

## RESEARCH ARTICLE

## Assessing the performance of genome-wide association studies for predicting disease risk

Jonas Patron<sup>1</sup>, Arnau Serra-Cayuela<sup>1</sup>, Beomsoo Han<sup>1</sup>, Carin Li<sup>1</sup>, David Scott Wishart<sup>1,2\*</sup><sup>1</sup> Department of Biological Sciences, University of Alberta, Edmonton, Canada, <sup>2</sup> Department of Computing Science, University of Alberta, Edmonton, Canada\* [dwishart@ualberta.ca](mailto:dwishart@ualberta.ca)

## OPEN ACCESS

**Citation:** Patron J, Serra-Cayuela A, Han B, Li C, Wishart DS (2019) Assessing the performance of genome-wide association studies for predicting disease risk. PLoS ONE 14(12): e0220215. <https://doi.org/10.1371/journal.pone.0220215>

**Editor:** Joseph Devaney, GeneDx, UNITED STATES

**Received:** July 6, 2019

**Accepted:** November 1, 2019

**Published:** December 5, 2019

**Peer Review History:** PLOS recognizes the benefits of transparency in the peer review process; therefore, we enable the publication of all of the content of peer review and author responses alongside final, published articles. The editorial history of this article is available here: <https://doi.org/10.1371/journal.pone.0220215>

**Copyright:** © 2019 Patron et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

**Data Availability Statement:** The code is freely available for download at <https://github.com/jonaspatronjp/GWIZ-Rscript/>, and is compatible with a wide variety of UNIX platforms, Mac OS and Windows operating systems. The GitHub download file also includes an example of the \*.csv input file,

## Abstract

To date more than 3700 genome-wide association studies (GWAS) have been published that look at the genetic contributions of single nucleotide polymorphisms (SNPs) to human conditions or human phenotypes. Through these studies many highly significant SNPs have been identified for hundreds of diseases or medical conditions. However, the extent to which GWAS-identified SNPs or combinations of SNP biomarkers can predict disease risk is not well known. One of the most commonly used approaches to assess the performance of predictive biomarkers is to determine the area under the receiver-operator characteristic curve (AUROC). We have developed an R package called G-WIZ to generate ROC curves and calculate the AUROC using summary-level GWAS data. We first tested the performance of G-WIZ by using AUROC values derived from patient-level SNP data, as well as literature-reported AUROC values. We found that G-WIZ predicts the AUROC with <3% error. Next, we used the summary level GWAS data from GWAS Central to determine the ROC curves and AUROC values for 569 different GWA studies spanning 219 different conditions. Using these data we found a small number of GWA studies with SNP-derived risk predictors that have very high AUROCs (>0.75). On the other hand, the average GWA study produces a multi-SNP risk predictor with an AUROC of 0.55. Detailed AUROC comparisons indicate that most SNP-derived risk predictions are not as good as clinically based disease risk predictors. All our calculations (ROC curves, AUROCs, explained heritability) are in a publicly accessible database called GWAS-ROCS (<http://gwasrocs.ca>). The G-WIZ code is freely available for download at <https://github.com/jonaspatronjp/GWIZ-Rscript/>.

## Introduction

A genome-wide association study (GWAS) is a comprehensive genetic analysis of the association between certain observable traits and specific genetic variations in the form of Single Nucleotide Polymorphisms (SNPs). The appeal of genome wide association (GWA) studies is that they provide a relatively facile approach for detecting potential genetic contributors to common, complex diseases (such as diabetes) or phenotypes (such as body mass index or hair color) using a simple case-control study model. The first GWA study was performed in 2005

as well as instructions on how to run the program. In addition, the complete set of simulated populations used in this study, along with their corresponding SNP profiles, disease status, calculated ROC curves and AUROC data is available for download at <http://gwasrocs.ca>.

**Funding:** Funding for this research has been provided by Genome Canada [<https://www.genomecanada.ca/>] (to DSW), Genome Alberta [<http://www.genomealberta.ca/>] (to DSW), the Canadian Institutes of Health Research [<http://www.cihr-irsc.gc.ca/e/193.html>] (to DSW), the Canada Foundation for Innovation [<https://www.innovation.ca/>] (to DSW) and the Natural Sciences and Engineering Research Council [[http://www.nserc-crsng.gc.ca/index\\_eng.asp](http://www.nserc-crsng.gc.ca/index_eng.asp)] (to DSW). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

**Competing interests:** The authors have declared that no competing interests exist.

[1]. This early work explored the association between certain SNPs and age-related macular degeneration in a study population of 146 individuals (96 cases, 50 controls). To date, thousands of GWA studies looking at almost an equal number of conditions or phenotypes, with study populations as large as 1.3 million have been published [2]. Many of these GWA studies are now archived in public databases such as GWAS Central [3] and the NHGRI-EBI GWAS Catalog [4].

Public databases such as GWAS Central contain summary level findings from GWA studies collected on humans [3]. GWAS Central currently houses data from more than 3300 publications corresponding to over 6100 GWA studies, and lists more than 21 million p-values, ranging between  $5 \times 10^{-2}$  and  $1 \times 10^{-584}$ . The p-values in a GWA study report the likelihood of the odd-ratios between two different alleles being statistically different than one. The typical threshold of significance for most published GWA studies is  $p = 5 \times 10^{-8}$ . The average odds ratio for a statistically significant SNP is 1.33 with very few SNPs having an odds ratio above 3.0 [5].

GWA studies differ from most other 'omics or clinical/epidemiological studies in their method of reporting of disease-associated or disease risk markers. In particular, relatively strong emphasis is placed on the reporting of p-values of a single SNP marker (with p-values of  $< 10^{-50}$  often being achieved) along with its associated odds ratio [6]. Interestingly, very few GWA studies combine multiple SNPs together to produce a multi-marker risk predictor. As a result, it is rare to find a GWA study that reports on the sensitivity, specificity, receiver operating characteristic (ROC) curve or the area under the ROC curve (AUROC or C-statistic) associated with a given SNP or combination of SNPs for disease risk or trait prediction. In contrast to most GWA studies, clinical, proteomic and metabolomic studies involving marker identification rarely achieve p-values of  $< 10^{-6}$  and they infrequently report the performance of their markers in terms of odds ratios. Instead, most marker-based clinical, proteomic or metabolomic studies tend to combine multiple clinical or 'omic measures to generate multi-marker risk predictors. In these studies, multi-marker sensitivity, specificity, ROC curves and/or AUROCs are routinely reported [7].

One way for the different 'omics communities to better understand the predictive or discriminative ability of GWAS data is to convert the reported SNP data into a more conventional biomarker reporting format. In particular, if multi-SNP biomarker data could be combined and converted into ROC or AUROC data, then a more direct comparison could be performed regarding the performance of SNPs for predicting disease risks relative to clinical, metabolite or protein biomarkers for the same conditions. The only precise way of generating ROC curves or calculating AUROC values from GWA studies is to have the complete patient SNP data set. Unfortunately, obtaining ethics approvals and permissions for access to the full patient data sets from thousands of individual GWA studies (found in GWAS Central, for example), conducted in multiple countries and spanning more than 15 years is currently impractical.

Over the past decade several approaches to calculating ROC curves and AUROC data from summary level SNP data have appeared [8–13]. Unfortunately, we found that these methods were not particularly accurate, very limited in their capabilities or required more information than what was available in standard GWAS Central summary data (see [Results](#) for more details). To overcome these issues we developed a novel approach to accurately generate ROC curves and to calculate the AUROC for different SNP combinations using the summary-level data that is standardly found in GWAS databases. Summary level data, which only contains study-wide averages, p-values, odds ratios, risk allele frequencies and other summary statistics for the entire study population, is the only GWAS data can be deposited in public repositories such as GWAS Central. Our new method (called G-WIZ) combines population modeling with logistic regression (for risk prediction) to generate study-specific ROC curves and AUROC

values. Our approach also enables the estimation of SNP heritability directly from summary level GWAS data. We have extensively tested G-WIZ by measuring its ability to predict ROC curves and AUROC values derived from a large collection of authentic individual patient SNP data (i.e. non-summary data), as well as to predict literature-reported AUROC values. We found that G-WIZ predicts the AUROC with <3% error, which is much better than any other method published to date. We then used G-WIZ to, calculate ROC curves and AUROC values for 569 GWA studies spanning 219 different conditions/phenotypes using summary level GWAS data from GWAS Central. Using these data we found at least five conditions/phenotypes that exhibit very high AUROCs (>0.75) and SNP-heritability values (>30%) using a multi-component SNP risk predictor. On the other hand, we found that the average AUROC value for GWAS risk predictors is only 0.55. In contrast, most predictive clinical, metabolite or protein-based risk predictors have AUROCs of >0.7 or 0.8 [7,14]. We also used the data from these calculations to derive a simple formula to estimate SNP heritability directly from the AUROC value. All of the G-WIZ calculations and accompanying analyses (ROC curves, AUROCs, explained heritability, etc.) have been deposited into a publicly accessible database called GWAS-ROCS (<http://gwasrocs.ca>). Having such a large, centralized database of SNP combinations, ROC curves, AUROC values and heritability estimates should open the door to performing more systematic comparisons of GWA studies or to identify new and unexpected trends or novel disease-gene relationships.

## Methods

### Ethics approvals

We received research ethics approval from the University of Alberta (REB—Pro00084706) and approval from the Wellcome Sanger Institute (request ID: 8303 and 7104) to obtain access to the Wellcome Trust Case Control Consortium (WTCCC) GWAS datasets [15].

### Selection of GWAS Central studies

An in-house Python script was created to screen-scrape the summary data in GWAS Central. A total of 3307 GWAS publications corresponding to 6137 GWAS summaries were collected in this manner. These GWA studies were further filtered based on the inclusion of an odds ratio (OR) and a risk allele frequency (RAF) for each reported SNP. Additionally, to improve AUROC value prediction consistency we discarded GWA studies with a sample size of less than 1000 cases or 1000 controls. We thought this was reasonable since the WTCCC datasets all had sample sizes of more than 1000 cases and more than 1000 controls. Moreover, an analysis of every study in GWAS Central indicated that roughly 70% of all published GWA studies (since 2009) have had sample sizes of greater than 1000 cases and 1000 controls. While choosing this threshold should have left us with about 4300 studies to analyze, we found that many studies in GWAS Central were missing either the SNP odds ratios or the risk allele frequencies—or both, which prevented their use in our calculations. Indeed, after applying these filters, we were left with a total of 569 GWAS Central studies, corresponding to 219 different conditions.

### Population simulation for ROC and AUROC calculations

The calculation of sensitivity, specificity, ROC and AUROC data from case/control studies normally requires a data set where the variables of interest (SNPs, proteins, metabolites, clinical measures, etc.) are assigned to individual patients along with their health status. Unfortunately, the GWAS data as found in GWAS Central or other public databases, is only available

as summary data. This means that this patient-specific information is not readily accessible. Indeed, the only information retrievable is the cohort size (cases and controls), the SNP identifiers, the p-values, ORs for each SNP and RAFs for each reported SNP. While it is possible to obtain detailed SNP profiles and health status information for each patient in each study, doing so for the entire GWAS Central Collection would have required extensive ethics reviews along with time and resources that were far beyond our means. This necessitated the development of a modeling program (called G-WIZ or Gwas WIZard) that would calculate ROC and AUROC data from summary-level information only. To do so we exploited that fact that almost no SNPs separated by more than 10 kb are in absolute linkage disequilibrium and that the vast majority of reported, disease-significant SNPs are in Hardy-Weinberg equilibrium [16–18]. As a result, we assumed the independence of SNPs to create simulated patient populations with specific SNP profiles and assigned health conditions from the publicly available OR and RAF data. These synthetic populations were designed to be sufficiently large (typically >30,000 individuals) so that statistically anomalies would be averaged out. To assign a single SNP to an individual in the simulated population the following methodology was used:

- Let  $H_s$  denote the number of risk alleles in the healthy group,
- $H_n$  denote the number of non-risk alleles in healthy group,
- $D_s$  denote the number of risk alleles in the diseased group,
- $D_n$  denote the number of non-risk alleles in diseased group,
- $RAF$  denote risk allele frequency,
- $OR$  denote odds ratio,
- $N_{control}$  denote the number of controls in the simulated dataset,
- $N_{case}$  denote the number of cases in the simulated dataset,
- $\lceil \cdot \rceil$  denote the ceiling function.

