

1 Skyhawk: An Artificial Neural Network-  
2 based discriminator for reviewing  
3 clinically significant genomic variants

4

5 **Author**

6 Ruibang Luo<sup>1,2,\*</sup>, Tak-Wah Lam<sup>1</sup>, Michael C. Schatz<sup>2</sup>

7

8 <sup>1</sup> Department of Computer Science, The University of Hong Kong, Hong Kong

9 <sup>2</sup> Department of Computer Science, Johns Hopkins University, Baltimore, MD, USA

10

11 \* Correspondence should be addressed to [rbluo@cs.hku.hk](mailto:rbluo@cs.hku.hk)

12

## 13 Abstract

14 **Motivation:** Many rare diseases and cancers are fundamentally diseases of the genome. In  
15 the past several years, genome sequencing has become one of the most important tools in  
16 clinical practice for rare disease diagnosis and targeted cancer therapy. However, variant  
17 interpretation remains the bottleneck as is not yet automated and may take a specialist several  
18 hours of work per patient. On average, one-fifth of this time is spent on visually confirming  
19 the authenticity of the candidate variants.

20 **Results:** We developed Skyhawk, an artificial neural network-based discriminator that  
21 mimics the process of expert review on clinically significant genomics variants. Skyhawk  
22 runs in less than one minute to review ten thousand variants, and about 30 minutes to review  
23 all variants in a typical whole-genome sequencing sample. Among the false positive  
24 singletons identified by GATK HaplotypeCaller, UnifiedGenotyper and 16GT in the HG005  
25 GIAB sample, 79.7% were rejected by Skyhawk. Worked on the Variants with Unknown  
26 Significance (VUS), Skyhawk marked most of the false positive variants for manual review  
27 and most of the true positive variants no need for review.

28  
29 **Availability:** Skyhawk is easy to use and freely available at  
30 <https://github.com/aquaskyline/Skyhawk>

31

## 32 Keywords

33 Variant Validation, Clinical Decision Support, Artificial Intelligence

34

## 35 Introduction

36 The dramatic reduction in the cost of whole genome, exome and amplicon sequencing has  
37 allowed these technologies to be increasingly accessible for genetic testing, opening the door  
38 to broad applications in Mendelian disorders, cancer diagnosis and personalized medicine [1].  
39 However, sequencing data include both systematic and random errors that hinder any of the  
40 current variant identification algorithms from working perfectly. Even using state-of-the-art  
41 approaches, typically 1-3% of the candidate variants are false positives with Illumina  
42 sequencing [2]. With the help of a genome browser such as IGV [3], or web applications such  
43 as VIPER [4], a specialist can visually inspect a graphical layout of the read alignments to

44 assess supporting and contradicting evidence to make an arbitration. Though necessary, this  
45 is a tedious and fallible procedure because of three major drawbacks. 1) It is time-consuming  
46 and empirical studies report it requires about one minute per variant, sometimes summing up  
47 to a few hours per patient [5]. 2) It is tedious, not infallible, and even experienced genetic-  
48 specialists might draw different conclusions for a candidate variant with limited or  
49 contradicting evidence. 3) There is no agreed standard between genetic-specialists to judge  
50 various types of variants, including SNPs (Single Nucleotide Polymorphisms) and Indels. A  
51 specialist might be more stringent on SNPs because there are more clinical assertions and  
52 fewer candidate SNPs will be less likely to get contradicting medical conclusions, whereas  
53 another specialist might be more demanding on indels because they are rarer and harder to be  
54 identified.

55

56 An efficient, accurate and consistent computational method is strongly needed that automates  
57 assessing the candidate variants as they would be visually validated. Importantly, the new  
58 validation method needs to be orthogonal, i.e., independent of the algorithms used to identify  
59 the candidate variants. The new validation method also needs to capture the complex non-  
60 linear relationship between read alignments and the authenticity of a variant from a limited  
61 amount of labeled training data. Variant validation is a task with a different nature from  
62 variant filtration. Our target is to indicate the need of a variant being manually reviewed, as  
63 opposed to a hard filter that removes a variant from consideration. To achieve our target,  
64 failing to flag a false positive variant for review is less favorable than flagging a true variant  
65 for manual review, i.e., as a validation method, the precision must be maximized, and false  
66 positives must be minimized. Consequently, instead of using hand-coded models or rule-  
67 based learning, a more powerful and agnostic machine learning approach such as an Artificial  
68 Neural Network (ANN) is needed.

