

RESEARCH

Open Access



# Validating and automating learning of cardiometabolic polygenic risk scores from direct-to-consumer genetic and phenotypic data: implications for scaling precision health research

Arturo Lopez-Pineda<sup>1,2†</sup>, Manvi Vernekar<sup>3,4†</sup>, Sonia Moreno-Grau<sup>1</sup>, Agustin Rojas-Muñoz<sup>1</sup>, Babak Moatamed<sup>1</sup>, Ming Ta Michael Lee<sup>1</sup>, Marco A. Nava-Aguilar<sup>1,2</sup>, Gilberto Gonzalez-Arroyo<sup>1,2</sup>, Kensuke Numakura<sup>3,4</sup>, Yuta Matsuda<sup>3,4</sup>, Alexander Ioannidis<sup>1,5</sup>, Nicholas Katsanis<sup>1</sup>, Tomohiro Takano<sup>3,4\*</sup> and Carlos D. Bustamante<sup>1,5,6\*</sup>

## Abstract

**Introduction:** A major challenge to enabling precision health at a global scale is the bias between those who enroll in state sponsored genomic research and those suffering from chronic disease. More than 30 million people have been genotyped by direct-to-consumer (DTC) companies such as 23andMe, Ancestry DNA, and MyHeritage, providing a potential mechanism for democratizing access to medical interventions and thus catalyzing improvements in patient outcomes as the cost of data acquisition drops. However, much of these data are sequestered in the initial provider network, without the ability for the scientific community to either access or validate. Here, we present a novel geno-pheno platform that integrates heterogeneous data sources and applies learnings to common chronic disease conditions including Type 2 diabetes (T2D) and hypertension.

**Methods:** We collected genotyped data from a novel DTC platform where participants upload their genotype data files and were invited to answer general health questionnaires regarding cardiometabolic traits over a period of 6 months. Quality control, imputation, and genome-wide association studies were performed on this dataset, and polygenic risk scores were built in a case-control setting using the BASIL algorithm.

**Results:** We collected data on  $N = 4,550$  (389 cases / 4,161 controls) who reported being affected or previously affected for T2D and  $N = 4,528$  (1,027 cases / 3,501 controls) for hypertension. We identified 164 out of 272 variants showing identical effect direction to previously reported genome-significant findings in Europeans. Performance metric of the PRS models was  $AUC = 0.68$ , which is comparable to previously published PRS models obtained with larger datasets including clinical biomarkers.

<sup>†</sup>Arturo Lopez-Pineda and Manvi Vernekar have contributed equal to this work

\*Correspondence: tomo@genomelink.io; carlos.bustamante@galatea.bio

<sup>1</sup> Galatea Bio, Inc., 975 W 22nd Street, Hialeah, Florida 33010, USA

<sup>3</sup> Genomelink, Inc., 2150 Shattuck Avenue, Berkeley, California 94704, USA  
Full list of author information is available at the end of the article



**Discussion:** DTC platforms have the potential of inverting research models of genome sequencing and phenotypic data acquisition. Quality control (QC) mechanisms proved to successfully enable traditional GWAS and PRS analyses. The direct participation of individuals has shown the potential to generate rich datasets enabling the creation of PRS cardiometabolic models. More importantly, federated learning of PRS from reuse of DTC data provides a mechanism for scaling precision health care delivery beyond the small number of countries who can afford to finance these efforts directly.

**Conclusions:** The genetics of T2D and hypertension have been studied extensively in controlled datasets, and various polygenic risk scores (PRS) have been developed. We developed predictive tools for both phenotypes trained with heterogeneous genotypic and phenotypic data generated outside of the clinical environment and show that our methods can recapitulate prior findings with fidelity. From these observations, we conclude that it is possible to leverage DTC genetic repositories to identify individuals at risk of debilitating diseases based on their unique genetic landscape so that informed, timely clinical interventions can be incorporated.

**Keywords:** Polygenic risk score, Genome-wide association study, Type 2 diabetes mellitus, Hypertension, Ancestry

## Background

Early diagnosis and prevention of chronic modern diseases, including type 2 diabetes (T2D) and hypertension, have the potential to make a significant impact in patient outcomes. However, the Centers for Disease Control (CDC) estimated that over 20% of T2D cases are undiagnosed [1] and that only 11% of the over 80 million US residents that suffer from prediabetes have been diagnosed (CDC National Diabetes Statistics Report 2017). Early diagnosis could allow for better allocation of intervention strategies known to be effective at reducing the risk of disease progression. According to medical practitioners, insufficient screening is lacking mainly due to the fact that chronic diseases tend to progress slowly until they manifest clinically later in life. One of the main barriers to effectively identifying individuals at risk is the lack of predictive tools trained on heterogeneous datasets that are able to predict susceptibility using historical data available outside of clinical and research settings.

The World Health Organization (WHO) reports a sustained increase in diabetes mellitus, with projections increasing to 3% of the world population by 2030, becoming the seventh leading cause of death globally [1]. A sedentary lifestyle and a diet pattern with high intake of foods rich in hydrogenated fat, refined grains, and red meat have contributed to the increase in overweight and obesity and led to the increased incidence of T2D [2]. An important challenge to this health crisis is to decrease mortality, especially at younger ages, and in low and low-middle countries [3], with more than 400 million people affected globally [4].

The overlap between T2D and hypertension is common among the population [5]. Hypertension alone affects more than 1.28 billion people worldwide [6]. T2D can lead to complications which can be exacerbated when the patient also presents hypertension, for example, in the progression of diabetic nephropathy [7]. Both T2D

and hypertension are risk factors associated with stroke and other serious and life-threatening events [8]. In fact, during the recent COVID-19 pandemic, outcomes of patients seem to be negatively affected by the presence of T2D, hypertension, and obesity [9].

Genetics of T2D has been extensively studied [10–15], with over 400 genetic variants found to be associated with the diseases [16]. In addition to individual studies focusing on defined ethnic groups like Hispanics [17, 18], there have been consortium efforts to investigate the genetic architecture of complex traits in diverse populations. Some of these consortia include a) the Population Architecture using Genomics and Epidemiology (PAGE) consortium [19]; b) the rich data offered by the UK Biobank allowing associations between complex traits, genetics, and lifestyle [20]; c) the Trans-Omics for Precision Medicine (TOPMed) Consortium [21] which improves imputation quality and detection of rare variant associations; and more recently d) the Meta-Analysis Biobank Initiative [22], a collaborative network of biobanks across the world representing millions of consented individuals.