Using  $N_{control}$ ,  $N_{case}$ ,  $OR$  and  $RAF$  as the pre-specified input values, we can calculate  $H_s$ ,  $H_n$ ,  $D_s$ ,  $D_n$  as follows:

$$H_s = \lceil 2 \cdot N_{control} \cdot RAF \rceil$$

$$H_n = 2 \cdot N_{control} - H_s$$

$$D_s = \left\lceil \frac{2 \cdot OR \cdot N_{case} \cdot H_s}{H_n + OR \cdot H_s} \right\rceil$$

$$D_n = 2 \cdot N_{case} - D_s$$

Under the assumption that each SNP is independent of the others we can repeat the above procedure to create a full SNP profile and an assigned health state for each member of the simulated population. More specifically, a G-WIZ simulation starts by creating a population of individuals assigned as cases or controls in accordance with the selected GWAS Central record. Next, by using the risk allele frequency for the controls and the odds ratio between the cases and controls G-WIZ calculates the risk allele frequency in the cases. Once the risk allele frequency in both the case and control groups is generated, G-WIZ can appropriately assign the SNP profiles to each group. All G-WIZ models were built using all the SNPs reported by each of the respective GWA studies. These SNPs from GWAS Central were reported on the basis of their significance as identified by the original depositors. However, we considered that it might still be possible that models created using only a subset (feature selection) of reported SNPs would perform better, as this would have controlled for over-parameterization. We tested for this by performing feature selection, using only the SNPs with lowest p-values,

however no improvements to the models' performance were found. In the end our SNP profiles used every reported SNP. On average, each study had a SNP profile consisting of 6 significant SNPs. Moreover, the maximum SNP p-value was  $9 \times 10^{-6}$  and the minimum SNP p-value was  $1 \times 10^{-295}$ , indicating that the reported SNPs are all highly significant. Further, each G-WIZ model had on average 34,491 simulated patients (cases and controls). The largest number of SNPs used in any given SNP profile was 50, and the largest synthetic population generated by G-WIZ consisted of 808,380 individuals.

### Statistical modeling for ROC curve generation

The creation of simulated populations consisting of full SNP profiles and assigned health states for each of the 569 condition/phenotype studies in GWAS Central allowed us to calculate the corresponding ROC curves and AUROC values. A common modeling method used to generate ROC curves for multi-marker data is logistic regression. Logistic regression is a statistical method for modeling multiple independent variables (e.g. SNPs) to explain two possible outcomes (e.g. healthy or diseased). Once constructed, a logistic regression model will return a risk score between 0 and 1. A cut-off value can then be chosen (e.g. 0.5), and any individual that has a risk score above it is classified as 'diseased', and any individual below it is classified as 'healthy'. A plot of the sensitivity against 1-specificity for all possible cut-off values is known as a receiver-operating characteristic (ROC) curve. The classification accuracy of the logistic regression model can then be measured by calculating the area under the curve of the ROC curve (AUROC). A perfect model would have an AUROC of 1, while a model with no classification accuracy would have an AUROC of 0.5 [7].

We performed logistic regression analysis because it is easy to perform and interpret [19]. Although, it is possible that better performing multi-SNP profiles could have been developed using advanced machine learning algorithms such as neural networks, decision trees, or support vector machines [20–26], it is also possible to over-train models with these very powerful pattern-finding tools. Indeed, it is not unusual, during validation studies, to see these models fall short when compared to logistic regression models [20–23]. These concerns regarding overfitting led us to limit our model complexity and to exclusively use logistic and ridge logistic regression to estimate the classification accuracy of these GWA studies.

Another common issue with regression models containing many explanatory variables is multicollinearity. Multicollinearity increases the variance of parameter estimates, which will affect confidence intervals and hypothesis tests. This can lead to incorrect inferences about relationships between explanatory and response variables [27]. With these issues in mind we tested for multicollinearity by estimating the variance inflation factor (VIF) prior to building each regression model. The VIF is the quotient of the variance from a model which regresses one of the predictor variables against all the others. Multicollinearity was determined to exist when at least two variables showed an inflated coefficient (i.e. when the VIF was infinity). We tried a wide range of other VIF cutoff values (less than infinity), however the differences in the AUC estimates were very small ( $<0.009$ ). Whenever multicollinearity was observed we used ridge logistic regression [28] to generate a biomarker model, otherwise we used a standard logistic regression model. Because standard logistic regression is more easily interpretable than its ridge regression counterpart, we found it appropriate to restrict the use of ridge regression only to models with extreme (i.e. divergent) VIF estimates. In total 566 standard logistic regression models and 3 ridge logistic regression models were constructed for all 569 GWA studies.

To assess the performance of each biomarker or ROC-generative model, the simulated data was randomly split into training and testing sets. In the training set, nested cross-validation

(outer 3-fold and inner 2-fold) was used to obtain an estimate of the classification accuracy [28]. Once the model was properly tuned, it was validated using the testing set.

### The G-WIZ program

G-WIZ is written in the R programming language [29]. It consists of several modules including a custom-written tool to generate patient populations, the MLR package [30] to build the logistic regression models and perform cross-validation, as well as the pROC package [31] to generate ROC curves and to calculate the AUROCs. To analyze a study with G-WIZ the user must create a \*.csv file with the study sample size, and the odds ratios and risk allele frequencies for each SNP in the study. A typical study (with a sample size of 50,000 individuals) can be simulated and analyzed via G-WIZ on a modern laptop computer in less than 30 seconds. For any given study, G-WIZ performs population simulation, regression modelling, and then estimates ROC curves and AUROCs. The output consists of a \*.csv file with the regression coefficients, a \*.csv file with the sensitivity, specificity and AUROC of the study, a \*.csv file with the simulated population, and a \*.png file with the ROC curve. The code is freely available for download at <https://github.com/jonaspatronjp/GWIZ-Rscript/>, and is compatible with a wide variety of UNIX platforms, Mac OS and Windows operating systems. The GitHub download file also includes an example of the \*.csv input file, as well as instructions on how to run the program. In addition, the complete set of simulated populations used in this study, along with their corresponding SNP profiles, disease status, calculated ROC curves and AUROC data is available for download at <http://gwasrocs.ca>.

### SNP-derived heritability calculations

For each study in GWAS Central we also estimated the total variance in disease liability explained (often referred to as the SNP heritability) using the following formula described by Pawitan et al. [32].

Let  $var(\cdot)$  denote the variance  
 $OR$  denote the odds ratio  
 $log(\cdot)$  denote the natural logarithm  
 $RAF$  denote the risk allele frequency  
 $\pi$  denote the mathematical constant pi  
 $h^2$  denote the heritability  
 $g$  denote the random genetic effect  
 Then

$$h^2 = \frac{var(g)}{var(g) + \pi^2/3}$$

$$var(g) = 2 \sum_{k=1}^n RAF_k(1 - RAF_k)(logOR_k)^2$$

Where  $n$  is the number of SNPs in a particular study, and  $RAF_k$  and  $OR_k$  are the risk allele frequency and odds ratio of the  $k^{\text{th}}$  SNP. We created an in-house R script to run this formula on all 569 GWA studies collected from GWAS Central. The results are shown in [S1 Table](#).

### Validation

To further ensure that our modeling methods and assumptions were correct, we validated our predictions in two different ways. In the first approach, we used real patient GWAS data from the Wellcome Trust Case Control Consortium (WTCCC). Using this data we ran logistic

regression analyses for 7 different conditions on 2 different control datasets, and calculated the true ROCs and AUROCs. The 7 conditions were bipolar disorder (BD), coronary artery disease (CAD), Crohn's disease (CD), hypertension (HT), rheumatoid arthritis (RA), type 1 diabetes (T1D), and type 2 diabetes (T2D)—spanning AUROC values from 0.51 to 0.72. The two control datasets were WTCCC1 project samples from UK National Blood Service and WTCCC1 project samples from 1958 British Birth Cohort. The SNP profiles and risk alleles we used are reported in Table 1, and are the same SNPs and alleles reported as being statistically significant by the WTCCC researchers [15]. For both control datasets we then applied the G-WIZ modelling method to the same set of SNPs and generated a synthetic population with SNP profiles by directly calculating the RAFs and number of cases and controls for each disease study, and by estimating the ORs from the logistic regression coefficients. The true AUROCs were then compared to our G-WIZ calculated AUROCs. Additionally, the shape of the ROC curves was also compared using Delong's test [33].

The second approach we used to validate our G-WIZ ROC and AUROC predictions involved comparing the G-WIZ predictions against 21 AUROC values provided by previously published GWA studies [34–54]. The GWA studies selected for this validation test covered a broad range of conditions/phenotypes, and a wide range of AUROC values (0.55–0.74). The same set of SNPs reported for these studies were used in the G-WIZ summary-level

**Table 1. SNP profiles and risk alleles reported by the WTCCC publishers.**

WTCCC Disease	SNP	Risk Allele	Minor allele
Bipolar disorder	rs420259	A	G
Coronary artery disease	rs1333049	C	C
Crohn's disease	rs11805303	T	T
Crohn's disease	rs10210302	T	C
Crohn's disease	rs9858542	A	A
Crohn's disease	rs17234657	G	G
Crohn's disease	rs1000113	T	T
Crohn's disease	rs10761659	G	A
Crohn's disease	rs10883365	G	G
Crohn's disease	rs17221417	G	G
Crohn's disease	rs2542151	G	G
Hypertension	rs2820037	T	T
Hypertension	rs6997709	G	T
Hypertension	rs7961152	A	A
Hypertension	rs11110912	G	G
Hypertension	rs1937506	G	A
Hypertension	rs2398162	A	G
Rheumatoid arthritis	rs6679677	A	A
Rheumatoid arthritis	rs6457617	T	T
Type 1 diabetes	rs6679677	A	A
Type 1 diabetes	rs9272346	A	G
Type 1 diabetes	rs11171739	C	C
Type 1 diabetes	rs17696736	G	G
Type 1 diabetes	rs12708716	A	G
Type 2 diabetes	rs9465871	C	C
Type 2 diabetes	rs4506565	T	T
Type 2 diabetes	rs9939609	A	A

<https://doi.org/10.1371/journal.pone.0220215.t001>

simulations. We used the published ORs, RAFs and the number of cases and controls reported in the 21 papers. Logistic regression or ridge regression (as required) was used to generate the ROC curves and to calculate the AUROCs for each of the 21 studies. The reported AUROCs were then compared to our calculated AUROCs. To compile these 21 studies, we created a list of all applicable papers in PubMed by searching the literature for all publications with genomic prediction models that reported an AUROC. This gave us a list of 112 studies (see [S2 Table](#)). Of these 112 studies, we further filtered the studies by study design, sample size requirements, and available SNP information. This left us with 21 studies that were both case-control studies, that had more than 1000 cases and 1000 controls, and that had the relevant SNP summary data to run G-WIZ.

## Results

### Validation testing

A crucial part of developing G-WIZ was to perform an extensive set of validation tests to ensure that our AUROC predictions were accurate. We tested the performance of G-WIZ with respect to four different SNP inheritance models (dominant, recessive, odd SNPs are dominant, and odd SNPs are recessive). This was important because the SNP inheritance model in published GWA studies may vary from SNP to SNP, and is not typically disclosed. Thus, for each of the 7 WTCCC disease studies we calculated the true AUROC under these four SNP inheritance schemes. The results are outlined in [Table 2](#). In total, we ran 56 different tests (spanning all 7 different disease studies, 2 control datasets and 4 SNP inheritance schemes). The average difference in our AUROC predictions was 0.026. The average p-value from the Delong’s test was 0.43, indicating that the predicted ROC curves did not differ significantly from the true ROC curves. Indeed, 52 out of 56 tests had ROC curves that did not differ significantly when tested at the  $p = 0.05$  level of significance.

