

1 Exploring Various Polygenic Risk Scores for Skin Cancer in the Phenomes
2 of the Michigan Genomics Initiative and the UK Biobank with a Visual
3 Catalog: *PRSWeb*

4 Lars G. Fritsche^{1,2,†*}, Lauren J. Beesley^{1,†}, Peter VandeHaar^{1,2}, Robert B. Peng¹,
5 Maxwell Salvatore¹, Matthew Zawistowski^{1,2}, Sarah A. Gagliano^{1,2}, Sayantan Das^{1,2},
6 Jonathon LeFaive^{1,2}, Erin O. Kaleba³, Thomas T. Klumpner^{3,4}, Stephanie E. Moser³,
7 Victoria M. Blanc⁵, Chad M. Brummett^{3,4}, Sachin Kheterpal^{3,4}, Gonçalo R. Abecasis^{1,2},
8 Stephen B. Gruber⁶, Bhramar Mukherjee^{1,2,7,8,9*}

9
10 ¹ Department of Biostatistics, University of Michigan School of Public Health, Ann
11 Arbor, Michigan, United States of America

12 ² Center for Statistical Genetics, University of Michigan School of Public Health, Ann
13 Arbor, Michigan, United States of America

14 ³ Division of Pain Medicine, Department of Anesthesiology, University of Michigan
15 Medical School, Ann Arbor, Michigan, United States of America

16 ⁴ Institute for Healthcare Policy and Innovation, University of Michigan, Ann Arbor,
17 Michigan, United States of America

18 ⁵ Central Biorepository, University of Michigan Medical School, Ann Arbor, Michigan,
19 United States of America

20 ⁶ USC Norris Comprehensive Cancer Center, University of Southern California, Los
21 Angeles, California, United States of America

22 ⁷ Michigan Institute for Data Science, University of Michigan, Ann Arbor, Michigan,
23 United States of America

24 ⁸ Department of Epidemiology, University of Michigan School of Public Health, Ann
25 Arbor, Michigan, United States of America

26 ⁹ University of Michigan Rogel Cancer Center, University of Michigan, Ann Arbor,
27 Michigan, United States of America

28

29 * Corresponding authors

30 E-mail: bhramar@umich.edu (BM), larsf@umich.edu (LGF)

31 [†] These authors contributed equally to this work

32 **Abstract**

33 Polygenic risk scores (PRS) are designed to serve as a single summary measure,
34 condensing information from a large number of genetic variants associated with a
35 disease. They have been used for stratification and prediction of disease risk. The
36 construction of a PRS often depends on the purpose of the study, the available
37 data/summary estimates, and the underlying genetic architecture of a disease. In this
38 paper, we consider several choices for constructing a PRS using summary data
39 obtained from various publicly-available sources including the UK Biobank and evaluate
40 their abilities to predict outcomes derived from electronic health records (EHR). We
41 examine the three most common skin cancer subtypes in the USA: basal cell
42 carcinoma, cutaneous squamous cell carcinoma, and melanoma. The genetic risk
43 profiles of subtypes may consist of both shared and unique elements and we construct
44 PRS to understand the common versus distinct etiology. This study is conducted using
45 data from 30,702 unrelated, genotyped patients of recent European descent from the
46 Michigan Genomics Initiative (MGI), a longitudinal biorepository effort within Michigan
47 Medicine. Using these PRS for various skin cancer subtypes, we conduct a phenome-
48 wide association study (PheWAS) within the MGI data to evaluate their association with
49 secondary traits. PheWAS results are then replicated using population-based UK
50 Biobank data. We develop an accompanying visual catalog called *PRSweb* that
51 provides detailed PheWAS results and allows users to directly compare different PRS
52 construction methods. The results of this study can provide guidance regarding PRS
53 construction in future PRS-PheWAS studies using EHR data involving disease
54 subtypes.

55 **Author summary**

56 In the study of genetically complex diseases, polygenic risk scores synthesize
57 information from multiple genetic risk factors to provide insight into a patient's risk of
58 developing a disease based on his/her genetic profile. These risk scores can be
59 explored in conjunction with health and disease information available in the electronic
60 medical records. They may be associated with diseases that may be related to or
61 precursors of the underlying disease of interest. Limited work is available guiding risk
62 score construction when the goal is to identify associations across the medical
63 phenome. In this paper, we compare different polygenic risk score construction methods
64 in terms of their relationships with the medical phenome. We further propose methods
65 for using these risk scores to decouple the shared and unique genetic profiles of related
66 diseases and to explore related diseases' shared and unique secondary associations.
67 Leveraging and harnessing the rich data resources of the Michigan Genomics Initiative,
68 a biorepository effort at Michigan Medicine, and the larger population-based UK
69 Biobank study, we investigated the performance of genetic risk profiling methods for the
70 three most common types of skin cancer: melanoma, basal cell carcinoma and
71 squamous cell carcinoma.

72 Introduction

73 The underlying risk factors of genetically complex diseases are numerous.
74 Genome-wide association studies (GWAS) on thousands of diseases and traits have
75 made great strides to uncover a vast array of genetic variants that contribute to genetic
76 predispositions to a disease [1]. In order to harness the information from a large number
77 of genetic variants, a popular approach is to summarize their contribution through
78 polygenic risk scores (PRS). While the performance of PRS to predict disease
79 outcomes at a population level has been modest for many diseases including most
80 cancers, PRS have successfully been applied for risk stratification of cohorts [2, 3] and
81 recently have been used to screen a multitude of clinical phenotypes (collectively called
82 the medical phenome) for secondary trait associations [4, 5]. The goal of these
83 phenome-wide screenings is to uncover phenotypes that share genetic components
84 with the primary trait that, if pre-symptomatic, could shed biological insights into the
85 disease pathway and inform early interventions or screening efforts for individuals at
86 risk. However, limited prior work is available guiding the choice of PRS construction for
87 testing associations across the medical phenome.

88 In the post-GWAS era and with the availability of large biobank data from multiple
89 sources, general guidance for constructing a PRS for a phenotype of interest is needed.
90 A PRS of the general form $\sum_{i=1}^K \hat{\beta}_i G_i$ requires specification of three things: a list of
91 markers G_1, G_2, \dots, G_K , the depth of the list or the number of markers (K), and the choice
92 of the weights $\hat{\beta}_i$. These choices can be based on information extracted from the latest
93 GWAS or GWAS meta-analysis (when available), the NHGRI-EBI GWAS catalog of
94 published results [1] (when available), or summary data for GWAS corresponding to

95 each phenotype, e.g., from efforts that comprehensively screened the UK Biobank
96 (UKB) phenome [6, 7]. While various methods of constructing PRS have been widely
97 studied for predicting the primary phenotype collected through population-based
98 sampling [8, 9], it is unknown how the different PRS will be associated with a multitude
99 of other diagnoses across the medical phenome. This study attempts to bridge this
100 knowledge gap.

101 In this paper, we first explore strategies for constructing a PRS using markers
102 and weights obtained from either the latest GWAS or the NHGRI-EBI GWAS catalog
103 that have reached genome-wide significance. We compare the PRS in terms of their
104 performance [10] for the three most common skin cancer subtypes in the USA: basal
105 cell carcinoma (MIM: 614740) [11], cutaneous squamous cell carcinoma [12] and
106 melanoma (MIM: 155601) [13]. We compare the two strategies using an independent
107 biobank of genetic, demographic, and phenotype data collected by the Michigan
108 Genomics Initiative (MGI), a longitudinal biorepository effort within Michigan Medicine
109 (University of Michigan) [4, 14]. Based on these results, we choose a PRS construction
110 strategy for each skin cancer subtype for further analysis.

111 For the chosen PRS corresponding to each skin cancer subtype, we perform a
112 phenome-wide association study (PheWAS) relating the PRS to the electronic health
113 record (EHR)-based phenome of MGI. We call such a study a PRS-PheWAS.⁴ PRS-
114 PheWAS results are then replicated using the population-based UK Biobank data. In
115 order to identify secondary associations that are not driven by the primary phenotype,
116 we perform an additional “exclusion” PRS-PheWAS for each skin cancer subtype in
117 which we exclude subjects with any type of observed skin cancer.⁴ These studies

118 demonstrate differences in PheWAS results for PRS constructed for particular disease
119 subtypes and the ability of such studies to reproduce known associations between
120 secondary phenotypes and particular disease subtypes.

121 We then describe an approach for using PRS to (1) understand the shared and
122 unique genetic architecture of disease subtypes and to (2) identify shared and unique
123 secondary phenotype associations related to this genetic architecture. We define a new
124 PRS for each skin cancer subtype using loci **unique** to that subtype's chosen PRS. We
125 further construct a composite PRS for general skin cancer consisting of loci **common**
126 among all subtypes' PRS. While merging distinct clinical entities into a compound PRS
127 may seem counterintuitive in terms of specificity, such an approach may increase power
128 to identify dermatological features through PheWAS that are shared by all three
129 subtypes, which may in turn provide guidance for general skin cancer screening efforts
130 and sun protection behavior.

131 The NHGRI-EBI GWAS Catalog and Latest GWAS PRS construction methods
132 are based on published GWAS studies, which only report risk variants that reached
133 genome-wide significance (usually defined by a P-value threshold of $P < 5 \times 10^{-8}$).
134 However, it is likely that there are additional risk variants below this threshold that could
135 be associated with the trait but have not reached statistical significance [15].
136 Incorporating non-significant variants may conceivably improve the predictive power of
137 a PRS but may also add additional random false positive signals, which in turn could
138 dilute the discriminatory power of the true risk variants and diminish any predictive gain
139 [8, 16]. To explore whether a PRS constructed using additional non-significant loci may
140 outperform a PRS using only loci reaching genome-wide significance, we evaluated a

141 PRS constructed using publicly available genome-wide summary statistics from the UK
142 Biobank at six different p-value thresholds both in terms of associations with skin cancer
143 phenotypes and in terms of secondary phenotype associations. There is an extensive
144 literature on constructing genome-wide PRS using random effects, shrinkage methods,
145 or thresholding (our focus) [17-19], but none of these methods have been evaluated in a
146 PheWAS setting.

147 In this paper, we focus our attention on skin cancer, but the approaches used in
148 this paper can be applied to study many other phenotypes. We chose to use skin
149 cancer as a demonstrative example for a variety of reasons. First, our discovery dataset
150 (MGI) is particularly enriched for skin cancer cases due to the strong skin cancer clinical
151 program at Michigan Medicine and due to the high rate of surgery for skin cancer
152 patients. MGI primarily recruits participants undergoing surgery and is therefore
153 enriched for cancers and other medical comorbidities when compared to a general
154 population [4]. Additionally, skin cancer has well-defined subtypes, which allows us to
155 explore subtype-specific PRS constructed for several related but distinct diseases in
156 terms of their performance for related skin cancer outcomes. Skin cancer also provides
157 a setting in which there may be genetic factors uniquely related to particular subtypes
158 as well as genetic factors that are shared risk factors for all skin cancer subtypes. The
159 various PRS construction methods explored in this paper delivered tools to explore
160 shared and subtype-specific phenotypes and may provide an enhanced understanding
161 of the genome x phenome landscape.