Direct to consumer platforms are novel sources of information that have expanded quickly during the past decade. The earliest example in 2010 was the use of web-based self-reported questionnaires with complementary genetic testing, leading to the creation of a research database [23] which has allowed for novel polygenic risk scores in complex traits [24] and subsequent FDA submissions of novel diagnostics [25]. The clinicogenomic database developed by a consortium is led by a large pharmaceutical company, alongside an electronic health record company focused on oncological practices, and a direct to consumer (DTC) genetic testing company, putting together a comprehensive database [26], allowing sophisticated analysis of including the selection of novel biomarkers [27], drug effectiveness studies [28], or automatic eligibility criteria selection [29].

The rapid development of polygenic risk scores (PRS) in recent years stresses the importance of accurately assessing the ancestry makeup of participants in biomedical studies to avoid potential selection biases [30]. PRS represents a measure of an individual's overall genetic liability to a trait or disease [31]. The European Bioinformatics Institute (EBI) has launched a PRS catalog database, allowing for reproducibility and standardization of reporting of PRS models [32]. As of Feb 16, 2022, the catalog includes 35 models related to the T2D from 15 peer reviewed publications [33–45], and six models related to hypertension across 3 publications [38, 43, 45]. However, these models mostly include single ancestry participants (typically European) which may not generalize across other ancestral groups.

Even though the need for ancestry-focused research has been highlighted by many [46, 47], the lack of diversity resulted in systemic biases that threaten to widen existing health disparities among minority and majority populations in most developed countries. The overrepresentation of European individuals in genetic studies represents a major issue, hampering the translation of PRS across populations. In this study, we hypothesize that PRS models can be improved by defining the genetic ancestry of participants. Embedding genetic ancestry as a covariate, or scaling PRS scores as part of post-processing step, would result in more accurate models than traditional filtering to only European-based PRS models.

Furthermore, the validation of a DTC framework for validating and extending PRS provides a cost-effective means of enrolling understudied populations in complex disease genetics as 23andMe has done with their “Roots into the Future” effort. We aim to build on such successes by powering a public–private partnership that is collecting DTC data to be included in the Biobank Meta-analysis Network. The notion that DTC provided a “straight to mobile instead of landlines” opportunity is important to validate as most low- and middle-income countries are considering how to harness advances in genomics for the study of their own populations while building state-sponsored capacity is a fundamental challenge to many efforts.

## Methods

In this case–control retrospective observational study, adult participants from an international genetic platform were invited to self-report their health status and metabolic traits. Their genetic information was also previously uploaded in the same platform, which allowed us to explore their genetic susceptibility and to build polygenic risk scores (PRS) regarding these traits. Finally, we calculated ancestry estimation using Neural ADMIXTURE for all individuals. We were interested in evaluating

ancestry-aware polygenic scores for type 2 diabetes and hypertension.

### Cohort and eligibility criteria

We used the research database where participants were drawn from Genomelink (genomelink.io) users, which offers a DTC genetic traits platform with more than 500,000 users globally. After uploading their genetic information, generated in other DTC platforms, users can be informed on their susceptibility to an extensive set of genetic traits. All participants created an account and agreed to a consent on the use of their data and legal agreement. Upon signing up, participants were invited to undertake a health online survey. Participants were redirected to the survey once they gave online consent to be a part of the research. The online consent is in compliance with the institutional review board (IRB) at WCG IRB (<https://www.wcgirb.com/>) under IRB tracking ID 20,201,332.

The online survey included questions about general conditions like diabetes, blood pressure, lipid profile, and medication intake. It also included COVID-19, influenza and common cold-related questions along with age, sex, weight, height, and pandemic behavior. Data were collected over a period of six months, from May 01, 2021 to October 06, 2021. Additional file 1 shows the online questionnaire.

Only the initial answers of each participant were included in the study, if genotype information was available, and if they answered the age and sex questions. Case–control groups were created following participants' answers to the general condition question: T2D and hypertension. Additionally, for T2D, participants were included if they reported to have high levels of sugar in their blood work or if they were taking antidiabetic medication. Participants were defined as controls if they did not report managing health conditions listed in the questionnaire survey, or if they were not managing any health condition. Also, for the T2D cohort specifically, participants who reported having normal sugar levels were included as controls. Participants with missing values were excluded.

### Genotype data: quality control, imputation, and GWAS

This study includes seven independent genotyping arrays, comprising a total of 12,424 unrelated individuals. Genotype-level data for each array were processed by applying identical quality control and imputation procedures. Briefly, variants with a call rate of <95% and palindromic markers (A/T, G/C, MAF >0.4) were excluded. We performed an exact test for Hardy–Weinberg equilibrium for individuals of the largest ancestral group ( $p < 1 \times 10^{-12}$ , globally). Individual quality control (QC)

includes genotype call rates >97%, matching between gender identification and chromosomal sex, and no excess ancestry-adjusted heterozygosity. Samples genetically related to other individuals in the cohort and duplicates were detected and removed, by applying the King algorithm (`-make-king, king estimate > 0.177`; PLINK 2). Principal component analysis was performed to identify global ancestry per individual using 1000 genomes as reference population with PLINK 2 [48]. Further information about the number of markers per genotyping array pre- and post-QC is available on Additional file 2.

Imputation was carried out using 1000 genomes as a reference panel with Beagle [49]. Next, we generated a merged dataset combining imputed genotypes ( $MAF > 0.01$ ; imputation quality  $R^2 > 0.30$ ) from available datasets. Imputed makers with call rate >0.95 in the merged data were selected for downstream analysis.

The GWAS was performed for T2D and hypertension phenotypes ( $N = 4,550$ ;  $N = 4,528$ , respectively) using an additive genetic model with PLINK 2 (`-glm`). We include the top ten principal components (PCs), age, sex, and the genotyping array as covariables in the model. The results were depicted using the qqman package in R.

