	European	1.55	9.46x10⁻⁸
PROX2	010	Tuberculosis	156	19,844	European	1.40	9.70x10 ⁻⁷
USP44	010	Tuberculosis	156	19,844	European	1.44	1.07x10 ⁻⁶
GPRIN1	010	Tuberculosis	156	19,844	European	1.49	2.55x10 ⁻⁶
NSUN5	010	Tuberculosis	156	19,844	European	0.75	2.76x10 ⁻⁶
C19orf55/PROSER3	010	Tuberculosis	156	19,844	European	1.51	3.30x10 ⁻⁶
DNAJC17	054	Herpes simplex	154	3,241	African	0.65	2.75x10 ⁻⁷
CHMP5	054	Herpes simplex	154	3,241	African	1.36	2.22x10 ⁻⁶

GNRH1	054	Herpes simplex	154	3,241	African	1.47	4.64x10 ⁻⁶
TXNL1	152	Sexually transmitted infections excluding HIV and hepatitis	152	25,643	European	0.64	3.92x10 ⁻⁷
LTBP4	152	Sexually transmitted infections excluding HIV and hepatitis	152	25,643	European	0.70	2.13x10 ⁻⁶
CRYL1	152	Sexually transmitted infections excluding HIV and hepatitis	152	25,643	European	0.70	2.23x10 ⁻⁶
LIMK2	041.8	Helicobacter Pylori	150	19,844	European	0.71	2.88x10 ⁻⁶
WFDC12	320	Meningitis	144	25,170	European	1.16	4.23x10 ⁻⁶
TNNC2	071	Human immunodeficiency virus	139	3,241	African	1.44	2.49x10 ⁻⁶
PRELID2	071	Human immunodeficiency virus	139	3,241	African	1.47	2.68x10 ⁻⁶
PTCD3	324	Other CNS infections and poliomyelitis	136	25,170	European	0.89	3.74x10⁻⁸
ATG9A	324	Other CNS infections and poliomyelitis	136	25,170	European	0.76	2.46x10 ⁻⁶
EIF3M	078	Viral warts and human papilloma virus	136	3,241	African	1.48	3.22x10 ⁻⁶
CDK5	038.2	Gram positive septicemia	114	3,337	African	0.65	3.64x10 ⁻⁶

984

985

986

987

988

989

990

991 **STAR★METHODS**

992 **CONTACT FOR REAGENT AND RESOURCE SHARING**

993 Further information and requests for resources and reagents should be directed to and will be
994 fulfilled by the Lead Contact, Eric R. Gamazon (eric.gamazon@vanderbilt.edu).

995

996 **EXPERIMENTAL MODEL AND SUBJECT DETAILS**

997 *BioVU*

998 BioVU, one of the largest DNA biobanks tied to an EHR database, is a subset of the
999 synthetic derivative (SD), a deidentified electronic health record, consisting of individuals with
1000 whole-genome genetic information. Detailed information on the construction, utilization, ethics,
1001 and policies of the BioVU resource is described elsewhere (Roden et al., 2008). ID traits were
1002 defined based on a hierarchical grouping of International Classification of Diseases, Ninth
1003 Revision (ICD-9) codes into phenotype codes (Phecodes) representing clinical traits, as
1004 previously described (Denny et al., 2013; Denny et al., 2010). (See below for a description of
1005 pathogen culture and viral test data in the BioVU individuals, including genera detected from
1006 different types of cultures.) We used version 1.2 of the Phecode Map containing 1,965
1007 Phecodes based on 20,203 ICD-9 codes, which substantially improves signal-to-noise and more
1008 accurately reflects the clinical trait. Phecodes may exclude related phenotypes (e.g., in the case
1009 of Gram negative septicemia (Phecode = 038.1), the range of Phecodes given by 010-041.99,
1010 involving bacterial infection) and, importantly, include the definition of the appropriate control
1011 group (Wei et al., 2017). Detailed description of Phecode trait maps can be found at
1012 pewascatalog.org. As an efficient and viable model for human genetics research, the Phecode
1013 system has been used to perform phenome-wide association studies (PheWAS) for validation of

1014 known genetic associations and discovery of new genetic disorders (Denny et al., 2013; Unlu et
1015 al., 2020).

