



# Impute.me: An Open-Source, Non-profit Tool for Using Data From Direct-to-Consumer Genetic Testing to Calculate and Interpret Polygenic Risk Scores

Lasse Folkersen<sup>1\*</sup>, Oliver Pain<sup>2</sup>, Andrés Ingason<sup>1</sup>, Thomas Werge<sup>1†</sup>, Cathryn M. Lewis<sup>2,3†</sup> and Jehannine Austin<sup>4,5†</sup>

<sup>1</sup> Institute of Biological Psychiatry, Mental Health Centre Sankt Hans, Copenhagen, Denmark, <sup>2</sup> Social, Genetic and Developmental Psychiatry Centre, Institute of Psychiatry, Psychology & Neuroscience, King's College London, London, United Kingdom, <sup>3</sup> Department of Medical & Molecular Genetics, Faculty of Life Sciences & Medicine, King's College London, London, United Kingdom, <sup>4</sup> Department of Psychiatry, University of British Columbia, Vancouver, BC, Canada, <sup>5</sup> Department of Medical Genetics, University of British Columbia, Vancouver, BC, Canada

## OPEN ACCESS

### Edited by:

Gustavo Glusman,  
Institute for Systems Biology (ISB),  
United States

### Reviewed by:

Hervé Michel Chneiweiss,  
Centre National de la Recherche  
Scientifique (CNRS), France  
S. Hong Lee,  
University of South Australia, Australia

### \*Correspondence:

Lasse Folkersen  
lasse.folkersen@regionh.dk;  
lassefolkersen@gmail.com

† These authors have contributed  
equally to this work

### Specialty section:

This article was submitted to  
ELSI in Science and Genetics,  
a section of the journal  
Frontiers in Genetics

**Received:** 29 November 2019

**Accepted:** 11 May 2020

**Published:** 30 June 2020

### Citation:

Folkersen L, Pain O, Ingason A,  
Werge T, Lewis CM and Austin J  
(2020) Impute.me: An Open-Source,  
Non-profit Tool for Using Data From  
Direct-to-Consumer Genetic Testing  
to Calculate and Interpret Polygenic  
Risk Scores. *Front. Genet.* 11:578.  
doi: 10.3389/fgene.2020.00578

To date, interpretation of genomic information has focused on single variants conferring disease risk, but most disorders of major public concern have a polygenic architecture. Polygenic risk scores (PRSs) give a single measure of disease liability by summarizing disease risk across hundreds of thousands of genetic variants. They can be calculated in any genome-wide genotype data-source, using a prediction model based on genome-wide summary statistics from external studies. As genome-wide association studies increase in power, the predictive ability for disease risk will also increase. Although PRSs are unlikely ever to be fully diagnostic, they may give valuable medical information for risk stratification, prognosis, or treatment response prediction. Public engagement is therefore becoming important on the potential use and acceptability of PRSs. However, the current public perception of genetics is that it provides “yes/no” answers about the presence/absence of a condition, or the potential for developing a condition, which is not the case for common, complex disorders with polygenic architecture. Meanwhile, unregulated third-party applications are being developed to satisfy consumer demand for information on the impact of lower-risk variants on common diseases that are highly polygenic. Often, applications report results from single-nucleotide polymorphisms (SNPs) and disregard effect size, which is highly inappropriate for common, complex disorders where everybody carries risk variants. Tools are therefore needed to communicate our understanding of genetic vulnerability as a continuous trait, where a genetic liability confers risk for disease. Impute.me is one such tool, whose focus is on education and information on common, complex disorders with polygenic architecture. Its research-focused open-source website allows users to upload consumer genetics data to obtain PRSs, with results reported on a population-level normal distribution. Diseases can only be browsed by *International Classification of Diseases*, 10th Revision (ICD-10) chapter–location or alphabetically, thus prompting the

user to consider genetic risk scores in a medical context of relevance to the individual. Here, we present an overview of the implementation of the impute.me site, along with analysis of typical usage patterns, which may advance public perception of genomic risk and precision medicine.

**Keywords:** genetics, polygenic risk scores, direct-to-consumer, personal genomes, risk prediction

## INTRODUCTION

In clinical genetics, testing for rare strong-effect causal variants is routinely performed in the health-care system to confirm a diagnosis or to evaluate individual risk suspected from anamnestic information (Baig et al., 2016), and in such instances, the use of genome sequencing is expanding (Byrjalsen et al., 2018). Meanwhile, outside of the health-care system, direct-to-consumer (DTC) genetics expands rapidly, providing the public with access to individual genetic data profiles and to interpretation of common genetic variants derived from genotyping microarrays (Kaye, 2008; Greshake et al., 2014). This is developing as a sprawling industry of consumer services with widely diverging standards, including third-party genome analysis services. These services typically provide individual results from analysis of common single-nucleotide polymorphisms (SNPs) with (at best) weak effects. They are therefore severely mis-aligned with current state-of-the-art, which at least for common, complex disease is to use polygenic risk scores (PRSs) to estimate the combined risk of common variation in the genome (Lee et al., 2008; Lewis and Vassos, 2017).

We believe that the goal of the academic genetics community should extend beyond theory. This means engaging with the public and assisting those who seek information, even when it means helping them to interpret their own genomic data. We therefore developed impute.me as an online web-app for analysis and education in personal genetic analysis. The web-app is illustrated in **Figure 1**. Using any major DTC vendor, a user can download their raw data and then upload it at impute.me. Uploaded files are checked and formatted according to procedures that have been developed to handle most types of microarray-based consumer genetics data, including an imputation step. These data are then further subjected to automated analysis scripts including PRS calculations. This includes more than 2,000 traits, browsable in different interface types (modules). Each module is designed with the goal of putting findings in as relevant a context as possible, prompting users to see common variant genetics as a support tool rather than a diagnosis finder. The aim is to provide information as broadly as possible to offer a real alternative to the widespread practice of reporting on weak SNP genotypes for any trait, even though that means generation of reports that are below any sensible threshold for clinical usability. We hope that having this as an open and accessible resource for everyone will be of help to the debate on what exactly constitutes clinical usability beyond high-risk pathogenic variants.

