

LETTER TO THE EDITOR

The *CYP2D6* VCF Translator

The Pharmacogenomics Journal advance online publication,
15 March 2016; doi:10.1038/tpj.2016.14

The cytochrome P450-2D6 (*CYP2D6*) hepatic oxidase is directly involved in the metabolism of ~25% of all commonly used drugs,^{1–3} making it one of the most well-studied genes in the field of pharmacogenetics. As is likely well known by the readership of *The Pharmacogenomics Journal*, the *CYP2D6* gene is highly polymorphic, with over 100 variant star (*) alleles currently cataloged by the Human Cytochrome P450 (CYP) Allele Nomenclature Committee.⁴ This committee was formed in 1999 to manage *CYP450* allele nomenclature using established sequence variation guidelines and to provide a summary of curated star (*) alleles and their effects when known;⁵ however, the field of pharmacogenetics has been transformed by the recent advances in high-throughput sequencing technologies. Although next-generation sequencing has greatly facilitated pharmacogenomic discovery,^{6–8} one of the resultant challenges is now reconciling the historical star (*) allele nomenclature with current genome reference assemblies and a more informed understanding of the extent of variation in the human genome.⁹

Interrogating *CYP2D6* by short-read sequencing technologies is challenging due to pseudogene sequence homology and the existence of copy number variant *CYP2D6* alleles.¹⁰ In addition, one of the most cumbersome informatics issues is converting identified *CYP2D6* sequence variants from the genome assembly used for sequence analysis (for example, GRCh37/hg19) to the M33388.1 or AY545216.1 GenBank reference sequences used to define *CYP2D6* star (*) alleles (<http://www.cypalleles.ki.se/cyp2d6.htm>). This tedious but necessary process typically is accomplished by manual curation, but recently assisted by the *CYP2D6* haplotype tables available at PharmGKB (<https://www.pharmgkb.org>).

In an effort to simplify and automate this important step, we developed a *CYP2D6* variant call format (VCF) Translator tool that takes a standard VCF file (.vcf) limited to the *CYP2D6* coordinates from human reference hg19, including flanking sequences (chr22:42522071–42528563), and converts identified sequence variant coordinates to the M33388.1 and the AY545216.1 reference coordinates used to define *CYP2D6* star (*) alleles. Importantly, the reference coordinates in the *CYP2D6* VCF Translator have already been modified to include the necessary nucleotide corrections indicated by the Nomenclature Committee. Moreover, the *CYP2D6* VCF Translator corrects for the fact that the hg19 reference sequence actually contains *CYP2D6**2 and other common variants (for example, intron 1 conversion, 1661G>C, 2850C>T, 4180G>C and so on), so these nucleotides are automatically interrogated as if being compared with wild-type *CYP2D6**1. In addition, it converts variants to their reverse complement (as *CYP2D6* is encoded on the negative strand), annotates variants with gene location (for example, exon and intron), reference/alternate alleles, dbSNP identifier (when available) and genotype calls based on VCF metrics (for example, wild-type, heterozygous and mutant), which can all be downloaded as a tab-delimited text (.txt.) file after executing the program (Figure 1).

The script was written in Python and implemented on a web server with PHP, and was designed to focus on the *CYP2D6* variants cataloged on the *CYP2D6* Nomenclature Committee website to facilitate subsequent star (*) allele conversion. To accomplish this, the *CYP2D6* VCF Translator extracts variant coordinates from an uploaded VCF and compares them to an annotated reference file of *CYP2D6* sequence variants with hg19, M33388.1 and AY545216.1 coordinates (ATG start codon = nucleotide 1). The *CYP2D6**2 correction is accomplished by reversing the *CYP2D6* nucleotides that are incorrectly variant in hg19 before making genotype calls from the VCF. Accordingly, any *CYP2D6**2 variants not called in a VCF because they are homozygous reference in hg19 are annotated in the output file as 'Mutant or Not Sequenced' (Figure 1). The *CYP2D6* VCF Translator also annotates novel sequence variants not listed on the Nomenclature Committee website.

The *CYP2D6* VCF Translator is run on the high-performance computer cluster at the Icahn School of Medicine at Mount Sinai (Minerva) and is freely available for investigators to use at <http://stuartscottlab.org/vcf>. Future iterations of the *CYP2D6* VCF Translator aim to infer star (*) allele diplotypes based on genotype data; however, in the interim we have found it to have substantial utility when converting high-throughput targeted *CYP2D6* sequencing data to star (*) alleles, as well as translating *CYP2D6* coordinates extracted from exome and whole-genome sequencing VCFs that were derived using hg19. As such, this letter is submitted to *The Pharmacogenomics Journal* to increase visibility of this online tool and, therefore, facilitate more simplified *CYP2D6* sequencing analyses for those investigators studying this very important pharmacogenetic gene.

