

RESEARCH ARTICLE

Transfer learning enables prediction of *CYP2D6* haplotype function

Gregory McInnes¹, Rachel Dalton^{2,3}, Katrin Sangkuhl⁴, Michelle Whirl-Carrillo⁴, Seung-been Lee⁵, Philip S. Tsao^{6,7}, Andrea Gaedigk^{8,9}, Russ B. Altman^{4,10*}, Erica L. Woodahl^{2*}

1 Biomedical Informatics Training Program, Stanford University, Stanford, California, United States of America, **2** Department of Biomedical and Pharmaceutical Sciences, University of Montana, Missoula, Montana, United States of America, **3** Department of Biomedical and Translational Research, University of Florida, Gainesville, Florida, United States of America, **4** Department of Biomedical Data Science, Stanford University, Stanford, California, United States of America, **5** Department of Genome Sciences, University of Washington, Seattle, Washington, United States of America, **6** VA Palo Alto Epidemiology Research and Information Center for Genomics, VAPAHCS, Palo Alto, California, United States of America, **7** Department of Medicine, Stanford University School of Medicine, Stanford, California, United States of America, **8** Division of Clinical Pharmacology, Toxicology, and Therapeutic Innovation, Children's Mercy Kansas City, Kansas City, Missouri, United States of America, **9** School of Medicine, University of Missouri-Kansas City, Kansas City, Missouri, United States of America, **10** Departments of Bioengineering, Genetics, and Medicine, Stanford University, Stanford, California, United States of America

* raltman@stanford.edu (RBA); erica.woodahl@umontana.edu (ELW)



OPEN ACCESS

Citation: McInnes G, Dalton R, Sangkuhl K, Whirl-Carrillo M, Lee S-b, Tsao PS, et al. (2020) Transfer learning enables prediction of *CYP2D6* haplotype function. *PLoS Comput Biol* 16(11): e1008399. <https://doi.org/10.1371/journal.pcbi.1008399>

Editor: Avner Schlessinger, Icahn School of Medicine at Mount Sinai, UNITED STATES

Received: March 1, 2020

Accepted: September 24, 2020

Published: November 2, 2020

Copyright: © 2020 McInnes et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: The authors confirm that all data underlying the findings are fully available without restriction. Hubble.2D6 is freely available to use on GitHub (<https://github.com/gregmcinnes/Hubble2D6>). Data used in this study is available in various forms. Three training stages were used, and a different dataset was used at each stage. • Pre-training data • The first stage of pre-training used simulated data. Scripts used to generate simulated data can be found on GitHub (https://github.com/gregmcinnes/cyp2d6_simulator) or the generated data can be found on zenodo ([10.5281/zenodo.3951095](https://doi.org/10.5281/zenodo.3951095)). • The second

Abstract

Cytochrome P450 2D6 (*CYP2D6*) is a highly polymorphic gene whose protein product metabolizes more than 20% of clinically used drugs. Genetic variations in *CYP2D6* are responsible for interindividual heterogeneity in drug response that can lead to drug toxicity and ineffective treatment, making *CYP2D6* one of the most important pharmacogenes. Prediction of *CYP2D6* phenotype relies on curation of literature-derived functional studies to assign a functional status to *CYP2D6* haplotypes. As the number of large-scale sequencing efforts grows, new haplotypes continue to be discovered, and assignment of function is challenging to maintain. To address this challenge, we have trained a convolutional neural network to predict functional status of *CYP2D6* haplotypes, called Hubble.2D6. Hubble.2D6 predicts haplotype function from sequence data and was trained using two pre-training steps with a combination of real and simulated data. We find that Hubble.2D6 predicts *CYP2D6* haplotype functional status with 88% accuracy in a held-out test set and explains 47.5% of the variance in *in vitro* functional data among star alleles with unknown function. Hubble.2D6 may be a useful tool for assigning function to haplotypes with uncurated function, and used for screening individuals who are at risk of being poor metabolizers.

Author summary

CYP2D6 is an enzyme expressed in the liver that is responsible for metabolizing more than 20% of clinically used drugs. The gene that encodes *CYP2D6* is highly polymorphic, with more than 130 haplotypes currently known. Genetic variants in *CYP2D6* can

stage of pre-training used sequence data from liver microsomes. Summary statistics for this data can be found on zenodo ([10.5281/zenodo.3951095](https://doi.org/10.5281/zenodo.3951095)). • The data to train the final model was downloaded from PharmVar as a VCF, version 4.1.1 (<https://www.pharmvar.org/gene/CYP2D6>). The VCFs were processed as described in the Methods.

Funding: G.M. is supported by the Big Data to Knowledge (BD2K) from the National Institutes of Health (T32 LM012409). E.L.W. and R.D. are supported by the Northwest Alaska-Pharmacogenomics Research Network (NWA-PGRN) (P01GM116691). K.S., M.W.C., and R.B.A. are supported by NIH/NHGRI PharmGKB resource, (U24 HG010615). R.B.A. is also supported by the Chan Zuckerberg Biohub and NIH GM102365. A.G. is supported by the National Institutes of Health for the Pharmacogene Variation Consortium (R24GM123930). P.S.T. is supported by 1HL-101388 (NIH-NHLBI). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors of this manuscript have the following competing interests: RBA is a stockholder in [Personalis.com](https://www.personalis.com) and [23andme.com](https://www.23andme.com).

drastically affect how an individual responds to pharmaceuticals metabolized by CYP2D6, affecting both drug safety and efficacy. Guidelines have been developed to help providers better prescribe drugs using *CYP2D6* genotype, but because of the numerous variants that have been identified in *CYP2D6* these guidelines aren't suitable for all patients. The guidelines rely on the fact that the metabolic function of *CYP2D6* haplotypes is known, however, this is not the case for more than half of all existing haplotypes. Currently, we rely on human curation of haplotypes to assess function, but as genome sequencing has become routine the number of haplotypes identified has grown rapidly, making human curation untenable. We have developed a system using machine learning, named Hubble.2D6, that seeks to address this problem by making predictions of metabolic function for *CYP2D6* haplotypes. We trained a deep learning model using very little data by first pretraining the model on simulated data, allowing for the use of a complex model in a domain that is typically data limited. We find that Hubble.2D6 predicts haplotype function with high accuracy and that the predictions correlate with *in vitro* functional characterization. We believe that this predictive model can be helpful in assessing function of *CYP2D6* haplotypes with unknown function, and that the method of pretraining on simulated data may be useful in other domains within genetics.

Introduction

Cytochrome P450 family 2, subfamily D, polypeptide 6 (*CYP2D6*), is one of the most important pharmacogenes. The protein, a hepatic enzyme, metabolizes more than 20% of clinically used drugs including antidepressants, antipsychotics, opioids, antiemetics, antiarrhythmics, β -blockers, and cancer chemotherapeutics [1–3]. *CYP2D6* is highly polymorphic, making it challenging to address clinically[4,5]. More than 130 haplotypes comprised of single nucleotide variants (SNVs), insertions and deletions (INDELs), and structural variants (SVs) have been discovered and catalogued in the Pharmacogene Variation Consortium (PharmVar; www.pharmvar.org), many of which are known to alter enzymatic activity and protein expression levels[6,7]. Individuals are grouped by their *CYP2D6* metabolic function and are typically classified into one of four metabolizer (or phenotype) groups: normal (NM), intermediate (IM), poor (PM), and ultrarapid metabolizers (UM). Phenotype frequencies vary widely among global populations with PMs ranging from 0.5% to 5.4%, IMs ranging from 2.8% to 11%, and UMs ranging from 1.4% to 21.2%[8].

Despite its highly polymorphic nature, *CYP2D6* is one of the most clinically actionable pharmacogenes. A standardized method to translate *CYP2D6* genotype to phenotype has been recommended by the Clinical Pharmacogenomics Implementation Consortium (CPIC) and the Dutch Pharmacogenetics Working Group (DPWG)[9]. Clinical guidelines providing dosage recommendations for different phenotype groups have been published by CPIC and DPWG for drugs metabolized by *CYP2D6*, including opioids, selective serotonin reuptake inhibitors, tricyclic antidepressants, the attention-deficit/hyperactivity disorder drug atomoxetine, the estrogen receptor modulator tamoxifen, among others[10–16]. Following the guidelines for these drugs could improve patient outcomes by decreasing adverse effects or increasing efficacy[17]. It has even been suggested that pharmacogenetics-guided opioid therapy could be part of a solution for combating the opioid epidemic[18]. The Centers for Medicare and Medicaid Services, in addition to a major insurance company, recently announced coverage for *CYP2D6* genetic testing—among other genes—for improved selection of

antidepressants and other drugs carrying CPIC recommendations, marking a major advancement in the incorporation of pharmacogenomics in patient care[19,20].

Pharmacogenetic dosing guidelines presume that the clinician has access to the patient's *CYP2D6* genotype and that the resulting phenotype can be predicted with accuracy. The system used to translate *CYP2D6* genotype into phenotype is known as the activity score (AS) system, which has been widely adopted and is utilized by CPIC and DPWG[9,21]. *CYP2D6* haplotypes are named using the star allele nomenclature curated by PharmVar, which defines the core variants for each star allele[5–7]. Core variants are typically coding variants, but can also be functionally important variants such as splice junction variants. The AS works by first assigning a value to each star allele (0 for no function alleles, 0.25 or 0.5 for decreased function alleles, and 1 for normal function alleles; gene duplications receive double the value of their single counterpart), then summing the values assigned to each allele. The resulting AS for the person's pair of haplotypes (or diplotype) is used to determine the *CYP2D6* phenotype, which can then be used to inform treatment decisions[9].

The scores used for a *CYP2D6* haplotype's contribution to an individual's AS rely heavily on the manual curation of star allele function through a review of the literature. Most often, *in vitro* experiments and *in vivo* phenotype measures are used to make a determination of star allele function. Even where *in vitro* functional studies exist, however, it can be difficult to assess haplotype function due to variability between expression systems and substrates.

