1

2

3

4

5

6

# A variant by any name: quantifying annotation discordance across tools and clinical databases

9

10

11

12

Jennifer Yen[1], jennifer.yen@personalis.com

Sarah Garcia[1,2], stkgarcia@gmail.com

Aldrin Montana[1], aldrin.montana@personalis.com

Jason Harris[1], jason.harris@personalis.com

Steven Chervitz[1], schervitz@personalis.com

John West[1], john.west@personalis.com

Richard Chen[1], richard.chen@personalis.com

Deanna M. Church[1,2] , deanna.church@gmail.com

21

22

[1]Personalis, 1330 O'Brien Drive, Menlo Park, CA 94025

[2]10X Genomics, 7068 Koll Center Pkwy #401, Pleasanton, CA 94566

26

27 **ABSTRACT**

28 **Background**

29 Clinical genomic testing is dependent on the robust identification and reporting of

30 variant-level information in relation to disease. With the shift to high-throughput

31 sequencing, a major challenge for clinical diagnostics is the cross-identification of

32 variants called on their genomic position to resources that rely on transcript- or protein-

33 based descriptions.

34

35 **Methods**

36 We evaluated the accuracy of three tools (SnpEff, Variant Effect Predictor and Variation

37 Reporter) that generate transcript and protein-based variant nomenclature from genomic

38 coordinates according to guidelines by the Human Genome Variation Society (HGVS).

39 Our evaluation was based on comparisons to a manually curated list of 127 test variants

40 of various types drawn from data sources, each with HGVS-compliant transcript and

41 protein descriptors. We further evaluated the concordance between annotations

42 generated by Snpeff and Variant Effect Predictor with those in major germline and

43 cancer databases: ClinVar and COSMIC, respectively.

44

45 **Results**

46 We find that there is substantial discordance between the annotation tools and

47 databases in the description of insertion and/or deletions. Accuracy based on our ground

48 truth set was between 80-90% for coding and 50-70% for protein variants, numbers that

49    are not adequate for clinical reporting. Exact concordance for SNV syntax was over

50    99.5% between ClinVar and Variant Effect Predictor (VEP) and SnpEff, but less than

51    90% for non-SNV variants. For COSMIC, exact concordance for coding and protein

52    SNVs were between 65 and 88%, and less than 15% for insertions. Across the tools and

53    datasets, there was a wide range of equivalent expressions describing protein variants.

54

55    **Conclusion**

56    Our results reveal significant inconsistency in variant representation across tools and

57    databases. These results highlight the urgent need for the adoption and adherence to

58    uniform standards in variant annotation, with consistent reporting on the genomic

59    reference, to enable accurate and efficient data-driven clinical care.

60

61    **KEYWORDS**

62    HGVS, clinical genetic testing, genomics, annotation, sequencing, syntax, precision

63    medicine, variant

64

65    **INTRODUCTION**

66    High-throughput sequencing has transformed the landscape of clinical genetic testing.

67    This strategy, combined with the completion of massive public profiling datasets (ExAc

68    [1], 1000 Genomes [2]), has dramatically changed our approach towards cancer

69    treatment and the diagnosis of inherited disease. A major challenge in the analysis of

70    this throughput and volume of data is integrating variant level information from the

3

71    wealth of clinical and biological insight accumulated over decades of research,

72    particularly those from recent, large sequencing studies. Describing a variant's location

73    is a fundamental part of a clinical assessment, yet the practice remains difficult,

74    inconsistent and evolving.

75

76    Specifically, the clinical genomics community faces an enormous hurdle, which is

77    integrating data generated prior to the availability of a robust human reference assembly

78    with that generated using modern methods. Standards and guidelines for describing

79    variants at the genomic, transcript (coding) and protein level, provided by the Human

80    Genome Variation Society (HGVS) [3], were developed when testing was largely

81    transcript rather than genome-based. As laboratories shifted to high-throughput

82    sequencing, variant analysis transitioned to the genome level, confounding comparisons

83    with reports generated from previous transcript-based assays.

84

85    Reconciling variant coordinates from the transcript to the genome, and vice versa, is not

86    an unambiguous task. Requisite information about the genomic and transcript sequence

87    accessions, their versions, and the alignments used to relate the two sequences, are not

88    always reported in publications (Figure 1a-b). Alignment of cDNA to the genome remains

89    challenging and can result in substantially different exon structures depending on the

90    alignment approach (Figure 1a) [4,5]. In addition, variant reporting standards for VCF, a

91    format designed to store genomic variation, are different from those for HGVS, a format

92    that describes transcript and protein variants. In the context of nucleotide repeats, VCF

93    shifts left with respect to the genome, while HGVS shifts right with respect to the gene or

94    transcript (Figure 1c). Variants can therefore have very different locations depending on

95    their accession, version and alignment.

96

97    Even in relation to the same transcript, a variant can have multiple representations.

98    HGVS expressions can have long and short forms, preferred and non-preferred syntax,

99    and describe amino acids by their triple (e.g. Glu) or a single letter designation (e.g. E)

100    (Figure 1d). In a survey by Deans et al. (2016) [6], 20 laboratories reported the HGVS

101    syntax for a single variant in 14 different ways.  An evaluation of over 140 molecular

102    pathology laboratories in Europe and the UK revealed substantial errors in reported

103    HGVS variant descriptions for the EGFR gene [7].

104

105    Currently there are many tools that automatically generate HGVS syntax, including

106    SnpEff [8], Variant Effect Predictor (VEP) [9], Annovar [10], Variation Reporter (VR) [11],

107    Mutalyzer [12] and packages developed by individual clinical laboratories such as Invitae

108    [5] and Counsyl [13].  While the performance of different genomic variant callers have

109    been well-studied [14,15], the accuracy and consistency of HGVS generation tools have

110    not yet been described.

111

112    Previous comparison of Annovar and VEP revealed substantial differences in annotation

113    based on choice of transcript [16]. This low concordance, combined with the increasing

114    demand for automated syntax generation, prompted our re-evaluation of the

115    performance of well-supported, open source tools. We considered only freely available

116    tools as they would have the largest reach. Additionally, we wished to focus on

117    annotation differences that can occur even when the same transcript is used. In this

118    paper, we compare the concordance of variant nomenclature generated by VEP [9],

119    SnpEff [8] and Variation Reporter, benchmarked by a curated 'truth' set and variant

120    annotations described in large public datasets for germline (ClinVar) and cancer

121    (COSMIC) variant descriptions. We find that while the tools SnpEff and Variant Effect

122    Predictor produce comparable results, there remains significant discordance in variant

123    annotation among tools, public resources, and literature.

124

125    **METHODS**

126    **Datasets**

127    We curated a test set of 127 variants to establish a ground-truth set with which we can

128    evaluate the accuracy of the tools. Fifty-one variants were selected from public

129    repositories: ClinVar, dbSNP, COSMIC, My Cancer Genome, Emory database and

130    Leiden (Additional file 1: Table S1). We added 76 synthetic variants to ensure

131    representation across variant types and genomic features. Genomic, coding and protein

132    nomenclature for all variants were generated using a combination of the Mutalyzer [17]

133    and Variation Viewer [18] webservice. Effect impact was determined based on the

134    protein syntax and sequence ontology (SO) [19].

