








**SPECIAL ARTICLE**

# Predicting venous thromboembolism risk from exomes in the Critical Assessment of Genome Interpretation (CAGI) challenges

Gregory McInnes<sup>1</sup>  | Roxana Daneshjou<sup>2</sup> | Panagiostis Katsonis<sup>3</sup>  |  
 Olivier Lichtarge<sup>3,4,5,6</sup> | Raj G. Srinivasan<sup>7</sup> | Sadhna Rana<sup>7</sup>  | Predrag Radivojac<sup>8</sup>  |  
 Sean D. Mooney<sup>9</sup>  | Kymberleigh A. Pagel<sup>10</sup> | Moses Stamboulia<sup>10</sup> |  
 Yuxiang Jiang<sup>10</sup> | Emidio Capriotti<sup>11</sup> | Yanran Wang<sup>12</sup> | Yana Bromberg<sup>12</sup> |  
 Samuele Bovo<sup>13</sup> | Castrense Savojardo<sup>13</sup> | Pier Luigi Martelli<sup>13</sup> | Rita Casadio<sup>13,14</sup> |  
 Lipika R. Pal<sup>15</sup> | John Moul<sup>15,16</sup>  | Steven Brenner<sup>17</sup>  | Russ Altman<sup>18</sup>

<sup>1</sup>Biomedical Informatics Training Program, Stanford University, Stanford, California

<sup>2</sup>Department of Dermatology, Stanford School of Medicine, Stanford, California

<sup>3</sup>Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, Texas

<sup>4</sup>Department of Biochemistry & Molecular Biology, Baylor College of Medicine, Houston, Texas

<sup>5</sup>Department of Pharmacology, Baylor College of Medicine, Houston, Texas

<sup>6</sup>Computational and Integrative Biomedical Research Center, Baylor College of Medicine, Houston, Texas

<sup>7</sup>Innovations Labs, Tata Consultancy Services, Hyderabad, India

<sup>8</sup>Khoury College of Computer and Information Sciences, Northeastern University, Boston, Massachusetts

<sup>9</sup>Department of Biomedical Informatics and Medical Education, University of Washington, Seattle, Washington

<sup>10</sup>Department of Computer Science and Informatics, Indiana University, Bloomington, Indiana

<sup>11</sup>BioFOLD Unit, Department of Pharmacy and Biotechnology (FaBiT), University of Bologna, Bologna, Italy

<sup>12</sup>Department of Biochemistry and Microbiology, Rutgers University, New Brunswick, New Jersey

<sup>13</sup>Department of Pharmacy and Biotechnology, Bologna Biocomputing Group, University of Bologna, Italy

<sup>14</sup>Institute of Biomembrane and Bioenergetics, Consiglio Nazionale delle Ricerche, Bari, Italy

<sup>15</sup>Institute for Bioscience and Biotechnology Research, University of Maryland, Rockville, Maryland

<sup>16</sup>Department of Cell Biology and Molecular Genetics, University of Maryland, College Park, Maryland

<sup>17</sup>Department of Plant and Microbial biology, University of California Berkeley, Berkeley, California

<sup>18</sup>Departments of Bioengineering, Biomedical Data Science, Genetics, and Medicine, Stanford University, Stanford, California

**Correspondence**

Gregory McInnes, Biomedical Informatics Training Program, Stanford University, Stanford, CA.  
 Email: gmcinnes@stanford.edu

**Funding information**

National Institute of Health (NIH), Grant/Award Number: T32LM012409; NIH Office of Extramural Research, Grant/Award Numbers: HG006650, HG007346; National Institute of General Medical Sciences, Grant/Award Number: GM115486

**Abstract**

Genetics play a key role in venous thromboembolism (VTE) risk, however established risk factors in European populations do not translate to individuals of African descent because of the differences in allele frequencies between populations. As part of the fifth iteration of the Critical Assessment of Genome Interpretation, participants were asked to predict VTE status in exome data from African American subjects. Participants were provided with 103 unlabeled exomes from patients treated with warfarin for non-VTE causes or VTE and asked to predict which disease each subject had been treated for. Given the lack of training data, many participants opted to use unsupervised machine learning methods, clustering the exomes by variation in genes

known to be associated with VTE. The best performing method using only VTE related genes achieved an area under the ROC curve of 0.65. Here, we discuss the range of methods used in the prediction of VTE from sequence data and explore some of the difficulties of conducting a challenge with known confounders. In addition, we show that an existing genetic risk score for VTE that was developed in European subjects works well in African Americans.