The current reliance on manual curation for scoring of *CYP2D6* haplotypes limits the ultimate utility of *CYP2D6* phenotype prediction for pharmacogenetic guidance of treatment decisions. It is estimated that individuals carrying *CYP2D6* haplotypes with an unknown, uncertain, or uncurated function (herein referred to collectively as uncurated) range from 2 to 9%, with a study of the UK Biobank finding that 3.4% of individuals carry haplotypes that cannot be mapped to a predefined function[22,23]. Individuals carrying a haplotype with uncurated function cannot be assigned to a distinct metabolizer group using the AS and are instead labeled as "Indeterminate". Therefore, for these patients, pharmacogenomic-guided therapy for drugs metabolized by *CYP2D6* (e.g. CPIC and DPWG dosing guidelines) cannot be used. In fact, there are currently over 70 star alleles in PharmVar with unknown, uncertain, or uncurated function. Additionally, the extent of the true population level variation in *CYP2D6* is likely far greater than that which can be explained by existing star alleles. In gnomAD, a large aggregate database of genomes, 544 nonsynonymous SNVs and INDELS are identified in *CYP2D6*, and only 98 of those are included in existing star alleles in PharmVar[24]. *CYP2D6* is known to harbor other rare variants that may be functionally important that are not yet included in star alleles[22,25,26]. Further work is required, however, to determine whether these SNVs can truly be ascribed to *CYP2D6* or whether they are misaligned from *CYP2D7*, a highly homologous pseudogene of *CYP2D6*.

With the ever-increasing amount of sequencing data being generated, the number of novel haplotypes of *CYP2D6* keeps increasing making it even more difficult to generate functional data that fulfil the criteria for function assignment recently described by CPIC and DPWG [27]. The current framework of manual curation of literature in order to assign function to newly discovered star alleles will be challenged to keep up with the rate at which star alleles are discovered, since there will be a lag between the discovery of the star alleles and the generation of functional data for rare alleles.

Methods for predicting variant deleteriousness *in silico* are abundant, but these methods are often developed to be of general purpose for functional prediction of single variants. Many methods have been developed for predicting whether a single variant is likely to damage protein function[28], and several have been developed specifically to predict the impact of single variants in pharmacogenes[29]. For predicting the function of *CYP2D6* star alleles, however,

impact of all of variants combined must be considered opposed to single variants. There are existing methods that can predict function for pairs of variants[30], but *CYP2D6* star alleles can have as many as ten core variants (e.g. *CYP2D6**57). To predict the function of *CYP2D6* haplotypes *in silico*, a purpose-built tool is needed. Tools have been developed to assign *CYP2D6* star alleles to sequence data, but these tools do not predict star allele function[31,32]. A recent publication highlighted a deep learning approach to predicting the function of *CYP2D6* drug metabolism using a fully connected neural network[33]. This method predicts metabolic function on a continuous scale using known *CYP2D6* variants as input, but cannot evaluate the impact of variants not captured by existing star alleles.

The rise of big data and machine learning, in particular deep learning, has revolutionized computer vision and is being successfully applied to applications in genetics[34]. Deep learning presents an attractive solution for making functional predictions about variation in highly polymorphic genes, like *CYP2D6*. Deep learning algorithms require massive amounts of data to be used effectively, however, and using deep learning to address all problems in genetics is not straightforward. Many domains in genetics are data limited and therefore challenging to address with complex algorithms. Data scarcity has been addressed in computer vision through the use of transfer learning, where a neural network is pretrained on a large corpus of data, then adapted to a new domain by transferring the coefficients learned by the network to a new network. This is particularly useful in cases where there are spatial or sequential motifs that are shared between the source and target tasks. A neural network can be pre-trained using real data from a related task, or simulated data labeled using an existing knowledge source.

Neural networks offer flexibility in the way data is represented. For computer vision tasks, data is represented using three stacked matrices, or channels, representing the RGB channels used to color pixels[35]. Genetic data is represented using a one-hot encoded 4xN matrix, where there are four channels for each of the possible nucleotides and N represents the sequence length[34]. It is possible to include additional information about each position in the sequence by adding additional channels to this input matrix. This flexibility allows us to include variant annotation information, such as whether a variant is in a coding region or whether it is predicted to be deleterious. A functional representation of genetic variants such as this may provide context for observed variants that the neural network can use to reason about newly observed variants.

Here we present, Hubble.2D6 which predicts a functional phenotype of *CYP2D6* haplotypes from DNA sequence data using a convolutional neural network (CNN). The model predicts whether a *CYP2D6* haplotype will have normal, decreased, or no function. We validated our model using *in vitro* studies from literature of 46 star alleles, of which 37 had not previously been seen by the model. We generated predictions for 71 *CYP2D6* star alleles without curated function assignments. We trained our model using transfer learning and a functional variant representation, and show that these features greatly improve the ability of the network to learn to predict function.

Results

We trained Hubble.2D6, a deep learning model that predicts star allele functional phenotype from DNA sequence, to classify *CYP2D6* star alleles as normal, decreased, or no function (Fig 1). Hubble.2D6 was trained to classify function on a set of star alleles for which a functional label has been assigned by curators. We used a training set of 31 star alleles and a separate set of 25 star alleles for model validation. All 56 star alleles used for training and model validation have a curated function assignment. Hubble.2D6 correctly predicts the function of 100% of the 31 star alleles used for training and 88.0% of the 25 star alleles used as a held-out validation set

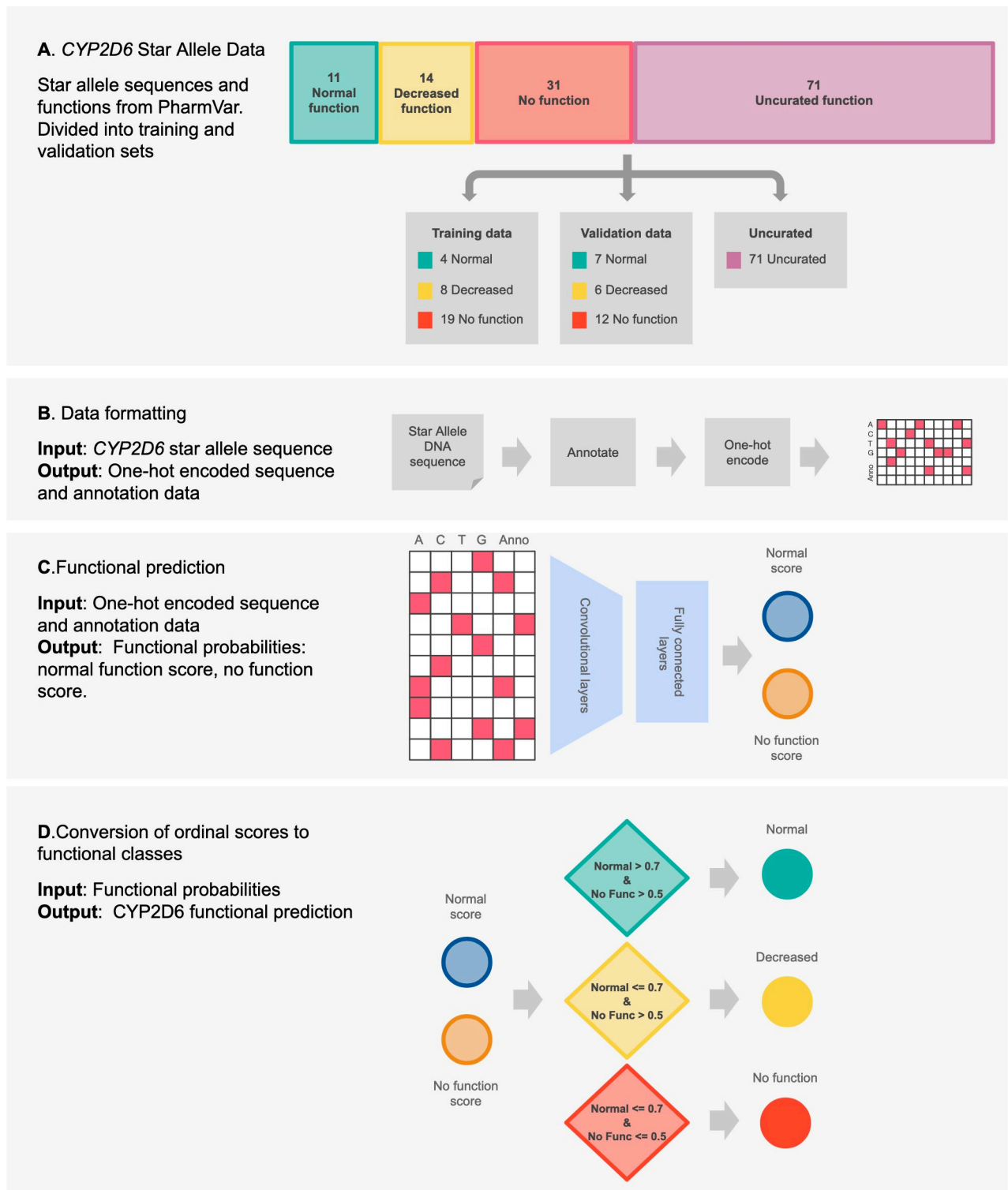


Fig 1. Schematic overview of the Hubble.2D6 workflow. (A) Sequences and functions for all existing star alleles in PharmVar were collected and divided into training and validation datasets. Star alleles with uncurated function were held from training. (B) Star allele sequences were annotated with functional annotations and one-hot encoded as preparation for input into the deep learning model. (C) One-hot encoded sequence and annotation data was read into a convolutional neural network that output scores for two classes: a score indicating a normal function allele, and a score indicating a no function allele. (D) The two score outputs from the model were transformed into one of the three functional classes using cutoffs that were set to optimize sensitivity and specificity in the training data.