69

## 70 **Implementation**

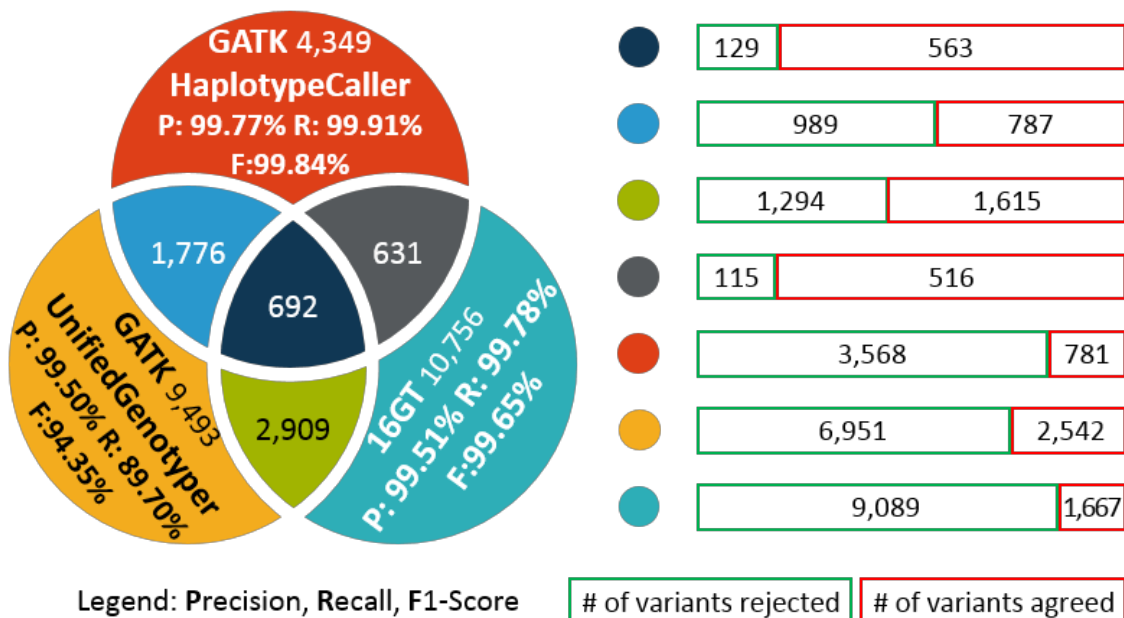
71 We implemented Skyhawk, a computational discriminator that is fast and accurate for  
72 validating candidate variants in clinical practice. Skyhawk mimics how a human visually  
73 identifies genomic features comprising a variant and decides whether the evidence supports  
74 or contradicts the sequencing read alignments. To reach this goal, we repurposed the network  
75 architecture we developed in a previous study named Clairvoyante [6]. The multi-task ANN  
76 was designed for variant calling in Single Molecule Sequencing, and the method is  
77 orthogonal to traditional variant callers using algorithms such as Bayesian or local-assembly.

78 In Skyhawk, we used a repurposed network to generate a probability of each possible option  
79 for multiple categories including 1) variant type, 2) alternative allele, 3) zygosity, and 4)  
80 indel-length. We then compare a candidate variant to Skyhawk's prediction on each category.  
81 Skyhawk will agree with a variant if all categories are matched but will reject and provide  
82 possible corrections if any category is unmatched. We have provided pre-trained models for  
83 Skyhawk on GitHub trained using the known variants and Illumina data of multiple human  
84 genomes, including sequencing libraries prepared by either the PCR or the PCR-free  
85 protocol. With a trained model, Skyhawk accepts a VCF input with candidate SNPs and  
86 Indels, and a BAM input with read alignments. Skyhawk outputs a judgment and a quality  
87 score on how confident the judgment was made for each candidate variant. Skyhawk was  
88 implemented in Python and Tensorflow and has been carefully tuned to maximize its speed.  
89

## 90 Results

91 Using four deeply Illumina sequenced genomes (HG001, HG002, HG003, and HG004) with  
92 13.5M known truth variants from the Genome In A Bottle (GIAB) project [2], we trained  
93 Skyhawk to recognize how the truth variants are different from another 20M non-variants we  
94 randomly sampled from the four genomes. The sample details and the commands used are in  
95 the **Supplementary Note**. For benchmarking and identifying the false positive variant calls,  
96 we used the known truth variants in HG005, which was not included in the model training. A  
97 false positive variant is defined as a variant called by a variant caller but cannot be found in  
98 the HG005 GIAB truth dataset and will be used for the subsequent analysis. We expect the  
99 false positive variants that are supported by only one variant caller, but not the other variant  
100 callers are very likely to be erroneous and should be marked for manual review (i.e., rejected  
101 by Skyhawk) [7]. Thus, we called variants using three different variant callers with different  
102 calculation models, including GATK HaplotypeCaller (HC) [8], GATK UnifiedGenotyper  
103 (UG) [8] and 16GT [9]. A Venn diagram of the variant set called by the three callers  
104 comprise seven different types of variant: 1) three types of singleton variant that have support  
105 from only one caller, 2) three types of doubleton variant that have support from two of the  
106 three callers, and 3) one type of tripleton variant that is supported by all three callers.  
107 Empirically, doubleton and especially tripleton variants are relatively less likely to be real  
108 false positives and should be less likely to be rejected by Skyhawk. Conversely, singletons  
109 called by only one caller are more likely to be genuine false positive and should be more  
110 likely to be rejected by Skyhawk. The results are shown in **Figure 1**. Only 18.64% of the  
111 tripleton variants were rejected while 79.70% of the singletons were rejected by Skyhawk.