#### KEYWORDS

exomes, machine learning, phenotype prediction, prediction challenge, venous thromboembolism

## 1 | INTRODUCTION

There are 300,000 to 900,000 cases of venous thromboembolism (VTE) a year in the United States alone (Beckman, Hooper, Critchley, & Ortel, 2010). VTE captures both deep vein thrombosis (DVT) and pulmonary embolism (PE). There are differences in the incidence of VTE based on ancestry; individuals of African ancestry have a 30–60% higher incidence of VTEs than people of European ancestry (Roberts, Patel, & Arya, 2009; Zakai & McClure, 2011). VTE risk is multifactorial, both environmental and genetic factors are involved (Feero, 2004). For individuals of European descent, the commonly seen VTE risk factors are *F5* R506Q (NC\_000001.11:g.169549811C>T; three- to five-fold increased risk of VTE in carriers) and *F2* G20210A (NC\_000011.10:g.46739505 G>A; two to three-fold increased risk of VTE in carriers; Middeldorp & van Hylckama Vlieg, 2008; Rosendaal & Reitsma, 2009).

However, the genetic variants that confer risk in populations of European descent differ in allele frequency between Europeans and Africans, and population-specific genetic factors influencing the higher VTE rate are not well characterized in African Americans (Dowling et al., 2003). Recent studies identified a population-specific genetic risk factor in African Americans, but much of the genetic risk is still undiscovered (Daneshjou et al., 2016; Hernandez et al., 2016). Previous work has been done to develop genetic risk models for VTE in European populations, but no such risk model exists for individuals of African descent and the existing models have not been tested in African populations (Soria et al., 2014).

The Critical Assessment of Genomic Interpretation (CAGI) aims to objectively assess the prediction of phenotypic impacts of genetic variation. In the fifth iteration of CAGI, participants were challenged to predict the VTE status of 103 African American individuals from exome data. This dataset was used as part of a warfarin dosage prediction challenge in CAGI 4, where participants were asked to predict the precise warfarin dosage of each individual (Daneshjou et al., 2017). VTE often requires long term use of anticoagulants. The dataset comprised 66 individuals with a VTE diagnosis and 37 individuals on warfarin for non-VTE causes (such as atrial fibrillation prophylaxis, AF). Thus, we were able to repurpose this data CAGI 5 and participants were asked to distinguish between individuals that

were prescribed warfarin for a clotting disorder versus those that were prescribed warfarin for non-VTE purposes.

The use of exome data presents a challenge for the prediction of complex disease. Exomes only capture the coding regions of the genome. Previous studies have shown that noncoding portions of the genome explain 79% of the heritability for complex traits, whereas coding regions explain less than 10% (Gusev et al., 2014). As a complex trait, we might expect VTE to have a similar amount of the heritability explained by the exome. Studies of VTE in Europeans in the UK Biobank have calculated the heritability on the liability scale in Europeans to be 0.14 and disease prevalence to be 2%, which would indicate that the theoretical maximum area under the ROC curve (AUC) that could be achieved in predicting VTE from coding regions is approximately 0.62 (Canela-Xandri, Rawlik, & Tenesa, 2018; Wray, Yang, Goddard, & Visscher, 2010). The theoretical maximum AUC may be different in African American populations, but because no large-scale studies have been performed to determine heritability estimates of VTE in African Americans we cannot know *a priori* what to expect. However, studies have shown that the contribution of rare variants explains a significant amount of the observed variance in some complex traits, which may limit estimates heritability derived from genotyping panels (De Rubeis et al., 2014; Marouli et al., 2017; Purcell et al., 2014; Simons, Turchin, Pritchard, & Sella, 2014). Exome sequencing data provides an opportunity to evaluate the predictive power of coding variants to the prediction of VTE.

## 2 | METHODS

### 2.1 | Data distribution

Participants were provided exome data for all 103 subjects in the VCF file format as well as corresponding covariate data. The covariate data included was subject age, height, weight, sex, and drug regimen (aspirin, amiodarone, and warfarin dose). Amiodarone is an antiarrhythmic drug used to treat atrial fibrillation, which could be a clear sign that the subject belonged in the AF group. However, only one subject was on amiodarone, so this conferred no predictive

advantage to the participants. Participants consented to the CAGI data use agreement.

## 2.2 | Predicting phenotypes

Participants were asked to make VTE status predictions for all 103 subjects in the provided data. No labeled training data was provided. Participants were required to return a text file with predicted disease status and confidence in the prediction for each subject. They were also provided with a validation script to check their output formatting. Participants were asked to provide a brief description of their prediction methods for each submission. The prediction results were presented at the CAGI 5 meeting.