### PRS analysis

The Batch Screening Iterative LASSO (BASIL) algorithm [50] is a meta-algorithm (algorithms that learn from the output of other algorithms), which employs a Lasso algorithm [51] and enhances this output with another layer for faster variable selection in ultra-high-dimensional problems. Similar to the Lasso algorithm, the purpose of BASIL is to find a parameter vector  $\beta$  whose components are the coefficients for the independent variable of the linear regression that approximates the solution of the problem.

BASIL solves the Lasso solution path in an iterative fashion, starting with a sequence of candidate parameters. From these candidate solutions, each iteration discards the ones that do not meet the requirements to be a suitable solution. The set of variables who make it into the final set for a viable solution are those who were also screened satisfying a desired threshold requirement, while the others are discarded (i.e., those solutions in which the coefficients in their positions inside the  $\beta$  parameter are meant to be 0). This process is repeated until the optimum parameter  $\lambda_0$  is found, which is the one that minimizes ( $\lambda_0$ ). The BASIL algorithm guarantees to find the exact solution and not only an approximation, via the Karush–Kuhn–Tucker condition (the first derivative necessary conditions for a solution to be optimal) [52] which is verified along each iteration. This condition is necessary and sufficient to prove it.

### Genetic ancestry

To address the confounding factor of population stratification in PRS estimations, various approaches have been taken [53, 54]. Here, we follow the convention of using the first 10 principal components of the PCA to the adjustment of the GWA study. For the correction of PRS models, we make use of estimates of global ancestry. For this purpose, we use Neural ADMIXTURE [55], a faster adaptation of the ADMIXTURE algorithm [56] with similar (or better) clustering results. Utilizing the Python implementation of Neural ADMIXTURE, we use data from the 1000 Genomes Project Consortium [57] for training a model in the supervised mode of Neural ADMIXTURE with the default parameters. We utilized the results of global ancestry inference as a covariate in the training of our PRS models.

### Statistical analysis

We first used descriptive statistics to understand the characteristics of this cohort, including the geographic distribution of participants, age, sex, and comorbidities. This can become useful in future replication studies. We built three PRS models using BASIL, a lasso-based linear model: a) one model included genotype features alone; b) another model used only the covariates of age, sex, and the first ten genetic principal components (PCs); c) and a final (full) model used both genotypes and covariates. We used the first 10 PCs to account for residual population micro-stratification as fixed effects. With the three models, we then calculate their predictive performances, as well as the performance of the full model against the covariate-only model (delta between models).

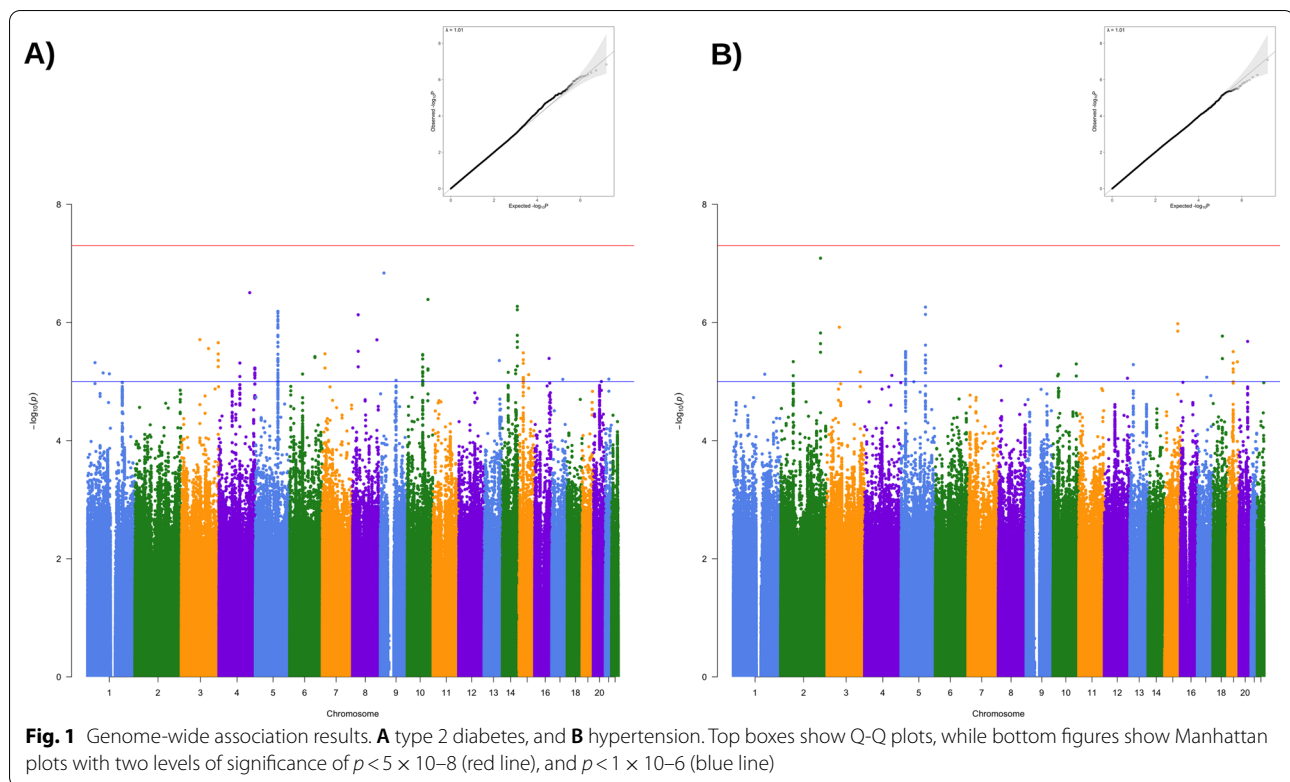
The predictive ability of these PRS models was evaluated using the area under the curve (AUC) receiver operating characteristic (ROC) curves using the pROC package in R [58]. To make comparisons between AUC curves from each model, we used the nonparametric method developed by DeLong et al. [59], which is a commonly used method.

## Results

### Genome-wide association results for diabetes type 2 and hypertension

We combined genome-wide association data for: 1) 389 T2D cases and 4,161 controls ( $N = 4,550$ ); and 2) 1,027 hypertension cases and 3,501 controls of European ancestry ( $N = 4,528$ ). We tested ~8 million variants for T2D and hypertension passing quality control and imputation filters ( $MAF > 0.01$ ,  $R^2 > 0.3$ ). Both results showed low inflation of test statistics ( $\lambda_{GC} = 1.01$ ;  $\lambda_{GC} = 1.01$ ) (Fig. 1).