**Table 2. Comparison of true versus G-WIZ predicted AUROCs for the WTCCC data.**

Cases Cohort	Controls Cohort: WTCCC1 UK National Blood Service		
	True AUROC*	Predicted AUROC*	Difference
WTCCC1 Bipolar disorder	0.52	0.52	0.00
WTCCC1 Coronary artery disease	0.53	0.56	0.03
WTCCC1 Chron’s disease	0.63	0.62	0.01
WTCCC1 Hypertension	0.57	0.57	0.00
WTCCC1 Rheumatoid arthritis	0.61	0.62	0.01
WTCCC1 Type 1 diabetes	0.68	0.70	0.02
WTCCC1 Type 2 diabetes	0.57	0.57	0.00
Cases Cohort	Controls Cohort: WTCCC1 1958 British Birth Cohort		
	True AUROC*	Predicted AUROC*	Difference
WTCCC1 Bipolar disorder	0.52	0.52	0.00
WTCCC1 Coronary artery disease	0.54	0.57	0.03
WTCCC1 Chron’s disease	0.64	0.60	0.04
WTCCC1 Hypertension	0.55	0.55	0.00
WTCCC1 Rheumatoid arthritis	0.61	0.61	0.00
WTCCC1 Type 1 diabetes	0.66	0.71	0.05
WTCCC1 Type 2 diabetes	0.56	0.56	0.00

\* Values are the average calculated over the 4 different SNP inheritance models dominant, recessive, dominant (odd)—recessive (even), dominant (even)—recessive (odd)

<https://doi.org/10.1371/journal.pone.0220215.t002>



**Table 3. Comparison of published versus G-WIZ predicted AUROCs for 21 publications.**

Phenotype/Condition	Published AUROC [Reference]	Predicted AUROC	Difference
Type 2 diabetes	0.58 [36]	0.59	0.01
Type 2 diabetes	0.60 [34]	0.62	0.02
Type 2 diabetes	0.63 [37]	0.62	0.01
Type 2 diabetes	0.60 [35]	0.58	0.02
Type 2 diabetes	0.62 [39]	0.58	0.04
Type 2 diabetes	0.63 [40]	0.60	0.03
Coronary artery disease	0.61 [41]	0.58	0.02
Psoriasis	0.72 [42]	0.67	0.05
Rheumatoid arthritis	0.59 [43]	0.58	0.01
Breast cancer	0.58 [44]	0.63	0.05
Colorectal cancer	0.57 [38]	0.60	0.03
Breast cancer	0.58 [45]	0.61	0.03
Prostate cancer	0.66 [46]	0.65	0.01
Lung cancer	0.55 [47]	0.52	0.04
Esophageal squamous-cell carcinoma	0.63 [48]	0.63	0.00
Rheumatoid arthritis	0.66 [49]	0.69	0.03
Venous thrombosis	0.66 [50]	0.64	0.02
Alzheimer's disease	0.70 [51]	0.77	0.07
Colorectal cancer	0.56 [52]	0.59	0.03
Leprosy	0.74 [53]	0.80	0.06
Glaucoma	0.62 [54]	0.57	0.05

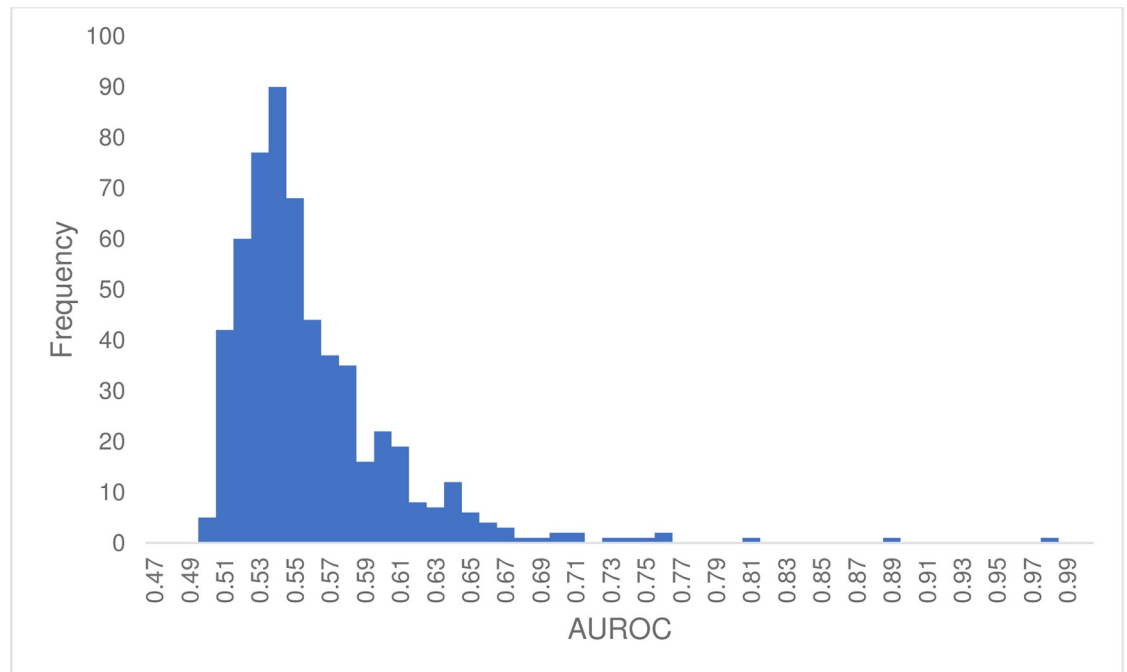
<https://doi.org/10.1371/journal.pone.0220215.t003>

Table 3 lists the AUROC values calculated by our logistic regression modeling technique along with the AUROCs from the 21 previously published GWAS papers. These published AUROC values were derived from GWA studies that used patient-identifiable, authentic clinical data. As can be seen from this table, our AUROC predictions had an average difference of 0.030 relative to the published AUROCs, with a maximum difference of 0.071 and a minimum difference of 0.004. This is remarkably good considering that our data consisted of simulated populations while the published studies were working with complete population data. Previously published studies on AUROC estimates have shown that the type of modeling method used to generate a multi-component biomarker signature can lead to AUROC variations of +/- 0.07 [20–23]. Similarly, we observed that the size of the sample population can also lead to differences in AUROC values (of +/- 0.04) with smaller populations or unbalanced numbers of cases and controls leading to larger differences.

The excellent performance of our G-WIZ AUROC estimates in this validation study (along with the robustness shown in the cross-validation tests) gave us sufficient confidence to apply our logistic modeling method to all 569 GWA studies compiled from GWAS Central and calculate their ROC curves and AUROCs.

### Analysis of GWAS Central studies

The complete set of ROCs, along with the calculated AUROCs for all 569 GWAS Central studies is available at <http://gwasrocs.ca>. These datasets were used to calculate the plot shown in Fig 1. This figure displays a histogram of the calculated AUROC distribution for all 569 GWAS Central studies. The average AUROC was 0.55 while the median AUROC was 0.54.



**Fig 1. Histogram of the G-WIZ predicted AUROCs of the 569 studies from GWAS Central.** Bin width of 0.01.

<https://doi.org/10.1371/journal.pone.0220215.g001>

Ten GWA studies (2%) had an AUROC greater than 0.70 and 73 (13%) had an AUROC above 0.60. The GWA study with the highest AUROC was “Black vs. non-black hair color” [55,56], which had a predicted AUROC of 0.98. The GWA study “Shingles” [57,58] had the lowest AUROC with a predicted AUROC of 0.50. The logistic regression models that we built used on average 6 SNPs. The largest number of SNPs used in a single model was 50. Six studies were modeled using 50 SNPs and 165 studies (29%) were modeled using only a single SNP (see S3 Table).

### GWAS vs. non-GWAS risk prediction performance

One of the motivating factors behind this study was to compare the performance of GWAS-derived or SNP-derived biomarker profiles for disease prediction with other predictive biomarker profiles derived from clinical, metabolomic and/or proteomic (i.e. non-GWAS) data. These data are presented in Table 4 along with types of biomarkers used in each non-GWAS biomarker profile. In compiling this table, we performed an extensive literature survey (of 30 publications) to identify a number of validated or widely used clinical and/or ‘omics’ biomarker sets for disease risk prediction in 12 different conditions/phenotypes (Alzheimer’s, type 2 diabetes, metabolic syndrome, prostate cancer, etc.) along with their reported AUROCs. About one third of the AUROC values quoted in this table used logistic regression while the others used either PLS-DA or random forest classifiers. While the choice of the classification model can affect the reported AUROC, the differences between different types of (good) classifiers are generally small (<5%) [20–23]. The average GWAS AUROC was 0.64 whereas the average non-GWAS AUROC was 0.81.

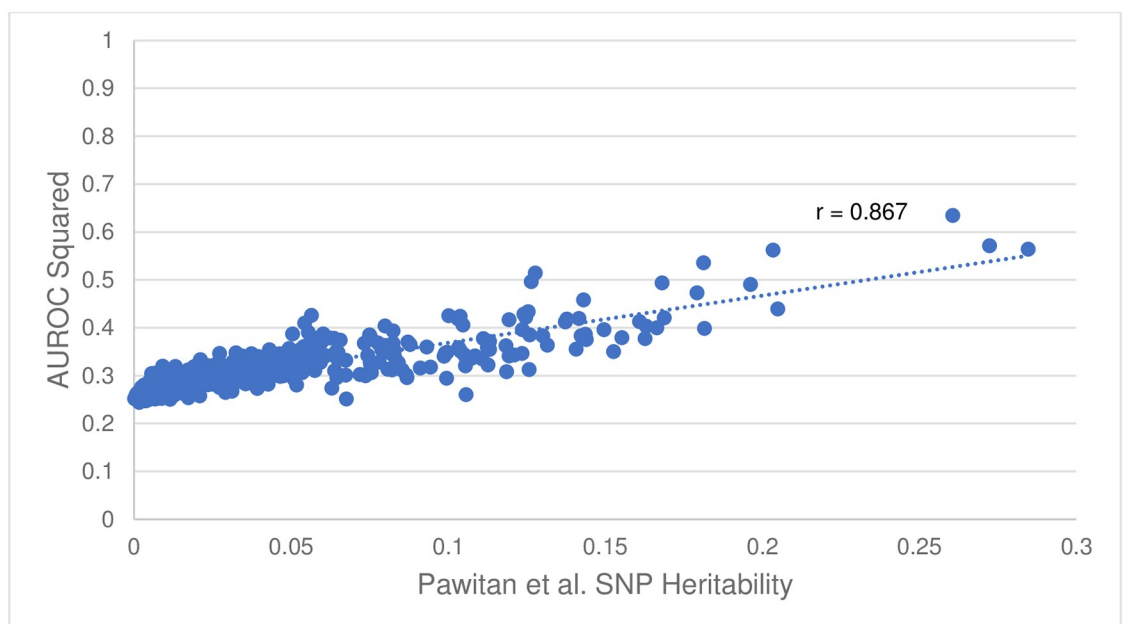
**Table 4. Comparison of typical metabolic, clinical, proteomic and genetic marker AUROCs for common conditions.**

