

Association of Polygenic Risk Scores for Multiple Cancers in a Phenome-wide Study: Results from The Michigan Genomics Initiative

Lars G. Fritsche,^{1,2,3} Stephen B. Gruber,⁴ Zhenke Wu,^{1,5} Ellen M. Schmidt,¹ Matthew Zawistowski^{1,2}, Stephanie E. Moser,⁶ Victoria M. Blanc,⁷ Chad M. Brummett,^{6,8} Sachin Kheterpal,^{6,8} Gonçalo R. Abecasis,^{1,2} Bhramar Mukherjee,^{1,2,5,9,10,11*}

1. Department of Biostatistics, University of Michigan School of Public Health, Ann Arbor, MI 48109, USA.
2. Center for Statistical Genetics, University of Michigan School of Public Health, Ann Arbor, MI 48109, USA.
3. K.G. Jebsen Center for Genetic Epidemiology, Department of Public Health and Nursing, Norwegian University of Science and Technology, 7491 Trondheim, Sør-Trøndelag, Norway.
4. USC Norris Comprehensive Cancer Center, Los Angeles, CA 90033, USA.
5. Michigan Institute for Data Science, University of Michigan, Ann Arbor, MI 48109, USA.
6. Division of Pain Medicine, Department of Anesthesiology, University of Michigan Medical School, Ann Arbor, MI 48109, USA.
7. Central Biorepository, University of Michigan Medical School, Ann Arbor, MI 48109, USA.
8. Institute for Healthcare Policy and Innovation, University of Michigan, Ann Arbor, MI 48109, USA.
9. Department of Epidemiology, University of Michigan School of Public Health, Ann Arbor, MI 48109, USA.
10. University of Michigan Comprehensive Cancer Center, University of Michigan, Ann Arbor, MI 48109, USA.
11. Present Address: Department of Biostatistics and Epidemiology, University of Michigan School of Public Health, 1415 Washington Heights, SPH 1 Room 4619, Ann Arbor, MI 48109, USA.

* bhramar@umich.edu

Abstract

Health systems are stewards of patient electronic health record (EHR) data with extraordinarily rich depth and breadth, reflecting thousands of diagnoses and exposures. Measures of genomic variation integrated with EHRs offer a potential strategy to accurately stratify patients for risk profiling and discover new relationships between diagnoses and genomes. The objective of this study was to evaluate whether Polygenic Risk Scores (PRS) for common cancers are associated with multiple phenotypes in a Phenome-wide Association Study (PheWAS) conducted in 28,260 unrelated, genotyped patients of recent European ancestry who consented to participate in the Michigan Genomics Initiative, a longitudinal biorepository effort within Michigan Medicine. PRS for 12 cancer traits were calculated using summary statistics from the NHGRI-EBI catalog. A total of 1,711 synthetic case-control studies was used for PheWAS analyses. There were 13,490 (47.7%) patients with at least one cancer diagnosis in this study sample. PRSs exhibited strong association for several cancer traits they were designed for including female breast cancer, prostate cancer, melanoma, basal cell carcinoma, squamous cell carcinoma and thyroid cancer. Phenome-wide significant associations were observed between PRS and many non-cancer diagnoses. To differentiate PRS associations driven by the primary trait from associations arising through shared genetic risk profiles, the idea of “exclusion PRS PheWAS” was introduced. This approach led to phenome-wide significant associations between a lower risk for hypothyroidism in patients with high thyroid cancer PRS and a higher risk for actinic keratosis in patients with high squamous cell carcinoma PRS after removing all cases of the primary cancer trait. Further analysis of temporal order of the diagnoses improved our understanding of

these secondary associations. This is the first comprehensive PheWAS study using PRS instead of a single variant.

Introduction

In the past decade, genome-wide association studies (GWAS) using single nucleotide polymorphisms (SNPs) led to discovery of many common disease susceptibility loci¹⁻³. An alternative agnostic way of exploring gene-disease association is through phenome-wide association studies (PheWAS)⁴⁻⁶. PheWAS enable simultaneous exploration of the association between genetic variants and a broad spectrum of physiological/clinical phenotypes. To explore the joint genome x phenome landscape, one needs access to both Electronic Health Records (EHRs) and GWAS data. The promise and potential of these studies have recently been illustrated by the electronic Medical Records and Genomics (eMERGE) network^{7; 8}. Beyond genetic associations, EHR has enabled discovery of new associations between disease and secondary effects of drugs or blood biomarker levels⁹⁻¹¹.

PheWAS have been used to both replicate known genetic-phenotypic associations and to discover new consequences for disease associated variants. PheWAS use computable phenotypes derived from EHR databases. Traditional PheWAS have used International Classification of Disease (ICD) codes to define a set of computable phenotypes or “PheWAS codes” defined and validated by experts using a combination of ICD codes¹². Standard PheWAS have primarily focused on correlating genetic variants, **one at a time**, to a spectrum of phenotypes. When each variant is associated with a small effect size, these studies can only provide limited insight. For this reason, many areas of genetics now use ensembles of variants that cumulatively explain substantial variation in disease

risk¹³⁻¹⁵. For example, PRS constructed from multiple GWAS identified loci that have been proposed for cancer screening, risk prediction and risk stratification¹⁶⁻¹⁹.

In this paper, we introduce the new concept of exploring PRS in a PheWAS setting instead of a traditional PheWAS that considers single variants, one at a time. We focus on cancer traits while constructing the PRS. We construct PRS for multiple cancers including some of the most common groupings of cancers in the United States: prostate cancer (PCa, MIM: 176807), breast cancer (MIM: 114480), colorectal cancer (MIM: 114500), lung cancer (MIM: 211980), melanoma of skin (MIM: 155601) and basal cell carcinoma (MIM: 614740) and correlate them with PheWAS codes.

Our study is based on the Michigan Genomics Initiative (MGI) launched in 2012, a biorepository effort to create a longitudinal cohort of participants in Michigan Medicine.

MGI enrolled participants undergoing anesthesia prior to a surgery or diagnostic procedure, creating a patient community with genome-wide data, electronic health information, and permission for follow-up and re-contact in future studies. Our current analysis of 28,260 patients in MGI indicates that 47.7% of these patients have at least one current or previous neoplasm diagnosis (excluding benign neoplasms). This presents a unique opportunity to study multiple cancer outcomes leveraging both EHR and genomic data in MGI.

At the same time, this enrichment of cancer patients in MGI highlights some of the special features of the sampling frame for the study and source population. Because of the self-selective, consent-based nature of MGI patient enrollment, the sample selection mechanism is non-probabilistic, that is, the probability of a sampling unit being included in the study is not pre-determined. The MGI sample is enriched for neoplasm diagnoses

which could be related to the fact that many surgeries are diagnostic procedures are specifically related to cancer treatment and screening (e.g., colonoscopy, skin biopsies). Cancer patients undergo surgery more frequently than the general population and frequently choose an academic medical center for diagnostic and/or interventional procedures. The analytic framework presented in this paper (conducts careful sensitivity analysis for protecting our inference against such selection biases, unbalanced case-control ratios, and phenotypic enrichment).

There are several innovative aspects to our study. Our study represents the first comprehensive PheWAS focused on using **PRS in a cancer-enriched cohort** accrued in an academic health center. Our study is also the first PheWAS focused **on cancer**. Our results demonstrate PRS, a summary score constructed based on results of large population-based GWAS, can be potentially useful for cancer risk stratification among patients in an academic medical center. We also note that when a PRS-based PheWAS leads to the association of a cancer-specific PRS (e.g., prostate cancer PRS) with other secondary related phenotypes (e.g., erectile dysfunction or urinary incontinence), these findings may require careful consideration. We observe that many of these secondary associations are often driven by the primary cancer diagnosis. We introduce the notion of “exclusion PRS PheWAS” to detect independent secondary associations that have shared genetic etiology. We extract the temporal order of diagnoses from the EHR to shed further insight into these secondary associations.

Subjects and Methods

MGI cohort

Participants were recruited through Michigan Medicine health system while awaiting diagnostic or interventional procedures either during a preoperative visit prior to the procedure or on the day of procedure that required anesthesia. Opt-in written informed consent is obtained. In addition to coded biosamples and protected secure health information, participants understand that all EHR, claims, and national data sources linkable to the participant may be incorporated into the MGI databank. Each participant donates a blood sample for genetic analysis, undergoes baseline vital signs and a comprehensive history and physical, and completes validated self-report measures of pain, mood and function, including NIH Patient Reported Outcomes Measurement Information System (PROMIS) measures. Data were collected according to Declaration of Helsinki principles. Study participants provided written informed consent, and protocols were reviewed and approved by local ethics committees (IRB ID HUM00099605). In the current study, we report results obtained from 28,260 genotyped samples of European ancestry with available integrated EHR data (see summary characteristics of the cohort in **Table 1**).

Genotyping and Sample Quality Control (QC)

DNA from 37,412 blood samples was genotyped on two batches of customized Illumina HumanCoreExome v12.1 bead arrays (“UM_HUNT_Biobank_11788091_A1” [N = 21,207] and “UM_HUNT_Biobank_v1-1_20006200_A” [N = 16,205]) that in addition to standard genome-wide tagging SNPs (~N=240,000) and exomic variants (N=~280,000) contained about 70,000 additional custom content variants, e.g. candidate variants from GWAS experiments, nonsense and missense variants from sequencing studies, ancestry informative markers, and Neanderthal variants. Genotype analysis was performed with

Illumina GenomeStudio (module 1.9.4, algorithm GenTrain 2.0). After initial clustering, variant cluster boundaries were re-defined in a second run using only individuals with call rate of at least 99% and genotyped the remaining samples afterwards.

We excluded samples with: (1) call rate <99%, (2) estimated contamination >2.5 % (BAF Regress)²⁰, (3) large chromosomal copy number variants (single chromosome with missingness \geq five times larger than other chromosomes), (4) lower call rate than its technical duplicate or twin, (5) gonosomal constellations other than XX and XY, or (6) whose inferred sex did not match the reported gender. We excluded variants if: (1) their probes could not be perfectly mapped or mapped perfectly to multiple position in the human genome assembly (Genome Reference Consortium Human genome build 37 and revised Cambridge Reference Sequence of the human mitochondrial DNA; BLAT)²¹, (2) they showed deviations from Hardy Weinberg equilibrium in European ancestry samples ($P < 0.0001$), (3) had a call rate <99%, (4) another variant with higher call rate assayed the same variant or (5) if the allele frequency differences between the two array versions within unrelated, European ancestry samples had a P-value < 0.001 (PLINK v1.90)²². After quality control, 392,323 polymorphic variants remained.

Before preparing the final analytical data set, we reduced the data to 33,028 samples for which we had complete age and ICD9 data available. Next, we estimated pairwise kinship with the software KING²³ and limited further analysis to a subset that contained no pairs of individuals with a 1st- or 2nd-degree relationship. We inferred recent ancestry by projecting all genotyped samples into the space of the principal components of the Human Genome Diversity Project reference panel using PLINK (938 unrelated individuals)^{24; 25}. We limited the principal component analysis to variants that were shared between the

HGDP reference and the MGI data, had a minor allele frequency >1%, and remained after LD pruning ($r^2 < 0.5$; PLINK). Samples of recent European ancestry (~90% of participants) were defined as samples that fell into a circle around the center of the European HGDP populations in the PC1 versus PC2 space, whereas the circle's radius was set to 1/8 of the distance between the center of the European HGDP populations and the centroid of the centers of the European, East Asian and Sub-Saharan populations (**Figure S1**). Principal components were stored and used for further association tests. After quality control, 28,260 unrelated, genotyped individuals of recent European ancestry with age and ICD9 data remained for further analysis.

Phasing and Genotype Imputation

We imputed genotypes of the Haplotype Reference Consortium using the Michigan Imputation Server²⁶ and filtered poorly imputed variants with $R^2 < 0.3$ and/or minor allele frequency (MAF) < 0.1% resulting in over 17 million imputed variants available after quality control and filtering. The obtained accuracy for imputed variants, i.e. the average empirical R^2 values for different MAF frequency bins, was: 0.89 ($0.1\% \leq \text{MAF} \leq 0.5\%$), 0.94 ($0.5\% < \text{MAF} \leq 5\%$), and 0.96 ($\text{MAF} > 5\%$).

Phenome Generation

We extracted the ICD9 data for 28,260 unrelated, genotyped individuals of recent European ancestry and mapped a total of 3.5 million ICD9 codes to PheWAS codes (PheWAS translation table version 1.2)¹². The ICD9 codes (10,322 unique ICD9 codes) were aggregated to PheWAS traits using the PheWAS R package¹². Cases for a given PheWAS code were defined if an individual had at least one assignment of that PheWAS code in their record. The remaining individuals that did not have overlapping PheWAS

codes that are a part of the exclusion criteria were considered as controls. A total of 1,857 case control studies were generated of which 1,711 with ≥ 20 cases were used for further analyses (see **Table S1**; phenotypes with < 50 cases were coded as “ <50 ”).

GWAS Catalog SNP Extraction and Construction of PRS

We downloaded previously reported GWAS variants from the NHGRI-EBI Catalog (file date: June 31, 2017) ^{27; 28}. None of the discovery studies included in the catalog used any subset of the MGI cohort. This is primarily because MGI started recruiting in 2012 and the genotype data only became available recently. Variant positions were converted to GRCh37 using variant IDs from dbSNP version 144 (UCSC genome browser) after updating outdated dbSNP IDs to their merged dbSNP IDs. Entries with missing risk alleles, risk allele frequencies, or odds ratios were excluded. We corrected alleles of non-ambiguous SNPs to the forward strand of the genomic reference sequence so that the reported risk allele matched one of the alleles found at the corresponding position in the 1000 Genomes Project genotype data. We only included entries with broad European ancestry (as reported by the NHGRI-EBI Catalog). To allow an additional quality control check, we compared the reported risk allele frequencies (RAF) in controls with the frequencies of the 503 European samples of the 1000 Genomes Project reference data (Phase 3, release 5) ²⁹. We then excluded entries whose RAF deviated more than 15% from the reference. This chosen threshold is subjective and was based on clear differentiation between correct and likely flipped alleles on the two diagonals (see **Figure S2**) as noted frequently in GWAS meta-analyses quality control procedures ³⁰. For each analyzed cancer type, we extracted overlapping GWAS hits in our genotype data and estimated pairwise LD (r^2) using the available allele dosages of the corresponding

controls. For pairwise correlated SNPs ($r^2 > 0.5$) or SNPs with multiple entries, we kept the SNP with the younger publication date (and smaller P value, if necessary) and excluded the other (**Figure S2 & Table S2**). Finally, we weighted the allele dosages of risk SNPs of the risk increasing alleles with their reported log odds ratios and calculated PRS as their sum. Namely, for subject j in MGI the PRS was of the form $PRS_j = \sum_i \beta_i G_{ij}$ where the sum extends over all included loci, β_i are the log odds ratios retrieved from the GWAS catalog for locus i and G_{ij} was the measured dosage data for the risk allele on locus i in subject j . This variable was created for each MGI participant and for each cancer separately.

Statistical Analysis

For the current study, we initially explored 30 cancer traits that had matching entries in the GWAS catalog (**Table S3**), and restricted our analysis to 12 cancer traits with at least 5 risk SNPs detected in the GWAS catalog after filtering that had relatively larger samples sizes in MGI (namely $N \geq 250$ cases) (**Table 2**). Logistic regression was used for all genetic association analysis. Firth's bias reduction method was applied to all single SNP and PRS models to resolve the problem of separation in logistic regression (Logistf in R package "EHR")³¹⁻³³, a common problem for binary or categorical outcome models when for a certain part of the covariate space there is only one observed value of the outcome which often leads to very large parameter estimates and standard errors. Firth's bias-reduction³⁴ is a penalized likelihood method that reduces the bias in such situations by adding a penalty term to the likelihood.