162 We develop an online visual web catalog called *PRSweb* that provides PRS-
163 PheWAS results for melanoma, basal cell carcinoma, and squamous cell carcinoma.

164 PheWAS results are available using three different PRS construction methods explored
165 in this paper: Latest GWAS, NHGRI-EBI GWAS Catalog, and the UK Biobank GWAS
166 summary statistics using different significance thresholds. The weights and the marker
167 list for each PRS method can be downloaded. Furthermore, PheWAS summary
168 statistics can be accessed from *PRSweb* (see **Web resources**), providing future
169 investigators with readily available and useful tools to perform further analyses.

170 Comprehensive phenome-wide and genome-wide analyses of large biobank
171 studies with publicly available summary statistics can be rich resources for PRS
172 construction, especially if the trait-of-interest's prevalence is high in the biobank. Using
173 PRS, we can synthesize complex genetic information that is then used to identify these
174 shared genetic components across phenotypes. Compared to prior and existing
175 literature, our contribution is new in four principal directions: (1) comparing various PRS
176 construction methods in terms of their relationships with related EHR-derived
177 phenotypes (2) comparing PRS associations with secondary phenotypes across the
178 phenome of MGI (academic medical center) and UK Biobank (population-based), (3)
179 developing PRS-based methods for understanding the shared and unique genetic
180 contribution across disease sub-types both in terms of disease biology and in terms of
181 secondary phenotype associations, and (4) introducing a publicly accessible online
182 visual catalog to visually represent the genome x phenome landscape and access
183 summary data from GWAS and PheWAS.

184

185 **Material and methods**

186

187 **Discovery and replication cohorts**

188 **MGI cohort (discovery cohort).** Participants were recruited through the
189 Michigan Medicine health system while awaiting diagnostic or interventional procedures
190 either during a preoperative visit prior to the procedure or on the day of the procedure
191 that required anaesthesia. Opt-in written informed consent was obtained. In addition to
192 coded biosamples and secure protected health information, participants understood that
193 all EHR, claims, and national data sources – linkable to the participant – may be
194 incorporated into the MGI databank. Each participant donated a blood sample for
195 genetic analysis, underwent baseline vital signs and a comprehensive history and
196 physical assessment. Data were collected according to Declaration of Helsinki
197 principles. Study participants' consent forms and protocols were reviewed and approved
198 by local ethics committees (IRB ID HUM00099605). In the current study, we report
199 results obtained from 30,702 unrelated, genotyped samples of recent European
200 ancestry with available integrated EHR data (~90 % of all MGI participants were inferred
201 to be of recent European ancestry) [4].

202 **UK Biobank cohort (replication cohort).** The UK Biobank is a population-
203 based cohort collected from multiple sites across the United Kingdom and includes over
204 500,000 participants aged between 40 and 69 years when recruited in 2006–2010 [20].
205 The open access UK Biobank data used in this study included genotypes, ICD9 and
206 ICD10 codes, inferred sex, inferred white British-European ancestry, kinship estimates

207 down to third degree, birthyear, genotype array, and precomputed principal components
208 of the genotypes.

209

210 **Genotyping, sample quality control and imputation**

211 **MGI.** DNA from 37,412 blood samples was genotyped on customized Illumina
212 Infinium CoreExome-24 bead arrays and subjected to various quality control filters that
213 resulted in a set of 392,323 polymorphic variants. Principal components and ancestry
214 were estimated by projecting all genotyped samples into the space of the principal
215 components of the Human Genome Diversity Project reference panel using PLINK (938
216 unrelated individuals) [21, 22]. Pairwise kinship was assessed with the software KING
217 [23], and the software fastindep was used to reduce the data to a maximal subset that
218 contained no pairs of individuals with 3rd-or closer degree relationship [24]. We also
219 removed patients not of recent European descent from the analysis, resulting in a final
220 sample of 30,702 unrelated subjects. Additional genotypes were obtained using the
221 Haplotype Reference Consortium using the Michigan Imputation Server [25] and
222 included over 17 million imputed variants with $R^2 < 0.3$ and/or minor allele frequency
223 (MAF) $< 0.1\%$. Genotyping, quality control and imputation are described in detail
224 elsewhere [4]. **Table 1** provides some descriptive statistics of the MGI and UK Biobank
225 samples.

226 **Table 1. Demographics and Clinical Characteristics of the Analytic Datasets**

Characteristic	MGI	UK Biobank*
n	30,702	408,961
Females, n (%)	16,297 (53.1%)	221,052 (54.1%)
Mean Age, years (S.D.)	54.2 (15.9)	57.7 (8.1)
Median number of visits per participant	27	n/a
Median days between first and last visit	1,469	n/a
Total number of ICD9 code days	3,459,331	49,085
Number of unique ICD9 codes	10,323	3,126
Median ICD9 code days per participant	58	2
Total number of ICD10 code days	1,311,264	2,764,868
Number of unique ICD10 codes	14,997	11,059
Median ICD10 code days per participant	27	6
Total number of PheWAS code days	6,367,117	3,679,624
Number of unique PheWAS codes	1,856	1,680
Median PheWAS code days per participant	94	8
n cases with Skin Cancer	4,503	13,782 (13,624***)
n cases with melanomas of skin	1,772	2,724 (2,718***)
n cases with epithelial skin cancer and others**	3,220	11,152 (11,030***)
n cases with basal cell carcinoma	1,303	Not available
n cases with squamous cell carcinoma	836	Not available

227
 228 * The provided characteristics are based a subset of white British subjects of the UK
 229 Biobank Study for which phenotype data and imputed data was available. To retain as
 230 many unrelated cases as possible for each trait, a maximal set of unrelated cases was
 231 identified before choosing controls from the pool of subjects unrelated to these cases or
 232 to each other.

233 ** Original PheWAS code “172.2” description "Other non-epithelial cancer of skin".

234 *** Unrelated cases

235 ICD9 and ICD10: International Statistical Classification of Diseases codes (9th and 10th
 236 revision)

237 **UK Biobank.** The UK Biobank is a population-based cohort collected from
238 multiple sites across the United Kingdom [20]. After quality control, we phased and
239 imputed the 487,409 UK Biobank genotyped samples against the Trans-Omics for
240 Precision Medicine (TOPMed) reference panel (see **Web resources**), which is
241 composed of 60,039 multi-ethnic samples and 239,756,147 SNP and indel variants
242 sequenced at high depth (30x). The phasing step was carried out on 81 chromosomal
243 chunks with around 10,000 genotyped variants in each chunk using the software Eagle
244 (with the “kbpwt” parameter set at 80,000) [26]. The imputation was carried out in 137
245 chromosomal chunks of around 20 Mbp in length with Mbp of total overlap on either
246 side using the imputation tool Minimac4 (see **Web resources**). To increase
247 computational efficiency, we imputed each of the chunks in batches of 10,000 samples
248 at a time and then merged them back using BCFtools. Since Minimac4 imputes each
249 sample independently, analyzing our samples in batches did not change their
250 imputation estimates. However, this sampling would result in different imputation quality
251 estimates for each batch, and thus we collapsed the estimates to generate imputation
252 quality estimates across all the study samples. After imputation, we filtered out variants
253 with estimated imputation accuracy of $R^2 < 0.1$, which left us with 177,895,992 variants.

254

255 **Phenome generation**

256 **MGI.** The MGI phenome was used as the discovery dataset and was based on
257 the Ninth and Tenth Revision of the International Statistical Classification of Diseases
258 (ICD9 and ICD10) code data for 30,702 unrelated, genotyped individuals of recent
259 European ancestry. These ICD9 and ICD10 codes were aggregated to form up to 1,857

260 PheWAS traits using the PheWAS R package (as described in detail elsewhere[4, 27]).
261 For each trait, we identified case and control samples. To minimize differences in age
262 and sex distributions or extreme case-control ratios as well as to reduce computational
263 burden, we matched up to 10 controls to each case using the R package “MatchIt” [28].
264 Nearest neighbor matching was applied for age and PC1-4 (using Mahalanobis-metric
265 matching; matching window caliper/width of 0.25 standard deviations) and exact
266 matching was applied for sex and genotyping array. A total of 1,578 case control studies
267 with >50 cases were used for our analyses of the MGI phenome.

268 **UK Biobank.** The UK Biobank phenome was used as a replication dataset and
269 was based on ICD9 and ICD10 code data of 408,961 white British [14], genotyped
270 individuals that were aggregated to PheWAS traits in a similar fashion (as described
271 elsewhere [7]). To remove related individuals and to retain larger sample sizes, we first
272 selected a maximal set of unrelated cases for each phenotype (defined as no pairwise
273 relationship of 3rd degree or closer [24, 29]) before selecting a maximal set of unrelated
274 controls unrelated to these cases. Similar to MGI, we matched up to 10 controls to each
275 case using the R package “MatchIt” [28]. Nearest neighbor matching was applied for
276 birthyear and PC1-4 (using Mahalanobis-metric matching; matching window
277 caliper/width of 0.25 standard deviations) and exact matching was applied for sex and
278 genotyping array. 1,366 case control studies with >50 cases each were used for our
279 analyses of the UK Biobank phenome.

280 Additional phenotype information for MGI and UK Biobank is included in **S1 Text**
281 **Fig B** and **S2 Text Tables F-H**.

282

283 **Risk SNP selection**

284 For each skin cancer subtype (melanoma, basal cell carcinoma, and squamous
285 cell carcinoma), we generated three different sets of PRS: (1) based on merged
286 summary statistics published in the NHGRI EBI GWAS catalog [1], (2) based on the
287 latest available GWAS meta-analysis [30-32] and (3) based on publicly available GWAS
288 summary statistics from the UK Biobank data [7].

289 **GWAS Catalog SNP selection.** We downloaded previously reported GWAS
290 variants from the NHGRI-EBI GWAS Catalog (file date: February 28, 2018) [1, 33].
291 None of the currently available skin cancer discovery studies included in the catalog
292 used any subset of the MGI cohort or data from the UK Biobank. Single nucleotide
293 polymorphism (SNP) positions were converted to GRCh37 using variant IDs from
294 dbSNP: build 150 (UCSC Genome Browser) after updating outdated dbSNP IDs to their
295 merged dbSNP IDs. Entries with missing risk alleles, risk allele frequencies, or odds
296 ratios were excluded. If a reported risk allele did not match any of the reported forward
297 strand alleles of a non-ambiguous SNP (not A/T or C/G) in the imputed genotype data
298 (which correspond to the alleles of the imputation reference panel), we assumed minus
299 strand designation and corrected the effect allele to its complementary base of the
300 forward strand. Entries with a reported risk allele that did not match any of the alleles of
301 an ambiguous SNP (A/T and C/G) in our data were excluded at this step. We only
302 included entries with broad European ancestry (as reported by the NHGRI-EBI GWAS
303 Catalog). As a quality control check, we compared the reported risk allele frequencies
304 (RAF) in controls with the RAF of 14,770 MGI individuals who had no cancer diagnosis
305 (for chromosome X variants, we calculated RAF in females only). We then excluded

306 entries whose RAF deviated more than 15%. This chosen threshold is subjective and
307 was based on clear differentiation between correct and likely flipped alleles on the two
308 diagonals (see **S1 Text Fig A**) as noted frequently in GWAS meta-analyses quality
309 control procedures [34]. For each analyzed cancer type, we extracted risk variants that
310 were also present in our genotype data and estimated pairwise linkage disequilibrium
311 (LD; correlation r^2) using the allele dosages of the corresponding controls. For pairwise
312 correlated SNPs ($r^2 > 0.1$) or SNPs with multiple entries, we kept the SNP with the most
313 recent publication date (and smaller P value, if necessary) and excluded the other (**S2**
314 **File Table I**).