Conditions/Phenotypes	Metabolic Marker AUROC [Reference]	Clinical Factor AUROC [Reference]	Proteomic Marker AUROC [Reference]	Genetic Marker AUROC [Reference]
Alzheimer's Disease	0.770 [59]	0.770 [60]	---	0.640 [61]
Type 2 Diabetes	0.849 [62]	0.870 [63]	---	0.641 [64]
Age Related Macular Degeneration	---	0.780 [65]	---	0.820 [66]
Metabolic Syndrome	0.820 [67]	0.810 [68]	---	0.640 [69]
Colorectal Cancer	0.895 [70]	---	0.885 [71]	0.570 [38]
Prostate Cancer	---	---	0.690 [72]	0.610 [73]
Pulmonary Tuberculosis	---	0.910 [74]	---	0.640 [75]
Cardiovascular Disease	---	0.770 [76]	---	0.600 [77]
Breast Cancer	---	0.690 [78]	0.856 [79]	0.638 [80]
Esophageal Squamous-cell Carcinoma	---	0.639 [48]	---	0.632 [48]
Leprosy	0.862 [81]	---	---	0.707 [53]
Lung Cancer	0.800 [82]	0.903 [83]	---	0.551 [47]
Average:	0.83	0.79	0.81	0.64

<https://doi.org/10.1371/journal.pone.0220215.t004>

### Calculating SNP heritability using the AUROC

While comparing the AUROC values and heritability estimates plotted in the GWAS-ROCS website, an interesting trend was noted. In particular, the heritability seemed to be well correlated with the square of the AUROC ( $r = 0.87$ , see Fig 2). This led to a more detailed investigation regarding the potential rationale for this observation. Upon further reading we found that the  $AUROC = (D + 1)/2$ , where D is the Somers' rank correlation [84] between risk profile and disease status (1 = diseased, 0 = not diseased). Note that the squared Somers' D rank correlation is in fact the proportion of explained variance [85], and that the definition of heritability is precisely the proportion of explained variance. Thus, in the context of a SNP-only model



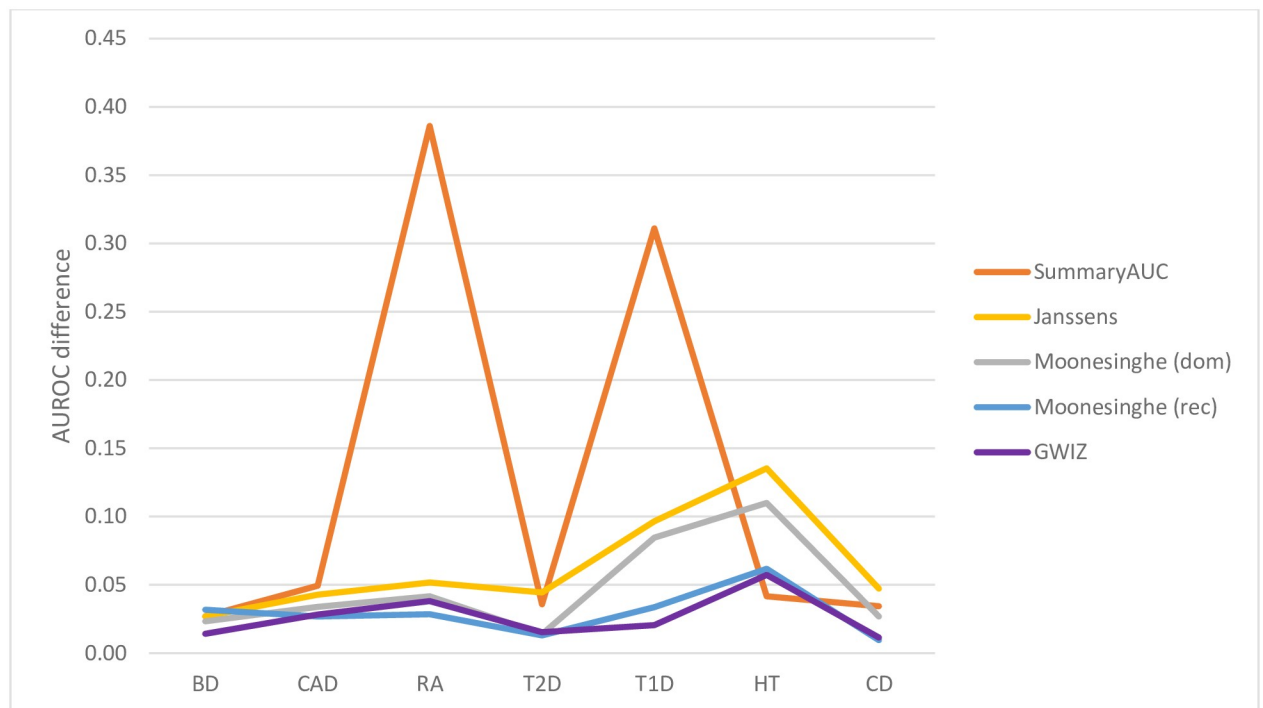
**Fig 2. Plot of the square of the AUROC versus the Pawitan et al. [32] SNP Heritability.** Correlation coefficient,  $r = 0.867$ .

<https://doi.org/10.1371/journal.pone.0220215.g002>

trying to predict disease status, the squared Somers' D rank correlation is, in fact, the SNP heritability. Rearranging for D and then squaring in the formula above, we find that  $h^2 = D^2 = (2 \cdot AUROC - 1)^2$ . This result highlights the utility of AUROC calculations for not only assessing the predictive performance of a multi-SNP panel but for also easily and rapidly calculating heritability of such a SNP panel.

### Comparison with competing methods

As noted earlier, several other methods have been described for estimating AUROC values [8–13] or generating ROC curves [11,12] from summary-level GWAS data. The methods are referred in this paper by the name of the first author. These include the methods by Lu [13], Moonesinghe [9] and Gail [10] which analytically determine AUROC estimates (but not ROC curves), Pepe [12] and Janssens [11] which use population simulation techniques to generate ROC curves from summary-level GWAS data and SummaryAUC [8], a recently released package that can be used to estimate AUROC values based on summary level GWAS data. As noted earlier the methods by Lu, Moonesinghe, Gail and Song do not generate ROC curves nor do they produce population level data. Likewise, the methods by Janssens and Moonesinghe require disease prevalence information to calculate ROC and AUROC data while the method by Janssens is not capable of processing single SNP data. Similarly, the methods by Gail and Lu cannot handle more than 14 SNPs. While SummaryAUC is ideal for calculating AUROCs for a large number of SNPs (100–2000 SNPs) this is not typical of most published GWA studies. As shown in Fig 3, when we compared the actual AUC values generated from the raw WTCCC1 data with the summary-level calculations generated by these programs the performance was disappointing. Indeed, G-WIZ outperforms all of the programs in this test



**Fig 3. Comparison of AUROC predictions accuracies made by G-WIZ and SummaryAUC where the difference between the actual AUROC and predicted AUROC for 7 different disorders is shown.** Bipolar disorder (BD), coronary artery disease (CAD), Crohn's disease (CD), hypertension (HT), rheumatoid arthritis (RA), type 1 diabetes (T1D), and type 2 diabetes (T2D).

<https://doi.org/10.1371/journal.pone.0220215.g003>

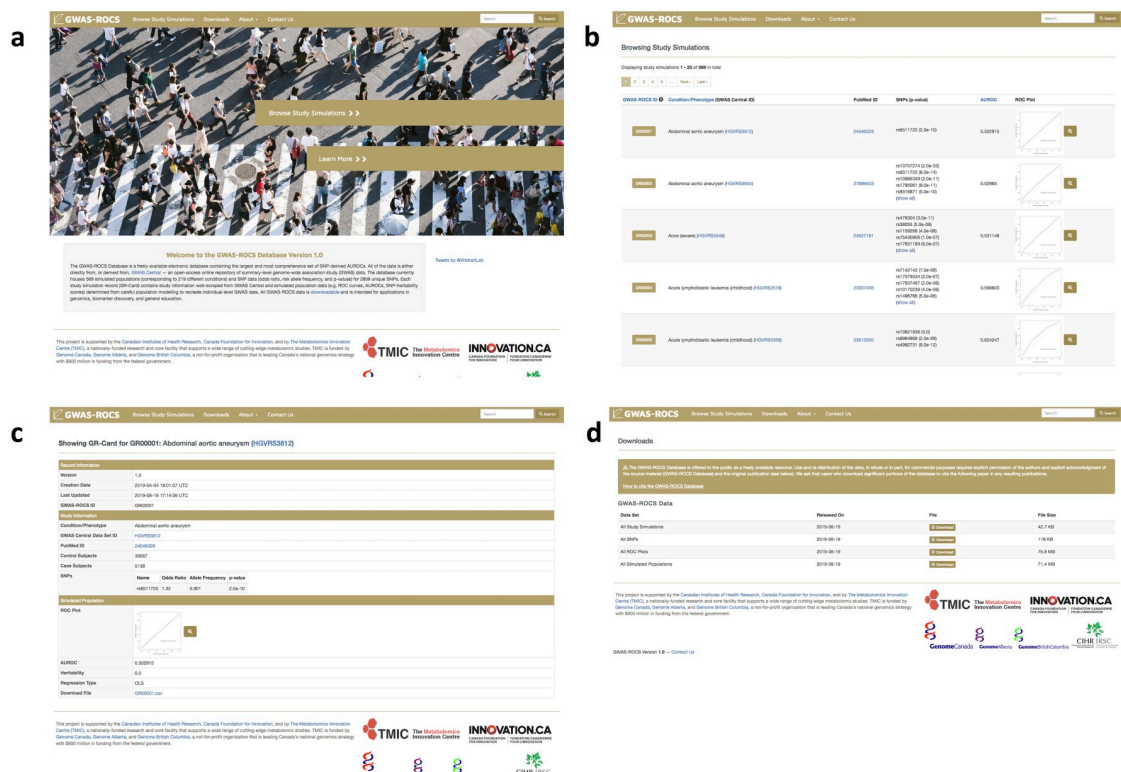
set. In particular, the average prediction difference was 0.026 for G-WIZ, versus 0.13 for SummaryAUC, 0.033 for Moonesinghe's (dominant inheritance) model, 0.029 for Moonesinghe's (recessive inheritance) model and 0.064 for Janssen's model. Based on these comparisons, we believe that G-WIZ is not only more accurate than other tools or methods so far described, but also requires less input data (no need for disease incidence or frequency). Furthermore, G-WIZ handles single SNP as well as multi-SNP profiles (up to 200 SNPs) and generates a much wider range of useful data including ROC curves, AUROC values, heritability estimates, predictive logistic regression models and simulated patient populations. This diverse output can be used for a much wide number of other research applications that go beyond merely knowing the AUROC.

### The GWAS-ROCS database

Along with the creation of G-WIZ, we also created the GWAS-ROCS Database—a freely available electronic database containing SNP-derived AUROCs calculated via G-WIZ. GWAS-ROCS currently contains all the ROC curves, AUROCs, SNP heritability calculations, simulated populations, references, and SNP meta-data generated in the course of producing this paper. The database currently houses 569 simulated populations (corresponding to 219 different conditions) and SNP data (odds ratios, risk allele frequencies, and p-values) for 2986 unique SNPs. Over the coming months, we intend to extend our web scraping to other databases such as NHGRI-EBI GWAS Catalog. In doing so we expect to be able to expand this database and add dozens more GWAS studies. All the data in the GWAS-ROCS database is downloadable, open-access and intended for applications in genomics, biomarker discovery, and general education.

The GWAS-ROCS database also contains numerous datasets with individual level GWAS simulated populations. A GWAS-ROCS simulated population is a \*.csv file with computer-generated “patients” who are marked as either cases or controls. Each individual is assigned data about the presence of risk alleles (1 having the risk allele, or 0 not having the risk allele) at SNPs previously identified as being significant. Users can download these simulated patient populations to develop and test their own genetic risk prediction models or to perform other kinds of synthetic population experiments.