Nineteen T2D variants displayed significant evidence of replication ( $p < 0.05$ ) in this dataset. Among them, we identified variants closely associated with genes which have been previously linked to type 2 diabetes susceptibility (e.g., *CDKALI*, *KCNQ1*), as well as variants in the *FTO* locus linked previously with both BMI and T2D. Overall, we identified 164 out of 272 variants showing identical effect direction to previously reported genome-significant findings in Europeans (Additional file 3) [35]. For the hypertension dataset, we replicated ten hypertension genetic markers and identified 230 out of 365 variants having identical effect direction [60] (Additional file 3; Fig. 1).

We validated our GWAS using independent GWAS meta-analysis datasets from Mahajan et al. 2018 (74,124 T2D cases, 824,006 controls) [35] and Evangelou et al. 2018 (757,201 individuals) [60]. We compared the  $p$ -values and the effect sizes for the variants assessed in both our studies that had identical chromosomal coordinates and alleles with the independent GWAS. The direction of the effect sizes (estimated as OR) was set to match the effect alleles in each study. We observed that the effect sizes of the genome-wide significant variants in the independent GWAS [35, 60] were concordant in directionality in both our T2D and hypertension GWAS. (Effect sizes had the same direction across both studies, Additional file 3.)

Our observations highlight how carefully curated DTC repositories with ever increasing sample sizes and variant diversity can replicate previous findings and hold the potential of delivering enhanced discovery and single-variant resolution of causal T2D and hypertension risk and protection alleles. Additionally, our findings confirm the potential impact of DTC resources on mechanistic insights and clinical translation efforts.

#### Estimation of cardiometabolic PRS models using SNPnet

PRS models were built for each phenotype using the BASIL algorithm [50]. The predictor variable was binary (presence or absence of diabetes or hypertension) as reported by participants. After tenfold cross-validation, the genotype-only models reported a predictive performance (AUC) of 0.56 for both diabetes and hypertension and increased to 0.68 for the full model (genotype and covariates together). Similarly, when filtering for participants of European ancestry, the genotype-only models reported a predictive performance of 0.57 and 0.53, increasing to 0.69 and 0.66 in the full model, respectively. We compared the performance of these models using DeLong's method [59] with no statistical significance (See Additional file 4 for additional metrics). The imbalance ratio of both cohorts was an important factor that impacted the accuracy of these models. Additionally, the majority proportion of participants of mostly European

ancestry also explains the small differences in performance between both types of models. Figure 2 shows the comparison between AUC curves in both phenotypes.

The European Bioinformatics Institute (EBI) developed the PGS Catalog [31], which is an open resource of published polygenic scores (including variants, alleles, and weights). We investigated those published PRS for T2D and hypertension. For those with reported AUC, we also obtained the number of variants, number of individuals whose data was used to train the model under various ancestry groups. See Additional file 5 for more information on the PGS Catalog reported scores.

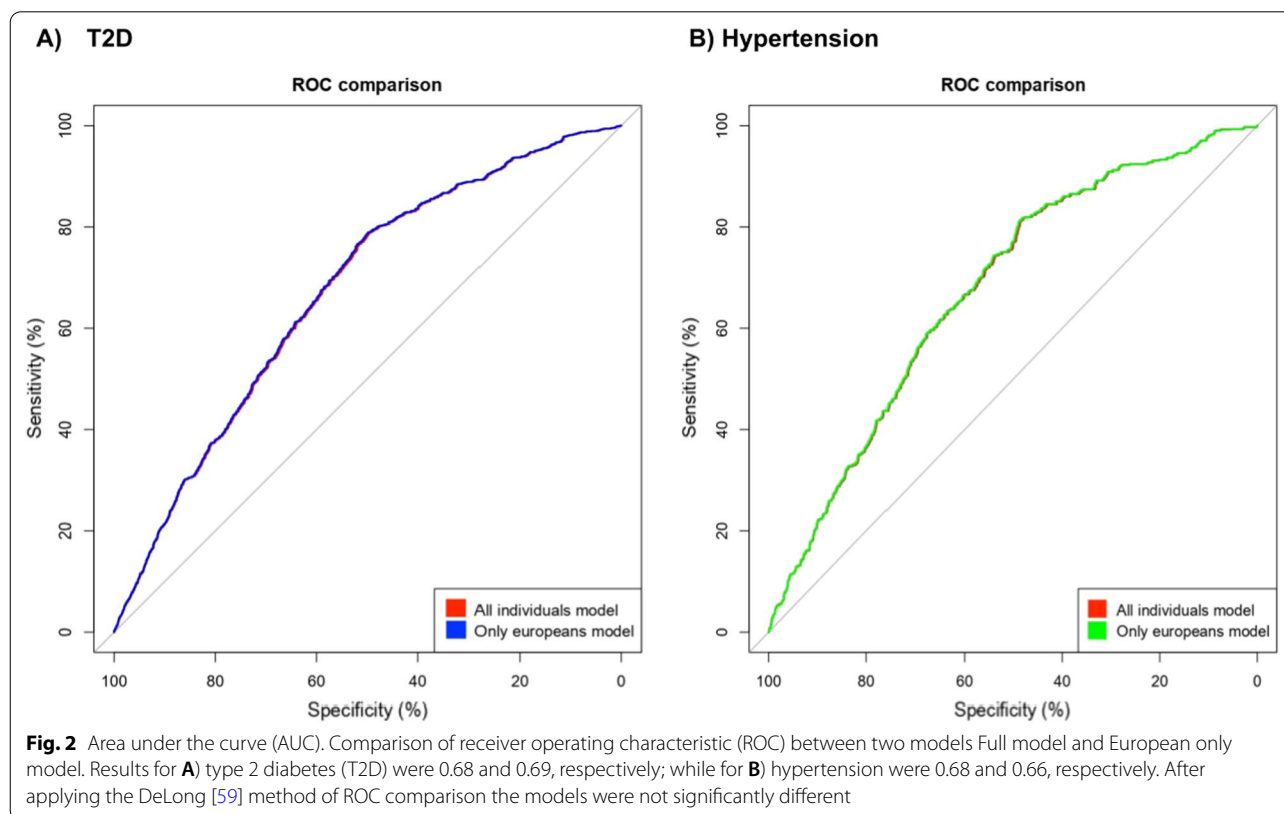
Our models are comparable to previously published PRS models. Figure 3 shows the comparison in reported AUC for those models in the EBI PGS catalog including our own models. The average AUC between these models was 0.70. The small number of variables used by our models (125 for T2D and 666 for hypertension) makes them comparable to those reported by Tanigawa et al. [43] who also used the BASIL algorithm. Likewise, the number of individuals whose data was used to train the models is modest in comparison with large academic and clinical databases. Nevertheless, the predictive performance does not seem to be overtly affected by the number of individuals in the study or the number of included variants highlighting that genetic array data from DTC