## 2.3 | Data quality

The data had previously undergone rigorous QC using ancestry informative markers to confirm self-reported ancestry and identity by state (IBS) analysis to ensure that samples were not related, as previously described (Daneshjou et al., 2014).

## 2.4 | Assessing predicted phenotypes

To assess the submissions of each group, several evaluation metrics were used. Predictions were evaluated using AUC, accuracy, sensitivity, specificity, and F1 scores. Some participants submitted binary class predictions rather than probabilities. To fairly evaluate predictions across all groups the predictions that were submitted as probabilities were binarized using a cutoff of 0.5, where a score greater than 0.5 indicates a VTE prediction. Accuracy, sensitivity, specificity, recall, and F1 scores were then computed with the binarized data, whereas AUC was calculated on the submitted scores. Confidence intervals for the AUC scores were calculated according to the method presented by DeLong, DeLong, and Clarke-Pearson (1988).

## 2.5 | Establishing a baseline

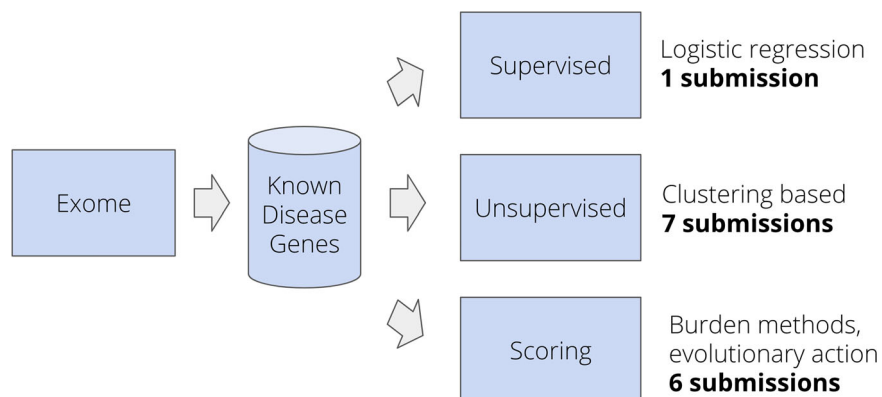
A baseline prediction score was calculated using the multilocus genetic risk score proposed by Soria et al. (2014). The proposed method uses a linear model of 17 loci across nine VTE related genes. To compute the scores the number of the alternate alleles at each site was multiplied by the corresponding coefficients proposed by Soria et al., (2014). As with the participant submitted scores, the genetic risk scores were binarized using a threshold of 0.5 before calculating accuracy, sensitivity, specificity, recall, and F1 scores, and the raw scores were used to compute AUC. Four sites included in the risk model could not be used for this data: Two sites that were outside the exome capture region, and two INDEL variants. INDELS were not called during variant calling of the data used for this study, as such ABO\*A1 haplotype could not be assigned to subjects and was excluded from the risk model calculations.

European and African allele frequencies for the SNPs used in the baseline model were determined by querying gnoMAD (Karczewski et al., 2019). Allele frequencies for the SNPs in the study population were calculated using VCFtools (Danecek et al., 2011).

## 3 | RESULTS

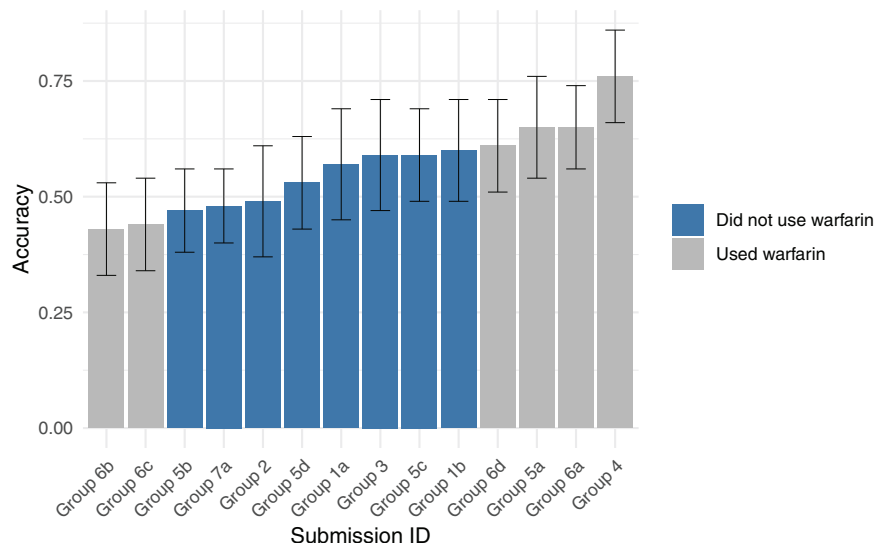
We assessed 14 submissions of phenotype predictions from seven groups. As no training dataset was provided, most participants chose to use unsupervised models trained on variants from genes previously reported to be associated with the phenotypes. Some groups used burden based scoring methods, scoring samples by the frequency of damaging variants in selected genes.

