

RESEARCH

Open Access



Genetic variation in human drug-related genes

Charlotta Pauline Irmgard Schärfe^{1,2,3}, Roman Tremmel⁴, Matthias Schwab^{4,5,6,7}, Oliver Kohlbacher^{2,3,8,9,10*} and Debora Susan Marks^{1*}

Abstract

Background: Variability in drug efficacy and adverse effects are observed in clinical practice. While the extent of genetic variability in classic pharmacokinetic genes is rather well understood, the role of genetic variation in drug targets is typically less studied.

Methods: Based on 60,706 human exomes from the ExAC dataset, we performed an in-depth computational analysis of the prevalence of functional variants in 806 drug-related genes, including 628 known drug targets. We further computed the likelihood of 1236 FDA-approved drugs to be affected by functional variants in their targets in the whole ExAC population as well as different geographic sub-populations.

Results: We find that most genetic variants in drug-related genes are very rare ($f < 0.1\%$) and thus will likely not be observed in clinical trials. Furthermore, we show that patient risk varies for many drugs and with respect to geographic ancestry. A focused analysis of oncological drug targets indicates that the probability of a patient carrying germline variants in oncological drug targets is, at 44%, high enough to suggest that not only somatic alterations but also germline variants carried over into the tumor genome could affect the response to antineoplastic agents.

Conclusions: This study indicates that even though many variants are very rare and thus likely not observed in clinical trials, four in five patients are likely to carry a variant with possibly functional effects in a target for commonly prescribed drugs. Such variants could potentially alter drug efficacy.

Keywords: Bioinformatics analysis, Exome sequence analysis, Pharmacogenomics

Background

About three in five Americans aged 20 years and above take prescription drugs every month [1] and many either encounter adverse drug reactions or reduced treatment efficacy [2]. The strong genetic component of altered drug response in patients is well known [3] and attributed to variants affecting drug pharmacokinetics (PK) and pharmacodynamics (PD) [4]. Methods to identify these genetic determinants have been developed in population-stratified [5–7] or individualized settings [4, 8]. Particularly, the vast amount of genetic information now available has opened up the possibility to systematically study

inter-individual differences in drug response using genome-wide association (GWA) studies [9, 10]. Results of these efforts have so far led to the pharmacogenomics labeling of 170 drugs by the Food and Drug Administration (FDA) [11] and the establishment of pharmacogenomics screening in many large hospitals in the US [12] and Europe [13].

However, typical pharmacogenomics GWA studies struggle with study sizes that are only large enough to detect common variants with an effect on the phenotype, but are unable to statistically pick up signals from rare variants with a functional effect [9, 10]. Thus, data from recent genetic population catalogs such as the 1000 Genomes project [14] and the NHLBI Exome Sequencing Project (ESP) have been used to determine the spectrum of variation in pharmacokinetics-related genes. Especially variants considered to be on the rare

* Correspondence: oliver.kohlbacher@uni-tuebingen.de; debbie@hms.harvard.edu

²Center for Bioinformatics, University of Tübingen, 72076 Tübingen, Germany

¹Department of Systems Biology, Harvard Medical School, Boston 02115, Massachusetts, USA

Full list of author information is available at the end of the article



end of the spectrum (minor allele frequency < 0.5%) were found to be abundant in genes associated with drug absorption, distribution, metabolism, and excretion (ADME) [15, 16] as well as in potential drug targets [17]. Based on these surveys, it was estimated that at least 97% of individuals carry actionable high-risk pharmacological variants affecting drug ADME in their genome [12, 18]. However, the role of genetic variation in pharmacologically established drug targets is less well studied.

The Exome Aggregation Consortium (ExAC) [19] has aggregated data from several large sequencing studies comprising exome sequencing data of 60,706 individuals—nearly an order of magnitude larger than the public population catalogs mentioned above. Using a cohort of this size, it now becomes possible to study even very rare variants in drug target and ADME genes and to calculate the overall risk of containing a functional variation for each patient. Furthermore, even though geographic ancestry is a known confounding factor for drug response and has been incorporated in clinical decision-making in the absence of individual genotype data [20], a comprehensive inventory of functional genetic variation in drug-associated genes across populations is still lacking. A cohort of the size of the ExAC catalog now allows determining the allele frequency (AF) of very rare variants in distinct population subgroups and comparing their prevalence.

In this study, we provide a comprehensive analysis of genetic variation predicted to result in altered protein function (“functional variants”) in 806 drug-related genes including 628 drug targets (163 targeted by cancer therapeutics). We further describe how this may affect the likelihood of 1236 FDA-approved drugs being affected by functional variants in their targets and how this likelihood varies between different populations. Even though variants in non-coding regions, copy number alterations, and chromosomal structural changes as well as epigenetics may further contribute to drug PK and PD variability [21], such alterations were not part of this study.

Methods

Data selection and handling

Known pharmacogenomics associations between drugs and genetic variants were retrieved from PharmGKB [22]. Data about drugs and drug-related genes were collated from DrugBank 5 [23]. Information about drug approval status, ATC code, and details about the drug–gene relationship (target, pharmacological action, and action type) were extracted from the xml file using python. We further obtained a list of the top 100 most prescribed drugs of 2013 from drugs.com [24] and the list of WHO essential medicines by parsing the Index of the 19th WHO Model List of Essential Medicines [25]. Drugs obtained from the top 100 list and WHO essential

medicines catalog were mapped to DrugBank compounds and those where this was not possible were excluded. Relations between hyaluronic acid and human gene targets as well as between dihydropyridines and skeletal *CACNA1S* were removed because the literature in the database entry did not support the pharmacological involvement of these pairs. We further removed ethanol from the list of WHO essential medicines because it is listed as a surface disinfectant and thus not dependent on the patient’s cellular targets.

Drug target genes were extracted from the drug–gene relationships in DrugBank by filtering this set for only those relations with an established pharmacological action flag and in which the gene is annotated as the drug target. Based on previous studies, a list of pharmacologically relevant cellular receptors, metabolic enzymes, and nuclear receptors was obtained from the supplemental material of recent pharmacogenomics surveys [15, 26], which comprise the set of ADME genes.

Genetic variant information including variant types, allele frequencies, and deleterious prediction scores were extracted from the ExAC VCF file (release 0.3) downloaded from the ExAC FTP server [19]. Multi-allelic variants were split using *vcflib breakmulti* (<https://github.com/vcflib/vcflib>) and synonymous variants were excluded. We then calculated for each variant the allele frequency (AF) in the full cohort ($n = 60,706$) as well as in each ExAC population separately by dividing the allele count (AC) by the allele number (AN). The following information about ancestry was used: AFR = African/African American ($n = 5203$), SAS = South Asian ($n = 8256$), EAS = East Asian ($n = 4327$), FIN = Finnish ($n = 3307$), NFE = Non-Finnish European ($n = 33,370$), AMR = admixed American/Latino ($n = 5789$), excluding OTH = other ($n = 545$) from the study. We further excluded variants whose loci were not observed at least once in every geographic population and in 50% of all possible samples (i.e., minimal allele number of 60,706). We removed duplicated variants using a unique identifier based on chromosome position, reference, and mutant allele.

Identifier mapping, filtering, and annotation were performed using the Konstanz Information Miner (KNIME) workflow system [27] and the Python programming language (Python Software Foundation, <https://www.python.org/>).

Variant subsets

To evaluate variants with functional effects in the ExAC catalog, we created subsets of variants with functional effects (“functional variants”): 1) loss-of-function variants affecting stop codons, splice sites, and shifts in the reading frame as annotated by the Loss-Of-Function Transcript Effect Estimator (LOFTEE) tool [28] in the ExAC VCF file; and 2) variants predicted to have a

damaging effect on the protein as predicted unanimously by PolyPhen-2 [29] (“possibly damaging” or “probably damaging”) and SIFT [30] (“deleterious”) as annotated in the ExAC VCF file. Functional variants with AFs above 0.5 were excluded from this set after observing annotation or reference genome mapping problems.

Deleteriousness predictors exhibit varying degrees of incorrect predictions, both false positive and false negative, but nevertheless show agreement on most predictions [19, 31–33]. To decrease the likelihood of false positive predictions, we combined SIFT and PolyPhen-2 into a consensus predictor as described above. In order to estimate the extent of false positive predictions of this consensus predictor, we calculated the intersection of the set of damaging variants with predictions by two independent predictors (CADD [34] and EVmutation [33]). Prediction scores for genes covered in ExAC were obtained from the tools’ websites (<http://cadd.gs.washington.edu/download>, https://marks.hms.harvard.edu/evmutation/human_proteins.html). For CADD, the threshold of a scaled score >20 for classifying a variant as damaging was chosen based on the conservative recommendation from the corresponding publication [34]. For each gene we calculated the fraction of common (AF ≥ 0.1%) and rare (AF < 0.1%) alleles.

Computation of cumulative probabilities for drugs and their related genes

To quantify the risk of an individual person in the population carrying functional variants in a particular gene, we define the “cumulative allele probability” (CAP) statistic, which captures both the number of functional variants and their allele frequencies per gene. Formally, this score is the probability for an individual to carry at least one variant allele a of the observed alleles A in a gene g :

$$CAP(g) = 1 - \prod_{a \in A} (1 - AF(a))^2$$

Two types of CAP scores were calculated, one for all functional variants in a drug-related gene and one based only on loss-of-function (LoF) variants.

To estimate how much each drug can be affected by functional variants in its target genes, we further define the drug-specific “drug risk probability” (DRP) score by combining the CAP scores for all drug target genes. Formally, the DRP score is defined as:

$$DRP(D) = 1 - \prod_{g \in G} \prod_{a \in A_g} (1 - AF(a))^2$$

Here G is the set of all target genes for drug D , as documented in DrugBank, and A_g the set of all variant alleles observed in gene g .