A screenshot montage illustrating the contents and design of the GWAS-ROCS database is shown in Fig 4. As can be seen from this figure the GWAS-ROCS website has a simple webpage layout (Fig 4a). There are four tabs at the top of the page: “Browse Study Simulations”, “Downloads”, “About” and “Contact Us”. Clicking on the “Browse Study Simulations” tab allows users to view a scrollable series of images where they can easily browse through the 569 simulated GWA studies produced in this paper (Fig 4b). Users can sort the GWA studies according to their ID number, the condition/phenotype and the AUROC value. Clicking on a study sends users to a webpage with more information about that specific study (Fig 4c). This includes information such as hyperlinks to the reference GWAS Central study and the original PubMed publication, the number of control and case subjects, the SNP accession IDs, the ORs and RAFs, and the simulated ROC plots, all of which can be found on this page. Additionally, a downloadable \*.csv file with the simulated population for that specific GWA study can be found on this page too. The “Downloads” tab gives users a way to quickly and efficiently download the \*.csv files with simulated populations and ROC plots, for every single study in GWAS-ROCS (Fig 4d). The “About” tab contains some documentation to help users navigate the site. And finally, the “Contact Us” tab gives users an easy way to contact the GWAS-ROCS team with any questions or concerns they may have.



**Fig 4. A screenshot montage of the GWAS-ROCS database at <http://gwasrocs.ca>. a) The landing page. b) Browse Study Simulations tab. c) Study simulation page. d) Downloads tab.**

<https://doi.org/10.1371/journal.pone.0220215.g004>

## Discussion

GWA studies have contributed significantly to our understanding of the genetic contributions to disease and disease risk. Hundreds of novel genes have been identified and implicated in various traits or conditions and many of these have led to new biological understandings and insights [86]. With continued improvements to GWA study designs (increased sample sizes, better population selection to remove confounders, more narrowly defined phenotypes) and GWAS analysis it is likely that many more important biological or genetic insights will be gained [86]. While the value of GWA studies is indisputable, there are still lingering concerns over the inability of SNPs to explain as much of the heritable variation as originally hoped (the missing heritability problem [87,88]) or as much of the disease risk as expected [88].

As remarked earlier, GWA studies that explore disease risk do not often adopt the convention used by most other multi-marker risk predictors to assess performance. In particular, the use of multi-component SNP models and the evaluation of ROC curves or AUROC values (C-statistics) is quite rare. Of the >3700 GWAS publications we evaluated, only 112 have published ROC curves or AUROC data. Of these, fewer than 30 provided sufficient data to independently validate their reported ROC or AUROC results. This has made it difficult to compare the performance of SNP-derived or GWAS-derived biomarkers in disease risk prediction with other types of disease-risk prediction biomarkers or models (clinical, metabolomic, proteomic, etc.). Furthermore, the difference in reporting methodologies between GWA studies (with an emphasis on p-values and odds ratios for individual SNPs) and non-GWA studies (with an emphasis on ROC curves and AUROCs calculated for multiple markers) has also led to an expectation by many non-GWAS specialists, or those with limited statistical

training, that the predictive performance of GWAS-derived biomarkers should be much better than non-GWAS derived biomarkers.

Because of this “cultural” difference we undertook this study to help standardize biomarker reporting between GWAS derived and non-GWAS derived biomarker profiles. In particular, we used logistic regression modeling and simulated patient data to generate a comprehensive and publicly available database of GWAS ROC curves, AUROCs and SNP-heritability scores for a large number of conditions (219) and a large number of GWAS studies (569). These data were placed into an open-access database called GWAS-ROCS, which is publicly available at <http://gwasrocs.ca>.

In creating the GWAS-ROCS database we hoped to accomplish several objectives. First, we wanted to compile and consolidate an accurate and comprehensive set of SNP-derived AUROCs into a single, open-access site. Second, we wanted to use this consolidated data to systematically analyze interesting trends or features in GWAS AUROC data. One of the trends we wanted to explore in more detail concerned the performance of SNP biomarker panels in disease or phenotype prediction. Our results indicate that the average AUROC for a typical GWAS-derived biomarker profile is low, just 0.55 with a standard deviation of 0.05. This is significantly lower than what we expected given that (the few) published AUROCs typically report a range between 0.62–0.67 (see [S2 Table](#), [88]). The fact that published GWAS AUROCs tend to be high (~0.65) and unpublished GWAS AUROCs tend to be low (~0.55), suggests that one reason for the paucity of published GWAS AUROCs is that many AUROCs for SNP biomarker profiles are either uninterestingly low (<0.55), or not statistically different from those generated by a random predictor.

Another aspect that we wanted to explore in more detail was the performance of GWAS-derived SNP profiles for disease prediction compared to non-GWAS profiles for predicting identical diseases. As noted previously, the fact that GWAS disease risk assessments are not typically presented or measured in the same way as non-GWAS disease risk assessments, has made this kind of comparison difficult [89]. As seen in [Table 4](#), we found that non-genetic factors were generally better at predicting disease than genetic factors. In particular, for those conditions where GWAS-derived AUROCs and non-GWAS (clinical/proteomic/metabolic) derived AUROCs could be compared, we found that a typical non-GWA study reported AUROCs closer to 0.81, which is significantly more than the average GWAS-derived biomarker profile of 0.64 (see [Table 4](#)). On the other hand, it is important to note that the predictive ability or disease risk scores of SNP-derived biomarker profiles can effectively occur at birth (many decades prior to the onset of disease) while the non-SNP-derived biomarkers are generally only useful a few months or at most a few years prior to the onset of the disease. In this regard, the utility of SNP-biomarker profiles for long-term disease prevention or disease prophylaxis, even if modest compared to non-SNP profiles, is still quite significant.

A third objective of this study was to identify those conditions that appear to exhibit the best AUROC performance with multi-component SNP data. These high AUROC conditions would be expected to have a relatively high “explainable” genetic component with regard to disease risk. From the data compiled in GWAS Central we identified 5 conditions/phenotypes that had an AUROC greater than 0.75 and an estimated heritability of >25%. As can be seen in [S4 Table](#), age-related macular degeneration (AUROC = 0.75), celiac disease (AUROC = 0.88), progressive supranuclear palsy (AUROC = 0.81), Craniofacial microsomia (AUROC = 0.76) and black hair colour (AUROC = 0.98) can be largely determined through SNP profiles. It is interesting to note that the very first SNP study ever recorded was one done on macular degeneration [1] and that macular degeneration has among the highest levels of heritability and among the highest AUROC values of all conditions we investigated. In many respects, macular degeneration was the equivalent of hitting the “mother lode” for GWA studies.

A fourth objective of this study was to identify those conditions where SNP information appears to be relatively uninformative with regard to disease risk prediction. Our data indicates that 192 conditions (out of 219) have AUROCs  $< 0.60$ , and 202 conditions (out of 219) have SNP heritability values of  $< 15\%$ . The majority of these conditions (such as type 1 diabetes, hypertension, amyotrophic lateral sclerosis and Parkinson's disease) are known to have a significant environmental or lifestyle contribution to disease risk. Indeed, non-GWAS derived risk scores (from clinical, proteomic or metabolomic studies) for several of these conditions have predictive AUROCs approaching 0.8 or 0.9. While many GWA studies have been undertaken due to moderately high heritability data estimated from twin studies, it appears that this heritability may be over-estimated due to small sample sizes or through undetected/unaccounted environmental (gut microbiome heritability) or epigenetic effects. We would suggest that prior to conducting large-scale GWA studies based on twin heritability estimates that some assessment regarding disease risk scores with (published) non-GWAS data should be conducted. Objectively assessing the contributions of clinical/phenotypic data, metabolomic data or proteomic data to disease risk prior to conducting a GWA study would certainly identify which kinds of conditions would most likely yield useful GWAS results or useful GWAS-derived biomarkers.

A fifth objective of this study was to demonstrate the utility of using simulated populations to model SNP distributions and to show how these populations, along with logistic regression modeling, could be used to create multi-marker SNP profiles from publicly available, summary-level GWAS data. Critical to this modeling is the assumption that all genotypes and allele proportions were in Hardy-Weinberg equilibrium with no linkage disequilibrium between genes. The method appears to be robust and particularly well-suited to handling most publicly available SNP data. It is also far faster and less resource intensive than having to apply for study access and research ethics approvals (for non-public data) each time a new GWA study is released or a new GWA study is deposited in GWAS Central. However, simulated populations are not "real" and if the exact population and SNP structure is needed to answer a specific question or if absolute precision is needed in determining some SNP-derived model, then there is obviously no substitution for actual clinical data.

A sixth objective of this study was to explore whether certain trends in AUROCs, disease types or heritability estimates could be discerned by analyzing a large AUROC data set. As noted earlier, we found that the SNP-heritability as determined by Pawitan et al. [32], seemed to be well correlated with the square of the AUROC ( $r = 0.87$ , see Fig 2). This led to the discovery that  $h^2 = (2 \cdot AUROC - 1)^2$ . We used this newly derived formula to estimate the SNP-heritability for all 569 studies from GWAS Central. These heritability estimates can be found at <http://gwasrocs.ca>. Moreover, we compared our heritability estimates against the heritability values for 10 different conditions reported in the literature. Table 5 shows that on average our estimates differed from the true values by just 0.013. The formula we derived suggests that an AUROC of approximately 0.85 is needed to explain 50% of the heritability.

Much is often made of the difference in p-values between GWA studies and the p-values reported for non-GWA studies. As noted before, many GWA studies have SNPs with p-values  $< 1 \times 10^{-50}$ , while it is rare for non-GWA studies to have clinical, protein or metabolite measures with p-values  $< 1 \times 10^{-6}$ . From our dataset one particular GWA study, "Alzheimer's disease" [97,98], had one SNP with the lowest observed p-value ( $1 \times 10^{-295}$ ) of all SNPs in our GWAS Central dataset. This corresponds to a SNP located close to the well-known Alzheimer's disease-associated gene ApoE4 [99]. Using this single SNP alone it was possible to create an Alzheimer's disease risk predictor with an AUROC of 0.62. The addition of 3 more SNPs with p-values of  $2 \times 10^{-10}$ ,  $4 \times 10^{-8}$ ,  $1 \times 10^{-7}$  led to an increase in the Alzheimer's disease risk prediction AUROC to 0.65 [100]. So, despite the extremely low p-values for these Alzheimer's



Table 5. Comparison of published heritability values versus those predicted using the AUROC for 10 conditions.

Condition	Published SNP Heritability	AUROC derived SNP Heritability	Difference
Common psoriasis	0.280 [90]	0.270*	0.010
Crohn's disease	0.053 [91]	0.068**	0.015
Type 1 diabetes	0.127 [92]	0.130**	0.003
Coronary artery disease	0.106 [93]	0.073***	0.033
Age related macular degeneration	0.272 [94]	0.252***	0.020
Lupus	0.209 [95]	0.213***	0.004
Alzheimer's disease	0.059 [95]	0.055***	0.004
Type 2 diabetes	0.082 [95]	0.078***	0.004
Inflammatory bowel disease	0.117 [95]	0.091***	0.026
Rheumatoid arthritis	0.047 [96]	0.062***	0.015
Average:	0.135	0.129	0.013

\* Using AUROC published in the same paper

\*\* Using the true AUROCs calculated from raw WTCCC1 data

\*\*\* Using the highest predicted AUROC from GWAS Central

<https://doi.org/10.1371/journal.pone.0220215.t005>

disease SNPs, the influence on the AUROC (and the heritability) was relatively modest. In another interesting example, the GWA study “Coronary Artery Disease” [101,102] had 50 SNPs, and a p-value of a staggering  $1 \times 10^{-101}$  for the most significant SNP. However, even with so many SNPs and the inclusion of SNPs with remarkably low p-values, the AUROC of this SNP-panel reached just 0.58. Overall, our results indicate that it is not the p-value, but rather the odds ratio (OR), in conjunction with the risk allele frequency (RAF), that are most important for determining biomarker performance in disease risk prediction.