Each of the participants formulated their own strategy for predicting phenotype from the exome data. Although each was unique, there were many similarities between the methods used (Figure 1). All submissions but one primarily used the genetic data, each group first selected genes related to the phenotypes of interest from a disease-gene database, then used the variants in those genes



**FIGURE 1** General participant workflow. Each group formed their own approach to predicting phenotypes of the exomes, but there were some similarities across all submissions. All groups subset the exome into genes known to be involved in the phenotypes of interest, then made predictions on the basis of the variants in those subset genes. Some groups generated scores for each individual on the basis of burden of variants of a certain class. Others clustered the genotypes alone and segmented the clusters into predicted phenotypes

Area under the ROC curve for participant submissions



**FIGURE 2** Area under the ROC (receiver operating characteristic) curve for all submissions. Submissions that used knowledge of warfarin confounding in the dataset (either by including the warfarin dose or including genes involved in warfarin pharmacogenetics) are shown in red, submissions that did not use the warfarin confounding in any way are shown in blue. The error bars indicate the 95% confidence interval of the AUC. AUC, area under the ROC curve

for downstream analysis. Half of the 14 submissions used an unsupervised approach, clustering the variants from the selected genes using a variety of approaches. Clustering methods used by participants included the principle component analysis, k-means clustering, and a single submission using a deep learning-based approach with autoencoders. Six of the groups used scoring-based methods to the variants within the selected genes to calculate an overall burden score for each subject. A single submission did not use the genotype data at all and trained a logistic regression classifier to predict VTE status based clinical covariates.

The dataset was originally collected to study the genetics of warfarin dosage and had been previously published on and the original publication reports that VTE status is significantly associated with warfarin dosage (Daneshjou et al., 2016, Figure S1). Warfarin dosage was provided to the participants as a covariate for each subject. The known relationship between VTE status and warfarin dose in the dataset was exploited by several groups in their predictions. The most extreme case classified individuals as patients with VTE if they were on a high warfarin dose, and classified individuals on a low warfarin dose as AF. This was the best performing method overall achieving 72% accuracy. Because warfarin dosage is largely influenced by genetics, several groups included genes involved in warfarin pharmacokinetics and pharmacodynamics in their models. Overall, five of the 14 submissions used knowledge that warfarin dosage is associated with VTE status in this dataset in some form.

Of the nine submissions that did not use warfarin dosage to inform their predictions, all utilized either an unsupervised, clustering-based, approach to distinguish the two classes, or used various methods to score variants on the basis of predicted deleteriousness. The top performing group that did not inform their predictions with warfarin dosage information achieved an AUC of 0.65. This method selected genes associated with VTE, PE, and deep vein thrombosis (25 genes total) from DisGeNET and performed k-means clustering on variants determined to be non-neutral by SNAP (Bromberg,

Yachdav, & Rost, 2008; Piñero et al., 2017). The distribution of AUC scores for all predictions can be seen in Figure 2 and a complete list of the scores for each submission is presented in Table 1. Details of all prediction methods can be found in the Supporting Information material.

A baseline prediction accuracy was generated using a linear model proposed by Soria et al. (2014) The baseline model outperformed all submissions that did not use warfarin, achieving a prediction accuracy of 67% and an AUC of 0.71 (Figure 3). Allele frequencies of the SNPs used in the baseline model are shown in Table 2 for individuals of European descent (on whom the genetic risk score was developed), African descent, and the allele frequency observed in the study population.

## 4 | DISCUSSION

This CAGI exome prediction challenge has yielded several insights into the genetics of VTE in African Americans as well as insights into the challenges conducting prediction challenges.

Predicting VTE risk from genetic sequence is a difficult task and the challenge offered in CAGI 5 was no exception. Participants were asked to differentiate exomes of individuals suffering from VTE and those who may be treated with warfarin for a different indication. This task was further complicated by the lack of training data to for participants to validate their proposed methods. This led most participants to develop methods using existing biological knowledge to perform feature selection.