Correlation analysis of the DRP scores with the number of targets was performed using linear regression with ordinary least squares fitting using the Python package statsmodels [35] to compute the coefficient of determination r^2 .

Statistical analysis of population differences

Population comparisons for CAP and DRP scores were performed using the absolute risk difference (RD) metric:

$$RD = |P(\text{event in group 2}) - P(\text{event in group 1})|$$

The RD for a drug was calculated by subtracting the score for the population with the smallest DRP score from the score for the population with the highest DRP. To identify for which drugs a population has above or below average risks, we further calculated all pairwise risk differences between populations from which we then computed the population-specific mean RDs.

Detailed variant analyses in case studies

Protein structures for the porcine *TUBB1* homologue (PDB IDs 1tub [36], 3j6g [37]), *ADRB2* (PDB ID 2rh1 [38]), *PTGS1* (PDB ID 3n8w [39]), and *NOS2* (PDB ID 4nos [40]), were obtained from the Protein Data Bank. Recently published homology models for *VKORC1* were downloaded from the supplement of the respective publications [41, 42]. Co-evolution analysis of residues was done using plmc-based EVcouplings [43] and based on jackhmmer [44] alignments created with the UniProt entries of the respective protein as queries against the Uniref100 database [45] (release 01/2017). Alignment columns with more than 70% gaps and sequences with more than 50% gaps were excluded from the model. Functional impact was predicted using EVmutation [33] and, in the case of *VKORC1*, compared to experimental warfarin binding data [42]. Protein structures were analyzed and rendered using the UCSF Chimera package from the Computer Graphics Laboratory, University of California, San Francisco [46].

Statistical analysis and code availability

Statistical analysis of the data set was performed in jupyter/IPython notebooks [47] using pandas [48] and other packages of the SciPy stack [49]. The code used to analyze the data set and produce the figures is available on github (https://github.com/debbiemarkslab/variants_pharmacogenes).

Results

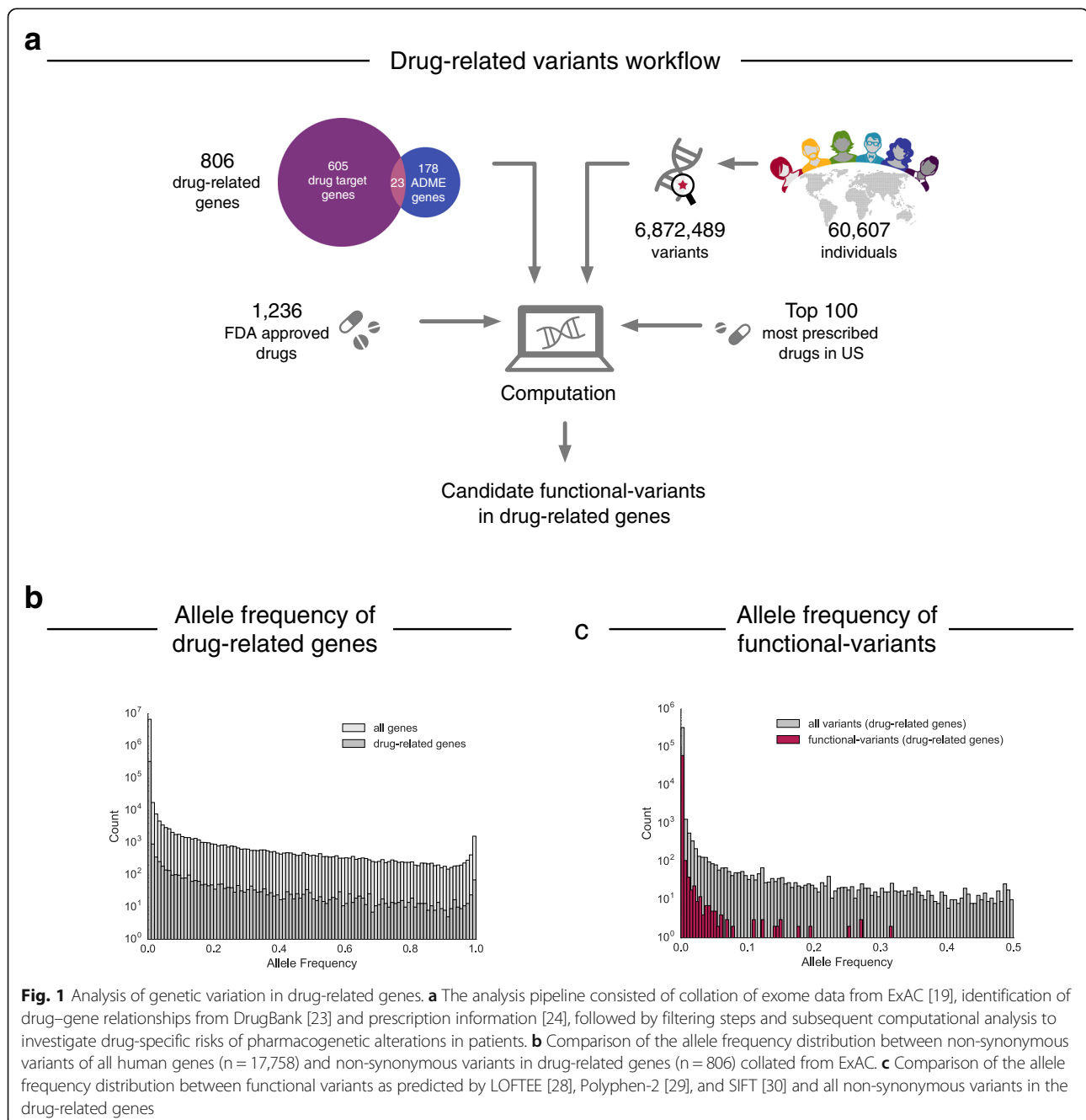
Drug-related genes show a high extent of genetic variability across 60 K individuals

To explore the extent of non-synonymous genetic variation in drug-related genes in the human populations, we analyzed single-nucleotide variants in 60,706 human

individual exomes from ExAC [19] in a set of 806 drug-related genes collated from DrugBank [23] and other sources [15, 26] (Fig. 1a; Supplementary Table found at doi: 10.6084/m9.figshare.5631751). The AF distribution of non-synonymous variants in drug-related genes is almost identical to that of all genes ($n = 17,758$) and 97.5% of observed non-synonymous variants have an allele frequency $< 0.1\%$ (sometimes termed a “rare variant” [19]) (Fig. 1b; Additional file 1: Figure S1). Of note, 71% of the variants in the human exome, including drug-related genes, have not been observed previously in

public repositories such as dbSNP and therefore can be considered novel (Additional file 1: Figure S1).

To identify variants that are most likely to affect the gene function (functional variants), we filtered the set of non-synonymous variants for those resulting in the complete loss of the protein’s primary biological function (LoF) by affecting splice sites or stop codons or resulting in frameshifts [19] or predicted to be “damaging” by PolyPhen-2 [29] and SIFT [30]. This resulted in 61,134 functional variants in 806 drug-related genes (of which 767 genes included at least one LoF variant). Variants



predicted to be damaging by SIFT and PolyPhen further agree with those predicted to have damaging effects with high sensitivity (90.4%) and specificity (68.8%) by the complementary prediction tools CADD [34] and EVmutation [33] (Additional file 1: Figure S2).

Not surprisingly, these functional variants tend to have lower AFs than all other non-synonymous variants (98.7% have an AF < 0.1%; Fig. 1c). Nevertheless, 43% of the drug-related genes with predicted functional variants have at least one functional variant with AF \geq 0.1%. The drug-related genes with the most frequent functional variants are membrane transporter genes related to drug >efflux and uptake such as *ABCB5* (three LoF, six damaging), *SLC22A1* (nine damaging), and *SLC22A14* (eight damaging). In the clinically highly important polymorphic cytochrome P450 enzyme *CYP2D6* eight damaging variants have been identified (Additional file 2: Table S1). Since the ExAC cohort contains an order of magnitude more individuals than previously available, it also allowed us to identify genes with many different functional variants even though each variant may be individually rare. The ADME genes with the most functional variants per residue reflect similar findings from smaller cohort studies and include the glutathione S-transferase sodium/bile transporter *SLC10A1* (0.36 variants/residue), *GSTA5* (0.31 variants/residue), and some cytochrome P450s such as *CYP1A1* (0.30 variants/residue) and *CYP2C19* (0.28 variants/residue) [15]. Furthermore, our analysis revealed drug target genes with comparable numbers of functional variants per residue, including the dofetilide target *KCNJ12* (0.31 variants/residue) and the target for the rheumatoid arthritis drug niflumic acid, *PLA2GLB* (0.30 variants/residue) (Additional file 2: Table S2).