315 **Selection of risk SNPs from largest GWAS.** In a similar fashion, we extracted
316 and filtered reported association signals from large GWAS meta-analyses on basal cell
317 carcinoma [31], cutaneous squamous cell carcinoma [30] and melanoma [32] (**S2 File**
318 **Table I**).

319 **Genome-wide SNP selection of UK-Biobank-based GWAS.** We obtained
320 GWAS summary statistics for the ICD9- and ICD10-based PheWAS codes “172” (skin
321 cancer; 13,752 cases versus 395,071 controls), “172.11” (melanoma; 2,691 cases
322 versus 395,071 controls), and “172.2” (non-epithelial skin cancer; 11,149 cases versus
323 395,071 controls) from a public download [7] (see **Web resources**). These GWAS
324 analyzed up to 408,961 white British European-ancestry samples with generalized
325 mixed model association tests that used the saddlepoint approximation to calibrate the
326 distribution of score test statistics and thus could control for unbalanced case-control
327 ratios and sample relatedness [7]. For each trait, we reduced these summary statistics
328 to SNPs that were reported with minor allele frequencies $> 0.5\%$ and were also

329 available for the MGI data. Next, we performed linkage LD clumping of all variants with
330 p-values $< 5 \times 10^{-4}$ using the imputed allele dosages to obtain independent risk SNPs (LD
331 threshold of $r^2 > 0.1$ and a maximal SNP distance of 1 Mb). We limited the LD
332 calculations to 10,000 randomly selected, unrelated, white British individuals to reduce
333 the computational burden. Finally, we created subsets of these independent SNPs with
334 p-values $< 5 \times 10^{-9}$, $< 5 \times 10^{-8}$, $< 5 \times 10^{-7}$, $< 5 \times 10^{-6}$, $< 5 \times 10^{-5}$, and $< 5 \times 10^{-4}$ (**S2 File Table J**).

335

336 **Construction of the polygenic risk scores**

337 For each of the obtained SNP sets for each trait, we constructed a PRS as the
338 sum of the allele dosages of risk increasing alleles of the SNPs weighted by their
339 reported log odds ratios. Restated, the PRS for subject j in MGI was of the form
340 $PRS_j = \sum_i \beta_i G_{ij}$ where i indexes the included loci for that trait, β_i is the log odds ratios
341 retrieved from the external GWAS summary statistics for locus i , and G_{ij} is a continuous
342 version of the measured dosage data for the risk allele on locus i in subject j . The PRS
343 variable was created for each MGI and UKB participant. For comparability of effect
344 sizes corresponding to the continuous PRS across cancer traits and PRS construction
345 methods, we transformed each PRS to the standard Normal distribution using
346 “ztransform” of the R package “GenABEL” [35].

347

348 **Statistical analysis**

349 In this study, we constructed PRS for three skin cancer subtypes using two
350 different PRS construction methods (using the Latest GWAS or the corresponding
351 entries of the GWAS Catalog). To compare the association between PRS and skin

352 cancer phenotypes across different PRS construction methods, we fit the following
353 model for each PRS and skin cancer phenotype:

354 $\text{logit} (P(\text{Phenotype is present} \mid \text{PRS, Age, Sex, Array, PC})) = \beta_0 + \beta_{PRS} \text{PRS} +$
355 $\beta_{Age} \text{Age} + \beta_{Sex} \text{Sex} + \beta_{Array} \text{Array} + \beta \text{PC}$, where the PCs were the first four principal
356 components obtained from the principal component analysis of the genotyped GWAS
357 markers and where “Array” represents the genotyping array. Our primary interest is in
358 β_{PRS} , while the other factors (Age, Sex and PC) were included to address potential
359 residual confounding and do not provide interpretable estimates due to the preceding
360 application of case control matching. Firth's bias reduction method was used to resolve
361 the problem of separation in logistic regression (Logistf in R package “EHR”) [36-38], a
362 common problem for binary or categorical outcome models when for a certain part of
363 the covariate space there is only one observed value of the outcome, which often leads
364 to very large parameter estimates and standard errors.

365 We then evaluated each PRS's (1) ability to discriminate between cases and
366 controls by determining the area under the receiver-operator characteristics (ROC)
367 curve (AUC) using R package “pROC” [39]; (2) calibration using Hosmer-Lemeshow
368 Goodness Of Fit test of the R package “ResourceSelection” [40, 41]; and (3) accuracy
369 with the Brier Score of R package “DescTools” [42]. These evaluations did not adjust for
370 additional covariates. We used these metrics and the logistic regression results to
371 choose a PRS construction method to use for each skin cancer subtype moving
372 forward. To explore the impact of incorporating non-significant loci into the PRS
373 construction, we further performed the above analyses with PRS constructed using UK
374 Biobank GWAS summary statistics with different p-value thresholds.

375 Using the chosen PRS for each subtype, we conducted two PheWAS to identify
376 other phenotypes associated with the PRS first for the 1,578 phenotypes in MGI and
377 then for the 1,366 phenotypes from UK Biobank. To evaluate PRS-phenotype
378 associations, we conducted Firth bias-corrected logistic regression by fitting a model of
379 the above form for each phenotype and data source. Age represents the birth year in
380 UK Biobank. To adjust for multiple testing, we applied the conservative phenome-wide
381 Bonferroni correction according to the analyzed PheWAS codes ($n_{\text{MGI}} = 1,578$ or n_{UK}
382 $n_{\text{Biobank}} = 1,366$). In Manhattan plots, we present $-\log_{10}(\textit{p}\text{-value})$ corresponding to tests
383 of $H_0: \beta_{\text{PRS}} = 0$. Directional triangles on the PheWAS plot indicate whether a phenome-
384 wide significant trait was positively (pointing up) or negatively (pointing down)
385 associated with the PRS.

386 To investigate the possibility of the secondary trait associations with PRS being
387 completely driven by the primary trait association, we performed a second set of
388 PheWAS after excluding individuals affected with the primary or related cancer traits for
389 which the PRS was constructed, referred to as “exclusion PRS PheWAS” as described
390 previously [4]. We then constructed new PRS scores representing shared and subsite-
391 unique genetic components and performed a PheWAS for each.

392 To evaluate how well prior presence of an identified secondary non-skin-cancer
393 diagnosis can identify subjects with increased risk of developing skin cancer, we
394 created a binary variable taking the value 1 if a given subject (1) was diagnosed with the
395 non-skin-cancer diagnosis and then diagnosed with skin cancer at least 365 days after
396 or (2) was diagnosed with the non-skin-cancer diagnosis and never diagnosed with skin
397 cancer. We then fit a Firth bias-corrected logistic regression of the following form:

398 $\text{logit}(P(\text{Primary phenotype is present} \mid \text{Predictor, Age, Sex, Array, PC}))$
399 $=\beta_0 + \beta_{PRS}I(\text{Secondary non skin cancer trait}) + \beta_{Age}Age + \beta_{Sex}Sex + \beta_{Array}Array + \beta PC$
400 where Array and PC were defined as before. Unless otherwise stated, analyses were
401 performed using R 3.4.4 [43].

402

403 **Development of an online visual catalog: *PRSweb***

404 The online open access online visual catalog *PRSweb* available at
405 <https://statgen.github.io/PRSweb> was implemented using “Pandas”, a Data Analysis
406 Library, which offers high level of performance for large data structures and data
407 analysis in the Python3 environment [44]. In combination with “Jinja2”, a templating
408 language for Python, and “Bootstrap”, a Cascading Style Sheets (CSS) framework (see
409 Web resources), static HTML files were compiled and allow easy and fast hosting of all
410 PRS-PheWAS results. The interactive plots are drawn with the JavaScript library
411 “LocusZoom.js” (see Web resources) offers dynamic plotting, automatic plot sizing and
412 label positioning.

413 Results

414 Assessing various PRS construction methods

415 We first explored the comparative performance of two PRS construction
 416 strategies in terms of the resulting PRS associations with related phenotypes in the skin
 417 cancer setting. **Table 2** provides the results.

418

419 **Table 2. Associations of constructed PRS with skin cancer traits in MGI**

PRS (Number of SNPs)		Skin Cancer n = 4,503	Melanoma n = 1,896	Basal Cell Carcinoma n = 1,303	Squamous Cell Carcinoma n = 836
<i>PRS based on GWAS Catalog</i>					
Melanoma (29)	PRS OR ^a P-value ^a AUC ^b HL χ^2 , P-value ^c Brier Score	1.41 (1.35,1.47) 2.7x10 ⁻⁵³ 0.57 (0.56,0.58) 10,0.24 0.14	1.68 (1.57,1.79) 1.3x10 ⁻⁵³ 0.61 (0.60,0.62) 5.3,0.72 0.09	1.42 (1.31,1.53) 7.3x10 ⁻¹⁹ 0.57 (0.56,0.59) 12,0.16 0.091	1.3 (1.19,1.44) 4.3x10 ⁻⁰⁸ 0.55 (0.53,0.57) 3.7,0.89 0.09
Basal cell carcinoma (32)	PRS OR ^a P-value ^a AUC ^b HL χ^2 , P-value ^c Brier Score	1.39 (1.33,1.44) 8x10 ⁻⁶⁰ 0.57 (0.56,0.58) 13,0.12 0.14	1.37 (1.29,1.45) 4.8x10 ⁻²⁵ 0.57 (0.56,0.58) 8.6,0.38 0.091	1.82 (1.70,1.95) 3.6x10 ⁻⁶⁵ 0.64 (0.62,0.65) 9.5,0.3 0.09	1.4 (1.28,1.52) 1.4x10 ⁻¹⁴ 0.57 (0.55,0.59) 13,0.11 0.09
Squamous cell carcinoma (18)	PRS OR ^a P-value ^a AUC ^b HL χ^2 , P-value ^c Brier Score	1.28 (1.24,1.33) 4.8x10 ⁻⁴² 0.56 (0.56,0.57) 4.7,0.79 0.14	1.35 (1.28,1.43) 2x10 ⁻²⁸ 0.58 (0.56,0.59) 5.3,0.72 0.091	1.39 (1.31,1.48) 7.9x10 ⁻²⁶ 0.59 (0.57,0.60) 5.1,0.75 0.091	1.29 (1.19,1.39) 1.8x10 ⁻¹⁰ 0.56 (0.54,0.59) 7.8,0.46 0.09
<i>PRS based on Latest GWAS</i>					
Melanoma (20)	PRS OR ^a P-value ^a AUC ^b HL χ^2 , P-value ^c Brier Score	1.48 (1.41,1.55) 3.5x10 ⁻⁵⁵ 0.57 (0.56,0.58) 3,0.93 0.14	1.78 (1.65,1.92) 7x10 ⁻⁵³ 0.61 (0.59,0.62) 6.7,0.56 0.09	1.60 (1.47,1.75) 7.9x10 ⁻²⁷ 0.59 (0.57,0.60) 2.5,0.96 0.091	1.38 (1.24,1.53) 4x10 ⁻⁰⁹ 0.56 (0.54,0.58) 4.3,0.83 0.09
Basal cell carcinoma (28)	PRS OR ^a P-value ^a AUC ^b HL χ^2 , P-value ^c Brier Score	1.42 (1.36,1.48) 5.8x10 ⁻⁶¹ 0.58 (0.57,0.58) 4.3,0.83 0.14	1.43 (1.34,1.52) 7x10 ⁻²⁹ 0.58 (0.56,0.59) 16,0.051 0.091	1.84 (1.71,1.97) 2.8x10 ⁻⁶⁰ 0.63 (0.62,0.65) 4,0.86 0.09	1.45 (1.32,1.58) 1.2x10 ⁻¹⁵ 0.57 (0.55,0.60) 17,0.035 0.09
Squamous cell carcinoma (10)	PRS OR ^a P-value ^a AUC ^b HL χ^2 , P-value ^c Brier Score	1.44 (1.38,1.5) 1.1x10 ⁻⁷⁰ 0.58 (0.57,0.59) 17,0.027 0.14	1.54 (1.45,1.64) 2.9x10 ⁻⁴⁶ 0.60 (0.58,0.61) 13,0.13 0.09	1.62 (1.52,1.73) 1.8x10 ⁻⁴³ 0.61 (0.60,0.63) 6,0.64 0.09	1.52 (1.39,1.65) 2.1x10 ⁻²¹ 0.59 (0.57,0.61) 4.9,0.76 0.09