One criticism of our approach to calculating AUROCs from GWAS data is that is computationally inefficient. In particular, we construct large, simulated patient populations, and then used those simulated patient/SNP populations to estimate the AUROCs. A more efficient approach would have been to use machine learning methods or statistical techniques [9,10] to predict the AUROC values directly. There are two reasons why we chose the population simulation approach. First, we believed it would be more useful for the scientific community to have access to simulated patient populations (with SNP data). This would allow others to perform their own statistical or modeling experiments. Furthermore, simulated SNP data can be used to create synthetic “patients” for electronic medical record (EMR) testing and training. Indeed, because of ethics and privacy restrictions, access to individual level GWAS data is often difficult, making generation of realistic genetic data for patients equally difficult. On the other hand, simulated individual level SNP data (and other ‘omics’ data) could be of great utility in the development and testing of “next-generation” EMR software and databases with realistic genetic data. As a result, G-WIZ was created as part of a separate EMR project to generate realistic “synthetic” patients with realistic conditions/phenotypes and correspondingly realistic genomic (SNP), metabolomic, proteomic and clinical profiles. In addition to the appeal of creating synthetic patient data, we also realized that by creating simulated populations we would be able to determine and plot ROC curves (with which we could determine the AUROC values). Having a calculated ROC curve would give us another set of data with which to compare and validate our results. Indeed, we used the G-WIZ generated ROC curves to visually validate a number of the early ROC results during the testing phase of the program.

## Conclusion

To summarize, we have created a software tool called G-WIZ to accurately predict GWAS ROC curves and AUROCs from summary level GWAS data. We subsequently compiled data from every sufficiently informative large-scale study in GWAS Central and calculated the corresponding ROC curves and AUROCs using G-WIZ. Using these calculated data, we conducted a number of comparisons to look for interesting results or unexpected trends. In particular, we compared these calculated GWAS AUROCs to typical AUROCs reported in other 'omics' studies and found some striking differences. We also derived a novel formula to calculate SNP-heritability and calculated the proportion of heritability explained by SNPs for all 569 GWAS Central studies that we analyzed. Through this assessment we were able to make some general suggestions regarding the evaluation and selection of medical conditions that should hopefully yield more significant and useful GWAS outcomes. The results of our G-WIZ calculations, along with other meta-data about each GWA study and the predicted heritability have been placed in an open-access database called GWAS-ROCS.

## Supporting information

**S1 Table. SNP heritability calculated using the method described by Pawitan et al. (2009) for every study in GWAS-ROCS.**

(PDF)

**S2 Table. List of every GWA study found in PubMed which reported an AUROC.**

(PDF)

**S3 Table. A breakdown of the number of SNPs reported by the 569 studies we scraped from GWAS Central.**

(PDF)

**S4 Table. List with the highest predicted AUROC for each of the 219 diseases analyzed from GWAS Central.**

(PDF)

## Acknowledgments

The authors thank Mr. John Bacon for his critical comments and suggested edits.

## Author Contributions

**Conceptualization:** Jonas Patron, Arnau Serra-Cayuela, Beomsoo Han, David Scott Wishart.

**Data curation:** Jonas Patron, Carin Li.

**Formal analysis:** Jonas Patron.

**Funding acquisition:** David Scott Wishart.

**Investigation:** Jonas Patron.

**Methodology:** Jonas Patron, Beomsoo Han, David Scott Wishart.

**Project administration:** David Scott Wishart.

**Resources:** David Scott Wishart.

**Software:** Jonas Patron, Arnau Serra-Cayuela, Beomsoo Han, Carin Li.

**Supervision:** David Scott Wishart.

**Validation:** Jonas Patron.

**Visualization:** Jonas Patron, Carin Li.

**Writing – original draft:** Jonas Patron, David Scott Wishart.

**Writing – review & editing:** Jonas Patron, Arnau Serra-Cayuela, Beomsoo Han, Carin Li, David Scott Wishart.

## References

1. Klein RJ, Zeiss C, Chew EY, Tsai J-Y, Sackler RS, Haynes C, et al. Complement factor H polymorphism in age-related macular degeneration. *Science*. 2005 Apr 15; 308(5720):385–9. <https://doi.org/10.1126/science.1109557> PMID: 15761122
2. Jansen PR, Watanabe K, Stringer S, Skene N, Bryois J, Hammerschlag AR, et al. Genome-wide Analysis of Insomnia (N = 1,331,010) Identifies Novel Loci and Functional Pathways. *bioRxiv*. 2018 Feb 1;214973.
3. Beck T, Hastings RK, Gollapudi S, Free RC, Brookes AJ. GWAS Central: a comprehensive resource for the comparison and interrogation of genome-wide association studies. *Eur J Hum Genet EJHG*. 2014 Jul; 22(7):949–52. <https://doi.org/10.1038/ejhg.2013.274> PMID: 24301061
4. MacArthur J, Bowler E, Cerezo M, Gil L, Hall P, Hastings E, et al. The new NHGRI-EBI Catalog of published genome-wide association studies (GWAS Catalog). *Nucleic Acids Res*. 2017 Jan 4; 45(D1): D896–901. <https://doi.org/10.1093/nar/gkw1133> PMID: 27899670
5. Manolio TA. Genomewide association studies and assessment of the risk of disease. *N Engl J Med*. 2010 Jul 8; 363(2):166–76. <https://doi.org/10.1056/NEJMra0905980> PMID: 20647212
6. Reed E, Nunez S, Kulp D, Qian J, Reilly MP, Foulkes AS. A guide to genome-wide association analysis and post-analytic interrogation. *Stat Med*. 2015 Dec 10; 34(28):3769–92. <https://doi.org/10.1002/sim.6605> PMID: 26343929
7. Xia J, Broadhurst DI, Wilson M, Wishart DS. Translational biomarker discovery in clinical metabolomics: an introductory tutorial. *Metabolomics*. 2013 Apr; 9(2):280–99. <https://doi.org/10.1007/s11306-012-0482-9> PMID: 23543913
8. Song L, Liu A, Consortium MGOS, Shi J, Gejman P, Sanders A, et al. SummaryAUC: a tool for evaluating the performance of polygenic risk prediction models in validation datasets with only summary level statistics. *Bioinformatics*. 2019 Mar 26. btz176. <https://doi.org/10.1093/bioinformatics/btz176> PMID: 30911754
9. Moonesinghe R, Liu T, Khoury MJ. Evaluation of the discriminative accuracy of genomic profiling in the prediction of common complex diseases. *Eur J Hum Genet EJHG*. 2010 Apr; 18(4):485–9. <https://doi.org/10.1038/ejhg.2009.209> PMID: 19935832
10. Gail MH. Discriminatory accuracy from single-nucleotide polymorphisms in models to predict breast cancer risk. *J Natl Cancer Inst*. 2008 Jul 16; 100(14):1037–41. <https://doi.org/10.1093/jnci/djn180> PMID: 18612136
11. Janssens ACJW, Aulchenko YS, Elefante S, Borsboom GJJM, Steyerberg EW, van Duijn CM. Predictive testing for complex diseases using multiple genes: fact or fiction? *Genet Med Off J Am Coll Med Genet*. 2006 Jul; 8(7):395–400.
12. Pepe MS, Gu JW, Morris DE. The potential of genes and other markers to inform about risk. *Cancer Epidemiol Biomark Prev Publ Am Assoc Cancer Res Cosponsored Am Soc Prev Oncol*. 2010 Mar; 19(3):655–65.
13. Lu Q, Elston RC. Using the optimal receiver operating characteristic curve to design a predictive genetic test, exemplified with type 2 diabetes. *Am J Hum Genet*. 2008 Mar; 82(3):641–51. <https://doi.org/10.1016/j.ajhg.2007.12.025> PMID: 18319073
14. Bitton A, Gaziano T. The Framingham Heart Study's Impact on Global Risk Assessment. *Prog Cardiovasc Dis*. 2010; 53(1):68–78. <https://doi.org/10.1016/j.pcad.2010.04.001> PMID: 20620429
15. Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature*. 2007 Jun 7; 447(7145):661–78. <https://doi.org/10.1038/nature05911> PMID: 17554300
16. Shifman S, Kuypers J, Kokoris M, Yakir B, Darvasi A. Linkage disequilibrium patterns of the human genome across populations. *Hum Mol Genet*. 2003 Apr 1; 12(7):771–6. <https://doi.org/10.1093/hmg/ddg088> PMID: 12651872
17. Turner S, Armstrong LL, Bradford Y, Carlson CS, Crawford DC, Crenshaw AT, et al. Quality Control Procedures for Genome Wide Association Studies. *Curr Protoc Hum Genet Editor Board Jonathan Haines Al*. 2011 Jan;CHAPTER:Unit1.19.

18. Marees AT, de Kluiver H, Stringer S, Vorspan F, Curis E, Marie-Claire C, et al. A tutorial on conducting genome-wide association studies: Quality control and statistical analysis. *Int J Methods Psychiatr Res.* 2018 Jun; 27(2).
19. Bush WS, Moore JH. Chapter 11: Genome-Wide Association Studies. *PLoS Comput Biol.* 2012 Dec 27; 8(12): e1002822. <https://doi.org/10.1371/journal.pcbi.1002822> PMID: 23300413
20. Forberg JL, Green M, Björk J, Ohlsson M, Edenbrandt L, Öhlin H, et al. In search of the best method to predict acute coronary syndrome using only the electrocardiogram from the emergency department. *J Electrocardiol.* 2009 Jan 1; 42(1):58–63. <https://doi.org/10.1016/j.jelectrocard.2008.07.010> PMID: 18804783
21. Gulkesen KH, Koksall IT, Bilge U, Saka O. Comparison of methods for prediction of prostate cancer in Turkish men with PSA levels of 0–10 ng/mL. *J Balk Union Oncol.* 2010 Sep; 15(3):537–42.
22. van der Ploeg T, Smits M, Dippel DW, Hunink M, Steyerberg EW. Prediction of intracranial findings on CT-scans by alternative modelling techniques. *BMC Med Res Methodol.* 2011 Oct 25; 11(1):143.
23. Wu J, Roy J, Stewart WF. Prediction Modeling Using EHR Data: Challenges, Strategies, and a Comparison of Machine Learning Approaches. *Med Care.* 2010 Jun; 48:S106–13. <https://doi.org/10.1097/MLR.0b013e3181de9e17> PMID: 20473190
24. Lee HJ, Hwang SI, Han S, Park SH, Kim SH, Cho JY, et al. Image-based clinical decision support for transrectal ultrasound in the diagnosis of prostate cancer: comparison of multiple logistic regression, artificial neural network, and support vector machine. *Eur Radiol.* 2010 Jun 1; 20(6):1476–84. <https://doi.org/10.1007/s00330-009-1686-x> PMID: 20016902
25. Muniz AMS, Liu H, Lyons KE, Pahwa R, Liu W, Nobre FF, et al. Comparison among probabilistic neural network, support vector machine and logistic regression for evaluating the effect of subthalamic stimulation in Parkinson disease on ground reaction force during gait. *J Biomech.* 2010 Mar 3; 43(4): 720–6. <https://doi.org/10.1016/j.jbiomech.2009.10.018> PMID: 19914622
26. Kim S, Kim W, Park RW. A Comparison of Intensive Care Unit Mortality Prediction Models through the Use of Data Mining Techniques. *Healthc Inform Res.* 2011 Dec 1; 17(4):232–43. <https://doi.org/10.4258/hir.2011.17.4.232> PMID: 22259725
27. Midi H, Sarkar SK, Rana S. Collinearity diagnostics of binary logistic regression model. *J Interdiscip Math.* 2010 Jun 1; 13(3):253–67.
28. McDonald GC. Ridge regression. *Wiley Interdiscip Rev Comput Stat.* 2009 Jul 1; 1(1):93–100.
29. R Core Team. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. [Internet]. 2018. <https://www.R-project.org/>.
30. Bischl B, Lang M, Kotthoff L, Schiffner J, Richter J, Studerus E, et al. mlr: Machine Learning in R. *J Mach Learn Res.* 2016; 17(170):1–5.
31. Robin X, Turck N, Hainard A, Tiberti N, Lisacek F, Sanchez J-C, et al. pROC: an open-source package for R and S+ to analyze and compare ROC curves. *BMC Bioinformatics.* 2011; 12(1):77.
32. Pawitan Y, Seng KC, Magnusson PKE. How Many Genetic Variants Remain to Be Discovered? *PLoS ONE.* 2009 Dec 2; 4(12): e7969. <https://doi.org/10.1371/journal.pone.0007969> PMID: 19956539
33. DeLong ER, DeLong DM, Clarke-Pearson DL. Comparing the areas under two or more correlated receiver operating characteristic curves: a nonparametric approach. *Biometrics.* 1988 Sep; 44(3): 837–45. PMID: 3203132
34. Lango H, Palmer CNA, Morris AD, Zeggini E, Hattersley AT, McCarthy MI, et al. Assessing the Combined Impact of 18 Common Genetic Variants of Modest Effect Sizes on Type 2 Diabetes Risk. *Diabetes.* 2008 Nov; 57(11):3129–35. <https://doi.org/10.2337/db08-0504> PMID: 18591388
35. Sparsø T, Grarup N, Andreasen C, Albrechtsen A, Holmkvist J, Andersen G, et al. Combined analysis of 19 common validated type 2 diabetes susceptibility gene variants shows moderate discriminative value and no evidence of gene–gene interaction. *Diabetologia.* 2009 Jul 1; 52(7):1308–14. <https://doi.org/10.1007/s00125-009-1362-3> PMID: 19404609
36. Weedon MN, McCarthy MI, Hitman G, Walker M, Groves CJ, Zeggini E, et al. Combining information from common type 2 diabetes risk polymorphisms improves disease prediction. *PLoS Med.* 2006 Oct; 3(10):e374. <https://doi.org/10.1371/journal.pmed.0030374> PMID: 17020404
37. Miyake K, Yang W, Hara K, Yasuda K, Horikawa Y, Osawa H, et al. Construction of a prediction model for type 2 diabetes mellitus in the Japanese population based on 11 genes with strong evidence of the association. *J Hum Genet.* 2009 Apr; 54(4):236–41. <https://doi.org/10.1038/jhg.2009.17> PMID: 19247372
38. Dunlop MG, Tenesa A, Farrington SM, Ballereau S, Brewster DH, Koessler T, et al. Cumulative impact of common genetic variants and other risk factors on colorectal cancer risk in 42,103 individuals. *Gut.* 2013 Jun; 62(6):871–81. <https://doi.org/10.1136/gutjnl-2011-300537> PMID: 22490517

