

1 Genetic ancestral patterns in *CYP2D6* alleles: structural variants, rare variants, and clinical 2 associations in 479,144 UK Biobank genomes

3
4 Xiao Jiang¹, Fengyuan Hu¹, Xueqing Zou¹, Ali Abbasi¹, Sri V. V. Deevi¹, Santosh S. Atanur¹,
5 Amanda O'Neill¹, Jen Harrow¹, Margarete Fabre^{1,2,3}, Quanli Wang⁴, Slavé Petrovski¹, William
6 Rae⁵, Oliver Burren¹, Katherine R. Smith^{1,*}

8 Abstract

9 Cytochrome P450 2D6 (*CYP2D6*) is involved in metabolising over 20% of clinical drugs, yet its
10 genetic variation across ancestries is underexplored in large-scale whole-genome sequencing
11 (WGS) datasets. We analysed WGS data from 479,144 UK Biobank participants, identifying 95
12 distinct *CYP2D6* star alleles across five biogeographic groups. Of these, 48 alleles had currently
13 unknown effects. These alleles were more prevalent in African, admixed American, and South
14 Asian groups (~5%) compared to European and East Asian groups (~2%), affecting the ability
15 to provide pharmacogenomics recommendations across ancestries. We identified 99,656
16 (20.8%) individuals carrying *CYP2D6* structural variations and predicted the *CYP2D6* ultra-
17 rapid metaboliser phenotype to be most common in Africans (4.5%) and rarest in East Asians
18 (0.32%). Less than half (45.7%) of rare protein-truncating variant carriers were categorised as
19 poor or intermediate metabolisers, indicating an underrepresentation of rare functional
20 variants in the current *CYP2D6* star allele evaluation. Phenome-wide association studies
21 confirmed links with narcotic allergies and found new associations with plasma BAFFR and
22 BAFF proteins, offering insights for the BAFF-targeted clinical therapy. Collectively, this largest
23 WGS study of *CYP2D6* to date highlights the importance of leveraging all genetic variations for
24 pharmacogenomic insights affecting therapeutic safety and development.

25

26 Introduction

27 Pharmacogenomics (PGx) explores the influence of the genome on an individual's response
28 to medication. Consequently, PGx plays a pivotal role in advancing the principles and practices
29 of precision medicine by utilising genetic characteristics to optimise drug therapy and enhance
30 treatment outcomes. Pharmacogenes, are genes heavily involved in drug metabolism. For
31 example, cytochrome P450 2D6 (*CYP2D6*) is responsible for metabolising over 20% of
32 currently clinically prescribed drugs, particularly those with psychoactive properties¹. *CYP2D6*
33 is highly polymorphic, with over 160 haplotype patterns documented in the
34 Pharmacogenomics Knowledgebase (PharmGKB)². Referred to as star alleles, these
35 encompass combinations of small nucleotide variants (SNVs), small in-frame insertions and
36 deletions (INDELs), and structural variants (SVs)³, including hybrid gene arrangements with its
37 upstream adjacent paralogous pseudogene, *CYP2D7*^{4,5}. Additionally, despite of limited cohort

¹Centre for Genomics Research, Discovery Sciences, BioPharmaceuticals R&D, AstraZeneca, Cambridge, UK; ²Department of Haematology, University of Cambridge, Cambridge, UK; ³Department of Haematology, Cambridge University Hospitals NHS Foundation Trust, Cambridge, UK; ⁴Centre for Genomics Research, Discovery Sciences, BioPharmaceuticals R&D, AstraZeneca, Waltham, Massachusetts, USA; ⁵Clinical Development – Neurology & Ophthalmology, Alexion Pharmaceuticals Inc, Boston, Massachusetts, USA.

*Correspondence: katherine.smith1@astrazeneca.com

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

1 sizes for non-European genetic ancestries, it has been shown that broad genetic ancestral
2 groups exhibit distinctive *CYP2D6* star allele frequencies^{6–10}.

3
4 Pharmacogenomic resources predict an individual's *CYP2D6* enzyme activity level, or gPheno
5 based on their star allele haplotypes, with five potential metaboliser classifications: poor (PM),
6 intermediate (IM), normal (NM), ultra-rapid (UM) or unknown¹¹. The *CYP2D6* enzyme activity
7 can be associated with either poor efficacy or adverse events when specific pharmacogene-
8 metabolised therapeutic treatments are administered¹². Therefore, gPheno prediction can be
9 critical in clinical decision-making for prescribing drugs known to be metabolised by *CYP2D6*.
10 The Clinical Pharmacology Implementation Consortium (CPIC)¹³ has published clinical
11 guidelines for several *CYP2D6*-metabolised drugs, including Atomoxetine¹⁴, Ondansetron and
12 Tropisetron¹⁵, Tamoxifen¹⁶, Tricyclic Antidepressants¹⁷, and certain types of Opioids¹⁸.
13 Recently, studies have evaluated *CYP2D6* genotypes for the UK Biobank (UKBB) participants,
14 integrating imputed genotypes with whole-exome-sequencing data^{19,20}. However, the lack of
15 SV calls in these studies meant that certain loss of function star alleles such as *5, *13, and
16 *36 could not be detected, impacting allele frequency estimation and therefore gPheno
17 prediction.

18
19 In this study, we used the DRAGEN (v3.7.8) *CYP2D6* genotype caller²¹ to identify the *CYP2D6*
20 star allele genotypes for 490,524 multi-ancestry UKBB participants using whole-genome-
21 sequencing (WGS) data²². We used these data to characterise the genetic architecture of
22 *CYP2D6* across five broad genetic ancestries and assess the associations of *CYP2D6* gPheno
23 predictions with 20,488 clinical and molecular phenotypes. This study, the largest WGS
24 analysis of *CYP2D6* to date, highlights the importance of leveraging all types of genetic
25 variation in a large population-based cohort to provide PGx insights that may impact
26 therapeutic safety and development.

27 28 **Methods**

29 **UKBB whole genome Sequencing processing**

30 Whole-genome-sequencing (WGS) data of the UKBB participants were generated by deCODE
31 Genetics and the Wellcome Trust Sanger Institute as part of a public-private partnership
32 involving AstraZeneca, Amgen, GlaxoSmithKline, Johnson & Johnson, Wellcome Trust Sanger,
33 UK Research and Innovation, and the UKBB. The WGS sequencing methods and QC have been
34 previously described^{22,23}. Briefly, genomic DNA underwent paired-end sequencing on Illumina
35 NovaSeq6000 instruments with a read length of 2×151 and an average coverage of 32.5x.
36 Conversion of sequencing data in BCL format to FASTQ format and the assignments of paired-
37 end sequence reads to samples were based on 10-base barcodes, using bcl2fastq v2.19.0.
38 Initial quality control was performed by deCODE and Wellcome Sanger, which included sex
39 discordance, contamination, unresolved duplicate sequences, and discordance with
40 microarray genotyping data checks.

41 42 **AstraZeneca Centre for Genomics Research (CGR) UKB WGS small variant, *CYP2D6* star allele, and 43 copy number variant calling**

44 UKBB genomes were processed at AstraZeneca CGR using the provided CRAM format files. A
45 custom-built Amazon Web Services (AWS) cloud compute platform running Illumina DRAGEN
46 Bio-IT Platform Germline Pipeline v3.7.8 (DRAGEN v3.7.8) was used to align the reads to the
47 GRCh38 genome reference and to call small variants. Small variants were annotated using

1 SnpEff v4.3²⁴ against Ensembl Build 38.92²⁵. We adopted the DRAGEN v3.7.8 in-build *CYP2D6*
2 Caller employing the method in Cyrius²¹ to identify the *CYP2D6* star allele diplotype for each
3 UKBB genome. We retained for downstream analysis calls from 482,033 individuals where the
4 caller could determine the two haplotypes with high confidence (status "PASS"), discarding
5 those where the variants could not be matched to star alleles ("No_call", n=3,594) or were
6 consistent with more than one possible genotype ("More_than_one_possible_genotype",
7 n=4,931). Furthermore, we used *Peddy*²⁶ and 1000 Genomes Project data^{27,28} to classify
8 participants (peddy probability ≥ 0.80) into broad genetic ancestries. Copy number variants of
9 UKBB WGS were called by DRAGEN v3.7.8 germline CNV caller. Post-hoc sample QC and variant
10 QC were applied to ensure a good quality CNV call set (median mendelian violation rate =
11 4.1%; heterozygous de novo rate = 1.9%). According to the overlapping between CNVs and
12 genes, we annotated the functional consequences of CNVs against 19,348 protein-coding
13 genes from Ensembl Build 38.92²⁵. All genes, exons, UTRs from Ensembl Build 38.92
14 annotation were included in the analyses. If a gene has multiple transcripts, we included all
15 transcripts in the analyses. UTRs were defined as the union regions of all transcripts.
16 Promoters were defined as the 1kb window before the transcription start site.

17

18 **Mendelian Consistency in the UKB trios and monozygotic twins**

19 We utilised KING v2.2.7 on the UKBB array genotypes with parameters "--build --degree 1"
20 and determined 1,047 sets of trios (consisting of two parents and an offspring) within the
21 UKBB WGS dataset. Among these, 1,014 trios had *CYP2D6* star allele genotypes called for all
22 three members. Mendelian consistency within trios was determined by confirming that the
23 offspring inherited one haplotype star allele from each parent. Additionally, we identified 362
24 pairs of monozygotic twins from the UKBB WGS data, with *CYP2D6* star allele genotypes called
25 for both individuals in 360 of these pairs. Mendelian consistency in monozygotic twins was
26 defined by the presence of identical *CYP2D6* star allele genotypes in both individuals.

27

28 ***CYP2D6* star allele haplotype frequency estimation**

29 The star allele frequency was assessed separately within each broad genetic ancestral group
30 in the UKBB dataset. Within each group, we initially tallied the presence of each haplotype.
31 Instances of identical gene duplications with different copies were evaluated individually,
32 except for those with three or more copies of *1, *2, *4, where counts were aggregated and
33 categorised as "*1x \geq 3", "*2x \geq 3", and "*4x \geq 3" to facilitate comparison with data in
34 PharmGKB²⁹. Subsequently, the frequency for each star allele haplotype was estimated by
35 dividing the haplotype count by the total number of haplotypes within the respective
36 ancestral groups (i.e., twice the sample size). The estimated frequencies of star allele
37 haplotypes in the AFR, AMR, EAS, EUR, and SAS groups were compared with those of "African
38 American/Afro-Caribbean", "American", "East Asian", "European", and "Central/South Asian"
39 in PharmGKB using the *cor.test()* function within the "stats" package (v4.3.3) in R.