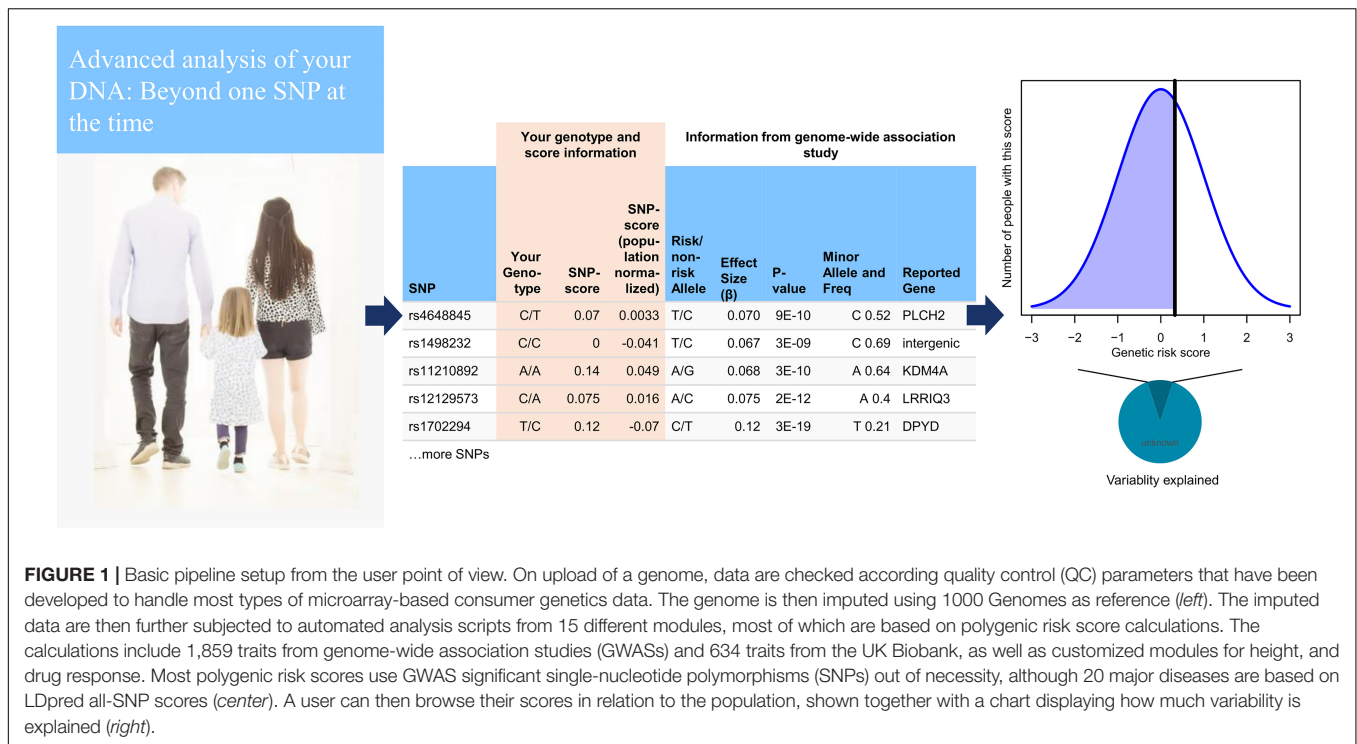
In this article, we will describe the (i) development and setup, (ii) validation and testing, (iii) evaluation of usage, (iv) communication of risk scores, and (v) ethics and implications.

In the section *Development and Setup*, we discuss some of the challenges faced when developing a full personal-genome scoring pipeline. The goal of this section is to motivate and explain the choices made in development. In the second section, *Validation and Testing*, we use public Biobank data from individuals who consented for genetic research to test the effect of the impute.me scores on known disease outcomes. The purpose of this section is to test and validate scores, as well as to investigate consequences of some of the challenges that were raised in the first section. In the third section, *Evaluation of Usage*, we evaluate usage metrics of impute.me users. The goal of this section is to shed light on behavioral patterns of individuals who use DTC genetics for health questions and to offer recommendations that may be of use in other personal-genome scoring pipelines. In the section *Communication of Risk Scores*, we discuss our views on future directions particularly with respect to improving how genetic findings are presented to people. Finally, in *Ethics and Implications*, we discuss the ethics of providing access to health-related interpretation of DNA data.

## DEVELOPMENT AND SETUP

The first challenge in development of personal genomic services is standardization. As the name impute.me implies, all genotype data are processed by imputation of genotype data (Howie et al., 2009; Delaneau et al., 2013). This procedure expands the data available into ungenotyped SNPs and increases overlap with public genome-wide association study (GWAS) summary statistics used to estimate risk. It also expands the SNP overlap between microarray types from the major vendors, such as 23andMe, MyHeritage, and Ancestry.com. Further, we have found that imputation helps in avoiding major errors, for example, strand-flip issues that arise from the dozens of different data formats. Eliminating such problems from further processing is one important step to minimize mis-interpretation of genome analysis. To ensure high standard of reported results, impute.me requires a fully completed imputation for continued analysis.

The second challenge is to estimate PRSs that are accurate and robust to heterogenous data sources. This is particularly important to an application utilized by people from around the world leveraging data from dozens of different vendors and data types. Importantly, PRSs calculated from GWAS of a population of (for example) European ancestry will perform better for individuals of the same ancestry, and the systematic shift (i.e., bias) in risk scores in individuals from other populations is a problem (Curtis, 2018). Because studies of all disease traits are not yet available for all non-European populations, the pragmatic solution has been to include a population-specific normalization



attempting to minimize the systematic shifts of scores for non-European ancestry users. Further, it is computationally and logistically easier to implement PRSs that use only the most (i.e., genome-wide) significant SNPs (often referred to as top SNPs), but the prediction strength is better when more SNPs are included (all-SNP), which, however, is more sensitive to ancestry biases (Lam et al., 2019). The impute.me pipelines calculate PRSs for each trait or disease on the basis of all-SNP-based PRS calculations if full genome-wide summary statistics are available and processed, and top-SNP-based PRS calculations if not.

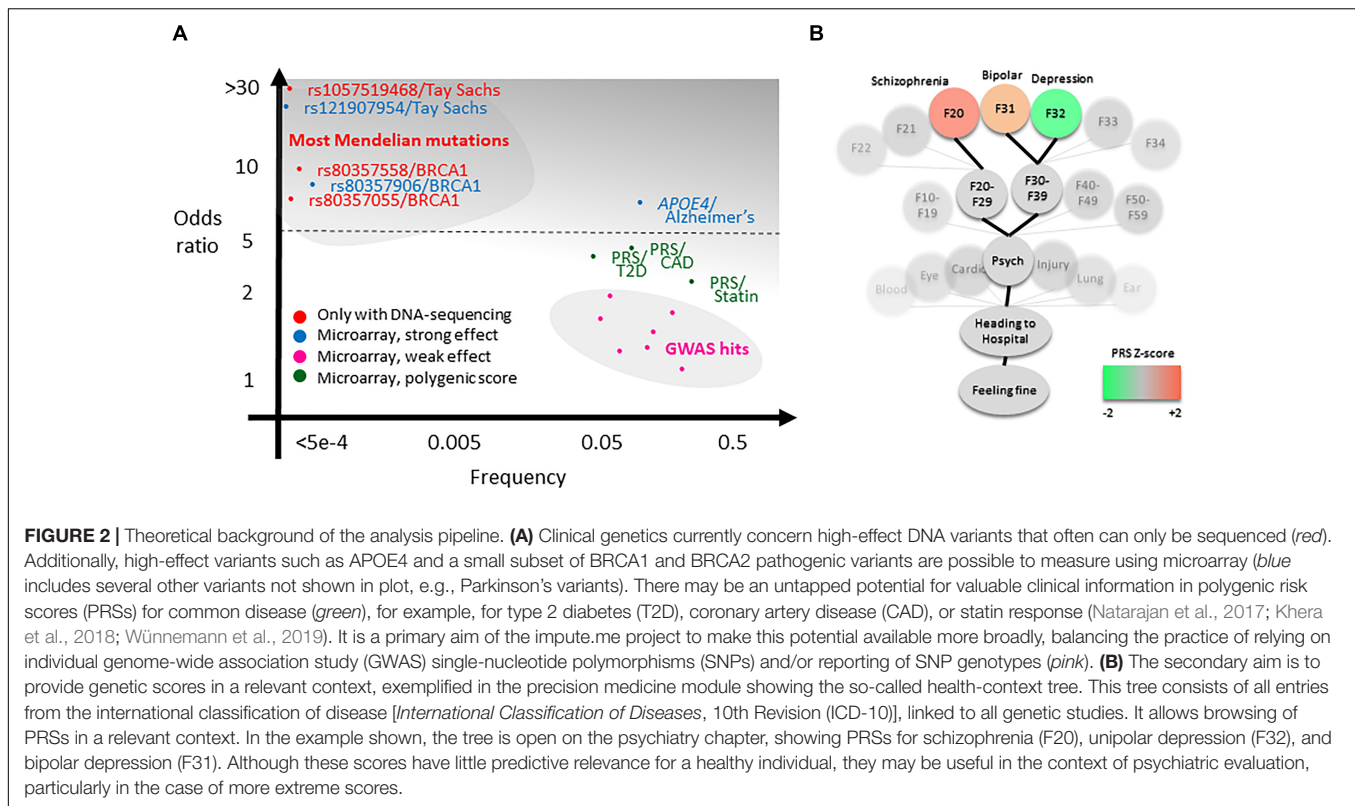
The third challenge is presentation. For a single rare large-effect variant, such as for the pathogenic variants in the *BRCA* genes conferring very high risk of cancers (odds ratio >10; **Figure 2A**, upper left), presentation focuses on absence versus presence (Maxwell et al., 2016). However, also, low-effect variants, for example, as in pharmacogenetics, impacting statin response, is considered as having potential clinical use (Natarajan et al., 2017; **Figure 2A**, lower right). This difference in effect magnitude is a major challenge in result presentation and understanding, particularly because a firm threshold is difficult to set: In the context of a drug-prescription situation or a question of which of two suspected disease risks is the most likely, it may be useful to know such scores. But in the context of an otherwise healthy individual, genetic risks are only relevant if we are very certain of them, they are serious, and preferably actionable [e.g., *BRCA* variants (Kalia et al., 2017)]. For this reason, we have made the design choice to avoid the use of lists sorted by risk score. Currently, scores are accessible through either an alphabetically sorted list or in a tree-like setup where genetic scores are reported in a health-context tree (**Figure 2B**). In this, all