1016

1017 *Pathogen culture and virology data linked to whole-genome genetic information*

1018 The SD consists of a wide range of clinical microbiological data. For individuals with
1019 whole-genome genetic information, we analyzed pathogen (bacterial, mycobacterial, and fungal)
1020 culture data derived from the following positive cultures for the indicated clinical samples: 1)
1021 blood (n = 7,699), 2) sputum (n = 2,478), 3) sinus/nasopharyngeal (n = 1,820), 4) bronchial-
1022 alveolar lavage (n = 1,265), and 5) tracheal sampling (n = 422). Furthermore, we analyzed a
1023 respiratory panel containing 28 viral strains from 2,890 individuals with whole-genome genetic
1024 information. Viral strains included the following: 1) Adenovirus, 2) Bocavirus, 3) Bordetella
1025 parapertussis, 4) Bordetella pertussis, 5) Chlamydia pneumoniae, 6) Coronavirus 229E, 7)
1026 Coronavirus HKU1, 8) Coronavirus NL63, 9) Coronavirus NOS, 10) Coronavirus OC43, 11)
1027 Enterovirus/Rhinovirus, 12) Human Metapneumovirus, 13) Influenza A, 14) Influenza A, H1, 15)
1028 Influenza A, H1N1, 16) Influenza A, H3, 17) Influenza B, 18) Mycoplasma pneumoniae, 19)
1029 Parainfluenza, 20) Parainfluenza 1, 21) Parainfluenza 2, 22) Parainfluenza 3, 23) Parainfluenza
1030 4, 24) Respiratory syncytial virus (RSV), 25) RSV, A, 26) RSV, B, and 27) Rhinovirus. The
1031 pathogen information for each individual in our study included: 1) Total number of cultures; 2)
1032 Number of negative cultures (i.e., no pathogen growth); 3) Number of ambiguous cultures (i.e.,
1033 normal upper respiratory bacteria or low level contamination); 4) Number of positive cultures
1034 (i.e., the number of cultures with growth consistent with clinical infection); 5) Genus or genera
1035 isolated (up to 96 unique genera per sample site), which ranged from zero to 10 per sample.

1036

1037 **METHODS DETAILS**

1038 *GWAS of ID traits*

1039 GWAS of the ID traits were performed on the 23,294 BioVU individuals of European
1040 ancestry. Quality control pre-processing and SNP-level imputation were conducted, as
1041 previously described (Unlu et al., 2020). Genomic ancestry was quantified using principal
1042 components analysis of the genotype data (Derks et al., 2017; Price et al., 2006). The
1043 association analysis was performed using age, gender, batch, and the first five principal
1044 components as covariates.

1045

1046 *Conditional SNP-level analysis*

1047 We performed conditional analysis on the top GWAS association with the ID trait (in this
1048 case, bacterial pneumonia) to determine whether it was driven by a related covariate (in this
1049 case, cystic fibrosis status). We used logistic regression to model the conditional probability of
1050 the infectious disease:

$$1051 \ln \frac{P(Y=1 | s)}{1 - (P(Y=1 | s))} = \beta_0 + \beta_1 s + \beta_2 (CF)$$

1052 where s is the genotype at the sentinel variant, Y is the disease (i.e., bacterial pneumonia)
1053 status, and CF is the covariate of interest (i.e., cystic fibrosis).

1054

1055 *Transcriptome-wide association studies (TWAS) using PrediXcan*

1056 We performed multi-tissue PrediXcan (Barbeira et al., 2019; Gamazon et al., 2018; Gamazon et al., 2015) in the
1057 23,294 BioVU subjects. Experiment-wide significance was determined using Bonferroni
1058 correction for the total number of genes tested ($n = 9,868$) across 35 phenotypes (i.e., $p <$
1059 1.4×10^{-7}). Trait-specific significance was determined using Bonferroni correction for the total
1060 number of genes tested ($n = 9,868$, $p < 5.07 \times 10^{-6}$). Genomic ancestry was quantified using
1061 principal components analysis (Derks et al., 2017; Price et al., 2006).