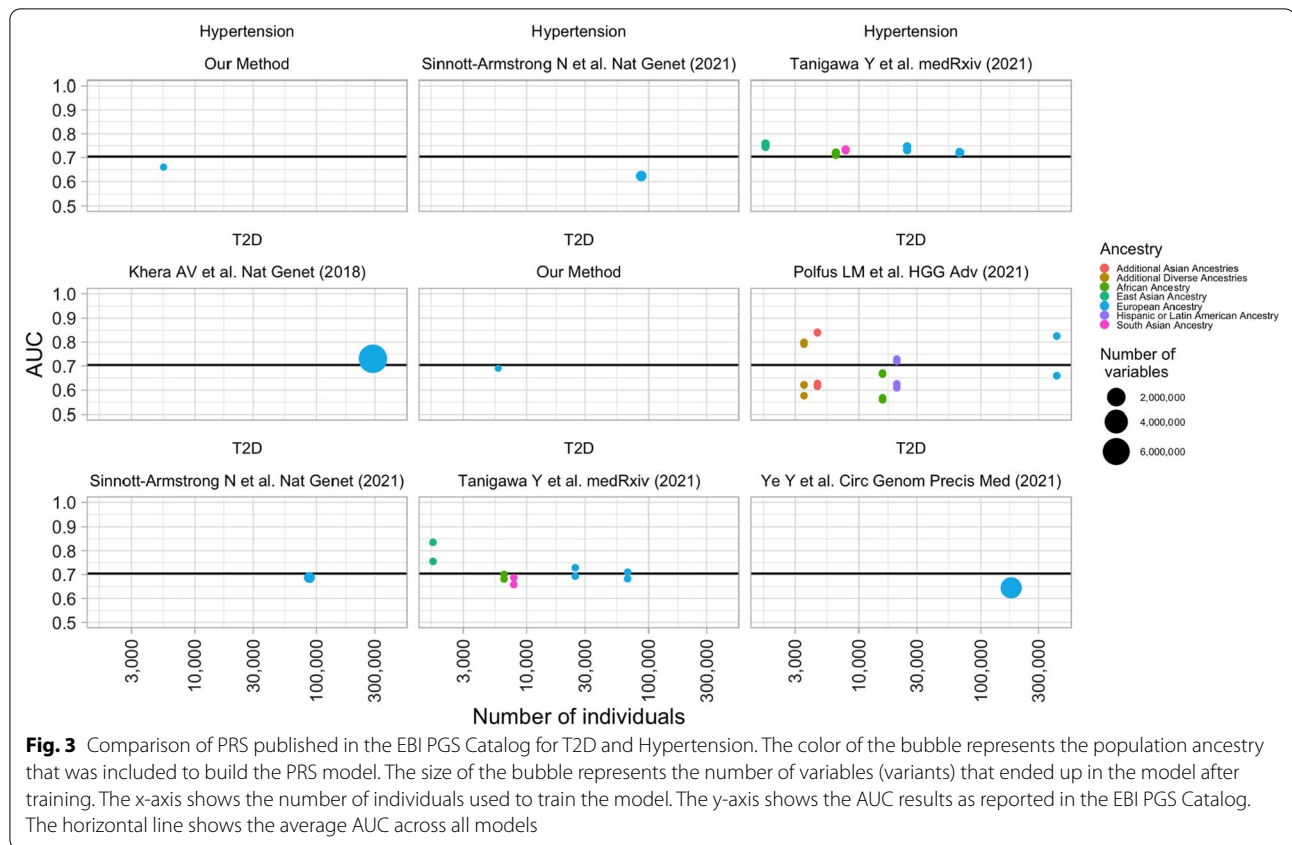
repositories carry immense promise for the development of PRS tools aimed at improving early detection and prevention of T2D and hypertension.

## Discussion

In this study we generated PRS for T2D and hypertension from a heterogeneous dataset housing a combination of genetic data and self-reported information from a DTC genetics company. See Additional file 6, for more information about this cohort. Despite a relatively modest predictive ability our PRS models are able to identify subsets of users at substantially increased risk of presenting T2D or hypertension. This finding is remarkable because it suggests that the ever-increasing availability of genetic data from DTC providers, most of it not annotated for traits of clinical relevance, can be leveraged to generate predictive tools able to improve diagnosis and prevention of diseases with genetic determinants.

Our study tested the possibility of inverting the model regarding genetic and phenotypic data acquisition in a research study. Individuals participating in our study shared their genetic array information from other DTC providers and were invited to take an online survey regarding their general health condition. We found no difference in predictive performance between our trained models that included respondents from all inferred





ancestries and those models with respondents from European heritage, due to the fact that 86% of our database was of European origin. The genetics of our PRS models for T2D and hypertension are supported by our ability to replicate known variants from publicly available independent GWAS studies.

Multiple array types were available in our database, and imputation across platforms (up-imputing) was necessary to harmonize these diverse datasets. The fact that individuals can self-report their genomic information could potentially corrupt the file being uploaded into our platform. However, applying appropriate quality control (QC) principles proved to successfully enable traditional GWAS and PRS analyses.