Most participants opted to use clustering-based approaches to predict VTE status. This was a prudent decision given the lack of training data and the stated goal of distinguishing two traits. The other common approach was to score variants within genes based on their predicted deleteriousness, then to create a final score for each individual based on the number of deleterious variants. Although the best performing method used a clustering approach (k-means), there

**TABLE 1** Evaluation metrics for all submissions and for the baseline method. The table is broken up by submissions that used the known warfarin confounding, those that did not, and the baseline method. Within each group scores are sorted by AUC. Accuracy, sensitivity, specificity, and F1 are calculated using a cutoff of 0.5 for all predictions

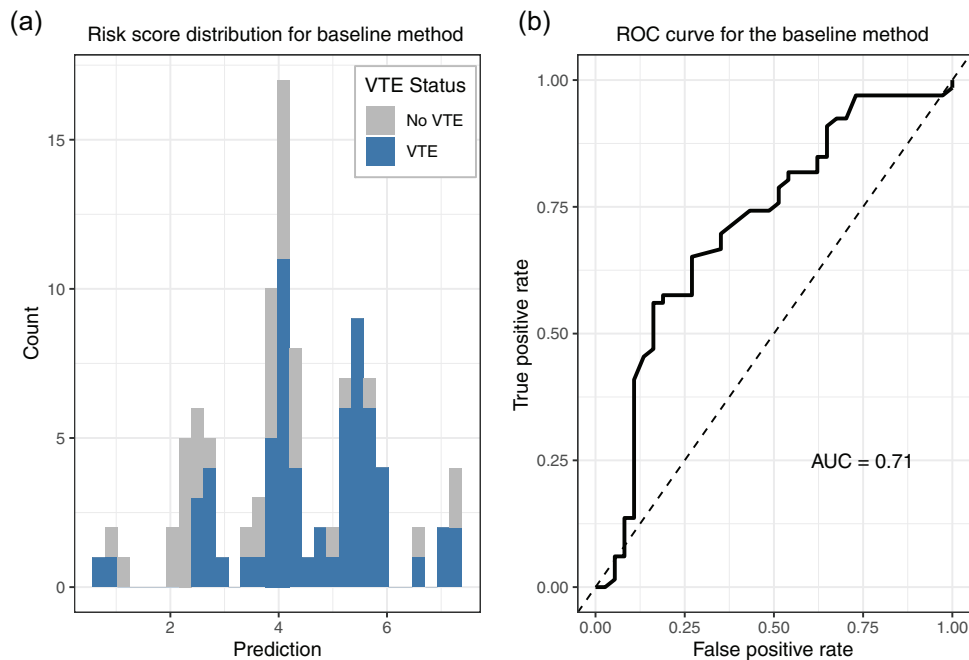
Description	Submission	Approach	AUC $\pm$ 95% CI	AUC	Accuracy	Sensitivity	Specificity	F1
Did not use warfarin in prediction	Group 5a	Scoring	0.65 $\pm$ 0.11	0.65	0.51	0.26	0.95	0.40
	Group 1b	Scoring	0.60 $\pm$ 0.11	0.60	0.60	0.59	0.59	0.65
	Group 5c	Unsupervised	0.59 $\pm$ 0.10	0.59	0.63	0.70	0.49	0.70
	Group 3	Scoring	0.59 $\pm$ 0.12	0.59	0.34	0.23	0.54	0.31
	Group 1a	Scoring	0.57 $\pm$ 0.12	0.57	0.47	0.30	0.76	0.42
	Group 5d	Unsupervised	0.53 $\pm$ 0.10	0.53	0.59	0.73	0.32	0.69
	Group 2	Scoring	0.49 $\pm$ 0.12	0.49	0.41	0.12	0.92	0.21
	Group 7a	Unsupervised	0.48 $\pm$ 0.08	0.48	0.41	0.21	0.76	0.31
	Group 5b	Unsupervised	0.47 $\pm$ 0.09	0.47	0.53	0.65	0.30	0.64
	Used warfarin in prediction	Group 4	Supervised	0.76 $\pm$ 0.10	0.76	0.70	0.71	0.65
Group 6a		Scoring	0.65 $\pm$ 0.09	0.65	0.72	0.85	0.46	0.79
Group 6d		Unsupervised	0.61 $\pm$ 0.10	0.61	0.64	0.70	0.51	0.71
Group 6c		Unsupervised	0.44 $\pm$ 0.10	0.44	0.47	0.53	0.35	0.56
Group 6b		Unsupervised	0.43 $\pm$ 0.10	0.43	0.47	0.56	0.30	0.57
Soria et al.	Baseline	Genetic risk score	0.71 $\pm$ 0.11	0.71	0.67	0.68	0.65	0.73

Abbreviation: AUC, area under the ROC curve.

was no clear advantage to using clustering methods over scoring methods.