While both metrics described above may be useful to evaluate the extent of genetic variation in the human population, they do not quantify the risk of an individual person in the population carrying functional variants in a particular gene. In order to estimate this risk, we define a statistic, the cumulative allele probability (CAP), which captures both the number of functional variants and their allele frequencies per gene (“Methods”; Additional file 2: Table S1). We want to emphasize that the CAP score of a gene does not necessarily reflect the extent to which the variants change the pharmacological behavior of the drug and therefore should be regarded as a score solely indicating a potential pharmacogenetic risk. Amongst the genes with the highest CAP scores—that is, the highest probability of being affected by a functional variant—are both ADME genes and drug targets. The ADME genes with the highest CAP scores include *NAT2* (81%, involved in metabolizing arylamine and hydrazine drugs), *CYP2D6* (59.6%, involved in the metabolism of 20% of most prescribed drugs in the US

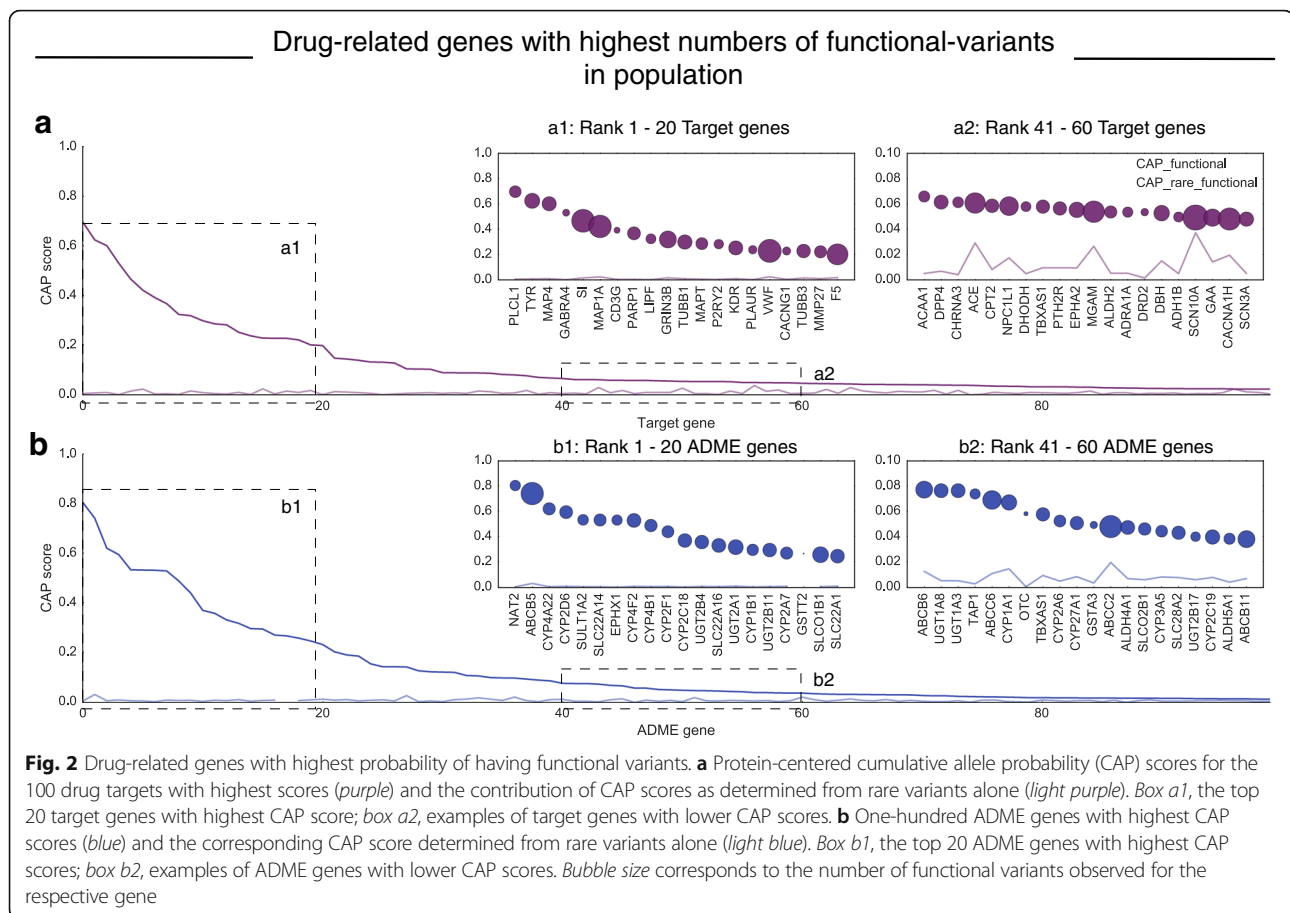
[50]), and the transporter gene *SLCO1B1* (26.0%, a high risk gene for simvastatin-related myopathy/rhabdomyolysis [51]). The drug target genes with comparable high CAP scores include tyrosinase (*TYR*; 62.4%, targeted by the acne drug azelaic acid), the alpha-4 subunit of the GABA_A receptor *GABRA4* (53%, targeted by benzodiazepines), and *F5* (20.1%, targeted by drotrecogin alpha, which was withdrawn from the market due to unacceptable high numbers of adverse drug reactions) (Fig. 2). We calculated an additional score, CAP_{LoF} based on LoF variants only. Again, genes with the highest CAP_{LoF} scores are ADME and target genes, including *CYP2F1* with CAP_{LoF} of 38.4%, *GSTT2* (26.9%), and *FCGR2A* (19.6%).

The major proportion of the CAP score for these highest “risk” genes derives from common genetic variants, many of which have been observed previously. Nevertheless, for many genes a non-negligible proportion of the score is contributed by rare functional variants, which were identified through the sufficiently large cohort size (see the lines in light purple and light blue in Fig. 2a and b, respectively, and Additional file 2: Table S1). In addition, we estimate that more than 60% of the drug-related genes in our set are putative novel candidates for pharmacogenomic research, so far missing relevant information from clinical studies (Additional file 1: Figure S3) [22].

Cancer drug target genes have many germline functional variants

Especially in cancer therapy, genetic variation in drug targets has been recognized to play a crucial role for treatment success [52, 53]. While some cancer drugs do not act in the tumor tissue, the cancer drug’s primary site of action usually is in the tumor, whose genome contains tumor-specific somatic variants as well as a subset of patient-specific germline variants [54]. Information on somatic variants from tumor samples is thus increasingly used to enable research on drug design and to implement stratified or personalized cancer therapy. However, the patient’s germline genome is routinely masked in these tumor sequencing analysis protocols [52, 53].

We thus wanted to assess whether target genes of drugs used in cancer therapy contain germline variants in the population that may affect the drug action and may be missed by current tumor sequencing analysis protocols. More than 15% of the drugs in this report (193 of the 1236) are used in oncology (as defined by the WHO ATC code [55]) and between them have 163 gene targets. Several of these targets have high probabilities of having a functional variant in the germline (Additional file 2: Table S1). For some of these targets the germline risk directly corresponds to potential altered treatment effects. This is the case for the kinase *KDR* (also known as *VEGFR2*; CAP = 25%), which is targeted by sorafenib and sunitinib to inhibit vascularization of the tumor site



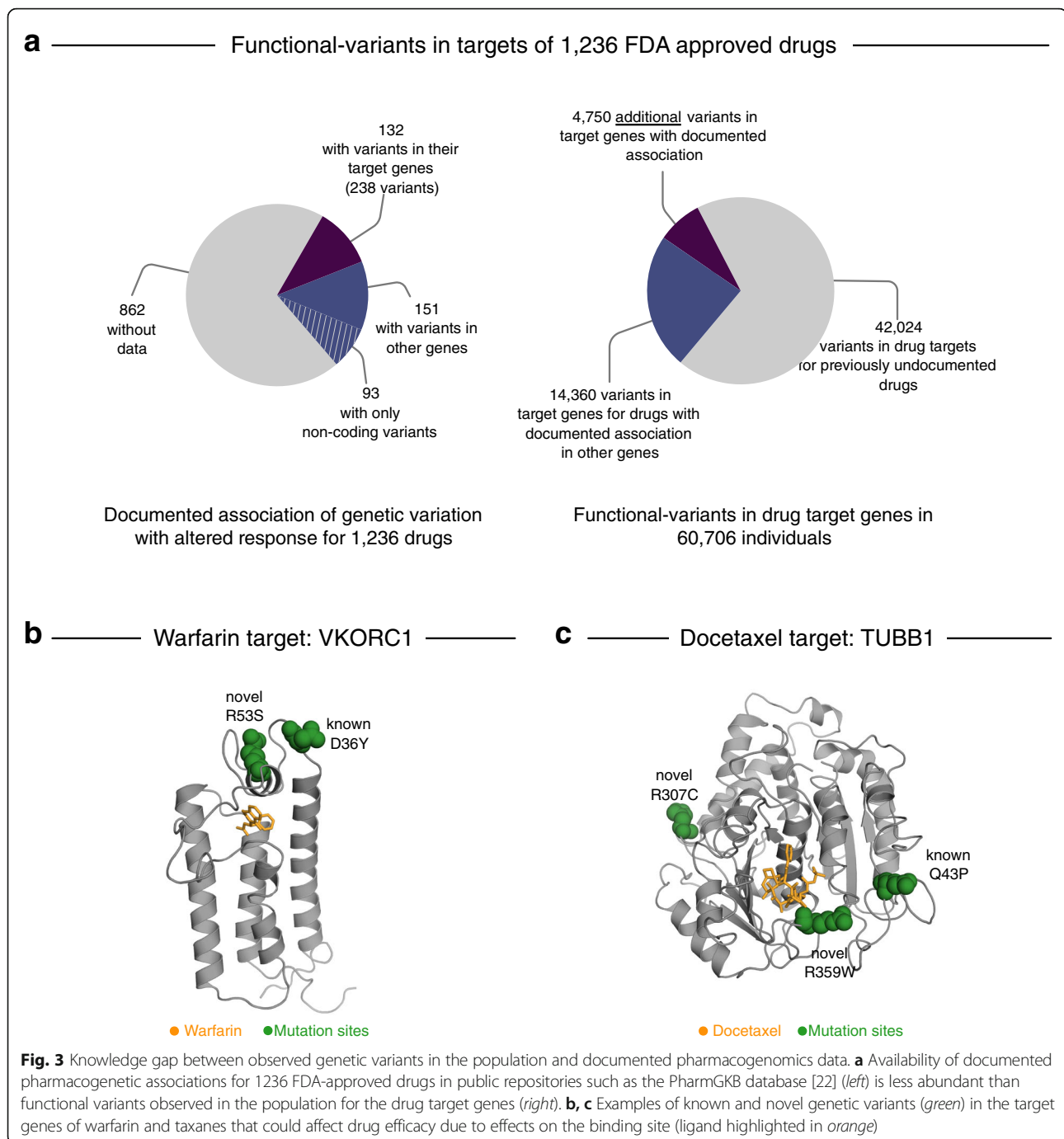
[56]. Other drug targets for cancer therapeutics with high CAP scores include *MAP4* (60%) and *TUBB1* (30%), which are targets of paclitaxel, *MAP1A* (42%), a target of estramustine, *CD3G* (39%), a target of munitomab, and *PARP1* (37%), a target of olaparib (Fig. 2). Overall, 40 cancer drug target genes, including 34 target genes with kinase domains, show CAP scores > 1%. For these examples, functional germline variants are only relevant for treatment response if the tumor genome also carries them. While there is not a complete overlap between both germline and tumor genome due to loss of heterozygosity and other alterations in carcinogenesis [54], our analysis suggests that a large percentage of the population may contain functional variants in cancer therapeutic targets in the germline that may carry over to the cancer genome and could be easily overlooked by current analysis protocols.