135

136    We used the ClinVar GRCh37 VCF and annotations from the tab separated file

137    downloaded from the FTP site [20] (January 5[th] 2016 release). We used the rsid and

138    alternative allele to connect variants between the two files. We obtained genomic

139    coordinates from the COSMIC GRCh37 VCF and connected them with the transcript and

140     nomenclature in the CosmicCompleteExport.tsv file from the COSMIC website [21] (v75)

141     with the COSMID.

142

143     **VCF normalization**

144     We used vt-normalize [22] to left-justify all variants in each of the dataset VCFs used. A

145     breakdown of insertions and deletions (indels) for each dataset and the number

146     normalized is represented in Additional file 1: Table S2.

147

148     **Tools used**

149     We ran SnpEff (v4.1L) [8,23], VEP (v82) [9] and VR [11] on our ground truth set, and

150     subsequently only SnpEff and VEP on the ClinVar and COSMIC datasets. The Snpeff

151     database was built using the NCBI GRCh37 GFF corresponding to the NCBI annotation

152     '*Homo sapiens* 105'. The Snpeff database for Ensembl transcripts was built using the

153     GRCh37 Ensembl transcript GFF [24]. We ran VEP with the corresponding RefSeq or

154     Ensembl cache (v83). For all tools we used NCBI GRCh37p13 as the input reference

155     genome.

156

157     **Assessment of syntax**

158     To assess the performance of the variant annotation tools, we performed string match

159     comparisons between the output and the reference syntax (Figure 2). Annotations were

160     evaluated according to the HGVS guidelines [25] [26]. Variant annotations were labeled

161     as 'exact' matches when the HGVS string and the query annotation matched as-is. If the

7

162    string did not match perfectly, but could be transformed to the query string by applying

163    HGVS recommendations, the tool's annotation was labeled 'equivalent'. For this study,

164    both 'exact' and 'equivalent' annotations are regarded as correct. The code and module

165    for performing the syntax assessment will be distributed on GitHub.

166

167

168    **RESULTS**

169    **Comparisons to a ground truth test set**

170    In order to assess the performance of different variant annotation tools against a ground

171    truth, we used a contrived test set of 127 manually curated variants (Figure 3a)

172    comprised of 52 previously reported variants in the literature or databases, and an

173    additional 76 synthetic variants targeting a spectrum of variants (Additional file 1:

174    Table S3). All annotations were reviewed manually using a combination of the

175    Mutalyzer and Variation Viewer web services. This structured test set would allow us

176    to deeply evaluate variants across different classes, effects, and genomic features.

177

178    Using the analysis flowchart summarized in Figure 4, we compared the annotations

179    generated by VR, VEP [9] and SnpEff [8] to the ground-truth test set (Additional File 1:

180    Table S4). VEP and SnpEff accept VCF as input files; at the time of analysis, the VR

181    API was limited in its functionality in processing large VCF files. Input HGVS

182    expressions were also required for Mutalyzer, but we did not assess this tool because

183    it was used to construct the ground truth set. We compared only annotations made on

184    the same RefSeq transcript version. Although the input transcript alignments for

185　　SnpEff and VEP were identical, the tools produced a different number of transcripts

186　　and annotations. For example, we could not extract the relevant transcript for 4

187　　variants in the SnpEff output and 5 from the VEP output, in addition to 5 variants

188　　absent from both tools. The importance of transcript collection was more pronounced

189　　for VR, which uses its own in-house alignments. As a result, 18% of the test variants

190　　could not be assessed by VR because NCBI carries only the most up-to-date

191　　transcripts. VR also frequently yielded multiple annotations for a single variant and

192　　transcript. In these cases, we chose the first variant in the output to evaluate in this test.

193　　In total, only 121 out of the 127 variants were annotated on the relevant transcript for

194　　any of the three tools.

195

196　　A major challenge in comparing nomenclature between tools was evaluating the

197　　equivalency of the many HGVS expressions for a given variant. Protein variant syntax

198　　was far more variable than coding variant syntax: between 14-20% of protein

199　　annotations were described with equivalent nomenclature across the three tools,

200　　compared to only 2.5% to 3.0% of coding syntax (Figure 4c-d). Each tool had distinct

201　　frameshift and synonymous annotations; frameshift has both long and short form

202　　alternatives, while synonymous variants can be described in several ways; e.g. 'p.(=)',

203　　'p.=', 'p.Thre258=', 'p.Thr258Thr' (PTV012, Additional file 1: Table S5). Exact

204　　concordance in annotation between SnpEff and VEP was higher at the coding level

205　　(77.6% of variants) than at the protein level (68.8.1%) (Figure 4b, right panel). Less

206　　agreement was observed between Variation Reporter and either VEP or SnpEff:

207　　approximately 50% with either tool for coding and protein syntax.

208

9

209    For variants in our ground-truth set, SnpEff and VEP exhibited comparable accuracy

210    and precision. At the coding level, SnpEff, VEP and VR annotated between 80% and

211    85% of substitutions correctly (out of 65 and 68), compared to 100% of substitutions

212    for VR (out of 55). For deletions and insertions, VR performed poorly largely due to

213    systematic errors in reporting. VR incorrectly described all but two deletions as indels.

214    The remaining two annotations diverged from HGVS guidelines by omitting the 'del'

215    designation altogether (e.g. c.2199-1301GA>A) (PTV062, PTV067, Additional file 1:

216    Table S5). Duplications were also annotated as indels, but with technically equivalent

217    (and redundant) nomenclature (c.1961dupG as c.1960delCinsCG). Such VR errors at

218    the coding level led to inaccurate protein syntax for 18 variants.

219

220    We tested the ability of the tools to discriminate between the genomic reference and

221    RefSeq transcript sequences, both of which are independently curated by the NCBI

222    [27]. Since RefSeq transcripts typically receive a high level of manual review, conflicts

223    between the RefSeq and genomic sequence often reveal an error in the latter. For this

224    reason, we included nine test instances of RefSeq-Genomic differences in our ground

225    truth set. Strikingly, none of the nine test examples of RefSeq-Genomic differences

226    were identified by either VEP or SnpEff (Additional file 1: Table S5), and were

227    erroneously reported as missense or deletion variants. While VR correctly identified 7

228    out of 9 RefSeq Genomic differences (the remaining two variants were not

229    annotated), it mistakenly called differences for an additional 22 variants, indicating a

230    poor precision for recognizing true differences. HGVS expressions should always

231    reflect the base on the relevant genomic or transcript sequence to avoid asserting

232    variants at positions where there is no change.

233

234  Both SnpEff and VEP correctly annotated the phased dinucleotide substitutions,

235  which are variants present in consecutive bases, also known as multinucleotide

236  variants (MNV) (Additional file 1: Table S6). Dinucleotide substitutions are highly

237  prevalent in cancers associated with clear mutagen exposures such as melanoma, lung

238  adenoma and lung squamous cell carcinoma [28]. Similarly, treatment by the

239  chemotherapeutic agents cisplatin and meclorethamine have also been shown to cause

240  dinucleotide substitutions at appreciable rates [28]. VR incorrectly annotated the phased

241  dinucleotide substitutions as frameshift variants (PTV044, PTV068). We found that

242  MNVs must be phased in the VCF, as the tools annotated adjacent but independent

243  substitutions in the VCF separately instead of as a pair. For example, two BRAF

244  variants (PTV045, PTV046) were incorrectly annotated as p.V600E and p.V600M, when

245  the combined result would be p.V600K. These results indicate that for cancers with a

246  high mutation load, prior phasing for dinucleotide pairs will be especially crucial to

247  circumvent potential clinical oversights [29].