All groups subset the exome to genes known to be associated with VTE or AF to use for downstream predictions. The groups with the top two highest scoring submissions both used DisGeNET to select phenotype associated genes. There is clear value in limiting the

search space of the genome and DisGeNET seems to be a useful asset for selecting phenotype associated genes. Groups 1 and 5 (which accounted for five of the top six submissions that did not use warfarin), both used DisGeNET to select genes for their predictions but then applied different methods to predicting subject status from those genes.



**FIGURE 3** Performance of the baseline method on prediction of venous thromboembolism. Here we show the risk scores and predictive performance of the genetic risk score developed by Soria et al. We show the distribution of risk scores across all patients with subjects with VTE shown in blue and those without VTE shown in gray (left). We also show a ROC curve to illustrate the predictive performance of the baseline method (right). AUC, area under the ROC curve; VTE, venous thromboembolism

**TABLE 2** Allele frequencies of variants used in Soria genetic risk model. Here we show allele frequencies for the 13 SNPs used in the genetic risk model used as a baseline. We show allele frequencies shown for European and African populations, which are derived from gnoMAD, and observed allele frequencies for each site in the study population of African Americans used for the CAGI challenge. Allele frequencies for sites outside the exome capture region are assigned a value of N/A

SNP	AF in Europeans	AF in Africans	AF in this study
rs6025 (NC_000001.11:g.169549811C>T)	0.97	0.99	0.99
rs118203905 (NC_000001.11:g.169555300T>C)	0.0	0.0	0.0
rs118203906 (NC_000001.11:g.169555299C>G)	$1.7 \times 10^{-4}$	0.0	0.0
rs1799963 (NC_000011.10:g.46739505 G>A)	$1.3 \times 10^{-2}$	$2.8 \times 10^{-3}$	N/A
rs8176719 (NC_000009.12:g.133257521_133257522insC)	0.37	0.31	N/A (INDELS not called)
rs1801020 (NC_000005.10:g.177409531A>G)	0.75	0.55	0.55
rs5985 (NC_000006.12:g.6318562C>T)	0.25	0.18	0.18
rs2232698 (NC_000014.9:g.94290332 G>A)	$6.5 \times 10^{-3}$	$2.0 \times 10^{-3}$	0.0
rs121909548 (NC_000001.11:g.173904038C>G)	$1.4 \times 10^{-3}$	$1.6 \times 10^{-4}$	0.0
rs2036914 (NC_000004.12:g.186271327T>C)	0.54	0.65	N/A
rs2066865 (NC_000004.12:g.154604124 G>A)	0.24	0.30	0.38
rs710446 (NC_000003.12:g.186742138T>C)	0.41	0.52	0.56
rs2289252 (NC_000004.12:g.186286227C>T)	0.40	0.26	N/A

Abbreviation: CAGI, critical assessment of genome interpretation.

One major point of contention in conducting this challenge was the confounding effect of warfarin on the prediction task. Participants were asked to make predictions about VTE status with the given dataset but were not advised to avoid using warfarin dosage which may be a confounding variable unique to this dataset and not broadly correlated with VTE risk. It is therefore reasonable that the participants sought the best performance possible and used the previously published correlation between VTE and warfarin dose in this dataset. This does, however, go against the spirit of the CAGI experiment, which is meant to better understand the impact of genetics on phenotype and this challenge in particular aimed to improve our collective ability to predict VTE from genetic data. For this reason, we have divided the prediction results between those who used warfarin dosage information, either directly (using the subjects warfarin dose) or indirectly (through the inclusion of warfarin dose related genes), in their models.

The point was raised at the CAGI conference that if this confounding factor was known, why give the warfarin dosage to participants at all? This was because of a miscommunication between the data providers and the conference organizers. However, it may have made little difference as the participants were provided with the entire exome and there is a strong genetic relationship between warfarin dosage and genetics. In hindsight, it may have been better for the challenge to not provide warfarin dosage to the participants and to remove genes related to warfarin pharmacokinetics and pharmacodynamics from the exomes. Alternatively, it may have been better to explicitly inform the participants to avoid using any knowledge about warfarin in their predictions.