scores are included, but scores that are less relevant to healthy individuals (i.e., most of them) are buried deeper into the health-context tree. As further discussed in the section *Future Challenges*, there are a lot of remaining challenges to solve in this question.

## VALIDATION AND TESTING

To evaluate pipelines on individuals with known disease outcomes, we investigated 242 samples from the CommonMind data set. The CommonMind data set includes patients with schizophrenia (SCZ), bipolar disorder, and controls, from European ancestry and from African ancestry. For each disorder and each ancestry group, the full impute.me pipelines were applied, including imputation and PRS calculation. Additionally, SNP sets corresponding to each of three major DTC companies were extracted and re-calculated. This was done to test the hypothesis that PRS calculation in mixed SNP sets poses particular challenges with regard to missing SNPs. Such sets of genotyped SNPs that are different in each sample are an unavoidable consequence of working with online data uploads.

We found that disease prediction strength, measured as variability explained, corresponded well to theoretical expectations of known SNP heritability (Lee et al., 2017; Li et al., 2017; Wünnemann et al., 2019). Secondly, we found that using all-SNP scores resulted in better prediction than top-SNP scores, which was as expected (Vilhjálmsón et al., 2015). Thirdly, we found that prediction was more accurate in individuals of European ancestry compared with individuals of African ancestry, which is concordant with the PRSs being developed



from European Ancestry GWAS (Schizophrenia Working Group of the Psychiatric Genomics Consortium, 2014; Hou et al., 2016; Lee et al., 2018). These observations match well with findings from studies of PRSs in much larger data sets. We caution that universally valid estimates of variability explained are better derived from larger studies that can consider the numerous issues such as balancing of cases and controls, realistic sampling conditions, and other inflations of effects. The intention here is to provide a specific test of impute-me pipelines and address DTC data-related questions.

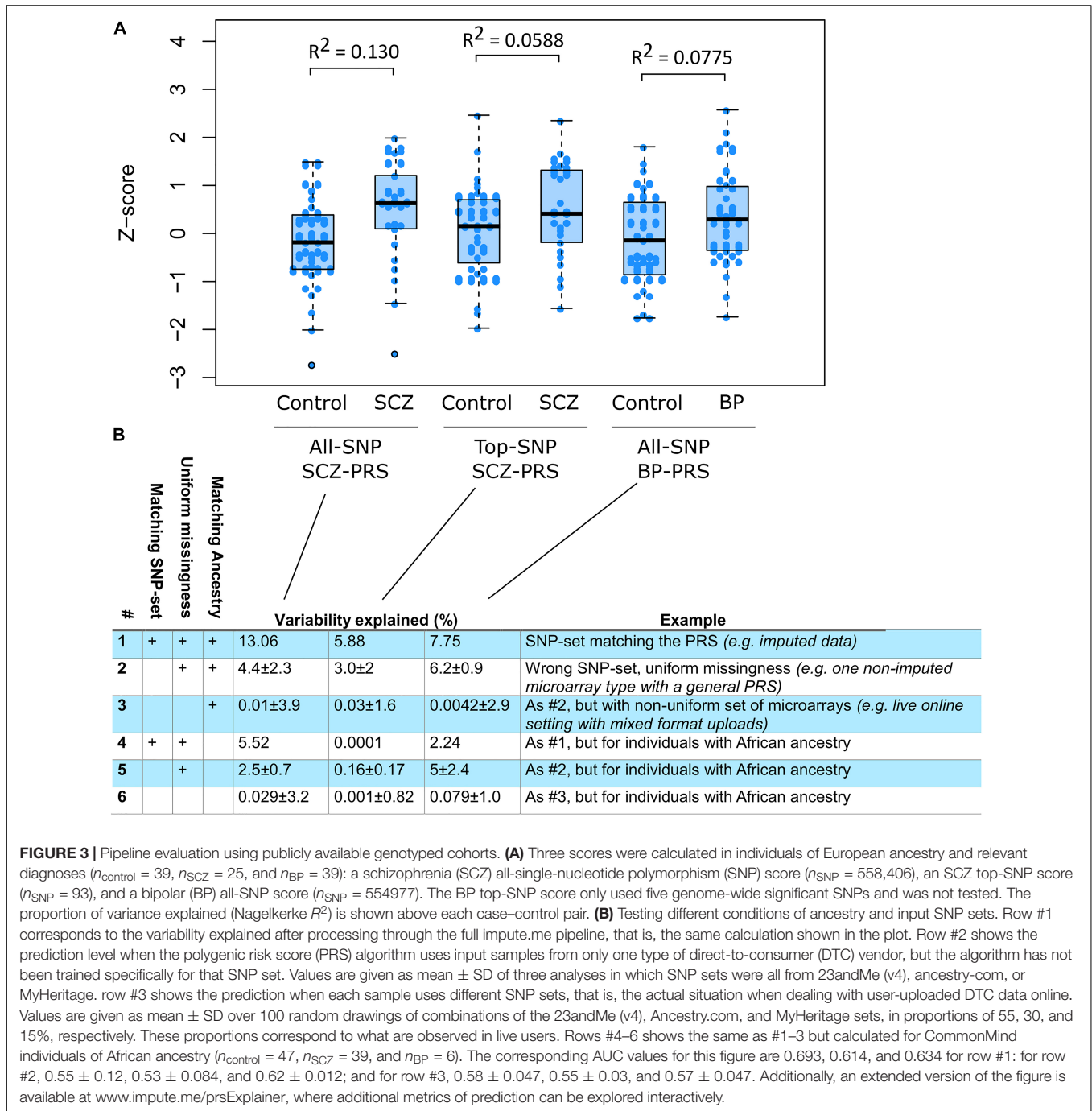
Of importance to this, we found that PRS prediction in mixed samples of non-imputed data causes severe problems. When training PRS algorithms, an SNP set is prespecified. The pipelines evaluated here were trained with HapMap3 as SNP set. Similar choices are made in other published PRSs. However, such SNP sets may not match with the SNPs available in downloadable raw data from DTC vendors. We therefore tested what prediction strength would be possible when using raw data directly from DTC vendors, both in a uniform setting (e.g., “all individuals use 23andMe v4 data”) and in a mixed setting (e.g., “individuals have data from different vendors”). We found that in the uniform setting, roughly half the predictive strength remained when using genotype data that are not imputed to match the HapMap3 SNP sets (Figure 3, rows 2 and 4). In the mixed setting, virtually no predictive strength remained (Figure 3, rows 3 and 6). The mixed setting is the reality that is faced, both for third-party analytical services and for DTC vendors with different chip versions. Imputation is therefore likely to be an essential requirement in such scenarios.