<https://doi.org/10.1371/journal.pcbi.1008399.g001>

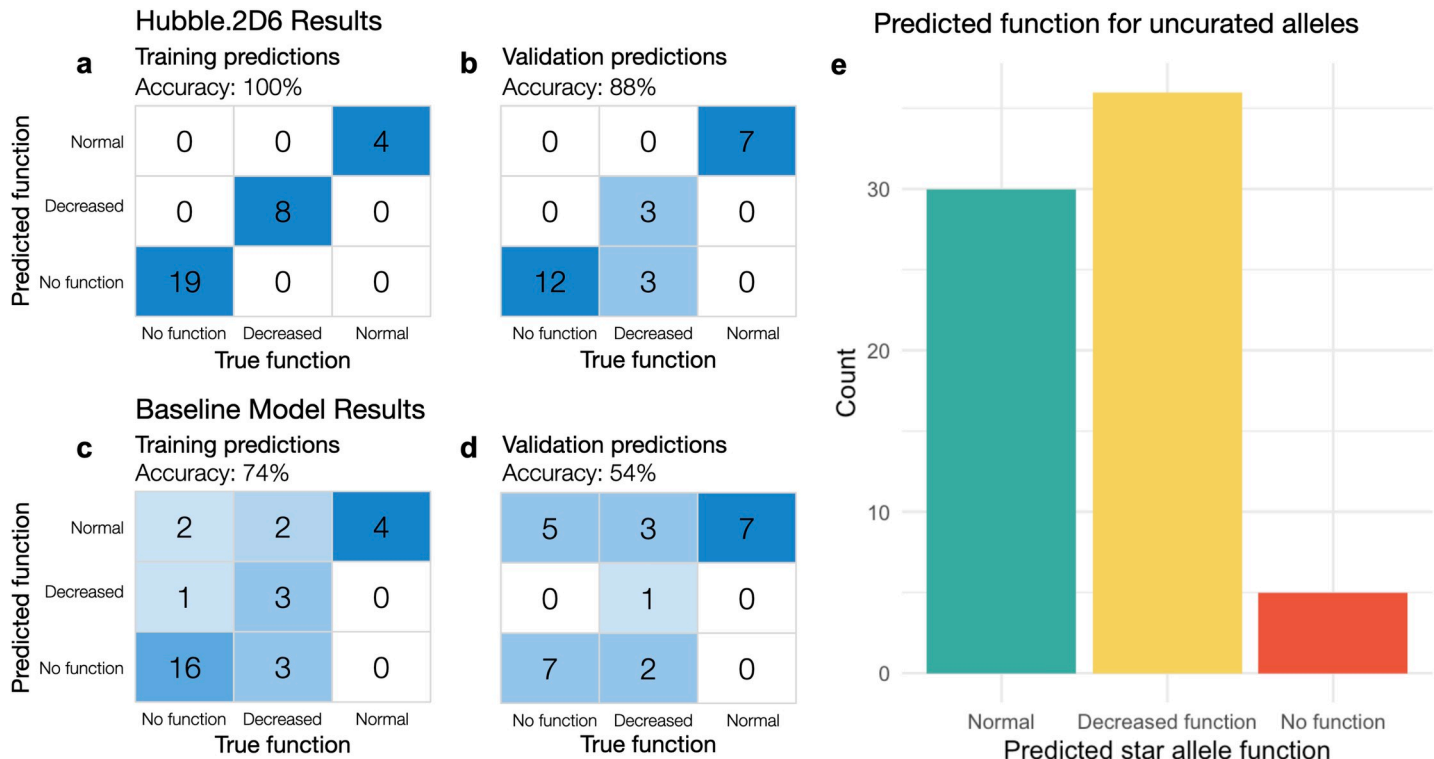


Fig 2. Star allele classification results. The figure depicts performance metrics for the prediction of star allele function in the training and validation sets; confusion matrices for class prediction in training and validation are shown in (a) and (b), for Hubble.2D6 and in (c) and (d) for the baseline model. (e) shows the frequency of predicted function for uncurated star alleles.

<https://doi.org/10.1371/journal.pcbi.1008399.g002>

(Fig 2A and 2B). The only misclassifications among samples with curated function are two decreased function alleles in the test set that were predicted to be no function alleles (i.e. *CYP2D6*14* and *CYP2D6*72*). Of the 71 star alleles with unknown, uncertain or not yet curated function, 30 were predicted to have normal function, 36 decreased function, and 5 were predicted to have no function, although the true function of these star alleles remains uncurated (Fig 2E). Predictions for all investigated star alleles, including those with uncurated function, are provided in S1 Table.

The weights in the convolutional and fully connected layers in Hubble.2D6 were trained using two pre-training steps (S1 Fig). Each pre-training step uses the same *CYP2D6* region as input and variant representation as the final model. First, we trained a model to classify the *CYP2D6* activity score using 50,000 simulated pairs of star alleles. This model predicted the activity score of an additional set of 10,000 simulated pairs of star alleles with 100% accuracy. The weights from the convolutional layers were then transferred to a new network, which was trained to predict the measured metabolic function for a *CYP2D6* substrate, dextromethorphan, in 249 liver microsomes[36], a regression problem. This model explained 61% of the variance in the training data and 70% of the variance in a held-out set of 59 samples. The weights from the convolutional and fully connected layers were then transferred to the final model and fine-tuned using the star allele data.

We created an ordinal logistic regression model as a baseline to compare Hubble.2D6 against. The baseline model was trained using counts of each annotation in the functional variant representation as input, such that it would have access to the same annotation data as Hubble.2D6. The baseline model correctly predicted the function of 74% of the 31 star alleles

used for training and 54% of the 25 star alleles used as a held-out validation set (Fig 2C and 2D).

We evaluated our predictions with *in vitro* data from a study describing functional characterization of 50 CYP2D6 star alleles[37]. Four star alleles were excluded from analysis, CYP2D6*61 and CYP2D6*63 were excluded because they are CYP2D6-2D7 hybrid genes, which are not analyzed by Hubble.2D6. We also excluded the allele formerly known as CYP2D6*14B. Finally, CYP2D6*53 was determined to be an extreme outlier and excluded from this analysis because it had a 28-fold increase in *in vitro* function compared to reference. Of the 46 remaining star alleles characterized by this study, 30 had a CPIC assigned clinical allele function label, and 16 were uncurated. Of the 30 alleles with a curated function and with *in vitro* functional data, 21 were held out from training to be used in the validation set for model evaluation. For these 21 star alleles used for evaluation, our predicted labels explained 71.0% of the variance, approximately equal to the variance explained by the CPIC assigned function labels, 71.1% (Fig 3A and 3B). We also assessed the function of 16 star alleles from this study that have not yet been assigned function by CPIC. For these uncurated alleles, two are predicted to have no function, nine decreased function, and five normal function. Our predicted labels explained 47.5% of the variance in the measured activity in the uncurated alleles, the mean measured *in vitro* activity of each predicted phenotype group for the uncurated alleles were significantly different as a result of a one-way ANOVA ($P = 0.014$, Fig 3C and 3D).

We interpreted the predictions made by Hubble.2D6 by calculating importance scores for each variant in each star allele sequence using DeepLIFT. This allowed us to see the relative

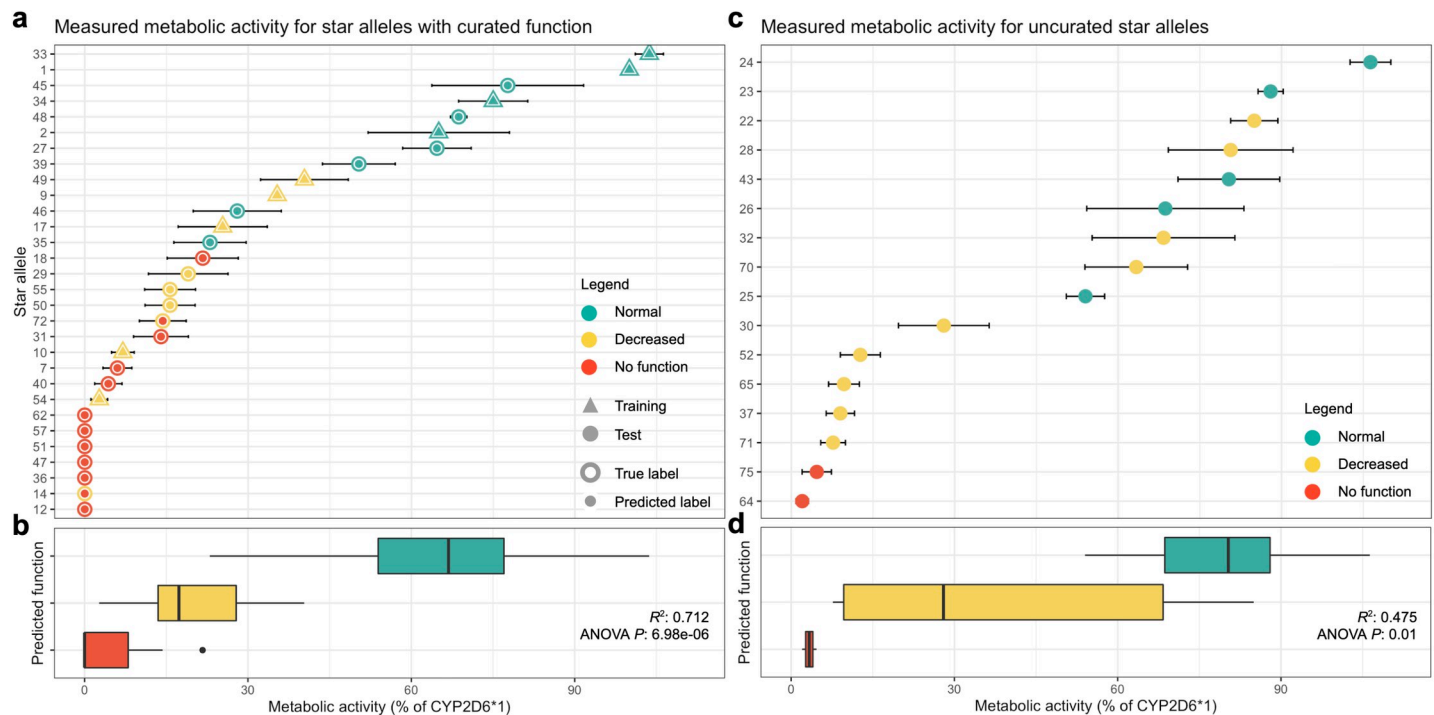


Fig 3. Prediction of star allele function with *in vitro* data. The figures summarize the distribution of metabolic activity measured *in vitro* for star alleles whose function was predicted by Hubble. The distribution of functional activity is shown in (a) and (b) for star alleles with CPIC-assigned clinical function assignments. (a) star alleles included in the training process are depicted with a triangle, and those held for testing are depicted with a circle. Error bars depict the standard error of the measured function. The outer edge of each point indicates the true, curator-assigned phenotype, while the inner color represents predicted function. (b) distribution of values for each predicted functional class for data shown in (a). (c) star alleles without assigned function status; colors represent the predicted function. (d) variance in measured activity of the star alleles for each predicted label for data shown in (c).