CONFLICT OF INTEREST

The authors declare no conflict of interest.

ACKNOWLEDGMENTS

This research was supported in part by the National Institute of General Medical Sciences (NIGMS) of the National Institutes of Health (NIH) through grant K23 GM104401 (SAS), and the computational resources and staff expertise provided by the Department of Scientific Computing at the Icahn School of Medicine at Mount Sinai.

W Qiao¹, J Wang^{1,2}, BS Pullman³, R Chen^{1,2}, Y Yang¹ and SA Scott¹
¹Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York, NY, USA;
²Icahn Institute for Genomics and Multiscale Biology, Icahn School of Medicine at Mount Sinai, New York, NY, USA and
³Department of Computer Science and Engineering, University of California San Diego, La Jolla, CA, USA
E-mail: stuart.scott@mssm.edu

REFERENCES

- Gonzalez FJ, Skoda RC, Kimura S, Umeno M, Zanger UM, Nebert DW *et al.* Characterization of the common genetic defect in humans deficient in debrisoquine metabolism. *Nature* 1988; **331**: 442–446.
- Gough AC, Miles JS, Spurr NK, Moss JE, Gaedigk A, Eichelbaum M *et al.* Identification of the primary gene defect at the cytochrome P450 CYP2D locus. *Nature* 1990; **347**: 773–776.

a

stuartscottlab.org/vcf

b

INPUT (.vcf)

| #CHROM | POS | ID | REF | ALT | QUAL | FILTER | INFO | FORMAT | SAMPLE |
|--------|----------|-------------|-----|-----|---------|--------|---------------------------------|--------|------------------------------|
| chr22 | 42522613 | rs1135840 | G | C | 767.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:25,26:51:99:796,0,677 |
| chr22 | 42523003 | rs116917064 | A | G | 93.03 | PASS | AC=2;AF=1;AN=2;DIGT:AD:DP:GQ:PL | | 1/1:0,4:4:12:121,12,0 |
| chr22 | 42523209 | rs28371730 | T | C | 2062.77 | PASS | AC=2;AF=1;AN=2;DIGT:AD:DP:GQ:PL | | 1/1:0,52:52:99:2091,156,0 |
| chr22 | 42523211 | rs2004511 | T | C | 967.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:20,33:53:99:996,0,477 |
| chr22 | 42523409 | rs1985942 | G | T | 779.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:109,36:145:99:808,0,3304 |
| chr22 | 42523528 | rs1058172 | C | T | 945.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:74,40:114:99:974,0,2004 |
| chr22 | 42523636 | rs3915951 | C | A | 173.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:67,13:80:99:202,0,2092 |
| chr22 | 42523943 | rs16947 | A | G | 1777.77 | PASS | AC=2;AF=1;AN=2;DIGT:AD:DP:GQ:PL | | 1/1:0,58:58:99:1806,173,0 |
| chr22 | 42524243 | rs35742686 | CT | C | 824.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:14,24:38:99:853,0,444 |
| chr22 | 42524696 | rs58440431 | T | C | 883.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:29,36:65:99:912,0,756 |
| chr22 | 42524708 | rs111564371 | T | C | 502.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:50,17:67:99:531,0,2002 |
| chr22 | 42524713 | rs112568578 | C | G | 566.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:46,19:65:99:595,0,1874 |
| chr22 | 42524743 | rs113889384 | G | A | 546.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:53,23:76:99:575,0,1440 |
| chr22 | 42524795 | A | G | A | 468.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:40,19:59:99:497,0,1123 |
| chr22 | 42524947 | rs3892097 | C | T | 347.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:24,15:39:99:376,0,743 |
| chr22 | 42525132 | rs1058164 | G | C | 671.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:22,24:46:99:700,0,640 |
| chr22 | 42525798 | rs28371705 | G | C | 53.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:17,5:22:82:82,0,489 |
| chr22 | 42525811 | rs28371704 | T | C | 121.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:15,4:19:99:150,0,859 |
| chr22 | 42525821 | rs28371703 | G | T | 121.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:14,6:20:99:150,0,854 |
| chr22 | 42525952 | rs1328650 | C | A | 192.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:6,10:16:99:221,0,163 |
| chr22 | 42526484 | rs28371699 | A | C | 727.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:47,33:80:99:756,0,1399 |
| chr22 | 42526549 | rs56011157 | C | T | 2715.77 | PASS | AC=2;AF=1;AN=2;DIGT:AD:DP:GQ:PL | | 1/1:0,81:81:99:2744,243,0 |
| chr22 | 42526561 | rs28695233 | G | T | 3412.77 | PASS | AC=2;AF=1;AN=2;DIGT:AD:DP:GQ:PL | | 1/1:0,75:75:99:3441,232,0 |
| chr22 | 42526562 | rs75276289 | G | C | 3412.77 | PASS | AC=2;AF=1;AN=2;DIGT:AD:DP:GQ:PL | | 1/1:0,74:74:99:3441,232,0 |
| chr22 | 42526567 | rs76312385 | G | A | 3256.77 | PASS | AC=2;AF=1;AN=2;DIGT:AD:DP:GQ:PL | | 1/1:0,71:71:99:3285,220,0 |
| chr22 | 42526571 | rs74644586 | C | G | 3166.77 | PASS | AC=2;AF=1;AN=2;DIGT:AD:DP:GQ:PL | | 1/1:0,69:69:99:3195,214,0 |
| chr22 | 42526573 | rs1080996 | T | G | 3179.76 | PASS | AC=2;AF=1;AN=2;DIGT:AD:DP:GQ:PL | | 1/1:0,68:68:99:3208,217,0 |
| chr22 | 42526580 | rs1080995 | G | C | 2988.77 | PASS | AC=2;AF=1;AN=2;DIGT:AD:DP:GQ:PL | | 1/1:0,66:66:99:3017,205,0 |
| chr22 | 42526694 | rs1065852 | G | A | 1249.77 | PASS | AC=1;AF=0.5;AN=2;GT:AD:DP:GQ:PL | | 0/1:25,43:68:99:1278,0,612 |