Aggregating risk for functional variants in targets by drug highlights drug candidates for future pharmacogenomics research

About 70% of the FDA-approved drugs analyzed here do not have any pharmacogenomics data associated with them in public repositories [22]. However, our analysis

shows that there are many functional variants in their target genes (Fig. 3a). To estimate how much each drug can be affected by functional variants in its target genes and to highlight possible candidates for future research, we computed the probability of containing a functional variant in any number of its reported targets in DrugBank [23] by combining the CAP scores of the drug's target genes to a "drug risk probability" (DRP; see "Methods" for details). For all FDA-approved drugs considered here ($n = 1236$), 43% have a DRP greater than 1% (Additional file 2: Table S3). The DRPs are weakly correlated to the number of targets (linear regression, $r^2 = 0.28$), leaving many drugs with few targets but higher than expected DRPs (determined by root mean square errors (RMSE), of the model; red circles in Additional file 1: Figure S4). For instance, one of the two human targets of azelaic acid, tyrosinase (*TYR*) is highly mutated in the population, causing a DRP of 62.5% for this drug, which results in an RMSE of 0.34.

Drugs with the top DRP scores are paclitaxel and docetaxel (82%), quinacrine (70%), azelaic acid (63%), triazolam, and other benzodiazepines (>50%) (Additional file 2: Table S3). This means that any individual in the population has a probability of more than 50% to carry a



functional variant that may affect the medication outcome of these drugs. Several of the drugs with high DRPs are considered “essential medicines” by the WHO [25]. In addition to paclitaxel and docetaxel, these include the opioid methadone (13.6%), the diuretic amiloride (11.7%), and the local anesthetic lidocaine (11.4%). For instance, the drug methadone targets the D- and M-type opioid receptors (*OPRD1*, *OPRM1*), and

whilst some non-coding variants and a single coding variant (rs1799971) have previously been associated with required dose adjustments and treatment response, we observe another 132 functional variants in these target genes, which could therefore be candidates for further testing. Since variants with predicted damaging effects dominate especially the rather high DRPs, we filtered the variants for only those resulting in LoF (DRP_{LoF}).

Restricting to these high-confidence variants, the DRP decreases below 10% and the drugs with the highest DRP_{LoF} include the anti-cancer drug marimastat ($DRP_{LoF} = 8.3\%$), the anti-ulcer medication sulfacrate ($DRP_{LoF} = 8.2\%$), the anti-flu drug oseltamivir ($DRP_{LoF} = 6.0\%$), which targets human *CESI* for activation, and several liptins used for diabetes that inhibit *DPP4* ($DRP_{LoF} = 5.6\%$) (Additional file 2: Table S3).

We then focused our analysis on the top 100 most prescribed medications in the US (from 2013 [24]), which resulted in a list of 77 unique drug compounds for further investigation. Of these drugs, 42% have a DRP for a functional variant of greater than 1% and the probability of an individual carrying a functional variant in any of the targets for these 77 top prescribed drugs is 81%. For some of these drugs it is already well established that there is some genetic component to drug response, even if the details are debated [57]. For instance, five of the top 15 most prescribed drugs in the US are asthma drugs (budesonide, salbutamol, salmeterol, fluticasone, and tiotropium). Whilst each of the DRPs is not particularly high (ranging from 0.06 to 0.25%), their widespread prescription rate (>100 million prescriptions in 2013) still results in thousands of individuals who may be affected by a functional variant. Similarly, statins (e.g., atorvastatin and rosuvastatin) are prescribed to nearly one in five adults in the US [1] and primarily target *HMGCR*. Due to genetic variation in this target gene, statins have a DRP of 0.18%. This means that of the 40 million individuals who are prescribed a statin in the US, more than 80,000 individuals could be at risk of altered PD of statin treatment due to a functional variant in the target *HMGCR*. This finding is underlined by previous pharmacogenetic studies showing that *HMGCR* is the most important polymorphic gene for treatment success of statins [58].

Overall, the genetic variability of drug targets of many of the top 100 prescribed drugs has not been systematically annotated so far (Additional file 1: Figure S5), including the Alzheimer's drug memantine (DRP = 7.2%), the pain-medication acetaminophen (DRP = 4.7%), and the proton-pump inhibitor esomeprazole (DRP = 3.1%), which all have high DRPs. While these drugs, to our knowledge, are not associated with functional variants in drug targets with regard to their action, clinical studies show that certain proportions of patients treated with them do not respond to treatment. The extent of this non-response is reflected by the number needed to treat (NNT) [59]. For instance, for every one patient successfully treated for Alzheimer's diseases with memantine, between two and seven patients do not respond to treatment [60] (NNT = 3 to 8). Similarly, the NNT for acetaminophen and its indication of pain is five [61] and for esomeprazole and reflux disease is 54 [62].

Drug-related genes show geographic difference in genetic variability

It is known that individuals with different geographic ancestry carry genetic variants with different frequencies [63]. The six populations differentiated in ExAC are of African, South Asian, East Asian, Finnish, Non-Finnish European, and Admixed American (Latino) ancestry [19]. About half of all functional variants in drug-related genes ($M = 54\%$, $SD = 15.2\%$) are unique to only one of the six populations and only 0.1% of functional variants occur with an $AF \geq 0.1\%$ across all populations. Consequently, this results in drug-related genes that have a high risk of functional variants depending on geographic ancestry.

For instance, using a cutoff of $CAP > 1\%$, we found that 231 drug-related genes have functional variants in the cohort of European ancestry compared to 298 genes with functional variants for the cohort of African ancestry.

Nevertheless, 114 drug-related genes showed a CAP score above 1% in each population, indicating that there are genes with a similar world-wide pharmacogenetic relevance.

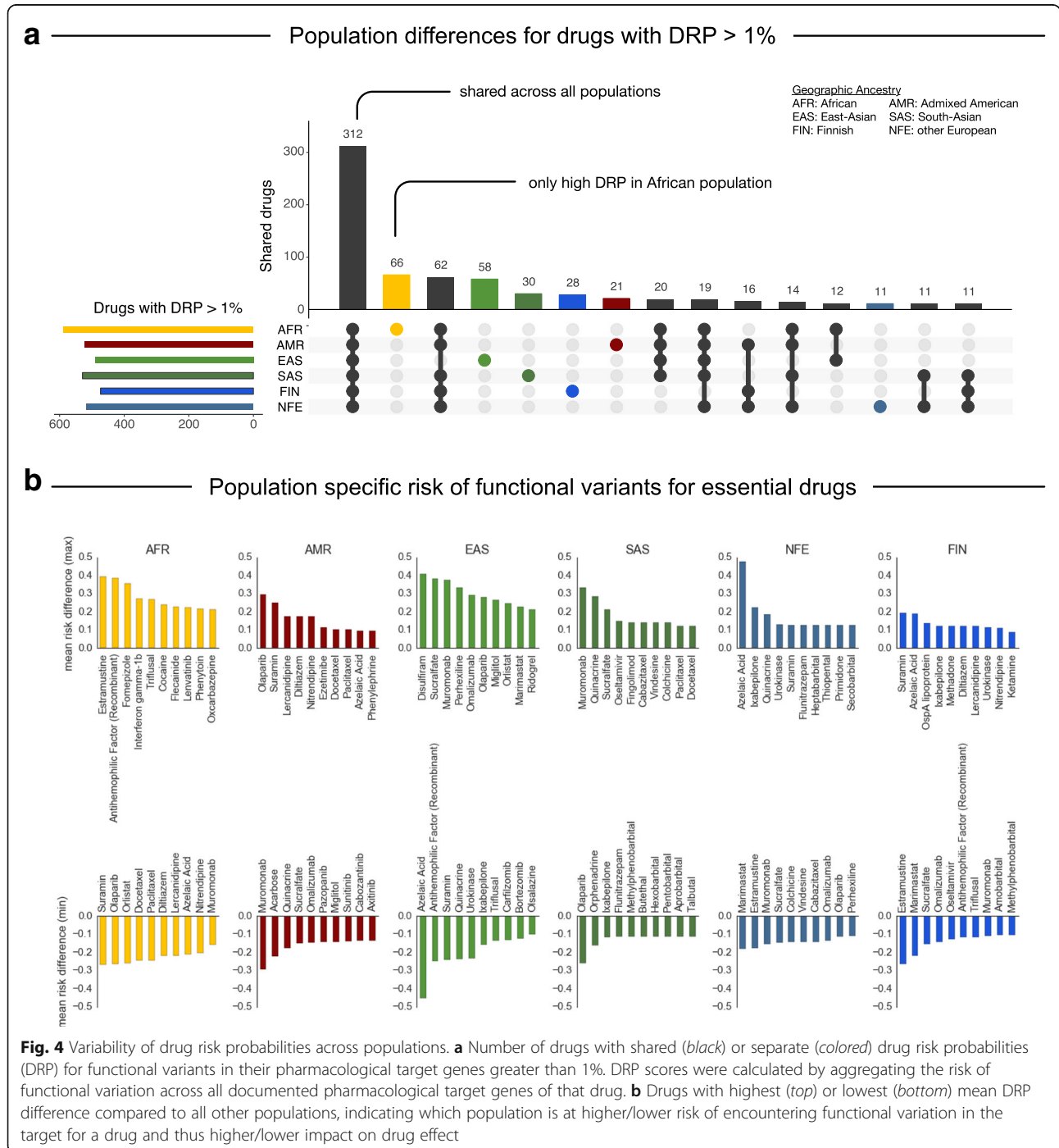
Not surprisingly, amongst those genes with the highest difference in CAP score between populations are many cytochrome P450s and phase II enzymes (Additional file 2: Table S4), as noted in previous studies with smaller population sizes [26]. Similarly, we observe drug target genes with markedly different CAP scores across populations. Among the target genes with the highest absolute CAP score difference are *VWF* (which is targeted by antihemophilic factor), *SIRT5* (targeted by suramin for treating sleeping sickness), and the gastric lipase *LIPF* (targeted by orlistat for obesity treatment). The latter has 65 functional variants and the most frequent variants differ especially between African and East Asian cohorts (CAP 8 vs 51%). Target genes with high subpopulation differences also include several targets for antineoplastic agents, such as the olaparib-target *PARP1*, for which the CAP score ranges from 10.2% in patients of African ancestry to 69.6% in Latino patients. While the efficacy of olaparib depends on the tumor genome and not the germline, the risk to carry germline-originated variants in the tumor should not be ignored. We also observed population differences in the nucleoside transporter *SLC28A1*. While the CAP score is 4% in Non-Finnish Europeans, individuals with an East Asian ancestry have a risk of 60%. Interestingly, several variants in *SLC28A1* have been associated with different outcomes in non-small cell lung cancer and breast cancer [64, 65] when treated with gemcitabine, suggesting that variant differences across the populations may be involved.