112 Those doubletons have an intermediate 45.11% rejected by Skyhawk. In the true positive  
 113 variants, only 1,879/3,232,539 (0.058%) in HC, 43/2,902,052 (0.0014%) in UG, and  
 114 124/3,228,537 (0.0038%) in 16GT were rejected. By deducting the rejected variants from  
 115 both the number of true positives and true negatives, the precision increased from 99.77% to  
 116 99.92% for HC, 99.50% to 99.58% for UG and 99.51% to 99.84% for 16GT.  
 117



118 **Figure 1.** The variant calling results of GATK HaplotypeCaller, GATK UnifiedGenotyper,  
 119 and 16GT. The Venn diagram on the left shows 1) the precision rate (P), recall rate (R) and  
 120 fl-score (F) of each variant caller on all variants of the entire HG005 genome, and 2) the  
 121 number of false positive variants produced by each variant caller. The bars on the right shows  
 122 the number of false positive variants rejected or agreed by Skyhawk. The bar length is  
 123 proportionate to the total number of false positive variants in that type.  
 124

125  
 126 Another experiment better mimics how medical doctors would use Skyhawk in clinical  
 127 diagnosis. Instead of fully removing manual review, which is impossible in a stringent  
 128 clinical context the emphasizes accountability, Skyhawk’s target is to help doctors to  
 129 prioritize which variants should they invest efforts in further investigation and lab validation.  
 130 In practice, those variants categorized as “Pathogenic” or “Likely Pathogenic” are rare and  
 131 should be given priority [10], thus all these variants are preferred to be manually reviewed.  
 132 “Benign”, and most of the time together with the “Likely Benign” category, suggest variants  
 133 without much value in clinical diagnosis and therapy, thus not requiring manual review. The  
 134 one category left, named Variant of Unknown Significance, or VUS, contains variants that  
 135 are potentially impactful, and requires doctors to sort through them. The number of VUS is

136 usually tens to even hundreds of time larger than the sum of other categories [11]. Thus,  
 137 Skyhawk will benefit the clinical doctors if it can significantly decrease the number VUS to  
 138 be manually reviewed. To assess the intended function, we firstly ran GATK  
 139 HaplotypeCaller on the HG002 sample. In total about 5M variants were called. Then we  
 140 annotated all variants using SeattleSeq version 151 (with dbSNP v151) [12]. We extracted  
 141 those variants that are 1) not in dbSNP (RSID tag equals to 0) and, 2) are in a human gene  
 142 (GL tag not empty). Finally, we ran Skyhawk on the extracted variants with a model trained  
 143 on four samples including HG001, HG003, HG004, and HG005, and annotated the variants  
 144 as either true positive (TP) or false positive (FP) against the HG002 GIAB truth dataset.  
 145 Skyhawk performed as expected, and the results are shown in **Table 1**. For SNPs, 53.4% of  
 146 the FPs are flagged for manual review, while only 0.3% of the TPs are flagged. For Indels,  
 147 78.3% of the FPs are flagged for manual review, while only 25.5% of the TPs are flagged. A  
 148 higher rate of TP Indels is flagged for manual review because longer Indels are usually more  
 149 error-prone and can lead to more several clinical consequences than SNPs, thus we required  
 150 all Indels  $\geq 4$ bp to be manually reviewed. Noteworthy, although an ideal percentage of FP  
 151 being marked for manual review is 100%, it is not yet achievable because as mentioned in the  
 152 previous paragraph, FP still have a chance to be an authentic variant especially when it is  
 153 supported by multiple variant callers. Nevertheless, the trend of having significantly more FP  
 154 variants marked for manual review than TP variants verified Skyhawk's effectiveness.