248

249  To complement the analysis of protein and coding annotations, we also assessed the

250  variant effects predicted by the tools. Predicted effect is commonly used for

251  evaluating pathogenicity during variant interpretation [30]. In instances where a

252  variant could be associated with two functional consequences (for example, as

253  intronic but also at a slice acceptor site), the annotation was considered to be correct

254  if one association was described. Overall, the accuracy of effect prediction correlated

255  highly with that of protein annotation (Additional file 1: Figure S1) even if they are

256  calculated independently [8]). Compared to coding and protein syntax, efforts among

257    tools to converge on a standardized set of variant effect annotations were far more

258    evident (Additional file 2: Figure S1).

259

260    **Comparison with the ClinVar dataset**

261    Having established baseline accuracy for automated syntax generation, we sought to

262    assess the syntax concordance of these tools with those in public datasets. We

263    started with ClinVar [31], a large public archive of variant and disease relationships

264    that is widely used for evaluating Mendelian disease. Of the 106,110 small variants in

265    the ClinVar VCF, the vast majority are SNVs (84%); the rest comprise a smaller

266    number of deletions (10%), duplications (3.3%), insertions (1%) and indels (1%)

267    (Figure 3b). We evaluated the performance of VEP and SnpEff on the ClinVar dataset

268    (Additional file 3); because of the limited functionality and long running time of the VR

269    tool, we did not include it in subsequent annotation assessments (Additional file 1:

270    Table S7).

271

272    Approximately 10% of transcripts in the ClinVar dataset had different versions from

273    those in our input transcript alignment file, which was used to build resources for both

274    VEP and SnpEff (Figure 5a). Approximately 1.8% of ClinVar transcript accessions

275    were not represented in the alignment input at all. Because of these discrepancies in

276    transcript accession and versions, we could not assess the SnpEff or VEP

277    annotations for 7 and 7.5% of ClinVar variants, again underscoring the importance of

278    the input transcript set.

279

12

280    Overall concordance for both SnpEff and VEP was remarkably high, which can be

281    attributed to the large proportion of SNVs (Figure 5b). At the coding level, both SnpEff

282    and VEP yielded nearly perfect concordance for SNVs, matching the exact ClinVar

283    nomenclature for over 99.9% of the SNVs (Figure 5a, Additional file 2: Figure S2b). In

284    rare instances of error, SnpEff and VEP were typically incorrect by one base

285    (Additional file 1: Table S8). Exact concordance was lowest for variants annotated as

286    insertions in ClinVar (approximately 75-80% for both tools), largely due to their correct

287    assertion as duplications by VEP and SnpEff. In contrast, concordance was slightly

288    higher for deletions and indels (between 86 and 88%). There were 25 instances in

289    which neither tool could have predicted correct coding HGVS syntax without prior

290    reports of the splicing product. A single nucleotide change at a splice site in the AGA

291    gene NC_000004.11:g.178354367C>A (NM_000027.3:c.940+1G>T) results in the

292    skipping of exon 8 and a final syntax of c.807_940del134 (Additional file 1: Table S9). In

293    five of these cases, this type of error also resulted in the incorrect protein syntax.

294

295    As with the ground truth test set, we observed greater variation in protein syntax

296    (Table 1). This was mostly evident for deletions, duplications, and insertions, where

297    between 16 to 78% of annotations were reported correctly but with alternative

298    nomenclature (Additional file 2: Figure S2b). Overall concordance was again high for

299    SNVs (99%), with 75% and 83% exact nomenclature for SnpEff and VEP. However,

300    neither VEP nor SnpEff performed as well on deletions, duplications and insertions

301    (between 76.3% and 94.4% overall concordance). For non-SNV variant types, our

302    results show that between 60-70% of annotations output by these tools do not match

303    the ClinVar HGVS, and between 5-20% of these annotations are completely

304    discordant.  Dinucleotide substitutions, which ClinVar reports as indels, were

13

305     annotated as independent substitutions for both SnpEff and VEP (Additional file 1:

306     Table S6). In a few cases at the boundaries between coding and non-coding regions,

307     VEP and ClinVar yielded no output while SnpEff reported the ambiguity as

308     'p.Thr662_Glu663delins???'. Even for substitutions, there were instances where all three

309     tools yielded distinct nomenclature for the same variant. For

310     NM_001126128.1:c.163delA, ClinVar, SnpEff and VEP output p.Ile55Terfs, p.Ile55fs,

311     and p.Ile55Ter respectively. The correct HGVS syntax for this variant is p.Ile55Ter.

312

313     Interestingly, agreement between SnpEff and VEP sometimes revealed errors or

314     inconsistencies in the ClinVar output. For example, for rs34618570, the TTN variant

315     NM_133378.4:c.10361-2293A>T is purported to be a missense variant (p.Ile3877Phe),

316     when the variant is intronic for that transcript. For rs398123611 (Additional file 1: Table

317     S8), ClinVar recognizes NM_133378.4:c.1138_1140dupGGC as a duplication in the

318     coding syntax but annotates the protein as an insertion (p.Gly380_Ala381insGly). At

319     least 626 variants in the ClinVar dataset were incorrectly annotated as both nonsense

320     and frameshift (p.Glu307Terfs), when the output should simply be nonsense. Together,

321     these results demonstrate that while there is near perfect consensus between the HGVS

322     tools and ClinVar annotations for SNVs, the uniformity and correctness for other variant

323     types, which are often of the most clinically relevant (e.g. frameshifts), still needs

324     improvement.

325

326     **Comparison with the COSMIC dataset**

327     Clinical cancer care is dependent on identifying relationships between tumor variants

328     and relevant information about their prognostic and therapeutic significance. We

14

329     investigated the consistency between annotation output by SnpEff and VEP with

330     COSMIC, currently the largest public resource of somatic mutations in human cancer

331     [32] that is also widely used by clinical laboratories. Again, we did not include VR in

332     our assessment because of its limited functionality and long running time. Because

333     COSMIC annotates variants in relation to Ensembl instead of NCBI RefSeq transcript

334     accessions, we built a second, separate database to run VEP and SnpEff according to

335     Ensembl transcript alignments.

336

337     We queried a total of 3,075,504 coding COSMIC variants. Following normalization and

338     de-duplication of the COSMIC VCF, there remained a set of 2,215,076 variants (Figure

339     3). Approximately 142,134 variants were insertions, deletions or indels, 19% of which

340     required left justification (Additional file 1: Table S2). We compared syntax

341     representations (Figure 5c, Additional file 4). Both SnpEff and VEP generated

342     annotations for approximately 90% of the COSMIC dataset. Because the cancer field

343     employs the convention of abbreviating amino acids to a single letter while the

344     annotation tools, and HGVS, all use the three-letter convention, we converted the

345     COSMIC annotations to three-letter amino acids to facilitate annotation comparison.

346

347     At the coding level, VEP recapitulated the exact syntax as COSMIC for 85.9% of the

348     total variants, compared to 76.8% of variants by SnpEff, with less than 1% of

349     equivalent syntax for both tools (Figure 5b). However, the majority of the COSMIC

350     dataset are SNVs (95%); for variant types other than SNVs, neither VEP nor SnpEff

351     achieved comparable concordance (Additional File 2: Figure S2b). Notable

352     differences in annotations include COSMIC's reporting of all duplications as

15

353    insertions, resulting in nearly complete discordance for variants of this type. We did

354    not assert the equivalency of multi-base insertions with duplications due to the

