

16GT: a fast and sensitive variant caller using a 16-genotype probabilistic model

Ruibang Luo^{1,2,*}, Michael C. Schatz^{1,2}, Steven L. Salzberg^{1,2,3}

¹Department of Computer Science, Johns Hopkins University

²Center for Computational Biology, McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University School of Medicine

³Departments of Biomedical Engineering and Biostatistics, Johns Hopkins University

Email addresses:

Ruibang Luo, rluo5@jhu.edu

Michael C. Schatz, mschatz@jhu.edu

Steven L. Salzberg, salzberg@jhu.edu

*Corresponding Author

28 **Abstract**

29 16GT is a variant caller for Illumina whole-genome and whole-exome sequencing data. It uses
30 a new 16-genotype probabilistic model to unify SNP and indel calling in a single variant calling
31 algorithm. In benchmark comparisons with five other widely used variant callers on a modern
32 36-core server, 16GT ran faster and demonstrated improved sensitivity in calling SNPs, and
33 it provided comparable sensitivity and accuracy for calling indels as compared to the GATK
34 HaplotypeCaller. 16GT is available at <https://github.com/aquaskyline/16GT>.

36 **Keywords**

37 Variant calling; Bayesian model; SNP calling; Indel calling

39 **Background**

40 Single nucleotide polymorphisms (SNPs) and insertions and deletions (indels) that occur at
41 a specific genome position are interdependent; i.e., evidence that elevates the probability of
42 one variant type should decrease the probability of other possible variant types, and the
43 probability of all possible alleles should sum to 1. However, widely-used tools such as
44 GATK's UnifiedGenotyper [1] and SAMtools [2] use separate models for SNP and indel
45 detection. The model for SNP calling in these two tools is nearly identical: both assume all
46 variants are biallelic (i.e., exactly two haplotypes are present) and use a probabilistic model
47 allowing for 10 genotypes: AA, AC, AG, AT, CC, CG, CT, GG, GT, TT. For indel calling, the
48 GATK UnifiedGenotyper uses a model from the Dindel's variant caller [3], while SAMtools'
49 model is from BAQ [4].

Findings

In order to detect SNPs and indels with a unified approach, we developed a new 16-genotype probabilistic model and its implementation named 16GT. Building on an idea first introduced in Luo et al. [5], 16GT uses an empirically improved model and is the first publicly available implementation. Using X and Y to denote the indels with the highest (X) and second highest (Y) support, we add 6 new genotypes (AX , CX , GX , TX , XX and XY) to the traditional 10-genotype probabilistic model. The six new genotypes include: 1) one homozygous indel (XX); 2) one reference allele plus one heterozygous indel (AX , CX , GX , TX); 3) one heterozygous SNP plus one heterozygous indel (AX , CX , GX , TX , reusing the genotypes in 2); and 4) two heterozygous indels (XY). We exclude the 5 possible combinations AY , CY , GY , TY , YY because X has higher support than Y . By unifying SNP and indel calling in a single variant calling algorithm, 16GT not only runs 4 times faster, but also demonstrates improved sensitivity in calling SNPs and comparable sensitivity in calling indels to the GATK HaplotypeCaller.

Posterior probabilities of these 16 genotypes are calculated using a Bayesian model $P(L|F) \propto P(F|L)P(L)$, where L is an assumed genotype. F refers to the observation of the 6 alleles (A , C , G , T , X , Y) at a given genome position. $P(L)$ is the prior probability of the genotype, $P(F|L)$ is the likelihood of the observed genotype. and $P(L|F)$ is the posterior probability of the genotype. The resulting genotype L_{max} is assigned to the genotype with the highest posterior probability. The distance between the highest posterior probability and the second highest posterior probability is used as a quality metric in 16GT, along with some other metrics introduced by GATK [1].

Calculating the probability of an observation F given the genotype L

To test how well an observation fits the expectation of different genotypes, we use a two-tailed Fisher's Exact Test P and use the resulting p -value as the goodness of fit. When

78 calculating the likelihood of a homozygous genotype, ideally we expect 100% single allele
 79 support from the observation. For example, consider genotype 'AA':

$$80 \quad P(F|'AA') = P_{hom}(F_A) \times P_e(F_C, F_G, F_T, F_X, F_Y)$$

81 where P_e is the probability of an erroneous base call.