40

41 ***CYP2D6* genetically predicted metaboliser phenotypes**

42 The activity values for the identified star alleles were sourced from the functionality table of
43 *CYP2D6* alleles in PharmGKB³⁰. To estimate the *CYP2D6* metaboliser activity score for each
44 individual, we summed the activity values corresponding to both their star allele haplotypes.
45 For star alleles with predicted duplications or higher-order multiplications, we multiplied the
46 activity value by the number of repetitions. Following the definition provided by CPIC^{11,31},
47 individuals were categorised as follows: 1) *CYP2D6* poor metaboliser (PM) if their activity score

1 equalled zero; 2) *CYP2D6* intermediate metaboliser (IM) if their activity score fell between 0
2 and 1.25 (exclusive); 3) *CYP2D6* normal metaboliser (NM) if their activity score ranged from
3 1.25 to 2.25 (inclusive); 4) *CYP2D6* ultra-rapid metaboliser (UM) if their activity score exceeded
4 2.25 (exclusive); 5) *CYP2D6* indeterminate if they carried a star allele of unknown function. In
5 addition, 23 individuals who carrying unknown effect star alleles but with an activity score
6 greater than 2.25 based on remaining alleles were classified as UM metabolisers.

7

8 ***CYP2D6* predicted protein-truncating variation definition and filtration criteria**

9 As described in the previous study³², we defined protein-truncating variants (PTVs) based on
10 SnpEff v4.3²⁴ annotations of variants as exon_loss_variant, frameshift_variant, start_lost,
11 stop_gained, stop_lost, splice_acceptor_variant, splice_donor_variant, gene_fusion,
12 bidirectional_gene_fusion, rare_amino_acid_variant, and transcript_ablation.

13

14 The applied quality control filters were as follows: a minor allele frequency (MAF) of $\leq 5\%$, a
15 minimum coverage of 10X, and annotation within CCDS transcripts (release 22; approximately
16 34 Mb). For homozygous genotypes, alternate reads were limited to 80%, while for
17 heterozygous variants, the proportion of alternate reads had to be between 0.25 and 0.8.
18 Additionally, a binomial test for the alternate allele proportion in the heterozygous state had
19 to show no significant departure from 50% ($P > 1 \times 10^{-6}$). Other criteria included a genotype
20 quality (GQ) of ≥ 20 , a Fisher's strand bias score (FS) of ≤ 200 for indels and ≤ 60 for SNVs,
21 a mapping quality (MQ) of ≥ 40 , a quality score (QUAL) of ≥ 30 , a read position rank sum
22 score (RPRS) of ≥ -2 , and a mapping quality rank sum (MQRS) of ≥ -8 . Variants had to pass
23 the DRAGEN variant status. The variant site must have at least 10X coverage in over 90% of
24 sequences and must not fail any of the specified quality control criteria in 5% or more of the
25 sequences. Furthermore, the variant site required tenfold coverage in at least 25% of gnomAD
26 exomes³³, and if present in gnomAD exomes, the variant needed an exome z-score of ≥ -2.0
27 and an exome MQ of ≥ 30 .

28

29 **Statistical Associations**

30 We conducted ancestry-specific association tests for 12 predefined contrasts of *CYP2D6*
31 metaboliser categories for 15,909 binary clinical outcomes and 1,656 quantitative traits
32 collected from the UKBB dataset. We tested 11 contrasts defined so that the predicted activity
33 score in cases was higher than controls, and one negative control contrast comparing extreme
34 (PM and UM) to central (IM and NM) categories. We adopted Fisher's exact test for the
35 association analyses against the binary clinical outcomes, and linear regression test for
36 quantitative traits, which were first transformed by rank-based inverse normalisation. The
37 linear regression was adjusted for the age at recruitment and the genetically predicted sex.
38 Subsequently, a meta-analysis was conducted to consolidate the results of the tested
39 phenotypes within each ancestral group. For combining the association results of binary
40 clinical outcomes, we employed the exact conditional test using the *mantelhaen.test()* within
41 the "stats" package (v4.3.3) in R, while the association results of quantitative traits were
42 combined using the fixed-effect inverse variance weighting method. Due to the limited
43 numbers of non-EUR participants with UKBB plasma protein abundance data currently
44 available, we performed a pan-ancestry analysis by regressing the *CYP2D6* metaboliser status
45 for 47,599 individuals against 2,941 plasma protein abundance levels corresponding to 2,923
46 proteins, adopting the same adjustments described previously³⁴. To address multiple testing,

1 we implemented the Bonferroni correction accounting for the number of tested phenotypes
2 and models, setting the corrected significance level at 2.0×10^{-7} .

3 4 **Test for Heterogeneity**

5 In comparison to individuals classified as *CYP2D6* NMs, those identified as PMs and IMs
6 exhibited significant associations with plasma protein abundance of POMC, BAFF and BAFFR.
7 We assessed evidence for heterogeneity in associations with these proteins for PMs vs NMs
8 and IMs vs NMs, or between PMs vs IMs and PMs vs NMs. For example, to test for
9 heterogeneity between PMs vs NMs and IMs vs NMs, we first ensured the independence of
10 associations between the two pairs of *CYP2D6* metaboliser comparisons by dividing *CYP2D6*
11 NM individuals randomly into two subgroups based on the ratio of PM and IM cohort sizes
12 (NM_P and NM_I). Subsequently, we conducted regression analyses for the newly formed
13 metaboliser pairs (PM vs. NM_P ; IM vs. NM_I) against POMC, BAFF and BAFFR, respectively,
14 employing the same adjustments as previously reported³⁴. We then evaluated the
15 heterogeneity of associations using the *rma()* function (method="REML") within the
16 "metafor" package (v4.4.0) in R. This process was repeated 1,000 times, and the median of
17 the P-values for Cochran's Q test was reported. We defined significance of heterogeneity as
18 Cochran's Q test P-value < 5%.

19 20 **Results**

21 ***CYP2D6* genotyping and phenotyping calling**

22 Excluding 36 individuals under consent withdrawals, for the remaining 490,524 UKBB
23 genomes, we were able to call *CYP2D6* star allele genotypes for 481,999 (98.3%) individuals,
24 with additional 1% ambiguous genotype calls ("More than one possible genotype", N=4,931)
25 and 0.7% "No Call" status (N=3,594), returned by DRAGEN v3.7.8. For successful calls, we
26 examined 1,014 trios and 360 monozygotic twins and found Mendelian consistency rates
27 (**Methods**) of 98.8% and 100% respectively, comparable with that reported previously from
28 1000 Genome (1KG) Project²¹. Notably, the 12 families that demonstrated Mendelian
29 inconsistency encompassed at least one *CYP2D6* SV carrier. We used 1KG data to classify
30 individuals in to five broad genetic ancestral groups: African (AFR; N=8,822), admixed
31 American (AMR; N=955), East Asian (EAS; N=2,470), European (EUR; N=456,778), and South
32 Asian (SAS; N=10,119) (**Methods**). Subsequently, we used star allele genotypes to predict an
33 individual's *CYP2D6* activity and overall *CYP2D6* metaboliser phenotype (gPheno). In total,
34 479,144 participants with *CYP2D6* gPhenos were available for downstream analyses (**Figure**
35 **1**).

36
37 Utilising a standalone copy number variant (CNV) calling pipeline (**Methods**), we identified
38 nine *CYP2D6* CNVs spanning the *CYP2D8P-CYP2D7-CYP2D6* genomic region (GRCh38
39 chr22:42,126,499-42,155,001), encompassing 10,373 participants in the UKBB. Among these,
40 9,956 individuals were involved into our study, with 9,940 (99.8%) being classified as SV
41 carriers based on the identified *CYP2D6* star alleles as well (**Supplementary Table 1**).

42 43 **Underrepresentation of diversity in the current evaluation of *CYP2D6* genotypes**

44 In total, we identified 95 distinct *CYP2D6* star alleles. These star alleles were unevenly
45 distributed among the five broad genetic ancestry groups (**Figure 2A, Supplementary Table**
46 **2**), with the proportions of both non-functional (activity value=0) and normal functioning
47 (activity value=1) star alleles being the smallest in the EUR individuals. Notably, of the 95 star

1 alleles, 48 (50.5%) lacked known effects on *CYP2D6* enzyme activities according to PharmGKB.
2 Overall, we found these constituted between 30% to 50% of star alleles identified across the
3 genetic ancestries examined (**Figure 2A**). Although indeterminate star alleles were only
4 observed in 1.9% of UKBB participants, they led to challenges in gPheno classification for, on
5 average 5% of AFR, AMR and SAS individuals in contrast to 2% of EUR and EAS.

6
7 In total, we were able to predict gPhenos for 424,682 individuals (88.6%) whose star allele
8 diplotype, activity scores, and the corresponding predicted phenotype were documented in
9 PharmGKB, where all predicted metaboliser phenotypes aligned consistently with the *CYP2D6*
10 gPhenos within our study. The remaining 54,462 individuals (11.4%) all carried non-identical
11 SVs with duplications or higher-order multiplications (e.g., *10+*36, *4+*68). We further
12 conducted genetic ancestry-specific analysis to assess the consistency of *CYP2D6* star allele
13 haplotype frequencies to those reported in PharmGKB (**Methods; Supplementary Figure 1;**
14 **Supplementary Table 3**). Overall, we observed a high degree of correlation in haplotype
15 frequencies among four genetic ancestries: AFR ($r=0.97$), AMR ($r=0.99$), EUR ($r=0.98$), and SAS
16 ($r=0.97$). Notably, the correlation in individuals with EAS was lower ($r=0.86$). This was primarily
17 due to the absence of frequency information for the *10+*36 star allele combination in
18 PharmGKB as we were able to improve correlation to 0.89 by incorporating information from
19 a Chinese- and Malay-focused cohort¹⁰. As previously reported¹⁰, we observed a lower
20 frequency of *10 within EAS individuals with respect that reported by the PharmGKB. We also
21 compared star allele frequencies of the UKBB genetic African ancestry with those of “African
22 American/Afro-Caribbean” and “Sub-Saharan African” documented in the PharmGKB
23 (**Supplementary Figure 2; Supplementary Table 4**).

24
25 We selected the five most common *CYP2D6* star allele haplotypes within each genetic
26 ancestry and then compared their prevalence across the remaining genetic ancestry groups
27 (**Figure 2B**). This highlighted ten distinct star alleles, including two structural variants
28 (*CYP2D6**5, *10+*36). Among these ten star alleles, three were normal-functioning alleles
29 (activity value=1) of which two were common across ancestries (*1, *2), with the other (*35),
30 found to be more common in EUR (5.1%) and AMR (3.2%) individuals. Five were reduced-
31 functional alleles (activity value=0.25 or 0.5) including *17 and *29 that were over nine-times
32 more frequent in AFR individuals than other genetic ancestries. The *10 and *10+*36 hybrid
33 arrangements were most frequent in EAS individuals (15.6% and 32.3%, respectively);
34 whereas the remaining *41 had a significant enrichment in the SAS individuals. Finally, the
35 two remaining star alleles were non-functional (activity value=0). These included the EUR-
36 enriched *4 (14.2%) and the whole gene deletion variation *5 most frequently found in AFR
37 and EAS individuals (**Figure 2B; Supplementary Table 3**).