355    involvement of verifying duplicated bases in the reference transcript. As a result, none

356    of the indel annotations were exact string matches. Additionally, in complete

357    departure from current HGVS standard, COSMIC reports indels as block substitutions

358    (c.569_570TC>AT vs c.569_570delTCinsAT, Table 1), which we assessed as

359    'equivalent'.  This format could be attributed to the historical representation of

360    dinucleotide variants [33], which remains popular despite the adoption of the HGVS

361    standard by most clinical resources (e.g. My Cancer Genome). By failing to

362    consistently right justify insertion and deletion positions, the concordance between

363    tools and COSMIC nomenclature for deletions was less than 50% (Additional file 1:

364    Table S10).

365

366    For protein variants, SnpEff reproduced the exact protein syntax for 75.8% of

367    COSMIC variants compared to 58.4% by VEP (Figure 5b). A large fraction of VEP

368    discordance could be attributed to VEP's annotation of all frameshifting indels as

369    nonsense variants (Additional file 1: Table S10, COSM1476431). Further, over 90% of

370    VEP alternative protein expressions were due to discrepant reporting of synonymous

371    variants as p.= compared to p.Gly35Gly by both COSMIC and SnpEff (Table 1).

372    Similar to coding deletions, nuances in nomenclature revealed distinct expressions of

373    frameshifts for COSMIC, VEP and SnpEff.

374

375    For the majority of discordant annotations, the agreement between SnpEff and VEP

376    syntax suggest that the COSMIC syntax is incorrect. To verify the HGVS nomenclature

16

377    of these variants, we mapped the Ensembl transcript to its approximate corresponding

378    RefSeq accession through its consensus coding sequence (CCDS), since a number of

379    tools, including Mutalyzer, do not support Ensembl identifiers. A mutation in *TP53* at

380    position chr17:7578525 (COSM1683507) is annotated in COSMIC as c.404_405insC.

381    Because of a sequence of 4 C's at this position, the standardized left shifted VCF

382    position should be at chr17:7578523 and right-shifted HGVS syntax as c.405_406insC,

383    or c.405dupC. In another example, a HER2 insertion variant is described in My Cancer

384    Genome as c.2339_2340ins (with no insertion bases or transcript as reference) and

385    G778_P780dup. The correct coding syntax by both SnpEff and VEP is c.2331_2339dup

386    while the correct protein syntax (output only by VEP) is p.Gly778_Pro780dup. COSMIC

387    annotated neither the coding or protein syntax correctly (Table 1). Based on the

388    agreement of VEP and SnpEff alone, our results suggest that at least 2.7% of COSMIC

389    variant annotations are incorrect (Table 1). This is not surprising given its recent

390    transition from a research repository to a major clinical resource, although efforts to

391    comply with genomic and HGVS standards are apparently underway.

392

393    **Clinical impact of discordant variant annotation**

394    Ultimately, we are concerned about the concordance of positional and syntax

395    expressions because of its impact on clinical interpretation. To illustrate this point, we

396    describe a frameshift variant in the *PROK2* gene, which was differentially classified as

397    an exercise by two curators in our laboratory - one classifying as likely pathogenic

398    and the other as pathogenic for Kallman syndrome. The difference in classification

399    stemmed from the use of different syntax in constructing the string-based search. The

400    variant was described as 'NM_001126128.1:c.297dupT (p.Gly100Trpfs*22)'. Because

17

401    of alternative transcripts and HGVS representations, this variant could be searched by

402    multiple expressions (Additional file 1: Figure S3a). In one route, searching 'PROK2

403    c.297_298insT' or 'PROK2 c.234_235insT' immediately retrieved the relevant

404    literature to classify this variant. However, searching 'PROK2 *297_298ins*', 'PROK2

405    *234_235ins*', or the correct HGVS syntax 'c.297dup' or 'c.234dup' did not return any

406    relevant results (Additional file 1: Figures S2b). Searching for 'PROK2 G100fsX121',

407    'PROK2 c.297_298insT' or 'PROK2 c.234_235insT' identifies a paper by Abreu et al.

408    [34], which leads to a thread of reports that supports a final variant classification of

409    'pathogenic' (Additional file 1: Figure S3b-c). Because of these multiple variant

410    representations, identifying relevant information can entail navigating a complex matrix

411    of HGVS expressions and web results.

412

413    As another example of the importance of accurate HGVS nomenclature for clinical

414    care, a variant in a patient's melanoma sample was annotated in our pipeline as

415    'NM_004333.4:c.1799T>A (p.V600E)'. During visual review we found that the variant

416    was part of a dinucleotide pair, with a combined syntax of c.1799_1800delTGinsAT

417    and protein syntax of p.V600D. Although p.V600D is sensitive to BRAF inhibitors, this

418    variant is not as well-studied and characterized with respect to drug response and

419    efficacy compared p.V600E. Further, while V600E confers sensitivity to MEK

420    inhibition, the sensitivity of p.V600D to MEK remains unclear.

421

422    **DISCUSSION**

423    We have described some of the remaining challenges of moving clinical sequencing

424    into a high-throughput environment. Consistent with findings by McCarthy et al. [16],

18

425 we find that the transcript collection has a significant impact on the yield of relevant

426 variant annotations. Our examination of automated syntax from HGVS tools and the

427 ClinVar or ground truth datasets reveal that approximately 10% of variants could not

428 be assessed due to discordant transcript accessions or versions. The fact that ClinVar

429 and COSMIC, the largest public repositories of germline and somatic data

430 respectively, do not share the same collection of transcript accessions reflects the

431 degree of harmonization and the need for a universal store of transcript to genome

432 alignments.

433

434 Importantly, although variant calling is performed almost exclusively on genomic data,

435 variants are still being primarily referenced with respect to the transcript. Recent

436 publications continue to describe variants according to their protein and/or coding

437 syntax [35-37], sometimes even without the transcript identifier [38,39]. In a survey by

438 the American Society of Molecular Pathologists, 50% of clinical cancer labs report

439 variants exclusively by coding and protein HGVS nomenclature but without

440 accompanying genomic coordinates. The same survey also found that 70% of clinical

441 cancer labs use as a resource MyCancerGenome.org, which references variants by

442 their popular single-letter amino acid or coding-level convention, again, without

443 transcript or genomic coordinates. As our analyses show, transforming genomic

444 positions to transcript loci is challenging and prone to error; ambiguity in

445 representation is best avoided by always referencing variants by their genomic

446 position and assembly version. For this reason, HGVS recommends reporting clinical

447 variants by their Locus Reference Genomic sequence (LRG), a system designed for

448 clinically relevant variants that is based on un-versioned and stabled accession

449 sequences [26,40].

450

451    Despite the precision achieved with generating syntax for SNVs, the positions of

452    insertions and/or deletions remain stubbornly difficult to annotate, regardless of the

453    VCF or HGVS genomic standard. The presence of duplicates in nearly one-fifth of the

454    COSMIC VCF highlights the importance of using tools for normalization to reconcile

455    the multiple possible positions to represent a single variant. At the level of HGVS, we

456    found that none of the non-SNV variant types were annotated with near 100%

457    accuracy or compliance with HGVS conventions for any of the tools or databases that

458    we queried. Given the rigorous reporting requirements of a clinical genetics lab, this is

459    concerning, and suggests that it remains critical to manually review the syntax when

460    reporting non-SNVs.