420

421 ^a Association of each cancer with continuous PRS that were transformed to standard normal
422 distribution. Point estimates, 95% confidence intervals and P- values are obtained by fitting Firth's
423 Bias-Corrected Logistic Regression.

424 ^b Area under the curve of the receiver operating characteristic (ROC) curve with 95% confidence
425 intervals.

426 ^c Hosmer-Lemeshow Goodness-of-Fit Test

427

428 **Comparisons within methods.** Using the GWAS Catalog construction method,
429 the melanoma PRS was more strongly associated with and had better discrimination for
430 the melanoma phenotype than the other skin cancer phenotypes. For the PRS based on
431 the GWAS Catalog, the odds ratio (OR) of the melanoma PRS was 1.68 (95% CI, [1.57,
432 1.79]). By “discrimination,” we refer to the ability of the PRS to distinguish melanoma
433 cases and controls, which is measured by AUC. The melanoma PRS AUC for the
434 melanoma phenotype is 0.61 (95 % CI, [0.60, 0.62]). Similarly, the basal cell carcinoma
435 PRS was most strongly associated with and had the best discrimination for the basal
436 cell carcinoma phenotype, with an OR of 1.82 (95% CI, [1.70, 1.95]) and an AUC of
437 0.64 (95% CI, [0.62, 0.65]). Unlike the other cancer subtypes, the squamous cell
438 carcinoma PRS did not appear to be most strongly associated with the squamous cell
439 carcinoma phenotype. Instead, it was most strongly associated with and most
440 discriminative for basal cell carcinoma. For all three skin cancer subtypes, the PRS
441 produced higher Brier scores for overall skin cancer, suggesting that the subtype-
442 defined PRS were less accurate for predicting skin cancer as a whole. We obtain similar
443 conclusions for the Latest GWAS method.

444 **Comparisons across methods.** For each cancer subtype, we compared the
445 PRS-subtype associations for the two PRS construction methods. **Melanoma:** For the
446 melanoma PRS, the GWAS Catalog method and the Latest GWAS method produced
447 similar performance in terms of AUC, OR, Hosmer-Lemeshow goodness of fit, and Brier
448 score. For example, the AUC for melanoma for the GWAS Catalog melanoma PRS was
449 0.61 (95% CI, [0.60, 0.62]). The corresponding AUC for the Latest GWAS method was
450 0.61 (95% CI, [0.59, 0.62]). **S1 Text Fig J** compares PRS weights to corresponding
451 SNP-melanoma associations in MGI and UK Biobank. **Basal Cell Carcinoma:** As with
452 melanoma, the basal cell carcinoma PRS produced similar results under the GWAS
453 Catalog and Latest GWAS construction methods. The basal cell carcinoma AUC under
454 the GWAS catalog method was 0.64 (95% CI, [0.62, 0.65]) and the AUC under the
455 Latest GWAS method was 0.63 (95% CI, [0.62, 0.65]). The OR values and Brier score
456 values were nearly identical, and neither approach produced evidence of lack of fit
457 based on the Hosmer-Lemeshow statistic. **Squamous Cell Carcinoma:** The squamous
458 cell carcinoma PRS was not more strongly associated with the squamous cell
459 carcinoma phenotype than the other phenotypes. However, we do observe that the
460 squamous cell carcinoma phenotype using the GWAS Catalog method (0.56, 95% CI
461 [0.54, 0.59]) produced a lower AUC compared to the Latest GWAS method (0.59, 95%
462 CI [0.57, 0.61]). While a difference of 0.03 may not seem like a large difference in AUC
463 in other applications, any improvement in AUC for PRS associations with observed
464 phenotypes may be considered appreciable [45]. These two methods produced identical
465 Brier scores, and the Latest GWAS method resulted in a stronger association between

466 the PRS and the squamous cell carcinoma phenotype (OR of 1.29, 95% CI [1.19, 1.39]
467 vs OR of 1.52, 95% CI [1.39, 1.65]).

468 Using the above comparisons between the two PRS construction methods, we
469 chose a single PRS construction method for each skin cancer subtype to use in
470 subsequent analyses. For melanoma and basal cell carcinoma, we chose the GWAS
471 Catalog method. While the GWAS Catalog and Latest GWAS methods were very
472 similar for these two subtypes, we chose to pursue the GWAS Catalog PRS for future
473 analysis due to the larger number of loci for these PRS (29 vs 20 for melanoma and 32
474 vs 28 for basal cell carcinoma). We choose the Latest GWAS method for squamous cell
475 carcinoma due to its improved AUC over the GWAS Catalog method. We will denote
476 the chosen PRS for melanoma, basal cell carcinoma, and squamous cell carcinoma as
477 mPRS, bPRS, and sPRS respectively.

478

479 **PheWAS using the chosen PRS in MGI**

480 Using each of the chosen PRS described above (mPRS, bPRS, and sPRS), we
481 tested the association between each PRS and each of the 1,578 constructed
482 phenotypes in MGI. For each PRS, the strongest associations were observed with
483 dermatologic neoplasms that included overall skin cancer, melanoma, “other non-
484 epithelial cancer of skin” (the PheWAS over-category of basal and squamous cell
485 carcinoma), and carcinoma in situ of skin. In addition, secondary dermatologic traits
486 such as actinic keratosis (with over-category “degenerative skin conditions and other
487 dermatoses”), chronic dermatitis due to solar radiation (with over-category “dermatitis
488 due to solar radiation”), and seborrheic keratosis were found to be associated with all

489 three PRS (**Fig 1** and **S2 File Table K**). mPRS was most strongly associated with the
490 melanoma phenotype (OR 1.67, 95% CI [1.56, 1.79]), while bPRS was most strongly
491 associated with carcinoma in situ of the skin (OR 1.51, 95% CI [1.39, 1.64]) followed
492 closely by “non-epithelial cancer of the skin” (OR 1.47, 95% CI [1.41, 1.54]). sPRS was
493 most strongly associated with carcinoma in situ of the skin (OR 1.79, 95% CI [1.65,
494 1.94]). The OR of all these phenotypes indicated an increased risk for primary and
495 secondary traits with increasing PRS.

496

497 **Validation of PRS-PheWAS in UK Biobank**

498 To substantiate the detected dermatologic associations, we reiterated the
499 association screen of the three PRS in the matched phenome of the population-based
500 UK Biobank data set (**Fig 1**). In general, stronger evidence for association was found in
501 UKB compared to MGI. This may be driven by the larger sample sizes, e.g. a total of
502 13,623 skin cancer cases versus 4,503 in MGI. In the UK Biobank phenome, the large
503 majority of the previous associations with dermatologic neoplasms were validated with
504 the exception of the trait “dermatitis due to solar radiation”, which had substantially
505 fewer cases in UKB compared to MGI (390 versus 2,959 cases). Unlike MGI, all three
506 PRS were significantly associated (at the phenome-wide level) with “cancer, suspected
507 or other” and “malignant neoplasm, other.”

508

509 **Exclusion PheWAS using the chosen PRS in MGI**

510 In order to explore whether the identified PRS-phenotype associations were
511 driven by the primary trait used to define the PRS (for example, as a side effect of

512 treatment given after diagnosis with the primary trait), we performed a PheWAS for
513 each PRS in which we excluded subjects who were cases for the primary trait or other
514 skin cancer subtypes [4]. Results are shown in **S2 File Table K** and **S1 Text Fig C**.
515 Actinic keratosis, a skin condition believed to be a precursor to non-melanoma skin
516 cancers, remained significantly associated with the squamous cell carcinoma PRS in
517 MGI and all three PRS in UK Biobank [46][47, 48]. No other phenotypes were significant
518 for MGI. “Sebaceous cyst” and its over-category “diseases of the sebaceous gland”
519 were significant in the main UK Biobank PheWAS and remained significantly associated
520 with basal cell carcinoma PRS and squamous cell carcinoma PRS in UK Biobank in the
521 Exclusion PheWAS.

522

523 **Sub-analysis of actinic keratosis as a predictor of future skin cancer**

524 Actinic keratosis (AK) is a rough, scaly patch of skin that usually develops after
525 years of cumulative skin exposure [49]. Previous research has identified actinic
526 keratosis as a common pre-malignant condition for squamous cell carcinoma (SCC)
527 [46]. Actinic keratosis has also been identified as a potential precursor to basal cell
528 carcinoma (BCC) [47, 48]. The availability of temporal information of diagnoses in the
529 MGI cohort offered the opportunity to explore actinic keratosis as a potential precursor
530 for development of skin cancer in MGI.

531 **Fig 2** shows the ROC curves and AUC values for diagnosis of actinic keratosis at
532 least one year before any skin cancer diagnosis and its association with future BCC or
533 SCC diagnosis. AK diagnosis alone has little discrimination abilities, with AUC values of
534 0.52 (95% CI [0.51, 0.53]) for BCC and 0.51 (95% CI [0.50, 0.61]) for SCC. The bPRS

535 and sPRS provide comparatively good discrimination SCC (AUC 0.63 [0.62, 0.65] for
536 BCC and 0.59 [0.57, 0.61] for SCC). The combination of prior AK diagnosis and bPRS
537 provided further improvement in discrimination, with an AUC of 0.65 (95% CI [0.64,
538 0.67]).