38 39 **Evaluation of *CYP2D6* structural variations**

40 *CYP2D6* SVs are essential for precisely determining an individual's *CYP2D6* gPheno, particularly
41 for the UM category. These *CYP2D6* SVs can include whole gene deletions (i.e., *5),
42 duplications or higher-order multiplications of identical gene copies (e.g., *1x2, *4x2) and
43 non-identical gene copies (e.g., *4+*68, *10+*36), singleton hybrid genes (e.g., *4.013, *13),
44 and combined arrangements (e.g., *4+*68x2, *10+*36x2). We found that 99,656 (20.8%) of
45 participants carried at least one *CYP2D6* SV, and that these individuals were spread unevenly
46 across the five broad genetic ancestry groups (**Figure 3A**). Notably, 63.6% EAS individuals
47 carried one or more *CYP2D6* SVs, a figure higher than a previous estimate of 55.6% in an

1 admixed Asian population from Singapore¹⁰. To investigate this, we randomly selected a
2 comparable cohort size from UKBB individuals (N=2,000) with proportions of EAS and SAS
3 individuals matched to those in the Singaporean study. By 1,000-time iterations, we observed
4 a 58.3% [IQR 57.9-58.7%] median SV carrying frequency. Therefore, we conclude that the
5 observed discrepancy (63.6% vs 55.6%) is most likely due to differences in admixture between
6 the two studies.

7

8 **Gene deletions** *CYP2D6**5 denotes a complete deletion of the *CYP2D6* gene, resulting in a
9 total loss of enzyme activity. Overall, it was one of the most common SVs with 31,388 (31.5%
10 of the total *CYP2D6* SV carriers) carriers and consistent with a previous report³⁵, it was most
11 frequent in AFR individuals. Furthermore, we discovered a more than two-fold higher
12 frequency of homozygous carriers (*5/*5) in AFR individuals (0.43%) compared to individuals
13 within the other groups (ranging from 0.11% AMR ancestry to 0.16% EAS ancestry). We
14 observed that whilst the frequency of the *5 allele in the East Asian cohort (5.2%) was
15 identical to that reported in PharmGKB²⁹, it was higher than the estimate from a recent
16 *CYP2D6* SV-focused study³.

17

18 **Identical gene duplications and multiplications** *CYP2D6* duplication and multiplication SV
19 events are key to predicting the UM category (**Supplementary Table 5**) and we identified
20 16,144 carriers (3.4% of the overall population). Among them, 181 were homozygous carriers,
21 with 170 heterozygous participants carrying alternative identical gene duplication or
22 multiplication haplotypes (e.g., *1x2/*2x2, *1x3/*2x3). Overall, we observed a higher
23 frequency of carriers in non-EUR individuals (**Figure 3A, Supplementary Table 6**), for example,
24 one in seven (14.5%) of AFR participants were identified as identical gene duplication or
25 multiplication carriers. Furthermore, we found that *10+*36xN exhibited the widest range of
26 and highest copy numbers in the EAS participants (N=1, 2, 3, 4, 5, 6, 22; **Supplementary Table**
27 **3**). Interestingly, we did not observe any multicopy carriers of the *4 allele in EAS participants,
28 which is the most frequent star allele contributing to non-functional *CYP2D6* enzyme in EUR
29 individuals (haplotype frequency 14%), further demonstrating the ancestry-specific
30 distribution of *CYP2D6* star alleles.

31

32 **Hybrid genes** Certain *CYP2D6* star alleles indicate unequal recombination between *CYP2D6*
33 and the adjacent homologous *CYP2D7* pseudogene that create hybrid genes. The naming
34 convention for these hybrid genes depends on the orientation of the sequence, with the gene
35 occupying the 5' portion mentioned first. In our cohort, we observed the currently only one
36 *CYP2D7-2D6* hybrid allele (*13) and three *CYP2D6-2D7* hybrid alleles (*4.013, *36, *68). We
37 found that singleton *13, *36, and *68 hybrids were rare (**Supplementary Table 3**), and we
38 observed only one unique gene recombination arrangement for *13 (*CYP2D6**2+*13), *4.013
39 (*CYP2D6**4+*4.013), and *36 (*CYP2D6**10+*36). While, for *68, we observed three non-
40 identical gene duplications: *1+*68, *4+*68, and *45+*68, which were most frequent in AMR
41 (2.1×10^{-3}), EUR (0.05), and AFR (6.8×10^{-4}) groups, respectively. Consistent with previous
42 findings³, 26.3% of *CYP2D6**4 European carriers had the hybrid arrangement pattern:
43 *4+*68xN (N=1,2).

44

45 In summary, we observed that 54.6% *CYP2D6* SV carriers consisted of hybrid genes with non-
46 identical gene duplications or multiplications. While Gustafson et. al. (2024)³⁶ stated it could
47 be challenging to identify both the gene deletion and the hybrid arrangements in the same

1 individual for *CYP2D6* in the short read WGS, in total, we found that 2,003 participants
2 carrying a whole gene deletion (*5) and another type of SV on different haplotypes (**Figure**
3 **3B**). Among them, 1,898 individuals had a combination of whole gene deletion (*5) and hybrid
4 genes with non-identical gene duplication or multiplication arrangements, including the
5 recently reported *5/*10+*36³⁶. Under this genotype setting, the most frequent was
6 *CYP2D6**4+*68/*5 carried by 1,533 individuals. Additionally, we found 84 individuals with the
7 genotype combining *CYP2D6**5 and identical gene duplications or multiplications, including
8 *5/*17x2, *5/*36x2, and *5/*43x2. Furthermore, we found 21 individuals had genotypes
9 combining *CYP2D6**5 and singleton hybrid genes, which were *5/*13, *5/*36, and *5/*68.

11 **Underrepresentation of predicted rare protein-truncating variants in the current evaluation** 12 **of *CYP2D6* genotypes**

13 We used the known effects of star alleles on *CYP2D6* enzyme activity, to classify individuals
14 into poor (PM), intermediate (IM), normal (NM), and ultra-rapid (UM)¹¹ metabolisers
15 (**Methods**). Cohort-wide, 49.9% of individuals were predicted as NM, 39.6% as IM, 7.0% as
16 PM and 1.6% as UM, with the remaining 1.9% as indeterminate because of the unknown effect
17 of one or more of the identified star alleles on enzyme activity status (**Figure 4A, Table 1**).
18 Across genetic ancestries, the prevalence of NM was higher in non-EUR participants,
19 especially in AMR and SAS genetic ancestries where it represented almost two thirds of
20 individuals.

21
22 We identified 78 predicted protein-truncating variants (PTVs) with a minor allele frequency
23 (MAF) \leq 5% in the *CYP2D6* coding regions (**Methods, Supplementary Table 7**), carried by a
24 total of 25,855 individuals. Overall, we observed a significant enrichment, with 97.1% of PTV
25 carriers showing reduced *CYP2D6* enzyme activity (i.e., categorised as PM or IM). However,
26 this enrichment was significantly lower in the AFR (76.6%), EAS (72.0%), and SAS (85.7%)
27 cohorts (**Figure 4B, Supplementary Table 8**). This variation may be attributed to different
28 patterns of linkage disequilibrium between *CYP2D6* star alleles and PTVs across different
29 genetic ancestral groups.

30
31 We found that 76 of the *CYP2D6* predicted PTVs were rare mutations with a MAF of < 0.1%.
32 65 of which were not included in the current evaluation of *CYP2D6* star alleles and were
33 carried by 269 individuals. Notably, less than half (45.7%) of these people were predicted with
34 reduced *CYP2D6* enzyme levels (i.e., PM or IM). The proportion was even lower in the AFR and
35 SAS cohorts, at 12.5% and 16.7% respectively (**Figure 4B, Supplementary Table 8**). This
36 indicates an underrepresentation of rare functional variants in the current *CYP2D6* star allele
37 definition, potentially affecting the accuracy of metaboliser phenotype predictions across
38 different genetic ancestral groups, disproportionately.

39 **Associations between predicted *CYP2D6* metabolisers and clinically relevant phenotypes**

40 We sought to examine whether there were associations between *CYP2D6* gPhenos and
41 UKBB³⁷ phenotypes by performing a phenome-wide association analysis, including the plasma
42 abundance of 2,923 proteins from the UKBB PPP data^{34,38} in a subset of 47,599 individuals.
43 Considering the uneven distribution of predicted *CYP2D6* gPhenos, we established 12 binary
44 pairs (**Methods, Supplementary Table 9**). These pairs involved comparisons among various
45 combinations of *CYP2D6* metaboliser groups, including a negative control that combines PMs
46 with UMs, contrasted with a comparison involving IMs and NMs.
47

1
2 Following the adjustment for multiple testing (corrected $P=2 \times 10^{-7}$; **Methods**), we uncovered
3 significant associations between one or more of our *CYP2D6* gPheno contrasts and two binary
4 clinical outcomes, six quantitative traits, and changes of three plasma protein levels (**Figure**
5 **5A&B, Supplementary Table 10-12**).

6
7 Importantly, we identified a significant association between *CYP2D6* NM and UM individuals
8 compared to PM and IM metabolisers with 'Allergy events related to narcotic agents' (ICD10:
9 Z88.5, OR=1.19, 95% CI: [1.12-1.27], $P=4.5 \times 10^{-9}$), serving as a positive control¹⁸. According to
10 the CPIC guidelines, individuals with *CYP2D6* UM status should avoid certain opioid drugs
11 (such as codeine and tramadol) because of the risk of severe toxicity, whereas PMs should
12 avoid these medications because of lack of efficacy¹⁸. Underscoring this when we limited
13 association analysis to UM vs PM metaboliser groups, we observed a higher point estimate of
14 the effect size (OR=1.87, 95% CI: [1.47, 2.37], $P=3.5 \times 10^{-7}$). Additionally, we identified an
15 association between increased risk of developing Calculus of the kidney in NM & UM gPhenos
16 compared to PM & IM (OR=1.19, 95% CI: [1.12, 1.27], $P=8.9 \times 10^{-8}$).

17
18 Quantitative trait associations included ankle spacing width (PM & IM vs NM; $\beta=0.02$, 95%
19 CI: [0.02, 0.03], $P=1.1 \times 10^{-11}$), serum creatinine (PM & IM vs NM & UM; $\beta=-0.01$, 95% CI: [-
20 0.02, -0.01], $P=2.5 \times 10^{-9}$), distance to coast (PM & IM vs NM; $\beta=-0.02$, 95% CI: [-0.02, -0.01],
21 $P=1.5 \times 10^{-8}$), spherical power (PM vs NM; $\beta=0.07$, 95% CI: [0.04, 0.09], $P=7.2 \times 10^{-8}$), avMSE
22 (PM vs NM; $\beta=0.06$, 95% CI: [0.04, 0.09], $P=1.6 \times 10^{-7}$), and records in HES inpatient main
23 dataset (PM vs NM; $\beta=0.03$, 95% CI: [0.02, 0.04], $P=3.5 \times 10^{-7}$), while the molecular
24 mechanisms underlying the identified associations remained further validation in other data
25 resources. It is possible that some of these results may represent noise or confounding by
26 fine-scale ancestry stratification.