<https://doi.org/10.1371/journal.pcbi.1008399.g003>

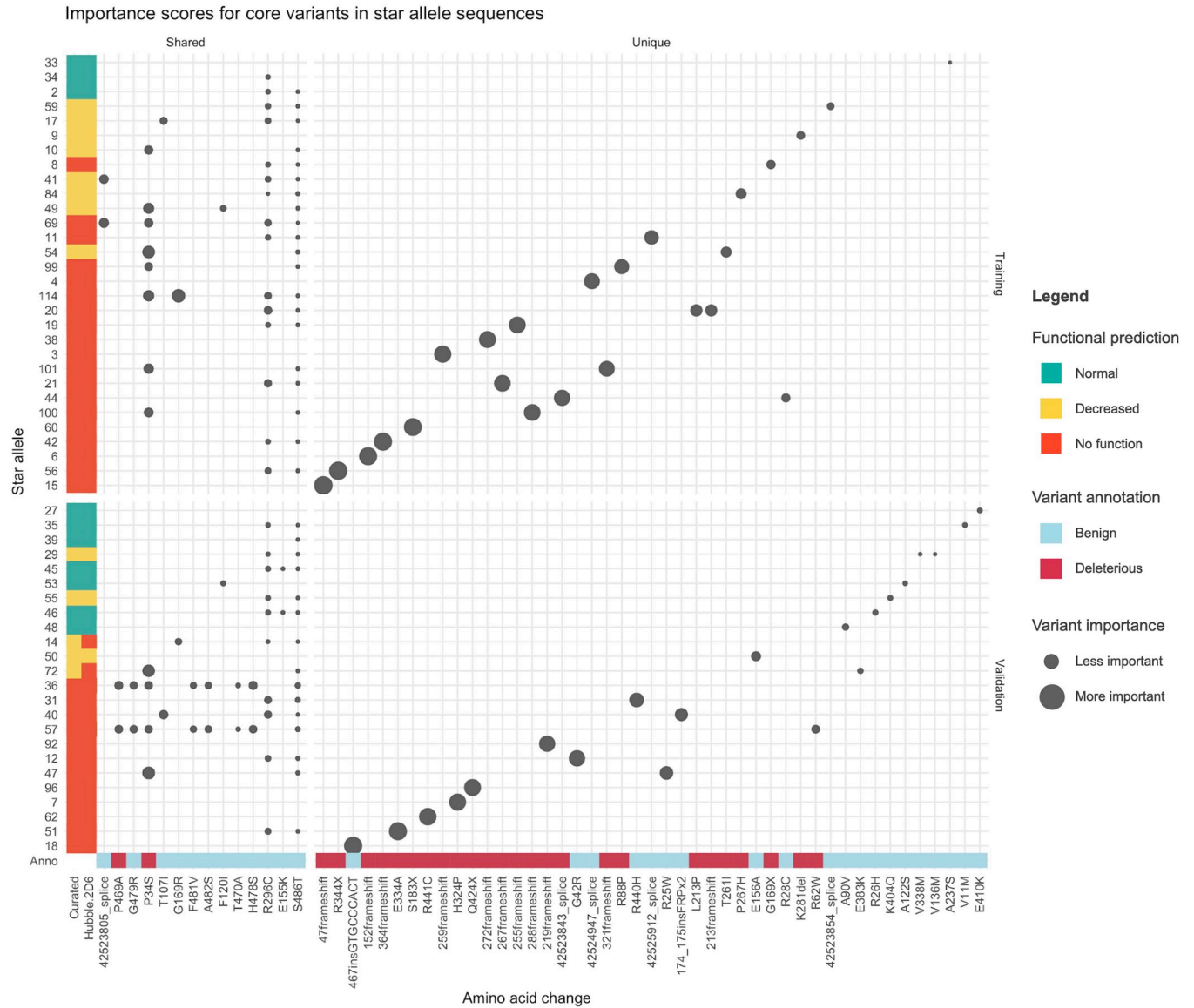


Fig 4. Importance scores for core variants in each star allele used for training and test of Hubble.2D6. Star alleles are along the y-axis and core variants (both amino acid changes and non-coding changes) are listed along the x-axis. Each dot represents the importance of the core variant to the final prediction as determined by DeepLIFT. The size of the dot represents the value of the importance score, with larger dots indicating variants with larger importance scores, typically associated with a negative impact on function. Star alleles are annotated with the curated function as well as the Hubble.2D6 predicted function. Star alleles are sorted by the sum of the importance scores, with those with the largest sums at the bottom. Core variants are divided along the x-axis by those that are uniquely in either the training or test samples (right), and those that are shared between star alleles in train and test (left). Core variants are sorted by their mean importance score across all star alleles. Core variants are annotated with the deleteriousness prediction used in the functional variant representation with red indicating a variant predicted to be deleterious and blue indicating a variant predicted to be benign (described in Methods).

<https://doi.org/10.1371/journal.pcbi.1008399.g004>

importance of each variant to the final prediction (Fig 4). Additionally, we wanted to understand whether the model was relying on core variants shared between star alleles to make predictions, or whether novel variants were driving the predictions. We found that, although there is some overlap in core variants between the train and test groups, most star alleles predicted to be of decreased or no function carried unique variants with large importance scores. Importance scores for uncurated star alleles are shown in S2 Fig.

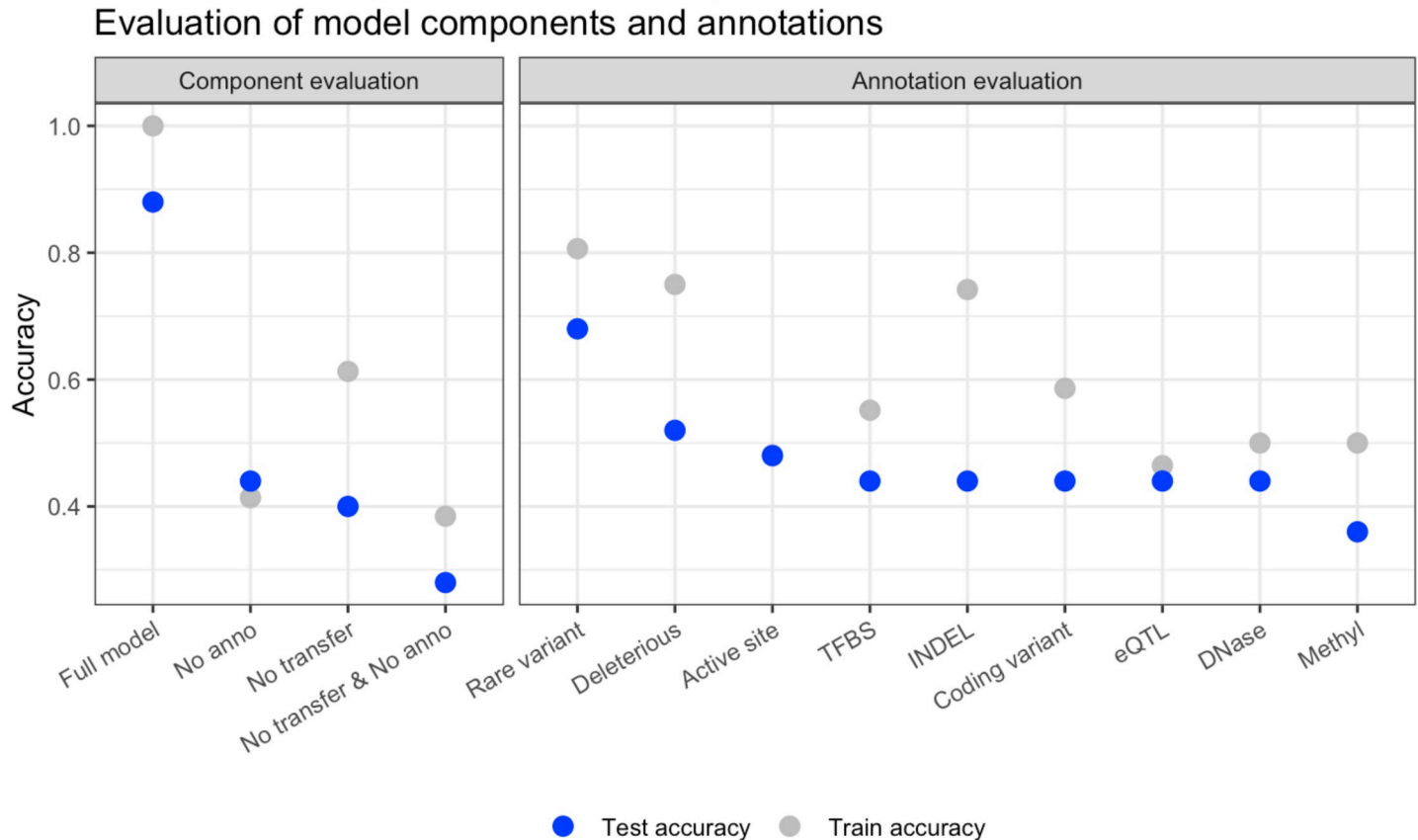


Fig 5. Evaluation of the contribution of deep learning model components. The figure depicts the training and test classification for models trained under various constraints. Under “Component evaluation”, we test the contribution of transfer learning and the inclusion of annotations in the variant encoding by training new models to predict star allele function. Each model is identical in every way to the full Hubble.2D6 model except for the stated difference. We tested the effect of including annotations and transfer learning individually, together, and one model was built with neither component. Under “Annotation evaluation” we depict classification accuracy for models trained with a single added annotation. Each point represents the accuracy of a model trained to predict star allele function using transfer learning with a one-hot encoding of the nucleotide sequence, but only the specified annotation was included in the encoding of the variant. The full model contained all listed annotations together.

<https://doi.org/10.1371/journal.pcbi.1008399.g005>

We evaluated the contribution of transfer learning and using a functional variant representation by training three new models: two leaving out a single component (one without transfer learning, one without the functional variant representation), and one without either component. We compared the classification accuracy of each of these models to the full model that was trained with both components (Fig 5). Test accuracy was 28% when excluding both components, 40% for excluding only transfer learning, and 44% excluding only the functional variant annotation, compared to 88% accuracy for the full model. We evaluated each included annotation in the functional variant representation in a similar way. The annotation that most improved the classification accuracy indicated whether a variant was rare in the population (68% test accuracy), followed by an annotation indicating whether a variant is predicted to be deleterious (52% test accuracy).

Discussion

Here we present Hubble.2D6, a model for predicting functional phenotypes of *CYP2D6* star alleles. Given a star allele sequence, Hubble.2D6 classifies the haplotype as having normal, decreased, or no function. Hubble.2D6 has 88.0% accuracy on a held-out validation set of

existing haplotypes with curated function. Additionally, we found that in star alleles with uncurated function, the functional predictions from Hubble.2D6 explained a significant amount of the variance in *in vitro* functional measurements, indicating that the Hubble.2D6 assigned labels correlate with actual function.