461

462    Our analyses further provide a glimpse into the diverse matrix of possible HGVS

463    representations for a given variant - a disturbing concept for attempts to mine and

464    exploit existing resources through string-based search. Internal efforts can be made to

465    standardize HGVS syntax within knowledge-bases and clinical enterprises; variants

466    can be transformed into a standard, minimal expression to enable a uniform query

467    across curated databases [41]. However, while this is useful for a limited set of data, it

468    is impractical for mining beyond internally curated information. The alternative is

469    exhaustive but impractical, requiring the search for every permutation of an HGVS

470    expression for a particular variant. A thoughtful discussion should be made about

471    asserting HGVS guidelines as rules to enforce a strict convergence across

472    laboratories, resources, and literature.

473

474    By design, the HGVS annotation system was not intended for mining large bodies of

475    genomic information, while approximations of syntax are not acceptable because of

476    their impact on clinical care. A means of clinical intervention in oncology is to directly

477    connect clinically actionable variants in patient tumor samples with relevant therapeutic

478    strategies, such as approved drugs or eligibility for clinical trials. In the ACMG guidelines

479    for the classification of germline variants, at least five categories of evidence require

480    interrogating variants from previous reports in reliable databases or the published

481    literature [30]. Already, studies have shown that there remains substantial heterogeneity

482    in the interpretation of genomic variants by clinical laboratories [6,7,42]. Imprecise

483    nomenclature can lead to variant misclassification and consequent misdiagnosis [29].

484    The applications of genomics in clinical care will require concerted efforts to converge on

485    standardized reporting mechanisms to enable data sharing and integration across

486    diverse datasets and resources. Reporting on the same genomic reference, according to

487    uniform variant syntax, will be one crucial step towards the achieving this aim and the

488    ultimate goal of precision medicine.

489

490    **DECLARATIONS**

491    **List of abbreviations**

492    HGVS, Human Genome Variation Society; VEP, Variation Effect Predictor; VR, Variation

493    Reporter; VCF, Variant Call Format

494

495    **Competing Interests**

496    JY, SG, AM, JH, SC, JW, RC, DMC were full time employees of Personalis at the time of

497    this study.

498

**Authors' contributions**

499

500    JY designed the study, performed the bioinformatics and data analysis, interpreted

501    results and wrote the manuscript. DMC conceived the idea and guided the study design,

502    results interpretation and manuscript preparation. SG contributed to the data analysis,

503    study design and results interpretation. AM and SC contributed to the bioinformatics

504    analysis. Both AM and JH provided guidance for the bioinformatics analysis. All authors

505    read and approved the final manuscript.

506

507    **Acknowledgements**

511

512

513    **ADDITIONAL FILES**

514    **Additional file 1**

515          **Table S1.** Ground Truth Set Variants

516          **Table S2.** Variants Normalized by Dataset

517          **Table S3.** Ground Truth Set Contents by Features

518          **Table S4.** Ground Truth Set Comparison Results. Exact matches between the

519          reference annotation in COSMIC and annotations provided by Snpeff and VEP

520          are noted as "yes", equivalent matches as "yes_m" ("yes modified") and not

521          equivalent annotations as "no".

522          **Table S5.** Example Nomenclature Discrepancies Ground Truth Set Variants

523          **Table S6.** Annotation of dinucleotide substitutions

524          **Table S7.** Tools Run Time

525          **Table S8.** Example Nomenclature Discrepancies from the ClinVar Dataset

526          **Table S9.** Examples of genomic SNVs resulting in deletions at the transcript level

527          **Table S10.** Example Nomenclature Discrepancies from the COSMIC Dataset

528

529    **Additional file 2.**

530          **Figure S1. Comparison of effect annotation between tools and the HGVS**

531          **test set.**

532     a) Concordance in effect nomenclature between the HGVS test set and SnpEff

533     and VEP by variant type.

534     **Figure S2.** Concordance in variant syntax by variant type between tools and a)

535     ClinVar or b) COSMIC datasets at the coding (upper panel) and protein (lower

536     panel) level. Bars represent fraction of exact (blue) and equivalent (orange)

537     matches. All duplications were marked as insertions in COSMIC.

538     **Figure S3.  Impact of HGVS nomenclature on clinical interpretation**

539     a) Transcripts and nomenclature associated with variant (chr3:g.71821968dupA).

540     b) PubMed and Google results from search strings.

541     c) From a single search string to evidence and classification.

542

543     **Additional file 3.** ClinVar Comparison Results. Exact matches between the reference

544     annotation in COSMIC and annotations provided by Snpeff and VEP are noted as "yes",

545     equivalent matches as "yes_m" ("yes modified") and not equivalent annotations as "no".

546

547     **Additional file 4.** COSMIC Comparison Results. . Exact matches between the reference

548     annotation in COSMIC and annotations provided by Snpeff and VEP are noted as "yes",

549     equivalent matches as "yes_m" ("yes modified") and not equivalent annotations as "no".

550

551

## REFERENCES

552

553    1. Lek M. ExAC_Main_Submission_151029. 2015 Oct pp. 1–26.

554    2. Altshuler DL, Abecasis GR, Chakravarti A, La Vega De FM, Donnelly P, Gibbs RA, et
555    al. A map of human genome variation from population-scale sequencing. Nature.
556    2010;467:1061–73.

557    3. Dunnen den JT, Antonarakis SE. Mutation nomenclature extensions and suggestions
558    to describe complex mutations: a discussion. Hum. Mutat. 2000;15:7–12.

559    4. Kapustin Y, Souvorov A, Tatusova T, Lipman D. Splign: algorithms for computing
560    spliced alignments with identification of paralogs. Biol Direct. 2008;3:20–13.

561    5. Hart RK, Rico R, Hare E, Garcia J, Westbrook J, Fusaro VA. A Python package for
562    parsing, validating, mapping and formatting sequence variants using HGVS
563    nomenclature. Bioinformatics. 2015;31:268–70.

564    6. Deans Z, Fairley JA, Dunnen den JT, Clark C. HGVS Nomenclature in Practice: An
565    Example from the United Kingdom National External Quality Assessment Scheme. Hum.
566    Mutat. 2016.

567    7. Tack V, Deans ZC, Wolstenholme N, Patton S, Dequeker EMC. What's in a Name? A
568    Coordinated Approach toward the Correct Use of a Uniform Nomenclature to Improve
569    Patient Reports and Databases. Human Mutation. 2016;37:570–5.

570    8. Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, et al. A program for
571    annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs
572    in the genome of Drosophila melanogaster strain w1118; iso-2; iso-3. Fly (Austin).
573    2012;6:80–92.

574    9. McLaren W, Pritchard B, Rios D, Chen Y, Flicek P, Cunningham F. Deriving the
575    consequences of genomic variants with the Ensembl API and SNP Effect Predictor.
576    Bioinformatics. 2010;26:2069–70.

577    10. Wang K, Li M, Hakonarson H. ANNOVAR: functional annotation of genetic variants
578    from high-throughput sequencing data. Nucleic Acids Research. 2010;38:e164–4.

579    11. Variation Reporter [Internet]. http://www.ncbi.nlm.nih.gov/variation/tools/reporter.
580    [cited 2016 May 16]. Available from: http://www.hgvs.org/mutnomen/

581    12. Wildeman M, van Ophuizen E, Dunnen den JT, Taschner PEM. Improving sequence
582    variant descriptions in mutation databases and literature using the Mutalyzer sequence
583    variation nomenclature checker. Hum. Mutat. 2008;29:6–13.