Analysis of the DRP score reveals a population-specific risk for several drugs

Of the 1236 FDA approved drugs considered, 241 have more than 10% absolute difference in DRP scores between at least two sub-population cohorts and 24 of these have more than 30% difference (Additional file 2: Tables S5 and S6). Out of this subset of drugs, 11 belong to the 100 most prescribed drugs in the US and 28 are recommended worldwide by the WHO for their therapeutic

use, including oxcarbazepine, amobarbital, and dolasetron. Of the 1236 drugs, 312 have a high risk (DRP >1%) in all six sub-populations (Fig. 4a; the DRP top 20 drugs stratified by population are illustrated in Fig. 4b).

Well-known differences, such as response to disulfiram (treatment for chronic alcoholism), are recapitulated in the data (Fig. 4b). Specifically, the genetic variant E487K in the disulfiram target *ALDH2* (rs671) is seen in the ExAC East Asian population at similarly high frequencies



as seen in previous genetic studies [66]. Similarly, population-specific AFs in ExAC significantly correlate with AFs described in CPIC guidelines for other well studied drug-related genes, such as *TPMT* and *CYP2D6* [67, 68] (Additional file 2: Table S7).

The different responses with the asthma-medication salbutamol and the blood-thinner warfarin have been attributed to variants in their respective drug targets, including R16G in *ADRB2* (rs1042713) for salbutamol [69] and 1639G > A (rs9923231) in *VKORC1* for warfarin [70]. Since the well-known response-altering variants were not annotated by mutation prediction software as functional variants, we did not expect to see the drugs appear high in our ranked list of risk differences across the populations (see “Discussion”). Nevertheless, our analysis shows that salbutamol still has a high risk ratio between populations, caused by 29 variants with a dominant contribution from one variant separating the individuals of Finnish ancestry from those of African ancestry (rs201257377, N69S, $AF_{FIN} = 0.01$). To our knowledge this variant has not been functionally characterized or previously associated with salbutamol response. Similarly, we observe 19 functional variants in the warfarin target *VKORC1* that are population-specific, including a functional variant observed most frequently in individuals of Non-Finnish European or Latino ancestry, (rs61742245, D36Y, $AF_{NFE} = 0.003$, $AF_{Latino} = 0.001$), which has been previously associated with predisposition for warfarin resistance [71]. However, 16 of the functional variants may be novel risk factors, including a functional variant primarily observed in individuals of East Asian ancestry (R53S, ENST00000394975.2:c.157C > A, $AF_{EAS} = 0.001$). Using a recent protein 3D model [41, 42] of *VKORC1*, we mapped the R53S variant to the putative warfarin binding pocket (Fig. 3b). Furthermore, analysis of coevolution in the protein using EVfold [43] shows that R53 is strongly coupled to other residues in the protein and changes in this site are predicted by EVmutation [33] to affect protein fitness due to epistatic variant effects (Additional file 1: Figure S6). Together, this suggests that this mutation might be negatively associated with warfarin binding.

Triflusal, a treatment for stroke re-occurrence, targets four genes (*PTGS1* (also known as Cox-1), *NOS2*, *NFKB1*, and *PDE10A*) that together have more functional variants in the African population than in any other population ($DRP_{AFR} = 37\%$; Fig. 4b). This difference between populations is mainly due to a SNP in *NOS2*, which occurs in the population of African ancestry with higher than average frequency (rs3730017, $AF_{AFR} = 19\%$ vs $AF_{global} = 4\%$) and while not functionally characterized, has been associated with protection against cerebral malaria [72]. In *PTGS1*, three functional variants have AFs above 0.1% in the cohort of African ancestry. The most frequent variant (rs5789, L237M,

$AF_{AFR} = 0.5\%$ vs $AF_{global} = 1.7\%$) lies on the dimer interface and has previously been associated with reduced metabolic activity of the enzyme [73]. A second variant is an indel which is predicted to result in the total loss of protein function ($AF_{AFR} = 0.3\%$ vs $AF_{global} = 0.02\%$). The effects of the third functional variant common in the African cohort (rs139956360, E259A, $AF_{AFR} = 0.2\%$ vs $AF_{global} = 0.02\%$) on enzyme activity or drug binding is less clear from the three-dimensional structure of the protein and would require further exploration. Since triflusal is prescribed for prophylactic use in the same way as aspirin for stroke prevention, it is clearly worth further investigating the effects of these observed functional variants.

Population differences in functional variants for cancer drugs

Our results also highlight a large DRP variability of cancer drugs between the populations. While for many of these drugs it is not the germline but the tumor genome that is relevant for drug action, germline DRPs of these drugs give an estimate of the population risk to possess potentially resistance-causing variants in the tumor and should be screened accordingly. For instance, the DRPs of taxanes (docetaxel, paclitaxel, and cabazitaxel) are 30 percentage points higher in the cohorts of South Asian and European ancestry compared to the cohort of African ancestry ($DRP_{SAS/NFE} = 85\%$ vs $DRP_{AFR} = 45\%$) due to functional variants in the four taxane targets, *TUBB1*, *MAP2*, *MAP4*, and *MAPT*. Among these are three distinct positions in *TUBB1* (Q43P/H, R307C, R359W) that occur with comparably high frequencies in the South-Asian population. While Q43P ($AF_{SAS} = 14\%$) has recently been associated with decreased progression-free survival in urothelial cell carcinoma when treated with cabazitaxel [74], less is known about the effects of the other two variants. Mapping the affected residues onto the 3D structure of docetaxel bound to tubulin (PDB ID 1tub [36]) shows that R359 interacts with the drug (Fig. 3c). The effect of R307C is less obvious from structural observations as it does not lie very close to the binding site or the interface between the monomers in the polymer (R307 to K124 < 15 Å, mapped on PDB ID 3j6g [37]).

Discussion

In this study, we analyzed the extent of functional genetic variation in drug-related genes and its implication for 1236 FDA-approved drugs in exome sequencing data of 60,706 individuals. We show that the risk of carrying functional variants not only in ADME-related genes but also in drug targets is high for an individual patient. For ADME genes this observation is in line with previous studies [12, 15, 18], but it is novel for drug-target genes.

We observed functional variants in 98% of the drug-related genes and at least one high confidence LoF variant in 93% of the genes. The prevalence of functional variants in drug-related genes is thus higher than previously shown [18]. When considering drug target genes for the 100 most prescribed medications in the US the probability of carrying at least one functional variant is above 80% for each patient. Together with the high risk for clinically actionable variants in ADME genes (98% [12]) these findings indicate that genetic variability may contribute significantly to observed differences in drug response between patients.

While individualized cancer therapies often focus on the somatic variants present only in tumor tissue, we can show that functional germline variants, which are routinely masked out in the analysis of somatic variants, are common in many cancer drug targets. By excluding germline variants that the tumor inherited from its progenitor cell from cancer genome analysis in the context of therapeutic decision-making may thus result in the oversight of important determinants for treatment response or resistance development. To what extent the tumor genome varies from the germline genome is dependent on patient and cancer type. Loss of heterozygosity, where the germline allele is lost in the disease progression, and copy number alterations can indeed result in drastic changes between genetic variants observed in the normal tissue of a patient and the cancer [54, 75]. The high prevalence of variants in systemic cancer therapy targets, such as *KDR* for sorafenib, further indicates that the germline variants of target genes in addition to ADME genes should be considered for clinical decision-making.

Geographic ancestry is a well-established confounding factor for drug response, but few drugs have been assessed for efficacy across global populations. Even where clinical trials have been carried out in different populations, particularly non-European and non-Asian individuals remain understudied. By calculating risk probabilities for drugs and different populations, we showed that the frequency of functional variants in drug-related genes varies widely across populations. Even for drugs where population differences in response are observed, additional patient groups may be at high risk of altered PD due to genetic variants in drug targets. Especially for drugs commonly used around the world, such as those on the WHO essential medicines list, this could result in large numbers of patients with reduced drug efficacy in some, but not all, of the populations where they are applied.