39. Hu C, Zhang R, Wang C, Wang J, Ma X, Lu J, et al. PPAR $\gamma$ , KCNJ11, CDKAL1, CDKN2A-CDKN2B, IDE-KIF11-HHEX, IGF2BP2 and SLC30A8 are associated with type 2 diabetes in a Chinese population. *PLoS One*. 2009 Oct 28; 4(10):e7643. <https://doi.org/10.1371/journal.pone.0007643> PMID: 19862325
40. Fontaine-Bisson B, Renström F, Rolandsson O, MAGIC, Payne F, Hallmans G, et al. Evaluating the discriminative power of multi-trait genetic risk scores for type 2 diabetes in a northern Swedish population. *Diabetologia*. 2010 Oct; 53(10):2155–62. <https://doi.org/10.1007/s00125-010-1792-y> PMID: 20571754
41. Davies RW, Dandona S, Stewart AFR, Chen L, Ellis SG, Tang WHW, et al. Improved prediction of cardiovascular disease based on a panel of single nucleotide polymorphisms identified through genome-wide association studies. *Circ Cardiovasc Genet*. 2010 Oct; 3(5):468–74. <https://doi.org/10.1161/CIRCGENETICS.110.946269> PMID: 20729558
42. Chen H, Poon A, Yeung C, Helms C, Pons J, Bowcock AM, et al. A genetic risk score combining ten psoriasis risk loci improves disease prediction. *PLoS One*. 2011 Apr 29; 6(4):e19454. <https://doi.org/10.1371/journal.pone.0019454> PMID: 21559375
43. Nakaoka H, Cui T, Tajima A, Oka A, Mitsunaga S, Kashiwase K, et al. A systems genetics approach provides a bridge from discovered genetic variants to biological pathways in rheumatoid arthritis. *PLoS One*. 2011; 6(9):e25389. <https://doi.org/10.1371/journal.pone.0025389> PMID: 21980439
44. Darabi H, Czene K, Zhao W, Liu J, Hall P, Humphreys K. Breast cancer risk prediction and individualised screening based on common genetic variation and breast density measurement. *Breast Cancer Res BCR*. 2012 Feb 7; 14(1):R25. <https://doi.org/10.1186/bcr3110> PMID: 22314178
45. Hüsing A, Canzian F, Beckmann L, Garcia-Closas M, Diver WR, Thun MJ, et al. Prediction of breast cancer risk by genetic risk factors, overall and by hormone receptor status. *J Med Genet*. 2012 Sep 1; 49(9):601–8. <https://doi.org/10.1136/jmedgenet-2011-100716> PMID: 22972951
46. Akamatsu S, Takahashi A, Takata R, Kubo M, Inoue T, Morizono T, et al. Reproducibility, performance, and clinical utility of a genetic risk prediction model for prostate cancer in Japanese. *PLoS One*. 2012; 7(10):e46454. <https://doi.org/10.1371/journal.pone.0046454> PMID: 23071574
47. Li H, Yang L, Zhao X, Wang J, Qian J, Chen H, et al. Prediction of lung cancer risk in a Chinese population using a multifactorial genetic model. *BMC Med Genet*. 2012 Dec 10; 13:118. <https://doi.org/10.1186/1471-2350-13-118> PMID: 23228068
48. Chang J, Huang Y, Wei L, Ma B, Miao X, Li Y, et al. Risk prediction of esophageal squamous-cell carcinoma with common genetic variants and lifestyle factors in Chinese population. *Carcinogenesis*. 2013 Aug; 34(8):1782–6. <https://doi.org/10.1093/carcin/bgt106> PMID: 23536576
49. Scott IC, Seegobin SD, Steer S, Tan R, Forabosco P, Hinks A, et al. Predicting the risk of rheumatoid arthritis and its age of onset through modelling genetic risk variants with smoking. *PLoS Genet*. 2013; 9(9):e1003808. <https://doi.org/10.1371/journal.pgen.1003808> PMID: 24068971
50. Bruzelius M, Bottai M, Sabater-Lleal M, Strawbridge RJ, Bergendal A, Silveira A, et al. Predicting venous thrombosis in women using a combination of genetic markers and clinical risk factors. *J Thromb Haemost JTH*. 2015 Feb; 13(2):219–27. <https://doi.org/10.1111/jth.12808> PMID: 25472531
51. Sleegers K, Bettens K, De Roeck A, Van Cauwenberghe C, Cuyvers E, Verheijen J, et al. A 22-single nucleotide polymorphism Alzheimer's disease risk score correlates with family history, onset age, and cerebrospinal fluid A $\beta$ 42. *Alzheimers Dement J Alzheimers Assoc*. 2015 Dec; 11(12):1452–60.
52. Ibáñez-Sanz G, Díez-Villanueva A, Alonso MH, Rodríguez-Moranta F, Pérez-Gómez B, Bustamante M, et al. Risk Model for Colorectal Cancer in Spanish Population Using Environmental and Genetic Factors: Results from the MCC-Spain study. *Sci Rep*. 2017 24; 7:43263. <https://doi.org/10.1038/srep43263> PMID: 28233817
53. Wang N, Wang Z, Wang C, Fu X, Yu G, Yue Z, et al. Prediction of leprosy in the Chinese population based on a weighted genetic risk score. *PLoS Negl Trop Dis*. 2018; 12(9):e0006789. <https://doi.org/10.1371/journal.pntd.0006789> PMID: 30231057
54. Taylor KD, Guo X, Zangwill LM, Liebmann JM, Girkin CA, Feldman RM, et al. Genetic Architecture of Primary Open-Angle Glaucoma in Individuals of African Descent: The African Descent and Glaucoma Evaluation Study III. *Ophthalmology*. 2019 Jan; 126(1):38–48. <https://doi.org/10.1016/j.ophtha.2018.10.031> PMID: 30352225
55. Lin BD, Mbarek H, Willemsen G, Dolan CV, Fedko IO, Abdellaoui A, et al. Heritability and Genome-Wide Association Studies for Hair Color in a Dutch Twin Family Based Sample. *Genes*. 2015 Jul 13; 6(3):559–76. <https://doi.org/10.3390/genes6030559> PMID: 26184321
56. GWAS-ROCS Database: Showing GR-Card for GR00070: Black vs. non-black hair color (HGVR4156) [Internet]. [cited 2019 Sept 13]. [http://gwasrocs.ca/study\\_simulations/GR00070](http://gwasrocs.ca/study_simulations/GR00070)

