

1 **Leveraging genome and phenome-wide association studies to investigate genetic risk of**
2 **acute lymphoblastic leukemia**

3
4 Eleanor C. Semmes^{1,2}, Jayaram Vijayakrishnan³, Chenan Zhang PhD⁴, Jillian H. Hurst PhD²,
5 Richard S. Houlston³, Kyle M. Walsh PhD^{2,4,5,6*}

6
7 ¹ Medical Scientist Training Program, Duke University, Durham, NC, United States of America

8
9 ² Children's Health and Discovery Initiative, Department of Pediatrics, Duke University, Durham,
10 NC, United States of America

11
12 ³ Division of Genetics and Epidemiology, The Institute of Cancer Research, Sutton, Surrey SM2
13 5NG, United Kingdom

14
15 ⁴ Department of Epidemiology and Biostatistics, University of California, San Francisco, CA,
16 United States of America

17
18 ⁵ Department of Neurosurgery, Duke University, Durham, NC, United States of America

19
20 ⁶ Duke Cancer Institute, Duke University, Durham, NC, United States of America

21
22 Running title: Leveraging GWAS and PheWAS in childhood leukemia

23
24 Keywords: acute lymphoblastic leukemia, GWAS, PheWAS, childhood cancer, cancer
25 epidemiology

26
27
28 *corresponding author
29 Email: kyle.walsh@duke.edu (KMW)
30 Mailing address: DUMC Box 3050, Durham, NC 27710
31 Phone: (919) 684-8732
32 Fax: (919) 684-5207

33
34 Conflict of interest disclosure: The authors declare no competing interests.

35 Word count: 250 (abstract), 4,000 (manuscript excluding abstract and references)

36 Number of elements: 4 tables, 2 figures

37

38 **Financial support:** This study was supported by R21CA242439-01 (KMW), Alex's Lemonade
39 Stand Foundation "A" Awards (KMW), The Children's Health and Discovery Initiative of
40 Translating Duke Health (ECS, JHH, KMW), and NIH T32CA151022-06 (CZ).
41

42 **ABSTRACT**

43 **Background:** Genome-wide association studies (GWAS) of childhood cancers remain limited,
44 highlighting the need for novel analytic strategies. We describe a hybrid GWAS and phenome-
45 wide association study (PheWAS) approach to uncover genotype-phenotype relationships and
46 candidate risk loci, applying it to acute lymphoblastic leukemia (ALL).

47 **Methods:** PheWAS was performed for 12 ALL SNPs identified by prior GWAS and two control
48 SNP-sets using UK Biobank data. PheWAS-traits significantly associated with ALL SNPs
49 compared to control SNPs were assessed for association with ALL risk (959 cases, 2624
50 controls) using polygenic score and Mendelian randomization analyses. Trait-associated SNPs
51 were tested for association with ALL risk in single-SNP analyses, with replication in an
52 independent case-control dataset (1618 cases, 9409 controls).

53 **Results:** Platelet count was the trait most enriched for association with known ALL risk loci. A
54 polygenic score for platelet count (223 SNPs) was not associated with ALL risk ($P=0.82$) and
55 Mendelian randomization did not suggest a causal relationship. However, twelve platelet count-
56 associated SNPs were nominally associated with ALL risk in COG data and 3 were replicated in
57 UK data (rs10058074, rs210142, rs2836441).

58 **Conclusions:** In our hybrid GWAS-PheWAS approach, we identify pleiotropic genetic variation
59 contributing to ALL risk and platelet count. Three SNPs known to influence platelet count were
60 reproducibly associated with ALL risk, implicating genomic regions containing *IRF1*, pro-
61 apoptotic protein *BAK1*, and *ERG* in platelet production and leukemogenesis.

62 **Impact:** Incorporating PheWAS data into association studies can leverage genetic pleiotropy to
63 identify cancer risk loci, highlighting the utility of our novel approach.

64

65

66

67

68 **Introduction**

69 Genome-wide association studies (GWAS) have greatly enhanced our understanding of
70 inherited genetic susceptibility to cancer (1), but GWAS of pediatric cancers remain limited due
71 to lower disease incidence (2). Because of limited sample size, GWAS of childhood
72 malignancies are often underpowered to detect variants of small-to-moderate effect size,
73 preventing potentially important risk loci from reaching genome-wide statistical significance (*i.e.*
74 $P < 5.0 \times 10^{-8}$) (2, 3). Novel analytic approaches are needed to investigate how germline genetic
75 variation contributes to childhood cancer risk. The incorporation of polygenic scores (4, 5),
76 Mendelian randomization (MR) analyses, gene-pathway analyses (6), and phenome-wide
77 association studies (PheWAS) can augment traditional GWAS approaches to expand our
78 understanding of the genetic etiology of pediatric malignancies and other rare diseases.

79 PheWAS have not been widely applied to childhood cancer etiology research, but
80 represent a promising approach to understanding genetic risk in childhood cancer (7, 8). While
81 GWAS examine millions of genetic loci and test for association with a single phenotype or
82 disease, PheWAS test hundreds or thousands of phenotypes for association with a single
83 genetic variant, essentially a reversal of the GWAS paradigm (9, 10). This methodology has
84 recently become feasible through large collaborative efforts linking electronic health records
85 (EHR) data with high-throughput genomic data (7). Using PheWAS to discover additional traits
86 associated with cancer risk variants can reveal “intermediate phenotypes” (*e.g.*, height, smoking
87 behaviors) (4) that may mediate the relationship between SNPs and cancer development. Trait-
88 disease relationships can be further investigated using polygenic scores and MR approaches.
89 PheWAS data can also be integrated into case-control studies to identify trait-associated
90 genetic variants, create empirical candidate-SNP lists, and test for association with cancer case-
91 control status. Thus, integrating PheWAS and GWAS approaches in analyses of case-control
92 datasets may enhance our understanding of pathways driving pediatric cancer predisposition.

93 Acute lymphoblastic leukemia (ALL) is the most common childhood malignancy,
94 accounting for nearly one-third of pediatric cancers (11). Its etiology is complex, but the disease
95 is likely initiated *in utero*, with driver pre-leukemic fusion genes arising in lymphoid progenitors.
96 ALL development is also influenced by pre/postnatal environmental exposures (e.g., infections,
97 ionizing radiation) (11-14) and by germline genetic variants. GWAS have uncovered important
98 inherited genetic risk loci for ALL in hematopoietic transcription factors (*IKZF1*, *CEBPE*,
99 *ARID5B*, *GATA3*, *ELK3*), cell cycle regulators (*CDKN2A/CDKN2B*, *SP4*), and chromatin
100 remodeling enzymes (*BMI1*), though the precise mechanisms by which these GWAS-identified
101 risk loci influence leukemogenesis are not completely understood (15-22).

102 We have developed an integrated GWAS-PheWAS approach to identify candidate traits
103 and trait-associated variants that may modify cancer risk. We apply this methodology to ALL,
104 uncovering novel phenotypes associated with known ALL risk variants and pleiotropic ALL risk
105 loci, which we successfully replicate in an independent dataset. Our findings suggest that this
106 hybrid GWAS-PheWAS methodology is a promising new approach for deciphering germline
107 genetic risk in rare diseases, such as childhood cancers, where GWAS power remains limited.

108

109 **Materials and Methods**

110 **Prior GWAS ALL risk loci.** We accessed the NHGRI-EBI GWAS Catalog
111 (<https://www.ebi.ac.uk/gwas/>) to compile a list of variants previously identified by GWAS as
112 associated with B-cell precursor ALL risk in European-ancestry populations at genome-wide
113 statistical significance (*i.e.* $P < 5.0 \times 10^{-8}$) (access date: November 27, 2018) (23). We pruned this
114 list of significant variants for linkage disequilibrium ($R^2 \leq 0.15$ in European-ancestry populations)
115 using LDlink (24) and cross-referenced recent reviews on ALL GWAS (3), identifying 12
116 genome-wide significant independent ALL risk SNPs, which were included in our ALL SNP-set.