To estimate the association of PRS with the primary cancer phenotype, we first determined the PRS quartiles using all control samples, categorized all samples according to these PRS quartiles and fitted Firth bias corrected logistic regression adjusting for age, sex, genotyping array, and the first four principal components. We report odds ratios corresponding to the top versus the bottom quartile PRS (reference), referred to as **PRS OR**. We also used continuous PRS instead of the categorized version as the covariate for enhanced power.

To compare reported associations of individual GWAS catalog SNPs with association observed in the MGI data set, we tested the association between reported GWAS hits and its corresponding trait using Firth bias-corrected logistic regression implemented in EPACTS (version 3.3, see **Web Resources**). Age, sex, genotyping array, and principal components 1–4 were included as covariates (see ***Kinship and Ancestry Inference***).

To determine the agreement of estimated effect sizes [estimated log(odds ratios)] between the MGI case-control studies and the published GWAS catalog hits, we estimated Pearson's correlation coefficient [$\hat{\rho}$] and Lin's concordance measure between the two sets of coefficients^{35; 36}. Towards more standard discovery type genome-wide association analysis with MGI data, we performed GWAS for the 9 cancer traits where the correspondence between the effect sizes were relatively strong [$\hat{\rho} \geq 0.6$], **Table 2**). For computationally efficient GWA analysis we used the score test-based saddle point approximation (SPA)³⁷ method adjusting for age, sex, genotyping array, the first four principal components. SPA was reported to provide accurate test statistics even for extremely unbalanced case-control ratios similar to Firth bias corrected logistic regression (see below) but was estimated to be 100 times faster than the latter³⁷.

For our primary PRS-PheWAS, for the six PRS in **Table 2** that showed strong and significant association, we conducted Firth bias-corrected logistic regression by fitting a model of the following form and repeated them for each of the 1,711 phenotypes.

$\text{logit}(P(\text{Disease}=1|\text{PRS}, \text{Age}, \text{Sex}, \text{Array}, \text{PC}))$

$$=\beta_0 + \beta_{PRS}PRS + \beta_{Age}Age + \beta_{Sex}Sex + \beta_{Array}Array + \beta PC$$

where the PCs were the first four principal components obtained from the principal component analysis of the genotyped GWAS markers and where “Array” represents the two genotyping array versions used in MGI accounting for potential batch effects. To adjust for multiple testing, we applied the conservative phenome-wide Bonferroni correction according to the 1,711 analyzed PheWAS codes (**Table S1**). Through a PheWAS plot, we present $-\log_{10}(P\text{-values})$ corresponding to each of the 1,711 association tests for $H_0: \beta_{PRS} = 0$. Directional arrows on the PheWAS plot indicate whether a phenome-wide significant trait was positively or negatively associated with the PRS.

Furthermore, our extensive sensitivity analyses included: (a) similar models adjusting for 20 PCs, (b) matching cases to control and conducting conditional logistic regression analysis and (c) using the unweighted risk allele counts as predictor. The reason for these three sensitivity analyses was to check (a) if the first 4 PCs were sufficient to control for population stratification, (b) if differences in age and sex distributions or extreme case-control ratios influenced the main analysis and (c) if ignoring effect sizes and using total risk allele count produced similar results. For (b), we matched cases and controls using the R package “MatchIt” and applied nearest neighbor matching for age, PC1-4 (using Mahalanobis-metric matching; matching window caliper/width of 0.25 standard

deviations) and exact matching for sex. We considered a varying set of case:control matching ratios from 1:1 to 1:10. We observed gain in precision with increasing the number of matched controls per case, but the gain in precision became negligible after 1:10 matching ratio (**Table S6**). Moreover, for some cancers with large number of cases we could not attain 1:10 matching ratio for all cases and ended up with varying number of controls per case. For example, for prostate cancer the average number of matched controls per case was around 5³⁸.

To investigate the possibility of the secondary trait associations with PRS being completely driven by the primary trait association, we performed a second set of PheWAS after excluding individuals affected with the primary cancer trait for which the PRS was constructed, referred to as “exclusion PRS PheWAS”. We applied “exclusion PRS PheWAS” instead of a PRS PheWAS that uses the primary cancer trait as covariate, because the control exclusion criteria implemented in the PheWAS phenotype construction pipeline will often eliminate these primary cancer cases from being eligible controls for some selected secondary related phenotypes and thus a logistic regression analysis will lead to complete separation¹². We also stratified the MGI data set (or the corresponding gender subset depending on cancer type) into ten groups of equal size by PRS deciles and determined the percentage of observed cases for secondary traits in each risk decile and conducted a test of significance in difference in proportions across the deciles before and after removing individuals affected with cancer traits related to the primary cancer trait. As a follow-up tool to understand the secondary associations we created a plot to display the temporal ordering of diseases plotted against time of diagnoses. If not stated otherwise, analyses were performed using R 3.4.1³⁹

Results

In the current study, we report results obtained from 28,260 genotyped and unrelated samples of inferred European ancestry with available integrated ICD9-based EHR data. The study sample contains 53.5% females and the mean age is 54 years (see **Table 1 for summary**). We conducted our initial analysis on 12 cancer traits that after quality control had at least five independent risk variants in the NHGRI EBI GWAS Catalog and more than 250 cases in our cohort (**Table 2, Table S2**). **Table 2** summarizes data on 8,423 distinct individuals that were affected by at least one of the 12 cancers. Of these patients, 6,398 had one cancer, 1,574 had two cancers, and 451 had more than two cancer sites involved.

Correspondence of MGI effect estimates with those reported in GWAS: To assess the calibration properties of the 12 ICD-9 based cancer case-control studies, we first compared the concordance of observed effect estimates (log odds ratios) from MGI with published effect estimates reported in the NHGRI EBI GWAS Catalog.

We found strong positive correlation (estimated Pearson's correlation coefficient [$\hat{\rho}$] > 0.6) between the MGI and GWAS reported estimates for 9 of the 12 cancers: female breast cancer (78 SNPs; [$\hat{\rho}$]=0.67 [95% CI: 0.53,0.78]), prostate cancer (PCa; 93 SNPs; [$\hat{\rho}$]=0.81 [0.73,0.87]), melanoma (16 SNPs; [$\hat{\rho}$]=0.92 [0.77,0.97]), basal cell carcinoma (19 SNPs; [$\hat{\rho}$]=0.88 [0.71,0.95]), bladder cancer (MIM: 109800; 16 SNPs; [$\hat{\rho}$]=0.65 [0.22,0.86]), squamous cell carcinoma (5 SNPs; [$\hat{\rho}$]=0.95 [0.39,1]), lung cancer (MIM: 211980; 9 SNPs; [$\hat{\rho}$]=0.90 [0.6,0.98]), thyroid cancer (9 SNPs; [$\hat{\rho}$]=0.79 [0.26,0.95]), and cancer of brain and nervous system (9 SNPs; [$\hat{\rho}$]=0.79 [0.26,0.95]) (**Table 2; Table S5; Figure 1; Figure S3**).

Cancer GWAS in MGI: After having established strong positive correlation for 9 of the 12 cancer traits and thus a phenotype quality that appears to be in line with their corresponding published GWAS, we performed for each of these 9 cancers a GWAS to explore our ability to replicate and/or uncover cancer risk variants in a genome-wide setting. For the 9 cancers, we could replicate a total 55 of the 253 included risk SNPs with consistent effect orientation with $P < 0.05$ after correcting for the number of SNPs per phenotype (**Table S7**). We found genome-wide significant signals ($P < 5 \times 10^{-8}$) for female breast cancer, melanoma of skin, basal cell carcinoma, squamous cell carcinoma and thyroid cancer. All but one of the genome-wide signals were found in loci already reported in the GWAS catalog for the corresponding cancer trait or related phenotypes. For instance, the four melanoma of skin loci with risk variants near *SLC45A2* (MIM 606202), *IRF4* (MIM 601900), *MC1R* (MIM 155555) and *ASIP/RALY* (MIM: 600201) were previously reported to be associated with melanoma, non-melanoma skin cancer, squamous cell carcinoma, or basal cell carcinoma⁴⁰⁻⁴⁴. Also, the two breast cancer risk loci near *FGFR2* (MIM: 176943) and *FGF3/FGF4* (MIM: 164950 / 164980) as well as the thyroid cancer risk loci near *NRG1* (MIM: 142445) and *FOXE1* (MIM: 602617) were previously described⁴⁵⁻⁴⁸. The only potentially novel finding was the SNP rs77909434 on chromosome 13 showing borderline genome-wide association with melanoma (MAF in cases = 5.3%; MAF in controls = 3.4%; $P = 1.5 \times 10^{-8}$) located 53 kb downstream of the Fibroblast Growth Factor 9 gene (*FGF9*, MIM: 600921) on chromosome 13. Since multiple phenotypes were involved in the genome-wide explorations, this SNP would not have passed the Bonferroni multiple testing correction for multiple GWAS. Further

exploration in larger studies are warranted to substantiate this suggestive finding. We present GWAS Manhattan and QQ plots for all nine cancer traits in **Figure S4**.

Owing to the smaller sample sizes compared to the studies included in the NHGRI-EBI GWAS Catalog, only 8 of 253 catalog SNPs exceeded the genome-wide significance (**Table S5**). However, we found catalogued risk SNPs in **Table 2** were markedly enriched in the top 1% of GWAS associations, especially for the larger case/control studies. For example, 27 out of the 93 GWAS Catalog PCa risk SNPs fall in top 1% of associated SNPs in the MGI GWAS (with $P < 0.0083$) (**Table S8**).

Replicability of PRS Primary Cancer association: PRS integrates multiple SNPs, weighted by prior effect estimates and is expected to substantially improve the power to detect an association compared to an analysis with individual variants. To evaluate the association of PRS with the primary cancer trait, we estimated the OR for patients in the top risk quartile compared to the bottom quartile (PRS OR). Six of the 12 cancer PRS revealed an at least twofold enrichment of cases with $P_{Q1vsQ4} < 2.0 \times 10^{-10}$; all of which also showed strong positive correlation between the MGI and GWAS reported estimates (see above): female breast cancer (PRS OR = 2.3 [95% CI: 2.0;2.7], $P_{Q1vsQ4} = 2.5 \times 10^{-29}$), prostate cancer (PRS OR = 3.3 [95% CI: 2.7;3.9], $P_{Q1vsQ4} = 3.7 \times 10^{-43}$), melanoma (PRS OR = 2.4 [95% CI: 2.0;2.8], $P_{Q1vsQ4} = 2.6 \times 10^{-31}$), basal cell carcinoma (PRS OR = 2.7 [95% CI: 2.2;3.2], $P_{Q1vsQ4} = 1.1 \times 10^{-27}$), squamous cell carcinoma (PRS OR = 2.0 [95% CI: 1.6;2.5], $P_{Q1vsQ4} = 2.0 \times 10^{-10}$), and thyroid cancer (PRS OR = 3.2 [95% CI: 2.5;4.5], $P_{Q1vsQ4} = 1.8 \times 10^{-18}$) (**Figure 1A-F, Table 2**).

The corresponding P-values obtained from Firth's bias-reduced logistic regression using continuous PRS were even stronger as expected and indicated that these six cancer traits

would withstand a Bonferroni multiple testing correction in a phenome setting (1,711 traits; $P_{PRS} < 2.9 \times 10^{-5}$): female breast cancer ($P_{PRS}=3.6 \times 10^{-37}$), prostate cancer ($P_{PRS}=3.8 \times 10^{-69}$), melanoma ($P_{PRS}=6.7 \times 10^{-36}$), basal cell carcinoma ($P_{PRS}=3.3 \times 10^{-44}$), cutaneous squamous cell carcinoma (SCC, $P_{PRS}=1.8 \times 10^{-18}$), thyroid cancer (MIM: 188550) $P_{PRS}=4.8 \times 10^{-19}$) (**Table 2, Table S3**). We excluded the remaining six cancer traits from further investigation, because in this initial analysis they showed only little or moderate association (PRS OR < 1.5) and consequently only modest power for subsequent exploration of phenome-wide associations (**Table 2**).

PRS PheWAS: Next, we evaluated each of the six remaining PRS that were strongly associated with the primary cancer trait across a collection of 1,711 EHR-derived phenotypes (not limited to cancer traits) with at least 20 cases each (**Table S1**). For each of the six cancer PRS, we found strongest associations with their primary traits, except for squamous cell carcinoma PRS which revealed its strongest association with the more general skin cancer trait definition ($P_{PRS} = 7.2 \times 10^{-61}$) (**Figure 2**). Overall, we found no or little sign for inflation in our PheWAS results (median chi-squared based Lambda ≤ 1.16). Notably, we observed deflation for some PRS PheWAS that might be caused by lack of power especially for the phenotypes with small number of cases (**Figure S5**). We displayed the results from the three type of sensitivity analyses PheWAS next to the original results: conditional logistic regression results from 1:10 matched (for age, sex and first four PCs) case-control studies; adjusting for 20 principal components; using unweighted sum of risk allele counts instead of weighted PRS (**Figures S6-11**). The results remained robust with respect to these design and analytic choices.

Secondary Associations: In addition, we identified for each PRS novel associations with secondary traits besides their primary traits (**Figure 2 A-F, Table S9**). For example, we observed associations of the three skin cancer PRSs (PRS for melanoma, basal cell carcinoma, and squamous cell carcinoma) with overall skin cancer and other skin cancer sub categories – expected due to their overlapping SNP sets – but also significantly associated with multiple dermatologic phenotypes, e.g. actinic keratosis ($P_{PRS} < 1.2 \times 10^{-10}$) and other degenerative skin conditions or disorders, all potential pre-cancer stages (**Figure 2C,D,E**).

Similarly, the female breast cancer PRS was associated not only with breast cancer ($P_{PRS} = 3.6 \times 10^{-37}$) but also with acquired absence of breast ($P_{PRS} = 2.4 \times 10^{-14}$), abnormal mammogram ($P_{PRS} = 1.3 \times 10^{-8}$), benign neoplasms of the breast ($P_{PRS} = 1.8 \times 10^{-7}$) and benign mammary dysplasias ($P_{PRS} = 1.2 \times 10^{-5}$) (**Figure 2A**). The PRS originally constructed for prostate cancer was associated with prostate cancer ($P_{PRS} = 3.8 \times 10^{-69}$), as expected, but also with four additional traits: elevated prostate specific antigen ($P_{PRS} = 9.3 \times 10^{-27}$), erectile dysfunction ($P_{PRS} = 6.3 \times 10^{-15}$), urinary incontinence ($P_{PRS} = 6.6 \times 10^{-11}$), frequency of urination and polyuria ($P_{PRS} = 2.9 \times 10^{-6}$), and hyperplasia of prostate ($P_{PRS} = 3.6 \times 10^{-6}$) (**Figure 2B, Table S9**).

While all of the above mentioned secondary trait associations were in the same effect orientation as their primary traits, i.e. increasing PRSs were associated with increased risk for the secondary trait, we observed an association of increasing thyroid cancer PRS with decreased risk for hypothyroidism ($P_{PRS} = 7.0 \times 10^{-10}$) (**Figure 2F**).

Exploring Secondary PRS PheWAS Associations via Exclusion PRS PheWAS:

Since we already applied exclusion criteria to the controls during our phenome

generation, e.g., individuals with elevated prostate specific antigen levels were excluded from being controls for prostate cancer and vice versa, we could not adjust for the primary cancer traits as a predictor in logistic regression models to identify independent secondary PRS PheWAS associations due to the issue of complete separation. To alternatively explore the secondary associations in PRS PheWAS (**Figure 2**), we proposed and performed “exclusion PRS PheWAS” by removing subjects affected with the cancer or related cancer traits for which the PRS was constructed. After removing all breast cancer cases (N = 1,894) no association with breast cancer PRS remained significant, e.g., acquired absence of breast ($P_{PRS} = 0.52$), abnormal mammogram ($P_{PRS} = 0.76$) or benign neoplasms of the breast ($P_{PRS} = 0.49$), indicating that the secondary trait associations were driven by the primary trait (**Figure S12A**). However, we noted that the majority of cases of the non-neoplasm phenotype “Acquired absence of breast” (>94.4%; 624 of 661) were removed in this step as they are highly correlated with breast cancer. We made similar observations for prostate cancer PRS where none of the previously detected secondary trait associations remained phenome-wide significant after removing all 1,425 prostate cancer cases (**Figure S12B**).