hg19 coordinates

OUTPUT (.txt)

| hg19_POS | REF | ALT | 2d6_POS(M33388 ATG=1) | Alleles | Gene Location | 2D6_POS(AY545216 ATG=1) | 2D6_POS(AY545216) | dbSNP ID | Genotype | |
|----------|-----|-----|-----------------------|-----------|---------------|-------------------------|-------------------|-------------|------------|----------------------|
| 42526694 | G | A | 100 | C>T | Exon1 | 100 | 4300 | rs1065852 | HET | |
| 42526580 | G | C | 214 | G>C | Intron1 | 214 | 4414 | rs1080995 | WT | |
| 42526573 | T | G | 221 | C>A | Intron1 | 221 | 4421 | rs1080996 | WT | |
| 42526571 | C | G | 223 | C>G | Intron1 | 223 | 4423 | rs74644586 | WT | |
| 42526567 | G | A | 227 | T>C | Intron1 | 227 | 4427 | rs76312385 | WT | |
| 42526562 | C | A | 232 | G>C | Intron1 | 232 | 4432 | rs75276289 | WT | |
| 42526561 | G | T | 233 | A>C | Intron1 | 233 | 4433 | rs28695233 | WT | |
| 42526549 | C | T | 245 | A>G | Intron1 | 245 | 4445 | rs56011157 | WT | |
| 42526484 | A | C | 310 | G>T | Intron1 | 310 | 4510 | rs28371699 | HET | |
| 42525952 | C | A | 843 | T>G | Intron1 | 842 | 3048 | rs28371702 | HET | |
| 42525821 | G | T | 974 | C>A | Exon2 | 973 | 5173 | rs28371703 | HET | |
| 42525811 | T | C | 984 | A>G | Exon2 | 983 | 5183 | rs28371704 | HET | |
| 42525798 | G | C | 997 | C>G | Exon2 | 996 | 5196 | rs28371705 | HET | |
| 42525132 | C | G | 1661 | G>C | Exon3 | 1662 | 3862 | rs1058164 | HET | |
| 42524947 | T | C | 1846 | G>A | Intron3 | 1847 | 6047 | rs3892097 | HET | |
| 42524795 | A | G | 1998 | | Exon4 | 1999 | 6199 | | HET | |
| 42524743 | G | A | 2050 | | Intron4 | 2051 | 6251 | rs113889384 | HET | |
| 42524713 | C | G | 2080 | | Intron4 | 2081 | 6281 | rs112568578 | HET | |
| 42524708 | T | C | 2085 | | Intron4 | 2086 | 6286 | rs111564371 | HET | |
| 42524696 | T | C | 2097 | A>G | Intron4 | 2098 | 6298 | rs58440431 | HET | |
| 42524243 | CT | C | 2550 | 2549A>del | Exon5 | 2550 | 6750 | rs35742686 | HET | |
| 42523943 | A | G | 2850 | C>T | Exon6 | 2851 | 7051 | rs16947 | WT | |
| 42523636 | C | A | 3157 | | Exon7 | 3158 | 7338 | rs3915951 | WT | |
| 42523528 | C | T | 3265 | | Exon7 | 3266 | 7466 | rs1058172 | HET | |
| 42523409 | G | T | 3384 | A>C | Intron7 | 3385 | 7585 | rs1985942 | HET | |
| 42523211 | T | C | 3582 | A>G | Intron7 | 3583 | 7783 | rs2004511 | HET | |
| 42523209 | T | C | 3584 | G>A | Intron7 | 3585 | 7785 | rs28371730 | WT | |
| 42523003 | A | G | 3790 | C>T | Intron7 | 3791 | 7991 | rs116917064 | WT | |
| 42522613 | G | C | 4180 | G>C | Exon9 | 4181 | 8381 | rs1135840 | HET | |
| No Call | | | -1584 | C>G | Upstream | | -1581 | 2617 | rs1080995 | MUT or Not Sequenced |
| No Call | | | 746 | C>G | Intron1 | | 745 | 4945 | rs28371701 | MUT or Not Sequenced |
| No Call | | | 4722 | T>G | Downstream | | 4723 | 8923 | rs35028622 | MUT or Not Sequenced |