Predictions of star allele function *in silico*, such as the ones output by Hubble.2D6, can be utilized in several applications in the current pharmacogenomic landscape. Functional predictions could serve as an additional source of evidence when assigning function to star alleles. In cases where *in vitro* data is limited or variable, predicted function from Hubble.2D6 may be used to predict a diplotype's function status, analogous to a computer-aided decision tool. High throughput assays to generate *in vitro* activity data for large numbers of variants have been suggested for pharmacogenomic applications but have yet to be run comprehensively for CYP2D6 [38]. Since Hubble.2D6 takes as input the full coding and non-coding sequence of the *CYP2D6* locus, any variant or combination of variants within the *CYP2D6* locus can be rapidly assessed.

Ultimately, functional predictions are anticipated to be part of the clinical prediction of haplotype function in the absence of other data. Currently, if a clinician encounters a patient with a haplotype of uncertain, unknown, or uncurated function in a patient, his/her metabolizer status remains indeterminate. Considering the number of star alleles in one of these categories, as many as 8% of patients will not receive a phenotype assignment. Rather than proceed without pharmacogenetic guidance, the patient's phenotype may be predicted using Hubble.2D6. For example, the patient may be more closely followed, drug level monitoring employed following the initiation of drug therapy, or a drug from a different class selected as a precaution to avoid adverse effects. Such a scenario may be preferable opposed to proceeding blindly without any indication of the patient's metabolizer status. As we strive to make personalized predictions of drug response a reality for every patient, it is important to be able to provide accurate predicted phenotypes for as many patients as possible.

There are four uncurated haplotypes with variants that likely obliterate function. *CYP2D6*81*, *CYP2D6*120*, and *CYP2D6*129* carry nonsense variants, and *CYP2D6*124* has a frameshift insertion. Of these four alleles, only one was predicted to be a no function allele (*CYP2D6*120*), and the other three were predicted to have decreased function, rather than no function. A prediction of "no function" would be more intuitive, as these types of variants are well known to lead to a non-functional protein. The model importance scores show that the model heavily weighted the variant leading to a premature stop in *CYP2D6*120*, however, the presumably loss-of-function variants in the other star alleles received considerably lower importance scores (S2 Fig). Each of the four loss-of-function variants is predicted to be deleterious by LOFTEE[39]. It seems that the model treats newly seen variants conservatively and outputs a prediction of decreased function rather than no function. Although this is an important distinction, providing an indication that the allele will metabolize drugs abnormally is beneficial. There are important cases, such as the conversion of a prodrug to its active form by CYP2D6 (e.g. codeine to morphine), in which the distinction between a decreased function allele and a no function allele is critical. More training data with more varied examples of loss-of-function variants will likely improve future versions of the model.

The variance explained in the *in vitro* data by Hubble.2D6 classifications is substantially lower among star alleles with unknown or uncertain function ($R^2 = 0.475$) compared to those with curator-assigned function ($R^2 = 0.712$). Star alleles with unknown or uncertain function are typically more difficult for curators to assess because of variability in published *in vitro* data, hence the lack of a functional assignment. The increased difficulty for human curators likely translates to a similarly increased challenge for a machine learning system to assess star alleles whose function is less obvious.

An important contribution of this work is the finding that deep learning models can be trained for small-scale problems in genetics. The highly polymorphic nature of *CYP2D6* makes it an attractive application for deep learning, but the limited number of haplotypes with known function to use as training data limits our ability to train large models effectively. We show that through the use of both transfer learning and the inclusion of variant annotations we can train a model to predict haplotype function with high accuracy. The intuition that motivates this is that, first, by pretraining the model on simulated and real data it has learned weights for motifs related to commonly observed variants. This prevents the model from having to learn all motifs from scratch during the training phase. Second, by including variant annotations in the variant encoding the model has been provided with additional useful information. With enough data a model may be able to learn that a certain nucleotide change would be likely to be deleterious, but in a space where data is scarce it may be advantageous to provide the model with that information upfront. As can be seen from the variant importance scores in Fig 4 many of the variants with large importance scores are unique to a single star allele, so it is important for the model to learn the principles that lead to a functional change rather than nucleotide motifs alone. We show that by combining these two methods the result is a more accurate predictor than either method alone.

There are several limitations to our *CYP2D6* predictive models. First, we do not include structural variants in our model. This prevents the model from being able to predict function of star alleles with structural variation including hybrid genes[40] including *CYP2D7-2D6* or *CYP2D6-2D7*, the *CYP2D6*5* gene deletion,[41] as well as gene duplications (e.g. *CYP2D6*1xN*, **2xN*, etc)[42]. Thus, Hubble.2D6 does not predict increased function at this point in time. While these are important types of alleles to consider, all known occurrences of hybrid genes result in a complete loss of function, so *in silico* prediction of function is not necessary, and changes in copy number are easily accounted for with the use of the activity score. In the current star allele definitions, the only source of increased function alleles is due to gene duplication events. Second, since we only considered a narrow range of upstream and downstream sequences of *CYP2D6*, we have not captured distal effects on gene expression, such as a long-distance ‘enhancer’ SNV that has been described to be associated with decreased function [43–45]. Distal regulatory effects are not included by the current model, which may limit our ability to fully explain the observed variability in enzymatic activity using *CYP2D6* sequence alone. However, these distal effects would not be captured by conventional *in vitro* expression studies either. Third, the validation performed used *in vitro* data for only 50 star alleles, and the data was generated using one expression system in a single laboratory. Substantial variance exists in measured functional activity of star alleles between laboratories. Further functional assessment of star alleles compared to Hubble.2D6 predictions could inspire greater confidence in our ability to predict function *in silico*. Fourth, there are known substrate specific effects on *CYP2D6* metabolic function, which are not considered by this approach. We aimed to assign function to star alleles to be in line with the CPIC clinical functional classifications used by the existing guidelines. The current dosing guidelines use star allele functional assignments that are substrate agnostic. This is an important consideration for future work and may improve overall performance given enough data. Finally, there may be factors impacting *CYP2D6* activity outside the gene sequence, such as variation in other genes or whether an individual is taking *CYP2D6* inhibitors such as selective serotonin reuptake inhibitors[46,47]. These factors are not included in our model as we only focus on the genetic variation in the *CYP2D6* locus.

In summary, we have created a model for the prediction of metabolic activity for *CYP2D6* haplotypes from sequence data. We find that our model has high accuracy predicting allelic function, and that our predicted function labels explain a significant amount of the variance

observed in CYP2D6 metabolic activity *in vitro*. This model may be useful to predict phenotype of patients carrying haplotypes, which are either uncurated or have unknown/uncertain function assignments

Methods

Model overview

We trained a multiclass convolutional neural network (CNN) classifier using supervised learning to predict CYP2D6 haplotype function from DNA sequence. The deep learning model used is a three layer CNN following the Basset architecture (S2 Fig)[48]. The output is one of three functional statuses: “No function”, “Decreased function”, or “Normal function”, following the clinical function assignments provided by CPIC and posted by PharmVar, excluding “Increased function” (Fig 1). Hubble.2D6 does not make predictions for increased function alleles because the only increased function alleles identified for CYP2D6 occur as a result of gene duplications; only SNVs and small INDELS are considered. Hubble.2D6 reads in a one-hot encoded matrix representing the full DNA sequence of the haplotype in addition to eight variant level functional annotations and outputs two scores: (1) the probability that the haplotype is a no function allele and (2) the probability that the haplotype is a normal function haplotype. The two scores are then transformed into one of the three functional classes (no function, decreased function, or normal function) using cutoffs that are defined to maximize sensitivity and specificity of the functional predictions in the training data. Specifically, we used the R package *cutpoint* to select the value for each of the two scores that maximized the sensitivity and specificity of the functional predictions of star alleles in the training data[49]. For example, a haplotype receiving a no function probability greater than 0.5 and a normal function probability less than 0.7 is classified as a decreased function allele.

Although only two scores are output they are combined to map the function to one of the three possible outcomes; this setup formats the prediction task as an ordinal regression problem such that the network can learn the ranking of the functional groups[50]. The first score, the no function score, differentiates between no function alleles and alleles with either decreased or normal function. No function alleles are indicated with a 0 as the first score while alleles with any other function are indicated with a 1. Likewise, the second score, the normal score, differentiates between alleles with normal function and alleles with either decreased or no function. Normal function alleles have a score of 1 for the second score while others have a score of 0. Leading to a scoring system where each of the three function classes can be yielded from only two scores. This scoring system is superior to a classification setup that assumes class independence because it allows the network to learn that no function alleles are more similar to decreased function alleles than they are to normal function alleles, and vice versa.

The input into the model is genetic sequence data. For each star allele the data is converted into a one-hot encoded matrix of DNA variation and variant-specific functional annotations for each position in the CYP2D6 locus (chromosome 22, 42,521,567–42,528,984, hg19). The interrogated region includes 2,103 bp of upstream and 934 bp of downstream sequence which corresponds to the region covered by the PGRNseq platform[51]. We used eight variant level annotations that may influence gene expression or protein function to create a functional variant representation. Each base in the capture window was annotated with a binarized annotation for the following characteristics:

1. If the variant is in a coding region, as defined by the RefSeq (NG_008376.4)[52]
2. If it is rare in the population. Defined as allele frequency among all populations in gnomAD < 0.05.

3. If it is deleterious. If it is a coding variant, we use an ADME optimized framework for predicting deleteriousness[29]. If it is non-coding we use a majority vote of CADD, DANN, and FATHMM[29,53–55]. LOFTEE predictions of deleteriousness supersede the other methods, if available[39]. LOFTEE is designed specifically to evaluate the functional consequence of stop gain, splice site, and frameshift variants based on sequence features such as position in transcript. Since LOFTEE does not evaluate all types of variants, we use the ADME optimized framework to predict the consequence of all remaining SNVs and INDELS.
4. If it is an INDEL of any length. INDELS are reduced to the first nucleotide and given the INDEL annotation, so as to keep the length of each sequence the same.
5. If it is in a methylation mark, as defined by UCSC Genome Browser tracks wgEncodeHaibMethyl450Gm12878SitesRep1 and wgEncodeHaibMethylRrbsGm12878HaibSitesRep1 [56,57].
6. If it is in a DNase hypersensitivity site. We use UCSC Genome Browser track wgEncodeAwgDnaseMasterSites.
7. If it is in a transcription factor binding site. We use UCSC Genome Browser tracks tfbsConsSites and wgEncodeRegTfbsClusteredV3.
8. If it is a known *CYP2D6* expression quantitative trait loci (eQTL) for any tissue in gTEX v6 [58].
9. If it codes for a residue in the *CYP2D6* active site where the substrate binds to the protein [59].