To compare these findings with approaches that look at one SNP at the time, we extracted the SNPedia/Promethease SNPs that were indicated as associated with SCZ (Cariaso and Lennon, 2012). All cases ( $n = 25$ ) and all controls ( $n = 39$ ) had at least one risk variant from at least one of the 139 SNPs that indicated SCZ association. When focusing on SNPs that had the SNPedia/Promethease-defined “magnitude”-level (*sic.*) at  $>1.5$ , we found that 80% of the SCZ cases (20 of 25) had at least one SNPedia/Promethease risk variant. Among the healthy controls, 84% (33 of 39) had at least one such risk variant ( $p = 0.9$  for difference in proportions). In other words, it is not very predictive to know if you have a SCZ SNP. This illustrates the importance of considering more than one SNP at the time.

Finally, we compared pipeline reproducibility using two genome-data files, one obtained from MyHeritage and one from Ancestry.com, but sampled from the same person. After processing through the impute-me pipelines, the correlation between PRS values over 1,468 traits was  $r = 0.933$  between the two samples. Traits that showed discrepancy between the two data files typically were based on only few SNPs, of which one did not meet imputation quality thresholds for one of the data files.

## EVALUATION OF USAGE

As of June 2019, a total of 28,651 genomes had been uploaded to impute.me, and a total of 3.1 million analytical queries had been performed (Figure 4A). The following additional

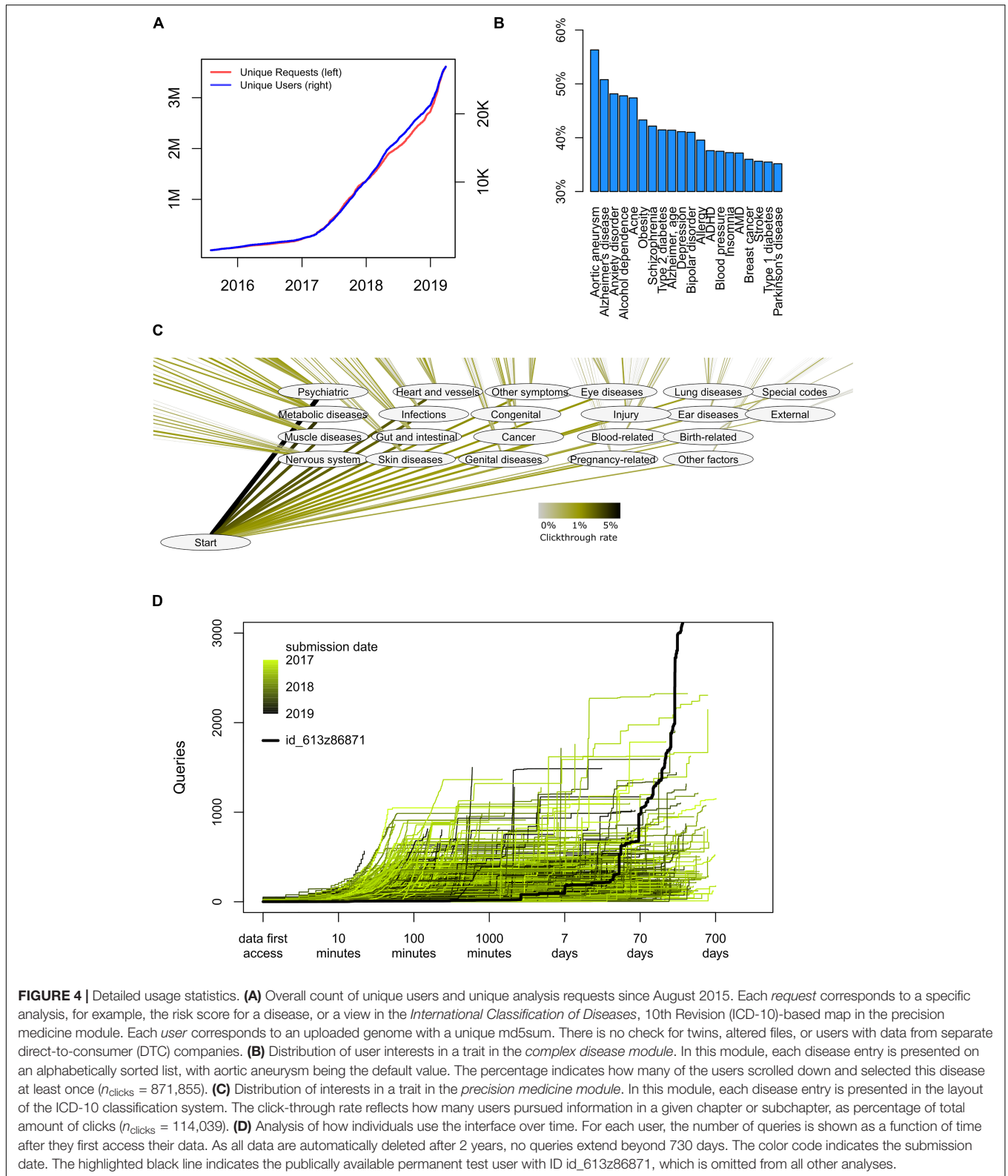


observations about user behavior may be of use to the genetics research community.

Common and well-known diseases are the most sought after. By overall click count and comparing over several different modules, there is no doubt that users are most interested in common disease types; diseases of the brain, heart, and metabolism are more requested. Interface design may of course play important roles in such choices. For example, the choice to serve disease traits as alphabetically sorted lists is likely to artificially inflate interest in, for example,

abdominal aneurysm (Figure 4B). However, the larger interest in psychiatry, cardiovascular, and metabolic disorders remains also in the precision medicine module, which is not presented as an alphabetically sorted list (Figure 4C). It is possible that greater scientific interest in PRSs in these fields also drives some of these effects, but we cannot explain why other fields where PRSs are actively discussed, such as cancer, are not attracting more attention.

Likewise, it seems that common disease (“complex disease module”) is more sought after than rare disease (“rare disease



**FIGURE 4 |** Detailed usage statistics. **(A)** Overall count of unique users and unique analysis requests since August 2015. Each *request* corresponds to a specific analysis, for example, the risk score for a disease, or a view in the *International Classification of Diseases*, 10th Revision (ICD-10)-based map in the precision medicine module. Each *user* corresponds to an uploaded genome with a unique md5sum. There is no check for twins, altered files, or users with data from separate direct-to-consumer (DTC) companies. **(B)** Distribution of user interests in a trait in the *complex disease module*. In this module, each disease entry is presented on an alphabetically sorted list, with aortic aneurysm being the default value. The percentage indicates how many of the users scrolled down and selected this disease at least once ( $n_{\text{clicks}} = 871,855$ ). **(C)** Distribution of interests in a trait in the *precision medicine module*. In this module, each disease entry is presented in the layout of the ICD-10 classification system. The click-through rate reflects how many users pursued information in a given chapter or subchapter, as percentage of total amount of clicks ( $n_{\text{clicks}} = 114,039$ ). **(D)** Analysis of how individuals use the interface over time. For each user, the number of queries is shown as a function of time after they first access their data. As all data are automatically deleted after 2 years, no queries extend beyond 730 days. The color code indicates the submission date. The highlighted black line indicates the publicly available permanent test user with ID id\_613z86871, which is omitted from all other analyses.

module”); 95% of all users visit the first, whereas only 70% visit the second. Again, interface design and project goals probably play a big role in this—the landing page headers says *Beyond*

*one SNP at the time*, and the rare disease module is found in the navigation bar only below seven other module entries. But it may also illustrate a central communication challenge for the field:

People are more interested in the genetics of common, complex diseases with small effect sizes (**Figure 2A**, lower right) but may interpret the results as if they were for rare diseases with large effect sizes (**Figure 2A**, upper left).