To assess the submitted predictions against an existing gold standard we calculated genetic risk scores for each exome using the method proposed by Soria et al. (2014) The genetic risk scores calculated using this method achieved an AUC of 0.71, greater than

that of any submitted method that did not use warfarin dose in their predictions. This method was developed using data from individuals of European descent and had not been previously validated in individuals of African descent. The AUC achieved by this method in African Americans exceeds the reported AUC in the original study population (0.677), and greatly exceeds the expected theoretical maximum AUC that can be achieved for predicting complex phenotypes from coding variants alone. This is particularly interesting because previous work has shown that genetic risk models developed in European perform on average 25% as well in African populations as they do in European populations (Martin et al., 2019). This suggests that the method proposed by Soria et al may be clinically useful in predicting VTE in African Americans, but would need to be evaluated in a larger cohort.

#### ACKNOWLEDGMENTS

We would like to thank and acknowledge the CAGI planning committee and participants. The CAGI experiment coordination is supported by NIH U41 HG007346 and the CAGI conference by NIH R13 HG006650. G.M. is supported by BD2K grant number T32 LM 012409. Y.B. and Y.W. were supported by the NIH U01 GM115486 grant.

#### CONFLICT OF INTEREST

The authors declare that there is no conflict of interest.

#### ORCID

Gregory McInnes  <http://orcid.org/0000-0002-8033-2499>

Panagiostis Katsonis  <http://orcid.org/0000-0002-7172-1644>

Sadhna Rana  <http://orcid.org/0000-0001-8241-057X>  
 Predrag Radivojac  <http://orcid.org/0000-0002-6769-0793>  
 Sean D. Mooney  <http://orcid.org/0000-0003-2654-0833>  
 John Moulton  <http://orcid.org/0000-0002-3012-2282>  
 Steven Brenner  <http://orcid.org/0000-0001-7559-6185>