117

118 **Control SNP Sets.** We compiled 2 comparison SNP-sets to serve as controls for PheWAS
119 analyses. A set of unlinked control SNPs (1000 Genomes Project) was generated using
120 SNPsnap (Broad Institute) (25). Four control SNPs were matched to the 12 ALL risk SNPs on:
121 minor allele frequency ($\pm 5\%$), surrounding gene density ($\pm 50\%$), distance to nearest gene
122 ($\pm 50\%$) and, as a proxy for haplotype block size, the number of other SNPs in LD at $R^2 \geq 0.50$
123 ($\pm 50\%$). For several ALL risk SNPs, we could not generate more than 4 control SNPs without
124 loosening our matching parameters, but the gain in statistical power achieved beyond a case-to-
125 control ratio of 1:4 is minimal (26, 27).

126 Because ALL risk SNPs are trait-associated variants that may be more likely to
127 associate with additional traits in PheWAS analyses, we identified a second control SNP-set by
128 querying the GWAS catalog for chronic lymphocytic leukemia (CLL) risk SNPs. We used the
129 same methodology as for ALL risk SNPs, yielding 31 unlinked CLL-associated variants used as
130 another control SNP-set.

131
132 **eQTL and *in silico* SNP functional analyses.** We characterized ALL risk SNPs and control
133 SNP-sets using HaploReg to annotate chromatin state and regulatory motifs surrounding each
134 SNP (28). We examined whether variants were expression quantitative trait loci (eQTLs),
135 protein-binding, located in DNase hypersensitive sites, promoter or enhancer histone marks, or
136 predicted to change transcription factor binding motifs.

137
138 **UK Biobank GeneATLAS and PheWAS analyses.** The UK Biobank atlas of genetic
139 associations (<http://geneatlas.roslin.ed.au.uk/>) was constructed by genotyping 452,264
140 European-ancestry individuals for 805,426 genetic variants, performing genome-wide SNP
141 imputation and quality-controls, and linking genetic data to EHR data (29). GeneATLAS
142 contains data for 778 traits (118 quantitative, 660 binary) and associations with 9,113,133
143 genetic variants (genotyped or imputed). GeneATLAS is searchable and can be queried for

144 genetic (e.g. SNPs) or phenotypic (e.g., height) data to assess genotype-phenotype
145 associations (see Canela-Xandri *et. al.* for additional details) (29).

146 We queried GeneATLAS for trait associations with 12 known ALL risk SNPs and two
147 control SNP-sets (31 CLL-associated SNPs, 48 matched SNPsnap controls). Summary
148 statistics for traits associated with each queried variant were downloaded from GeneATLAS for
149 downstream analyses. Significant SNP-trait associations ($P < 0.01$) were carried forward in
150 subsequent SNP-set analyses. Although a more stringent p-value threshold for carrying SNP-
151 trait associations forward was considered (e.g. 0.05/778), this was determined to be too
152 conservative because many of the 778 traits in GeneATLAS have high genetic correlations with
153 each other (e.g. weight and hip circumference, 0.909; reticulocyte percentage and reticulocyte
154 count, 0.952). Additionally, these individual SNP-trait associations were carried forward for
155 SNP-set enrichment comparisons between ALL-associated SNPs and control SNP-sets, and as
156 such the PheWAS significance threshold is somewhat arbitrary so long as it is the same
157 threshold across all SNP-sets. PheWAS results for the 12 ALL risk SNPs and 778 traits were
158 compared to results for the two control SNP-sets using the R Statistical Programming
159 Environment (<http://www.R-project.org/>, version 3.5.2). Using Fisher's exact tests, we compared
160 PheWAS traits associated with >1 ALL SNP between the ALL and control SNP-sets to
161 determine if traits were enriched for association with known ALL risk variants.

162

163 **ALL case-control discovery cohort.** We included 959 European-ancestry ALL cases from the
164 Children's Oncology Group (COG) in our discovery dataset (16). Genotype data were
165 downloaded from dbGaP study accession phs000638.v1.p1, including ALL patients from COG
166 protocols 9904 and 9905 for whom DNA was obtained from remission blood samples (30).
167 Controls included 2624 European-ancestry subjects from the Wellcome Trust Case-Control
168 Consortium (<http://www.wtccc.org.uk/>) (31). Cases and controls were genotyped on the
169 Affymetrix 6.0 array. As described previously, genotyping quality-control (QC) measures were

170 implemented for cases and controls (16). We excluded samples or SNPs with genotyping call
171 rates <98%, individuals with suggested non-European-ancestry, IBD proportion >0.20, or with
172 discrepant sex between genotype and clinical report.

173

174 **Genotype imputation.** ALL case-control SNP data underwent genome-wide imputation as
175 previously described (5). Haplotype phasing was performed with SHAPEIT (version 2.790) (32),
176 and whole-genome imputation was performed using Minimac3 software (33) with 64,976 human
177 haplotypes from the Haplotype Reference Consortium (2016 release) as the reference panel
178 (34). SNPs with imputation quality (info) scores <0.60 or posterior probabilities <0.90 were
179 excluded (16).

180

181 **Platelet count polygenic score and single-SNP associations.** We constructed a polygenic
182 score for platelet count using 287 independent genetic variants associated with platelet count in
183 a prior GWAS of blood cell trait indices (223 were included after QC filtering) (35) (**S7 Table**).
184 The polygenic score for each individual in the ALL case-control dataset was determined based
185 on signed, weighted beta estimates for each platelet count-associated variant, as reported in
186 Astle *et. al.* (35) and calculated using the PLINK toolkit (36). We performed logistic regression
187 for the platelet count polygenic score, adjusting for 10 principal components (PCs). We also
188 tested platelet count-associated SNPs for association with ALL case-control status via single-
189 SNP association analyses.

190

191 **Mendelian randomization analyses.** To assess for a causal relationship between platelet
192 count and ALL risk, we performed formal MR analyses with the R package
193 “MendelianRandomization”(37, 38). Using summary statistics of SNP-exposure (*i.e.*, platelet
194 count) and SNP-outcome (*i.e.*, ALL) associations, we used the (1) inverse-variance weighted

195 (IVW), (2) MR-Egger, and (3) weighted median methods to test for a causal relationship
196 between platelet count and ALL risk in our case-control dataset (39, 40).

197

198 **ALL replication study.** The ALL replication dataset was a meta-analysis of two prior published
199 GWAS of B-cell precursor ALL, including German GWAS (834 cases, 2024 controls) (19) and
200 UK GWAS II (784 cases, 7,385 controls) (20). German cases were genotyped using Illumina
201 Human OmniExpress-12v1,0 arrays and controls were genotyped using the same platform or
202 Illumina-HumanOmni1-Quad1_v1. UK GWAS II cases and controls were genotyped using an
203 Illumina Infinium OncoArray-500K. Fixed-effects meta-analysis was used to estimate beta
204 values, standard errors, and p-values for queried risk loci in this combined GWAS meta-analysis
205 (1618 ALL cases, 9409 controls). For additional information on the GWAS meta-analysis used
206 for replication, see Vijayakrishnan *et. al.* (21).