584    13. Counsyl HGVS variant name parsing and generation [Internet]. [cited 2016 May 16].
585    Available from: https://github.com/counsyl/hgvs

586    14. Hwang S, Kim E, Lee I, Marcotte EM. Systematic comparison of variant calling
587    pipelines using gold standard personal exome variants. Scientific Reports. The Author(s)
588    SN ; 2015;5:17875.

589   15. Cornish A, Guda C. A Comparison of Variant Calling Pipelines Using Genome in a
590   Bottle as a Reference. BioMed Research International. 2015;2015:1–11.

591   16. McCarthy DJ, Humburg P, Kanapin A, Rivas MA, Gaulton K, Cazier J-B, et al.
592   Choice of transcripts and software has a large effect on variant annotation. Genome
593   Med. 2014;6:26.

594   17. Taschner PEM, Dunnen den JT. Describing structural changes by extending HGVS
595   sequence variation nomenclature. Lindblom A, Robinson PN, editors. Human Mutation
596   [Internet]. 2011;32:507–11. Available from: https://mutalyzer.nl/

597   18. Variation Viewer [Internet]. [cited 2016 May 16]. Available from:
598   http://www.ncbi.nlm.nih.gov/variation/view/

599   19. Eilbeck K, Lewis SE, Mungall CJ, Yandell M, Stein L, Durbin R, et al. The Sequence
600   Ontology: a tool for the unification of genome annotations. Genome Biology.
601   2005;6:R44–12.

602   20. ClinVar FTP Site [Internet]. [cited 2016 May 16]. Available from:
603   ftp://ftp.ncbi.nlm.nih.gov/pub/clinvar/

604   21. COSMIC [Internet]. [cited 2016 May 16]. Available from:
605   http://cancer.sanger.ac.uk/cosmic/download)

606   22. Tan A, Abecasis GR, Kang HM. Unified representation of genetic variants.
607   Bioinformatics. 2015;31:2202–4.

608   23. SnpEff [Internet]. [cited 2016 May 16]. Available from: http://snpeff.sourceforge.net/

609   24. Ensembl FTP Fownload [Internet]. [cited 2016 May 16]. Available from:
610   http://www.ensembl.org/info/data/ftp/index.html

611   25. HGVS Website [Internet]. [cited 2016 May 16]. Available from:
612   http://www.hgvs.org/mutnomen

613   26. Dunnen den JT, Dalgleish R, Maglott DR, Hart RK, Greenblatt MS, McGowan-
614   Jordan J, et al. HGVS Recommendations for the Description of Sequence Variants:
615   2016 Update. Human Mutation. 2016;:n/a–n/a.

616   27. Pruitt KD, Brown GR, Hiatt SM, Thibaud-Nissen F, Astashyn A, Ermolaeva O, et al.
617   RefSeq: an update on mammalian reference sequences. Nucleic Acids Research.
618   2014;42:D756–63.

619   28. Meier B, Cooke SL, Weiss J, Bailly AP, Alexandrov LB, Marshall J, et al. C. elegans
620   whole-genome sequencing reveals mutational signatures related to carcinogens and
621   DNA repair deficiency. Genome Research. Cold Spring Harbor Lab; 2014;24:1624–36.

622   29. Varga E, Chao EC, Yeager ND. The importance of proper bioinformatics analysis
623   and clinical interpretation of tumor genomic profiling: a case study of undifferentiated
624   sarcoma and a constitutional pathogenic BRCA2 mutation and an MLH1 variant of
625   uncertain significance. Familial Cancer. Springer Netherlands; 2015;:1–5.

626 30. Richards S, Aziz N, Bale S, Bick D, Das S, Gastier-Foster J, et al. Standards and
627 guidelines for the interpretation of sequence variants: a joint consensus recommendation
628 of the American College of Medical Genetics and Genomics and the Association for
629 Molecular Pathology. Genet Med. 2015;17:405–23.

630 31. Landrum MJ, Lee JM, Benson M, Brown G, Chao C, Chitipiralla S, et al. ClinVar:
631 public archive of interpretations of clinically relevant variants. Nucleic Acids Research.
632 Oxford University Press; 2016;44:D862–8.

633 32. Forbes SA, Beare D, Gunasekaran P, Leung K, Bindal N, Boutselakis H, et al.
634 COSMIC: exploring the world's knowledge of somatic mutations in human cancer.
635 Nucleic Acids Research. 2015;43:D805–11.

636 33. Pleasance ED, Cheetham RK, Stephens PJ, McBride DJ, Humphray SJ, Greenman
637 CD, et al. A comprehensive catalogue of somatic mutations from a human cancer
638 genome. Nature. Macmillan Publishers Limited. All rights reserved; 2009;463:191–6.

639 34. Abreu AP, Trarbach EB, de Castro M, Frade Costa EM, Versiani B, Matias Baptista
640 MT, et al. Loss-of-function mutations in the genes encoding prokineticin-2 or prokineticin
641 receptor-2 cause autosomal recessive Kallmann syndrome. J Clin Endocrinol Metab.
642 2008;93:4113–8.

643 35. Zhu X, Petrovski S, Xie P, Ruzzo EK, Lu Y-F, McSweeney KM, et al. Whole-exome
644 sequencing in undiagnosed genetic diseases: interpreting 119 trios. Genet Med.
645 2015;17:774–81.

646 36. Helbig KL, Farwell Hagman KD, Shinde DN, Mroske C, Powis Z, Li S, et al.
647 Diagnostic exome sequencing provides a molecular diagnosis for a significant proportion
648 of patients with epilepsy. Genet Med. 2016;:1–8.

649 37. Yang Y, Muzny DM, Reid JG, Bainbridge MN, Willis A, Ward PA, et al. Clinical
650 Whole-Exome Sequencing for the Diagnosis of Mendelian Disorders. New England
651 Journal of Medicine. 2013;369:1502–11.

652 38. Soden SE, Saunders CJ, Willig LK, Farrow EG, Smith LD, Petrikin JE, et al.
653 Effectiveness of exome and genome sequencing guided by acuity of illness for diagnosis
654 of neurodevelopmental disorders. Sci Transl Med. 2014;6:265ra168.

655 39. Pansuriya TC, van Eijk R, d'Adamo P, van Ruler MAJH, Kuijjer ML, Oosting J, et al.
656 Somatic mosaic IDH1 and IDH2 mutations are associated with enchondroma and
657 spindle cell hemangioma in Ollier disease and Maffucci syndrome. Nat Genet.
658 2011;43:1256–61.

659 40. MacArthur JAL, Morales J, Tully RE, Astashyn A, Gil L, Bruford EA, et al. Locus
660 Reference Genomic: reference sequences for the reporting of clinically relevant
661 sequence variants. Nucleic Acids Research. 2014;42:D873–8.

662 41. Patterson SE, Liu R, Statz CM, Durkin D, Lakshminarayana A, Mockus SM. The
663 clinical trial landscape in oncology and connectivity of somatic mutational profiles to
664 targeted therapies. Human Genomics. Human Genomics; 2016;:1–13.
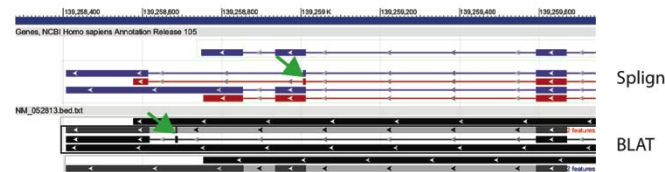
665    42. Pepin MG, Murray ML, Bailey S, Leistritz-Kessler D, Schwarze U, Byers PH. The
666    challenge of comprehensive and consistent sequence variant interpretation between
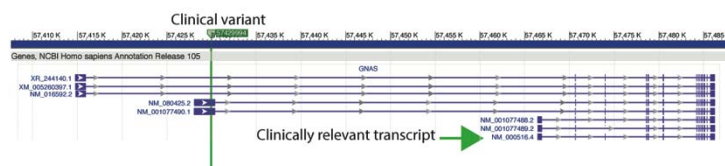667    clinical laboratories. Genet Med. 2015;18:20–4.