Finally, we have observed that usage of health genetic data surprisingly often is not just a test-and-forget event. When plotting query count as a function of time from first data access, we find an expected pattern of intense browsing the hours and days after first data access (**Figure 4D**). However, many users revisit their data even months and years after first data access, perhaps implying that results are considered and saved and then revisited at a later time in a different context.

## COMMUNICATION OF RISK SCORES

Generation of the PRS data presents one set of challenges, but communicating them to people in such a way as to make it both comprehensible and useful presents another (Lipkus and Hollands, 1999; Naik et al., 2012). We believe that this is a crucial unmet need in current genetics research, because presenting PRS data in a way that is useful requires an understanding of people's motivations for accessing them in the first place.

To date, studies of PRSs have focused on providing people with PRS information in relation to specific conditions [e.g., cancer (Bancroft et al., 2014, 2015; Smit et al., 2018; Young et al., 2018)] for which participants have an indicated risk and exploring understanding and reactions. No studies have examined what motivates people to seek out and access their own PRSs for common complex conditions, and little is known about how people understand or respond to the data they receive.

Polygenic risk scores information is inherently probabilistic in nature, which is well known to be difficult for people to understand (Hallowell et al., 1998; Smerecnik et al., 2009), and receiving information about genetic risk is not necessarily benign. When people receive genetic test results that they perceive to reflect high risk for a condition, this can have negative impact on outcomes like self-perception and affect, and in the case of receiving high-risk test results for Alzheimer's disease—can actually impact objective measures of cognitive performance (Wilhelm et al., 2009; Dar-Nimrod et al., 2013; Lineweaver et al., 2014; Lebowitz and Ahn, 2017; Turnwald et al., 2019). Therefore, how information about genetic risk is communicated matters.

The literature suggests that when communicating risk, the most useful and effective strategy is to use absolute risks (Lipkus and Hollands, 1999; Reyna et al., 2009; Naik et al., 2012). In the case of PRSs with modest predictive power, however, this may simply result in restating the population prevalence of a disease for everyone (Janssens, 2019). It is therefore important that the predictive strength is also included in this communication; that is how much the genetic component potentially could alter the absolute risk. The genetic component corresponds to the SNP heritability, and we are therefore exploring how to best include this information (e.g., **Figure 1**, right). Currently, we have registered the SNP heritability for 294 of the reported traits, available as an experimental option called “plot heritability.” We

believe that a main future direction is to experiment and expand on how to best communicate this to people.

It will therefore be useful to have a constant flow of people that are interested in interpreting their genetics and expose them to various modes of presentation. Some could involve statistically advanced concepts, like the area under the receiver operating characteristic curve (AUC) and SNP heritability, but others may take simpler approaches, such as the explanatory jar model pioneered for talking with families about genetics (Peay and Austin, 2011; Austin, 2019). One may even imagine layered models of increasing complexity. This should be followed up with questionnaires probing the level of understanding and general impact on users, something that is possible using the impute.me platform.

## ETHICS AND IMPLICATIONS

Using genetics to maximize the benefits and minimize the harms to individuals and society requires the effective management of the ethical, legal, and social implications of genetics. Researchers have a responsibility to ensure that the technology and the knowledge developed through genetic research are used responsibly, in light of the bioethical principles of beneficence, non-maleficence, justice, and autonomy (Lázaro-Muñoz et al., 2019). Given that for most complex disorders there is currently a lack of data regarding the harms or benefits of accessing PRS information, the fundamental principle in favor of making PRSs available to the public is that of autonomy—in the context of genetic testing, this refers to “the right of persons to make an informed, independent judgment about whether they wish to be tested and then whether they wish to know the details of the outcome of the testing” (Andrews et al., 1994). Accordingly, currently, DTC users can access health information through portals of DTC providers and through third-party applications (Kalokairinou et al., 2018; Tiller and Lacaze, 2018; Ahmed and Shabani, 2019). The problem is that many popular websites do not communicate high-quality genetic knowledge, in part possibly owing to the lack of engagement by the research communities (Badalato et al., 2017). One solution to this problem is to call for regulation and to ban such sites. Alternatively, as we propose here, it is possible to meet user demands and strive to do so as ethically as possible.

To exemplify this, as researchers, we have a choice in whether to provide access to a state-of-the-art PRS for a disease or not. We know that this PRS does not explain everything about the disease, does not account for all the genetic information, and is not part of today's clinical guidelines. However, we also know that users are already accessing information about disease through DTC genetics. These users may get their information from flawed assumptions of SNP effect sizes or from commercial platforms with little interest in explaining the limitations of the score. We argue that the choice that maximizes the potential for benefits to individuals is to provide the score and to provide it in a setting that puts its consequence in perspective.

An example of such perspective is that of giving reports by disease score, and not by individual risk variant as is currently

the case in most third-party analytics apps. Many people carry the high-risk allele for a common variant, but fewer people have a high PRS, which is the sum of all such risk variants. An example of this is the 84% frequency of SCZ risk variants in healthy users according to SNPedia, as reported above. This means that for those autonomously seeking information on health genetics data, the use of PRSs has the potential to decrease the level of induced worry in people in comparison with the current levels. Similarly, smart interface design can actively steer people toward browsing results by indication, and away from the pervasive practice of reporting the worst genetic scores for any disease first. This too may serve to reduce induced worry, in alignment with the general approach of testing only on indication to limit false-positive rates. Finally, of course, adaptive warnings based on risk levels, including referral to resources such as [findageneticcounselor.com](http://findageneticcounselor.com), is something we continuously strive to optimize.

## CONCLUSION

In summary, we present impute.me as a fully operational General Data Protection Regulation (GDPR)-compliant genetic analysis engine covering a very broad range of health-related traits, specifically focusing on optimizing possibilities from microarray-based DNA measurements. The challenges, their solutions, and the curation work behind them are highly relevant today in a setting of highly varying quality in interpretation of personal consumer genetics. In the future, we can expect that PRS predictiveness will increase. This will mean a continued and increasing relevance of the platform, even more so as the number of individuals doing genetic testing increases. With a directed push toward responsible use of genetics, this may even prove to be an overall clinical benefit.

## METHODS

### Data Privacy and Security

On data submission, each personal genome is assigned a nine-digit alphanumeric unique identifier (“uniqueID”). This uniqueID is used as login and identifier throughout all downstream processes because it has no information that is personally linkable, as opposed to, for example, an email address. The uniqueID is initially linked to two types of data: those that can be traced back to individual that submitted the genome and those that cannot. Genomic data, filename of submitted data, and email address are of the first type: genomic data because it can be used with software such as *gedmatch* to trace family patterns, filename because it often contains the name of the submitter (e.g., 23andMe data use full name as standard), and email for obvious reasons. Data of the traceable type are deleted 14 days after processing, which is the period in which users are able to download their full imputed data sets. The exception is email addresses, which are not deleted but instead unlinked from the uniqueID and kept elsewhere for the purpose of follow-up studies. Either way, this means that

14 days after processing, there exists nothing on the servers that can link results (designated with a uniqueID) with the person who submitted the data (any of the three traceable data types). Thus, even if the database is leaked or lost, it is not possible to link the data to an actual person. After 2 years, the remaining non-traceable data, for example, the derived calculations, the risk scores, and the genotypes of SNPs of specific interest, are all completely deleted. All ingoing and outgoing data transfers are encrypted using Transport Layer Security (TLS 1.3). All storage is encrypted using the AES-256 standard.

This means that all data are collected for specified, explicit, and legitimate purposes in a transparent manner and kept in a form that permits identification of data subjects for no longer than is necessary for the purposes for which the personal data are processed. We therefore consider that these measures both provide adequate security and privacy protection and are in accordance with the GDPR.