In contrast, we found a markedly stronger association between hypothyroidism and thyroid cancer PRS after removing 472 thyroid cancer cases ($P_{PRS} = 4.7 \times 10^{-19}$) compared to the full analysis ($P_{PRS} = 7.0 \times 10^{-10}$) which is consistent with the effect orientations between thyroid cancer PRS and hypothyroidism (**Figure 2F**).

To account for the substantial overlap between skin cancer sub types, e.g. 253 of the 1,404 individuals affected with melanoma are also affected by basal and/or squamous cell carcinoma (**Figure S13**) and to account for the likely intensified skin cancer screening

of individuals that were diagnosed with skin cancer once in their life time, we excluded any type of skin cancer (N = 3,910) and repeated the PheWAS for melanoma, basal cell carcinoma, and squamous cell carcinoma PRS. After doing so, only actinic keratosis remained statistically associated with squamous cell carcinoma PRS while all of the previously observed associations mainly driven by skin cancer diagnoses disappeared (**Figure 2C-E** and **Figure S12C-E**). The association between squamous cell carcinoma PRS and actinic keratosis was less pronounced after excluding skin cancer cases but still remained phenome-wide significant ($P_{PRS} = 2.3 \times 10^{-36}$ versus $P_{PRS} = 1.1 \times 10^{-12}$).

To further understand the discovered secondary associations in the PRS PheWAS analyses (**Figure 2** and **Figure S12**), we conducted a simple follow-up analysis by stratifying the data into PRS deciles. We only discuss selected secondary trait associations for the prostate cancer (PCa), squamous cell carcinoma (SCC) and thyroid example in the main text and relegate their comprehensive analysis and a similar analysis of breast cancer, melanoma and basal cell carcinoma PRS to the supplemental material (**Table S10**). For prostate cancer, we stratified a total of 12,026 male individuals in MGI with age ≥ 30 years into deciles of PCa PRS. The observed PCa PRS associations in the PheWAS analysis are further supported by their respective increasing trait prevalences that are observed across 10 PCa PRS decile-stratified strata (**Figure 3A**; **Table S10**). These strata were not adjusted for confounders, but it is less likely that PRS is strongly associated with other covariates. A striking observation is that the proportion of PCa cases in lowest versus highest decile of PCa PRS is 5.4% versus 23.4% ($\Delta=18.0\%$ [95% CI, 15.2 to 20.7%]; $P=8.3 \times 10^{-36}$) emphasizing that the PRS can distinguish well between high and low risk individuals in a realistic academic medical center population.

Focusing on the secondary traits that reached genome-wide significance with the PCa PRS, all these traits are known to be associated with PCa: erectile dysfunction (ED), and urinary incontinence (UI) – which commonly follows invasive surgical removal of the prostate – and elevated prostate specific antigen levels (ePSA) – which is a known biomarker for an increased PCa risk being closely monitored after prostatectomy. For example, when comparing the lowest versus the highest PRS risk decile, we found significant differences for ePSA (3.9% versus 11.2%; $\Delta=7.3\%$ [95% CI, 5.1 to 9.5%]; $P=2.0\times 10^{-36}$), ED (9.7% versus 17.1%; $\Delta=7.4\%$ [95% CI, 4.6 to 10.2%]; $P=1.4\times 10^{-7}$), and UI (4.7% versus 11.9%; $\Delta=7.2\%$ [95% CI, 5.0 to 9.5%]; $P=2.0\times 10^{-10}$). To test whether these associations are early indicator for PCa or whether they are driven by the fact that subjects affected by these secondary traits are also PCa cases (perhaps as a side effect of PCa treatment), we removed PCa cases and evaluated secondary disease prevalence across PCa PRS deciles. By doing so, prevalence of all secondary traits became roughly constant across PRS strata (**Figure 3A; Table S10**) and can be illustrated by the comparison of the proportions of the lowest versus the highest PRS risk decile: ePSA (2.8% versus 2.2%, $\Delta=-0.6\%$ [95% CI, -1.9 to 0.8%]; $P=0.44$), ED (7.6% versus 6.5%, $\Delta=-1.1\%$ [95% CI, -3.2 to 1.0%]; $P=0.34$), and UI (3.0% versus 2.8%, $\Delta=0.2\%$ [95% CI, -1.6 to 1.3%]; $P=0.90$). Based on these observations, we hypothesize that the association of PCa PRS on the secondary traits ePSA, ED, and UI were driven by the PCa diagnosis, through either prior symptoms of PCa or prescribed medication, chemotherapy or surgical procedures for prostate removal (**Table S10**).

For SCC PRS stratification, there was a gradual increase of individuals affected with SCC with increasing PRS risk deciles, a trend that was also noted for actinic keratosis,

dermatitis due to solar radiation, and seborrheic keratosis (**Figure 3B**). However, when excluding cases with skin cancer, the upward trend for the latter two phenotypes disappeared. The previously observed difference between the top and bottom PRS risk decile of individuals affected with actinic keratosis (5.5% versus 13.2% ($\Delta=7.7\%$ [95% CI, 6.1 to 9.3%]; $P=4.0 \times 10^{-21}$) was markedly reduced after excluding skin cancer cases but still remained significant (2.8% versus 5.1% ($\Delta=2.4\%$ [95% CI, 1.3 to 3.5%]; $P=1.6 \times 10^{-5}$) (**Table S10**) suggesting the potential for common genetic risk profiles between SCC and actinic keratosis. Since actinic keratosis is a known precursor for squamous cell carcinoma⁴⁹, our approach indicated that it is possible to identify phenotypic risk factors through genome-wide association scans and careful follow-up investigation of primary and secondary diagnoses.

Finally, we found an attenuated association between increasing thyroid cancer PRS and reduced risk for hypothyroidism: within all 25,681 samples ≥ 30 years of age the difference between bottom and top decile was $\Delta=-3.5\%$ ([95% CI, -5.4 to -1.6%]; 15.1% versus 11.5%; $P=2.5 \times 10^{-4}$) and after excluding 452 thyroid cancer cases it increased to $\Delta=-5.3\%$ ([95% CI, -7.1 to -3.5%]; 14.4% versus 9.1%; $P=4.5 \times 10^{-9}$) (**Table S10**). Several studies previously reported genetic overlap of a subset of thyroid cancer risk variants and variants associated with serum levels of thyroid stimulating hormone (TSH) which matches the current observed association between thyroid cancer risk and risk for hypothyroidism^{47; 48}.

To further our understanding of the observed secondary associations, we take advantage of the temporally resolved electronic health records data and explore the temporal order in which the diagnoses appear. **Figure 4** shows that actinic keratosis diagnosis mostly

precedes the diagnosis of squamous cell carcinoma and sometimes by even 10 years. Erectile dysfunction or hypothyroidism, known side-effects of treatment of prostate and thyroid cancer (respectively), are mostly identified within a short timeframe of primary cancer diagnosis. Whereas elevated PSA, used as a screening tool for prostate cancer with known shared genetic correlation is observed mostly prior to a prostate cancer diagnosis and also after treatment as a prognostic marker. Having access to the electronic health records enables us to explore these temporally ordered data patterns and understand the explanation behind these secondary associations.

Discussion

Integration of large-scale biorepositories such as genetic data with EHRs are becoming increasingly common and indispensable for next-generation etiology studies. In this paper, we proposed, demonstrated and tested trait-specific PRS that summarize the results of large population-based GWAS studies towards cancer risk prediction in an actual academic medical center population managed by Michigan Medicine. Data repositories like MGI, allow us to explore many traits simultaneously whereas population-based case-control studies focus on one specific trait. It is indeed encouraging that the results of population-based studies corroborate with the phenotypes computed from EHR data. We found improved trait prediction power of the composite PRS compared to single-SNP analyses. We also replicated catalogued associations of SNPs for some cancer traits, observed excellent correspondence of effect estimates and discovered novel secondary trait associations with cancer PRS that were not driven by the primary cancer

diagnoses. To our knowledge this is the first comprehensive PRS PheWAS study and the first PheWAS study focused on cancer.

We have introduced several novel analytic strategies in this paper. We presented a principled framework and quality-control pipeline to create a PRS from a large curated, public database and to perform PRS PheWAS in a potentially biased sample. We introduced a primary PheWAS using Firth's bias-reduced logistic regression which has the advantage of resolving the problem of separation in logistic regression and providing well-controlled type I error rates for unbalanced case control studies with relatively small sample counts^{31; 32; 50}. These issues are often present in large EHR-based phenomes where controls are frequently hundredfold more abundant than cases. In addition, we conducted thorough sensitivity analyses to check the robustness of our findings by using PheWAS with unweighted risk allele counts, adjusting for 20 PCs and PRS PheWAS based on matched controls. All our reported results remained robust under these sensitivity analyses.

To distinguish PRS-trait associations that truly derive from a shared genetic risk profile from secondary associations that are potentially driven by the primary trait (for example urinary incontinence or erectile dysfunction following prostate cancer treatment), we further introduced a modified PRS PheWAS approach that excludes the PRS's underlying cancer traits. While reducing overall sample size, this "exclusion PRS PheWAS" approach is statistically preferable in contrast to a PRS PheWAS that conditions on the primary cancer trait. A conditional PheWAS approach is often affected by unilaterally applied exclusion criteria of controls that occur during the phenome construction, e.g. PCa cases were excluded from being eligible controls for elevated PSA levels and vice versa. Our

approach could directly discard trait associations driven by the primary cancer diagnosis and has the potential to identify clinically useful diagnostic traits among many that are conveniently measured in panel tests of biomarkers. When an association with a secondary trait disappears by removing the primary cancer cases in an exclusion PheWAS, there can be several alternative explanations: truly shared genetic correlation, intensified screening/examination due to detection of an initial cancer, a screening biomarker/pre-cancer phenotype or simply post treatment effects. We used the temporal ordering of the diagnoses to understand which of the above explanations appear plausible for a given scenario. Further exploration of our findings in larger biobank studies, like the UK Biobank study, is warranted and will empower a deeper understanding of relevant pre-cancer traits ⁵¹.

There are several limitations to the current study. We decided to rely on the associations reported in the NHGRI-EBI GWAS Catalog instead of focusing on the latest and largest GWAS study specific for each cancer trait. Our rationale for choosing the NHGRI-EBI GWAS Catalog as our source for extracting summary statistics were primarily three-fold: (1) Data quality: Summary statistics in the GWAS Catalog underwent a detailed expert curation and harmonization ^{28; 52} that avoids redundancy, allows reliable SNP position extraction, and most importantly ancestry matching; We wanted to use a database that is publicly accessible and applies the same set of criteria to update reported results across a wide variety of phenotypes. (2) Reproducibility: We provided detailed instructions on how to extract and filter GWAS Catalog summary statistics to construct PRS. This will allow interested readers to easily apply our approach to the regularly updated GWAS catalog versions or to a different ancestry group and/or broad set of disease categories

without requiring detailed and deep literature searches that could be somewhat subjective. (3) Scalability to Multiple Phenotypes: One can construct PRS for specific cancers of primary interest from the latest GWAS meta-analyses following the same prescriptions we provided. Using the latest GWAS result is likely to enhance power of a PRS PheWAS. Similarly using a PRS that is based on a truly polygenic model with many more variant (or the entire genome) instead of considering the GWAS hits may reveal new associations.

We restricted our analysis to GWAS results from studies of broad European ancestry to match them to our cohort of predominantly European ancestry and to allow an extra filtering of potential swaps in directionality of risk allele in published GWAS studies that otherwise could have negatively affected the correlative properties of our constructed PRS. One could modify or extend construction of PRS based on global ancestry, functionality of the variant and use other weighting schemes. Stratifying the present analysis by young onset cancers, metastatic/aggressive cancers or tumor subtype will shed further insight into cancer biology, cancer genetics and specificity of the PRS-cancer association. We have mostly ignored the temporal ordering in the diagnoses codes by defining dichotomous phenotypes of interest. Exploring the time-stamped data in greater detail may be instrumental in understanding the secondary associations like the negative association between hypothyroidism and thyroid cancer PRS.

Though we note some very encouraging and promising results for the cancer traits with modest number of cases and controls and with a larger number of variants reported in the NHGRI-EBI GWAS catalog, we also note that the correlation of effect estimates or the PRS-cancer association was not very strong for some cancers (**Table 2**). This could

be due to limited sample size/power, heterogeneity in the definition of the cancer phenotype, incomplete specification of PRS, differences in allele frequencies in the MGI population, or misclassification of ICD-9 codes. To address concerns with misclassification we conducted detailed chart review of 50 randomly sampled cases with at least one cancer PheWAS code and verified their primary and secondary cancer diagnosis. We could verify 149 of the 151 diagnoses and found 49/50 patients to have accurate record of their cancer diagnosis. Based on this we conclude that the rate of misclassification will likely be low for ICD 9 codes associated with cancer.

In this paper, we have focused on cancer traits. The low misclassification rate of cancer traits, typical within academic health and cancer centers, along with effective sensitivity analyses partly protect the results against imprecise case definitions and confounding. For non-cancer disease traits, more stringent ICD-9 defined cases, e.g., by repeated ICD-9 diagnoses, of adequate sample sizes might alleviate the biases from case misclassification. Future analysis will need to control for potentially different levels of misclassification error across phenotypes.

Our phenome comprised a total of 1,711 ICD9-based phenotypes and by its implemented design of hierarchical phenotypes with different levels of specificity induce a certain degree of redundancy. While we applied the multiple testing correction for 1,711 performed tests, we acknowledge that this threshold might be too conservative. For examples, we estimated a maximal set of 1,452 phenotypes with all pair-wise correlations $r^2 < 0.5$ before applying any exclusion criteria to the controls. In addition, the PheWAS approach often applies similar exclusion criteria to related phenotypes and thereby further

reduces the observable independence of case-control studies. Future studies are needed to determine the effective number of independent tests in such a phenome-wide analysis. Besides the ICD-9 codes used for case and control definitions, EHR databases generally contain vast amount of additional patient information including ICD-10 codes, temporal laboratory tests, drug prescriptions, inpatient and outpatient records, etc. Future analyses that leverage these heterogeneous data sources that might be predictive of disease outcomes could further improve disease risk predictions. It will be interesting to study whether PRS for cancer risk behaviors like smoking, alcohol and obesity predict cancer phenotypes. Tailored and validated models capable of integrating multiple sources of molecular and environmental data data for predicting risks of disease will be crucial.

Supplemental Data

Supplemental Data include 13 figures and 10 tables.

Conflicts of Interest

The authors declare no competing financial interest.

Acknowledgments

The authors acknowledge the University of Michigan Medical School Central Biorepository for providing biospecimen storage, management, and distribution services in support of the research reported in this publication.

Web Resources

Michigan Genomics Initiative, <https://www.michigangenomics.org>

OMIM, <http://www.omim.org>

BAF Regress, <http://genome.sph.umich.edu/wiki/BAFRegress>

PLINK v1.90, <https://www.cog-genomics.org/plink2>

KING v2.0, <http://people.virginia.edu/~wc9c/KING>

Human Genome Diversity Project reference panel,
<http://csg.sph.umich.edu/chaolong/LASER>

PheWAS R package, <https://github.com/PheWAS/PheWAS>

NHGRI-EBI GWAS Catalog, <https://www.ebi.ac.uk/gwas>

UCSC genome browser, <http://genome.ucsc.edu>

EPACTS, <http://genome.sph.umich.edu/wiki>

Figure Titles and Legends

Figure 1. Calibration of association parameters

Calibration of association parameters between the MGI-GWAS and NHGRI-EBI GWAS Catalog derived effect estimates [$\log(\text{OR})$] for (A) breast cancer (females only), (B) cancer of prostate, (C) melanoma, (D) basal cell carcinoma, (E) squamous cell carcinoma, and (F) thyroid cancer. The agreement of two sets of SNP-specific beta coefficients (non-reference allele is the effect allele), their Pearson Correlation (Coefficient $\hat{\rho}$, incl. 95% confidence interval and P) and Lin's correspondence correlation (coefficient CCC; incl. 95% confidence interval) are shown; dashed line: perfect concordance; solid line: fitted line.