## REFERENCES

- Beckman, M. G., Hooper, W. C., Critchley, S. E., & Ortel, T. L. (2010). Venous Thromboembolism. *American Journal of Preventive Medicine*, 38(4), S495–S501. <https://doi.org/10.1016/j.amepre.2009.12.017>
- Bromberg, Y., Yachdav, G., & Rost, B. (2008). SNAP predicts effect of mutations on protein function. *Bioinformatics*, 24, 2397–2398. <https://doi.org/10.1093/bioinformatics/btn435>
- Canela-Xandri, O., Rawlik, K., & Tenesa, A. (2018). An atlas of genetic associations in UK Biobank. *Nature Genetics*, 50, 1593–1599. <https://doi.org/10.1038/s41588-018-0248-z>
- Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M. A., ... Durbin, R. (2011). The variant call format and VCFtools. *Bioinformatics*, 27, 2156–2158. <https://doi.org/10.1093/bioinformatics/btr330>
- Daneshjou, R., Cavallari, L. H., Weeke, P. E., Karczewski, K. J., Drozda, K., Perera, M. A., ... Altman, R. B. (2016). Population-specific single-nucleotide polymorphism confers increased risk of venous thromboembolism in African Americans. *Molecular Genetics & Genomic Medicine*, 4, 513–520. <https://doi.org/10.1002/mgg3.226>
- Daneshjou, R., Gamazon, E. R., Burkley, B., Cavallari, L. H., Johnson, J. A., Klein, T. E., ... Perera, M. A. (2014). Genetic variant in folate homeostasis is associated with lower warfarin dose in African Americans. *Blood*, 124, 2298–2305. <https://doi.org/10.1182/blood-2014-04-568436>
- Daneshjou, R., Wang, Y., Bromberg, Y., Bovo, S., Martelli, P. L., Babbi, G., ... Morgan, A. A. (2017). Working toward precision medicine: Predicting phenotypes from exomes in the Critical Assessment of Genome Interpretation (CAGI) challenges. *Human Mutation*, 38, 1182–1192. <https://doi.org/10.1002/humu.23280>
- DeLong, E. R., DeLong, D. M., & Clarke-Pearson, D. L. (1988). Comparing the areas under two or more correlated receiver operating characteristic curves: A nonparametric approach. *Biometrics*, 44, 837–845. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/3203132>
- Dowling, N. F., Austin, H., Dille, A., Whitsett, C., Evatt, B. L., & Hooper, W. C. (2003). The epidemiology of venous thromboembolism in Caucasians and African-Americans: The GATE Study. *Journal of Thrombosis and Haemostasis: JTH*, 1, 80–87. Retrieved from <http://www.ncbi.nlm.nih.gov/pubmed/12871543>
- Feero, W. G. (2004). Genetic thrombophilia. *Primary Care*, 31, 685–709. <https://doi.org/10.1016/j.pop.2004.04.014>
- Gusev, A., Lee, S. H., Trynka, G., Finucane, H., Vilhjálmsson, B. J., Xu, H., ... Consortium, S. W. E. -S. C. Z. (2014). Partitioning heritability of regulatory and cell-type-specific variants across 11 common diseases. *American Journal of Human Genetics*, 95, 535–552. <https://doi.org/10.1016/j.ajhg.2014.10.004>
- Hernandez, W., Gamazon, E. R., Smithberger, E., O'Brien, T. J., Harralson, A. F., Tuck, M., ... Perera, M. A. (2016). Novel genetic predictors of venous thromboembolism risk in African Americans. *Blood*, 127, 1923–1929. <https://doi.org/10.1182/blood-2015-09-668525>
- Karczewski, K. J., Francioli, L. C., Tiao, G., Cummings, B. B., Alföldi, J., Wang, Q., ... MacArthur, D. G. (2019). Variation across 141,456 human exomes and genomes reveals the spectrum of loss-of-function intolerance across human protein-coding genes. *BioRxiv*, 531210. <https://doi.org/10.1101/531210>
- Marouli, E., Graff, M., Medina-Gomez, C., Lo, K. S., Wood, A. R., Kjaer, T. R., ... Lettre, G. (2017). Rare and low-frequency coding variants alter human adult height. *Nature*, 542, 186–190. <https://doi.org/10.1038/nature21039>
- Martin, A. R., Kanai, M., Kamatani, Y., Okada, Y., Neale, B. M., & Daly, M. J. (2019). Clinical use of current polygenic risk scores may exacerbate health disparities. *Nature Genetics*, 51, 584–591. <https://doi.org/10.1038/s41588-019-0379-x>
- Middeldorp, S., & vanHylckama Vlieg, A. (2008). Does thrombophilia testing help in the clinical management of patients. *British Journal of Haematology*, 143, 321–335. <https://doi.org/10.1111/j.1365-2141.2008.07339.x>
- Piñero, J., Bravo, À., Queralt-Rosinach, N., Gutiérrez-Sacristán, A., Deu-Pons, J., Centeno, E., ... Furlong, L. I. (2017). DisGeNET: A comprehensive platform integrating information on human disease-associated genes and variants. *Nucleic Acids Research*, 45, D833–D839. <https://doi.org/10.1093/nar/gkw943>
- Purcell, S. M., Moran, J. L., Fromer, M., Ruderfer, D., Solovieff, N., Roussos, P., ... Sklar, P. (2014). A polygenic burden of rare disruptive mutations in schizophrenia. *Nature*, 506, 185–190. <https://doi.org/10.1038/nature12975>
- Roberts, L. N., Patel, R. K., & Arya, R. (2009). Venous thromboembolism and ethnicity. *British Journal of Haematology*, 146, 369–383. <https://doi.org/10.1111/j.1365-2141.2009.07786.x>
- Rosendaal, F. R., & Reitsma, P. H. (2009). Genetics of venous thrombosis. *Journal of Thrombosis and Haemostasis*, 7, 301–304. <https://doi.org/10.1111/j.1538-7836.2009.03394.x>
- DeRubeis, S., He, X., Goldberg, A. P., Poultney, C. S., Samocha, K., Ercument Cicek, A., ... Buxbaum, J. D. (2014). Synaptic, transcriptional and chromatin genes disrupted in autism. *Nature*, 515, 209–215. <https://doi.org/10.1038/nature13772>
- Simons, Y. B., Turchin, M. C., Pritchard, J. K., & Sella, G. (2014). The deleterious mutation load is insensitive to recent population history. *Nature Genetics*, 46, 220–224. <https://doi.org/10.1038/ng.2896>
- Soria, J. M., Morange, P., Vila, J., Souto, J. C., Moyano, M., Tréguoët, D., ... Elosua, R. (2014). Multilocus genetic risk scores for venous thromboembolism risk assessment. *Journal of the American Heart Association*, 3, e001060. <https://doi.org/10.1161/JAHA.114.001060>
- Wray, N. R., Yang, J., Goddard, M. E., & Visscher, P. M. (2010). The genetic interpretation of area under the ROC curve in genomic profiling. *PLoS Genetics*, 6, e1000864. <https://doi.org/10.1371/journal.pgen.1000864>
- Zakai, N. A., & McClure, L. A. (2011). Racial differences in venous thromboembolism. *Journal of Thrombosis and Haemostasis*, 9, 1877–1882. <https://doi.org/10.1111/j.1538-7836.2011.04443.x>

## SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section.

**How to cite this article:** McInnes G, Daneshjou R, Katsonis P, et al. Predicting venous thromboembolism risk from exomes in the Critical Assessment of Genome Interpretation (CAGI) challenges. *Human Mutation*. 2019;1–7. <https://doi.org/10.1002/humu.23825>