### Preprocessing and Bioinformatics

After submission of data, a comprehensive bioinformatic processing of the genotype data takes place. This is done in order of free computing nodes becoming available, consisting of several support programs; first, a *shapeit* call is made to phase the data correctly (Delaneau et al., 2013), and then an *impute2* call is made with 1000 Genomes version 3 as reference (Howie et al., 2009; 1000 Genomes Project Consortium et al., 2015). Although the pipelines are not guaranteed to handle any format they receive, they currently operate with less than 1% processing failures, meaning uploads that cannot proceed through the full quality control and imputation pipelines. The failures are typically due to file formatting errors, missing chromosomes, or any number of other odd data corruptions that real-world data exchange suffers from.

Several customizations have been made with the goal of minimizing memory footprint and thereby allowing running in a clustered fashion on a series of small cloud computers. This allows for relatively easy scaling of capacity: one simple setup (“hub-only”), where calculations are run on the same computer as the website interface. Another is a hub + node-setup, where a central hub server stores data and shows the website, while a scalable number of node-servers perform all computationally heavy calculations. After preprocessing is finished, two new files are created: a *.gen* file with probabilistic information from imputation calls and a *simple format* file with best guess genotypes, called at a 0.9 *impute2* INFO threshold. All further calculations are based on these files. A mail with download links to these two files is returned to the user, along with a JSON-formatted file containing a machine-readable summary of all calculations, as well as links with guidance to obtain more in-depth information on personal DNA interpretation (Folkersen, 2018).

### Polygenic Risk Score Calculation

From the preprocessed data, a modular set of trait predictor algorithms is applied. For many of the modules, the calculations



are trivial. For example, this could be the reporting of presence and/or absence of a specific genotype, such as ACTN3 and ACE-gene SNPs known to be (weakly) associated with athletic performance. These are included mostly because users expect them to be. For others, we rely heavily on PRSs.

An important distinguishing factor between different PRS algorithms is how risk alleles are selected. A commonly used approach includes variants based on whether they surpass a given  $p$ -value threshold in the GWAS, retaining only linkage disequilibrium (LD)-independent variants using LD-based clumping, often with a  $p$ -value threshold of genome-wide significance ( $p < 5e^{-8}$ ). Herein, we refer to this approach as the “top-SNP” approach. The top-SNP approach has the advantage that it is simple to explain, is easy to obtain for many GWAS, and has a light computational burden (e.g., Buniello et al., 2019; Lambert et al., 2019; Patron et al., 2019; Watanabe et al., 2019). However, research has repeatedly shown that the inclusion of variants that do not achieve genome-wide significance improves the variance explained by PRSs, with PRSs including all variants often explaining the most variance. PRSs based on GWAS effect sizes that have undergone shrinkage to account for the LD between variants have been shown to explain more variance than PRSs that account for LD via LD-based clumping (Vilhjálmsdóttir et al., 2015). Herein, we refer to this approach as the all-SNP approach. It is more computationally and practically intensive to implement at scale. Consequently, within impute.me, each trait or disease reported shows all-SNP-based PRS calculations if such is available, and top-SNP-based PRS calculations if not.

In the top-SNP calculation mode, the results are scaled such that the mean of a population is zero and the standard deviation (SD) is 1, according to the relevant 1000 Genomes super-population: African, admixed American, East Asian, European, or South Asian.

$$\text{Population-score}_{\text{SNP}} = \text{frequency}_{\text{SNP}} \times 2 \times \text{beta}_{\text{SNP}}$$

$$\begin{aligned} \text{Zero-centered-score} = \\ \sum \text{Beta}_{\text{SNP}} \times \text{Effect-allele-count}_{\text{SNP}} \\ - \text{Population-score}_{\text{SNP}} \end{aligned}$$

$$\text{Z-score} = \frac{\text{Zero-centered-score}}{\text{Standard-deviation}_{\text{population}}}$$

where beta [or log(odds ratio)] is the reported effect size for the SNP effect allele,  $\text{frequency}_{\text{SNP}}$  is the allele frequency for the effect allele, and the  $\text{Effect-allele-count}_{\text{SNP}}$  is the allele count from genotype data (0, 1, or 2).

In the all-SNP calculation, the scaling is similar but done empirically, that is, based on previous impute.me users of matching ethnicity. This mode of scaling is also available as an optional functionality in the top-SNP calculations and generally seems to match well with the default 1000 Genomes super-population scaling.

The all-SNP scores were derived using weightings from the LDpred algorithm (Vilhjálmsdóttir et al., 2015). This algorithm adjusts the effect of each SNP allele for those of other SNP alleles in LD with it and also takes into account the likelihood of a given allele to have a true effect according to a user-defined parameter, which here was taken as  $w1$ , that is, the full set of SNPs. The algorithm was directed to use hapmap3 SNPs that had a minor allele frequency  $> 0.05$ , Hardy-Weinberg equilibrium  $p > 1e^{-05}$ , and genotype yield  $> 0.95$ , consistent with our expectation that these would be the best imputed SNPs after full pipeline processing.

## Pipeline Testing

To test the pipelines described herein, the CommonMind genotypes measured with the microarray of the type H1M were downloaded along with phenotypic information. Each sample was processed through the impute.me pipelines, using the batch upload functionality. Reported ethnicity was compared with pipeline (genotype) assigned ethnicity and found to be concordant.

After pipeline completion, we extracted three PRSs for each sample, corresponding to SCZ all-SNP, SCZ top-SNP, and BP all-SNP. In the github repository for impute.me, these three correspond to the scores labeled *SCZ\_2014\_PGC\_EXCL\_DK.EurUnrel.hapmap3.all.ldpred.effects*, *schizophrenia\_25056061*, and *BIP\_2016\_PGC.All.hapmap3.all.ldpred.effects* trait IDs (Schizophrenia Working Group of the Psychiatric Genomics Consortium, 2014; Hou et al., 2016). These extracted scores formed the basis of the row #1 and #4 calculations in **Figure 3**. The remaining rows were created by subsetting the best guess imputed genotypes into new sets of users, corresponding to each of three major DTC vendors and then re-running the scoring algorithms with either uniform data or mixed data. Uniform data are here defined as all 195 samples having the same set of SNPs available, corresponding to one of three DTC vendors in each run. Mixed data are defined as samples having different sets of SNPs available, a set corresponding to actual distributions of customers from different DTC vendors, with distributions redrawn 100 times. We estimated the predictive ability of the PRSs using Nagelkerke's  $R^2$  and AUC.

## Usage Evaluation

A log data freeze was performed on June 8, 2019 by making a copy of all usage log files and then removing the uniqueID of each user. This was done to prevent it from being linked with the genetic data of that user. The exception was the publicly available permanent test user with ID *id\_613z86871*, which was lifted out before analysis and is not included in other summary statistics. Generally, a user corresponds to an uploaded genome with a unique md5sum. Click-through rates were calculated as fraction of users that performed any query in the module in question; for example, the precision medicine module was only launched in September 2018 and, therefore, only counts clicks from people

who have used it. Plots were generated using base-R version 3.4.2 and cytoscape version 3.71.

## URLS

Code repository: <https://github.com/lassefolkersen/impute-me>

Web resource: <https://www.impute.me/>

## DATA AVAILABILITY STATEMENT

Publicly available datasets were analyzed in this study. This data can be found here: CommonMind data doi: 10.7303/syn2759792.