27
28 Looking at the subset of 47,599 individuals with plasma proteomics data we found that
29 compared to normal *CYP2D6* metabolisers (NM), PM and IM individuals had significantly
30 elevated plasma levels of B-cell-activating factor receptor (BAFFR; encoded by *TNFRSF13C* –
31 PM & IM (0) vs NM (1): $\beta=-0.11$, 95% CI: [-0.13, -0.09], $P=5.2 \times 10^{-34}$), and decreased levels
32 of its ligand BAFF (encoded by *TNFSF13B* - PM & IM (0) vs NM (1): $\beta=0.07$, 95% CI: [0.05,
33 0.09], $P=5.2 \times 10^{-15}$). For both proteins, effect sizes estimated for PM vs NM compared to those
34 for IM vs NM were significantly different (**Methods**; Cochran's Q test median $P_{\text{BAFFR}} = 1.9 \times 10^{-4}$,
35 $P_{\text{BAFF}} = 0.043$), with the former almost twice of the effect size as that of the latter (**Figure 5C,**
36 **Supplementary Table 12**). In addition, the estimated effect size for PM vs IM was significantly
37 different from that for PM vs NM in the association with plasma BAFFR levels, while at the
38 border line of significance for that with plasma BAFF levels (**Methods**; Cochran's Q test median
39 $P_{\text{BAFFR}}=4.8 \times 10^{-3}$, $P_{\text{BAFF}}=0.06$). Given the role BAFF and BAFFR play in the long-term survival of B
40 lymphocytes^{39,40}, we assessed *CYP2D6* metaboliser status associations with autoimmune
41 diseases and cancers. We found that PM and IM individuals had a 30% and 40% lower risk in
42 developing chronic myeloid leukaemia and follicular non-Hodgkin's lymphoma compared to
43 PMs at the significance level of 0.1% and 0.9%, respectively. In addition, *CYP2D6* gPhenos were
44 not associated with the plasma expressions of other measured TNF family members relevant
45 to B lymphocyte development (**Figure 5C**). We found that the previously reported³⁸ variant
46 rs763882049, a pQTL for both BAFFR and BAFF, was enriched in the *CYP2D6* NMs (OR [95% CI]
47 = 1.86 [1.84, 1.90]; $P=0$), with a consistent direction of effects for BAFFR and BAFF. Four

1 *CYP2D6* star alleles were enriched for rs763882049 (**Supplementary Table 13**) including *2
2 where we observed a three-fold increase of prevalence for rs763882049 carriers. Analyses
3 conditioning on the rs763882049 genotype resulted in an attenuated association of *CYP2D6*
4 PM/IM vs NM with BAFFR (beta=-0.02, 95% CI: [-0.04, -5x10⁻³], P=0.012). Additionally, we
5 observed that NM individuals were more likely to carry rarer protein coding pQTLs³⁴ in BAFFR
6 (*TNFRSF13C*), associated with lower abundance (**Supplementary Table 14**).

7
8 Furthermore, when contrasting PM and IM, we identified *CYP2D6* NM was negatively
9 associated with the plasma protein abundance of Proopiomelanocortin (POMC; PM & IM vs
10 NM; beta=-0.05, 95% CI: [-0.07, -0.03], P=8.5x10⁻⁹), a precursor polypeptide producing various
11 hormones, including β-endorphin. In addition, the estimated effect size for PM vs NM was
12 significantly different from that for IM vs NM (**Methods**; Cochran's Q test median P-
13 value=0.041). However, the estimated effect size for PM vs IM was not significantly different
14 from that for PM vs NM (**Methods**; Cochran's Q test median P-value=0.19). The average
15 plasma POMC abundance in the *CYP2D6* UM individuals was less than that of *CYP2D6* PM
16 population at the nominal statistical significance level (P=0.02) and was not significantly
17 different from that of *CYP2D6* IM or NM individuals (**Supplementary Table 15**).

18 19 **Discussion**

20 Pharmacogenomics research is pivotal to the development of personalised medicine, with
21 *CYP2D6* emerging as a critical component, given its role in metabolising over one-fifth of the
22 clinically prescribed drugs¹. This study of 479,144 whole-genome sequenced UKBB individuals
23 presents the most extensive multi-ancestry study of *CYP2D6* genetic variation to date,
24 allowing the identification of 99,656 (20.8%) carriers of *CYP2D6* SVs.

25
26 Across the five broad genetic ancestry groups, we identified 95 distinct *CYP2D6* star alleles,
27 with the highest number observed in EUR individuals, contrary to expectation given higher
28 AFR genetic diversity⁴¹. We attributed this to the limited cohort size of non-EUR genetic
29 ancestries in the UKBB, which may have restricted our ability to identify less common and
30 ancestry-specific star alleles. We were unable to assign function and therefore *CYP2D6*
31 metaboliser for 48 of the star alleles (50%), however, collectively these were found in just 1.9%
32 of individuals. Notably, these individuals were most prevalent in AFR, AMR, and SAS
33 individuals (4-5%) highlighting the need to identify and functionally characterise *CYP2D6*
34 variation in these ancestries to facilitate equitable capacity to provide pharmacogenomic
35 advice.

36
37 Comparing our findings with the published *CYP2D6* frequencies from PharmGKB, we observed
38 high concordance in star allele haplotype frequencies for AFR, AMR, EUR, and SAS broad
39 genetic ancestry groups ($r_{\min}=0.97$). However, concordance was lower ($r=0.89$) for EAS,
40 primarily due to a discrepancy in the estimated allele frequency for *10 that might be due to
41 the incomplete resolution of the *10+*36 haplotype in PharmGKB²⁹. We also replicated an
42 underestimation of the allele frequency of *1 for SAS genetic ancestries in PharmGKB¹⁰, as
43 well as that for AFR and EUR.

44
45 *CYP2D6* SVs were observed across broad genetic ancestry groups at different frequencies, with
46 highest rates (63.6%) observed in EAS. Overall, the most frequent (31.5%) SV was found in the
47 *CYP2D6**5 (whole gene deletion haplotype), highlighting its prevalence in the population.

1 Whilst, ultra-rapid metaboliser (UM) phenotype invariably involved identical gene
2 duplications or multiplications, we found that other metaboliser phenotypes could also
3 feature such duplications, (e.g. *10+*36xN, activity value=0.25). Additionally, we found that
4 non-identical duplications often comprised hybrid genes with ancestry-informative
5 characteristics. Overall, these observations emphasise the need to detect *CYP2D6* SVs in
6 clinical testing to ensure accurate therapeutic interventions and minimise potential adverse
7 reactions especially in underrepresented genetic ancestries.

8
9 In a previous study that utilised established qualifying models developed for phenome-wide
10 association studies (PheWAS)³², we identified individuals carrying computationally predicted
11 deleterious genetic mutations affecting the functionality of their protein-coding genes. Nearly
12 all (99%) cis-pQTL signals from the *ptv* model were linked to decreased plasma levels,
13 indicating the reliability of variation annotations³⁴. The strong enrichment (97.1%) of more
14 common PTV variant (MAF≤5%) carriers in reduced enzyme activity groups (i.e., IM or PM)
15 corroborated with predictions based on called *CYP2D6* star alleles. However, 65 (85.5%) rare
16 PTVs (MAF≤0.1%) were not considered in the current *CYP2D6* star allele definition. Less than
17 half of these rare PTV carriers (45.7%) were predicted to have reduced *CYP2D6* enzyme activity
18 levels (i.e., PM or IM) based on the current star allele evaluation, with further reduction of the
19 proportion in AFR (12.5%) and SAS (16.7%) cohorts, indicating the underrepresentation of the
20 *CYP2D6* rare functional variants in the current star allele definition. Although the contribution
21 may be marginal, incorporating further *CYP2D6* rare functional variants into star allele
22 definitions could improve the prediction of the *CYP2D6* metaboliser activity for some
23 individuals, especially in non-European ancestry populations. Further functional analyses are
24 warranted to integrate newly identified rare functional variations into the expanding
25 repertoire of *CYP2D6* star alleles.

26
27 A Phenome-wide analysis revealed the increased risk of *CYP2D6* NM and UM individuals of
28 allergic reactions to narcotic agents, such as opioids like morphine and oxycodone, serving as
29 a positive control. Using plasma proteome data available for a subset of approximately 50,000
30 individuals, we observed a significant association of *CYP2D6* PMs and IMs with the elevated
31 level of plasma POMC, compared to NMs. POMC is a precursor polypeptide that produces β-
32 endorphin, an endogenous opioid that binds to the μ-opioid receptor⁴². This receptor is also
33 targeted by exogenous opioids like morphine⁴³ and fentanyl⁴⁴. *CYP2D6* is crucial for
34 metabolising the prodrug codeine into morphine, facilitating pain relief. However, *CYP2D6* PM
35 individuals do not effectively convert codeine, leading to reduced effect in pain relief⁴⁵. An
36 increase in β-endorphin among these individuals may serve as a natural compensatory
37 mechanism. Conversely, no significant differences were found in plasma POMC expression
38 between PM/IM/NM and UM individuals, which might be because of the considerable
39 variation in the plasma POMC expression due to the limited size of the UM cohort. *CYP2D6*
40 UMs, due to their heightened enzyme activity, may produce excessive morphine, potentially
41 leading to toxic effects after codeine ingestion.

42
43 Furthermore, studying the UKBB PPP data enabled us to identify significant *cis*- and *trans*-
44 associations between *CYP2D6* metaboliser phenotypes and plasma levels of BAFF-R and BAFF,
45 respectively. The *cis*-association with BAFFR might be attributed to the linkage between the
46 previous identified³⁸ BAFFR *cis*-pQTL (rs763882049) with certain *CYP2D6* star alleles. BAFF is
47 a well-known target for belimumab, which has been approved for the clinical treatment of

1 systemic lupus erythematosus⁴⁶. Although no direct associations were found between *CYP2D6*
2 metaboliser phenotypes and common autoimmune diseases in the UKBB dataset, our findings
3 suggest a potential gene-drug relationship between *CYP2D6* and BAFF-targeted autoimmune
4 clinical treatments, warranting further dosage optimisation experiments tailored to different
5 *CYP2D6* metaboliser phenotypes, and highlighting the potential need of *CYP2D6* metaboliser
6 level assessment prior to the BAFF-targeted clinical practice.

7
8 We acknowledge several limitations to our study. Firstly, the UKBB dataset is predominantly
9 enriched in EUR participants (95.3% in our study), limiting our ability to fully survey ancestry
10 informative *CYP2D6* genetic features for non-EUR cohorts. However, our study highlights the
11 opportunities to deliver more equitable access to medicines through the growing availability
12 of large scale whole genome sequencing data of non-European broad genetic ancestry⁴⁷. We
13 chose not to study the relationship between *CYP2D6* genetic variation and clinical prescription
14 data in the UKBB as these data were only available for a subset of individuals (45.2%), there is
15 an average five-year interval between phenotype measurements, and there is an absence of
16 data concerning detailed consequences for changes in drug types and dosages, which
17 complicates modelling the efficacy of medication and dosage changes, as well as the
18 associated occurrences of side effect over time. Finally, functional analyses will be required to
19 validate the impact on enzyme function for the novel *CYP2D6* variants that we describe and
20 how these might affect clinically relevant phenotypes.