1062

1063 *GWAS and TWAS Replication in the UK Biobank and FinnGen consortia*

1064 Replication of GWAS and TWAS was performed in the UK Biobank (Bycroft et al., 2018)
1065 and FinnGen consortia (Locke et al., 2019). We used the UK Biobank
1066 (<http://www.nealelab.is/uk-biobank>) and the FinnGen (https://www.finnngen.fi/en/access_results)
1067 summary results to generate the gene-level associations.

1068
1069 *Classification of pathogen infection based on serology and culture data using several classifiers*

1070 Let X be a classifier (e.g., the Phecode or a logistic regression classifier) of serology and
1071 culture data based infection for a given pathogen, with probability density $\varphi_+(x)$ for positive
1072 instances and probability density $\varphi_-(x)$ for negative instances. The ROC curve plots the
1073 specificity (SP) and sensitivity (SN) at various thresholds:

1074
$$SN(T) = \int_T^{\infty} \varphi_+(x) dx$$

1075
$$SP(T) = 1 - \int_T^{\infty} \varphi_-(x) dx$$

1076 The area Ω under the curve (AUC) is given by:

1077
$$\Omega = \int_{-\infty}^{\infty} SN(T) SP'(T) dT = \int_{-\infty}^{\infty} \int_T^{\infty} \varphi_+(x) \varphi_-(T) dx dT = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} I(x > T) \varphi_+(x) \varphi_-(T) dx dT$$

1078 where $I(A)$ is the indicator function, i.e., equal to one if $(x, T) \in A$ and zero otherwise. The last
1079 equals the probability that the classifier X ranks a randomly chosen positive instance (of culture
1080 data based infection) higher than a randomly chosen negative instance. We estimated the
1081 sampling distribution of Ω , using bootstrapping ($n = 100$) (Efron, 1979). We used the pROC
1082 package for visualization.

1083
1084 *Causal inference by Mendelian Randomization*

1085 To infer causality between the infectious diseases and potential complications, we
1086 performed Mendelian Randomization (MR, (Davey Smith and Hemani, 2014; Lawlor et al.,

1087 2008)) in 23,294 individuals of European ancestry in BioVU. To define instrumental variables
1088 (IVs), we clumped the exposure-associated SNPs with high linkage disequilibrium (LD) using
1089 Plink1.9 ($p < 1 \times 10^{-5}$, $r^2 = 0.01$). Only biallelic non-palindromic variants were considered as IVs.
1090 Considering the pervasive horizontal pleiotropy in human genetic variation (Jordan et al., 2019),
1091 we applied summary statistics based MR-Egger regression (Bowden et al., 2015). MR-Egger
1092 regression generalizes the inverse-variance weighted method, where the intercept is assumed
1093 to be zero. We also used the weighted-median estimator (Bowden et al., 2016) to test the
1094 causal effect of the exposure trait on the outcome. We leveraged the R package
1095 'MendelianRandomization'.