The analysis in this study relied on external data for drug variant annotation and drug–gene associations. Even though it was possible to estimate the burden of functional variation in drug-related genes and quantify

to what extent individual drugs may be affected, there remain certain limitations. First of all, even manually curated drug–target associations and pharmacogenomics data are susceptible to spurious annotations. For example, some subunits of the GABA receptors, including *GABRA4*, are generally thought to give rise to receptors resistant to classic benzodiazepines such as diazepam [76] but have been annotated as targets for some benzodiazepines. Comparison to a different, independently curated set of drug–target associations [77] further shows that annotation of drug–target pairs does not always agree. Furthermore, to quantify the real risk for a drug, drug-specific ADME gene relations should be incorporated into the DRP calculation. For example, optimal warfarin dosing is known to be dependent on variants in *CYP2C9* in addition to *VKORC1* [78] and variants in the ADME gene *UGT1A1* are documented to contribute to different responses to the cancer drug irinotecan around the globe [79]. Unfortunately, comprehensive inclusion of ADME genes in the DRP calculations is currently not possible because sufficient data for ADME genes is lacking for most FDA-approved drugs, including the relative contribution of each enzyme. Our DRP estimates thus probably still underestimate the drug-specific risk of functional variation as well as population differences.

The vast majority of variants in drug-related genes considered in this study have not been seen previously and we thus lack validated knowledge about their functional impact on drug efficacy. We therefore had to rely on predictions of their impact on protein function. The probabilities presented are based on the assumption that the functional classification is correct and represents enzyme activity or drug efficacy. The relative risk between genes is based on the assumption that there has not been a significant bias in assessment when genes already have known deleterious mutations. That these assumptions are not always correct follows from the fact that variant classification tools are not exact, are often trained on disease-causing variant sets only, have issues with circularity in the classifier training data, and fail to sub-classify mutations [31]. Overall, variant effect predictors such as SIFT and PolyPhen lean towards over-prediction of deleteriousness [19, 31–33]. In contrast, there are also examples of false negatives, such as the well-studied pharmacogenetic variants in the anti-asthmatic target *ADRB2* (R16G/rs1042713, Q27E/rs1042714, and T164I/rs1800888) [69, 80] that are all misclassified as benign. Although we validated our consensus predictor approach using the meta-predictor CADD [34] and the independent predictor EVmutation [33], the impact of misclassified variants on the CAP and DRP scores should be considered in subsequent interpretation of the results presented here. Furthermore, the field of variant effect prediction is rapidly evolving, and especially

the ability to distinguish between activating and deactivating effects could prove to be crucial for predicting the downstream effects of variants on therapy response.

To increase predictor sensitivity, one could include additional prediction algorithms, which comes with the risk of reduced specificity (in some cases more than half of all non-synonymous variants were classified as functional [15]) as all currently available methods have their individual drawbacks [81]. Reliable computational classification methods for variant effects on drug response remain scarce due to insufficient training data [81], but may arise in the future if efforts are increased to create such data—for example, using novel high throughput methods such as deep mutational scans [82, 83]. For the present study we chose a conservative approach to variant annotation that requires the complete loss of the protein product—which should have a marked impact on the drug—or the consensus prediction of two independent prediction tools at the expense of missing some known variants (Fig. 3a). It is thus not unlikely that the effect of the functional variants is still in part underestimated in our study.

Sequencing data

The use of whole exome sequencing data comes with the intrinsic limitation that only variants in protein coding regions can be detected, potentially missing pharmacologically relevant non-coding variants [84]. Moreover, we only considered non-synonymous variation in relevant genes, thus excluding additional variant types that are known to have an effect on drugs, such as pseudogenes, epigenetics, structural variants, and copy number alterations. These variants are known to have an effect on certain drug-related gene families such as CYP450 [26] and their exclusion may thus result in an underestimation of the pharmacogenomic variability. Furthermore, even at low false-positive rates many called variants can be inaccurate [85] and several pharmacologically relevant gene families—namely CYPs, HLA, and UGTs—are at high risk for variant calling errors due to the complex genetic structure of their loci [86, 87]. While members of the cytochrome P450 family have indeed been found to be problematic in short-read sequencing [26], this does not apply to most other drug-related genes [15, 18]. To reduce the false-positive variant calls in our survey, we included only variants of sufficient locus coverage and high quality.

Homozygous occurrence of variants as well as combinations of variants in the drug-related gene may be required to noticeably alter the drug response in an individual [88]. While homozygous variant counts were reported in the ExAC dataset, and were consistently low, the aggregated format of the data set did not permit the study of particular haplotypes. Predicted effects of

heterozygous variants may thus be compensated for in an affected individual and their effect on the phenotype could be overestimated.

Furthermore, the ExAC cohort is very large in total, but not all populations are represented equally (Additional file 2: Figure S8) [19]. The power to detect very rare variants thus differs by an order of magnitude between the individual populations (from 0.01% AF for the Finnish and East Asian populations to 0.001% for Non-Finnish European). Due to legal restrictions in the underlying exome sequencing projects, sample-specific data, including haplotype phase, are missing in ExAC. Epistatic effects of variants could thus not be investigated, even though they are known to exist. For example, while the single variant rs12248560 (CYP2C18*17) results in increased CYP2C19 activity, the combination with another variant (rs28399504) is associated with LoF of the protein (CYP2C19*4B) [15]. Analysis of such haplotype patterns and comparison of their frequencies in the ExAC cohort to those in previous sequencing studies was not possible.

Implications

Many major medical institutions have started implementing genotyping protocols for preemptive pharmacogenetic testing [89–91]. However, these usually focus on a small number of ADME genes [12] and often only test a subset of established actionable variants using microarrays [92]. While these arrays facilitate fast and cheap screening, we show here that the vast majority of variants in drug-related genes seen in the human population is not covered. We further want to stress that the number of genes with pharmacogenomic variants should systematically include genes implicated in drug mechanism even though only very few examples have yet been characterized well enough to be part of a dosing guideline. Furthermore, with allele frequencies below 0.1%, many functional variants in drug-related genes are so rare that they cannot be observed in clinical trial cohorts, but may contribute to adverse events or diffuse lack of efficacy post-marketing. In the future, this information should be considered in all phases of clinical drug development and the effects of genetic variants in genes associated with PD and PK of the drug candidate should be systematically characterized.

Conclusions

Large-scale sequencing efforts can be used to identify and quantify the extent of genetic variation in genes relevant for drug action and metabolism. Identification of such variants is only the first step towards better treatment decisions. Newly identified variants of pharmacogenomic importance require validation and ultimately updated dosing guidelines. The development of quality-controlled and patient-centered software

solutions to combine available knowledge of pharmacologically actionable variants with a patient's genome as well as fast and accurate approaches (experimental and computational) to functionally classify novel variants will thus be of high importance for the future of personalized medicine.

Additional files

Additional file 1: All supplemental figures cited in the text. Details about each figure are provided in the figure legends below the figures. (PDF 2264 kb)

Additional file 2: All supplemental tables cited in the text. Enclosed data include data set meta-information, CAP scores for all drug-related genes, DRP scores for all drugs, CAP and DRP differences between populations, and a comparison between allele frequencies in the studied data set and CPIC guidelines. (XLSX 776 kb)

Abbreviations

AC: Allele count; ADME: Absorption, distribution, metabolism, and excretion; AF: Allele frequency; ADME: Absorption, distribution, metabolism, and excretion; AF: Allele frequency; CAP: Cumulative allele probability; DRP: Drug risk probability; ExAC: Exome Aggregation Consortium; FDA: Food and Drug Administration; GWA: Genome-wide association; LoF: Loss-of-function; PD: pharmacodynamics; PK: pharmacokinetics; RD: Risk difference; RMSE: Root mean square error; WHO: World Health Organisation

Acknowledgements

We would like to thank Ruomu Jiang for initial help with handling genetic variation data sets, Benjamin Schubert, Fabian Aichler, and Ulrich Mansmann for helpful discussions about the statistical analysis performed in the paper, and Thomas Hopf for support in using the EVmutation toolbox.

Funding

This work was also supported in part by the Robert Bosch Foundation, Stuttgart, Germany and the European Commission Horizon 2020 UPGx grant (668353).

Availability of data and materials

The raw dataset analyzed during the current study are available in the ExAC repository (ftp://ftp.broadinstitute.org/pub/ExAC_release). The subset of functional variants specifically analyzed in this study is further provided as a supplemental file on figshare (<https://doi.org/10.6084/m9.figshare.5631751>). Code is available on github (https://github.com/debbiemarkslab/variants_pharmacogenes).

Authors' contributions

CPS, DSM, and OK designed the study, CPS analyzed the data, DSM and OK helped to analyze the data, RT and MS provided expertise on pharmacogenetics and genomics and contributed in interpretation of the data, and CPS and DSM wrote the manuscript. All authors contributed to editing the manuscript.

Ethics approval and consent to participate

Not applicable. No human subjects or animals were involved in this study. All analyses were performed on aggregated, de-identified data obtained from the public ExAC database.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Publisher's Note

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

Author details

¹Department of Systems Biology, Harvard Medical School, Boston 02115, Massachusetts, USA. ²Center for Bioinformatics, University of Tübingen, 72076 Tübingen, Germany. ³Applied Bioinformatics, Department of Computer Science, 72076 Tübingen, Germany. ⁴Dr. Margarete Fischer-Bosch Institute of Clinical Pharmacology, 70376 Stuttgart, Germany. ⁵Department of Clinical Pharmacology, University Hospital Tübingen, 72076 Tübingen, Germany. ⁶Department of Pharmacy and Biochemistry, University of Tübingen, 72076 Tübingen, Germany. ⁷German Cancer Consortium (DKTK) and German Cancer Research Center (DKFZ), 69120 Heidelberg, Germany. ⁸Quantitative Biology Center, 72076 Tübingen, Germany. ⁹Faculty of Medicine, University of Tübingen, 72076 Tübingen, Germany. ¹⁰Biomolecular Interactions, Max Planck Institute for Developmental Biology, 72076 Tübingen, Germany.