21
22 In summary, our study, is to our knowledge, the largest to date to employ WGS data to analyse
23 *CYP2D6* pharmacogenomics across multiple ancestries. We highlight the lack of genetic
24 ancestral diversity in the current evaluation of *CYP2D6* star alleles and emphasise the
25 importance of including structural variations, available through WGS, to accurately predict
26 *CYP2D6* metaboliser phenotypes. Our results reveal clinically significant associations that
27 enhance the understanding of molecular mechanisms underpinning previously identified drug
28 responses and suggest potential dosage optimisation for BAFF-targeted treatments.

29 30 **Data and code availability**

31 The UK Biobank data were available via the registration for access procedure described at
32 <https://www.ukbiobank.ac.uk/enable-your-research>. Association tests described in this study
33 were performed using a custom framework, PEACOK (PEACOK 1.0.7), that is available via
34 GitHub (<https://github.com/astrazeneca-cgr-publications/PEACOK/>).

35 36 **Acknowledgments**

37 We thank the participants and investigators in the UKB study who made this work possible
38 (Resource Application Number 26041). We are grateful to the research and development
39 leadership teams at the 13 participating UKB-PPP member companies (AInylam
40 Pharmaceuticals, Amgen, AstraZeneca, Biogen, Bristol-Myers Squibb, Calico, Genentech,
41 Glaxo Smith Klein, Janssen Pharmaceuticals, Novo Nordisk, Pfizer, Regeneron and Takeda) for
42 funding the study.

43 44 **Author contributions**

45 X.J. and K.R.S conceptualised and designed this study. X.J and F.H. performed analyses and
46 statistical interpretation. X.Z.Z. and S.S.A. conducted the structural variation analyses using
47 an independent calling pipeline. S.V.V.D did the bioinformatics processing. A.A., G.A., A.O.,

1 J.H., and M.F. contributed to biological interpretation. X.J., O.B., and K.R.S. drafted the main
2 text and supplementary materials. Q.W., S.P., and W.R. supervised the study. All authors
3 read, commented on, and agreed upon the submitted manuscript.

4 5 **Declaration of interests**

6 X.J., F.H., X.Z.Z., A.A., S.V.V.D., S.S.A., G.A., A.O., J.H., M.F., Q.W., S.P., O.B., K.R.S. are current
7 employees and/or stockholders of AstraZeneca. W.R. is a current employee and stockholder
8 of Alexion Pharmaceuticals Inc.

9 10 **Reference**

- 11 1. Zanger, U.M., and Schwab, M. (2013). Cytochrome P450 enzymes in drug metabolism:
12 regulation of gene expression, enzyme activities, and impact of genetic variation.
13 *Pharmacol. Ther.* *138*, 103–141. <https://doi.org/10.1016/j.pharmthera.2012.12.007>.
- 14 2. Whirl-Carrillo, M., Huddart, R., Gong, L., Sangkuhl, K., Thorn, C.F., Whaley, R., and Klein,
15 T.E. (2021). An Evidence-Based Framework for Evaluating Pharmacogenomics Knowledge
16 for Personalized Medicine. *Clin. Pharmacol. Ther.* *110*, 563–572.
17 <https://doi.org/10.1002/cpt.2350>.
- 18 3. Turner, A.J., Nofziger, C., Ramey, B.E., Ly, R.C., Bousman, C.A., Agúndez, J.A.G., Sangkuhl,
19 K., Whirl-Carrillo, M., Vanoni, S., Dunnenberger, H.M., et al. (2023). PharmVar Tutorial on
20 CYP2D6 Structural Variation Testing and Recommendations on Reporting. *Clin.*
21 *Pharmacol. Ther.* *114*, 1220–1237. <https://doi.org/10.1002/cpt.3044>.
- 22 4. Charnaud, S., Munro, J.E., Semene, L., Mazhari, R., Brewster, J., Bourke, C., Ruybal-
23 Pesántez, S., James, R., Lautu-Gumal, D., Karunajeewa, H., et al. (2022). PacBio long-read
24 amplicon sequencing enables scalable high-resolution population allele typing of the
25 complex CYP2D6 locus. *Commun. Biol.* *5*, 1–10. [https://doi.org/10.1038/s42003-022-](https://doi.org/10.1038/s42003-022-03102-8)
26 [03102-8](https://doi.org/10.1038/s42003-022-03102-8).
- 27 5. Turner, A.J., Derezinski, A.D., Gaedigk, A., Berres, M.E., Gregornik, D.B., Brown, K.,
28 Broeckel, U., and Scharer, G. (2023). Characterization of complex structural variation in
29 the CYP2D6-CYP2D7-CYP2D8 gene loci using single-molecule long-read sequencing.
30 *Front. Pharmacol.* *14*. <https://doi.org/10.3389/fphar.2023.1195778>
- 31 6. Del Tredici, A.L., Malhotra, A., Dedek, M., Espin, F., Roach, D., Zhu, G., Volland, J., and
32 Moreno, T.A. (2018). Frequency of CYP2D6 Alleles Including Structural Variants in the
33 United States. *Front. Pharmacol.* *9*. <https://doi.org/10.3389/fphar.2018.00305>
- 34 7. Paradkar, M.U., Shah, S.A.V., Dherai, A.J., Shetty, D., and Ashavaid, T.F. (2018).
35 Distribution of CYP2D6 genotypes in the Indian population – preliminary report. *Drug*
36 *Metab. Pers. Ther.* *33*, 141–151. <https://doi.org/10.1515/dmpt-2018-0011>.
- 37 8. Koopmans, A.B., Braakman, M.H., Vinkers, D.J., Hoek, H.W., and van Harten, P.N. (2021).
38 Meta-analysis of probability estimates of worldwide variation of CYP2D6 and CYP2C19.
39 *Transl. Psychiatry* *11*, 1–16. <https://doi.org/10.1038/s41398-020-01129-1>.
- 40 9. Wang, W.Y., Twesigomwe, D., Nofziger, C., Turner, A.J., Helmecke, L.-S., Broeckel, U.,
41 Derezinski, A.D., Hazelhurst, S., and Gaedigk, A. (2022). Characterization of Novel

- 1 CYP2D6 Alleles across Sub-Saharan African Populations. *J. Pers. Med.* *12*, 1575.
2 <https://doi.org/10.3390/jpm12101575>.
- 3 10. Maulana, Y., Jimenez, R.T., Twesigomwe, D., Sani, L., Irwanto, A., Bertin, N., and
4 Gonzalez-Porta, M. (2024). The variation landscape of CYP2D6 in a multi-ethnic Asian
5 population. Preprint at bioRxiv, <https://doi.org/10.1101/2024.01.20.576401>
6 <https://doi.org/10.1101/2024.01.20.576401>.
- 7 11. Caudle, K.E., Sangkuhl, K., Whirl-Carrillo, M., Swen, J.J., Haidar, C.E., Klein, T.E., Gammal,
8 R.S., Relling, M.V., Scott, S.A., Hertz, D.L., et al. (2020). Standardizing CYP2D6 Genotype
9 to Phenotype Translation: Consensus Recommendations from the Clinical
10 Pharmacogenetics Implementation Consortium and Dutch Pharmacogenetics Working
11 Group. *Clin. Transl. Sci.* *13*, 116–124. <https://doi.org/10.1111/cts.12692>.
- 12 12. Swen, J.J., Wouden, C.H. van der, Manson, L.E., Abdullah-Koolmees, H., Blagec, K.,
13 Blagus, T., Böhringer, S., Cambon-Thomsen, A., Cecchin, E., Cheung, K.-C., et al. (2023). A
14 12-gene pharmacogenetic panel to prevent adverse drug reactions: an open-label,
15 multicentre, controlled, cluster-randomised crossover implementation study. *The Lancet*
16 *401*, 347–356. [https://doi.org/10.1016/S0140-6736\(22\)01841-4](https://doi.org/10.1016/S0140-6736(22)01841-4).
- 17 13. Relling, M.V., and Klein, T.E. (2011). CPIC: Clinical Pharmacogenetics Implementation
18 Consortium of the Pharmacogenomics Research Network. *Clin. Pharmacol. Ther.* *89*,
19 464–467. <https://doi.org/10.1038/clpt.2010.279>.
- 20 14. Brown, J.T., Bishop, J.R., Sangkuhl, K., Nurmi, E.L., Mueller, D.J., Dinh, J.C., Gaedigk, A.,
21 Klein, T.E., Caudle, K.E., McCracken, J.T., et al. (2019). Clinical Pharmacogenetics
22 Implementation Consortium Guideline for Cytochrome P450 (CYP)2D6 Genotype and
23 Atomoxetine Therapy. *Clin. Pharmacol. Ther.* *106*, 94–102.
24 <https://doi.org/10.1002/cpt.1409>.
- 25 15. Bell, G.C., Caudle, K.E., Whirl-Carrillo, M., Gordon, R.J., Hikino, K., Prows, C.A., Gaedigk,
26 A., Agundez, J., Sadhasivam, S., Klein, T.E., et al. (2017). Clinical Pharmacogenetics
27 Implementation Consortium (CPIC) guideline for CYP2D6 genotype and use of
28 ondansetron and tropisetron. *Clin. Pharmacol. Ther.* *102*, 213–218.
29 <https://doi.org/10.1002/cpt.598>.
- 30 16. Goetz, M.P., Sangkuhl, K., Guchelaar, H.-J., Schwab, M., Province, M., Whirl-Carrillo, M.,
31 Symmans, W.F., McLeod, H.L., Ratain, M.J., Zembutsu, H., et al. (2018). Clinical
32 Pharmacogenetics Implementation Consortium (CPIC) Guideline for CYP2D6 and
33 Tamoxifen Therapy. *Clin. Pharmacol. Ther.* *103*, 770–777.
34 <https://doi.org/10.1002/cpt.1007>.
- 35 17. Hicks, J.K., Sangkuhl, K., Swen, J.J., Ellingrod, V.L., Müller, D.J., Shimoda, K., Bishop, J.R.,
36 Kharasch, E.D., Skaar, T.C., Gaedigk, A., et al. (2017). Clinical pharmacogenetics
37 implementation consortium guideline (CPIC) for CYP2D6 and CYP2C19 genotypes and
38 dosing of tricyclic antidepressants: 2016 update. *Clin. Pharmacol. Ther.* *102*, 37–44.
39 <https://doi.org/10.1002/cpt.597>.