539 **S1 Text Tables A and B** provides odds ratio estimates relating AK and the PRS
540 to future BCC and SCC diagnosis. In unadjusted models, the odds of BCC diagnosis
541 were significantly higher in subjects with a prior actinic keratosis diagnosis (OR 1.46,
542 95% CI [1.18, 1.80]). Notably, when we adjust for both bPRS and AK diagnosis, the
543 unadjusted and adjusted effects of both variables are similar, suggesting that AK
544 diagnosis may be an independent predictor of future BCC diagnosis. In contrast, AK
545 diagnosis was **not** an independent predictor of SCC diagnosis. **S1 File Fig G** shows the
546 timing of an AK diagnosis relative to a skin cancer diagnosis for patients with both
547 diagnoses. For subjects with basal cell carcinoma or squamous cell carcinoma, AK
548 diagnoses tended to occur prior to the skin cancer diagnosis (often within 8 years).

549

550 **PRS-PheWAS for shared and unique loci across skin cancer subtypes**

551 In the PRS-PheWAS analyses, we note a striking overlap in the secondary
552 dermatological traits significantly associated with each of the three PRS (mPRS, bPRS,
553 sPRS). One potential explanation for this is that subjects may have more screening
554 after an initial skin cancer diagnosis. Indeed, many subjects have multiple skin cancer
555 diagnoses (**S1 Text Fig D**). **Fig 3** shows the number of risk loci shared by different
556 PRS. Six risk loci are shared between the mPRS, bPRS, and sPRS.

557 This observation inspired follow-up exploration in which we defined a PRS for
558 each cancer subtype using the loci unique to that subtype's chosen PRS. We call these
559 new PRS scores mPRS-u, bPRS-u, and sPRS-u, which reflect the unique loci in the
560 PRS for melanoma, basal cell carcinoma, and squamous cell carcinoma respectively.
561 We also define a PRS consisting of all loci shared across the three skin cancer
562 subtypes, which we call the shared PRS.

563 **S1 Text Table C** shows the association between the various constructed PRS
564 and the skin cancer phenotypes. As with mPRS, mPRS-u was most strongly associated
565 with the melanoma phenotype and is not significantly associated with the other skin
566 cancer subtypes. The bPRS-u score was similarly most strongly associated with basal
567 cell carcinoma and not significantly associated with the other subtypes. We note that the
568 melanoma AUC for the mPRS score was 0.61 (95% CI, [0.60, 0.62]) and is only 0.55
569 (95% CI, [0.53, 0.55]) for the mPRS-u score. Similarly, the basal cell carcinoma AUC for
570 the bPRS score was 0.64 (95% CI, [0.62, 0.65]) and is only 0.57 (95% CI, [0.55, 0.58])
571 for the bPRS-u score. The sPRS-u score is not more strongly associated with the
572 squamous cell carcinoma phenotype than the other skin cancer subtypes. For this
573 reason, we do not include this PRS in further analyses. The shared PRS constructed as
574 the unweighted sum of risk alleles of loci present in all three PRS scores (mPRS, bPRS,
575 and sPRS) is more strongly associated with all three subtype phenotypes than the
576 overall skin cancer phenotype, and the overall skin cancer phenotype also has the
577 lowest AUC and highest Brier score.

578 **S1 Text Fig E** shows PRS-PheWAS results using mPRS-u and bPRS-u. The
579 scores again reveal their subtype specificity in both phenomes, while no secondary

580 associations were observed. Although not shown here, additional exploration into the
581 loci identified uniquely for each subtype may provide some insight into subtype-specific
582 biological mechanisms. **S1 Text Fig F** shows PRS-PheWAS results for the shared PRS.
583 Most strikingly, the shared skin cancer PRS was associated with the top skin cancer
584 and dermatologic traits that were previously found to be associated with the three
585 partially overlapping PRS constructs, suggesting that a shared genetic risk may be
586 driving many of these secondary associations. These six underlying loci (*HERC2* [MIM
587 605837] / *OCA2* [MIM 611409], *IRF4* [MIM 601900], *MC1R* [MIM 155555], *RALY* [MIM
588 614663], *SLC45A2* [MIM 606202] and *TYR* [MIM 606933]) were previously found to be
589 associated not only with skin cancer traits, but also with pigmentation traits of skin, eyes
590 and hair (**Fig 3**; MIM 266300) [31, 50-69].

591 One of these pigmentation traits, skin tanning ability, the tendency of skin to
592 sunburn rather than to suntan, is a well-known risk factor for all skin cancer traits [69,
593 70]. A PRS based on the independent risk variants of a recent GWAS meta-analysis on
594 skin tanning ability [70] was strongly associated with overall skin cancer, melanoma,
595 basal cell carcinoma, and squamous cell carcinoma and even outperformed the
596 constructed PRS of the former two traits (**S1 Text Table C**). Furthermore, the skin
597 tanning ability PRS PheWAS identified a very similar set of traits as the shared skin
598 cancer PRS but in general revealed stronger associations (**S1 Text Fig F**).

599

600 **PRS construction based on UK Biobank summary statistics**

601 The NHGRI-EBI GWAS Catalog and Latest GWAS PRS construction methods
602 are based on published GWAS studies, which often only report risk variants that

603 reached genome-wide significance, but we may believe that incorporating additional risk
604 variance below this threshold may improve predictive power of a PRS. To explore
605 whether a PRS incorporating non-significant loci will outperform a PRS incorporating
606 only significant loci, we constructed PRS using loci related to the phenotype at six
607 different p-value thresholds based on publicly available GWAS summary statistics from
608 the UK Biobank. Larger p-values indicate greater SNP depth (with more SNPs being
609 incorporated into the PRS).

610 The collection of UK Biobank GWAS results did not include basal cell carcinoma
611 or squamous cell carcinoma subtypes; rather, it included only the merged trait ‘non-
612 epithelial cancer of skin’ (**S1 Text Fig B**). Thus, we limited our assessment of the
613 summary statistics to the overall skin cancer GWAS (UKB PheWAS code “172”: 13,752
614 skin cancer cases versus 395,071 controls) and the melanoma GWAS (UKB PheWAS
615 code “172.11”: 2,691 melanoma cases versus 395,071 controls) (**S2 File Table J**).

616 **S1 Text Table D** provides the results. As with the other PRS construction
617 methods, the melanoma PRS was most strongly associated with and discriminative for
618 the melanoma phenotype for all p-value cutoffs except 5×10^{-4} . For this p-value cutoff,
619 the melanoma PRS had similar AUC and OR for the melanoma and basal cell
620 carcinoma phenotypes. This p-value cutoff represents the least conservative inclusion
621 cutoff with 1,193 included loci, and its results indicated that inclusion of too many
622 suggestive SNPs at lower thresholds may reduce PRS performance. However, we also
623 note that the most conservative cutoff (5×10^{-9}) produced a PRS with only six loci and a
624 weaker OR and AUC compare to other PRS created with less stringent cutoffs. Like the
625 other PRS construction methods, the melanoma PRS was less accurate for predicting

626 overall skin cancer compared to the individual skin cancer subtypes. The best
627 performance in terms of AUC and OR relating to the melanoma phenotype were
628 observed for p-value thresholds 5×10^{-7} and 5×10^{-8} , which included 13 and 9 loci
629 respectively. The small number of loci identified by this method at more conservative p-
630 value cutoffs may be driven by the lower sample size for melanoma in the UK Biobank
631 compared to the published melanoma GWAS meta-analyses (n cases = 2,691 and n
632 cases = 6,628, respectively). We note that the melanoma PRS constructed using the
633 UK Biobank summary statistics produced lower AUC across all p-value thresholds than
634 was seen for the Latest GWAS and GWAS Catalog PRS construction methods.

635 The PRS constructed for overall skin cancer was most strongly associated with
636 and discriminative for basal cell carcinoma across all p-value thresholds, with AUCs
637 ranging from 0.59 (95% CI [0.57, 0.60]) to 0.64 (95% CI [0.62, 0.66]) and odds ratios
638 ranging from 1.42 (95% CI [1.33, 1.51]) to 1.73 (95% CI [1.63, 1.84]). The overall skin
639 cancer PRS had the highest Brier score for overall skin cancer, indicating that the
640 overall skin cancer PRS was more accurate at predicting the skin cancer subtypes
641 compared to overall skin cancer. The overall skin cancer PRS had very similar
642 association with and discrimination abilities for the overall skin cancer phenotype across
643 all p-value thresholds except the least conservative ($p = 5 \times 10^{-4}$), for which the AUC and
644 odds ratio were smaller. Overall, the highest AUCs and strongest OR signals for both
645 PRS and all skin cancer phenotypes were found at depths of 5×10^{-7} and 5×10^{-8} .

646 In addition to associations with the primary phenotype, we explored associations
647 between PRS constructed at various UK Biobank summary statistic depths and
648 secondary phenotypes. **Fig 3** (melanoma) and **S1 Text Fig H** (overall skin cancer) show

649 PRS-PheWAS results in MGI using PRS constructed at different depths. As shown in
650 **S1 Text Table E and Fig I**, depths of 5×10^{-7} and 5×10^{-8} produced very similar results,
651 and other depths identified fewer phenotypes associated with the corresponding PRS.
652 Phenotypes that were associated with the PRS at other depths had weaker associations
653 than those observed at 5×10^{-7} and 5×10^{-8} .

654

655 **Online visual catalog: *PRSweb***

656 For comparison of the aforementioned PRS-PheWAS results and to provide
657 researchers with resources for future PRS-based analyses, we developed an open
658 access, online visual catalog *PRSweb* available at <https://statgen.github.io/PRSweb> that
659 enables interactive exploration of the PheWAS results for each of the skin cancer
660 subtypes under each of three different PRS construction methods explored in this
661 paper, for both the MGI and UK Biobank phenomes. *PRSweb* shows PRS-PheWAS
662 plots with various choices of PRS in the drop-down menu (example screenshot in **Fig 5**)
663 and offers downloadable PRS constructs (list of independent risk variants with
664 corresponding weights). Mouse-over boxes offer detailed information about top results if
665 needed, without impeding the overall user experience (grey box in **Fig 5**). Enrichment of
666 cases in the tail of the PRS distribution are presented in interactive forest plots.

667

668 **Discussion**

669 PRS combine information from a large number of genetic variants to stratify
670 subjects in terms of their risk for developing a particular disease. However, there are
671 currently no general guidelines for how to construct a PRS for a given *EHR-derived*
672 phenotype. In this paper, we explore strategies for constructing a PRS using markers

673 and weights obtained from various publicly-available sources. First, we consider PRS
674 constructed using markers and weights identified in either (1) the latest GWAS or
675 GWAS meta-analysis or (2) the NHGRI-EBI GWAS Catalog. We compare these two
676 PRS construction methods in terms of their associations with EHR-derived phenotypes
677 for the three most common skin cancer subtypes in the USA: basal cell carcinoma,
678 cutaneous squamous cell carcinoma, and melanoma.

679 A priori, we may have some belief that the latest (and often the largest) GWAS
680 may provide a better source of evidence to use for PRS construction due to larger
681 sample sizes and (potentially) more carefully curated data. The Latest GWAS and
682 GWAS Catalog methods produced PRS with similar performance in terms of their
683 associations with and discrimination for the primary phenotype used to construct the
684 PRS for both basal cell carcinoma and melanoma. Generally, PRS constructed for
685 melanoma and basal cell carcinoma were most strongly associated with and
686 discriminative for their target phenotypes, indicating that both PRS construction
687 methods were able to provide a higher degree of specificity for the intended skin cancer
688 subtype. In contrast, the PRS for squamous cell carcinoma were not more strongly
689 associated with the squamous cell carcinoma phenotype compared to other skin cancer
690 phenotypes. This may suggest a need for further exploration into genetic factors
691 uniquely related to the squamous cell carcinoma subtype.