Figure 2. PRS PheWAS plots

PRS PheWAS plots for (A) breast cancer (females only), (B) cancer of prostate, (C) melanoma, (D) basal cell carcinoma, (E) squamous cell carcinoma, and (F) thyroid cancer. 1,711 traits are grouped into 16 color-coded categories as shown on the horizontal axis; the p-values for testing the associations of PRS with the traits are minus log-base-10-transformed and shown on the vertical axis. Triangles indicate phenome-wide significant associations with their effect orientation (up-pointing = risk increasing; down-pointing = risk decreasing). PRS upon multiplicity adjustment (see **Methods**). The solid horizontal line for $P=2.9 \times 10^{-5}$ cut-off.

Figure 3. Proportion of primary and secondary traits stratified by PRS deciles

Percentage of primary and selected secondary traits in each cancer PRS category for (A) prostate cancer, (B) squamous cell carcinomas, and (C) thyroid cancer. Observed percentages in the MGI database as represented by the height of bars for each of 10 increasing decile-stratified PRS strata from left to right. The PRS's underlying trait is shown on top and secondary traits below with (blue) and without (green) overlapping relevant cancer cases. Only individuals with age \geq 30 years were included in each analysis, with the prostate cancer PRS example only includes male individuals (see **Table S10** for detailed sample sizes and percentages).

Figure 4. Temporal order of diagnoses: (A) elevated PSA levels (ePSA) and PCa in 452 individuals with PCa and ePSA; (B) erectile dysfunction (ED) and prostate cancer (PCa) in 575 individuals with ED and PCa; (C) actinic keratosis (AK) and squamous cell carcinoma (SCC) in 286 individuals with AK and SCC; and (D) hypothyroidism (HT) and thyroid cancer (TCa) in 298 individuals with HT and TC. The time of the first non-cancer diagnosis relative to the cancer diagnosis is shown in weeks; before (blue) and after (red) the cancer diagnosis.

Tables

Table 1 Demographics and clinical characteristics of the final analytic data set

Characteristic	Analytic Data Set
N	28,260
Females, N (%)	15,113 (53.5%)
Mean Age, years (S.D.)	54.1 (15.9)
Total number of ICD9 code days	3.5 million
Number of unique ICD9 codes	10,322
Median number of visits per participant	23
Median days between first and last visit	1,265
Median ICD9 code days per participant	28

Table 2. Association analysis of cancer traits with at least five NHGRI EBI GWAS Catalog risk SNPs and more than 250 cases in MGI.

Cancer Trait ^a	N Cases	N Controls	N Risk SNPs used for PRS ^b	Effect Size Correlation [$\hat{\rho}$] between GWAS Catalog and MGI [95% CI]	Effect Size Correspondence (Lin's CCC) MGI versus GWAS Catalog [95% CI]	Estimated PRS-Cancer Association	
						PRS Odds Ratio [95% CI] ^c	Continuous PRS (normalized by IQR) Point Estimate [95% CI] ^d
<i>Breast cancer [female]</i> ^e	1,827	11,073	78	0.67 [0.53,0.78]	0.64 [0.51, 0.75]	2.3 [2.0,2.7] 2.5×10^{-29}	1.6 [1.5,1.7] 3.6×10^{-37}
<i>Cancer of prostate</i> ^e	1,425	9,793	93	0.81 [0.73,0.87]	0.74 [0.66, 0.81]	3.3 [2.7,3.9] 3.7×10^{-43}	2.0 [1.9,2.2] 3.8×10^{-69}
<i>Melanomas of skin</i> ^e	1,404	23,798	16	0.92 [0.77,0.97]	0.91 [0.77, 0.97]	2.4 [2.0,2.8] 2.6×10^{-31}	1.6 [1.5,1.7] 6.7×10^{-36}
<i>Basal cell carcinoma</i> ^e	1,124	23,798	19	0.88 [0.71,0.95]	0.85 [0.68, 0.94]	2.7 [2.2,3.2] 1.1×10^{-27}	1.8 [1.6,1.9] 3.3×10^{-44}
Cancer of bladder	978	26,748	16	0.65 [0.22,0.86]	0.57 [0.22, 0.79]	1.4 [1.2,1.7] 0.00018	1.2 [1.1,1.3] 4.9×10^{-6}
Non-Hodgkins lymphoma	878	26,794	18	0.51 [0.05, 0.79]	0.24 [0.028,0.43]	1.3 [1.1,1.6] 0.0063	1.1 [1.0,1.3] 0.0029
Colorectal cancer	718	22,183	42	0.48 [0.21,0.68]	0.39 [0.17,0.58]	1.3 [1.1,1.6] 0.011	1.2 [1.1,1.3] 0.00078
<i>Squamous cell carcinoma</i> ^e	703	23,798	5	0.95 [0.39,1.00]	0.92 [0.57, 0.99]	2.0 [1.6,2.5] 2.0×10^{-10}	1.6 [1.5,1.8] 1.8×10^{-18}
Malignant neoplasm of kidney, except pelvis	613	26,748	7	0.33 [-0.57,87]	0.053 [-0.094, 0.20]	0.98 [0.77,1.3] 0.89	0.99 [0.89,1.1] 0.86
Cancer of bronchus, lung	570	27,596	9	0.90 [0.60,0.98]	0.82 [0.53, 0.93]	1.2 [0.91,1.6] 0.13	1.1 [0.99,1.2] 0.091
<i>Thyroid cancer</i> ^e	472	26,692	8	0.97 [0.82,0.99]	0.94 [0.82, 0.98]	3.2 [2.5,4.5] 1.8×10^{-18}	1.7 [1.4,1.9] 4.8×10^{-19}
Cancer of brain and nervous system	321	27,069	9	0.79 [0.26,0.95]	0.66 [0.25,0.87]	1.3 [0.92,1.7] 0.13	1.2 [1.1,1.4] 0.042

^a underlying ICD9 codes are listed in **Table S4**

^b GWAS Catalog SNPs after quality control; corresponding summary statistics are listed in **Table S2, and Table S5**

^c Odds ratio for each cancer with top PRS quartile compared to bottom PRS quartile. Point estimates, confidence intervals and P- values are obtained by fitting Firth's Bias-Corrected Logistic Regression.

^d Association of each cancer with continuous PRS that were normalized by their interquartile ranges. Point estimates, confidence intervals and P- values are obtained by fitting Firth's Bias-Corrected Logistic Regression.

^e $OR > 1.5$, $P_{\text{continuous PRS}} < 2.9 \times 10^{-5}$, selected for PRS-PheWAS analysis

Notes: PRS = polygenic risk score, CI = confidence interval, IQR = interquartile range.