DTC platforms can offer a wide range of information about personal wellness, ancestry, physical characteristics, and traits. Advances in genomic research have led the DTC genomics industry to flourish and make accurate yet easy to interpret genomic results. Strict privacy policies of many companies disallow them to share customers' data without their consent. These platforms can serve as informative repositories giving actionable insights that aid traditional clinical approaches. The approach of subject recruitment for various complex phenotypes via online surveys is opening up multiple

avenues to complement conventional research and clinical strategies. DTC platforms also provide convenience along with a wider reach to recruit participants from various locations. They surpass barriers of single-point data collection centers to language restrictions thus allowing the aggregation of data from places with different ancestries and demographics. Democratizing the access to these genetic platforms and prediction tools will likely boost progress in precision medicine. In the future, we plan to investigate how federated learning approaches can further improve the possibility to increase the power of studies in DTC genomic analysis, but also how meta-analysis can be done in combination with academic and clinical datasets (including those from large consortiums).

We have shown that our DTC platform and research strategy has the potential to replicate the previously reported results with a very fast turnaround time. The participation of individual customers in our platform allowed the generation of a rich dataset that enabled the creation of PRS cardiometabolic models. The comparable predictive performance of our models also is a great indication of how we can quickly contribute more PRS models to the larger scientific community. As it stands right now, publicly available PRS models for T2D and

hypertension have an AUC of 0.7 on average as shown in Fig. 3. This is still a low accuracy, and it is even lower when compared to the small difference between the full and genotype-only models. However, even in our heterogeneous DTC platform, we have been able to replicate the findings seen in academic and government-funded biobanks.

T2D and hypertension are multifactorial diseases that are impacted by genetic and environmental determinants, including lifestyle factors like nutrition and exercise habits. Nevertheless, the inherent limitation of PRS models to provide accurate disease predictions compels the need to interpret these findings with caution, especially when they come from DTC genetic services. The clinical actionability of PRS models has yet to be determined through pragmatic trials involving real-world data. We hope to provide a novel source of information that can shed light on this important issue. Therefore, providing personalized information about T2D and hypertension predisposition is poised to improve early diagnosis and prevention bringing precision medicine at scale for all.

## Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s40246-022-00406-y>.

**Additional file 1.** Survey questionnaire.

**Additional file 2.** Number of variants and Genotyping Call Rate.

**Additional file 3.** Summary statistics for known D2T GWS signals in our dataset.

**Additional file 4.** Predictive performance results by model.

**Additional file 5.** PGS Catalog models for T2D and Hypertension.

**Additional file 6.** PGS Catalog models for T2D and Hypertension.

## Acknowledgements

We are grateful to all our participants for being a part of this study.

## Author contributions

ALP, MV, SMG, and ARM designed the study. SMG, MNA, GGA, ARM, and performed analysis of data. BM, MTML, KN, YM, AI, NK, TT, and CDB provided interpretation of the results. ALP, ARM, SMG, and CDB drafted the manuscript, and all authors contributed critically, read, revised and approved the final version.

## Funding

This research is based on results obtained from a project, JPNP19001, subsidized by the New Energy and Industrial Technology Development Organization (NEDO).

## Availability of data and materials

The data that supports the findings of this study is available for qualified researchers at non-profit institutions upon entering into an agreement with Genomelink. All information will be shared subject to the above criterion upon request via [info@genomelink.io](mailto:info@genomelink.io).

## Declarations

### Ethics approval and consent to participate

This study was approved by the institutional review board (IRB) at WCG IRB (<https://www.wcgirb.com/>) under IRB tracking number protocol number 20201332.

### Competing interests

ARM, SMG, MTML, ALP, CDB, AI, MNA, and NK are employees of or consultants to Galatea Bio. MV, KN, YM, and TT are employees of Genomelink. CDB, IA, and NK are shareholders of Galatea Bio stock. CDB, KN, YM, and TT are shareholders of Genomelink stock. The remaining authors declare that there is no conflict of interest regarding the publication of this article.

### Author details

<sup>1</sup>Galatea Bio, Inc., 975 W 22nd Street, Hialeah, Florida 33010, USA. <sup>2</sup>Amphora Health, Batallon Independencia 80, Morelia, Michoacan 58260, Mexico.

<sup>3</sup>Genomelink, Inc., 2150 Shattuck Avenue, Berkeley, California 94704, USA.

<sup>4</sup>Awakens Japan K.K., 2-11-3 Meguro, Meguro-ku, Tokyo 1530063, Japan.

<sup>5</sup>Department of Biomedical Data Science, Stanford University School of Medicine, 1265 Welch Road, Stanford, California 94305, USA. <sup>6</sup>Chan Zuckerberg Biohub, 499 Illinois Street, San Francisco, California 94158, USA.

Received: 8 March 2022 Accepted: 6 August 2022

Published online: 08 September 2022

## References

- Mathers CD, Loncar D. Projections of global mortality and burden of disease from 2002 to 2030. *PLoS Med*. 2006;3(11):442–520.
- Psaltopoulou T, Ilias I, Alevizaki M. The role of diet and lifestyle in primary, secondary, and tertiary diabetes prevention: a review of meta-analyses. *Rev Diabet Stud*. 2010;7(1):26–35.
- Cousin E, Duncan BB, Stein C, Ong KL, Vos T, Abbafati C, Haque S. Diabetes mortality and trends before 25 years of age: an analysis of the Global Burden of Disease Study 2019. *Lancet Diabetes Endocrinol*. 2022. [https://doi.org/10.1016/S2213-8587\(21\)00349-1](https://doi.org/10.1016/S2213-8587(21)00349-1).
- World Health Organization. (2022a). Diabetes. World Health Organization. Retrieved February 15, 2022, from <https://www.who.int/news-room/fact-sheets/detail/diabetes>
- Tsimihodimos V, Gonzalez-Villalpando C, Meigs JB, Ferrannini E. Hypertension and diabetes mellitus: coprediction and time trajectories. *Hypertension*. 2018;71(3):422–8.
- World Health Organization. (2022b). Hypertension. World Health Organization. Retrieved February 15, 2022, from <https://www.who.int/news-room/fact-sheets/detail/hypertension>
- Yamazaki D, Hitomi H, Nishiyama A. Hypertension with diabetes mellitus complications. *Hypertens Res*. 2018;41(3):147–56.
- Wang C, Du Z, Ye N, Shi C, Liu S, Geng D, Sun Y. Hyperlipidemia and hypertension have synergistic interaction on ischemic stroke: insights from a general population survey in China. *BMC Cardiovasc Disord*. 2022;22(1):47. <https://doi.org/10.1186/s12872-022-02491-2>.
- Buscemi S, Corleo D, Randazzo C. Risk Factors for COVID-19: diabetes, hypertension, and obesity. In: *Coronavirus therapeutics—volume II*. Springer, Cham; 2021. pp. 115–129.
- Sanghera DK, Blackett PR. Type 2 diabetes genetics: beyond GWAS. *J Diabetes Metab*. 2012. <https://doi.org/10.4172/2155-6156.1000198>.
- Hindy G, Dornbos P, Chaffin MD, Liu DJ, Wang M, Selvaraj MS, So WY. Rare coding variants in 35 genes associate with circulating lipid levels—A multi-ancestry analysis of 170,000 exomes. *Am J Hum Genet*. 2022;109(1):81–96.
- Patxot M, Banos DT, Kousathanas A, Orlic EJ, Ojavee SE, Moser G, Robinson MR. Probabilistic inference of the genetic architecture underlying functional enrichment of complex traits. *Nat Commun*. 2021;12(1):1–16.
- Rusu V, Hoch E, Mercader JM, Tenen DE, Gymrek M, Hartigan CR, Lander ES. Type 2 diabetes variants disrupt function of SLC16A11 through two distinct mechanisms. *Cell*. 2017;170(1):199–212.
- Burns SM, Vetere A, Walpita D, Dančik V, Khodier C, Perez J, Altshuler D. High-throughput luminescent reporter of insulin secretion for