M33388.1 coordinates Gene location Genotype calls

Figure 1. An overview of the *CYP2D6* variant call format (VCF) Translator with (a) a screenshot of the online homepage (<http://stuartscottlab.org/vcf>) and (b) example input (.vcf) and output (.txt) files. Highlighted in red boxes are some of the key components of the input and output files of an example VCF using DNA from the NA12878 HapMap cell line with a *CYP2D6**3/*4 diplotype sequenced through chr22:42522044–42527019 (hg19), or – 225 to 4752 using M33388.1 coordinates. Note the variants in the output file annotated as ‘Mutant or Not Sequenced’: – 1584C > G was not called in the VCF as it is outside of the sequenced region, and the common intronic 746C > G and downstream 4722T > G variants were not called in the VCF as they are found on both the *3 and *4 haplotypes, as well as in hg19. As such, these three variants were classified as ‘Mutant or Not Sequenced’ by the *CYP2D6* VCF Translator (blue asterisk). Similarly, the original VCF shows homozygous variants for the *CYP2D6* intron 1 conversion, 2850C > T, 3584G > A and 3790C > T; however, this is actually a reflection of those variants being present in hg19 and not in NA12878. As such, these variants were corrected to wild-type (WT) by the *CYP2D6* VCF Translator (green asterisks).

- 3 Owen RP, Sangkuhl K, Klein TE, Altman RB. Cytochrome P450 2D6. *Pharmacogenet Genomics* 2009; **19**: 559–562.
- 4 Sim SC, Ingelman-Sundberg M. The Human Cytochrome P450 (CYP) Allele Nomenclature website: a peer-reviewed database of CYP variants and their associated effects. *Hum Genomics* 2010; **4**: 278–281.
- 5 Nebert DW. Suggestions for the nomenclature of human alleles: relevance to ecogenetics, pharmacogenetics and molecular epidemiology. *Pharmacogenetics* 2000; **10**: 279–290.
- 6 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**: 100–104.
- 7 Daneshjou R, Gamazon ER, Burkley B, Cavallari LH, Johnson JA, Klein TE *et al.* Genetic variant in folate homeostasis is associated with lower warfarin dose in African Americans. *Blood* 2014; **124**: 2298–2305.
- 8 Gordon AS, Tabor HK, Johnson AD, Snively BM, Assimes TL, Auer PL *et al.* Quantifying rare, deleterious variation in 12 human cytochrome P450 drug-metabolism genes in a large-scale exome dataset. *Hum Mol Genet* 2014; **23**: 1957–1963.
- 9 Robarge JD, Li L, Desta Z, Nguyen A, Flockhart DA. The star-allele nomenclature: retooling for translational genomics. *Clin Pharmacol Ther* 2007; **82**: 244–248.
- 10 Gaedigk A. Complexities of CYP2D6 gene analysis and interpretation. *Int Rev Psychiatry* 2013; **25**: 534–553.