1096

1097 *High-throughput Human in vitro Susceptibility Testing (Hi-HOST)*

1098 We generated an atlas of TWAS associations with 79 pathogen-induced cellular traits –
1099 including infectivity and replication, cytokine levels, and host cell death (Wang et al., 2018)
1100 using the Hi-HOST platform (Ko et al., 2012; Ko et al., 2009). A list of populations, pathogens
1101 and project description may be found at <http://h2p2.oit.duke.edu/About/>, and phenotype
1102 definitions and family-based GWAS of the Hi-HOST Phenome Project were previously
1103 described (Wang et al., 2018). Briefly, lymphoblastoid cell lines (LCLs) from the 1000 Genomes
1104 Consortium (Auton et al., 2015) were obtained from the Coriell Institute. The LCLs represented
1105 diverse populations, including ESN (Esan in Nigeria), GWD (Gambians in Western Divisions in
1106 the Gambia), IBS (Iberian Population in Spain), and KHV (Kinh in Ho Chi Minh City, Vietnam).
1107 LCLs were cultured in RPMI 1640 media containing 10% fetal bovine serum, 2 mM glutamine,
1108 100 U/ml of penicillin-G, and 100 mg/ml streptomycin for 8 days prior to experimental use, as
1109 previously described (Wang et al., 2018). *Chlamydia trachomatis* infection of LCLs was
1110 performed using *C. trachomatis* LGV-L2 Rif^R pGFP::SW2 (Saka et al., 2011). *Salmonella*
1111 infection was performed using pMMB67GFP (Pujol and Bliska, 2003), and *sifA* deletion was
1112 constructed using lambda red and validated using PCR (Datsenko and Wanner, 2000; Ko et al.,

1113 2009). *Candida albicans* SC5314 infection was performed as previously described (Odds et al.,
1114 2004) and levels of fibroblast growth factor 2 were measured using enzyme linked
1115 immunosorbent assays. *Staphylococcus aureus* toxin (alpha-hemolysin) was obtained from
1116 Sigma and applied to LCLs at a concentration of 1 μ g/ml for 23 hours. Cell death was measured
1117 using 7-AAD staining and flow cytometry. Additional experimental details can be found at
1118 <http://h2p2.oit.duke.edu/About/>.

1119 We estimated the gene-level effect size on the Hi-HOST phenotypes, using GWAS
1120 summary statistics (Barbeira et al., 2018) in each of the 44 GTEx tissues (version 6p) (Battle et
1121 al., 2017). The gene expression prediction model was trained using GTEx as the reference
1122 dataset (<https://zenodo.org/record/3572842/files/GTEx-V6p-HapMap-2016-09-08.tar.gz>). The
1123 gene-level effect size was estimated using S-PrediXcan after allele harmonization (Barbeira et
1124 al., 2018). We also applied MultiXcan to improve the ability to identify potential target genes
1125 (Barbeira et al., 2019). In brief, MultiXcan regresses the cellular trait on the principal
1126 components of the predicted expression data across all the available tissues. For each gene,
1127 MultiXcan yields a joint effect estimate across the 44 tissues. We applied the summary-statistic
1128 based version (S-MultiXcan) and followed the guides from the tool's webpage
1129 <https://github.com/hakyimlab/MetaXcan>.

1130

1131

1132 **DATA AND SOFTWARE AVAILABILITY**

1133 All code is available at the project's github page:

1134 <https://github.com/gamazonlab/infectiousDiseaseResource>. All trait-level GWAS,

1135 PrediXcan, and Hi-HOST TWAS results are available at www.phewascatalog.org.

1136

1137

1138 **KEY RESOURCES TABLE**

REAGENT or RESOURCE	SOURCE	IDENTIFIER
PrediXcan genetic associations for 35 ID traits	This paper	Available in Supplementary Materials.
BioVU	(Denny et al., 2010; Roden et al., 2008)	https://victr.vanderbilt.edu/pub/biovu/
PrediXcan	(Gamazon et al., 2018; Gamazon et al., 2015)	https://github.com/hakyimlab/PrediXcan
GTEx	(2015; Battle et al., 2017; Consortium, 2013)	https://gtexportal.org/home/
Gene Set Enrichment Analysis (GSEA)	(Subramanian et al., 2005)	http://software.broadinstitute.org/gsea/index.jsp
Mendelian Randomization software package	(Yavorska and Burgess, 2017)	https://cran.r-project.org/web/packages/MendelianRandomization/MendelianRandomization.pdf
Hi-HOST GWAS	(Ko et al., 2012; Wang et al., 2018)	http://h2p2.oit.duke.edu/About/
Hi-HOST TWAS	This paper.	Available in Supplementary Materials.

1139

1140