- discovering regulators of pancreatic Beta-cell function. *Cell Metab.* 2015;21(1):126–37.
15. Dai N, Zhao L, Wrighting D, Krämer D, Majithia A, Wang Y, Avruch J. IGF2BP2/IMP2-deficient mice resist obesity through enhanced translation of Ucp1 mRNA and other mRNAs encoding mitochondrial proteins. *Cell Metab.* 2015;21(4):609–21.
  16. Mahajan A, Taliun D, Thurner M, Robertson NR, Torres JM, Rayner NW, McCarthy ML. Fine-mapping type 2 diabetes loci to single-variant resolution using high-density imputation and islet-specific epigenome maps. *Nat Genet.* 2018;50(11):1505–13.
  17. Williams AL, Jacobs SB, Moreno-Macias H, Huerta-Chagoya A, Churchhouse C, Márquez-Luna C, Altshuler D. Sequence variants in SLC16A11 are a common risk factor for type 2 diabetes in Mexico. *Nature.* 2014;506(7486):97.
  18. Estrada K, Aukrust I, Bjørkhaug L, Burt NP, Mercader JM, Garcia-Ortiz H, SIGMA Type 2 Diabetes Consortium. Association of a low-frequency variant in HNF1A with type 2 diabetes in a Latino population. *JAMA.* 2014;311(22):2305–14.
  19. Wojcik GL, Graff M, Nishimura KK, Tao R, Haessler J, Gignoux CR, Carlson CS. Genetic analyses of diverse populations improves discovery for complex traits. *Nature.* 2019;570(7762):514–8.
  20. Said MA, Verweij N, van der Harst P. Associations of combined genetic and lifestyle risks with incident cardiovascular disease and diabetes in the UK Biobank Study. *JAMA Cardiol.* 2018;3(8):693–702.
  21. Kowalski MH, Qian H, Hou Z, Rosen JD, Tapia AL, Shan Y, Li Y. Use of >100,000 NHLBI trans-omics for precision medicine (TOPMed) Consortium whole genome sequences improves imputation quality and detection of rare variant associations in admixed African and Hispanic/Latino populations. *PLoS Genet.* 2019;15(12): e1008500.
  22. Zhou W, Global Biobank Meta-analysis Initiative. In: Global Biobank Meta-analysis Initiative: Powering genetic discovery across human diseases; 2021. medRxiv.
  23. Eriksson N, Macpherson JM, Tung JY, Hon LS, Naughton B, Saxonov S, Mountain J. Web-based, participant-driven studies yield novel genetic associations for common traits. *PLoS Genet.* 2010;6(6): e1000993.
  24. Becker J. et al. Resource profile and user guide of the Polygenic Index Repository. *Nat Hum Behav.* 5.12 ; 2021: 1744–1758.
  25. Food and Drug Administration. FDA allows marketing of first direct-to-consumer tests that provide genetic risk information for certain conditions. U.S. Food and Drug Administration.; 2021. Retrieved February 16, 2022, from <https://www.fda.gov/news-events/press-announcements/fda-allows-marketing-first-direct-consumer-tests-provide-genetic-risk-information-certain-conditions>
  26. Singal G, Miller PG, Agarwala V, Li G, Kaushik G, Backenroth D, Miller VA. Association of patient characteristics and tumor genomics with clinical outcomes among patients with non-small cell lung cancer using a clinicogenomic database. *JAMA.* 2019;321(14):1391–9.
  27. Lee JK, Madison R, Classon A, Gjoerup O, Rosenzweig M, Frampton GM, Schrock AB. Characterization of non-small-cell lung cancers with MET Exon 14 skipping alterations detected in tissue or liquid: clinicogenomics and real-world treatment patterns. *JCO Precis Oncol.* 2021;5:1354–76.
  28. Turner S, Chia S, Kanakamedala H, Hsu WC, Park J, Chandiwana D, Rugo HS. Effectiveness of alpelisib+ fulvestrant compared with real-world standard treatment among patients with HR+, HER2-, PIK3CA-mutated breast cancer. *Oncologist.* 2021;26(7):e1133–42.
  29. Liu R, Rizzo S, Whipple S, Pal N, Pineda AL, Lu M, Zou J. Evaluating eligibility criteria of oncology trials using real-world data and AI. *Nature.* 2021;592(7855):629–33.
  30. Francisco M, Bustamante CD. Polygenic risk scores: a biased prediction? *Genome Med.* 2018;10(1):1–3.
  31. Choi SW, Mak TSH, O'Reilly PF. Tutorial: a guide to performing polygenic risk score analyses. *Nat Protoc.* 2020;15(9):2759–72.
  32. Lambert SA, Gil L, Jupp S, Ritchie SC, Xu Y, Buniello A, Inouye M. The polygenic score catalog as an open database for reproducibility and systematic evaluation. *Nat Genet.* 2021;53(4):420–5.
  33. Av K, Chaffin M, Aragam KG, Me H, Roselli C, Choi SH, Natarajan P, Lander ES, Lubitz SA, Ellinor PT, Kathiresan S. Genome-wide polygenic scores for common diseases identify individuals with risk equivalent to monogenic mutations. *Nat Genet.* 2018;50:1219–24.
  34. Läll K, Mägi R, Morris A, Metspalu A, Fischer K. Personalized risk prediction for type 2 diabetes: the potential of genetic risk scores. *Genet Med.* 2017;19(3):322–9.
  35. Vassy JL, Hivert MF, Porneala B, Dauriz M, Florez JC, Dupuis J, Meigs JB. Polygenic type 2 diabetes prediction at the limit of common variant detection. *Diabetes.* 2014;63(6):2172–82.