Received: 11 July 2017 Accepted: 24 November 2017

Published online: 22 December 2017

References

- Kantor ED, Rehm CD, Haas JS, Chan AT, Giovannucci EL. Trends in prescription drug use among adults in the United States From 1999-2012. *JAMA*. 2015;314:1818–30.
- Schork NJ. Time for one-person trials. *Nature*. 2015;520:609–11.
- Madian AG, Wheeler HE, Jones RB, Dolan ME. Relating human genetic variation to variation in drug responses. *Trends Genet*. 2012;28:487–95.
- Pirmohamed M. Personalized pharmacogenomics: predicting efficacy and adverse drug reactions. *Annu Rev Genomics Hum Genet*. 2014;15:349–70.
- Mette L, Mitropoulos K, Vozikis A, Patrinos GP. Pharmacogenomics and public health: implementing “populationalized” medicine. *Pharmacogenomics*. 2012;13:803–13.
- O'Donnell PH, Dolan ME. Cancer pharmacogenetics: ethnic differences in susceptibility to the effects of chemotherapy. *Clin Cancer Res*. 2009;15:4806–14.
- Yasuda SU, Zhang L, Huang SM. The role of ethnicity in variability in response to drugs: focus on clinical pharmacology studies. *Clin Pharmacol Ther*. 2008;84(3):417–23.
- Ma Q, Lu AYH. Pharmacogenetics, pharmacogenomics, and individualized medicine. *Pharmacol Rev*. 2011;63:437–59.
- Motsinger-Reif AA, Jorgenson E, Relling MV, Kroetz DL, Weinshilboum R, Cox NJ, et al. Genome-wide association studies in pharmacogenomics: successes and lessons. *Pharmacogenet Genomics*. 2013;23:383–94.
- Daly AK. Genome-wide association studies in pharmacogenomics. *Nat Rev Genet*. 2010;11:241–6.
- PharmGKB. Drug Labels. <https://www.pharmgkb.org/labels>. Accessed 14 Mar 2017.
- Dunnenberger HM, Crews KR, Hoffman JM, Caudle KE, Broeckel U, Howard SC, et al. Preemptive clinical pharmacogenetics implementation: current programs in five United States medical centers. *Annu Rev Pharmacol Toxicol*. 2015;55:89–106.
- van der Wouden CH, Cambon-Thomsen A, Cecchin E, Cheung KC, Dávila-Fajardo CL, Deneer VH, et al. Implementing pharmacogenomics in Europe: design and implementation strategy of the Ubiquitous Pharmacogenomics Consortium. *Clin Pharmacol Ther*. 2017;101:341–58.
- Consortium T1GP. A global reference for human genetic variation. *Nature*. 2015;526:68–74.
- 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*. 2016;19(1):20–29.
- Bush WS, Crosslin DR, Owusu Obeng A, Wallace J, Almoguera B, Basford MA, et al. Genetic variation among 82 pharmacogenes: The PGRNseq data from the eMERGE network. *Clin Pharmacol Ther*. 2016;100:160–9.
- Nelson MR, Wegmann D, Ehm MG, Kessner D, St Jean P, Verzilli C, et al. An abundance of rare functional variants in 202 drug target genes sequenced in 14,002 people. *Science*. 2012;337:100–4.
- Wright GEB, Carleton B, Hayden MR, Ross CJD. The global spectrum of protein-coding pharmacogenomic diversity. *Pharmacogenomics J*. 2016.
- Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, et al. Analysis of protein-coding genetic variation in 60,706 humans. *Nature*. 2016;536:285–91.

20. Ramos E, Doumatey A, Elkahoul AG, Shriner D, Huang H, Chen G, et al. Pharmacogenomics, ancestry and clinical decision making for global populations. *Pharmacogenomics J*. 2014;14:217–22.
21. He Y, Hoskins JM, McLeod HL. Copy number variants in pharmacogenetic genes. *Trends Mol Med*. 2011;17:244–51.
22. Whirl-Carrillo M, McDonagh EM, Hebert JM, Gong L, Sangkuhl K, Thorn CF, et al. Pharmacogenomics knowledge for personalized medicine. *Clin Pharmacol Ther*. 2012;92:414–7.
23. Law V, Knox C, Djoumbou Y, Jewison T, Guo AC, Liu Y, et al. DrugBank 4.0: shedding new light on drug metabolism. *Nucleic Acids Res*. 2014;42:D1091–7.
24. Top 100 Drugs for 2013 by Units—U.S. Pharmaceutical Statistics. <https://www.drugs.com/stats/top100/2013/units>. Accessed 18 Dec 2017.
25. Selection WECOT, Medicines UOE. WHO Model List of Essential Medicines. WHO Technical Report Series. The World Health Organisation; November 2015. <http://www.who.int/medicines/publications/essentialmedicines/en/>.
26. Fujikura K, Ingelman-Sundberg M, Lauschke VM. Genetic variation in the human cytochrome P450 supergene family. *Pharmacogenet Genomics*. 2015;25:584–94.
27. Bertold MR, Cebon N, Dill F, Gabriel TR, Kötter T, Meinel T, et al. KNIME: The Konstanz information miner. In: Preisach C, Burkhardt H, Schmidt-Thieme L, Decker R, editors. *Annual Conference of the German Classification Society*. Berlin, Heidelberg: Springer Berlin Heidelberg; 2008. p. 319–26.
28. MacArthur DG, Balasubramanian S, Frankish A, Huang N, Morris J, Walter K, et al. A systematic survey of loss-of-function variants in human protein-coding genes. *Science*. 2012;335:823–8.
29. Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, et al. A method and server for predicting damaging missense mutations. *Nat Methods*. 2010;7:248–9.
30. Ng PC, Henikoff S. SIFT: predicting amino acid changes that affect protein function. *Nucleic Acids Res*. 2003;31(13):3812–4.
31. Grimm DG, Azencott C-A, Aicheler F, Gieraths U, MacArthur DG, Samocha KE, et al. The evaluation of tools used to predict the impact of missense variants is hindered by two types of circularity. *Hum Mutat*. 2015;36:513–23.
32. van der Velde KJ, de Boer EN, van Diemen CC, Sikkema-Raddatz B, Abbott KM, Knoppers A, et al. GAVIN: Gene-Aware Variant Interpretation for medical sequencing. *Genome Biol*. 2017;18:6.
33. Hopf TA, Ingraham JB, Poelwijk FJ, Schärfe CPI, Springer M, Sander C, et al. Mutation effects predicted from sequence co-variation. *Nat Biotechnol*. 2017;35(2):128–35.
34. Kircher M, Witten DM, Jain P, O’Roak BJ, Cooper GM, Shendure J. A general framework for estimating the relative pathogenicity of human genetic variants. *Nat Genet*. 2014;46:310.
35. Seabold S, Perktold J. Statsmodels: Econometric and statistical modeling with python. *Proceedings of the 9th Python in Science Conference (SciPy 2010)*. 2010.
36. Nogales E, Wolf SG, Downing KH. Structure of the alpha beta tubulin dimer by electron crystallography. *Nature*. 1998;391:199–203.
37. Alushin GM, Lander GC, Kellogg EH, Zhang R, Baker D, Nogales E. High-resolution microtubule structures reveal the structural transitions in α -tubulin upon GTP hydrolysis. *Cell*. 2014;157:1117–29.
38. Cherezov V, Rosenbaum DM, Hanson MA, Rasmussen SGF, Thian FS, Kobilka TS, et al. High-resolution crystal structure of an engineered human beta2-adrenergic G protein-coupled receptor. *Science*. 2007;318:1258–65.
39. Sidhu RS, Lee JY, Yuan C, Smith WL. Comparison of cyclooxygenase-1 crystal structures: cross-talk between monomers comprising cyclooxygenase-1 homodimers. *Biochemistry*. 2010;49:7069–79.
40. Fischmann TO, Hruza A, Niu XD, Fossetta JD, Lunn CA, Dolphin E, et al. Structural characterization of nitric oxide synthase isoforms reveals striking active-site conservation. *Nat Struct Biol*. 1999;6:233–42.
41. Czogalla KJ, Biswas A, Höning K, Hornung V, Liphardt K, Watzka M, et al. Warfarin and vitamin K compete for binding to Phe55 in human VKOR. *Nat Struct Mol Biol*. 2017;24:77–85.
42. Shen G, Cui W, Zhang H, Zhou F, Huang W, Liu Q, et al. Warfarin traps human vitamin K epoxide reductase in an intermediate state during electron transfer. *Nat Struct Mol Biol*. 2017;24:69–76.
43. Marks DS, Colwell LJ, Sheridan R, Hopf TA, Pagnani A, Zecchina R, et al. Protein 3D structure computed from evolutionary sequence variation. *PLoS One*. 2011;6:e28766–17.
44. Johnson LS, Eddy SR, Portugaly E. Hidden Markov model speed heuristic and iterative HMM search procedure. *BMC Bioinformatics*. 2010;11:1.
45. Suzek BE, Wang Y, Huang H, McGarvey PB, Wu CH, UniProt Consortium. UniRef clusters: a comprehensive and scalable alternative for improving sequence similarity searches. *Bioinformatics*. 2015;31:926–32.
46. Pettersen EF, Goddard TD, Huang CC, Couch GS, Greenblatt DM, Meng EC, et al. UCSF Chimera—a visualization system for exploratory research and analysis. *J Comput Chem*. 2004;25:1605–12.
47. Perez F, Granger BE. IPython: a system for interactive scientific computing. *Comput Sci Eng IEEE*. 2007;9:21–9.
48. McKinney W. Data structures for statistical computing in python. *Proceedings of the 9th Python in Science Conference (SciPy 2010)*. 2010.
49. Jones E, Oliphant T, Peterson P. SciPy: Open source scientific tools for Python. 2001. <http://www.scipy.org>.
50. Zanger UM, Schwab M. Cytochrome P450 enzymes in drug metabolism: regulation of gene expression, enzyme activities, and impact of genetic variation. *Pharmacol Ther*. 2013;138:103–41.
51. Mosshammer D, Schaeffeler E, Schwab M, Moerike K. Mechanisms and assessment of statin-related muscular adverse effects. *Br J Clin Pharmacol*. 2014;78:454–66.
52. Rubio-Perez C, Tamborero D, Schroeder MP, Antolin AA, Deu-Pons J, Perez-Llamas C, et al. In silico prescription of anticancer drugs to cohorts of 28 tumor types reveals targeting opportunities. *Cancer Cell*. 2015;27:382–96.
53. Iorio F, Knijnenburg TA, Vis DJ, Bignell GR, Menden MP, Schubert M, et al. A landscape of pharmacogenomic interactions in cancer. *Cell*. 2016;166:740–54.
54. Stratton MR, Campbell PJ, Futreal PA. The cancer genome. *Nature*. 2009;458:719–24.
55. World Health Organization. ATC - Structure and principles. ATC classification and DDD. 2009. <http://www.fhi.no/en/hn/drug/who-collaborating-centre-for-drug-statistics-methodology/>. Accessed Jan 30 2017.
56. Adnane L, Trail PA, Taylor I, Wilhelm SM. Sorafenib (BAY 43-9006, Nexavar (R)), a dual-action inhibitor that targets RAF/MEK/ERK pathway in tumor cells and tyrosine kinases VEGFR/PDGFR in tumor vasculature. *Meth Enzymol*. 2006;407:597.
57. Blake K, Lima J. Pharmacogenomics of long-acting β 2-agonists. *Expert Opin Drug Metab Toxicol*. 2015;11:1733–51.
58. Chasman DI, Posada D, Subrahmanyam L, Cook NR, Stanton VP, Ridker PM. Pharmacogenetic study of statin therapy and cholesterol reduction. *JAMA*. 2004;291:2821–7.
59. Walter SD. Number needed to treat (NNT): estimation of a measure of clinical benefit. *Stat Med*. 2001;20:3947–62.
60. Livingston G, Katona C. The place of memantine in the treatment of Alzheimer’s disease: a number needed to treat analysis. *Int J Geriatr Psychiatry*. 2004;19:919–25.
61. Moore A, Collins S, Carroll D, McQuay H, Edwards J. Single dose paracetamol (acetaminophen), with and without codeine, for postoperative pain. *Cochrane Database Syst Rev*. 1996.
62. Gatta L, Vaira D, Sorrenti G, Zucchini S, Sama C, Vakili N. Meta-analysis: the efficacy of proton pump inhibitors for laryngeal symptoms attributed to gastro-oesophageal reflux disease. *Aliment Pharmacol Ther*. 2007;25:385–92.
63. Henn BM, Botigüé LR, Peischl S, Dupanloup I, Lipatov M, Maples BK, et al. Distance from sub-Saharan Africa predicts mutational load in diverse human genomes. *Proc Natl Acad Sci U S A*. 2016;113:E440–9.
64. Soo RA, Wang LZ, Ng SS, Chong PY, Yong WP, Lee SC, et al. Distribution of gemcitabine pathway genotypes in ethnic Asians and their association with outcome in non-small cell lung cancer patients. *Lung Cancer*. 2009;63:121–7.
65. Wong AL-A, Yap H-L, Yeo W-L, Soong R, Ng SS, Wang LZ, et al. Gemcitabine and platinum pathway pharmacogenetics in Asian breast cancer patients. *Cancer Genomics Proteomics*. 2011;8:255–9.
66. Eng MY, Luczak SE, Wall TL. ALDH2, ADH1B, and ADH1C genotypes in Asians: a literature review. *Alcohol Res Health*. 2007;30:22–7.
67. Relling MV, Gardner EE, Sandborn WJ, Schmiegelow K, Pui C-H, Yee SW, et al. Clinical Pharmacogenetics Implementation Consortium guidelines for thiopurine methyltransferase genotype and thiopurine dosing. *Clin Pharmacol Ther*. 2011;89:387–91.
68. Bell GC, Caudle KE, Whirl-Carrillo M, Gordon RJ, Hikino K, Prows CA, et al. Clinical Pharmacogenetics Implementation Consortium (CPIC) guideline for CYP2D6 genotype and use of ondansetron and tropisetron. *Clin Pharmacol Ther*. 2017;102:213–8.
69. Litonjua AA, Gong L, Duan QL, Shin J, Moore MJ, Weiss ST, et al. Very important pharmacogene summary ADRB2. *Pharmacogenet Genomics*. 2010;20:64–9.