Variants are annotated using a custom pipeline that includes annovar, VEP, and a script for binarizing the annotations[60,61]. The final dimensions of the input haplotype matrix are 7417x12 (the length of the locus window x the nucleotide vector). Every base in the sequence window, coding and non-coding, is annotated and converted into a 1x12 vector (four possible nucleotides, eight annotations).

Training procedure

The number of existing star alleles with curated function was small for typical deep learning applications, thus transfer learning was used to reduce the amount of training data required to create a robust model. Hubble.2D6 was trained in a stepwise process with two pre-training steps, first with simulated data and second with real data (described in detail in the subsequent paragraph and S1 Fig). Each step in the training procedure used a CNN with identical number of convolutional layers and filter shapes, although the model output varied at each step. This allowed the weights of the convolutional layers learned at each step to be transferred to the next stage in training, iteratively updating the weights with different datasets.

First, a CNN classifier was trained to predict the activity score of 50,000 simulated pairs of haplotypes (or diplotypes), thereby creating a neural network representation of the activity score (the simulation procedure is described in the following section). We simulate diplotypes rather than haplotypes in order to be able to further train the model using functional activity data collected from liver microsomes. In the second stage of training, a regression model was trained to predict measured metabolic activity of 314 liver microsome samples using dextromethorphan as substrate[36]. The weights derived from the convolutional layers of the first model were transferred to the second model and the fully connected portion of the model was retrained using randomly initialized weights. The input to this second model was the pair of

haplotypes identified through sequencing of the *CYP2D6* loci encoded in the same manner described previously (these data are described in the following section). The final model was created by removing the final output layer of the network trained on the liver microsome data and adding a new output layer with two neurons with sigmoid activations, corresponding to the two outputs described previously (one representing no function status, another representing normal function status). The new network was created by creating a new neural network with an identical architecture, except for the last layer. Then, the weights from the pretrained network are copied to the new network for each layer, except for the final output layer. The final layer with the two output neurons will then be trained from randomly initialized weights in a final training phase using the star allele haplotype sequences. Since the liver microsome model was trained on pairs of haplotypes, the final model input consisted of two identical copies of the input haplotype matrix for star allele classification. The network was then retrained to classify star allele functional status with the starting weights initialized from the liver microsome model. Ten models were trained, and their outputs averaged to form an ensemble.

Training data

Simulations. Simulated diplotypes used in pre-training were created by randomly selecting a pair of *CYP2D6* star alleles with curated function (normal, decreased, or no function haplotypes) that do not have any structural variants and constructing haplotypes with the variants associated with the star alleles. Star allele definitions were downloaded from PharmVar (v4.1.1). To introduce additional diversity to the training data, alternate alleles (both SNVs and INDELS) were sampled for variant sites not associated with any star allele following a uniform distribution with the probability of an alternate allele occurring equal to the population level alternate allele frequency published in gnomAD[24]. It is possible that rare, deleterious variants not currently represented in any star allele were added during this process. However, noisy pretraining data has been shown not to negatively impact the final model[62]. A total of 10,000 genotypes were selected for each AS (0, 0.5, 1, 1.5, 2), for a total of 50,000 simulated samples used in training, and an additional 10,000 total genotypes to use as a test set.

Liver microsome data. Activity data were available for 314 liver microsome samples from a prior study as a second pre-training step for Hubble.2D6[36]. The liver microsome data was collected from two sites, 249 samples from St. Jude's Children's Research Hospital (SJCRH), and 65 samples from University of Washington (UW). All samples were sequenced with the PGRNseq panel[51]. Metabolic activity was measured using two substrates, dextromethorphan and metoprolol, but only dextromethorphan activity from SJCRH was used for pre-training, because dextromethorphan is known CYP2D6 'probe drug' that is primarily metabolized by CYP2D6, whereas other enzymes can contribute to e.g. metoprolol metabolism[63]. The UW samples were used for validation of this intermediate model. Star alleles and structural variants were called for each sample using Stargazer, which uses statistical phasing performed by Beagle to predict haplotypes[31,64].

Star allele data. The sequences used to train the final model were constructed based on star allele definitions by PharmVar (version 4.1.1.1, downloaded 10/25/2019). We selected 31 star alleles (as per the core allele definitions) and each of their suballeles for training. Suballeles are versions of star alleles that have been identified that carry additional variants (typically non-coding variants) other than the core variants that define the star allele (e.g. *CYP2D6*1.002*). The model was evaluated with 24 randomly selected star alleles. Star alleles with no curated function (defined by CPIC and posted by PharmVar as "Uncertain", "Unknown", or "not available") were excluded from the training procedure. *CYP2D6* star allele functions are regularly updated, so official functional designations may have changed since

this work was completed. During the training of each model 10% of the samples from each functional class (no function, decreased function, and normal function) were randomly held out to be used to check for overfitting in the training process.

Baseline model

A baseline model was developed to compare performance of Hubble.2D6 against a more conventional model. No existing model takes in any arbitrary set of variants from the *CYP2D6* locus so we trained an ordinal logistic regression model to predict star allele function from the sums of the annotation vectors in the functional variant representation. Specifically, we input a single 9x1 vector representing each of the nine annotations where each value is the count of the variants in the star allele sequence that were found to have the corresponding annotation. This model then had access to the same annotations used to train Hubble.2D6, but lost the local context of each variant that is learned by the CNN. The baseline model was trained and evaluated on the same set of star alleles described in the previous section (star allele data). We trained an ordinal logistic regression model using the `polr` function in the R package MASS [65]. The baseline model was evaluated by calculating the accuracy on the training and validation sets and comparing values to those determined by Hubble.2D6.

Validation

The model was evaluated by predicting the function of the 25 star alleles with curated function that were excluded from the training process. The area under the receiver operator characteristic curve (AUROC) was calculated for the training and test groups for the two scores output by the model. In addition, the function of 71 star alleles with uncurated function was predicted.

In order to further validate our model, we used *in vitro* data from a study [37] that measured the activity of 46 star alleles (including wild-type) using three substrates. Of these 46 star alleles, 30 have curated function while 16 do not have a CPIC assigned function. The metabolic activity used for each star allele was the percent activity of the reference, taking the mean activity across all three substrates. We calculated the variance explained by the predicted function using a linear model and assessed the heterogeneity in the measured activity of each functional group (no function, decreased function, and normal function) using a one-way ANOVA. This was done separately for samples with curated function and those with uncurated function. *CYP2D6*53* was excluded from this analysis because it had a 28-fold increase in *in vitro* function compared to reference, and deemed to be an extreme outlier. Additionally, *CYP2D6*61* and *CYP2D6*63* were excluded because they are *CYP2D6-CYP2D7* hybrids, a result of structural variation which is not included in the Hubble.2D6 framework.

Model interpretation

We interpreted the model and the relative importance of the variants in each star allele for the predicted function using DeepLIFT from the DeepExplain package [66,67]. DeepLIFT compares the activations of a neural network for a given sample against a reference sample, and outputs importance scores for each input feature. We ran DeepLIFT on each star allele sequence with a *CYP2D6*1* reference sequence. This yielded importance scores for each variant in each star allele that were different from the variants in *CYP2D6*1*.

Model evaluation

We evaluated the added components of the model (transfer learning and the functional variant representation) by training new models with each component removed. Each newly trained

model was identical to the final Hubble.2D6 model in every way (e.g. the input data, predicting star allele function, etc.) except for the component or annotation under evaluation. Specifically, to evaluate the contribution of transfer learning to the final model, we trained a new model identical to the final model except the weights were randomly initialized rather than transferred from a pretrained network. To evaluate the contribution of the functional variant representation we trained a new model and input only the one-hot encoding of the nucleotide sequence, no annotations included. For this model weights were transferred from a pretrained network. Finally, we trained a third model that did not include any variant annotations and had randomly initialized weights. For each test case, we followed the same procedure to generate predictions as with the full model: an ensemble of seven models was trained and the average score for each of the two output scores taken across all seven models which was then converted into a single functional class. Classification accuracy was calculated for star alleles in training and test.

We also evaluated each of the annotations included in the functional variant representation. This was done by again training ensemble models with transfer learning, but in this case a single annotation was included in addition to the one-hot encoding of the nucleotide sequence. This was done for each of the eight annotations included in the final model. Again, the mean scores from each ensemble were taken and converted into predictions as previously described and then calculated classification accuracy.

Supporting information

S1 Fig. Model architecture and training procedure. Here we show the architecture of each model that contributes to the final model as well as the logic of the learned weight transfers between models. We train three models total, including two pre-training steps and a final model that predicts star allele function (Hubble.2D6). The first model predicts the *CYP2D6* Activity Score for simulated diplotype sequence data (a classification problem). The second predicts the measured metabolic activity for liver microsome data (regression). The final model predicts categorical function of a star allele (classification). Layer colors represent which model the learned weights were initialized in. The numbers in parentheses (e.g. 19x13) represent the dimensions of the filter for that layer. Layers with colors indicate layers with learnable weights, colors indicate which training step the weights were learned in. Asterisks indicate layer weights were transferred then fine-tuned in the corresponding model. (PNG)

S2 Fig. Importance scores for core variants in each star allele used for training and test of Hubble.2D6, as well as uncurated star alleles. Star alleles are along the y-axis and core variants (both amino acid changes and non-coding changes) are listed along the x-axis. Each dot represents the importance of the core variant to the final prediction as determined by DeepLIFT. The size of the dot represents the value of the importance score, with larger dots indicating variants with larger importance scores, typically associated with a negative impact on function. Star alleles are annotated with the curated function as well as the Hubble.2D6 predicted function. Star alleles are divided along the y-axis between star alleles that were included in the training data (top) and those used as test samples (bottom). Star alleles are sorted by the sum of the importance scores, with those with the largest sums at the bottom. Core variants are divided along the x-axis by those that are uniquely in either the training or test samples (right), and those that are shared between star alleles in train and test (left). Core variants are sorted by their mean importance score across all star alleles. Core variants are annotated with the deleteriousness annotation used in the functional variant representation. (PNG)

S1 Table. CYP2D6 Star allele function predictions. “Curated Function” refers to the function determined by human curators listed on PharmVar (<https://www.pharmvar.org/gene/CYP2D6>). Star alleles with an “uncurated” function refers to star alleles with either unknown or uncertain function. Curated functions were retrieved from PharmVar version 4.11.1, downloaded 10/25/2019. “Hubble.2D6 Predicted Function” is the star allele function predicted by Hubble.2D6.
(XLSX)

Acknowledgments

We thank Kenneth Thummel at the University of Washington and Erin Schuetz at St. Jude Children’s Research Hospital for providing the liver bank resources to conduct this work.