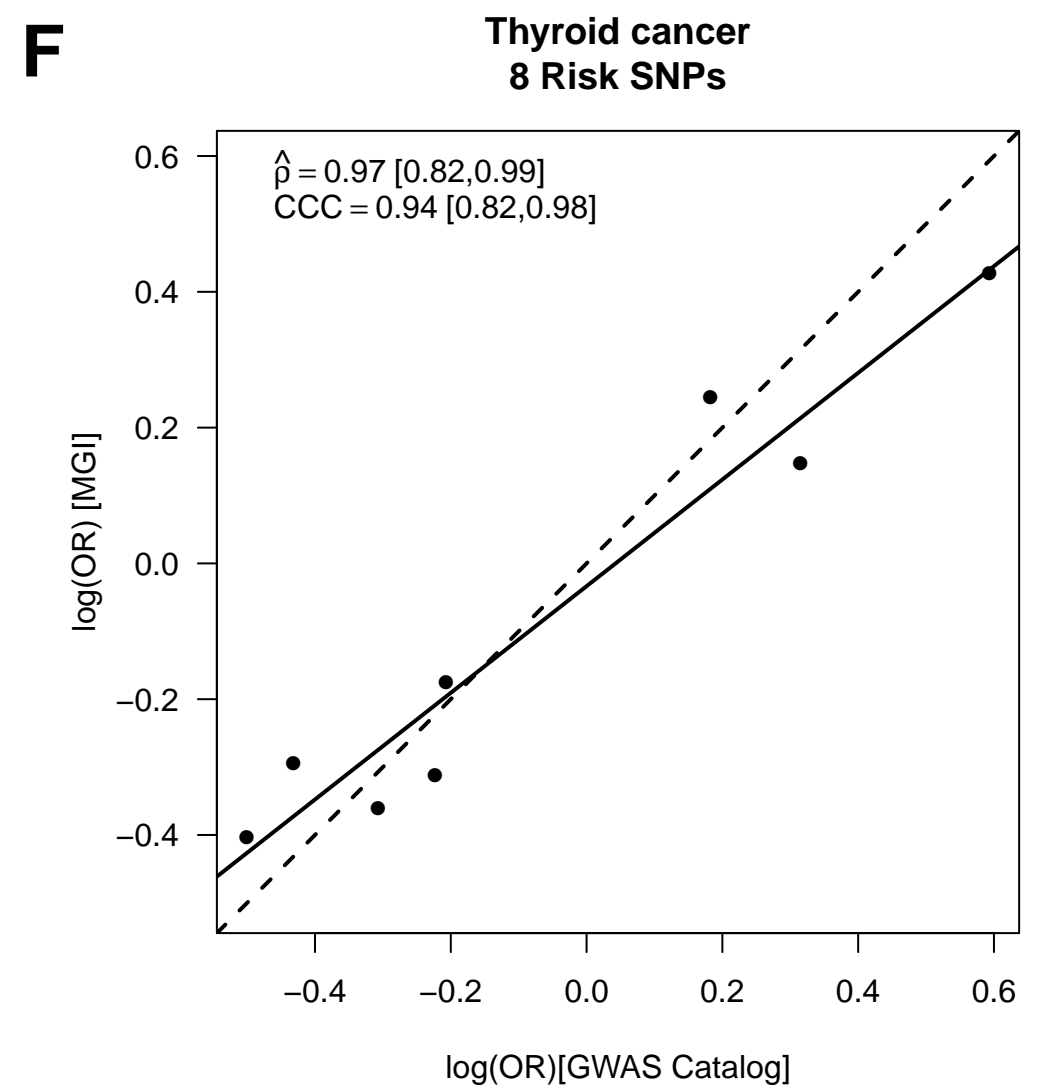
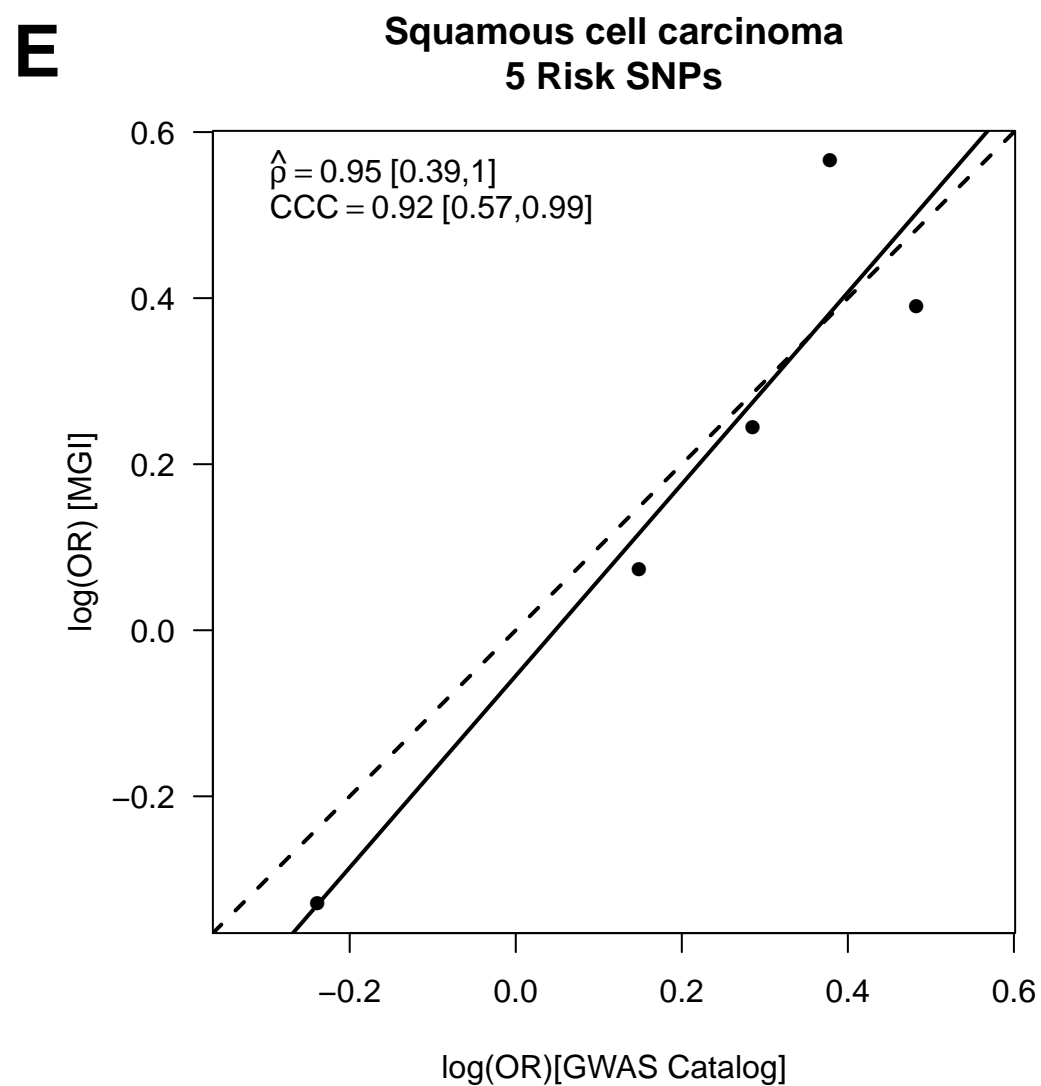
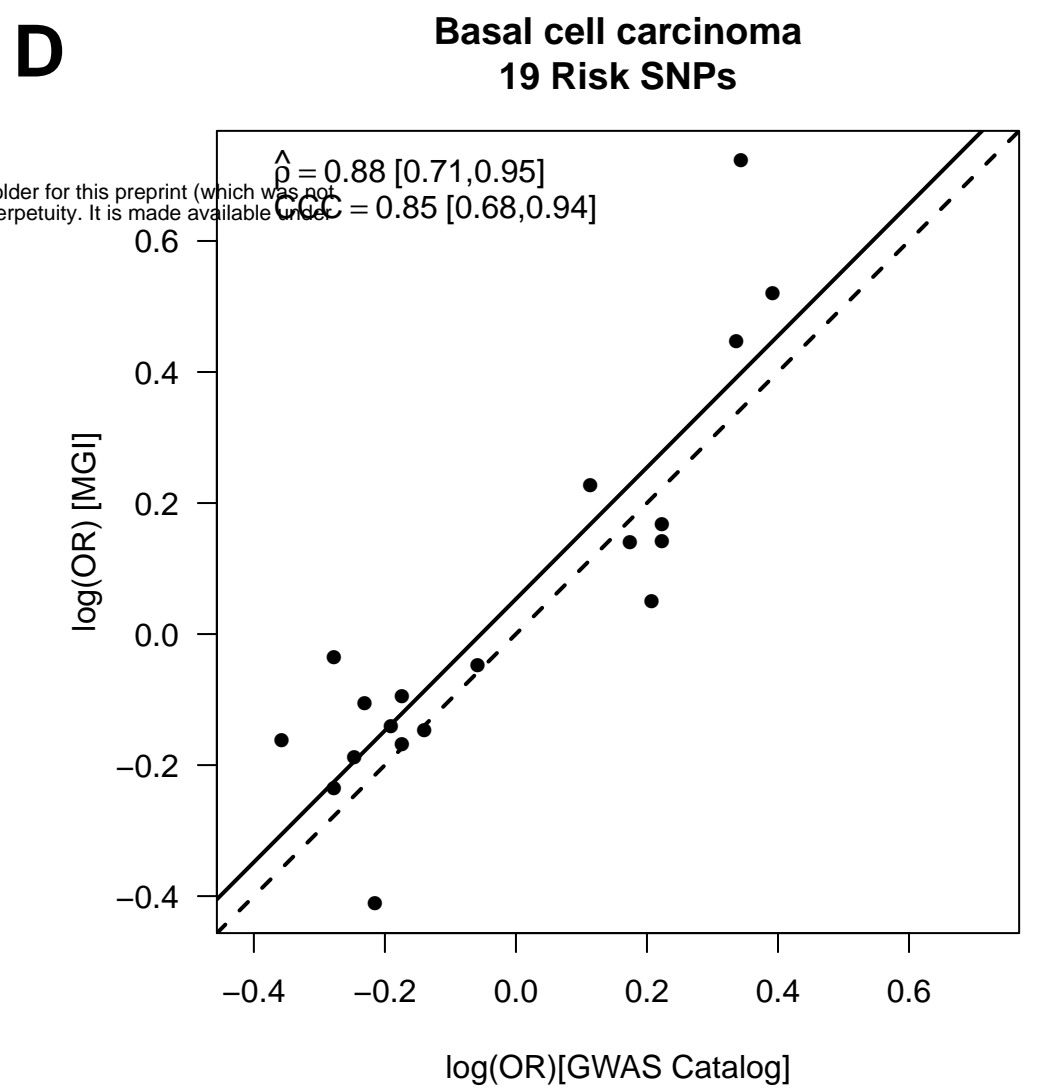
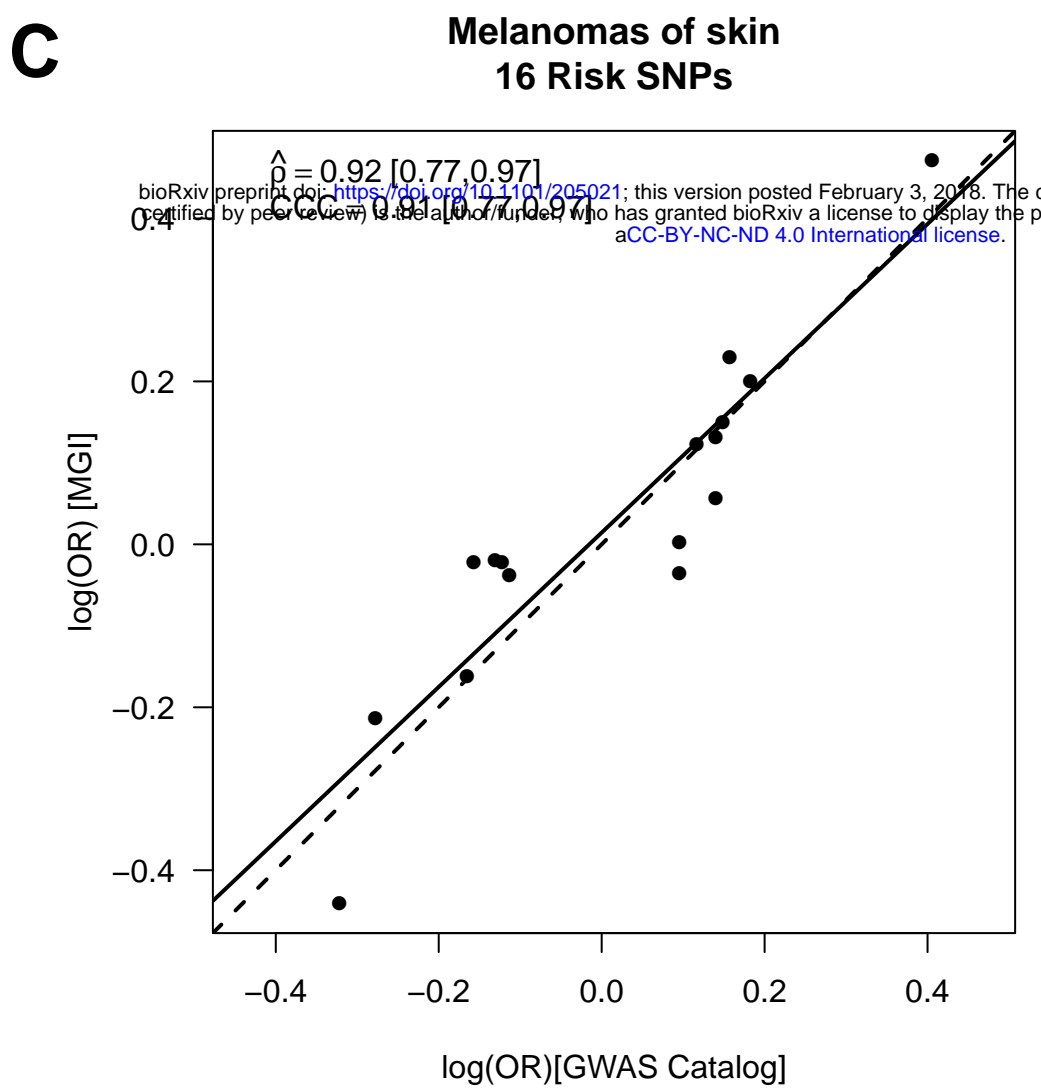
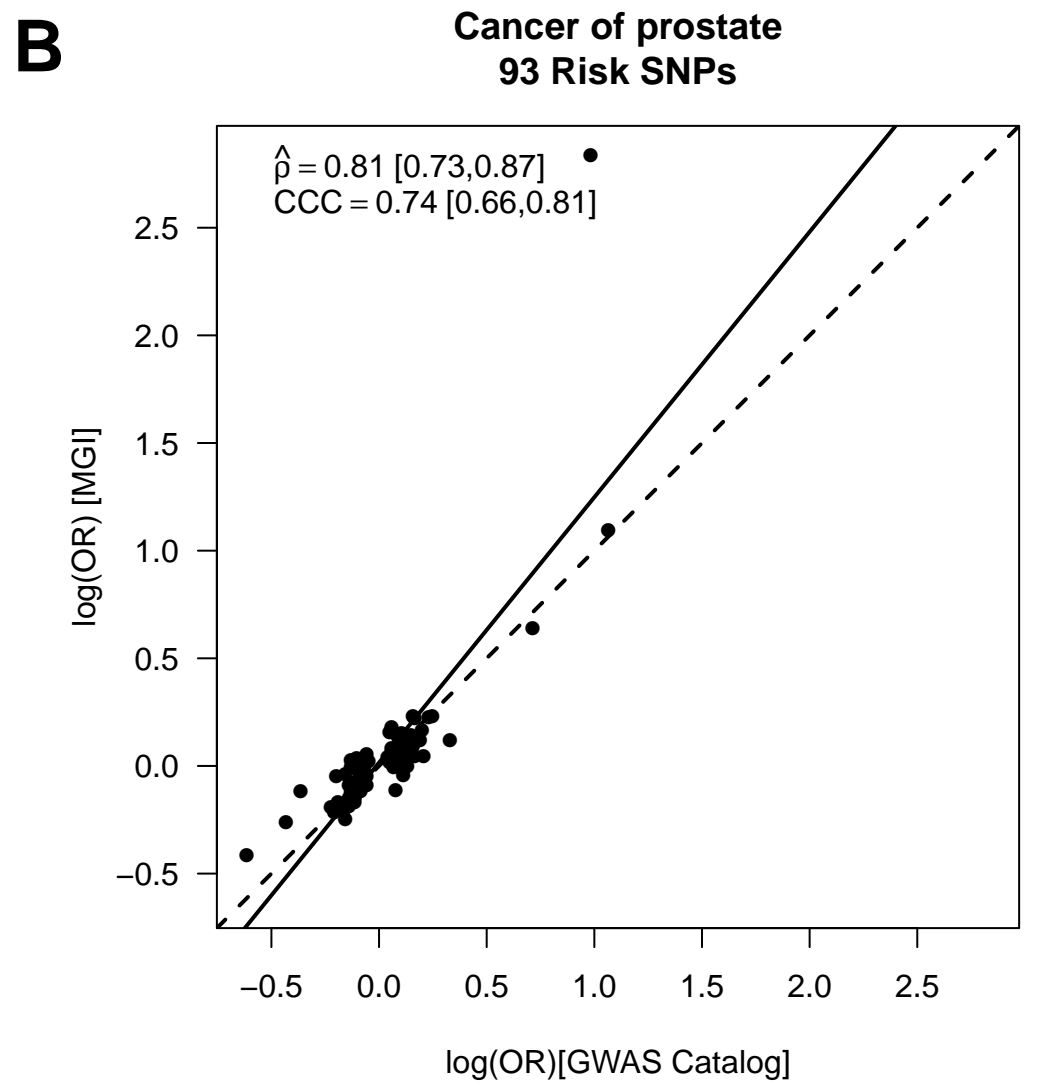
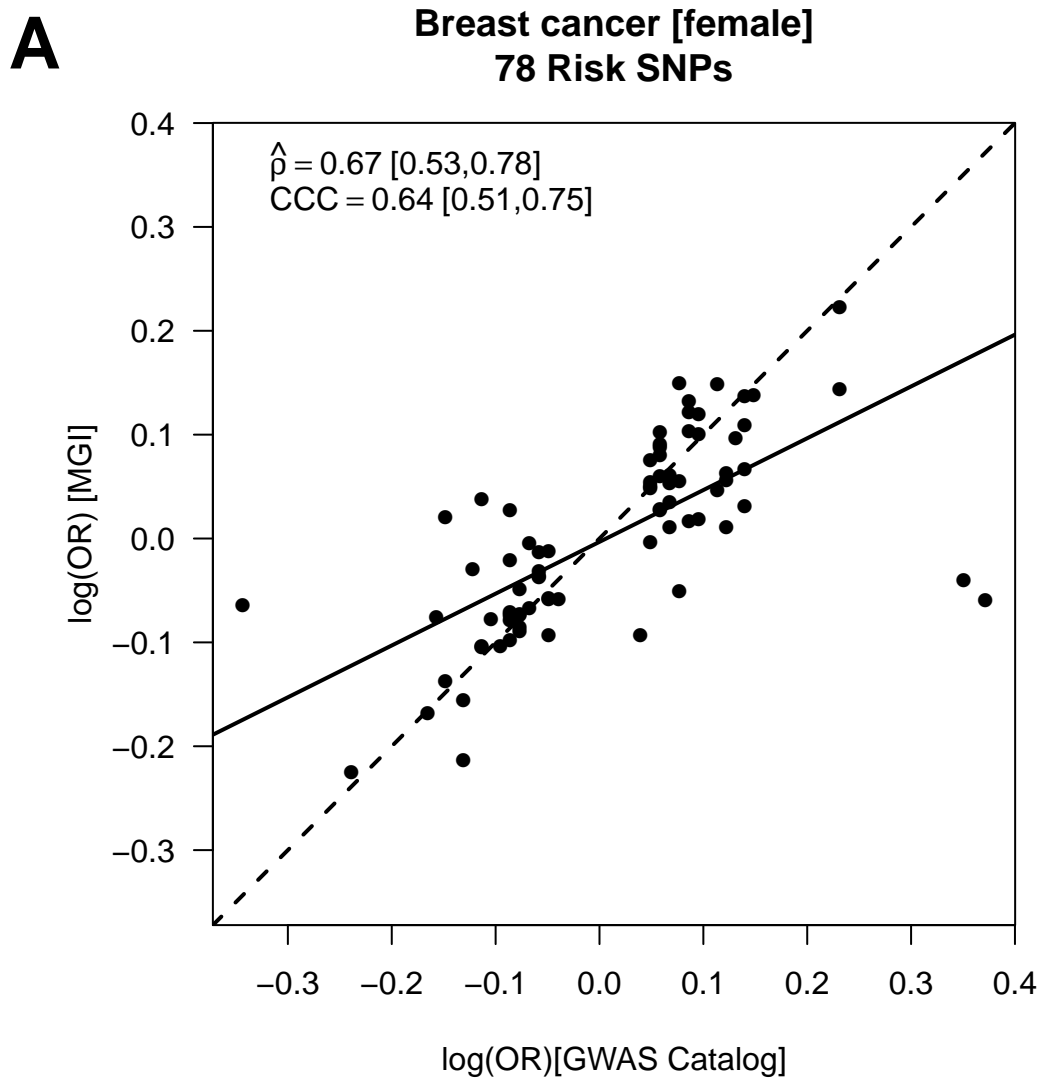
References

1. Witte, J.S. (2010). Genome-wide association studies and beyond. *Annu Rev Public Health* 31, 9-20 24 p following 20.
2. Manolio, T.A. (2010). Genomewide association studies and assessment of the risk of disease. *N Engl J Med* 363, 166-176.
3. Raychaudhuri, S. (2011). Mapping rare and common causal alleles for complex human diseases. *Cell* 147, 57-69.
4. Denny, J.C., Ritchie, M.D., Basford, M.A., Pulley, J.M., Bastarache, L., Brown-Gentry, K., Wang, D., Masys, D.R., Roden, D.M., and Crawford, D.C. (2010). PheWAS: demonstrating the feasibility of a phenome-wide scan to discover gene-disease associations. *Bioinformatics* 26, 1205-1210.
5. Denny, J.C., Bastarache, L., Ritchie, M.D., Carroll, R.J., Zink, R., Mosley, J.D., Field, J.R., Pulley, J.M., Ramirez, A.H., Bowton, E., et al. (2013). Systematic comparison of phenome-wide association study of electronic medical record data and genome-wide association study data. *Nat Biotechnol* 31, 1102-1110.
6. Bush, W.S., Oetjens, M.T., and Crawford, D.C. (2016). Unravelling the human genome-phenome relationship using phenome-wide association studies. *Nat Rev Genet* 17, 129-145.
7. Verma, A., Verma, S.S., Pendergrass, S.A., Crawford, D.C., Crosslin, D.R., Kuivaniemi, H., Bush, W.S., Bradford, Y., Kullo, I., Bielinski, S.J., et al. (2016). eMERGE Phenome-Wide Association Study (PheWAS) identifies clinical associations and pleiotropy for stop-gain variants. *BMC Med Genomics* 9 Suppl 1, 32.
8. Rasmussen, L.V., Overby, C.L., Connolly, J., Chute, C.G., Denny, J.C., Freimuth, R., Hartzler, A.L., Holm, I.A., Manzi, S., Pathak, J., et al. (2016). Practical considerations for implementing genomic information resources. Experiences from eMERGE and CSER. *Appl Clin Inform* 7, 870-882.
9. Ananthakrishnan, A.N., Cagan, A., Cai, T., Gainer, V.S., Shaw, S.Y., Churchill, S., Karlson, E.W., Murphy, S.N., Liao, K.P., and Kohane, I. (2016). Statin Use Is Associated With Reduced Risk of Colorectal Cancer in Patients With Inflammatory Bowel Diseases. *Clin Gastroenterol Hepatol* 14, 973-979.
10. Ananthakrishnan, A.N., Cagan, A., Cai, T., Gainer, V.S., Shaw, S.Y., Churchill, S., Karlson, E.W., Murphy, S.N., Kohane, I., Liao, K.P., et al. (2015). Common Genetic Variants Influence Circulating Vitamin D Levels in Inflammatory Bowel Diseases. *Inflamm Bowel Dis* 21, 2507-2514.
11. Ananthakrishnan, A.N., Cagan, A., Cai, T., Gainer, V.S., Shaw, S.Y., Savova, G., Churchill, S., Karlson, E.W., Murphy, S.N., Liao, K.P., et al. (2016). Identification of Nonresponse to Treatment Using Narrative Data in an Electronic Health Record Inflammatory Bowel Disease Cohort. *Inflamm Bowel Dis* 22, 151-158.
12. Carroll, R.J., Bastarache, L., and Denny, J.C. (2014). R PheWAS: data analysis and plotting tools for phenome-wide association studies in the R environment. *Bioinformatics* 30, 2375-2376.
13. Mavaddat, N., Pharoah, P.D., Michailidou, K., Tyrer, J., Brook, M.N., Bolla, M.K., Wang, Q., Dennis, J., Dunning, A.M., Shah, M., et al. (2015). Prediction of breast