207

208 **Results**

209 **Overview of methods.** An overview of the methodology applied in our study is displayed in
210 **Figure 1.** We used the GWAS catalog and a thorough literature review to identify known ALL
211 risk variants from GWAS of European-ancestry populations. PheWAS analyses were then
212 performed with the UK Biobank GeneATLAS database to test each ALL-associated variant and
213 control variant for association with 778 traits in the UK Biobank. After determining which traits
214 were enriched for association with the ALL SNP-set compared to control SNP-sets, we returned
215 to the GWAS catalog to identify SNPs associated with these traits. Using polygenic score, MR,
216 and candidate SNP approaches, we examined whether PheWAS-identified traits or trait-
217 associated variants conferred ALL risk, and replicated single-SNP associations in an
218 independent ALL case-control dataset.

219

220 **Risk variants for PheWAS analysis.** Using the GWAS catalog, we identified 12 independent
221 ($R^2 \leq 0.15$), genome-wide significant ($P < 5.0 \times 10^{-8}$) ALL risk SNPs (**Table 1**), which all previously
222 replicated in independent cohorts. Two SNP-sets served as controls for our PheWAS analyses,
223 including 31 SNPs previously associated with chronic lymphocytic leukemia (CLL) and 48
224 control SNPs matched to ALL risk SNPs on minor allele frequency, gene density, distance to
225 nearest gene, and number of SNPs in LD (25). Functional annotation and *in silico* analysis of
226 the ALL SNP-set and control SNP-sets demonstrated similar characteristics in terms of impact
227 on chromatin structure, including promoter and enhancer histone marks, DNase
228 hypersensitivity, and impact on regulatory motifs; however, ALL-associated and CLL-associated
229 SNPs were likelier to be eQTLs (**S1 Table**).

230

231 **UK Biobank PheWAS analyses.** We utilized the UK Biobank GeneAtlas database to conduct
232 a PheWAS for 12 ALL-associated SNPs (**S2 Table**), 31 CLL-associated SNPs (**S3 Table**), and
233 48 matched control SNPs (**S4 Table**) to test for association with 778 traits. We used the same
234 PheWAS approach and nominal significance threshold to identify SNP-trait associations (*i.e.*
235 $P < 0.01$) for ALL and control SNPs. The proportion of variants in each SNP-set (12 ALL-
236 associated, 31 CLL-associated, 48 matched control) that was associated with a particular
237 PheWAS trait was compared across groups to ascertain phenotypes enriched for association
238 with ALL SNPs compared with control SNPs (**S5 Table**).

239 We determined that 76 of the 778 traits in the database were nominally associated
240 ($P < 0.01$) with >1 of the 12 ALL risk SNPs. PheWAS traits significantly associated with >1 ALL
241 risk SNP were carried forward for enrichment comparisons between ALL and control SNP-sets
242 (**S6 Table**). All 76 PheWAS traits compared between SNP-sets are depicted in **Figure 2**
243 showing the relative proportion of significant SNP-trait associations in each SNP-set. Platelet
244 count was the phenotype most enriched for association with ALL risk variants. Specifically, 9 out
245 of 12 (75%) ALL SNPs were nominally associated with platelet count, compared to 11 of 31

246 (35.5%) CLL SNPs ($P = 0.047$) and 6 of 48 (12.5%) control SNPs ($P < 0.001$) (**Table 2**).

247 Notably, many of the PheWAS-identified traits were enriched for association in the ALL SNPs
248 compared to the control SNPs, but only 5 traits were significantly enriched for association with
249 ALL SNPs compared to both control SNPs *and* CLL SNPs, and platelet count was associated
250 with the highest proportion of ALL SNPs (**Table 2**).

251

252 **Platelet count polygenic score analyses.** Given that platelet count was the trait most enriched
253 for association with ALL SNPs in PheWAS analyses, we constructed a polygenic score for
254 platelet count using 287 previously-published variants from a recent GWAS on blood cell indices
255 (35) (**S7 Table**). Of these, 223 SNPs were successfully imputed (info score ≥ 0.60 , posterior
256 probability ≥ 0.90) in our ALL case-control dataset (959 cases, 2624 controls) and used in
257 polygenic score construction. The polygenic score for platelet count was not associated with
258 ALL case-control status in a logistic regression model adjusting for sex and 10 PCs ($P=0.819$).

259

260 **Mendelian randomization analyses.** To test for a causal relationship between platelet count
261 and ALL risk, we used several MR analytical approaches wherein genetic variants are used as
262 instrumental variables to assess causality in exposure/risk factor associations. Estimates from
263 IVW ($P_{IVW}=0.948$), MR-Egger ($P_{MR-Egger}=0.857$, $P_{MR-intercept}=0.912$), and median-weighted (P_{MR-}
264 $_{median}=0.857$) MR methods were non-significant and consistent with the null polygenic score
265 results. These MR results suggest that platelet count does not mediate ALL risk and that there
266 is no causal relationship between these two traits.

267

268 **Platelet count-associated SNPs as candidate ALL risk loci.** To examine whether individual
269 platelet count-associated variants might have pleiotropic effects on ALL risk, we performed
270 single-SNP association analyses for 223 platelet count-associated SNPs in 959 ALL cases and
271 2624 controls (**S8 Table**). Twelve SNPs were nominally associated ($P<0.05$) with ALL case-

272 control status (notably, not more than expected by chance) after adjusting for sex and 10 PCs
273 (**Table 3**). The directional effect of platelet count-associated alleles (*i.e.*, increased versus
274 decreased platelet count) did not correlate with the direction of effect on ALL susceptibility (*i.e.*,
275 protection versus risk).

276 These 12 candidate SNPs were carried forward for evaluation in an independent UK ALL
277 case-control dataset (1618 cases, 9409 controls). Nine SNPs were successfully genotyped or
278 imputed in this dataset, of which three associations were successfully replicated at $P < 5.6 \times 10^{-3}$
279 (*i.e.* 0.05/9) (**Table 4**). SNPs had similar magnitudes of effect in the UK ALL case-control and
280 discovery data. The replicated variants map to 3 distinct genomic loci on 5q31.1, 6p21.31, and
281 21q22.2 (**Table 4**). To interrogate these risk loci further, we identified the genes in which these
282 variants resided and associated genes for which these variants were expression quantitative
283 trait loci (eQTLs) (28). We found that the 5q31.1 region was adjacent to *IRF1*, a gene encoding
284 interferon regulatory factor 1, which regulates host immune responses, including interferon
285 signaling. The 6p21.31 region includes *BAK1*, a pro-apoptotic protein known to be disrupted in
286 adult-onset malignancies. Finally, the 21q22.2 region encodes the hematopoietic transcription
287 factor *ERG*, known to be associated with ALL risk in Hispanics and children with Trisomy 21 (*i.e.*
288 Down Syndrome) (41, 42).

289

290 Discussion

291 We provide a novel framework (**Figure 1**) for leveraging existing GWAS and PheWAS
292 data to uncover traits associated with known disease risk variants and to identify trait-associated
293 variants as possible candidate risk loci. We apply this framework to an investigation of ALL
294 predisposition that combines rich genotype-phenotype data available from the UK Biobank with
295 ALL case-control analyses. We first identify SNPs associated with ALL using the GWAS
296 catalog. We then perform PheWAS on these SNPs and control SNP-sets using the UK Biobank
297 GeneATLAS, identifying platelet count as the trait most enriched for association with ALL risk

298 loci. Returning to the GWAS catalog, we identify genetic determinants of platelet count. We then
299 use a two-stage case-control design (43, 44) to examine whether SNPs associated with platelet
300 count modify ALL risk, confirming three risk loci near *IRF1*, *BAK1*, and *ERG*.

301 Potential modifications to this hybrid GWAS-PheWAS approach could be implemented in
302 future applications based on features of the cancer undergoing analysis and the datasets
303 available. For cancers with many known GWAS hits (e.g. breast cancer), it may be preferable to
304 use a more stringent p-value threshold for the PheWAS analysis to streamline subsequent SNP-
305 set enrichment comparisons. Similarly, trait-associated SNPs could be evaluated for their
306 association with cancer using a more stringent p-value threshold in a one-stage case-control
307 design when sample sizes are large or when replication sets are unavailable.