668

669

670    **FIGURES**



671

672    **Figure 1. Factors affecting HGVS syntax generation.**

673    a) Transcript alignment approach can impact the transcript exon structure.  Alignment of cDNA
674    sequence by Splign and BLAT to the genome results in a 10kb difference in an exon positioning
675    in the CARD9 gene (green arrow).

676    b) Transcript accession can impact the variant association and HGVS syntax. Here, the identified
677    GNAS variant is outside the clinically relevant transcript. Small changes in versions may also
678    impact the coding sequence.

679    c) In the context of nucleotide repeats, variant justification can affect the variant's position.

29

680    d) Transcript annotation directly impacts its translation to a protein expression. Incorrect transcript
681    annotation can lead to incorrect protein syntax.

682    e) Representing the variant in a particular expression.  There are different ways of expressing the
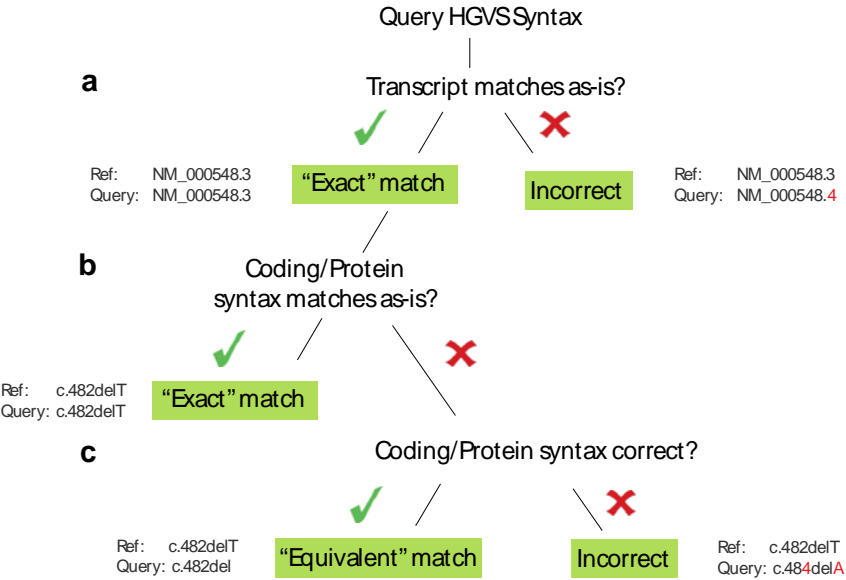683    same coding or protein variant.

684

685

686

687

**Figure 2. Methodology of HGVS syntax comparison.**

To compare two HGVS expressions in our dataset, we applied the following assessments.

a) The query transcript must match the reference transcript. If the accession or version

does not match, the variant is not assessed.

b) If the syntax for both expressions correspond as-is, the match is 'exact'.

c) If the syntax for both expressions are equivalent, the match is 'equivalent',  If the syntax is not
an alternative expression of the other HGVS variant, the match is 'incorrect'.

695

696

697

698

699

700

701

702

703

704

705



706

707

708

**Figure 3. Datasets by composition.**

a) Number of variants evaluated in the Ground Truth, ClinVar and COSMIC dataset. Note that the number of variants assessed may be less than the number of variants in the input set.

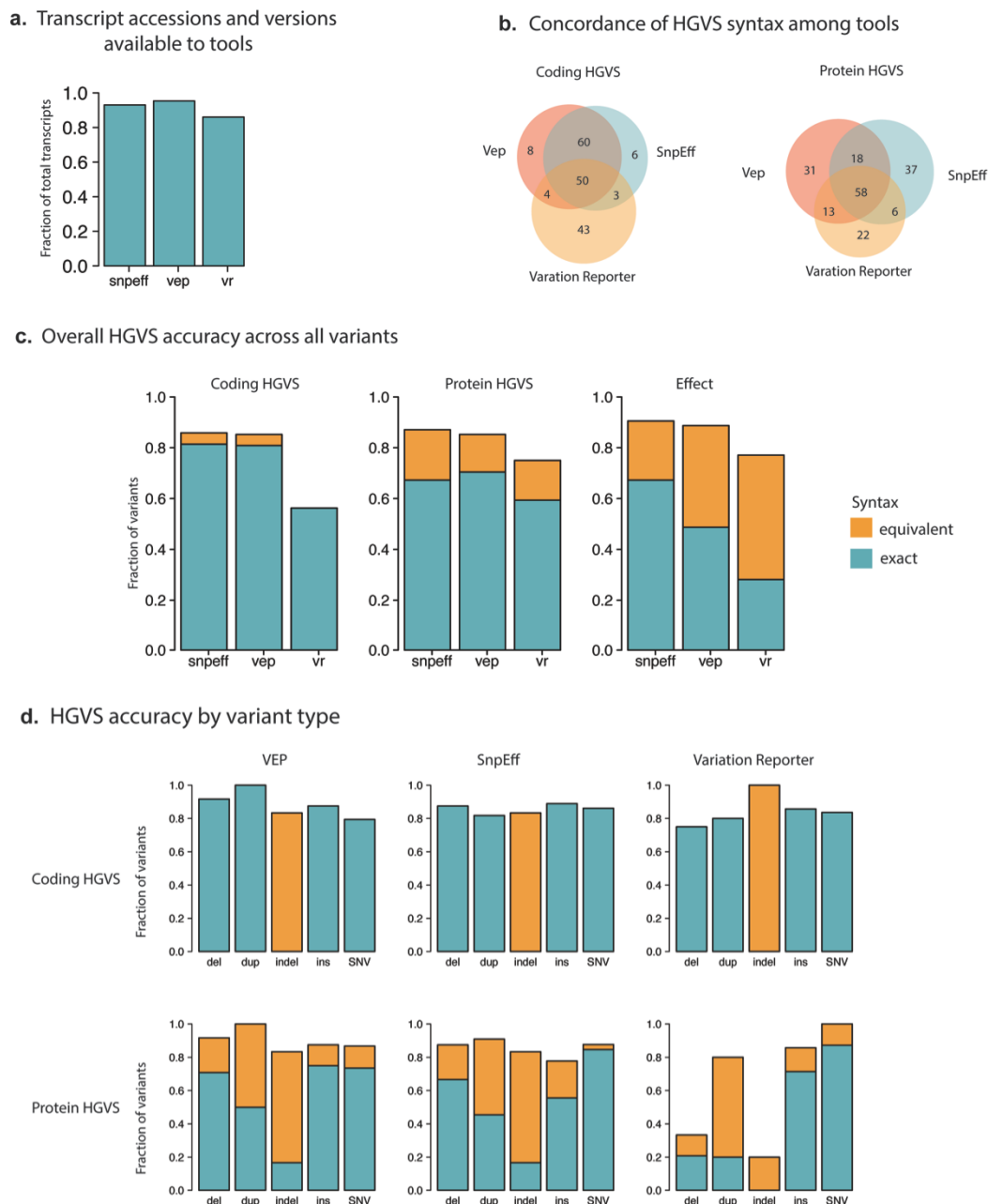b) Distribution of variant types for each dataset.  Duplications are included under insertions.