- cancer risk based on profiling with common genetic variants. *J Natl Cancer Inst* 107.
14. Frampton, M., and Houlston, R.S. (2017). Modeling the prevention of colorectal cancer from the combined impact of host and behavioral risk factors. *Genet Med* 19, 314-321.
 15. Szulkin, R., Whittington, T., Eklund, M., Aly, M., Eeles, R.A., Easton, D., Kote-Jarai, Z.S., Amin Al Olama, A., Benlloch, S., Muir, K., et al. (2015). Prediction of individual genetic risk to prostate cancer using a polygenic score. *Prostate* 75, 1467-1474.
 16. Maas, P., Barrdahl, M., Joshi, A.D., Auer, P.L., Gaudet, M.M., Milne, R.L., Schumacher, F.R., Anderson, W.F., Check, D., Chattopadhyay, S., et al. (2016). Breast Cancer Risk From Modifiable and Nonmodifiable Risk Factors Among White Women in the United States. *JAMA Oncol* 2, 1295-1302.
 17. Chatterjee, N., Shi, J., and Garcia-Closas, M. (2016). Developing and evaluating polygenic risk prediction models for stratified disease prevention. *Nat Rev Genet* 17, 392-406.
 18. Garcia-Closas, M., Rothman, N., Figueroa, J.D., Prokunina-Olsson, L., Han, S.S., Baris, D., Jacobs, E.J., Malats, N., De Vivo, I., Albanes, D., et al. (2013). Common genetic polymorphisms modify the effect of smoking on absolute risk of bladder cancer. *Cancer Res* 73, 2211-2220.
 19. Hsu, L., Jeon, J., Brenner, H., Gruber, S.B., Schoen, R.E., Berndt, S.I., Chan, A.T., Chang-Claude, J., Du, M., Gong, J., et al. (2015). A model to determine colorectal cancer risk using common genetic susceptibility loci. *Gastroenterology* 148, 1330-1339 e1314.
 20. Jun, G., Flickinger, M., Hetrick, K.N., Romm, J.M., Doheny, K.F., Abecasis, G.R., Boehnke, M., and Kang, H.M. (2012). Detecting and estimating contamination of human DNA samples in sequencing and array-based genotype data. *Am J Hum Genet* 91, 839-848.
 21. Kent, W.J. (2002). BLAT--the BLAST-like alignment tool. *Genome Res* 12, 656-664.
 22. Chang, C.C., Chow, C.C., Tellier, L.C., Vattikuti, S., Purcell, S.M., and Lee, J.J. (2015). Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* 4, 7.
 23. Manichaikul, A., Mychaleckyj, J.C., Rich, S.S., Daly, K., Sale, M., and Chen, W.M. (2010). Robust relationship inference in genome-wide association studies. *Bioinformatics* 26, 2867-2873.
 24. Wang, C., Zhan, X., Bragg-Gresham, J., Kang, H.M., Stambolian, D., Chew, E.Y., Branham, K.E., Heckenlively, J., Study, F., Fulton, R., et al. (2014). Ancestry estimation and control of population stratification for sequence-based association studies. *Nat Genet* 46, 409-415.
 25. Li, J.Z., Absher, D.M., Tang, H., Southwick, A.M., Casto, A.M., Ramachandran, S., Cann, H.M., Barsh, G.S., Feldman, M., Cavalli-Sforza, L.L., et al. (2008). Worldwide human relationships inferred from genome-wide patterns of variation. *Science* 319, 1100-1104.
 26. McCarthy, S., Das, S., Kretzschmar, W., Delaneau, O., Wood, A.R., Teumer, A., Kang, H.M., Fuchsberger, C., Danecek, P., Sharp, K., et al. (2016). A reference panel of 64,976 haplotypes for genotype imputation. *Nat Genet* 48, 1279-1283.

27. Welter, D., MacArthur, J., Morales, J., Burdett, T., Hall, P., Junkins, H., Klemm, A., Flicek, P., Manolio, T., Hindorff, L., et al. (2014). The NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Res* 42, D1001-1006.
28. MacArthur, J., Bowler, E., Cerezo, M., Gil, L., Hall, P., Hastings, E., Junkins, H., McMahon, A., Milano, A., Morales, J., et al. (2017). The new NHGRI-EBI Catalog of published genome-wide association studies (GWAS Catalog). *Nucleic Acids Res* 45, D896-D901.
29. The 1000 Genomes Project Consortium, Auton, A., Brooks, L.D., Durbin, R.M., Garrison, E.P., Kang, H.M., Korbel, J.O., Marchini, J.L., McCarthy, S., McVean, G.A., et al. (2015). A global reference for human genetic variation. *Nature* 526, 68-74.
30. Winkler, T.W., Day, F.R., Croteau-Chonka, D.C., Wood, A.R., Locke, A.E., Magi, R., Ferreira, T., Fall, T., Graff, M., Justice, A.E., et al. (2014). Quality control and conduct of genome-wide association meta-analyses. *Nat Protoc* 9, 1192-1212.
31. Heinze, G. (2006). A comparative investigation of methods for logistic regression with separated or nearly separated data. *Stat Med* 25, 4216-4226.
32. Heinze, G., Ploner, M., Dunkler, D., and Southworth, H. (2013). logistf: Firth's bias reduced logistic regression. In. (
33. Choi, L., and Beck, C. (2017). EHR: Electronic Health Record (EHR) Data Processing and Analysis Tool. In. (
34. Firth, D. (1993). Bias Reduction of Maximum Likelihood Estimates. *Biometrika* 80, 27-38.
35. Lin, L.I. (2000). A note on the concordance correlation coefficient. *Biometrics* 56, 324 - 325.
36. Lin, L.I. (1989). A concordance correlation coefficient to evaluate reproducibility. *Biometrics* 45, 255-268.
37. Dey, R., Schmidt, E.M., Abecasis, G.R., and Lee, S. (2017). A Fast and Accurate Algorithm to Test for Binary Phenotypes and Its Application to PheWAS. *Am J Hum Genet* 101, 37-49.
38. Ho, D.E., Imai, K., King, G., and Stuart, E.A. (2011). MatchIt : Nonparametric Preprocessing for Parametric Causal Inference. *Journal of Statistical Software* 42, 1-28.
39. R Core Team. (2016). R: A Language and Environment for Statistical Computing. In. (R Foundation for Statistical Computing, Vienna, Austria.
40. Nan, H., Xu, M., Kraft, P., Qureshi, A.A., Chen, C., Guo, Q., Hu, F.B., Curhan, G., Amos, C.I., Wang, L.E., et al. (2011). Genome-wide association study identifies novel alleles associated with risk of cutaneous basal cell carcinoma and squamous cell carcinoma. *Hum Mol Genet* 20, 3718-3724.
41. Barrett, J.H., Iles, M.M., Harland, M., Taylor, J.C., Aitken, J.F., Andresen, P.A., Akslen, L.A., Armstrong, B.K., Avril, M.F., Azizi, E., et al. (2011). Genome-wide association study identifies three new melanoma susceptibility loci. *Nat Genet* 43, 1108-1113.
42. Bishop, D.T., Demenais, F., Iles, M.M., Harland, M., Taylor, J.C., Corda, E., Randerson-Moor, J., Aitken, J.F., Avril, M.F., Azizi, E., et al. (2009). Genome-wide association study identifies three loci associated with melanoma risk. *Nat Genet* 41, 920-925.

43. Brown, K.M., Macgregor, S., Montgomery, G.W., Craig, D.W., Zhao, Z.Z., Iyadurai, K., Henders, A.K., Homer, N., Campbell, M.J., Stark, M., et al. (2008). Common sequence variants on 20q11.22 confer melanoma susceptibility. *Nat Genet* 40, 838-840.
44. Asgari, M.M., Wang, W., Ioannidis, N.M., Itnyre, J., Hoffmann, T., Jorgenson, E., and Whittemore, A.S. (2016). Identification of Susceptibility Loci for Cutaneous Squamous Cell Carcinoma. *J Invest Dermatol* 136, 930-937.
45. Hunter, D.J., Kraft, P., Jacobs, K.B., Cox, D.G., Yeager, M., Hankinson, S.E., Wacholder, S., Wang, Z., Welch, R., Hutchinson, A., et al. (2007). A genome-wide association study identifies alleles in FGFR2 associated with risk of sporadic postmenopausal breast cancer. *Nat Genet* 39, 870-874.
46. Turnbull, C., Ahmed, S., Morrison, J., Pernet, D., Renwick, A., Maranian, M., Seal, S., Ghoussaini, M., Hines, S., Healey, C.S., et al. (2010). Genome-wide association study identifies five new breast cancer susceptibility loci. *Nat Genet* 42, 504-507.
47. Gudmundsson, J., Sulem, P., Gudbjartsson, D.F., Jonasson, J.G., Masson, G., He, H., Jonasdottir, A., Sigurdsson, A., Stacey, S.N., Johannsdottir, H., et al. (2012). Discovery of common variants associated with low TSH levels and thyroid cancer risk. *Nat Genet* 44, 319-322.
48. Gudmundsson, J., Sulem, P., Gudbjartsson, D.F., Jonasson, J.G., Sigurdsson, A., Bergthorsson, J.T., He, H., Blondal, T., Geller, F., Jakobsdottir, M., et al. (2009). Common variants on 9q22.33 and 14q13.3 predispose to thyroid cancer in European populations. *Nat Genet* 41, 460-464.
49. Werner, R.N., Sammain, A., Erdmann, R., Hartmann, V., Stockfleth, E., and Nast, A. (2013). The natural history of actinic keratosis: a systematic review. *Br J Dermatol* 169, 502-518.
50. Ma, C., Blackwell, T., Boehnke, M., Scott, L.J., and Go, T.D.i. (2013). Recommended joint and meta-analysis strategies for case-control association testing of single low-count variants. *Genet Epidemiol* 37, 539-550.
51. Sudlow, C., Gallacher, J., Allen, N., Beral, V., Burton, P., Danesh, J., Downey, P., Elliott, P., Green, J., Landray, M., et al. (2015). UK biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS Med* 12, e1001779.
52. Morales, J., Bowler, E.H., Buniello, A., Cerezo, M., Hall, P., Harris, L.W., Hastings, E., Junkins, H.A., Malangone, C., McMahon, A.C., et al. (2017). A standardized framework for representation of ancestry data in genomics studies, with application to the NHGRI-EBI GWAS Catalog. *bioRxiv*.