692 For each skin cancer subtype, we performed a PRS-PheWAS to identify
693 secondary phenotypes that are associated with the corresponding PRS. We generally
694 identified many dermatological features in addition to the primary phenotype, indicating
695 the ability of PRS to reproduce associations with the primary phenotype even after

696 multiple testing corrections and covariate adjustment. The majority of these associations
697 were replicated in a PRS-PheWAS performed for the UK Biobank phenome. Our
698 analyses identified actinic keratosis, which is believed to be a precursor to squamous
699 cell and basal cell carcinoma, as an independent predictor of basal cell and squamous
700 cell carcinoma, and we demonstrated that incorporating the PRS in addition to clinical
701 information improved discrimination for future skin cancer diagnoses [46-48].

702 In an additional analysis, we identified loci that were shared among all three skin
703 cancer subtypes' PRS. Loci overlap between the PRS for the three subtypes may
704 indicate factors related to common biology between the subtypes. We noted that all
705 shared loci (*HERC2/OCA2*, *IRF4*, *MC1R*, *RALY*, *SLC45A2* and *TYR*) were also loci that
706 had been associated with human pigmentation traits and/or harbor key genes of the
707 biochemical pathway of melanogenesis [50, 54-62, 64, 67-71]. We constructed PRS
708 using SNPs shared by all three skin cancer subtypes and a PRS for skin tanning ability
709 using results from a recent GWAS meta-analysis.[70] The skin tanning ability PRS
710 PheWAS identified a very similar set of traits to the shared PRS PheWAS, suggesting
711 that the shared genetic component may in part represent genetic factors influencing the
712 skin pigmentation and the skin reaction to sun exposure. However, the PRS that are
713 unique to subtypes did not show such a common pathway or mechanism.

714 The Latest GWAS and the GWAS Catalog methods for constructing the PRS
715 involve incorporating only loci that reached genome-wide significance for at least one
716 study, as non-significant loci are usually not reported. However, incorporating non-
717 significant loci that are associated with the primary phenotype may help improve the
718 predictive ability of the PRS [8, 16]. We found that incorporating additional loci that

719 would not reach genome-wide significance did improve the PRS' ability to discriminate
720 cases from controls for the primary phenotype up to a point. In particular, PRS
721 constructed using SNPs with p-values less than 5×10^{-8} or 5×10^{-7} resulted in the best
722 performance, but further increasing the p-value threshold resulted in reduced
723 performance. Crucially, we also observed stronger associations between the PRS and
724 secondary phenotypes for PRS constructed using depths of 5×10^{-8} and 5×10^{-7} . These
725 results suggest that some benefit may be seen by incorporating loci that do not reach
726 significance into the PRS construction but incorporating too many loci with larger p-
727 values may not improve the predictive ability of the PRS (for both primary and
728 secondary phenotypes). However, this gain or reduction in PRS performance may
729 depend on the phenotype of interest and on the prevalence of the phenotype in the
730 analytical sample.

731 As a product of this study, we provide an online visual catalog *PRSweb* that
732 provides PRS-PheWAS results for the various skin cancer phenotypes for PRS
733 constructed using the different methods explored in this paper. *PRSweb* will provide a
734 routine way to compare different PRS construction methods and to explore PRS-
735 PheWAS results in detail. Additionally, *PRSweb* provides the PRS construction details,
736 which researchers can download and use in their own analyses. In the future, we plan
737 to extend this online platform to include PheWAS for many other cancer phenotypes,
738 which will make this online platform a general tool for identifying phenotypes related to
739 particular types of cancer.

740 One limitation of the generalizability of this study comes from the homogeneous
741 race profile of MGI and UK Biobank. UK Biobank consists of subjects of primarily

742 European descent, and we restricted our analyses to subjects of European descent in
743 MGI (excluding about 10% of the subjects in MGI) in order to ensure greater
744 comparability between the two datasets. Additionally, many of the existing GWAS were
745 conducted on European populations, and we wanted to consider similar samples when
746 comparing the performance of PRS constructed using summary statistics from
747 European populations. Unlike UK Biobank, MGI is not a population-based sample;
748 rather, it is a sample of patients recruited from a large academic medical center.
749 Patients were recruited prior to surgery through the anesthesiology department, and
750 therefore they may present a potential for selection bias. Additionally, the comparative
751 performance of the PRS across construction methods will depend on the phenotype of
752 interest. In spite of these limitations, a principled comparison of the various methods
753 explored in this paper may provide researchers with a sense of the robustness of their
754 PheWAS inference to the PRS construction method and an analytical framework for
755 exploration of shared genetic architecture of related traits.

756

757 **Acknowledgement**

758 The authors acknowledge the University of Michigan Medical School Central
759 Biorepository for providing biospecimen storage, management, and distribution services
760 in support of the research reported in this publication. The presented research was
761 funded by the National Cancer Institute at the National Institutes of Health (P30
762 CA046592 [LGF, LJB, MS, BM] and T32 CA83654 [RBP]). Part of this research has
763 been conducted using the UK Biobank Resource under application number 24460. This
764 material is based upon work supported by the National Science Foundation under Grant
765 No. (NSF DMS 1712933). Any opinions, findings, and conclusions or recommendations
766 expressed in this material are those of the author(s) and do not necessarily reflect the
767 views of the National Science Foundation.

768

769 **Web resources**

770 University of Michigan Medical School Central Biorepository;
771 <https://research.medicine.umich.edu/our-units/central-biorepository>
772 UK Biobank; <http://www.ukbiobank.ac.uk/>
773 UK Biobank GWAS summary statistics; <https://tinyurl.com/UKB-SAIGE>
774 TOPMed variant browser, <https://bravo.sph.umich.edu/freeze5/hg38/>
775 TOPMed program, [https://www.nhlbi.nih.gov/science/trans-omics-precision-medicine-](https://www.nhlbi.nih.gov/science/trans-omics-precision-medicine-topmed-program)
776 [topmed-program](https://www.nhlbi.nih.gov/science/trans-omics-precision-medicine-topmed-program)
777 Minimac4; <https://genome.sph.umich.edu/wiki/Minimac4>
778 BCFtools; <https://samtools.github.io/bcftools/bcftools.html>
779 KING; <http://people.virginia.edu/~wc9c/KING/>
780

- 781 FASTINDEP; <https://github.com/endrebak/fastindep>
- 782 PLINK; <https://www.cog-genomics.org/plink2/>
- 783 Eagle; <https://data.broadinstitute.org/alkesgroup/Eagle/>
- 784 UCSC Genome Browser; <http://genome.ucsc.edu/>
- 785 R; <https://cran.r-project.org/>
- 786 NHGRI-EBI GWAS Catalog; <https://www.ebi.ac.uk/gwas/>
- 787 dbSNP; <https://www.ncbi.nlm.nih.gov/projects/SNP/>
- 788 Imputation server; <https://imputationserver.sph.umich.edu/>
- 789 Jinja, <https://github.com/pallets/jinja>
- 790 Locuszoom, <https://github.com/statgen/locuszoom>
- 791 PRSweb; <https://statgen.github.io/PRSweb>

792 References

793

- 794 1. MacArthur J, Bowler E, Cerezo M, Gil L, Hall P, Hastings E, et al. The new
795 NHGRI-EBI Catalog of published genome-wide association studies (GWAS Catalog).
796 *Nucleic Acids Res.* 2017;45(D1):D896-D901.
- 797 2. Wray NR, Goddard ME, Visscher PM. Prediction of individual genetic risk to
798 disease from genome-wide association studies. *Genome Res.* 2007;17(10):1520-8.
- 799 3. Khera AV, Chaffin M, Aragam K, Emdin CA, Klarin D, Haas M, et al. Genome-
800 wide polygenic score to identify a monogenic risk-equivalent for coronary disease.
801 *bioRxiv.* 2017.
- 802 4. Fritsche LG, Gruber SB, Wu Z, Schmidt EM, Zawistowski M, Moser SE, et al.
803 Association of Polygenic Risk Scores for Multiple Cancers in a Phenome-wide Study:
804 Results from The Michigan Genomics Initiative. *Am J Hum Genet.* 2018;102(6):1048-
805 61.
- 806 5. Docherty AR, Moscati A, Dick D, Savage JE, Salvatore JE, Cooke M, et al.
807 Polygenic prediction of the phenome, across ancestry, in emerging adulthood. *Psychol*
808 *Med.* 2017:1-10.
- 809 6. Millard LAC, Davies NM, Gaunt TR, Davey Smith G, Tilling K. Software
810 Application Profile: PHESANT: a tool for performing automated phenome scans in UK
811 Biobank. *Int J Epidemiol.* 2017.
- 812 7. Zhou W, Nielsen JB, Fritsche LG, Dey R, Elvestad MB, Wolford BN, et al.
813 Efficiently controlling for case-control imbalance and sample relatedness in large-scale
814 genetic association studies. *bioRxiv.* 2017.
- 815 8. Dudbridge F. Power and predictive accuracy of polygenic risk scores. *PLoS*
816 *Genet.* 2013;9(3):e1003348.
- 817 9. Torkamani A, Wineinger NE, Topol EJ. The personal and clinical utility of
818 polygenic risk scores. *Nat Rev Genet.* 2018.
- 819 10. Steyerberg EW, Vickers AJ, Cook NR, Gerds T, Gonen M, Obuchowski N, et al.
820 Assessing the performance of prediction models: a framework for traditional and novel
821 measures. *Epidemiology.* 2010;21(1):128-38.
- 822 11. Kauvar AN, Cronin T, Jr., Roenigk R, Hruza G, Bennett R, American Society for
823 Dermatologic S. Consensus for nonmelanoma skin cancer treatment: basal cell
824 carcinoma, including a cost analysis of treatment methods. *Dermatol Surg.*
825 2015;41(5):550-71.
- 826 12. Kallini JR, Hamed N, Khachemoune A. Squamous cell carcinoma of the skin:
827 epidemiology, classification, management, and novel trends. *Int J Dermatol.*
828 2015;54(2):130-40.
- 829 13. Berwick M, Buller DB, Cust A, Gallagher R, Lee TK, Meyskens F, et al.
830 Melanoma Epidemiology and Prevention. *Cancer Treat Res.* 2016;167:17-49.
- 831 14. Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, et al. Genome-
832 wide genetic data on ~500,000 UK Biobank participants. *bioRxiv.* 2017.
- 833 15. Reisberg S, Iljasenko T, Lall K, Fischer K, Vilo J. Comparing distributions of
834 polygenic risk scores of type 2 diabetes and coronary heart disease within different
835 populations. *PLoS One.* 2017;12(7):e0179238.