- 1 18. Crews, K.R., Monte, A.A., Huddart, R., Caudle, K.E., Kharasch, E.D., Gaedigk, A.,
2 Dunnenberger, H.M., Leeder, J.S., Callaghan, J.T., Samer, C.F., et al. (2021). Clinical
3 Pharmacogenetics Implementation Consortium Guideline for CYP2D6, OPRM1, and
4 COMT Genotypes and Select Opioid Therapy. *Clin. Pharmacol. Ther.* *110*, 888–896.
5 <https://doi.org/10.1002/cpt.2149>.
- 6 19. McInnes, G., Lavertu, A., Sangkuhl, K., Klein, T.E., Whirl-Carrillo, M., and Altman, R.B.
7 (2021). Pharmacogenetics at Scale: An Analysis of the UK Biobank. *Clin. Pharmacol. Ther.*
8 *109*, 1528–1537. <https://doi.org/10.1002/cpt.2122>.
- 9 20. Li, B., Sangkuhl, K., Whaley, R., Woon, M., Keat, K., Whirl-Carrillo, M., Ritchie, M.D., and
10 Klein, T.E. (2023). Frequencies of pharmacogenomic alleles across biogeographic groups
11 in a large-scale biobank. *Am. J. Hum. Genet.* *110*, 1628–1647.
12 <https://doi.org/10.1016/j.ajhg.2023.09.001>.
- 13 21. Chen, X., Shen, F., Gonzaludo, N., Malhotra, A., Rogert, C., Taft, R.J., Bentley, D.R., and
14 Eberle, M.A. (2021). Cyrius: accurate CYP2D6 genotyping using whole-genome
15 sequencing data. *Pharmacogenomics J.* *21*, 251–261. [https://doi.org/10.1038/s41397-](https://doi.org/10.1038/s41397-020-00205-5)
16 [020-00205-5](https://doi.org/10.1038/s41397-020-00205-5).
- 17 22. Li, S., Carss, K.J., Halldorsson, B.V., Cortes, A., and Consortium, U.B.W.-G.S. (2023).
18 Whole-genome sequencing of half-a-million UK Biobank participants. Preprint at
19 medRxiv, <https://doi.org/10.1101/2023.12.06.23299426>
20 <https://doi.org/10.1101/2023.12.06.23299426>.
- 21 23. Halldorsson, B.V., Eggertsson, H.P., Moore, K.H.S., Hauswedell, H., Eiriksson, O.,
22 Ulfarsson, M.O., Palsson, G., Hardarson, M.T., Oddsson, A., Jensson, B.O., et al. (2022).
23 The sequences of 150,119 genomes in the UK Biobank. *Nature* *607*, 732–740.
24 <https://doi.org/10.1038/s41586-022-04965-x>.
- 25 24. Cingolani, P., Platts, A., Wang, L.L., Coon, M., Nguyen, T., Wang, L., Land, S.J., Lu, X., and
26 Ruden, D.M. (2012). A program for annotating and predicting the effects of single
27 nucleotide polymorphisms, SnpEff. *Fly (Austin)* *6*, 80–92.
28 <https://doi.org/10.4161/fly.19695>.
- 29 25. Zerbino, D.R., Achuthan, P., Akanni, W., Amode, M.R., Barrell, D., Bhai, J., Billis, K.,
30 Cummins, C., Gall, A., Girón, C.G., et al. (2018). Ensembl 2018. *Nucleic Acids Res.* *46*,
31 D754–D761. <https://doi.org/10.1093/nar/gkx1098>.
- 32 26. Pedersen, B.S., and Quinlan, A.R. (2017). Who’s Who? Detecting and Resolving Sample
33 Anomalies in Human DNA Sequencing Studies with Peddy. *Am. J. Hum. Genet.* *100*, 406–
34 413. <https://doi.org/10.1016/j.ajhg.2017.01.017>.
- 35 27. Auton, A., Abecasis, G.R., Altshuler, D.M., Durbin, R.M., Abecasis, G.R., Bentley, D.R.,
36 Chakravarti, A., Clark, A.G., Donnelly, P., Eichler, E.E., et al. (2015). A global reference for
37 human genetic variation. *Nature* *526*, 68–74. <https://doi.org/10.1038/nature15393>.
- 38 28. Sudmant, P.H., Rausch, T., Gardner, E.J., Handsaker, R.E., Abyzov, A., Huddleston, J.,
39 Zhang, Y., Ye, K., Jun, G., Hsi-Yang Fritz, M., et al. (2015). An integrated map of structural

- 1 variation in 2,504 human genomes. *Nature* 526, 75–81.
2 <https://doi.org/10.1038/nature15394>.
- 3 29. PharmGKB. CYP2D6 Frequency Table.
4 <https://www.pharmgkb.org/page/cyp2d6RefMaterials>.
- 5 30. PharmGKB: CYP2D6 Allele Functionality Table
6 <https://www.pharmgkb.org/page/cyp2d6RefMaterials>.
- 7 31. Relling, M.V., Klein, T.E., Gammal, R.S., Whirl-Carrillo, M., Hoffman, J.M., and Caudle, K.E.
8 (2020). The Clinical Pharmacogenetics Implementation Consortium: 10 Years Later. *Clin.*
9 *Pharmacol. Ther.* 107, 171–175. <https://doi.org/10.1002/cpt.1651>.
- 10 32. Wang, Q., Dhindsa, R.S., Carss, K., Harper, A.R., Nag, A., Tachmazidou, I., Vitsios, D.,
11 Deevi, S.V.V., Mackay, A., Muthas, D., et al. (2021). Rare variant contribution to human
12 disease in 281,104 UK Biobank exomes. *Nature* 597, 527–532.
13 <https://doi.org/10.1038/s41586-021-03855-y>.
- 14 33. Karczewski, K.J., Francioli, L.C., Tiao, G., Cummings, B.B., Alföldi, J., Wang, Q., Collins,
15 R.L., Laricchia, K.M., Ganna, A., Birnbaum, D.P., et al. (2020). The mutational constraint
16 spectrum quantified from variation in 141,456 humans. *Nature* 581, 434–443.
17 <https://doi.org/10.1038/s41586-020-2308-7>.
- 18 34. Dhindsa, R.S., Burren, O.S., Sun, B.B., Prins, B.P., Matelska, D., Wheeler, E., Mitchell, J.,
19 Oerton, E., Hristova, V.A., Smith, K.R., et al. (2023). Rare variant associations with plasma
20 protein levels in the UK Biobank. *Nature* 622, 339–347. [https://doi.org/10.1038/s41586-](https://doi.org/10.1038/s41586-023-06547-x)
21 [023-06547-x](https://doi.org/10.1038/s41586-023-06547-x).
- 22 35. Kane, M. (2021). CYP2D6 Overview: Allele and Phenotype Frequencies. In *Medical*
23 *Genetics Summaries* [Internet] (National Center for Biotechnology Information (US)).
- 24 36. Gustafson, J.A., Gibson, S.B., Damaraju, N., Zalusky, M.P., Hoekzema, K., Twesigomwe, D.,
25 Yang, L., Snead, A.A., Richmond, P.A., Coster, W.D., et al. (2024). Nanopore sequencing of
26 1000 Genomes Project samples to build a comprehensive catalog of human genetic
27 variation. Preprint at medRxiv, <https://doi.org/10.1101/2024.03.05.24303792>
28 <https://doi.org/10.1101/2024.03.05.24303792>.
- 29 37. Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L.T., Sharp, K., Motyer, A., Vukcevic,
30 D., Delaneau, O., O'Connell, J., et al. (2018). The UK Biobank resource with deep
31 phenotyping and genomic data. *Nature* 562, 203–209. [https://doi.org/10.1038/s41586-](https://doi.org/10.1038/s41586-018-0579-z)
32 [018-0579-z](https://doi.org/10.1038/s41586-018-0579-z).
- 33 38. Sun, B.B., Chiou, J., Traylor, M., Benner, C., Hsu, Y.-H., Richardson, T.G., Surendran, P.,
34 Mahajan, A., Robins, C., Vasquez-Grinnell, S.G., et al. (2023). Plasma proteomic
35 associations with genetics and health in the UK Biobank. *Nature* 622, 329–338.
36 <https://doi.org/10.1038/s41586-023-06592-6>.

- 1 39. Haiat, S., Billard, C., Quiney, C., Ajchenbaum-Cymbalista, F., and Kolb, J.-P. (2006). Role of
2 BAFF and APRIL in human B-cell chronic lymphocytic leukaemia. *Immunology* *118*, 281–
3 292. <https://doi.org/10.1111/j.1365-2567.2006.02377.x>.
- 4 40. Smulski, C.R., and Eibel, H. (2018). BAFF and BAFF-Receptor in B Cell Selection and
5 Survival. *Front. Immunol.* *9*. <https://doi.org/10.3389/fimmu.2018.02285>.
- 6 41. Rotimi, C.N., and Jorde, L.B. (2010). Ancestry and Disease in the Age of Genomic
7 Medicine. *N. Engl. J. Med.* *363*, 1551–1558. <https://doi.org/10.1056/NEJMra0911564>.
- 8 42. Smyth, D.G. (2016). 60 YEARS OF POMC: Lipotropin and beta-endorphin: a perspective. *J.*
9 *Mol. Endocrinol.* *56*, T13–T25. <https://doi.org/10.1530/JME-16-0033>.
- 10 43. Pasternak, G.W., and Pan, Y.-X. (2013). Mu Opioids and Their Receptors: Evolution of a
11 Concept. *Pharmacol. Rev.* *65*, 1257–1317. <https://doi.org/10.1124/pr.112.007138>.
- 12 44. Vo, Q.N., Mahinthichaichan, P., Shen, J., and Ellis, C.R. (2021). How μ -opioid receptor
13 recognizes fentanyl. *Nat. Commun.* *12*, 984. [https://doi.org/10.1038/s41467-021-21262-](https://doi.org/10.1038/s41467-021-21262-9)
14 [9](https://doi.org/10.1038/s41467-021-21262-9).
- 15 45. Armstrong, S.C., and Cozza, K.L. (2003). Pharmacokinetic Drug Interactions of Morphine,
16 Codeine, and Their Derivatives: Theory and Clinical Reality, Part II. *Psychosomatics* *44*,
17 515–520. <https://doi.org/10.1176/appi.psy.44.6.515>.
- 18 46. Dubey, A.K., Handu, S.S., Dubey, S., Sharma, P., Sharma, K.K., and Ahmed, Q.M. (2011).
19 Belimumab: First targeted biological treatment for systemic lupus erythematosus. *J.*
20 *Pharmacol. Pharmacother.* *2*, 317–319. <https://doi.org/10.4103/0976-500X.85930>.
- 21 47. Taliun, D., Harris, D.N., Kessler, M.D., Carlson, J., Szpiech, Z.A., Torres, R., Taliun, S.A.G.,
22 Corvelo, A., Gogarten, S.M., Kang, H.M., et al. (2021). Sequencing of 53,831 diverse
23 genomes from the NHLBI TOPMed Program. *Nature* *590*, 290–299.
24 <https://doi.org/10.1038/s41586-021-03205-y>.