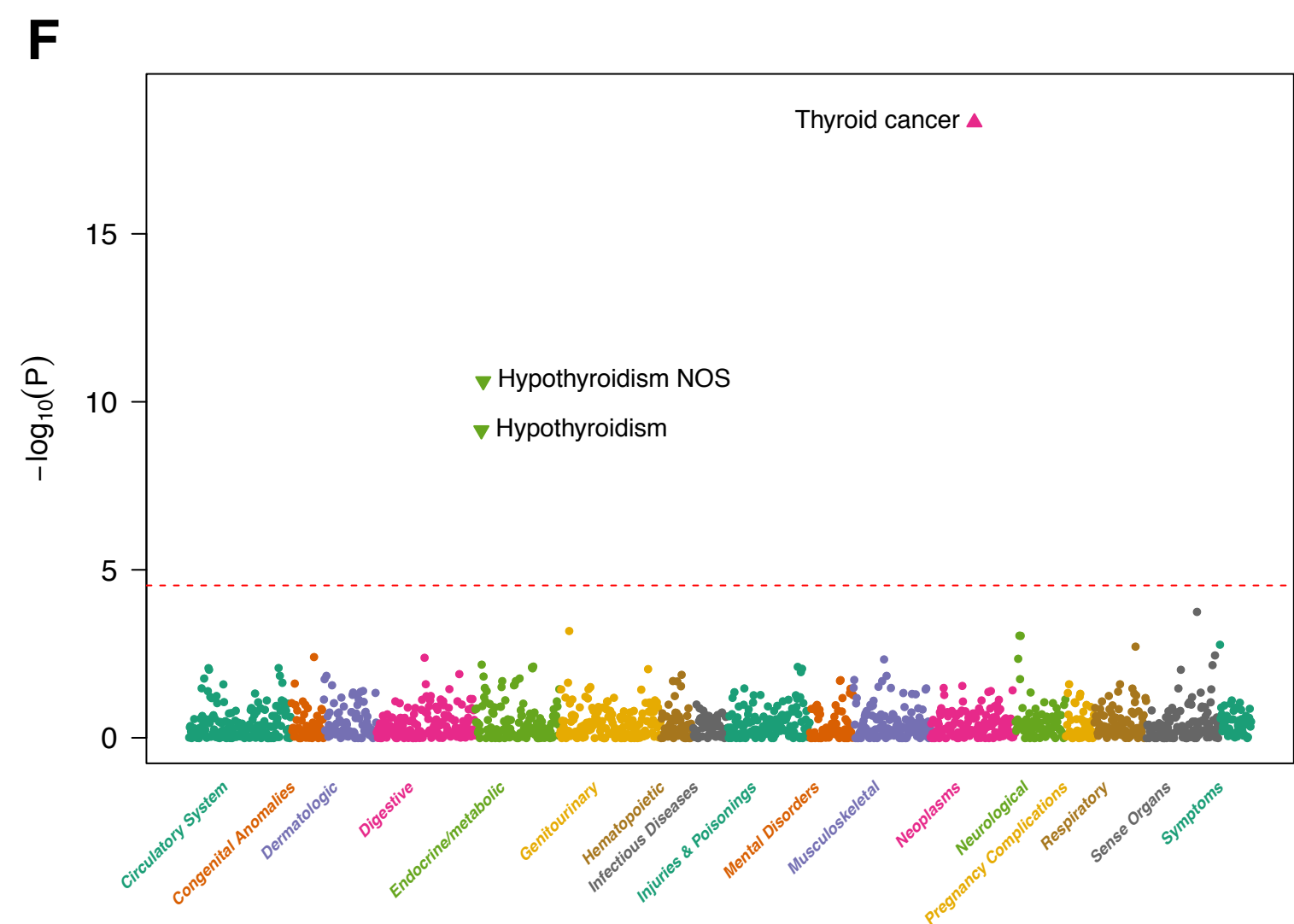
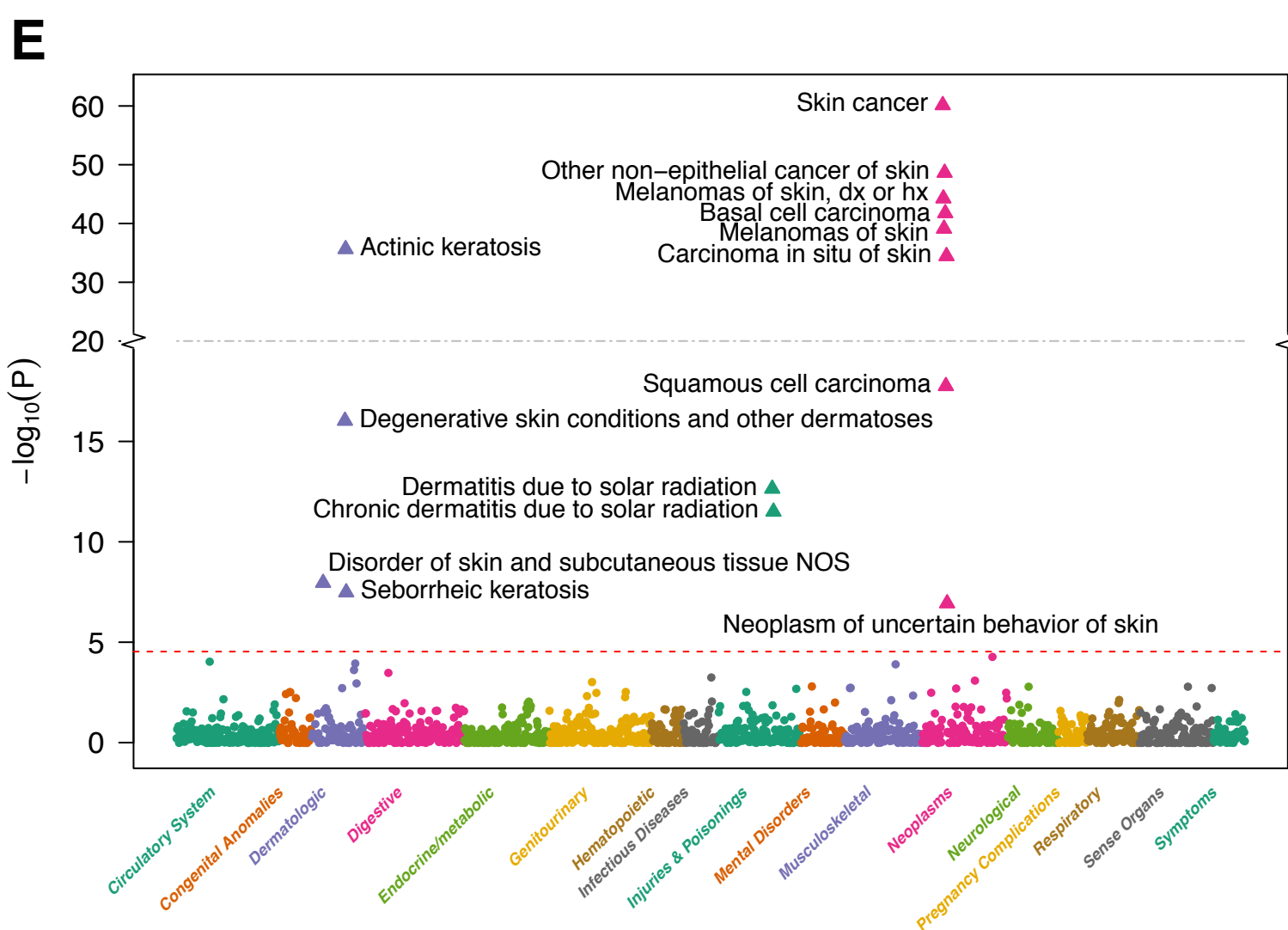
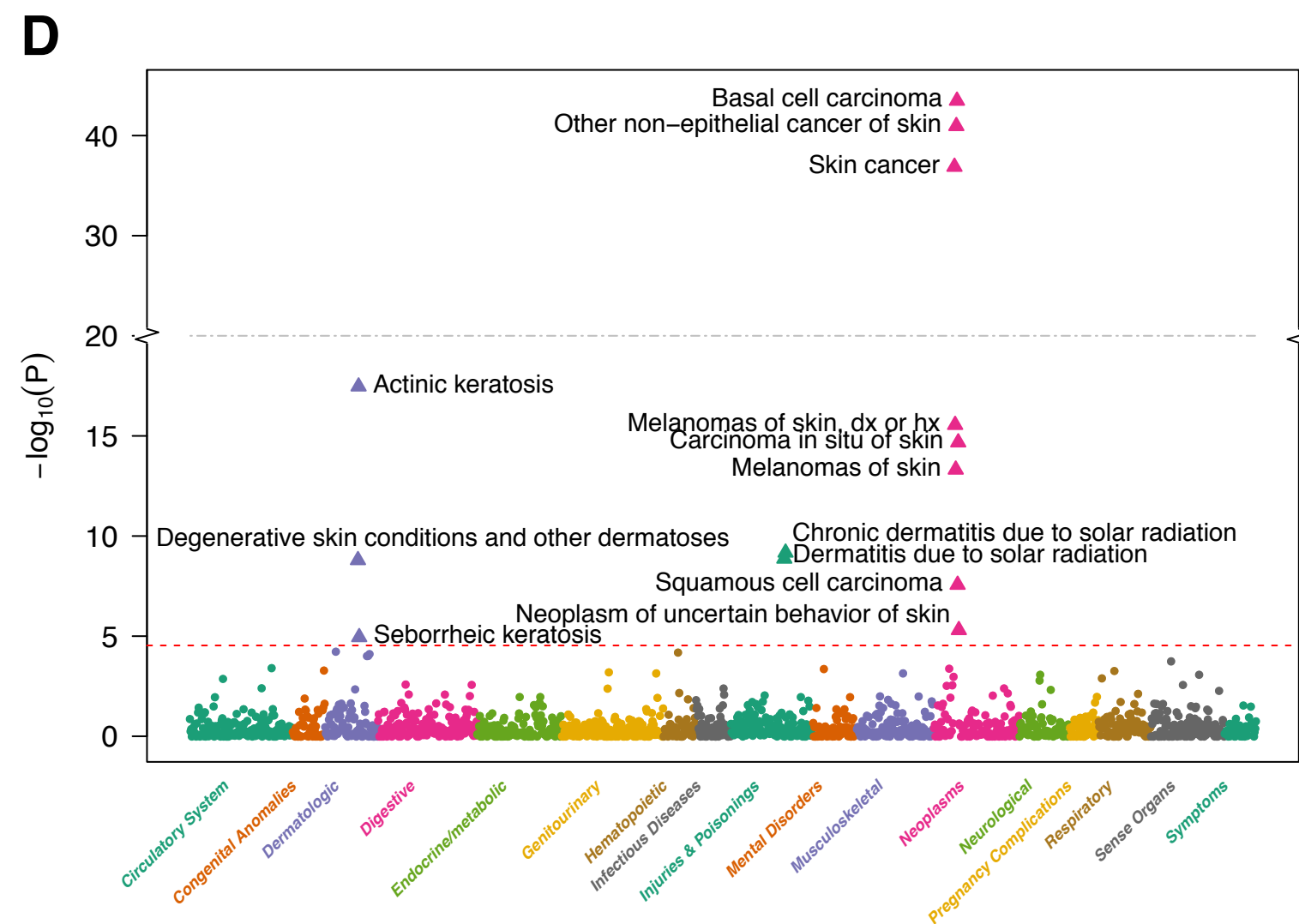
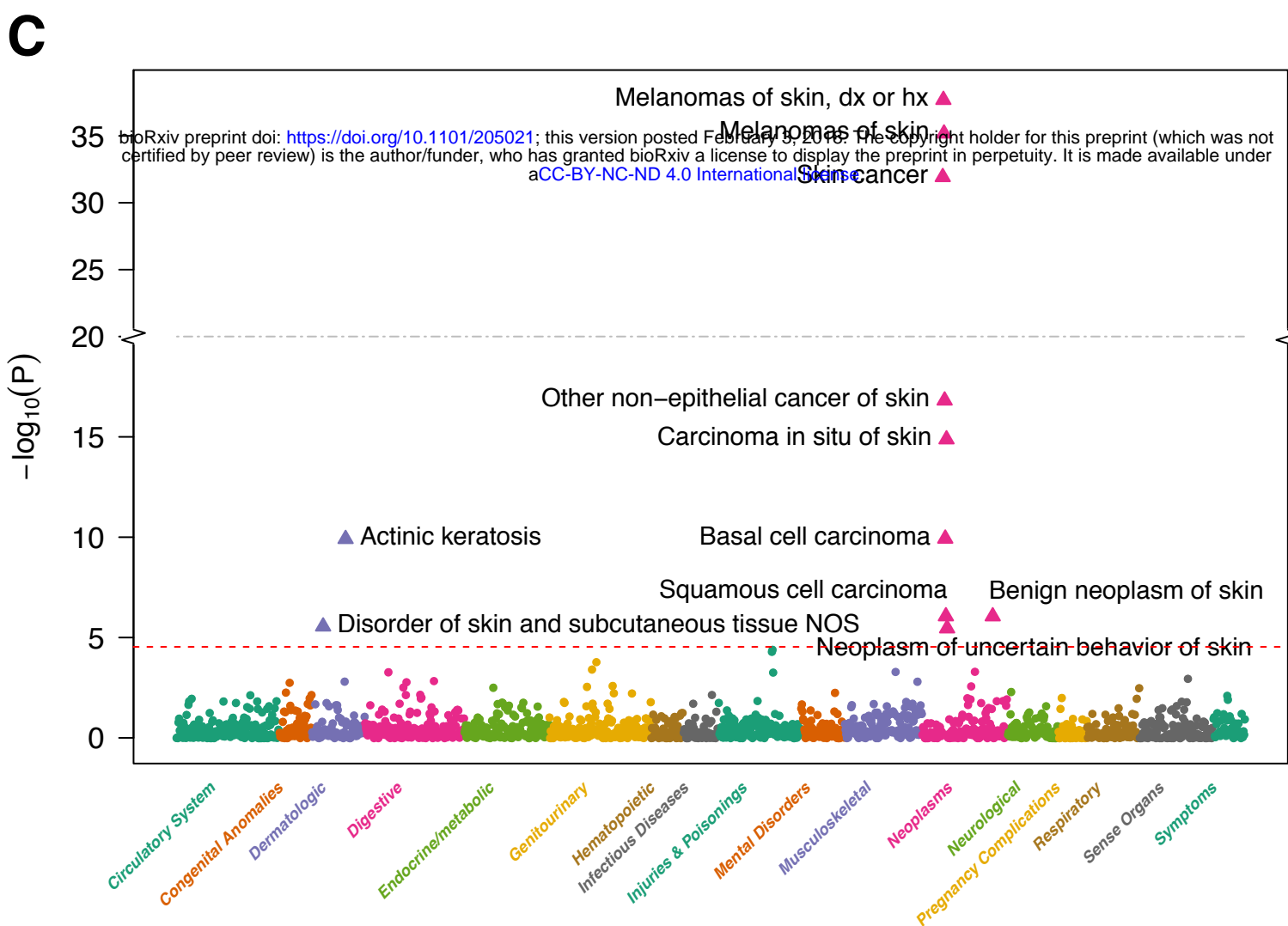
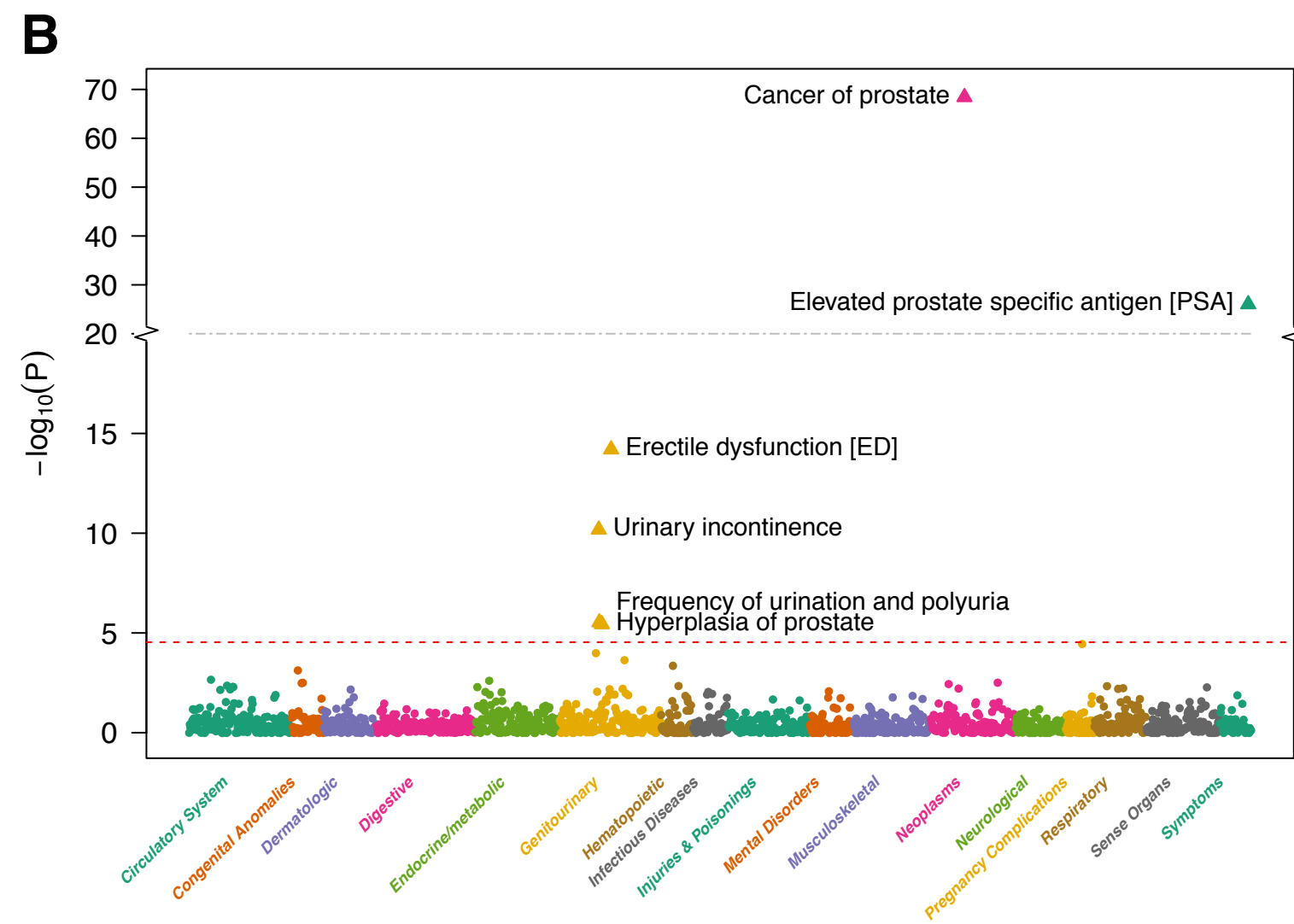
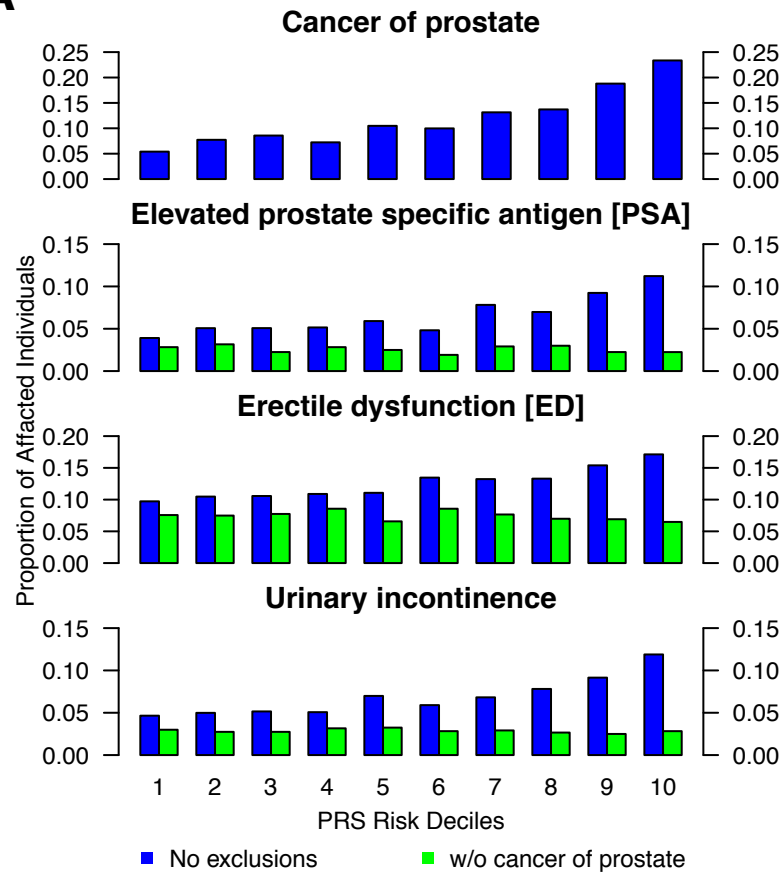
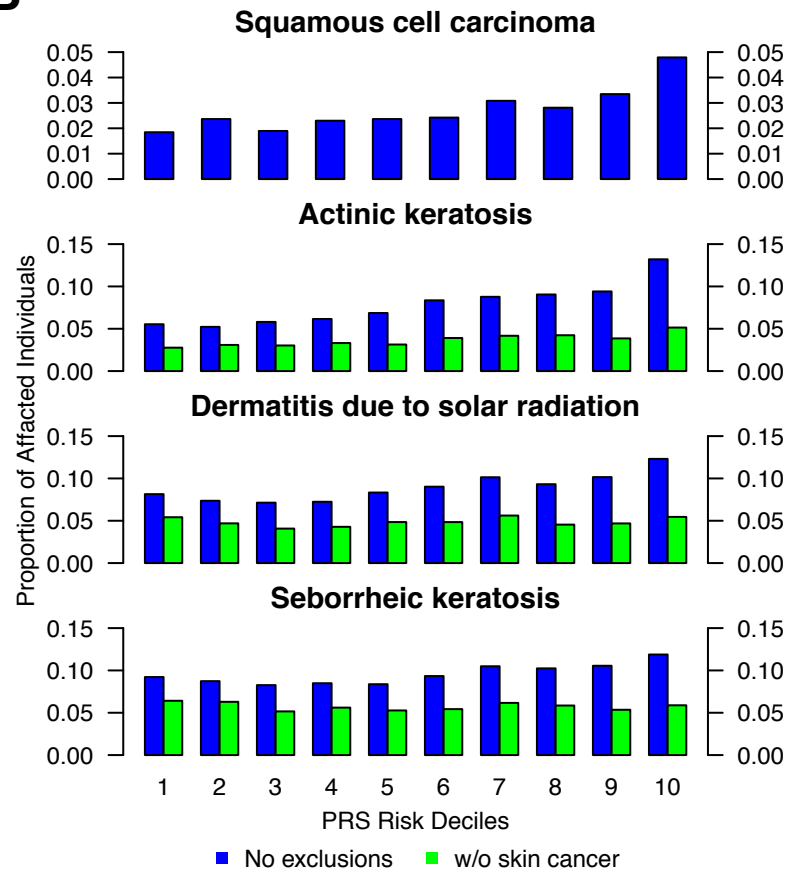