Author Contributions

Conceptualization: Gregory McInnes.

Data curation: Gregory McInnes, Rachel Dalton.

Formal analysis: Gregory McInnes.

Funding acquisition: Russ B. Altman, Erica L. Woodahl.

Investigation: Gregory McInnes.

Methodology: Gregory McInnes, Katrin Sangkuhl, Michelle Whirl-Carrillo, Andrea Gaedigk.

Project administration: Russ B. Altman, Erica L. Woodahl.

Resources: Philip S. Tsao, Russ B. Altman, Erica L. Woodahl.

Software: Seung-been Lee.

Supervision: Russ B. Altman, Erica L. Woodahl.

Validation: Gregory McInnes.

Visualization: Gregory McInnes.

Writing – original draft: Gregory McInnes.

Writing – review & editing: Gregory McInnes, Rachel Dalton, Katrin Sangkuhl, Michelle Whirl-Carrillo, Seung-been Lee, Philip S. Tsao, Andrea Gaedigk, Russ B. Altman, Erica L. Woodahl.

References

1. Zhou S-F. Polymorphism of human cytochrome P450 2D6 and its clinical significance: Part I. Clin Pharmacokinet. 2009; 48:689–723. <https://doi.org/10.2165/11318030-000000000-00000> PMID: 19817501
2. Zhou S-F. Polymorphism of human cytochrome P450 2D6 and its clinical significance: part II. Clin Pharmacokinet. 2009; 48:761–804. <https://doi.org/10.2165/11318070-000000000-00000> PMID: 19902987
3. Saravanakumar A, Sadighi A, Ryu R, Akhlaghi F. Physicochemical Properties, Biotransformation, and Transport Pathways of Established and Newly Approved Medications: A Systematic Review of the Top 200 Most Prescribed Drugs vs. the FDA-Approved Drugs Between 2005 and 2016. Clin Pharmacokinet. 2019. <https://doi.org/10.1007/s40262-019-00750-8> PMID: 30972694
4. Gaedigk A. Complexities of CYP2D6 gene analysis and interpretation. Int Rev Psychiatry. 2013; 25:534–553. <https://doi.org/10.3109/09540261.2013.825581> PMID: 24151800
5. Nofziger C, Turner AJ, Sangkuhl K, Whirl-Carrillo M, Agúndez JAG, Black JL, et al. PharmVar GeneFocus: CYP2D6. Clin Pharmacol Ther. 2019. <https://doi.org/10.1002/cpt.1643> PMID: 31544239

6. Gaedigk A, Ingelman-Sundberg M, Miller NA, Leeder JS, Whirl-Carrillo M, Klein TE, et al. The Pharmacogene Variation (PharmVar) Consortium: Incorporation of the Human Cytochrome P450 (CYP) Allele Nomenclature Database. *Clin Pharmacol Ther.* 2018; 103:399–401. <https://doi.org/10.1002/cpt.910> PMID: 29134625
7. Gaedigk A, Sangkuhl K, Whirl-Carrillo M, Twist GP, Klein TE, Miller NA, et al. The Evolution of PharmVar. *Clin Pharmacol Ther.* 2019; 105:29–32. <https://doi.org/10.1002/cpt.1275> PMID: 30536702
8. Gaedigk A, Sangkuhl K, Whirl-Carrillo M, Klein T, Leeder JS. Prediction of CYP2D6 phenotype from genotype across world populations. *Genet Med.* 2017; 19:69–76. <https://doi.org/10.1038/gim.2016.80> PMID: 27388693
9. Caudle KE, Sangkuhl K, Whirl-Carrillo M, Swen JJ, Haidar CE, Klein TE, et al. Standardizing CYP2D6 Genotype to Phenotype Translation: Consensus Recommendations from the Clinical Pharmacogenetics Implementation Consortium and Dutch Pharmacogenetics Working Group. *Clin Transl Sci.* 2019. <https://doi.org/10.1111/cts.12692> PMID: 31647186
10. Crews KR, Gaedigk A, Dunnenberger HM, Leeder JS, Klein TE, Caudle KE, et al. Clinical Pharmacogenetics Implementation Consortium guidelines for cytochrome P450 2D6 genotype and codeine therapy: 2014 update. *Clin Pharmacol Ther.* 2014; 95:376–382. <https://doi.org/10.1038/clpt.2013.254> PMID: 24458010
11. Hicks JK, Sangkuhl K, Swen JJ, Ellingrod VL, Müller DJ, Shimoda K, et al. Clinical pharmacogenetics implementation consortium guideline (CPIC) for CYP2D6 and CYP2C19 genotypes and dosing of tricyclic antidepressants: 2016 update. *Clin Pharmacol Ther.* 2017; 102:37–44. <https://doi.org/10.1002/cpt.597> PMID: 27997040
12. Hicks JK, Bishop JR, Sangkuhl K, Müller DJ, Ji Y, Leckband SG, et al. Clinical Pharmacogenetics Implementation Consortium (CPIC) Guideline for CYP2D6 and CYP2C19 Genotypes and Dosing of Selective Serotonin Reuptake Inhibitors. *Clin Pharmacol Ther.* 2015; 98:127–134. <https://doi.org/10.1002/cpt.147> PMID: 25974703
13. Brown JT, Bishop JR, Sangkuhl K, Nurmi EL, Mueller DJ, Dinh JC, et al. Clinical Pharmacogenetics Implementation Consortium Guideline for Cytochrome P450 (CYP)2D6 Genotype and Atomoxetine Therapy. *Clin Pharmacol Ther.* 2019; 106:94–102. <https://doi.org/10.1002/cpt.1409> PMID: 30801677
14. Goetz MP, Sangkuhl K, Guchelaar H-J, Schwab M, Province M, Whirl-Carrillo M, et al. Clinical Pharmacogenetics Implementation Consortium (CPIC) Guideline for CYP2D6 and Tamoxifen Therapy. *Clin Pharmacol Ther.* 2018; 103:770–777. <https://doi.org/10.1002/cpt.1007> PMID: 29385237
15. Swen JJ, Nijenhuis M, de Boer A, Grandia L, Maitland-van der Zee AH, Mulder H, et al. Pharmacogenetics: from bench to byte—an update of guidelines. *Clin Pharmacol Ther.* 2011; 89:662–673. <https://doi.org/10.1038/clpt.2011.34> PMID: 21412232
16. pharmacogenetic-recommendations-may-2020.pdf. Available from: <https://www.knmp.nl/downloads/pharmacogenetic-recommendations-may-2020.pdf>.
17. Schroth W, Goetz MP, Hamann U, Fasching PA, Schmidt M, Winter S, et al. Association between CYP2D6 polymorphisms and outcomes among women with early stage breast cancer treated with tamoxifen. *JAMA.* 2009; 302:1429–1436. <https://doi.org/10.1001/jama.2009.1420> PMID: 19809024
18. J Marcalus S, Bristow-Marcalus S. Combating opioid addiction and abuse—2 ways to effectively intervene in the cycle of addiction through pharmacogenomics. *J Am Pharm Assoc.* 2019. <https://doi.org/10.1016/j.japh.2019.04.016> PMID: 31126828
19. UnitedHealthcare Pharmacogenetic Testing. 2019. Available from: <https://www.uhcprovider.com/content/dam/provider/docs/public/policies/comm-medical-drug/pharmacogenetic-testing.pdf>.
20. FUTURE Local Coverage Determination for MoIDX: Pharmacogenomics Testing (L38294). [cited 29 Jun 2020]. Available from: <https://www.cms.gov/medicare-coverage-database/details/lcd-details.aspx?LCDId=38294&ver=16&DocID=L38294&SearchType=Advanced&bc=EAAAAAgAAAA&>.
21. Gaedigk A, Simon SD, Pearce RE, Bradford LD, Kennedy MJ, Leeder JS. The CYP2D6 activity score: translating genotype information into a qualitative measure of phenotype. *Clin Pharmacol Ther.* 2008; 83:234–242. <https://doi.org/10.1038/sj.cpt.6100406> PMID: 17971818
22. McInnes G, Lavertu A, Sangkuhl K, Klein TE, Whirl-Carrillo M, Altman RB. Pharmacogenetics at scale: An analysis of the UK Biobank. *bioRxiv.* 2020. p. 2020.05.30.125583. <https://doi.org/10.1101/2020.05.30.125583>
23. Whirl-Carrillo M, McDonagh EM, Hebert JM, Gong L, Sangkuhl K, Thorn CF, et al. Pharmacogenomics knowledge for personalized medicine. *Clin Pharmacol Ther.* 2012; 92:414–417. <https://doi.org/10.1038/clpt.2012.96> PMID: 22992668
24. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alföldi J, Wang Q, et al. Variation across 141,456 human exomes and genomes reveals the spectrum of loss-of-function intolerance across human protein-coding genes. *bioRxiv.* 2019. p. 531210. <https://doi.org/10.1101/531210>