25

Figure 1. Overview of *CYP2D6* star allele calling workflow and sample filtrations in the UKBB

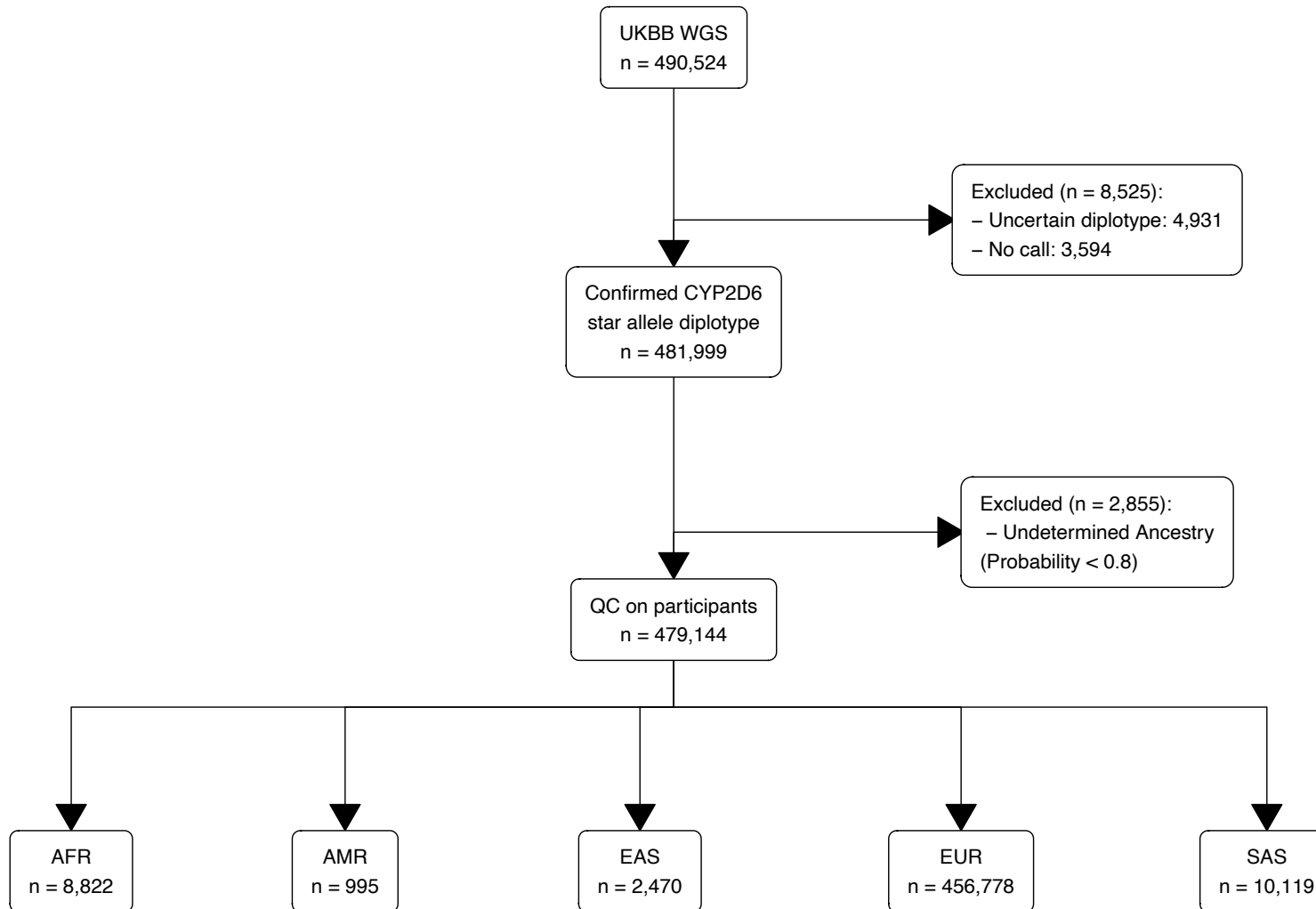
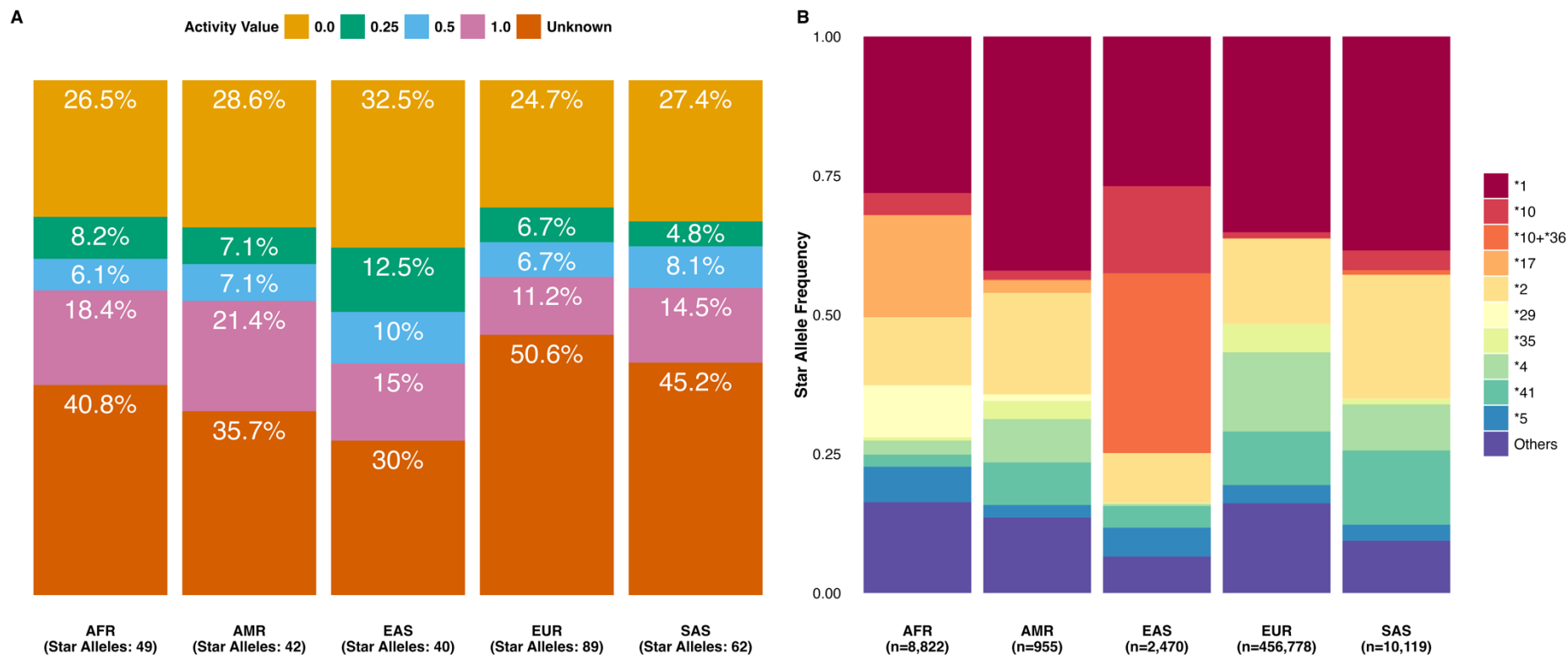
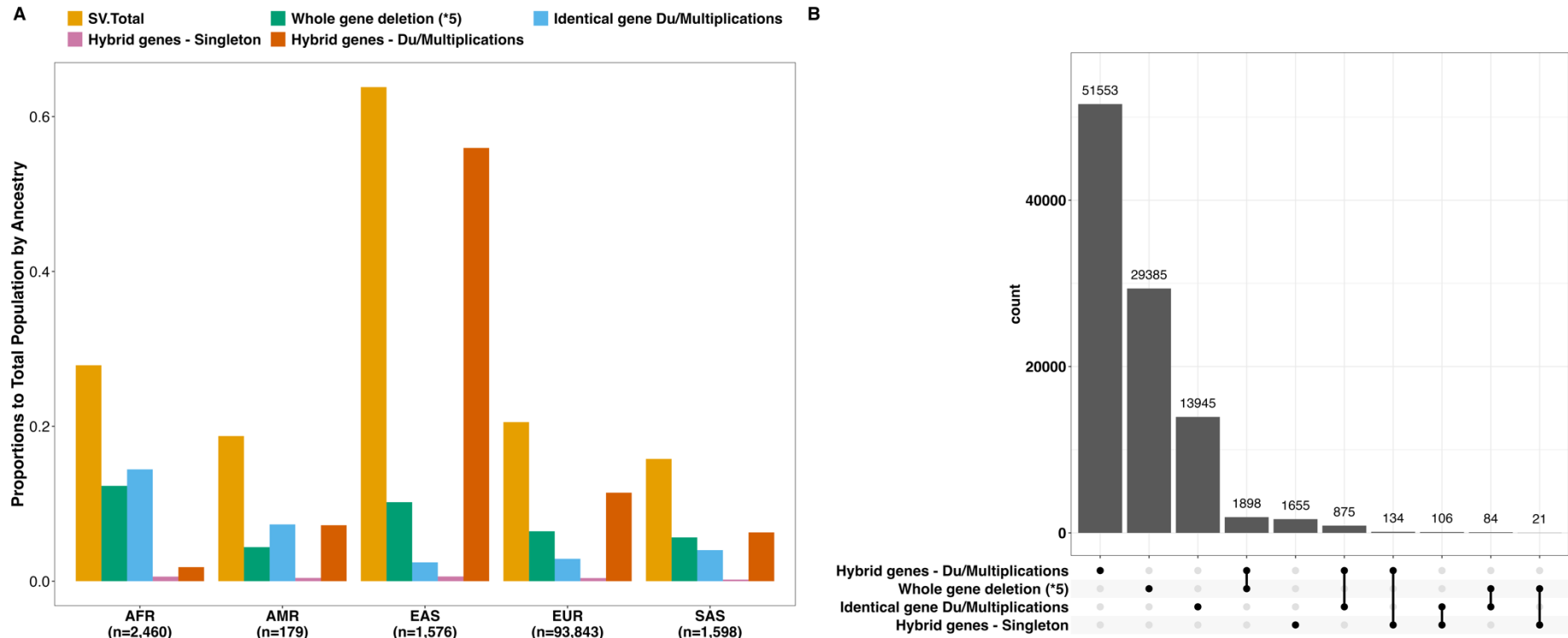