70. Owen RP, Gong L, Sagreya H, Klein TE, Altman RB. VKORC1 pharmacogenomics summary. *Pharmacogenet Genomics*. 2010;20:642–4.
71. Loebstein R, Dvoskin I, Halkin H, Vecsler M, Lubetsky A, Rechavi G, et al. A coding VKORC1 Asp36Tyr polymorphism predisposes to warfarin resistance. *Blood*. 2007;109:2477–80.
72. Trovoada Mde J, Martins M, Ben Mansour R, Sambo MDR, Fernandes AB, Antunes Gonçalves L, et al. NOS2 variants reveal a dual genetic control of nitric oxide levels, susceptibility to Plasmodium infection, and cerebral malaria. *Infect Immun*. 2014;82:1287–95.
73. Lee CR, Bottone FG, Krahn JM, Li L, Mohrenweiser HW, Cook ME, et al. Identification and functional characterization of polymorphisms in human cyclooxygenase-1 (PTGS1). *Pharmacogenet Genomics*. 2007;17:145–60.
74. Duran I, Hagen C, Arranz JA, Apellaniz-Ruiz M, Pérez-Valderrama B, Sala N, et al. SNPs associated with activity and toxicity of cabazitaxel in patients with advanced urothelial cell carcinoma. *Pharmacogenomics*. 2016;17:463–71.
75. Lu C, Xie M, Wendl MC, Wang J, McLellan MD, Leiserson MDM, et al. Patterns and functional implications of rare germline variants across 12 cancer types. *Nature*. 2015;6. <https://www.nature.com/articles/ncomms10086>.
76. Möhler H, Fritschy JM, Rudolph U. A new benzodiazepine pharmacology. *J Pharmacol Exp Ther*. 2002;300:2–8.
77. Santos R, Ursu O, Gaulton A, Bento AP, Donadi RS, Bologa CG, et al. A comprehensive map of molecular drug targets. *Nat Rev Drug Discov*. 2016;16:19–34.
78. Johnson JA, Caudle KE, Gong L, Whirl-Carrillo M, Stein CM, Scott SA, et al. Clinical Pharmacogenetics Implementation Consortium (CPIC) guideline for pharmacogenetics-guided warfarin dosing: 2017 update. *Clin Pharmacol Ther*. 2017;102:397–404.
79. Maitland ML, DiRienzo A, Ratain MJ. Interpreting disparate responses to cancer therapy: the role of human population genetics. *J Clin Oncol*. 2016;24:2151–7.
80. Ortega VE, Meyers DA. Pharmacogenetics: implications of race and ethnicity on defining genetic profiles for personalized medicine. *J Allergy Clin Immunol*. 2014;133:16–26.
81. Han SM, Park J, Lee JH, Lee SS, Kim H, Han H, et al. Targeted next-generation sequencing for comprehensive genetic profiling of pharmacogenes. *Clin Pharmacol Ther*. 2017;101:396–405.
82. Fowler DM, Fields S. Deep mutational scanning: a new style of protein science. *Nat Methods*. 2014;11:801–7.
83. Melnikov A, Rogov P, Wang L, Gnirke A, Mikkelsen TS. Comprehensive mutational scanning of a kinase in vivo reveals substrate-dependent fitness landscapes. *Nucleic Acids Res*. 2014;42:e112.
84. Hanson C, Cairns J, Wang L, Sinha S. Computational discovery of transcription factors associated with drug response. *Pharmacogenomics J*. 2016;16:573–82.
85. Shigemizu D, Fujimoto A, Akiyama S, Abe T, Nakano K, Boroevich KA, et al. A practical method to detect SNVs and indels from whole genome and exome sequencing data. *Sci Rep*. 2013;3:2161.
86. Droegemoeller BI, Wright GEB, Niehaus DJH, Emsley R, Warnich L. Next-generation sequencing of pharmacogenes: a critical analysis focusing on schizophrenia treatment. *Pharmacogenet Genomics*. 2013;23:666–74.
87. Tourancheau A, Margaillan G, Rouleau M, Gilbert I, Villeneuve L, Lévesque E, et al. Unravelling the transcriptomic landscape of the major phase II UDP-glucuronosyltransferase drug metabolizing pathway using targeted RNA sequencing. *Pharmacogenomics J*. 2016;16:60–70.
88. Roden DM, George AL. The genetic basis of variability in drug responses. *Nat Rev Drug Discov*. 2002;1:37–44.
89. Relling MV, Evans WE. Pharmacogenomics in the clinic. *Nature*. 2015;526:343–50.
90. Abbasi J. Getting pharmacogenomics into the clinic. *JAMA*. 2016;316:1533–5.
91. Drew L. Pharmacogenetics: the right drug for you. *Nature*. 2016;537:560–2.
92. Shahandeh A, Johnstone DM, Atkins JR, Sontag J-M, Heidari M, Daneshi N, et al. Advantages of array-based technologies for pre-emptive pharmacogenomics testing. *Microarrays (Basel)*. 2016;5:12.

Submit your next manuscript to BioMed Central and we will help you at every step:

- We accept pre-submission inquiries
- Our selector tool helps you to find the most relevant journal
- We provide round the clock customer support
- Convenient online submission
- Thorough peer review
- Inclusion in PubMed and all major indexing services
- Maximum visibility for your research

Submit your manuscript at
www.biomedcentral.com/submit