25. Kozyra M, Ingelman-Sundberg M, Lauschke VM. Rare genetic variants in cellular transporters, metabolic enzymes, and nuclear receptors can be important determinants of interindividual differences in drug response. *Genet Med*. 2017; 19:20–29. <https://doi.org/10.1038/gim.2016.33> PMID: 27101133
26. Ingelman-Sundberg M, Mkrтчian S, Zhou Y, Lauschke VM. Integrating rare genetic variants into pharmacogenetic drug response predictions. *Hum Genomics*. 2018; 12:26. <https://doi.org/10.1186/s40246-018-0157-3> PMID: 29793534
27. CPIC SOP for Assigning Allele Function. [cited 23 Jan 2020]. Available from: <https://cpicpgx.org/resources/cpic-draft-allele-function-sop/>.
28. Zhou Y, Fujikura K, Mkrтчian S, Lauschke VM. Computational Methods for the Pharmacogenetic Interpretation of Next Generation Sequencing Data. *Front Pharmacol*. 2018; 9:1437. <https://doi.org/10.3389/fphar.2018.01437> PMID: 30564131
29. Zhou Y, Mkrтчian S, Kumondai M, Hiratsuka M, Lauschke VM. An optimized prediction framework to assess the functional impact of pharmacogenetic variants. *Pharmacogenomics J*. 2019; 19:115–126. <https://doi.org/10.1038/s41397-018-0044-2> PMID: 30206299
30. Papadimitriou S, Gazzo A, Versbraegen N, Nachtegaele C, Aerts J, Moreau Y, et al. Predicting disease-causing variant combinations. *Proc Natl Acad Sci U S A*. 2019; 116:11878–11887. <https://doi.org/10.1073/pnas.1815601116> PMID: 31127050
31. Lee S-B, Wheeler MM, Patterson K, McGee S, Dalton R, Woodahl EL, et al. Stargazer: a software tool for calling star alleles from next-generation sequencing data using CYP2D6 as a model. *Genet Med*. 2019; 21:361–372. <https://doi.org/10.1038/s41436-018-0054-0> PMID: 29875422
32. Twist GP, Gaedigk A, Miller NA, Farrow EG, Willig LK, Dinwiddie DL, et al. Constellation: a tool for rapid, automated phenotype assignment of a highly polymorphic pharmacogene, CYP2D6, from whole-genome sequences. *NPJ Genom Med*. 2016; 1:15007. <https://doi.org/10.1038/nnpjgenmed.2015.7> PMID: 29263805
33. van der Lee M, Allard WG, Rolf H A, Baak-Pablo RF, Menafrá R, Birgit A L, et al. A unifying model to predict variable drug response for personalised medicine. *bioRxiv*. 2020. p. 2020.03.02.967554. <https://doi.org/10.1101/2020.03.02.967554>
34. Zou J, Huss M, Abid A, Mohammadi P, Torkamani A, Telenti A. A primer on deep learning in genomics. *Nat Genet*. 2019; 51:12–18. <https://doi.org/10.1038/s41588-018-0295-5> PMID: 30478442
35. Zhao Z-Q, Zheng P, Xu S-T, Wu X. Object Detection With Deep Learning: A Review. *IEEE Trans Neural Netw Learn Syst*. 2019; 30:3212–3232. <https://doi.org/10.1109/TNNLS.2018.2876865> PMID: 30703038
36. Dalton R, Lee S-B, Claw KG, Prasad B, Phillips BR, Shen DD, et al. Interrogation of CYP2D6 structural variant alleles improves the correlation between CYP2D6 genotype and CYP2D6-mediated metabolic activity. *Clin Transl Sci*. 2019. <https://doi.org/10.1111/cts.12695> PMID: 31536170
37. Muroi Y, Saito T, Takahashi M, Sakuyama K, Niinuma Y, Ito M, et al. Functional characterization of wild-type and 49 CYP2D6 allelic variants for N-desmethyltamoxifen 4-hydroxylation activity. *Drug Metab Pharmacokinet*. 2014; 29:360–366. <https://doi.org/10.2133/dmpk.dmpk-14-rg-014> PMID: 24647041
38. Chiasson M, Dunham MJ, Rettie AE, Fowler DM. Applying Multiplex Assays to Understand Variation in Pharmacogenes. *Clin Pharmacol Ther*. 2019; 106:290–294. <https://doi.org/10.1002/cpt.1468> PMID: 31145826
39. Karczewski KJ. LOFTEE (Loss-Of-Function Transcript Effect Estimator). 2015.
40. Gaedigk A, Fuhr U, Johnson C, Bérard LA, Bradford D, Leeder JS. CYP2D7-2D6 hybrid tandems: identification of novel CYP2D6 duplication arrangements and implications for phenotype prediction. *Pharmacogenomics*. 2010; 11:43–53. <https://doi.org/10.2217/pgs.09.133> PMID: 20017671
41. Gaedigk A, Blum M, Gaedigk R, Eichelbaum M, Meyer UA. Deletion of the entire cytochrome P450 CYP2D6 gene as a cause of impaired drug metabolism in poor metabolizers of the debrisoquine/sparteine polymorphism. *Am J Hum Genet*. 1991; 48:943–950. PMID: 1673290
42. Gaedigk A, Twist GP, Leeder JS. CYP2D6, SUL1A1 and UGT2B17 copy number variation: quantitative detection by multiplex PCR. *Pharmacogenomics*. 2012; 13:91–111. <https://doi.org/10.2217/pgs.11.135> PMID: 22111604
43. Wang D, Poi MJ, Sun X, Gaedigk A, Leeder JS, Sadee W. Common CYP2D6 polymorphisms affecting alternative splicing and transcription: long-range haplotypes with two regulatory variants modulate CYP2D6 activity. *Hum Mol Genet*. 2014; 23:268–278. <https://doi.org/10.1093/hmg/ddt417> PMID: 23985325
44. Ray B, Ozcagli E, Sadee W, Wang D. CYP2D6 haplotypes with enhancer single-nucleotide polymorphism rs5758550 and rs16947 (*2 allele): implications for CYP2D6 genotyping panels. *Pharmacogenet Genomics*. 2019; 29:39–47. <https://doi.org/10.1097/FPC.0000000000000363> PMID: 30520769

45. Wang D, Papp AC, Sun X. Functional characterization of CYP2D6 enhancer polymorphisms. *Hum Mol Genet.* 2015; 24:1556–1562. <https://doi.org/10.1093/hmg/ddu566> PMID: 25381333
46. Sandee D, Morrissey K, Agrawal V, Tam HK, Kramer MA, Tracy TS, et al. Effects of genetic variants of human P450 oxidoreductase on catalysis by CYP2D6 in vitro. *Pharmacogenet Genomics.* 2010; 20:677–686. <https://doi.org/10.1097/FPC.0b013e32833f4f9b> PMID: 20940534
47. Crewe HK, Lennard MS, Tucker GT, Woods FR, Haddock RE. The effect of selective serotonin reuptake inhibitors on cytochrome P4502D6 (CYP2D6) activity in human liver microsomes. *Br J Clin Pharmacol.* 2004; 58:S744–S747. <https://doi.org/10.1111/j.1365-2125.2004.02282.x> PMID: 15595963
48. Kelley DR, Snoek J, Rinn JL. Basset: learning the regulatory code of the accessible genome with deep convolutional neural networks. *Genome Res.* 2016; 26:990–999. <https://doi.org/10.1101/gr.200535.115> PMID: 27197224
49. Thiele C, Hirschfeld G. cutpointr: Improved Estimation and Validation of Optimal Cutpoints in R. *arXiv [stat.CO]*. 2020. Available from: <http://arxiv.org/abs/2002.09209>.
50. Cheng Jianlin, Wang Zheng, Pollastri G. A neural network approach to ordinal regression. 2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence). 2008. pp. 1279–1284.
51. Gordon AS, Fulton RS, Qin X, Mardis ER, Nickerson DA, Scherer S. PGRNseq: a targeted capture sequencing panel for pharmacogenetic research and implementation. *Pharmacogenet Genomics.* 2016; 26:161–168. <https://doi.org/10.1097/FPC.0000000000000202> PMID: 26736087
52. O'Leary NA, Wright MW, Brister JR, Ciufu S, Haddad D, McVeigh R, et al. Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res.* 2016; 44:D733–45. <https://doi.org/10.1093/nar/gkv1189> PMID: 26553804
53. Rentzsch P, Witten D, Cooper GM, Shendure J, Kircher M. CADD: predicting the deleteriousness of variants throughout the human genome. *Nucleic Acids Res.* 2019; 47:D886–D894. <https://doi.org/10.1093/nar/gky1016> PMID: 30371827
54. Quang D, Chen Y, Xie X. DANN: a deep learning approach for annotating the pathogenicity of genetic variants. *Bioinformatics.* 2015; 31:761–763. <https://doi.org/10.1093/bioinformatics/btu703> PMID: 25338716
55. Shihab HA, Rogers MF, Gough J, Mort M, Cooper DN, Day INM, et al. An integrative approach to predicting the functional effects of non-coding and coding sequence variation. *Bioinformatics.* 2015; 31:1536–1543. <https://doi.org/10.1093/bioinformatics/btv009> PMID: 25583119
56. Rosenbloom KR, Sloan CA, Malladi VS, Dreszer TR, Learned K, Kirkup VM, et al. ENCODE data in the UCSC Genome Browser: year 5 update. *Nucleic Acids Res.* 2013; 41:D56–63. <https://doi.org/10.1093/nar/gks1172> PMID: 23193274
57. Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, et al. The human genome browser at UCSC. *Genome Res.* 2002; 12:996–1006. <https://doi.org/10.1101/gr.229102> PMID: 12045153
58. GTEx Consortium, Laboratory, Data Analysis & Coordinating Center (LDACC)—Analysis Working Group, Statistical Methods groups—Analysis Working Group, Enhancing GTEx (eGTEx) groups, NIH Common Fund, NIH/NCI, et al. Genetic effects on gene expression across human tissues. *Nature.* 2017; 550:204–213. <https://doi.org/10.1038/nature24277> PMID: 29022597
59. Zhou S. *Cytochrome P450 2D6: Structure, Function, Regulation and Polymorphism.* CRC Press; 2016.
60. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data. *Nucleic Acids Res.* 2010; 38:e164. <https://doi.org/10.1093/nar/gkq603> PMID: 20601685
61. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GRS, Thormann A, et al. The Ensembl Variant Effect Predictor. *Genome Biol.* 2016; 17:122.
62. Mahajan D, Girshick R, Ramanathan V, He K, Paluri M, Li Y, et al. Exploring the Limits of Weakly Supervised Pretraining. *arXiv [cs.CV]*. 2018. Available from: <http://arxiv.org/abs/1805.00932>.
63. Berger B, Bachmann F, Duthaler U, Krähenbühl S, Haschke M. Cytochrome P450 Enzymes Involved in Metoprolol Metabolism and Use of Metoprolol as a CYP2D6 Phenotyping Probe Drug. *Front Pharmacol.* 2018; 9:774.
64. Browning SR, Browning BL. Rapid and accurate haplotype phasing and missing-data inference for whole-genome association studies by use of localized haplotype clustering. *Am J Hum Genet.* 2007; 81:1084–1097. <https://doi.org/10.1086/521987> PMID: 17924348
65. Ripley B, Venables B, Bates DM, Hornik K, Gebhardt A, Firth D, et al. Package “mass.” *Cran R.* 2013; 538. Available from: <ftp://192.218.129.11/pub/CRAN/web/packages/MASS/MASS.pdf>.
66. Shrikumar A, Greenside P, Kundaje A. Learning Important Features Through Propagating Activation Differences. *arXiv [cs.CV]*. 2017. Available from: <http://arxiv.org/abs/1704.02685>

67. Ancona M, Ceolini E, Öztireli C, Gross M. Towards better understanding of gradient-based attribution methods for Deep Neural Networks. arXiv [cs.LG]. 2017. Available from: <http://arxiv.org/abs/1711.06104>.