308 We identified platelet count as significantly enriched for association with ALL risk SNPs;
309 however, our results did not suggest a direct role for platelet count in mediating ALL risk, as the
310 polygenic score for platelet count was not associated with ALL case-control status. Null results
311 from MR analyses also support the conclusion that there is no causal relationship between
312 platelet count and ALL. This indicates that platelet count and ALL may have overlapping genetic
313 architecture due to pleiotropic loci independently influencing both traits, which appears
314 reasonable since regulatory variants in hematopoietic transcription factors could influence each
315 phenotype. This interpretation is supported by our single-SNP association results, identifying
316 and replicating 3 ALL risk loci using platelet count-associated variants as candidate SNPs. Two
317 of these ALL risk alleles were associated with higher platelet count (rs10058074 near *IRF1*,
318 rs210142 in *BAK1*), whereas one ALL risk allele was associated with reduced platelet count
319 (rs2836441 in *ERG*). In addition to hematopoietic transcription factor genes, pleiotropic variants
320 in cell-cycle regulators are also candidate modifiers of both platelet count and ALL risk, as
321 supported by our identification of a shared locus in pro-apoptotic protein *BAK1*.

322 The ALL risk SNP that we identify at 5q31.1 (rs10058074) is intronic, but has suggestive
323 functional significance as a *cis*-acting eQTL for *IRF1*, a master transcriptional regulator of

324 immune response and oncogenesis (45, 46), as well as for *PDLIM4*, an F-actin-binding protein
325 that influences T cell trafficking (47). This is one of the first ALL risk loci found that is related to
326 “immune response gene elements”, long posited to be important based on a hypothesized
327 infectious etiology for ALL (11). Interferon immune responses are particularly important for
328 controlling viral pathogens, which is notable since congenital cytomegalovirus infection was
329 recently associated with ALL risk (48, 49). Two associated variants at 6p21.31 (rs210142,
330 rs75080135) are also intronic, but both are *cis*-acting eQTLs for *BAK1* (BCL2 antagonist killer 1)
331 (50), which encodes a pro-apoptotic protein that is a known CLL GWAS hit (51) and important
332 for B cell homeostasis (52). Located 6kb apart and both associated with ALL risk in our
333 discovery analysis, rs210142 and rs75080135 are in only weak LD in 1000 Genomes
334 Europeans ($R^2=0.11$) and were both associated with ALL risk in UK replication data, although
335 only the rs210142 association survived Bonferroni correction. The associated SNP at 21q22.2,
336 rs2836441, is located in the 5' untranslated region of *ERG*, a transcription factor from the
337 erythroblast transformation-specific family that is frequently deleted or alternatively spliced in the
338 DUX4-rearranged ALL subtype (53).

339 Since our analyses were completed, the largest ALL GWAS to-date has been published
340 by Vijayakrishnan and colleagues (54). This meta-analysis of four GWAS totaling 5,321 cases
341 and 16,666 controls identified 4 novel B-ALL risk loci reaching genome-wide significance (54).
342 Of these 4 loci, one (9q21.31) was significant for B-ALL risk overall, two (5q31.1, 6p21.31) for
343 the high-hyperdiploid subtype, and one (17q21.32) for the ETV6-RUNX fusion subtype. Notably,
344 their 5q31.1 and 6p21.31 risk loci overlap substantially with those identified through our hybrid
345 GWAS-PheWAS approach. Their lead SNP at 5q31.1 (rs886285), located in C5orf56, is in weak
346 LD ($R^2=0.17$) with the SNP (rs10058074) discovered through our approach, yet both variants
347 appear to modulate expression of the master transcription factor *IRF1*. Compellingly, their lead
348 SNP at 6p21.31 (rs210143) is only ~100 bases away ($R^2=0.95$) from the SNP identified through
349 our approach (rs210142), and they too detected multiple signals in *BAK1* that implicate

350 decreased expression of this pro-apoptotic protein as an important hallmark of leukemogenesis.
351 The GWAS meta-analysis from Vijayakrishnan and colleagues also confirmed *ERG* (21q22.2)
352 as an ALL risk locus in European-ancestry populations, which we and others had previously
353 identified as a GWAS hit for ALL in Hispanic populations (41, 42), but had been unable to
354 replicate in European-ancestry populations.

355 While risk loci at 5q31.1, 6p21.31, and 21q22.2 were very recently associated with the
356 high-hyperdiploid subtype of ALL at genome-wide significance, our results suggest that these
357 susceptibility loci may influence ALL risk overall – not just subtype-specific risk – and may also
358 be broadly involved in non-malignant hematopoiesis. The fact that the same risk loci identified
359 through a largescale collaborative GWAS and a recent Hispanic ALL GWAS were uncovered
360 through our combined GWAS-PheWAS methodology, despite our limited sample size, confirms
361 the utility of the approach we have developed. These results provide further evidence of the
362 importance of these loci in B-cell ALL and suggests our approach has applicability to the study
363 of rare malignancies, including childhood cancers.

364 There are several limitations to our study and valid concerns of our hybrid GWAS-
365 PheWAS approach. One limitation of using the UK Biobank for this study investigating genetic
366 risk in a pediatric cancer is that the UK Biobank GeneATLAS PheWAS database was
367 constructed using genetic and EHR data from adults 40-69 years of age. Thus, applying this
368 approach to rare adult-onset diseases may be more appropriate than for pediatric diseases, as
369 the overlap between these adult traits and pediatric phenotypes are largely unknown. For the
370 PheWAS analyses, we used a p-value threshold of <0.01 to carry a SNP-trait association
371 forward to enrichment analyses, rather than a Bonferroni-corrected threshold (*i.e.* $0.05/778$
372 traits). While many of the traits in the GeneATLAS are highly correlated (*e.g.* standing height
373 and sitting height, BMI and waist-to-hip ratio), cancers that have more than just 12 GWAS hits to
374 evaluate via PheWAS may benefit from a more stringent threshold. Another significant limitation
375 of our study was the limited case sample size in our discovery dataset. Because of our limited

376 sample size, we implemented a two-stage study design, first screening for nominally-associated
377 SNPs in our discovery dataset ($P < 0.05$), and then attempting replication in an independent
378 sample. Despite these limitations, interest in our hybrid GWAS-PheWAS approach for
379 investigating inherited genetic risk in rare diseases, where traditional approaches remain limited,
380 appears warranted.

381 Although multiple GWAS in the past decade have contributed to our understanding of
382 inherited susceptibility to ALL, there remains significant missing heritability (55). The most
383 recent and largest ALL GWAS determined that known risk alleles accounted for 31% of the total
384 variance in genetic risk of ALL (54); thus, there is a need for additional studies investigating ALL
385 genetic risk loci. A recent review on the benefits and pitfalls GWAS emphasized the need for
386 novel analytic approaches to enhance our understanding of genotype-phenotype associations in
387 the post-GWAS era and the utility of large biorepository databases linking EHR and genotyping
388 data, polygenic scores, and innovative study designs (8). This review also highlighted that, while
389 increasing GWAS sample size may reveal more associations, new methods for analyzing the
390 wealth of existing data are essential (8).

391 One opportunity would be to leverage LD score regression to identify traits associated
392 with a cancer of interest. Although the sample-size limitations that apply to our study would also
393 apply to analyses using LD score regression, replacing the PheWAS portion of our methodology
394 with LD score regression is an intriguing approach for identifying traits with shared genetic
395 determinants in future applications. In summary, our novel hybrid application of PheWAS
396 represents a promising approach to investigate inherited genetic risk, especially in childhood
397 cancers where GWAS remain underpowered and where innovative analytic strategies can help
398 to decipher complex etiology and guide future prevention and screening strategies.