57. Crosslin DR, Carrell DS, Burt A, Kim DS, Underwood JG, Hanna DS, et al. Genetic variation in the HLA region is associated with susceptibility to herpes zoster. *Genes Immun*. 2015 Feb; 16(1):1–7. <https://doi.org/10.1038/gene.2014.51> PMID: 25297839
58. GWAS-ROCS Database: Showing GR-Card for GR00478: Shingles (HGVR5860) [Internet]. [cited 2019 Sept 13]. [http://gwasrocs.ca/study\\_simulations/GR00478](http://gwasrocs.ca/study_simulations/GR00478)
59. Mapstone M, Cheema AK, Fiandaca MS, Zhong X, Mhyre TR, MacArthur LH, et al. Plasma phospholipids identify antecedent memory impairment in older adults. *Nat Med*. 2014 Apr; 20(4):415–8. <https://doi.org/10.1038/nm.3466> PMID: 24608097
60. Barnes DE, Covinsky KE, Whitmer RA, Kuller LH, Lopez OL, Yaffe K. Dementia Risk Indices: A Framework for Identifying Individuals with a High Dementia Risk. *Alzheimers Dement J Alzheimers Assoc*. 2010 Mar; 6(2):138–41.
61. Xiao Q, Liu Z-J, Tao S, Sun Y-M, Jiang D, Li H-L, et al. Risk prediction for sporadic Alzheimer's disease using genetic risk score in the Han Chinese population. *Oncotarget*. 2015 Nov 10; 6(35):36955–64. <https://doi.org/10.18632/oncotarget.6271> PMID: 26543236
62. Floegel A, Stefan N, Yu Z, Mühlenbruch K, Drogan D, Joost H-G, et al. Identification of Serum Metabolites Associated With Risk of Type 2 Diabetes Using a Targeted Metabolomic Approach. *Diabetes*. 2013 Feb 1; 62(2):639–48. <https://doi.org/10.2337/db12-0495> PMID: 23043162
63. Paprott R, Mühlenbruch K, Mensink GBM, Thiele S, Schulze MB, Scheidt-Nave C, et al. Validation of the German Diabetes Risk Score among the general adult population: findings from the German Health Interview and Examination Surveys. *BMJ Open Diabetes Res Care*. 2016 Nov 1; 4(1):e000280. <https://doi.org/10.1136/bmjdr-2016-000280> PMID: 27933187
64. Walford GA, Porneala BC, Dauriz M, Vassy JL, Cheng S, Rhee EP, et al. Metabolite Traits and Genetic Risk Provide Complementary Information for the Prediction of Future Type 2 Diabetes. *Diabetes Care*. 2014 Sep; 37(9):2508–14. <https://doi.org/10.2337/dc14-0560> PMID: 24947790
65. Buitendijk GHS, Rohtchina E, Myers C, van Duijn CM, Lee KE, Klein BEK, et al. Prediction of age-related macular degeneration in the general population: the Three Continent AMD Consortium. *Ophthalmology*. 2013 Dec; 120(12):2644–55. <https://doi.org/10.1016/j.ophtha.2013.07.053> PMID: 24120328
66. Grassmann F, Fritsche LG, Keilhauer CN, Heid IM, Weber BHF. Modelling the Genetic Risk in Age-Related Macular Degeneration. *PLoS ONE*. 2012 May 30; 7(5): e37979. <https://doi.org/10.1371/journal.pone.0037979> PMID: 22666427
67. Pujos-Guillot E, Brandolini M, Pétéra M, Grissa D, Joly C, Lyan B, et al. Systems Metabolomics for Prediction of Metabolic Syndrome. *J Proteome Res*. 2017 2; 16(6):2262–72. <https://doi.org/10.1021/acs.jproteome.7b00116> PMID: 28440083
68. Bener A, Yousafzai MT, Darwish S, Al-Hamaq AOAA, Nasralla EA, Abdul-Ghani M. Obesity index that better predict metabolic syndrome: body mass index, waist circumference, waist hip ratio, or waist height ratio. *J Obes*. 2013; 2013:269038. <https://doi.org/10.1155/2013/269038> PMID: 24000310
69. van Pelt ED, Mescheriakova JY, Makhani N, Ketelslegers IA, Neuteboom RF, Kundu S, et al. Risk genes associated with pediatric-onset MS but not with monophasic acquired CNS demyelination. *Neurology*. 2013 Dec 3; 81(23):1996–2001. <https://doi.org/10.1212/01.wnl.0000436934.40034eb> PMID: 24198294
70. Qiu Y, Cai G, Zhou B, Li D, Zhao A, Xie G, et al. A distinct metabolic signature of human colorectal cancer with prognostic potential. *Clin Cancer Res Off J Am Assoc Cancer Res*. 2014 Apr 15; 20(8): 2136–46.
71. Qiao J, Fang C-Y, Chen S-X, Wang X-Q, Cui S-J, Liu X-H, et al. Stroma derived COL6A3 is a potential prognosis marker of colorectal carcinoma revealed by quantitative proteomics. *Oncotarget*. 2015 Oct 6; 6(30):29929–46. <https://doi.org/10.18632/oncotarget.4966> PMID: 26338966
72. Kim Y, Jeon J, Mejia S, Yao CQ, Ignatchenko V, Nyalwidhe JO, et al. Targeted proteomics identifies liquid-biopsy signatures for extracapsular prostate cancer. *Nat Commun*. 2016 28; 7:11906. <https://doi.org/10.1038/ncomms11906> PMID: 27350604
73. Helfand BT, Fought AJ, Loeb S, Meeks JJ, Kan D, Catalona WJ. Genetic prostate cancer risk assessment: common variants in 9 genomic regions are associated with cumulative risk. *J Urol*. 2010 Aug; 184(2):501–5. <https://doi.org/10.1016/j.juro.2010.04.032> PMID: 20620408
74. Heffernan C, Doroshenko A, Egedahl ML, Barrie J, Senthilvelan A, Long R. Predicting pulmonary tuberculosis in immigrants: a retrospective cohort study. *ERJ Open Res*. 2018 Apr; 4(2).
75. Hong EP, Go MJ, Kim H-L, Park JW. Risk prediction of pulmonary tuberculosis using genetic and conventional risk factors in adult Korean population. *PloS One*. 2017; 12(3):e0174642. <https://doi.org/10.1371/journal.pone.0174642> PMID: 28355295

76. Gander J. Factors Related to Coronary Heart Disease Risk Among Men: Validation of the Framingham Risk Score. *Prev Chronic Dis*. 2014 Aug 14; 11:E140. <https://doi.org/10.5888/pcd11.140045> PMID: 25121352
77. Kukava NG, Titov BV, Osmak GJ, Matveeva NA, Kulakova OG, Favorov AV, et al. Multilocus Analysis of Genetic Susceptibility to Myocardial Infarction in Russians: Replication Study. *Acta Naturae*. 2017; 9(4):74–83. PMID: 29340220
78. Winkel RR, von Euler-Chelpin M, Nielsen M, Petersen K, Lillholm M, Nielsen MB, et al. Mammographic density and structural features can individually and jointly contribute to breast cancer risk assessment in mammography screening: a case–control study. *BMC Cancer*. 2016 Jul 7; 16(1):414.
79. Chae YK, Gonzalez-Angulo AM. Implications of functional proteomics in breast cancer. *The Oncologist*. 2014 Apr; 19(4):328–35. <https://doi.org/10.1634/theoncologist.2013-0437> PMID: 24664486
80. Hsieh Y-C, Tu S-H, Su C-T, Cho E-C, Wu C-H, Hsieh M-C, et al. A polygenic risk score for breast cancer risk in a Taiwanese population. *Breast Cancer Res Treat*. 2017 May; 163(1):131–8. <https://doi.org/10.1007/s10549-017-4144-5> PMID: 28205043
81. Al-Mubarak R, Vander Heiden J, Broeckling CD, Balagon M, Brennan PJ, Vissa VD. Serum metabolomics reveals higher levels of polyunsaturated fatty acids in lepromatous leprosy: potential markers for susceptibility and pathogenesis. *PLoS Negl Trop Dis*. 2011 Sep; 5(9):e1303. <https://doi.org/10.1371/journal.pntd.0001303> PMID: 21909445
82. Haznadar M, Cai Q, Krausz KW, Bowman ED, Margono E, Noro R, et al. Urinary Metabolite Risk Biomarkers of Lung Cancer: A Prospective Cohort Study. *Cancer Epidemiol Biomark Prev Oncol*. 2016; 25(6):978–86.
83. Raghu VK, Zhao W, Pu J, Leader JK, Wang R, Herman J, et al. Feasibility of lung cancer prediction from low-dose CT scan and smoking factors using causal models. *Thorax*. 2019 Jul 1; 74(7):643–9. <https://doi.org/10.1136/thoraxjnl-2018-212638> PMID: 30862725
84. Somers RH. A New Asymmetric Measure of Association for Ordinal Variables. *Am Sociol Rev*. 1962; 27(6):799–811.
85. Mittlböck M, Schemper M. Explained variation for logistic regression. *Stat Med*. 1996 Oct 15; 15(19):1987–97. [https://doi.org/10.1002/\(SICI\)1097-0258\(19961015\)15:19<1987::AID-SIM318>3.0.CO;2-9](https://doi.org/10.1002/(SICI)1097-0258(19961015)15:19<1987::AID-SIM318>3.0.CO;2-9) PMID: 8896134
86. Visscher PM, Wray NR, Zhang Q, Sklar P, McCarthy MI, Brown MA, et al. 10 Years of GWAS Discovery: Biology, Function, and Translation. *Am J Hum Genet*. 2017 Jul 6; 101(1):5–22. <https://doi.org/10.1016/j.ajhg.2017.06.005> PMID: 28686856
87. Sandoval-Motta S, Aldana M, Martínez-Romero E, Frank A. The Human Microbiome and the Missing Heritability Problem. *Front Genet*. 2017; 8:80. <https://doi.org/10.3389/fgene.2017.00080> PMID: 28659968
88. Visscher PM, Brown MA, McCarthy MI, Yang J. Five Years of GWAS Discovery. *Am J Hum Genet*. 2012 Jan 13; 90(1):7–24. <https://doi.org/10.1016/j.ajhg.2011.11.029> PMID: 22243964
89. Kundu S, Mihaescu R, Meijer CMC, Bakker R, Janssens ACJW. Estimating the predictive ability of genetic risk models in simulated data based on published results from genome-wide association studies. *Front Genet*. 2014; 5: 179. <https://doi.org/10.3389/fgene.2014.00179> PMID: 24982668
90. Tsoi LC, Stuart PE, Tian C, Gudjonsson JE, Das S, Zawistowski M, et al. Large scale meta-analysis characterizes genetic architecture for common psoriasis associated variants. *Nat Commun*. 2017 May 24; 8:15382. <https://doi.org/10.1038/ncomms15382> PMID: 28537254
91. Yang S-K, Hong M, Zhao W, Jung Y, Baek J, Tayebi N, et al. Genome-wide association study of Crohn's disease in Koreans revealed three new susceptibility loci and common attributes of genetic susceptibility across ethnic populations. *Gut*. 2014 Jan 1; 63(1):80–7. <https://doi.org/10.1136/gutjnl-2013-305193> PMID: 23850713
92. Fan Y, Song Y-Q. Finding the Missing Heritability of Genome-wide Association Study Using Genotype Imputation. *Sci Matters*. 2016 May 4. <https://doi.org/10.19185/matters.201604000013>
93. The CARDIoGRAMplusC4D Consortium, Deloukas P, Kanoni S, Willenborg C, Farrall M, Assimes TL, et al. Large-scale association analysis identifies new risk loci for coronary artery disease. *Nat Genet*. 2013 Jan; 45(1):25–33. <https://doi.org/10.1038/ng.2480> PMID: 23202125
94. Fritsche LG, Igl W, Bailey JNC, Grassmann F, Sengupta S, Bragg-Gresham JL, et al. A large genome-wide association study of age-related macular degeneration highlights contributions of rare and common variants. *Nat Genet*. 2016 Feb; 48(2):134–43. <https://doi.org/10.1038/ng.3448> PMID: 26691988
95. Watanabe K, Stringer S, Frei O, Mirkov MU, Polderman TJC, van der Sluis S, et al. A global overview of pleiotropy and genetic architecture in complex traits. *bioRxiv*. 2018 Dec 19;500090.

96. Okada Y, Wu D, Trynka G, Raj T, Terao C, Ikari K, et al. Genetics of rheumatoid arthritis contributes to biology and drug discovery. *Nature*. 2014 Feb; 506(7488):376–81. <https://doi.org/10.1038/nature12873> PMID: 24390342
97. Seshadri S, Fitzpatrick AL, Ikram MA, DeStefano AL, Gudnason V, Boada M, et al. Genome-wide analysis of genetic loci associated with Alzheimer disease. *JAMA*. 2010 May 12; 303(18):1832–40. <https://doi.org/10.1001/jama.2010.574> PMID: 20460622
98. GWAS-ROCS Database: Showing GR-Card for GR00020: Alzheimer's disease (HGVR1094) [Internet]. [cited 2019 Sept 13]. [http://gwasrocs.ca/study\\_simulations/GR00020](http://gwasrocs.ca/study_simulations/GR00020)
99. Yu C-E, Seltman H, Peskind ER, Galloway N, Zhou PX, Rosenthal E, et al. Comprehensive Analysis of APOE and Selected Proximate Markers for Late-onset Alzheimer Disease: Pattern of Linkage Disequilibrium and Disease/Marker Association. *Genomics*. 2007 Jun; 89(6):655–65. <https://doi.org/10.1016/j.ygeno.2007.02.002> PMID: 17434289
100. GWAS-ROCS Database: Showing GR-Card for GR00299: Late onset Alzheimer's disease (HGVR1241) [Internet]. [cited 2019 Sept 13]. [http://gwasrocs.ca/study\\_simulations/GR00299](http://gwasrocs.ca/study_simulations/GR00299)
101. Nelson CP, Goel A, Butterworth AS, Kanoni S, Webb TR, Marouli E, et al. Association analyses based on false discovery rate implicate new loci for coronary artery disease. *Nat Genet*. 2017 Sep; 49(9): 1385–91. <https://doi.org/10.1038/ng.3913> PMID: 28714975
102. GWAS-ROCS Database: Showing GR-Card for GR00171: Coronary artery disease (HGVR4076) [Internet]. [cited 2019 Sept 13]. [http://gwasrocs.ca/study\\_simulations/GR00171](http://gwasrocs.ca/study_simulations/GR00171)