## ETHICS STATEMENT

The studies involving genotypes of human participants were reviewed and approved by the CommonMind Consortium. This data is generated from postmortem human brain specimens originating from tissue collections at the Mount Sinai NIH Brain Bank and Tissue Repository, University of Pennsylvania Brain Bank of Psychiatric illnesses and Alzheimer's Disease Core Center, The University of Pittsburgh NIH NeuroBioBank Brain and Tissue Repository, and the NIMH Human Brain Collection Core. Written informed consent for participation was not required for this study in accordance with the national legislation and the institutional requirements.

## REFERENCES

- 1000 Genomes Project Consortium, Auton, A., Brooks, L. D., Durbin, R. M., Garrison, E. P., Kang, H. M., et al. (2015). A global reference for human genetic variation. *Nature* 526, 68–74. doi: 10.1038/nature15393
- Ahmed, E., and Shabani, M. (2019). DNA data marketplace: an analysis of the ethical concerns regarding the participation of the individuals. *Front Genet.* 10:1107. doi: 10.3389/fgene.2019.01107
- Andrews, L. B., Fullerton, J. E., Holtzman, N. A., and Motulsky, A. G. (1994). *Assessing Genetic Risks - Implications for Health and Social Policy*. Washington, DC: National Academies Press.
- Austin, J. C. (2019). Evidence-based genetic counseling for psychiatric disorders: a road map. *Cold Spring Harb. Perspect. Med.* 9:a036608. doi: 10.1101/cshperspect.a036608
- Badalato, L., Kalokairinou, L., and Borry, P. (2017). Third party interpretation of raw genetic data: an ethical exploration. *Eur. J. Hum. Genet.* 25, 1189–1194. doi: 10.1038/ejhg.2017.126
- Baig, S. S., Strong, M., Rosser, E., Taverner, N. V., Glew, R., Miedzybrodzka, Z., et al. (2016). UK Huntington's disease prediction consortium, quarrell OW. 22 years of predictive testing for Huntington's disease: the experience of the UK Huntington's prediction consortium. *Eur. J. Hum. Genet.* 24, 1396–1402. doi: 10.1038/ejhg.2016.36
- Bancroft, E. K., Castro, E., Ardern-Jones, A., Moynihan, C., Page, E., Taylor, N., et al. (2014). It's all very well reading the letters in the genome, but it's a long way to being able to write: men's interpretations of undergoing genetic profiling to determine future risk of prostate cancer. *Fam. Cancer* 13, 625–635. doi: 10.1007/s10689-014-9734-3

## AUTHOR CONTRIBUTIONS

LF coded the code. All authors contributed to interpretation, drafting the work, critical revision for important intellectual content, and final approval of the manuscript.

## FUNDING

CL and OP were supported by UK Medical Research Council grant N015746 and by the National Institute for Health Research (NIHR) Biomedical Research Centre at South London and Maudsley NHS Foundation Trust and King's College London. CommonMind is supported by funding from Takeda Pharmaceuticals Company Limited, F. Hoffman-La Roche Ltd., and NIH grants R01MH085542, R01MH093725, P50MH066392, P50MH080405, R01MH097276, RO1-MH-075916, P50M096891, P50MH084053S1, R37MH057881, AG02219, AG05138, MH06692, R01MH110921, R01MH109677, R01MH109897, U01MH103392, and contract HHSN271201300031C through IRP NIMH. The funder was not involved in the study design, collection, analysis, interpretation of data, the writing of this article or the decision to submit it for publication.

## ACKNOWLEDGMENTS

In addition to the crucial resources cited in the text, we wish to thank the CommonMind Consortium for availability of testing data. We also wish to thank Ron Nudel for discussions on the ethics section.

- Bancroft, E. K., Castro, E., Bancroft, G. A., Ardern-Jones, A., Moynihan, C., Page, E., et al. (2015). The psychological impact of undergoing genetic-risk profiling in men with a family history of prostate cancer. *Psychooncology* 24, 1492–1499. doi: 10.1002/pon.3814
- Buniello, A., MacArthur, J. A. L., Cerezo, M., Harris, L. W., Hayhurst, J., Malangone, C., et al. (2019). The NHGRI-EBI GWAS catalog of published genome-wide association studies, targeted arrays and summary statistics. *Nucleic Acids Res.* 47, D1005–D1012. doi: 10.1093/nar/gky1120
- Byrjalsen, A., Stoltze, U., Wadt, K., Hjalgrim, L. L., Gerdes, A. M., Schmiegelow, K., et al. (2018). Pediatric cancer families' participation in whole-genome sequencing research in Denmark: parent perspectives. *Eur. J. Cancer Care* 27:e12877. doi: 10.1111/ecc.12877
- Cariaso, M., and Lennon, G. (2012). SNPedia: a wiki supporting personal genome annotation, interpretation and analysis. *Nucleic Acids Res.* 40, D1308–D1312. doi: 10.1093/nar/gkr798
- Curtis, D. (2018). Polygenic risk score for schizophrenia is more strongly associated with ancestry than with schizophrenia. *Psychiatr. Genet.* 28, 85–89. doi: 10.1097/YPG.0000000000000206
- Dar-Nimrod, I., Zuckerman, M., and Duberstein, P. R. (2013). The effects of learning about one's own genetic susceptibility to alcoholism: a randomized experiment. *Genet. Med.* 15, 132–138. doi: 10.1038/gim.2012.111
- Delaneau, O., Howie, B., Cox, A. J., Zagury, J. F., and Marchini, J. (2013). Haplotype estimation using sequencing reads. *Am. J. Hum. Genet.* 93, 687–696. doi: 10.1016/j.ajhg.2013.09.002
- Folkersen, L. (2018). *Understand Your DNA - A Guide*. Singapore: World Scientific Publishing.