399

400 **ACKNOWLEDGEMENTS**

401 The ALL Relapse GWAS dataset was generated at St. Jude Children's Research
402 Hospital and by the Children's Oncology Group, supported by NIH grants CA142665, CA21765,
403 CA158568, CA156449, CA36401, CA98543, CA114766, CA140729, and U01GM92666, Jeffrey
404 Pride Foundation, the National Childhood Cancer Foundation, and by ALSAC. Funding for the
405 project was provided by the Wellcome Trust under award 076113 and 085475.
406

407 **References**

- 408 1. Sud A, Kinnersley B, Houlston RS. Genome-wide association studies of cancer: current
409 insights and future perspectives. *Nature reviews Cancer*. 2017;17(11):692-704.
- 410 2. Enciso-Mora V, Hosking FJ, Sheridan E, Kinsey SE, Lightfoot T, Roman E, et al.
411 Common genetic variation contributes significantly to the risk of childhood B-cell precursor acute
412 lymphoblastic leukemia. *Leukemia*. 2012;26(10):2212-5.
- 413 3. Plon SE, Lupo PJ. Genetic Predisposition to Childhood Cancer in the Genomic Era.
414 *Annu Rev Genomics Hum Genet*. 2019;20:241-63.
- 415 4. Semmes EC, Zhang C, Walsh KM. Intermediate phenotypes underlying osteosarcoma
416 risk. *Oncotarget*. 2018;9(100):37345-6.
- 417 5. Zhang C, Morimoto LM, de Smith AJ, Hansen HM, Gonzalez-Maya J, Endicott AA, et al.
418 Genetic determinants of childhood and adult height associated with osteosarcoma risk. *Cancer*.
419 2018;124(18):3742-52.
- 420 6. Wang K, Li M, Hakonarson H. Analysing biological pathways in genome-wide
421 association studies. *Nature reviews Genetics*. 2010;11(12):843-54.
- 422 7. Bush WS, Oetjens MT, Crawford DC. Unravelling the human genome-phenome
423 relationship using phenome-wide association studies. *Nature reviews Genetics*. 2016;17(3):129-
424 45.
- 425 8. Tam V, Patel N, Turcotte M, Bosse Y, Pare G, Meyre D. Benefits and limitations of
426 genome-wide association studies. *Nature reviews Genetics*. 2019.
- 427 9. Denny JC, Bastarache L, Ritchie MD, Carroll RJ, Zink R, Mosley JD, et al. Systematic
428 comparison of phenome-wide association study of electronic medical record data and genome-
429 wide association study data. *Nat Biotechnol*. 2013;31(12):1102-10.
- 430 10. Denny JC, Bastarache L, Roden DM. Phenome-Wide Association Studies as a Tool to
431 Advance Precision Medicine. *Annu Rev Genomics Hum Genet*. 2016;17:353-73.
- 432 11. Greaves M. A causal mechanism for childhood acute lymphoblastic leukaemia. *Nature*
433 *reviews Cancer*. 2018;18(8):471-84.
- 434 12. Inaba H, Greaves M, Mullighan CG. Acute lymphoblastic leukaemia. *Lancet (London,*
435 *England)*. 2013;381(9881):1943-55.
- 436 13. Mori H, Colman SM, Xiao Z, Ford AM, Healy LE, Donaldson C, et al. Chromosome
437 translocations and covert leukemic clones are generated during normal fetal development.
438 *Proceedings of the National Academy of Sciences of the United States of America*.
439 2002;99(12):8242-7.
- 440 14. Iacobucci I, Mullighan CG. Genetic Basis of Acute Lymphoblastic Leukemia. *Journal of*
441 *clinical oncology : official journal of the American Society of Clinical Oncology*. 2017;35(9):975-
442 83.
- 443 15. Trevino LR, Yang W, French D, Hunger SP, Carroll WL, Devidas M, et al. Germline
444 genomic variants associated with childhood acute lymphoblastic leukemia. *Nature genetics*.
445 2009;41(9):1001-5.
- 446 16. Wiemels JL, Walsh KM, de Smith AJ, Metayer C, Gonseth S, Hansen HM, et al. GWAS
447 in childhood acute lymphoblastic leukemia reveals novel genetic associations at chromosomes
448 17q12 and 8q24.21. *Nature communications*. 2018;9(1):286.
- 449 17. de Smith AJ, Walsh KM, Francis SS, Zhang C, Hansen HM, Smirnov I, et al. BMI1
450 enhancer polymorphism underlies chromosome 10p12.31 association with childhood acute
451 lymphoblastic leukemia. *International journal of cancer*. 2018;143(11):2647-58.
- 452 18. Xu H, Zhang H, Yang W, Yadav R, Morrison AC, Qian M, et al. Inherited coding variants
453 at the CDKN2A locus influence susceptibility to acute lymphoblastic leukaemia in children.
454 *Nature communications*. 2015;6:7553.

- 455 19. Migliorini G, Fiege B, Hosking FJ, Ma Y, Kumar R, Sherborne AL, et al. Variation at
456 10p12.2 and 10p14 influences risk of childhood B-cell acute lymphoblastic leukemia and
457 phenotype. *Blood*. 2013;122(19):3298-307.
- 458 20. Papaemmanuil E, Hosking FJ, Vijayakrishnan J, Price A, Olver B, Sheridan E, et al. Loci
459 on 7p12.2, 10q21.2 and 14q11.2 are associated with risk of childhood acute lymphoblastic
460 leukemia. *Nature genetics*. 2009;41(9):1006-10.
- 461 21. Vijayakrishnan J, Studd J, Broderick P, Kinnersley B, Holroyd A, Law PJ, et al. Genome-
462 wide association study identifies susceptibility loci for B-cell childhood acute lymphoblastic
463 leukemia. *Nature communications*. 2018;9(1):1340.
- 464 22. Vijayakrishnan J, Kumar R, Henrion MY, Moorman AV, Rachakonda PS, Hosen I, et al.
465 A genome-wide association study identifies risk loci for childhood acute lymphoblastic leukemia
466 at 10q26.13 and 12q23.1. *Leukemia*. 2017;31(3):573-9.
- 467 23. MacArthur J, Bowler E, Cerezo M, Gil L, Hall P, Hastings E, et al. The new NHGRI-EBI
468 Catalog of published genome-wide association studies (GWAS Catalog). *Nucleic acids
469 research*. 2017;45(D1):D896-d901.
- 470 24. Machiela MJ, Chanock SJ. LDlink: a web-based application for exploring population-
471 specific haplotype structure and linking correlated alleles of possible functional variants.
472 *Bioinformatics (Oxford, England)*. 2015;31(21):3555-7.
- 473 25. Pers TH, Timshel P, Hirschhorn JN. SNPsnap: a Web-based tool for identification and
474 annotation of matched SNPs. *Bioinformatics (Oxford, England)*. 2015;31(3):418-20.
- 475 26. Rothman K. *Modern Epidemiology*. Boston, MA: Little, Brown and Company; 1986.
- 476 27. Hennessy S, Bilker WB, Berlin JA, Strom BL. Factors influencing the optimal control-to-
477 case ratio in matched case-control studies. *American journal of epidemiology*. 1999;149(2):195-
478 7.
- 479 28. Ward LD, Kellis M. HaploReg: a resource for exploring chromatin states, conservation,
480 and regulatory motif alterations within sets of genetically linked variants. *Nucleic acids research*.
481 2012;40(Database issue):D930-4.
- 482 29. Canela-Xandri O, Rawlik K, Tenesa A. An atlas of genetic associations in UK Biobank.
483 *Nat Genet*. 2018;50(11):1593-9.
- 484 30. Yang JJ, Cheng C, Devidas M, Cao X, Campana D, Yang W, et al. Genome-wide
485 association study identifies germline polymorphisms associated with relapse of childhood acute
486 lymphoblastic leukemia. *Blood*. 2012;120(20):4197-204.
- 487 31. Genome-wide association study of 14,000 cases of seven common diseases and 3,000
488 shared controls. *Nature*. 2007;447(7145):661-78.
- 489 32. O'Connell J, Gurdasani D, Delaneau O, Pirastu N, Ulivi S, Cocca M, et al. A general
490 approach for haplotype phasing across the full spectrum of relatedness. *PLoS genetics*.
491 2014;10(4):e1004234.
- 492 33. Das S, Forer L, Schonherr S, Sidore C, Locke AE, Kwong A, et al. Next-generation
493 genotype imputation service and methods. *Nature genetics*. 2016;48(10):1284-7.
- 494 34. McCarthy S, Das S, Kretzschmar W, Delaneau O, Wood AR, Teumer A, et al. A
495 reference panel of 64,976 haplotypes for genotype imputation. *Nature genetics*.
496 2016;48(10):1279-83.
- 497 35. Astle WJ, Elding H, Jiang T, Allen D, Ruklisa D, Mann AL, et al. The Allelic Landscape of
498 Human Blood Cell Trait Variation and Links to Common Complex Disease. *Cell*.
499 2016;167(5):1415-29.e19.
- 500 36. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, et al. PLINK: a
501 tool set for whole-genome association and population-based linkage analyses. *American journal
502 of human genetics*. 2007;81(3):559-75.
- 503 37. Yavorska OO, Burgess S. MendelianRandomization: an R package for performing
504 Mendelian randomization analyses using summarized data. *Int J Epidemiol*. 2017;46(6):1734-9.