  36. Qi Q, Stilp AM, Sofer T, Moon JY, Hidalgo B, Szpiro AA, Kaplan RC. Genetics of type 2 diabetes in US Hispanic/Latino individuals: results from the Hispanic Community Health Study/Study of Latinos (HCHS/SOL). *Diabetes.* 2017;66(5):1419–25.
  37. Mars N, Koskela JT, Ripatti P, Kiiskinen TT, Havulinna AS, Lindbohm JV, Ripatti S. Polygenic and clinical risk scores and their impact on age at onset and prediction of cardiometabolic diseases and common cancers. *Nat Med.* 2020;26(4):549–57.
  38. Sinnott-Armstrong N, Tanigawa Y, Amar D, Mars N, Benner C, Aguirre M, Rivas MA. Genetics of 35 blood and urine biomarkers in the UK Biobank. *Nat Genet.* 2021;53(2):185–94.
  39. Ritchie SC, Lambert SA, Arnold M, Teo SM, Lim S, Scepanovic P, Inouye M. Integrative analysis of the plasma proteome and polygenic risk of cardiometabolic diseases. *Nat Metab.* 2021;3(11):1476–83.
  40. Polfus LM, Darst BF, Highland H, Sheng X, Ng MC, Below JE, DIAMANTE Hispanic/Latino Consortium. Genetic discovery and risk characterization in type 2 diabetes across diverse populations. *Hum Genet Genom Adv.* 2021;2(2):100029.
  41. Mansour Aly D, Dwivedi OP, Prasad RB, Kärjämäki A, Hjort R, Thangam M, Ahlqvist E. Genome-wide association analyses highlight etiological differences underlying newly defined subtypes of diabetes. *Nat Genet.* 2021;53(11):1534–42.
  42. Aksit MA, Pace RG, Vecchio-Pagán B, Ling H, Rommens JM, Boelle PY, Blackman SM. Genetic modifiers of cystic fibrosis-related diabetes have extensive overlap with type 2 diabetes and related traits. *J Clin Endocrinol Metab.* 2020;105(5):1401–15.
  43. Tanigawa Y, Qian J, Venkataraman GR, Justesen JM, Li R, Tibshirani R, Rivas MA. Significant Sparse Polygenic Risk Scores across 428 traits in UK Biobank; 2021. medRxiv.
  44. Ye Y, Chen X, Han J, Jiang W, Natarajan P, Zhao H. Interactions between enhanced polygenic risk scores and lifestyle for cardiovascular disease, diabetes, and lipid levels. *Circ Genom Precis Med.* 2021;14(1):003128.
  45. Privé F, Aschard H, Carmi S, Folkersen L, Hoggart C, O'Reilly PF, Vilhjálmsson BJ. Portability of 245 polygenic scores when derived from the UK Biobank and applied to 9 ancestry groups from the same cohort. *Am J Hum Genet.* 2022;109(1):12–23.
  46. Popejoy AB, Fullerton SM. Genomics is failing on diversity. *Nature News.* 2016;538(7624):161.
  47. Hindorf LA, Bonham VL, Brody LC, Ginoza ME, Hutter CM, Manolio TA, Green ED. Prioritizing diversity in human genomics research. *Nat Rev Genet.* 2018;19(3):175.
  48. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, Sham PC. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet.* 2007;81(3):559–75.
  49. Ayres DL, Darling A, Zwickl DJ, Beerli P, Holder MT, Lewis PO, Suchard MA. BEAGLE: an application programming interface and high-performance computing library for statistical phylogenetics. *Syst Biol.* 2012;61(1):170–3.
  50. Qian J, Tanigawa Y, Du W, Aguirre M, Chang C, Tibshirani R, Hastie T. A fast and scalable framework for large-scale and ultrahigh-dimensional sparse regression with application to the UK Biobank. *PLoS Genet.* 2020;16(10): e1009141.
  51. Tibshirani R. Regression shrinkage and selection via the lasso. *J Roy Stat Soc Ser B Methodol.* 1996;58(1):267–88.
  52. Boyd S, Boyd SP, Vandenberghe L. Convex optimization. Cambridge University Press; 2004.
  53. Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. Principal components analysis corrects for stratification in genome-wide association studies. *Nat Genet.* 2006;38(8):904–9. <https://doi.org/10.1038/ng1847>.
  54. Price AL, Zaitlen NA, Reich D, Patterson N. New approaches to population stratification in genome-wide association studies. *Nat Rev Genet.* 2010;11(7):459–63.
  55. Mantes AD, Montserrat DM, Bustamante CD, Giró-i-Nieto X, Ioannidis AG. Neural ADMIXTURE: rapid population clustering with autoencoders; 2021. bioRxiv.

56. Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in unrelated individuals. *Genome Res.* 2009;19(9):1655–64. <https://doi.org/10.1101/gr.094052.109>.
57. Genomes Project Consortium. A global reference for human genetic variation. *Nature.* 2015;526(7571):68.
58. Robin X, Turck N, Hainard A, Tiberti N, Lisacek F, Sanchez JC, Müller M. pROC: an open-source package for R and S+ to analyze and compare ROC curves. *BMC Bioinform.* 2011;12(1):1–8.
59. DeLong ER, DeLong DM, Clarke-Pearson DL. Comparing the areas under two or more correlated receiver operating characteristic curves: a non-parametric approach. *Biometrics.* 1988. <https://doi.org/10.2307/2531595>.
60. Evangelou E, Warren HR, et al. Genetic analysis of over 1 million people identifies 535 new loci associated with blood pressure traits. *Nat Genet.* 2018;50(10):1412–25.

### Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