- Greshake, B., Bayer, P. E., Rausch, H., and Reda, J. (2014). openSNP—a crowdsourced web resource for personal genomics. *PLoS One* 9:e89204. doi: 10.1371/journal.pone.0089204
- Hallowell, N., Statham, H., and Murton, F. (1998). Women's understanding of their risk of developing breast/ovarian cancer before and after genetic counseling. *J. Genet. Couns.* 7, 345–364. doi: 10.1023/A:1022072017436
- Hou, L., Bergen, S. E., Akula, N., Song, J., Hultman, C. M., Landén, M., et al. (2016). Genome-wide association study of 40,000 individuals identifies two novel loci associated with bipolar disorder. *Hum. Mol. Genet.* 25, 3383–3394. doi: 10.1093/hmg/ddw181
- Howe, B. N., Donnelly, P., and Marchini, J. (2009). A flexible and accurate genotype imputation method for the next generation of genome-wide association studies. *PLoS Genet.* 5:e1000529. doi: 10.1371/journal.pgen.1000529
- Janssens, A. C. J. W. (2019). Proprietary algorithms for polygenic risk: protecting scientific innovation or hiding the lack of it? *Genes* 10:E448.
- Kalia, S. S., Adelman, K., Bale, S. J., Chung, W. K., Eng, C., Evans, J. P., et al. (2017). Recommendations for reporting of secondary findings in clinical exome and genome sequencing, update (ACMG SF v2.0): a policy statement of the American college of medical genetics and genomics. *Genet. Med.* 19, 249–255. doi: 10.1038/gim.2016.190
- Kalokairinou, L., Howard, H. C., Slokenberga, S., Fisher, E., Flatscher-Thöni, M., Hartlev, M., et al. (2018). Legislation of direct-to-consumer genetic testing in Europe: a fragmented regulatory landscape. *J. Commun. Genet.* 9, 117–132. doi: 10.1007/s12687-017-0344-2
- Kaye, J. (2008). The regulation of direct-to-consumer genetic tests. *Hum. Mol. Genet.* 17, R180–R183.
- Khera, A. V., Chaffin, M., Aragam, K. G., Haas, M. E., Roselli, C., Choi, S. H., et al. (2018). Genome-wide polygenic scores for common diseases identify individuals with risk equivalent to monogenic mutations. *Nat. Genet.* 50, 1219–1224. doi: 10.1038/s41588-018-0183-z
- Lam, M., Chen, C. Y., Li, Z., Martin, A. R., Bryois, J., Ma, X., et al. (2019). Comparative genetic architectures of schizophrenia in East Asian and European populations. *Nat. Genet.* 51, 1670–1678. doi: 10.1038/s41588-019-0512-x
- Lambert, S., Gil, L., Jupp, S., Chapman, M., Parkinson, H., Danesh, J., et al. (2019). *The Polygenic Score (PGS) Catalog: An Open Database To Enable Reproducibility And Systematic Evaluation*. Available at: www.pgscatalog.org (accessed November 2019).
- Lázaro-Muñoz, G., Sabatello, M., Huckins, L., Peay, H., Degenhardt, F., Meiser, B., et al. (2019). ISPG Ethics Committee. International society of psychiatric genetics ethics committee: issues facing us. *Am. J. Med. Genet. B Neuropsychiatr. Genet.* 180, 543–554. doi: 10.1002/ajmg.b.32736
- Lebowitz, M. S., and Ahn, W. K. (2017). Testing positive for a genetic predisposition to depression magnifies retrospective memory for depressive symptoms. *J. Consult. Clin. Psychol.* 85, 1052–1063. doi: 10.1037/ccp0000254
- Lee, S. H., Clark, S., and van der Werf, J. H. J. (2018). Estimation of genomic prediction accuracy from reference populations with varying degrees of relationship. *PLoS One* 12:e0189775. doi: 10.1371/journal.pone.0189775
- Lee, S. H., van der Werf, J. H., Hayes, B. J., Goddard, M. E., and Visscher, P. M. (2008). Predicting unobserved phenotypes for complex traits from whole-genome SNP data. *PLoS Genet.* 4:e1000231. doi: 10.1371/journal.pgen.1000231
- Lee, S. H., Weerasinghe, W. M., Wray, N. R., Goddard, M. E., and van der Werf, J. H. (2017). Using information of relatives in genomic prediction to apply effective stratified medicine. *Sci. Rep.* 7:42091. doi: 10.1038/srep42091
- Lewis, C. M., and Vassos, E. (2017). Prospects for using risk scores in polygenic medicine. *Genome Med.* 9:96. doi: 10.1186/s13073-017-0489-y
- Li, Z., Chen, J., Yu, H., He, L., Xu, Y., Zhang, D., et al. (2017). Genome-wide association analysis identifies 30 new susceptibility loci for schizophrenia. *Nat. Genet.* 49, 1576–1583.
- Lineweaver, T. T., Bondi, M. W., Galasko, D., and Salmon, D. P. (2014). Effect of knowledge of APOE genotype on subjective and objective memory performance in healthy older adults. *Am. J. Psychiatry* 171, 201–208. doi: 10.1176/appi.ajp.2013.12121590
- Lipkus, I. M., and Hollands, J. G. (1999). The visual communication of risk. *J. Natl. Cancer Inst.* 25, 149–163.
- Maxwell, K. N., Domchek, S. M., Nathanson, K. L., and Robson, M. E. (2016). Population frequency of germline BRCA1/2 mutations. *J. Clin. Oncol.* 34, 4183–4185. doi: 10.1200/jco.2016.67.0554
- Naik, G., Ahmed, H., and Edwards, A. G. (2012). Communicating risk to patients and the public. *Br. J. Gen. Pract.* 62, 213–216.
- Natarajan, P., Young, R., Stitzel, N. O., Padmanabhan, S., Baber, U., Mehran, R., et al. (2017). Polygenic risk score identifies subgroup with higher burden of atherosclerosis and greater relative benefit from statin therapy in the primary prevention setting. *Circulation* 135, 2091–2101. doi: 10.1161/circulationaha.116.024436
- Patron, J., Serra-Cayuela, A., Han, B., Li, C., and Wishart, D. S. (2019). Assessing the performance of genome-wide association studies for predicting disease risk. *PLoS One* 14:e0220215. doi: 10.1371/journal.pone.0220215
- Peay, H. L., and Austin, J. (2011). *How To Talk With Families About Genetics And Psychiatric Illness*. New York, NY: W. W. Norton & Company.
- Reyna, V. F., Nelson, W. L., Han, P. K., and Dieckmann, N. F. (2009). How numeracy influences risk comprehension and medical decision making. *Psychol. Bull.* 135, 943–973. doi: 10.1037/a0017327
- Schizophrenia Working Group of the Psychiatric Genomics Consortium, (2014). Biological insights from 108 schizophrenia-associated genetic loci. *Nature* 511, 421–427. doi: 10.1038/nature13595
- Smerecnik, C. M., Mesters, I., Verweij, E., de Vries, N. K., and de Vries, H. (2009). A systematic review of the impact of genetic counseling on risk perception accuracy. *J. Genet. Couns.* 18, 217–228. doi: 10.1007/s10897-008-9210-z
- Smit, A. K., Newson, A. J., Best, M., Badcock, C. A., Butow, P. N., Kirk, J., et al. (2018). Distress, uncertainty, and positive experiences associated with receiving information on personal genomic risk of melanoma. *Eur. J. Hum. Genet.* 26, 1094–1100. doi: 10.1038/s41431-018-0145-z
- Tiller, J., and Lacaze, P. (2018). Regulation of internet-based genetic testing: challenges for australia and other jurisdictions. *Front. Public Health* 6:24. doi: 10.3389/fgene.2019.00024
- Turnwald, B. P., Goyer, J. P., Boles, D. Z., Silder, A., Delp, S. L., and Crum, A. J. (2019). Learning one's genetic risk changes physiology independent of actual genetic risk. *Nat. Hum. Behav.* 3, 48–56. doi: 10.1038/s41562-018-0483-4
- Vilhjálmsdóttir, B. J., Yang, J., Finucane, H. K., Gusev, A., Lindström, S., Ripke, S., et al. (2015). Modeling linkage disequilibrium increases accuracy of polygenic risk scores. *Am. J. Hum. Genet.* 97, 576–592.
- Watanabe, K., Stringer, S., Frei, O., Umiaevia, Mirkov, M., de Leeuw, C., Polderman, T. J. C., et al. (2019). A global overview of pleiotropy and genetic architecture in complex traits. *Nat. Genet.* 51, 1339–1348. doi: 10.1038/s41588-019-0481-0
- Wilhelm, K., Meiser, B., Mitchell, P. B., Finch, A. W., Siegel, J. E., Parker, G., et al. (2009). Issues concerning feedback about genetic testing and risk of depression. *Br. J. Psychiatry* 194, 404–410. doi: 10.1192/bjp.bp.107.047514
- Wünnemann, F., Sin Lo, K., Langford-Avelar, A., Busseuil, D., Dubé, M. P., Tardif, J. C., et al. (2019). Validation of genome-wide polygenic risk scores for coronary artery disease in french Canadians. *Circ. Genom. Precis. Med.* 12:e002481.
- Young, M. A., Forrest, L. E., Rasmussen, V. M., James, P., Mitchell, G., Sawyer, S. D., et al. (2018). Making sense of SNPs: women's understanding and experiences of receiving a personalized profile of their breast cancer risks. *J. Genet. Couns.* 27, 702–708. doi: 10.1007/s10897-017-0162-z

**Disclaimer:** The views expressed are those of the author(s) and not necessarily those of the NHS, the NIHR or the Department of Health and Social Care.

**Conflict of Interest:** The authors declare that the voluntary donations received by impute.me go to a registered company, from where all of it is used to pay for server-costs. The company is a Danish-law IVS company with ID 37918806, financially audited under Danish tax law.

Copyright © 2020 Folkersen, Pain, Ingason, Werge, Lewis and Austin. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.