- 505 38. Burgess S, Thompson SG. Use of allele scores as instrumental variables for Mendelian
506 randomization. *Int J Epidemiol*. 2013;42(4):1134-44.
- 507 39. Bowden J, Davey Smith G, Haycock PC, Burgess S. Consistent Estimation in Mendelian
508 Randomization with Some Invalid Instruments Using a Weighted Median Estimator. *Genet*
509 *Epidemiol*. 2016;40(4):304-14.
- 510 40. Bowden J, Davey Smith G, Burgess S. Mendelian randomization with invalid
511 instruments: effect estimation and bias detection through Egger regression. *Int J Epidemiol*.
512 2015;44(2):512-25.
- 513 41. de Smith AJ, Walsh KM, Morimoto LM, Francis SS, Hansen HM, Jeon S, et al. Heritable
514 variation at the chromosome 21 gene *ERG* is associated with acute lymphoblastic leukemia risk
515 in children with and without Down syndrome. *Leukemia*. 2019;33(11):2746-51.
- 516 42. Qian M, Xu H, Perez-Andreu V, Roberts KG, Zhang H, Yang W, et al. Novel
517 susceptibility variants at the *ERG* locus for childhood acute lymphoblastic leukemia in
518 Hispanics. *Blood*. 2018.
- 519 43. Stanhope SA, Skol AD. Improved minimum cost and maximum power two stage
520 genome-wide association study designs. *PloS one*. 2012;7(9):e42367.
- 521 44. Wason JM, Dudbridge F. A general framework for two-stage analysis of genome-wide
522 association studies and its application to case-control studies. *American journal of human*
523 *genetics*. 2012;90(5):760-73.
- 524 45. Taniguchi T, Lamphier MS, Tanaka N. IRF-1: the transcription factor linking the
525 interferon response and oncogenesis. *Biochim Biophys Acta*. 1997;1333(1):M9-17.
- 526 46. Willman CL, Sever CE, Pallavicini MG, Harada H, Tanaka N, Slovak ML, et al. Deletion
527 of IRF-1, mapping to chromosome 5q31.1, in human leukemia and preleukemic myelodysplasia.
528 *Science (New York, NY)*. 1993;259(5097):968-71.
- 529 47. Fu C, Li Q, Zou J, Xing C, Luo M, Yin B, et al. JMJD3 regulates CD4 T cell trafficking by
530 targeting actin cytoskeleton regulatory gene *Pdlim4*. *The Journal of clinical investigation*.
531 2019;130:4745-57.
- 532 48. Wiemels JL, Talback M, Francis SS, Feychting M. Early infection with cytomegalovirus
533 and risk of childhood hematological malignancies. *Cancer epidemiology, biomarkers &*
534 *prevention : a publication of the American Association for Cancer Research, cosponsored by*
535 *the American Society of Preventive Oncology*. 2019.
- 536 49. Francis SS, Wallace AD, Wendt GA, Li L, Liu F, Riley LW, et al. In utero cytomegalovirus
537 infection and development of childhood acute lymphoblastic leukemia. *Blood*.
538 2017;129(12):1680-4.
- 539 50. Chittenden T, Harrington EA, O'Connor R, Flemington C, Lutz RJ, Evan GI, et al.
540 Induction of apoptosis by the Bcl-2 homologue Bak. *Nature*. 1995;374(6524):733-6.
- 541 51. Slager SL, Skibola CF, Di Bernardo MC, Conde L, Broderick P, McDonnell SK, et al.
542 Common variation at 6p21.31 (*BAK1*) influences the risk of chronic lymphocytic leukemia.
543 *Blood*. 2012;120(4):843-6.
- 544 52. Takeuchi O, Fisher J, Suh H, Harada H, Malynn BA, Korsmeyer SJ. Essential role of
545 BAX, BAK in B cell homeostasis and prevention of autoimmune disease. *Proceedings of the*
546 *National Academy of Sciences of the United States of America*. 2005;102(32):11272-7.
- 547 53. Zhang J, McCastlain K, Yoshihara H, Xu B, Chang Y, Churchman ML, et al.
548 Deregulation of DUX4 and *ERG* in acute lymphoblastic leukemia. *Nature genetics*.
549 2016;48(12):1481-9.
- 550 54. Vijayakrishnan J, Qian M, Studd JB, Yang W, Kinnersley B, Law PJ, et al. Identification
551 of four novel associations for B-cell acute lymphoblastic leukaemia risk. *Nature*
552 *communications*. 2019;10(1):5348.
- 553 55. Blanco-Gomez A, Castillo-Lluva S, Del Mar Saez-Freire M, Hontecillas-Prieto L, Mao JH,
554 Castellanos-Martin A, et al. Missing heritability of complex diseases: Enlightenment by genetic

555 variants from intermediate phenotypes. *BioEssays* : news and reviews in molecular, cellular and
556 developmental biology. 2016;38(7):664-73.
557 56. Buniello A, MacArthur JAL, Cerezo M, Harris LW, Hayhurst J, Malangone C, et al. The
558 NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and
559 summary statistics 2019. *Nucleic acids research*. 2019;47(D1):D1005-d12.