- 836 16. Evans DM, Visscher PM, Wray NR. Harnessing the information contained within
837 genome-wide association studies to improve individual prediction of complex disease
838 risk. *Hum Mol Genet.* 2009;18(18):3525-31.
- 839 17. Euesden J, Lewis CM, O'Reilly PF. PRSice: Polygenic Risk Score software.
840 *Bioinformatics.* 2015;31(9):1466-8.
- 841 18. Vilhjalmsson BJ, Yang J, Finucane HK, Gusev A, Lindstrom S, Ripke S, et al.
842 Modeling Linkage Disequilibrium Increases Accuracy of Polygenic Risk Scores. *Am J*
843 *Hum Genet.* 2015;97(4):576-92.
- 844 19. So HC, Sham PC. Improving polygenic risk prediction from summary statistics by
845 an empirical Bayes approach. *Sci Rep.* 2017;7:41262.
- 846 20. Sudlow C, Gallacher J, Allen N, Beral V, Burton P, Danesh J, et al. UK biobank:
847 an open access resource for identifying the causes of a wide range of complex diseases
848 of middle and old age. *PLoS Med.* 2015;12(3):e1001779.
- 849 21. Wang C, Zhan X, Bragg-Gresham J, Kang HM, Stambolian D, Chew EY, et al.
850 Ancestry estimation and control of population stratification for sequence-based
851 association studies. *Nat Genet.* 2014;46(4):409-15.
- 852 22. Li JZ, Absher DM, Tang H, Southwick AM, Casto AM, Ramachandran S, et al.
853 Worldwide human relationships inferred from genome-wide patterns of variation.
854 *Science.* 2008;319(5866):1100-4.
- 855 23. Manichaikul A, Mychaleckyj JC, Rich SS, Daly K, Sale M, Chen WM. Robust
856 relationship inference in genome-wide association studies. *Bioinformatics.*
857 2010;26(22):2867-73.
- 858 24. Abraham KJ, Diaz C. Identifying large sets of unrelated individuals and unrelated
859 markers. *Source Code Biol Med.* 2014;9(1):6.
- 860 25. McCarthy S, Das S, Kretzschmar W, Delaneau O, Wood AR, Teumer A, et al. A
861 reference panel of 64,976 haplotypes for genotype imputation. *Nat Genet.*
862 2016;48(10):1279-83.
- 863 26. Loh PR, Danecek P, Palamara PF, Fuchsberger C, Y AR, H KF, et al.
864 Reference-based phasing using the Haplotype Reference Consortium panel. *Nat Genet.*
865 2016;48(11):1443-8.
- 866 27. Carroll RJ, Bastarache L, Denny JC. R PheWAS: data analysis and plotting tools
867 for phenome-wide association studies in the R environment. *Bioinformatics.*
868 2014;30(16):2375-6.
- 869 28. Ho DE, Imai K, King G, Stuart EA. MatchIt: Nonparametric Preprocessing for
870 Parametric Causal Inference. *Journal of Statistical Software.* 2011;42(8):1-28.
- 871 29. Michailidou K, Lindstrom S, Dennis J, Beesley J, Hui S, Kar S, et al. Association
872 analysis identifies 65 new breast cancer risk loci. *Nature.* 2017;551(7678):92-4.
- 873 30. Chahal HS, Lin Y, Ransohoff KJ, Hinds DA, Wu W, Dai HJ, et al. Genome-wide
874 association study identifies novel susceptibility loci for cutaneous squamous cell
875 carcinoma. *Nat Commun.* 2016;7:12048.
- 876 31. Chahal HS, Wu W, Ransohoff KJ, Yang L, Hedlin H, Desai M, et al. Genome-
877 wide association study identifies 14 novel risk alleles associated with basal cell
878 carcinoma. *Nat Commun.* 2016;7:12510.
- 879 32. Ransohoff KJ, Wu W, Cho HG, Chahal HC, Lin Y, Dai HJ, et al. Two-stage
880 genome-wide association study identifies a novel susceptibility locus associated with
881 melanoma. *Oncotarget.* 2017;8(11):17586-92.

- 882 33. Welter D, MacArthur J, Morales J, Burdett T, Hall P, Junkins H, et al. The NHGRI
883 GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Res.*
884 2014;42(Database issue):D1001-6.
- 885 34. Winkler TW, Day FR, Croteau-Chonka DC, Wood AR, Locke AE, Magi R, et al.
886 Quality control and conduct of genome-wide association meta-analyses. *Nat Protoc.*
887 2014;9(5):1192-212.
- 888 35. GenABEL project developers. GenABEL: genome-wide SNP association
889 analysis. 2013.
- 890 36. Heinze G. A comparative investigation of methods for logistic regression with
891 separated or nearly separated data. *Stat Med.* 2006;25(24):4216-26.
- 892 37. Heinze G, Ploner M, Dunkler D, Southworth H. *logistf: Firth's bias reduced*
893 *logistic regression.* 2013.
- 894 38. Choi L, Beck C. *EHR: Electronic Health Record (EHR) Data Processing and*
895 *Analysis Tool.* 2017.
- 896 39. Robin X, Turck N, Hainard A, Tiberti N, Lisacek F, Sanchez JC, et al. pROC: an
897 open-source package for R and S+ to analyze and compare ROC curves. *BMC*
898 *Bioinformatics.* 2011;12:77.
- 899 40. Hosmer DW, Lemeshow S. *Applied Logistic Regression.* New York, USA: John
900 Wiley and Sons; 2010.
- 901 41. Lele S, R., Keim JL, Solymos P. *ResourceSelection: Resource Selection*
902 *(Probability) Functions for Use-Availability Data.* 2017.
- 903 42. Signorell A. *DescTools: Tools for Descriptive Statistics.* 2018.
- 904 43. R Core Team. *R: A Language and Environment for Statistical Computing.* R
905 Foundation for Statistical Computing, Vienna, Austria; 2016.
- 906 44. McKinney W. *Data Structures for Statistical Computing in Python*2010.
- 907 45. Baker SG, Schuit E, Steyerberg EW, Pencina MJ, Vickers A, Moons KG, et al.
908 How to interpret a small increase in AUC with an additional risk prediction marker:
909 decision analysis comes through. *Stat Med.* 2014;33(22):3946-59.
- 910 46. Fuchs A, Marmur E. The kinetics of skin cancer: progression of actinic keratosis
911 to squamous cell carcinoma. *Dermatol Surg.* 2007;33(9):1099-101.
- 912 47. Cohen JL. Actinic keratosis treatment as a key component of preventive
913 strategies for nonmelanoma skin cancer. *J Clin Aesthet Dermatol.* 2010;3(6):39-44.
- 914 48. Jacobs RJ, Phillips G. Basal cell carcinoma mistaken for actinic keratosis. *Clin*
915 *Exp Optom.* 2006;89(3):171-5.
- 916 49. Ko CJ. Actinic keratosis: facts and controversies. *Clin Dermatol.* 2010;28(3):249-
917 53.
- 918 50. Adhikari K, Fontanil T, Cal S, Mendoza-Revilla J, Fuentes-Guajardo M, Chacon-
919 Duque JC, et al. A genome-wide association scan in admixed Latin Americans identifies
920 loci influencing facial and scalp hair features. *Nat Commun.* 2016;7:10815.
- 921 51. Asgari MM, Wang W, Ioannidis NM, Itnyre J, Hoffmann T, Jorgenson E, et al.
922 Identification of Susceptibility Loci for Cutaneous Squamous Cell Carcinoma. *J Invest*
923 *Dermatol.* 2016;136(5):930-7.
- 924 52. Barrett JH, Iles MM, Harland M, Taylor JC, Aitken JF, Andresen PA, et al.
925 Genome-wide association study identifies three new melanoma susceptibility loci. *Nat*
926 *Genet.* 2011;43(11):1108-13.

- 927 53. Bishop DT, Demenais F, Iles MM, Harland M, Taylor JC, Corda E, et al.
928 Genome-wide association study identifies three loci associated with melanoma risk. *Nat*
929 *Genet.* 2009;41(8):920-5.
- 930 54. Eriksson N, Macpherson JM, Tung JY, Hon LS, Naughton B, Saxonov S, et al.
931 Web-based, participant-driven studies yield novel genetic associations for common
932 traits. *PLoS Genet.* 2010;6(6):e1000993.
- 933 55. Han J, Kraft P, Nan H, Guo Q, Chen C, Qureshi A, et al. A genome-wide
934 association study identifies novel alleles associated with hair color and skin
935 pigmentation. *PLoS Genet.* 2008;4(5):e1000074.
- 936 56. Hernandez-Pacheco N, Flores C, Alonso S, Eng C, Mak AC, Hunstman S, et al.
937 Identification of a novel locus associated with skin colour in African-admixed
938 populations. *Sci Rep.* 2017;7:44548.
- 939 57. Jacobs LC, Hamer MA, Gunn DA, Deelen J, Lall JS, van Heemst D, et al. A
940 Genome-Wide Association Study Identifies the Skin Color Genes IRF4, MC1R, ASIP,
941 and BNC2 Influencing Facial Pigmented Spots. *J Invest Dermatol.* 2015;135(7):1735-
942 42.
- 943 58. Law MH, Medland SE, Zhu G, Yazar S, Vinuela A, Wallace L, et al. Genome-
944 Wide Association Shows that Pigmentation Genes Play a Role in Skin Aging. *J Invest*
945 *Dermatol.* 2017;137(9):1887-94.
- 946 59. Lin BD, Mbarek H, Willemsen G, Dolan CV, Fedko IO, Abdellaoui A, et al.
947 Heritability and Genome-Wide Association Studies for Hair Color in a Dutch Twin Family
948 Based Sample. *Genes (Basel).* 2015;6(3):559-76.
- 949 60. Liu F, Wollstein A, Hysi PG, Ankra-Badu GA, Spector TD, Park D, et al. Digital
950 quantification of human eye color highlights genetic association of three new loci. *PLoS*
951 *Genet.* 2010;6(5):e1000934.
- 952 61. Liu F, Visser M, Duffy DL, Hysi PG, Jacobs LC, Lao O, et al. Genetics of skin
953 color variation in Europeans: genome-wide association studies with functional follow-up.
954 *Hum Genet.* 2015;134(8):823-35.
- 955 62. Nan H, Kraft P, Qureshi AA, Guo Q, Chen C, Hankinson SE, et al. Genome-wide
956 association study of tanning phenotype in a population of European ancestry. *J Invest*
957 *Dermatol.* 2009;129(9):2250-7.
- 958 63. Nan H, Xu M, Kraft P, Qureshi AA, Chen C, Guo Q, et al. Genome-wide
959 association study identifies novel alleles associated with risk of cutaneous basal cell
960 carcinoma and squamous cell carcinoma. *Hum Mol Genet.* 2011;20(18):3718-24.
- 961 64. Rawofi L, Edwards M, Krithika S, Le P, Cha D, Yang Z, et al. Genome-wide
962 association study of pigmentary traits (skin and iris color) in individuals of East Asian
963 ancestry. *PeerJ.* 2017;5:e3951.
- 964 65. Siiskonen SJ, Zhang M, Li WQ, Liang L, Kraft P, Nijsten T, et al. A Genome-Wide
965 Association Study of Cutaneous Squamous Cell Carcinoma among European
966 Descendants. *Cancer Epidemiol Biomarkers Prev.* 2016;25(4):714-20.
- 967 66. Song F, Amos CI, Lee JE, Lian CG, Fang S, Liu H, et al. Identification of a
968 melanoma susceptibility locus and somatic mutation in TET2. *Carcinogenesis.*
969 2014;35(9):2097-101.
- 970 67. Stokowski RP, Pant PV, Dadd T, Fereday A, Hinds DA, Jarman C, et al. A
971 genomewide association study of skin pigmentation in a South Asian population. *Am J*
972 *Hum Genet.* 2007;81(6):1119-32.