Figure 3

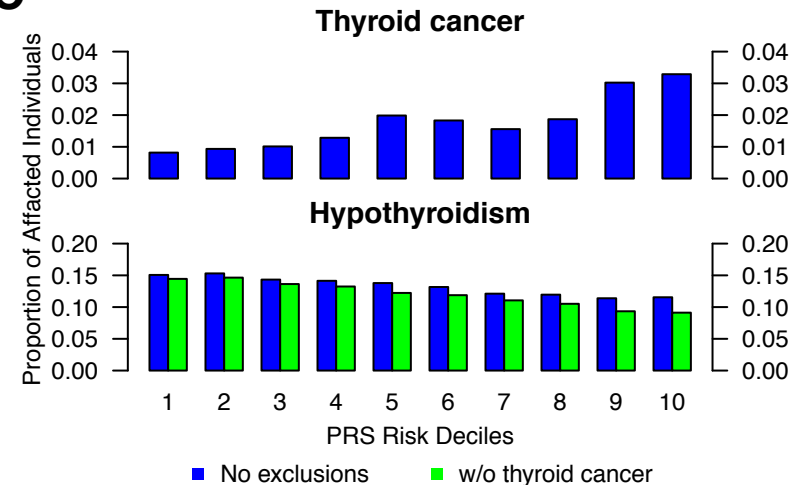
A

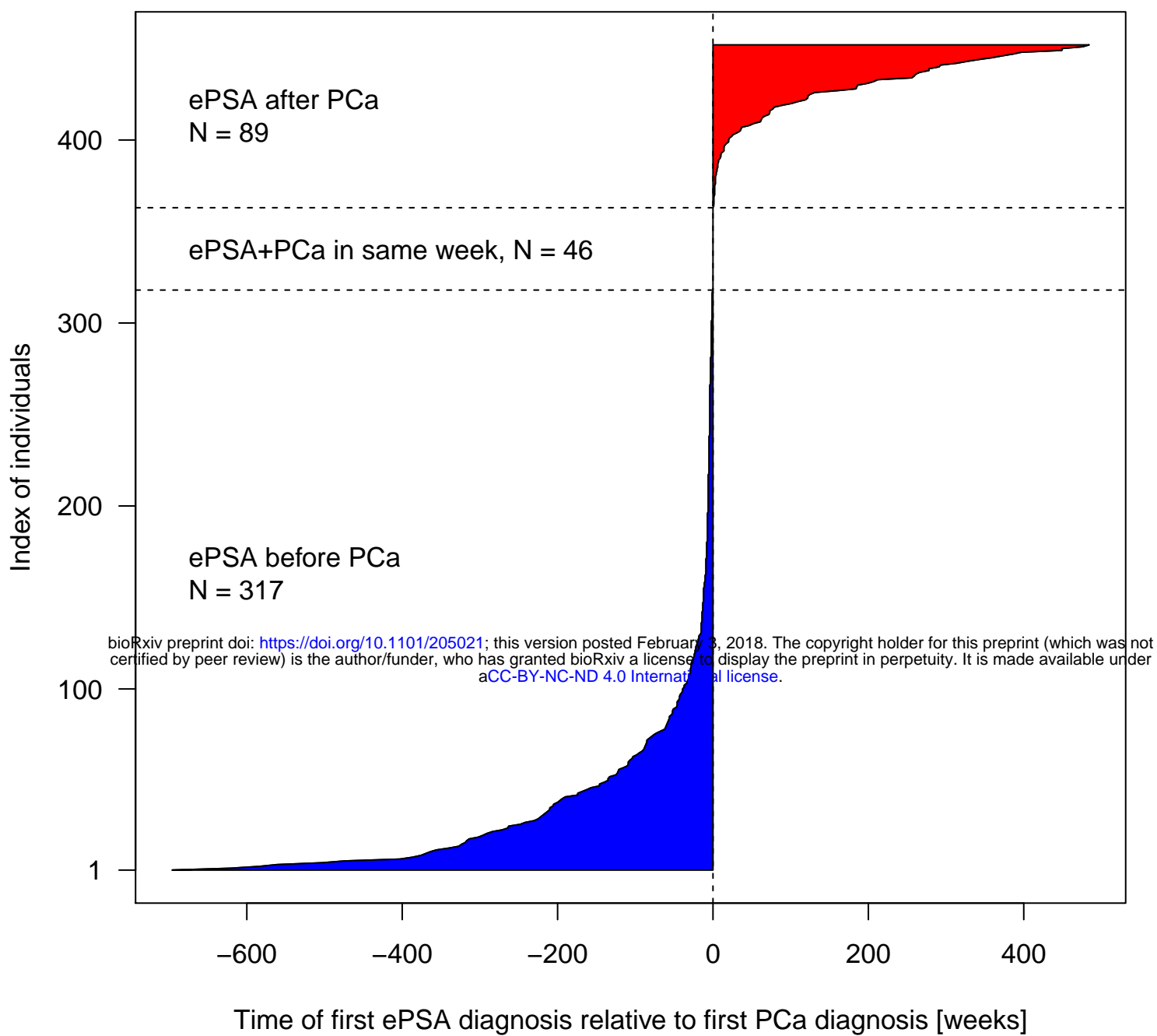
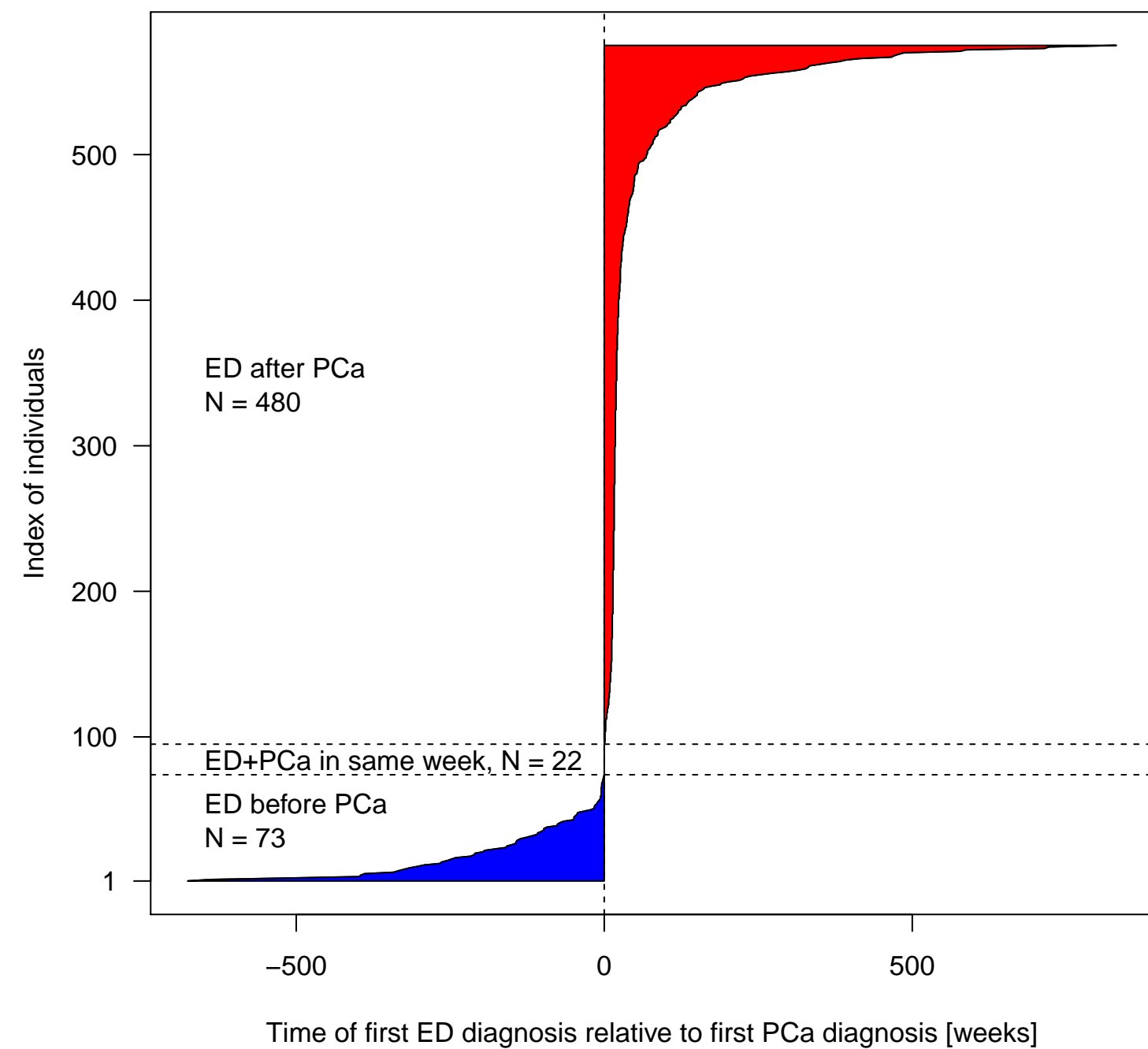
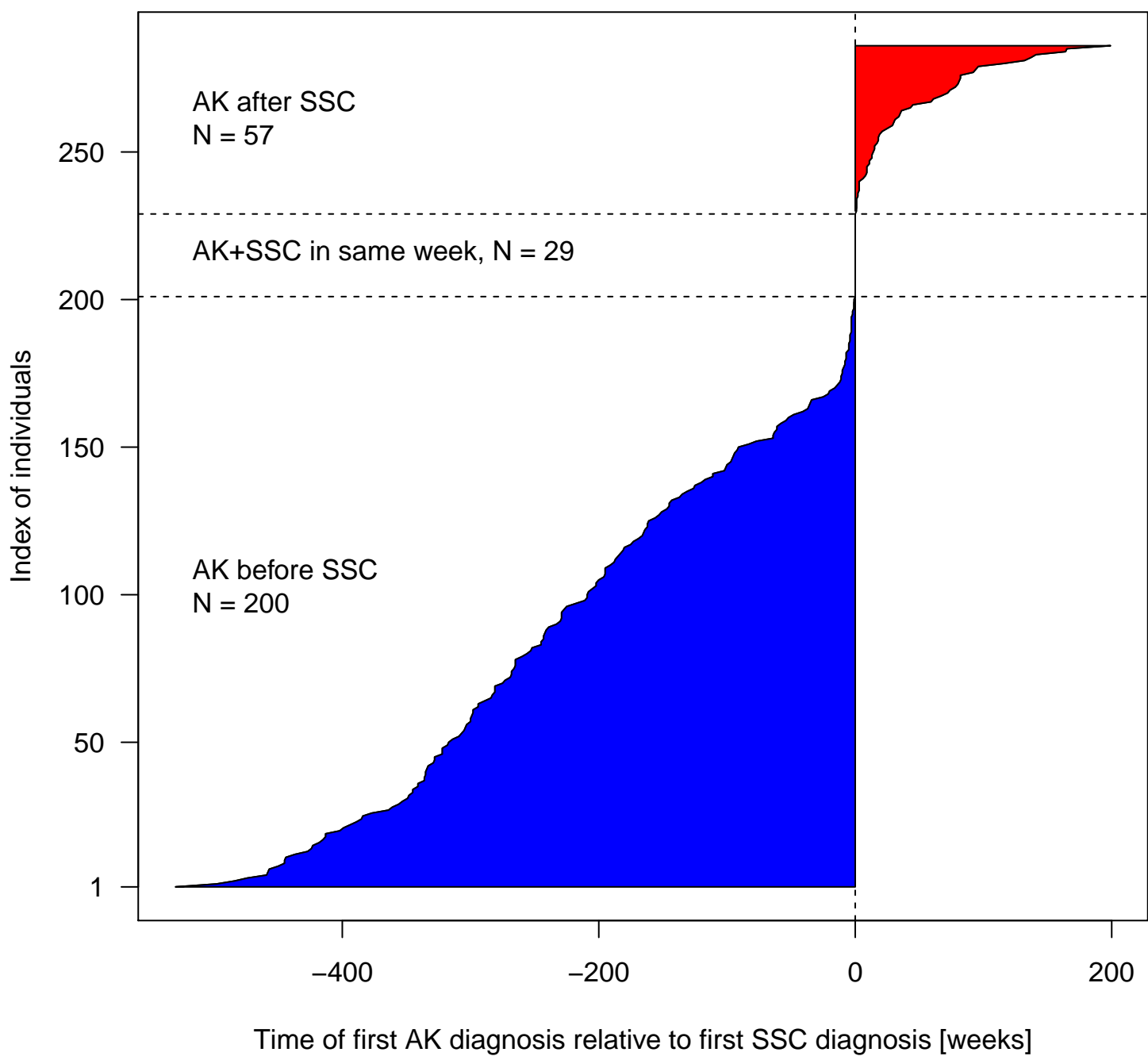


B



C



A**B****C****D**