560 **Table 1. Summary of previously published genome-wide significant risk loci for B-cell ALL**

561

Author (year reported) in ALL GWAS	Locus	ALL risk SNP	Gene	OR (95% CI)
Trevino LR <i>et al.</i> (2009)	7p12.2	rs11978267	<i>IKZF1</i>	1.69 (1.40-1.90)
Wiemels JL <i>et al.</i> (2018)	7p15.3	rs2390536	<i>SP4</i>	1.18 (1.11-1.24)
Wiemels JL <i>et al.</i> (2018)	8q24.21	rs4617118	Intergenic	1.34 (1.21-1.47)
Xu H <i>et al.</i> (2015)	9p21.3	rs3731249	<i>CDKN2A, CDKN2B</i>	1.63 (1.18-1.56)
de Smith AJ <i>et al.</i> (2018)	10p12.2	rs10741006	<i>PIP4K2A</i>	1.40 (1.40-1.53)
de Smith AJ <i>et al.</i> (2018)	10p12.31	rs12769953	<i>BMI1</i>	1.27 (1.20-1.35)
Migliorini G <i>et al.</i> (2013)	10p14	rs3824662	<i>GATA3</i>	1.31 (1.21-1.41)
Papaemmanuil E <i>et al.</i> (2009)	10q21.2	rs7089424	<i>ARID5B</i>	1.65 (1.54-1.36)
Wiemels JL <i>et al.</i> (2018)	10q26.13	rs3740540	<i>LHPP</i>	1.20 (1.15-1.28)
Vijayakrishnan J <i>et al.</i> (2017)	12q23.1	rs4762284	<i>ELK3</i>	1.19 (1.12-1.26)
Papaemmanuil E <i>et al.</i> (2009)	14q11.2	rs2239633	<i>CEBPE</i>	1.34 (1.22-1.41)
Wiemels JL <i>et al.</i> (2018)	17q21.1	rs2290400	<i>IKZF3</i>	1.18 (1.11-1.25)

562

563 Abbreviations: ALL, acute lymphoblastic leukemia; SNP, single nucleotide polymorphism; GWAS, genome-wide association study; OR,
 564 odds ratio; 95% CI, 95% confidence interval
 565 rsIDs from GRCh37/hg19 build

566

567

568 **Table 2. Selected PheWAS traits compared between ALL SNP-set and control SNP-sets for enrichment^a**
 569

UK Biobank PheWAS Trait	ALL vs. CLL SNP-set	ALL vs. matched control SNP-set ^b	UK Biobank PheWAS Trait	ALL vs. CLL SNP-set	ALL vs. matched control SNP-set
	<i>P</i> ^c	<i>P</i> ^c		<i>P</i> ^c	<i>P</i> ^c
platelet count	0.047	<0.001	weight	0.737	0.035
lymphocyte count	0.191	<0.001	white blood cell count	0.865	0.035
monocyte percentage	0.033	<0.001	whole body fat mass	0.568	0.004
monocyte count	0.309	0.002	whole body fat-free mass	0.737	0.035
neutrophil percentage	0.531	0.001	whole body water mass	0.737	0.035
platelet crit	0.664	<0.001	asthma	0.815	0.393
eosinophil count	0.892	0.002	body fat percentage	0.815	0.080
impedance of whole body	0.103	0.012	body mass index	0.615	0.393
standing height	0.248	0.039	melanoma/malignant skin neoplasms	1.000	0.028
basophil count	0.048	<0.001	hematocrit percentage	0.815	0.080
comparative height (age 10)	0.169	<0.001	hemoglobin concentration	1.000	0.080
eosinophil percentage	0.956	0.006	hip circumference	1.000	0.080
neutrophil count	1.000	0.015	chronic rheumatic heart diseases	0.026	0.028
basal metabolic rate	0.737	0.013	multiple valve diseases	0.026	0.028
basophil percentage	0.073	0.013	mean reticulocyte volume	0.387	0.393
lymphocyte percentage	0.583	0.071	mean sphered cell volume	0.156	0.162
mean platelet volume	1.000	0.124	oily fish intake	0.418	0.080
platelet distribution width	1.000	0.071	trunk fat mass	1.000	0.080
sitting height	1.000	0.035	waist circumference	1.000	0.028
trunk fat-free mass	1.000	0.124	water intake	0.105	0.005
trunk predicted mass	1.000	0.124	alcohol intake frequency	0.376	0.747

570 Abbreviations: ALL, acute lymphoblastic leukemia; SNPs, single nucleotide polymorphisms; CLL, chronic lymphocytic leukemia
 571 Bold values indicate nominal significance ($P < 0.05$)
 572

573
 574 ^a Individual SNP-trait associations available in S2 Table (ALL SNP-set), S3 Table (CLL SNP-set) and S4 Table (SNPsnap SNP-set)
 575 ^b Control SNPs generated using SNPsnap matched to ALL SNPs based on minor allele frequency ($\pm 5\%$), surrounding gene density
 576 ($\pm 50\%$), distance to nearest gene ($\pm 50\%$), and linkage disequilibrium at $R^2 \geq 0.50$ ($\pm 50\%$).
 577 ^c *P* value calculated with fisher's exact test, summary of all 76 PheWAS traits tested for enrichment in S6 Table

578
579

Table 3. Multivariate logistic regression of platelet count-associated variants and ALL risk in discovery case-control cohort^a

Locus	SNP rsID	Effect allele ^b	EAF ^c	Gene	OR ^e (95% CI)	P
2q32.3	rs7585866	G	0.37	<i>SDPR</i>	0.87 (0.77-0.99)	0.041
4q24	rs4699154	C	0.72	near <i>TET2</i> ^d	1.22 (1.06-1.39)	3.80 x 10⁻³
5q31.1	rs10058074	G	0.57	near <i>IRF1</i> ^d	1.15 (1.02-1.30)	0.023
6p21.31	rs210142	C	0.73	<i>BAK1</i>	1.17 (1.02-1.33)	0.021
6p21.31	rs75080135	C	0.24	<i>GGNBP1</i>	1.20 (1.02-1.40)	0.024
6q23.3	rs1331308	C	0.51	<i>HBS1L</i>	1.17 (1.04-1.32)	0.011
6q23.3	rs7776054	G	0.27	<i>HBS1L</i>	0.84 (0.73-0.97)	0.016
7q32.2	rs11556924	T	0.38	<i>ZC3HC1</i>	0.88 (0.77-0.99)	0.034
12q24.21	rs35427	T	0.60	intergenic	0.87 (0.76-0.98)	0.026
18q12.3	rs16977972	T	0.17	<i>SETBP1</i>	1.20 (1.01-1.41)	0.034
21q22.2	rs2836441	G	0.11	<i>ERG</i>	0.81 (0.67-0.96)	0.019
22q11.21	rs1059196	C	0.66	<i>SEPT5, GP1BB</i>	1.15 (1.00-1.32)	0.044

580

581 Abbreviations: SNP, single nucleotide polymorphism; EAF, effect allele frequency; OR, odds ratio; 95% CI, 95% confidence interval
582 Bold values indicate nominal significance ($P < 0.05$)

583

584 ^a Multivariate logistic regression adjusted for sex and top 10 ancestry-informative principal components

585 ^b Effect allele coded as allele previously associated with increased platelet count from Astle *et. al.* (35)

586 ^c Effect allele frequency in European-ancestry individuals from 1000 Genomes Project

587 ^d Neighboring gene located on UCSC Genome Browser

588 ^e Odds of ALL associated with each additional copy of the effect allele

589 Sample size in discovery cohort (959 Children's Oncology Group cases, 2624 controls); rsIDs from GRCh37/hg19 build

590 **Table 4. Independent replication of ALL risk loci in combined meta-analysis of UK GWAS II and German GWAS**
 591