155

156 **Table 1.** Skyhawk's performance on Variants of Unknown Significance (VUS)

		PASS		CHECK	
		#	%	#	%
TP	SNP	4,837	99.7%	14	0.3%
	Indel	7,126	74.5%	2,434	25.5%
FP	SNP	117	46.6%	134	53.4%
	Indel	41	21.7%	148	78.3%

157

158

## 159 Discussion and Conclusions

160 Skyhawk aims to relieve users from a heavy manual review workload without compromising  
 161 the accuracy. Instead of taking over the review of all variants, Skyhawk was configured to  
 162 review only 1) SNPs with a single alternative allele, and 2) Indels  $\leq 4$ bp. Skyhawk also  
 163 outputs a quality score ranging from 0 to 999 to indicate how confident a judgment is.  
 164 Among the false positive singletons, 27.46% of the judgments were with a quality score  
 165 lower than 150. Reviewing these variants manually shows that these variants were often

166 located in genome regions with homopolymer runs or very low depth. We suggest users to  
167 rely on Skyhawk only when the quality score of judgment is high and to manually review  
168 when the quality score falls below 150, or higher if the workload allows. Skyhawk requires  
169 less than a gigabyte of memory and less than a minute on one CPU core to review ten  
170 thousand variants, thus can be easily integrated into existing manual review workflows, such  
171 as VIPER [4] with minimal computational burden. Using 24 CPU cores, Skyhawk was able  
172 to review all five million whole genome sequencing variants of the HG002 sample in 30  
173 minutes. Overall, Skyhawk greatly reduces the workload on reviewing variants, and we  
174 believe Skyhawk will immediately increase the productivity of genetic-specialists in clinical  
175 practice.

176

## 177 Acknowledgment

178 R.L. was supported by the General Research Fund No. 27204518. T. L. was partially  
179 supported by Innovative and Technology Fund ITS/331/17FP from the Innovation and  
180 Technology Commission, HKSAR. This work was also supported, in part, by awards from  
181 the National Science Foundation (DBI-1350041) and the National Institutes of Health (R01-  
182 HG006677).

183

## 184 References

- 185 1. Katsanis SH, Katsanis N: **Molecular genetic testing and the future of clinical**  
186 **genomics.** *Nat Rev Genet* 2013, **14**:415-426.
- 187 2. Zook JM, Catoe D, McDaniel J, Vang L, Spies N, Sidow A, Weng Z, Liu Y, Mason  
188 CE, Alexander N, et al: **Extensive sequencing of seven human genomes to**  
189 **characterize benchmark reference materials.** *Sci Data* 2016, **3**:160025.
- 190 3. Robinson JT, Thorvaldsdottir H, Wenger AM, Zehir A, Mesirov JP: **Variant Review**  
191 **with the Integrative Genomics Viewer.** *Cancer Res* 2017, **77**:e31-e34.
- 192 4. Woste M, Dugas M: **VIPER: a web application for rapid expert review of variant**  
193 **calls.** *Bioinformatics* 2018.
- 194 5. Uzilov AV, Ding W, Fink MY, Antipin Y, Brohl AS, Davis C, Lau CY, Pandya C,  
195 Shah H, Kasai Y: **Development and clinical application of an integrative genomic**  
196 **approach to personalized cancer therapy.** *Genome medicine* 2016, **8**:62.
- 197 6. Luo R, Sedlazeck FJ, Lam T-W, Schatz M: **Clairvoyante: a multi-task**  
198 **convolutional deep neural network for variant calling in Single Molecule**  
199 **Sequencing.** *bioRxiv* 2018.
- 200 7. O'Rawe J, Jiang T, Sun G, Wu Y, Wang W, Hu J, Bodily P, Tian L, Hakonarson H,  
201 Johnson WE: **Low concordance of multiple variant-calling pipelines: practical**  
202 **implications for exome and genome sequencing.** *Genome medicine* 2013, **5**:28.
- 203 8. Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, Del Angel G, Levy-  
204 Moonshine A, Jordan T, Shakir K, Roazen D, Thibault J, et al: **From FastQ data to**

- 205 **high confidence variant calls: the Genome Analysis Toolkit best practices**  
206 **pipeline.** *Curr Protoc Bioinformatics* 2013, **43**:11 10 11-33.
- 207 9. Luo R, Schatz MC, Salzberg SL: **16GT: a fast and sensitive variant caller using a**  
208 **16-genotype probabilistic model.** *GigaScience* 2017, **6**:1-4.
- 209 10. Rehm HL, Bale SJ, Bayrak-Toydemir P, Berg JS, Brown KK, Deignan JL, Friez MJ,  
210 Funke BH, Hegde MR, Lyon E: **ACMG clinical laboratory standards for next-**  
211 **generation sequencing.** *Genetics in medicine* 2013, **15**:733.
- 212 11. Hoffman-Andrews L: **The known unknown: the challenges of genetic variants of**  
213 **uncertain significance in clinical practice.** *Journal of Law and the Biosciences*  
214 2017, **4**:648.
- 215 12. Ng SB, Turner EH, Robertson PD, Flygare SD, Bigham AW, Lee C, Shaffer T, Wong  
216 M, Bhattacharjee A, Eichler EE: **Targeted capture and massively parallel**  
217 **sequencing of 12 human exomes.** *Nature* 2009, **461**:272.
- 218