- 973 68. Sulem P, Gudbjartsson DF, Stacey SN, Helgason A, Rafnar T, Magnusson KP,
974 et al. Genetic determinants of hair, eye and skin pigmentation in Europeans. *Nat Genet.*
975 2007;39(12):1443-52.
- 976 69. Zhang M, Song F, Liang L, Nan H, Zhang J, Liu H, et al. Genome-wide
977 association studies identify several new loci associated with pigmentation traits and skin
978 cancer risk in European Americans. *Hum Mol Genet.* 2013;22(14):2948-59.
- 979 70. Visconti A, Duffy DL, Liu F, Zhu G, Wu W, Chen Y, et al. Genome-wide
980 association study in 176,678 Europeans reveals genetic loci for tanning response to sun
981 exposure. *Nat Commun.* 2018;9(1):1684.
- 982 71. Sturm RA, Duffy DL. Human pigmentation genes under environmental selection.
983 *Genome Biol.* 2012;13(9):248.
984

985 **Supporting information**

986 **S1 File. Supporting Material.** This file contains supporting Figures A-J and Tables A-E.

987 **S2 File. Supporting Tables.** This Excel file contains the following tables

- 988 ● Sheet 1: Table F, ICD9 codes to PheWAS code translations
- 989 ● Sheet 2: Table G, ICD10 codes to PheWAS code translations
- 990 ● Sheet 3: Table H, MGI and UK Biobank Phenome
- 991 ● Sheet 4: Table I, Risk SNP Selection
- 992 ● Sheet 5: Table J, Risk SNP Selection Depth
- 993 ● Sheet 6: Table K, Omnibus Table Significant Results

994

995 **Figure Titles and Legends**

996 **Fig 1. PRS-PheWAS in MGI and UKB phenomes.**

997 The horizontal line indicates phenome-wide significance.

998

999 **Fig 2 Comparison of predictors.**

1000 Actinic keratosis (AK), at least 365 prior to any skin cancer diagnosis as predictor for
1001 basal cell carcinoma (BCC) (A and B) and squamous cell carcinoma (SCC) (C and D).
1002 The PRS for BCC and SCC as well as the combined predictors are shown for
1003 comparison.

1004

1005 **Fig 3 Overlap between the three skin cancer trait loci.**

1006 Reported risk SNPs within 1 Mb were merged into the same locus. Loci that were also
1007 reported to be associated with skin tanning ability are highlighted in bold. Loci were
1008 named according to the closest RefSeq genes (except *M1CR* a 385 kb locus with 16
1009 RefSeq genes and *HV745896* named after a nearby, uncurated mRNA sequence).

1010

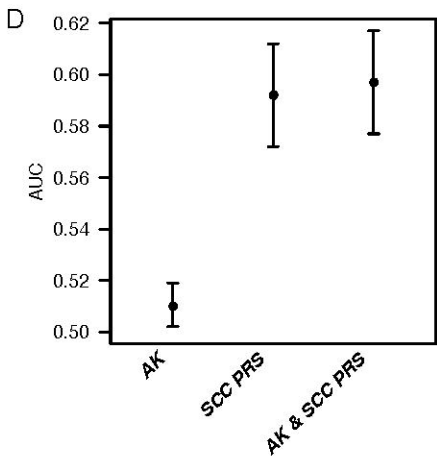
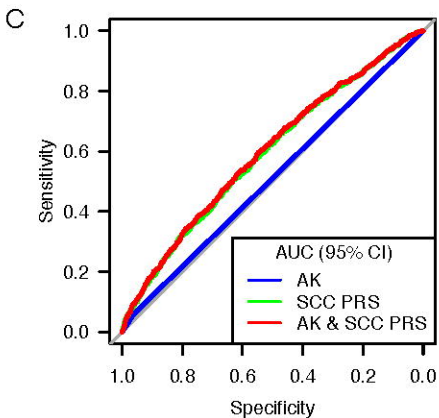
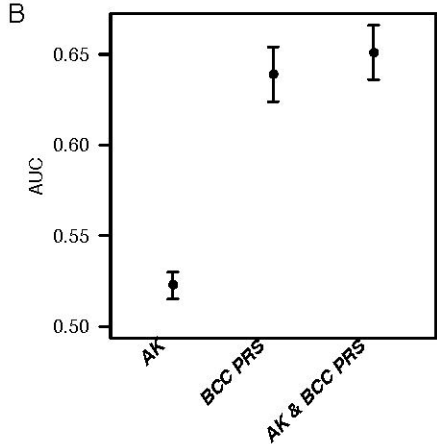
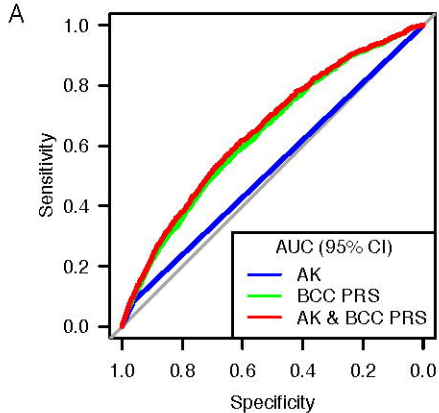
1011 **Fig 4 PheWAS on melanoma PRS constructed using UK Biobank statistics at**
1012 **different depths.**

1013 Results are shown with increasing depth from (A – F): $P \leq 5 \times 10^{-9}$, 5×10^{-8} , 5×10^{-7} , 5×10^{-6} ,
1014 5×10^{-5} , 5×10^{-4} .

1015

1016 **Fig 5. Example view from PRSweb.** A selection menu on top allows selection of PRS
1017 constructs and phenome while interactive plots with “PheWAS results”, “Exclusion

1018 PheWAS results”, and “Associations between PRS and Selected Phenotype” are
1019 generated after selection.



Melanoma

CASP8, ALS2CR12
CDKAL1, SOX4, LINC00340
CLPTM1L
LOC646329, FLJ43663
MTAP, CDKN2A, CDKN2B-AS1
MX2, MX1, TMPRSS2
OBFC1

ARNT
CDH1
FTO
GPR37
KLF4
LOC154092
LOC401177
MIR614
PARP1
RMDN2
STMN3
TERC
TMEM38B, ZNF462
TPCN2, MYEOV

AGR

HERC2, OCA2
IRF4
MC1R
RALY
SLC45A2
TYR

CUX1
FOXP1
HV745896
KRT5
MIR3939
FAM49A
NEU1
PADI6
PLIN3
RGS22

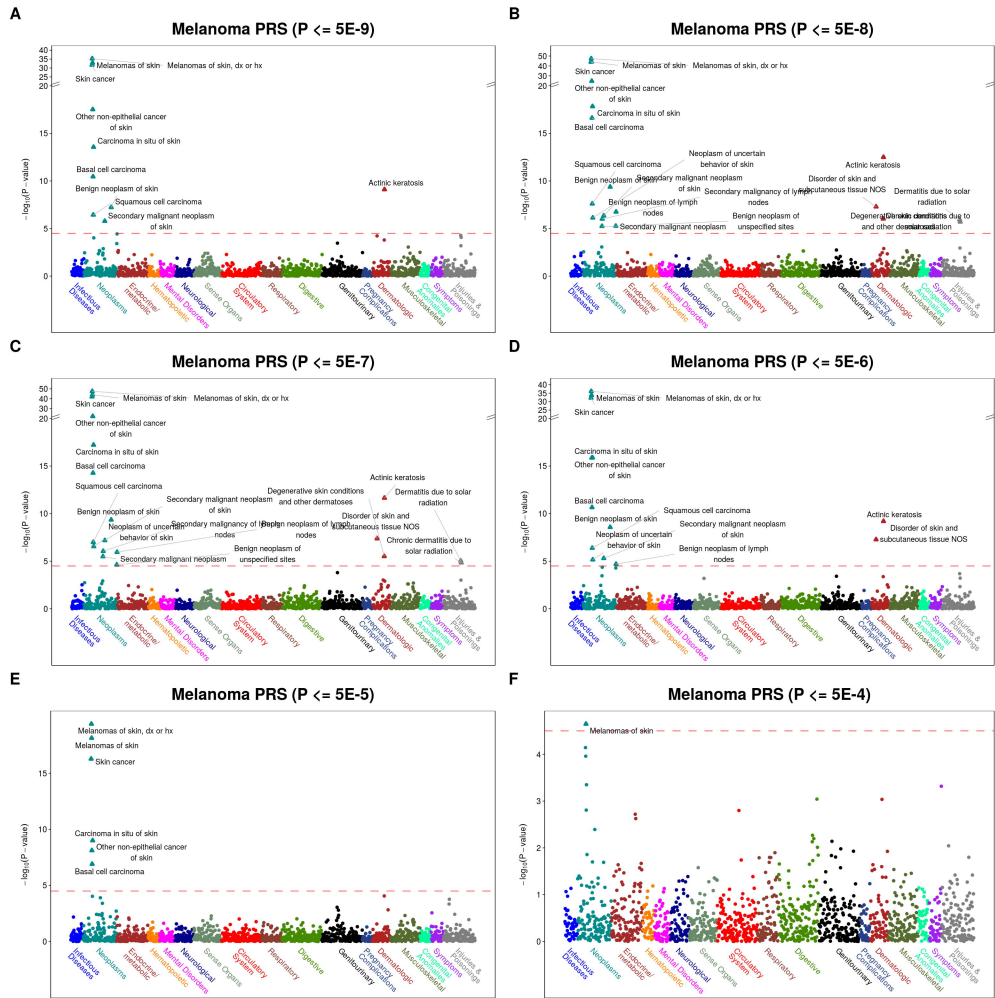
RHOA
TGM3
TNS3
TP53
UBAC2
ZBTB10
ZFHX4-AS1
ZNF365

BNC2

LINC00900
SEC16A

Basal Cell
Carcinoma

Squamous
Cell Carcinoma



PRS trait

Squamous Cell Carcinoma (172.22) ▾

PRS weights

GWAS-Catalog ▾

Study

MGI ▾

Figure 1: Squamous Cell Carcinoma PRS (172.22)

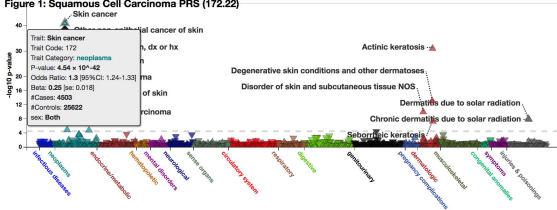


Figure 2: Squamous Cell Carcinoma PRS (172.22) (exclusion)

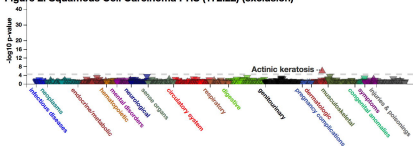
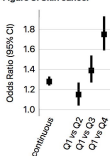


Figure 3: Skin cancer


[Download weights for GRCh37](#)
[Download weights for GRCh38](#)