Locus	SNP rsID	Effect allele	Gene	OR (95% CI) ^c	P	P _{heterogeneity}
2q32.3	rs7585866 ^a	G	<i>SDPR</i>	-	-	-
4q24	rs4699154	C	near <i>TET2</i> ^b	0.98 (0.89-1.06)	0.621	0.300
5q31.1	rs10058074	G	near <i>IRF1</i> ^b	1.15 (1.05-1.26)	8.46 x 10⁻⁴	0.625
6p21.31	rs210142	C	<i>BAK1</i>	1.19 (1.10-1.28)	1.20 x 10⁻⁴	0.780
6p21.31	rs75080135	C	<i>GGNBP1</i>	1.15 (1.05-1.25)	6.80 x 10 ⁻³	0.038
6q23.3	rs1331308	C	<i>HBS1L</i>	0.96 (0.88-1.04)	0.264	0.008
6q23.3	rs7776054	G	<i>HBS1L</i>	1.00 (0.91-1.09)	0.941	0.165
7q32.2	rs11556924	T	<i>ZC3HC1</i>	0.99 (0.90-1.07)	0.749	0.381
12q24.21	rs35427	T	intergenic	1.03 (0.95-1.13)	0.508	0.902
18q12.3	rs16977972 ^a	T	<i>SETBP1</i>	-	-	-
21q22.2	rs2836441	G	<i>ERG</i>	0.85 (0.77-0.94)	5.13 x 10⁻³	0.733
22q11.21	rs1059196 ^a	C	<i>SEPT5, GP1BB</i>	-	-	-

592
 593 Abbreviations: SNP, single nucleotide polymorphism; OR, odds ratio; 95% CI, 95% confidence interval
 594 Bold values indicate Bonferroni-corrected significance ($P < 0.05/9$) with concordant direction of effect in replication analyses
 595

596 ^a Data missing since SNPs did not pass quality control filtering in the replication cohort

597 ^b Neighboring gene located on UCSC Genome Browser

598 ^c Odds of ALL associated with each additional copy of the effect allele, estimates were determined using a fixed-effects model using beta
 599 values and standard errors

600 Sample size in replication cohort; combined UK GWAS II and German GWAS (1618 cases, 9409 controls) (21); rsIDs from GRCh37/hg19
 601 build

602 **FIGURE LEGENDS:**

603 **Figure 1. Methodology for hybrid analysis of GWAS and PheWAS data.** This figure illustrates our
604 approach for investigating phenotype associations with known disease risk variants in order to identify novel
605 candidate risk loci and/or intermediate phenotypes for subsequent analysis in case-control cohorts.
606 Specifically, this figure depicts our application of this approach to acute lymphoblastic leukemia (ALL),
607 which identified platelet count as a phenotype enriched for association with ALL GWAS hits and
608 downstream analysis of the role of platelet count-associated variants in relation to ALL risk in a case-control
609 cohort. Created with Biorender. NHGRI-EBI GWAS catalog diagram attributable to Buniello *et. al.* (56)

610 ^a GWAS catalog (NHGRI-EBI) - <https://www.ebi.ac.uk/gwas/> (23, 56)

611 ^b PheWAS catalog (UK Biobank GeneAtlas) - <http://geneatlas.roslin.ed.ac.uk/phewas/> (29)

612 ^c SNPsnap controls (Broad Institute) - [https://data.broadinstitute.org/mpg/snpsnap/\(25\)](https://data.broadinstitute.org/mpg/snpsnap/(25))

613 ^d PLINK (genome association analysis toolkit) - <https://www.cog-genomics.org/plink2> (36)

614 Abbreviations: GWAS, genome-wide association study; PheWAS, phenome-wide association study; ALL,
615 acute lymphoblastic leukemia; SNPs, single nucleotide polymorphisms.

616

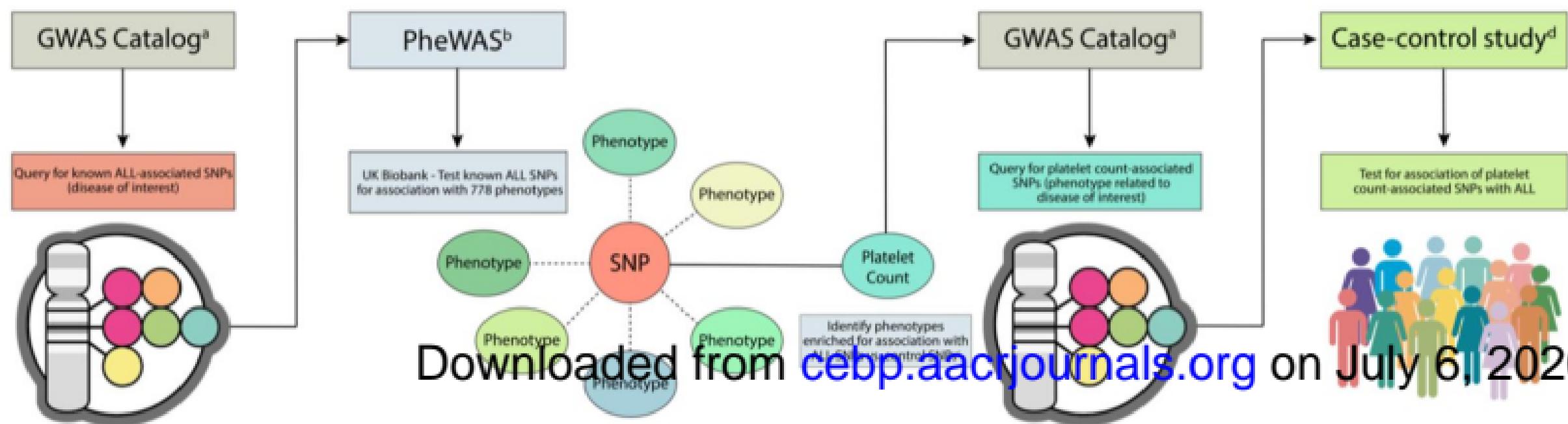
617 **Figure 2. UK Biobank PheWAS traits in ALL SNP-set versus control SNP-sets.** This figure shows the
618 percentage of ALL-associated (12 SNPs total), CLL-associated (31 SNPs total), and matched control SNPs
619 (48 SNPs total) that were significantly associated ($P < 0.01$) with PheWAS traits in the UK Biobank Traits
620 are depicted in descending order of percentage/proportion associated with ALL SNPs from left to right then
621 top to bottom. (A) shows the first subset of 38 traits and (B) shows the second subset of 38 traits, since 76
622 traits total were significantly associated with >1 SNP in the ALL SNP-set, and thus were carried forward for
623 statistical analysis and enrichment comparisons across SNP-sets (see **S5 Table and S6 Table** for full
624 results of proportions and of trait enrichment comparisons between SNP sets).

625 Abbreviations: PheWAS, phenome-wide association study; ALL, acute lymphoblastic leukemia; CLL,
626 chronic lymphocytic leukemia; SNP, single nucleotide polymorphism

627

628

Figure 1



Cancer Epidemiology, Biomarkers & Prevention

AACR American Association
for Cancer Research

Leveraging genome and phenome-wide association studies to investigate genetic risk of acute lymphoblastic leukemia

Eleanor C Semmes, Jayaram Vijaykrishnan, Chenan Zhang, et al.

Cancer Epidemiol Biomarkers Prev Published OnlineFirst May 28, 2020.

Updated version	Access the most recent version of this article at: doi: 10.1158/1055-9965.EPI-20-0113
Supplementary Material	Access the most recent supplemental material at: http://cebp.aacrjournals.org/content/suppl/2020/05/28/1055-9965.EPI-20-0113.DC1
Author Manuscript	Author manuscripts have been peer reviewed and accepted for publication but have not yet been edited.

E-mail alerts	Sign up to receive free email-alerts related to this article or journal.
Reprints and Subscriptions	To order reprints of this article or to subscribe to the journal, contact the AACR Publications Department at pubs@aacr.org .
Permissions	To request permission to re-use all or part of this article, use this link http://cebp.aacrjournals.org/content/early/2020/05/28/1055-9965.EPI-20-0113 . Click on "Request Permissions" which will take you to the Copyright Clearance Center's (CCC) Rightslink site.