Figure 2. Overview of distributions of *CYP2D6* star allele called in the UKBB



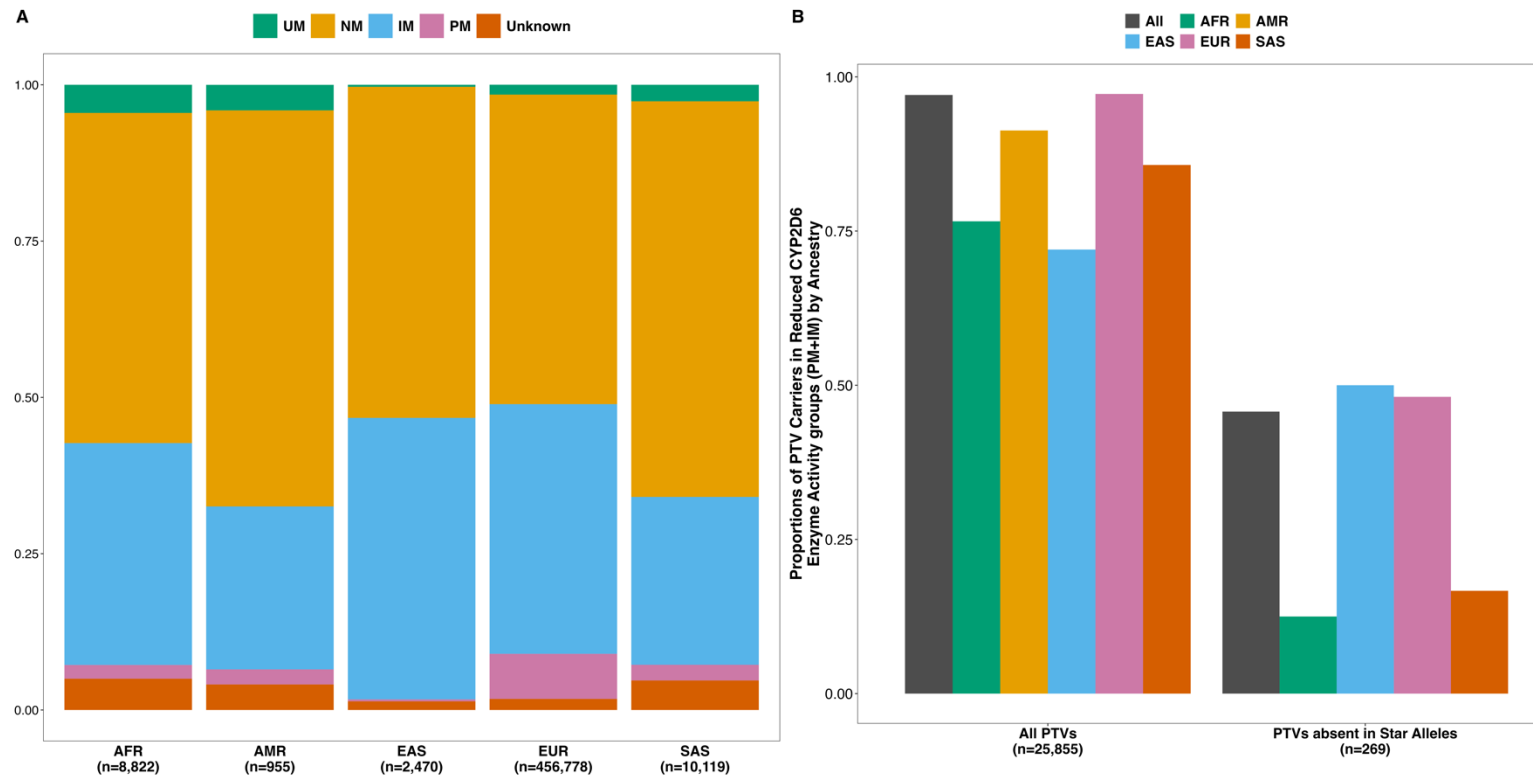
A) Distribution of functionality values of the identified *CYP2D6* star alleles in each genetically predicted ancestral group. The x-axis includes the count of distinct star alleles identified within each ancestral group. Each bar was segmented and coloured according to the distribution of star alleles based on their respective functionality values. **B) Distribution of the most prevalent (top 5) star alleles of each ancestry and their frequencies in other ancestries.** The x-axis provides the sample size for each ancestral group. Ten star alleles are highlighted, of which, two were SVs (*5, *10+*36). There were three normal-function alleles (activity value=1; *1, *2, *35), five decreased-function alleles (activity value= 0.25 or 0.5; *10, *10+*36, *17, *29, *41), and two non-functional alleles (activity value=0; *4, *5).

Figure 3. Overview of distributions of *CYP2D6* SVs in the UKBB



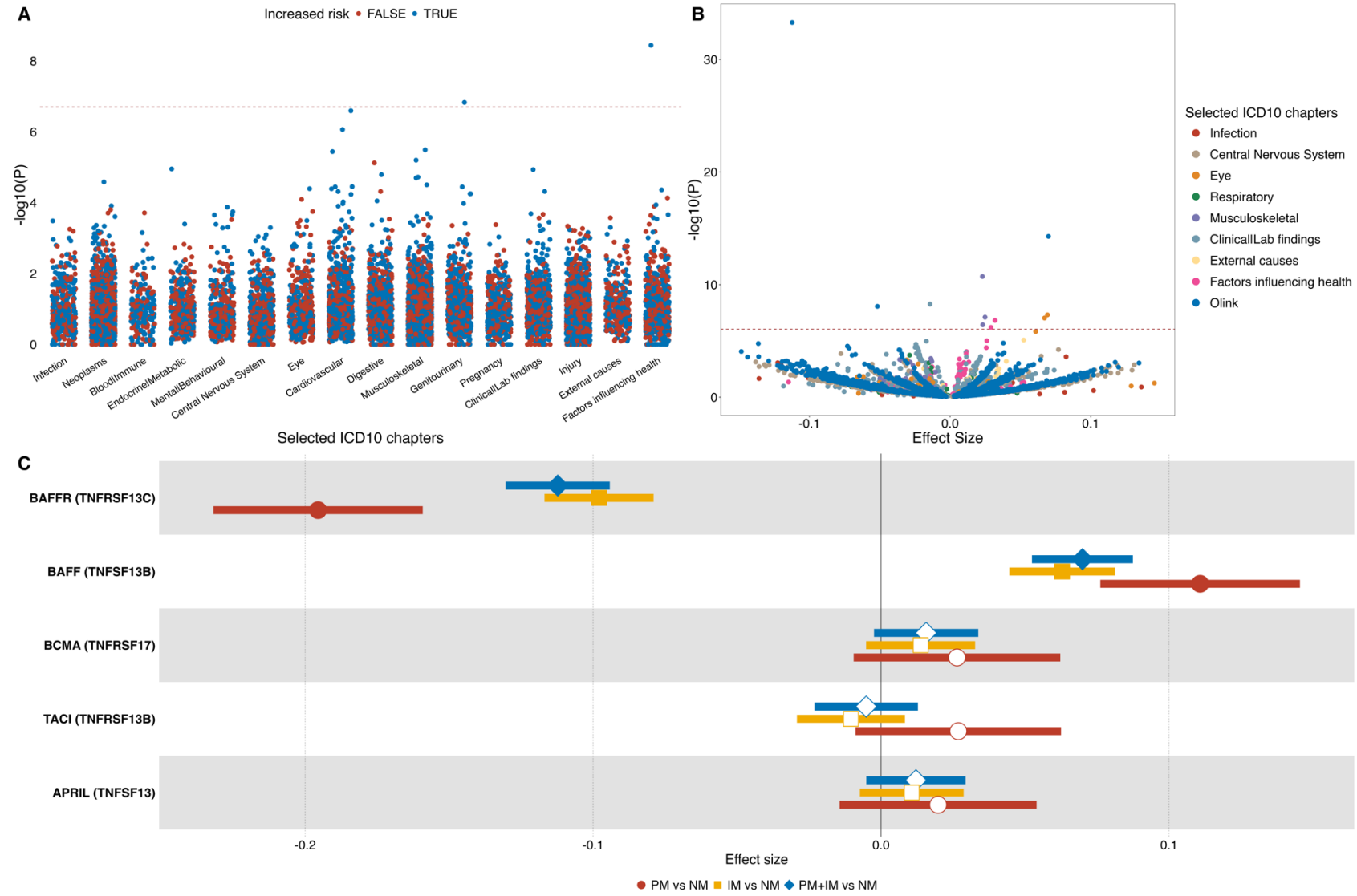
A) Distributions of *CYP2D6* SVs and the subtypes across genetic ancestries. The x-axis displays the count of *CYP2D6* SV carriers within each ancestral group. The proportions on the y-axis were determined by dividing the carrier count for each SV type by the total sample size of the respective ancestral population in the UKBB. **B) Distribution of carriers of each *CYP2D6* SV type and intersection of subtypes.** Over half of *CYP2D6* SV carriers (54.6%) were carrying the non-identical gene duplications or higher-order multiplications. Notably, there were 2,003 individuals carrying both the whole gene deletion (*5) and another *CYP2D6* SV type.

Figure 4. Overview of the distribution of *CYP2D6* gPhenos in the UKBB, and the comparison with CGR QV carriers in selected collapsing models



A) Distributions of *CYP2D6* gPhenos for each ancestral group in the UKBB. The bars were coloured according to the proportion of each *CYP2D6* gPheno carriers in each ancestry. Sizes of the ancestral population were included in the x-axis. **B) Proportions of PTV carriers in the reduced *CYP2D6* function gPheno categories (PM or IM), coloured by genetic ancestral groups.** The x-axis displayed the count of carriers of the overall *CYP2D6* PTVs (MAF≤5%) and those absent in the current *CYP2D6* star allele evaluation, respectively. The proportion on the y-axis was determined by dividing the number of PTV carriers in the decreased *CYP2D6* function gPheno categories (PM or IM) by the total number of identified carriers in the corresponding group.

Figure 5. Associations of *CYP2D6* metabolisers with clinical outcomes, biomarkers, and plasma protein abundance in the UKBB



A) Manhattan plot illustrating the phenome-wide association study of binary clinical outcomes with the defined *CYP2D6* gPheno category contrasts. Each point on the plot represents a binary clinical outcome test along with the most significantly associated *CYP2D6* gPheno comparison pair. If the estimated odds ratio exceeds one (indicating individuals in the case group of *CYP2D6* gPheno(s) had higher odds of the risk), the association is marked as 'increased risk=TRUE'. The selection of the ICD10 category for presentation was based on having at least one association with a significance level below 5% (i.e., $P < 0.05$). The red dotted horizontal line reflects the minus log₁₀ transformed significance threshold after the correction for multiple testing ($P=2 \times 10^{-7}$). **B) Volcano plot illustrating the phenome-wide association studies of quantitative traits with the defined *CYP2D6* gPheno comparison pairs.** Every point depicted on the plot symbolises the evaluation of quantitative traits along with the most significantly associated *CYP2D6* gPheno comparison pair. The ICD10 chapter was chosen for presentation if it had at least one association with a significance level below 5% (i.e., $P < 0.05$). The red dotted horizontal line reflects the minus log₁₀ transformed significance threshold after the correction for multiple testing ($P=2 \times 10^{-7}$). **C) Forest plot illustrating the associations between selected *CYP2D6* gPheno comparison pairs and TNF family members relevant to B lymphocyte development.** The plasma protein abundances underwent a rank-based inverse normalisation transformation. The effect size estimate displayed on the x-axis indicates the change per standard deviation of the plasma protein abundance when comparing between the control and case *CYP2D6* gPheno groups. Dots were filled if the association remained significant after the correction for multiple testing ($P=2 \times 10^{-7}$).

Table 1. Overview of *CYP2D6* genetically predicted metaboliser phenotypes

	AFR	AMR	EAS	EUR	SAS
PM	194 (2.2%)	23 (2.4%)	8 (0.3%)	32,913 (7.2%)	254 (2.5%)
IM	3,131 (35.5%)	249 (26.1%)	1,112 (45%)	182,434 (39.9%)	2,717 (26.9%)
NM	4,659 (52.8%)	605 (63.4%)	1,308 (53%)	226,143 (49.5%)	6,403 (63.3%)
UM	397 (4.5%)	39 (4.1%)	8 (0.3%)	7,222 (1.6%)	267 (2.6%)
Unknown	441 (5%)	39 (4.1%)	34 (1.4%)	8,066 (1.8%)	478 (4.7%)

PM – poor metaboliser; IM – intermediate metaboliser; NM – normal metaboliser; UM – ultra-rapid metaboliser; Unknown – individuals carrying star allele with currently unknown-effect. Number for each cell represents the count of individuals who were predicted to belong to the corresponding *CYP2D6* metaboliser category. The percentage in the bracket represents the proportion of the ancestral population in each *CYP2D6* metaboliser category.