**Figure 4. Summary of ground truth set HGVS syntax assessment.**

a) Fraction of unique transcript accessions and versions in the ground truth set that were available to the tools SnpEff (snpeff), VEP (vep), and Variation Reporter (vr). If a transcript was not accessible to the tool, the variant could not be annotated with respect to that transcript.

b) Exact concordance of HGVS syntax at the coding (left) and protein (right) level among the tools.
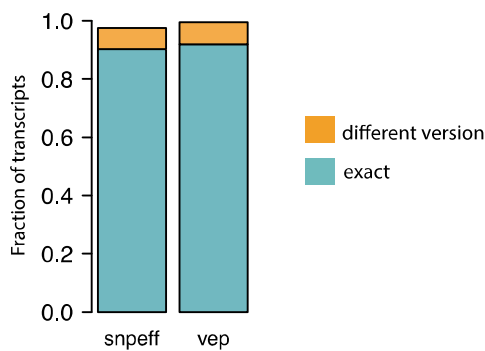
721     c) Accuracy of annotation across variants (n=121) described as exact (blue) and equivalent
722     (orange). Fraction shown is with respect to number of annotations on the relevant transcript on
723     the test set.

724     d) Accuracy of annotation for each variant type across the tools. Variant types evaluated were:
725     deletions (del), indels (delins), duplications (dup), insertions (ins) and single nucleotide variatnts
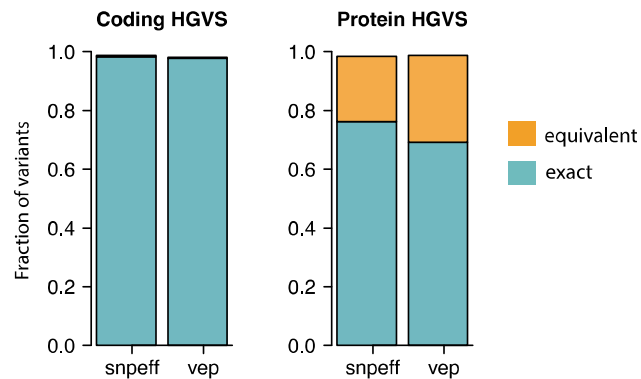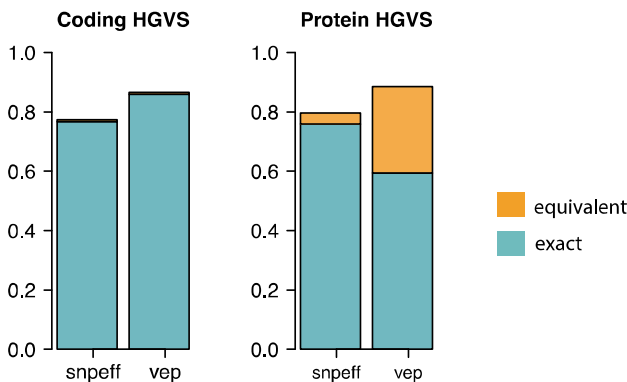726     (SNVs).

727

728



729

730    **Figure 5.  ClinVar and COSMIC HGVS syntax assessment.**

731    a) ClinVar transcript accessions and versions available to tools. Transcripts available to tools that
732    matched the ClinVar reference transcript are marked in blue; transcripts that were different
733    versions from the ClinVar transcript are marked in orange. Only unique transcripts were
734    considered.

735    b) Overall concordance in variant syntax across all variants between tools and ClinVar at the
736    coding (upper panel) and protein (lower panel) level. Bars represent fraction of exact (blue) and
737    equivalent (orange) matches.

738    c) Overall concordance in variant syntax across all variants between tools and COSMIC at the
739    coding (upper panel) and protein (lower panel) level. Bars represent fraction of exact (blue) and
740    equivalent (orange) matches. All duplications were considered insertions in COSMIC.

741

**Table 1. Exemplar variants demonstrating nomenclature discrepancies**

**Coding HGVS**

| variant type | ClinVar | COSMIC | SnpEff | Vep | VR | Reference ID |
|---|---|---|---|---|---|---|
| insertion | - | c.2262_2263ins14 | - | c.2262_2263insGGCATCTCAGCATC | - | COSM5254274 |
| duplication | - | c.422_423insA | c.428dupA | c.428dupA | - | COSM4719972 |
| duplication | c.567_568dup | - | c.567_568dupTT | c.567_568dupTT | - | rs137854332 |
| indel | - | c.3141_3142GA>TT | c.3141_3142delGAinsTT | c.3141_3142delGAinsTT | - | COSM4387531 |
| indel | c.68-5_68-3delinsTT | - | c.68-5_68-3delCTCinsTT | c.68-5_68-3delCTCinsTT | c.68-5_68-3delCTCinsTT | rs397516362 |
| deletion | c.562_563delCA | - | c.562_563delCA | c.562_563delCA | c.564delGinsCAG | PTV003 |
| insertion | | c.2339_2340insGGGCTCCCC | c.2331_2339dupGGGCTCCCC | c.2331_2339dupGGGCTCCCC | | COSM12555 * |

**Protein HGVS**

| effect | ClinVar | COSMIC | SnpEff | Vep | VR | Reference ID |
|---|---|---|---|---|---|---|
| synonymous variant | p.Arg317= | - | p.Arg317Arg | p.= | p.Arg317= | rs111033272 |
| synonymous variant | - | p.*1143* | p.Ter1143Ter | p.= | - | COSM3558732 |
| stop gained | p.Gln100Ter | - | p.Gln100* | p.Gln100Ter | - | rs119103276 |
| extension variant | - | p.*1133L | p.Ter1133Leuext*? | p.Ter1133LeuextTer22 | - | COSM1569676 |
| inframe insertion | - | - | p.Arg309_Arg310insArgArg | p.Arg310_Arg311dup | p.Arg311_Lys312insArgArg | PTV111 |
| inframe insertion | - | p.T502_His505delTTGH | p.Thr502_His505del | p.Thr502_His505del | - | COSM1163654 |
| inframe deletion | - | - | p.Ala1111_Ala1119del | p.Ala1111_Ala1119del | p.Ala1119_Gly1120insAlaAlaAlaAlaAlaAlaAlaAlaAla | PTV021 |
| inframe deletion | | p.N442delN | p.Asn442del | p.Asn442del | - | COSM5074446 |
| frameshift variant | p.Arg227Lysfs | - | p.Arg227fs | p.Arg227LysfsTer31 | - | rs80356649 |
| frameshift variant | - | p.P1176fs*>46 | p.Pro1176fs | p.Pro1176AlafsTer117 | - | COSM5196763 |
| frameshift variant | - | p.R613fs*15 | p.Arg613fs | p.Arg613AlafsTer15 | - | COSM5193613 |
| frameshift variant | - | - | p.Glu238fs | p.Glu238ProfsTer9 | p.Phe237_Glu238insPro | PTV008 |
| inframe insertion | - | p.Pro780_Tyr781insGlySerPro | p.Pro780_Tyr781insGlySerPro | p.Gly778_Pro780dup | - | COSM12555 * |

* known in My Cancer Genome as "c.2339_2340ins (G778_P780dup)"