82 For a heterozygous genotype, 50% support is expected for each allele in the genotype, for
 83 example consider 'CG':

$$84 \quad P(F|'CG') = P_{het}(F_C, F_G) \times P_e(F_A, F_T, F_X, F_Y)$$

85 where

$$86 \quad P_{hom}(F_A) = P \left(\begin{array}{c} F_A \\ (1 - P_{err})F \end{array} \quad \begin{array}{c} F \\ F \end{array} \right)$$

$$87 \quad P_{het}(F_C, F_G) = \sqrt{\prod_{i=C,G} P \left(\begin{array}{c} F_i \\ (0.5 - P_{err})F \end{array} \quad \begin{array}{c} F \\ F \end{array} \right)}$$

$$88 \quad P_e(F_A, F_T, F_X, F_Y) = P \left(\begin{array}{c} F_A + F_T + F_X + F_Y \\ P_{err} \times F \end{array} \quad \begin{array}{c} F \\ F \end{array} \right)$$

$$89 \quad F_s = \sum_{i=1}^n f(Q_i, M_i, s) \quad s \in \{A, C, G, T, X, Y\}$$

90 where s is the allele type, n is the number of reads supporting allele s , Q_i is the base quality,
 91 and M_i is the mapping quality. f is a function describing how s , Q_i and M_i change the
 92 observation:

$$93 \quad f(Q_i, M_i, s) = \alpha \times \beta \times \gamma \left\{ \begin{array}{l} \alpha = 0 \text{ if } M_i = 0 \\ \alpha = 1 \text{ if } M_i \neq 0 \\ \beta = 0 \text{ if } Q_i < 10 \\ \beta = 1 \text{ if } 10 \leq Q_i < 13 \\ \beta = 2 \text{ if } 13 \leq Q_i < 17 \\ \beta = 3 \text{ if } 17 \leq Q_i < 20 \\ \beta = 4 \text{ if } Q_i \geq 20 \\ \gamma = 1 \text{ if } s \in \{A, C, G, T\} \\ \gamma = 1.375 \text{ if } s \in \{X, Y\} \end{array} \right.$$

94
 95 The possible reasons for an observation that does not match the reference genome are: 1) a
 96 true variant; 2) an error generated in library construction; 3) a base calling error; 4) a
 97 mapping error; and 5) an error in the reference genome. Reasons 3 and 4 are explicitly

98 captured in our model. For reasons 2 and 5, we include two error probabilities, P_s for SNP
99 error and P_d for indel error. We define $P_{err}=P_s+P_d$, where P_s and P_d are set to 0.01 and 0.005,
100 respectively. These two values were set empirically based on the observation that SNP
101 errors are more common than indel errors in library construction and in the reference
102 genome.

103
104 In addition, most short read aligners use a dynamic programming algorithm to enable
105 gapped alignment, using a scoring scheme that usually penalizes gap opening and
106 extension more than mismatch. Consequently, authentic gaps that occur at an end of a read
107 are more likely to be substituted by a set of false SNPs or alternatively to get trimmed or
108 clipped. Thus, we applied a coefficient γ to weight indel observations more than SNPs, in
109 order to increase the sensitivity on indels.

110 111 ***Calculating the probability of the genotype L***

112 Given 1) a known rate of single nucleotide differences between two unrelated haplotypes; 2)
113 a known rate of single indel differences between two unrelated haplotypes; and 3) a known
114 Transitions to Transversions ratio (Ti/Tv), the 16GT model's prior probabilities are calculated
115 as shown in **Table 1**.

116
117
118
119
120
121
122
123
124

Table 1. $P(L)$, Genotype prior probabilities for a reference allele 'A'.

Hom.: homozygous. Het.: heterozygous.

L	Zygoty	Number of SNPs	Number of Indels	Number of Transversions	Prior Probability $P(L)$
AA	Hom.	-	-	0	1
GG	Hom.	1	0	2	$\theta/2*\epsilon*\epsilon$
CC, TT	Hom.	1	0	0	$\theta/2$
AG	Het.	1	0	1	$\theta*\epsilon$
AC, AT	Het.	1	0	0	θ
CG, GT	Het.	2	0	1	$\theta*\theta/2*\epsilon$
CT	Het.	2	0	0	$\theta*\theta/2$
AX	Het.	0	1	0	ω
GX	Het.	1	1	1	$\omega*\theta/2*\epsilon$
CX, TX	Het.	1	1	0	$\omega*\theta/2$
XX	Hom.	0	1	0	$\omega/2$
XY	Het.	0	2	0	$\omega*\omega/2$

Given 1) a known rate θ of single nucleotide differences between two unrelated haplotypes; 2) a known rate ω of single indel differences between two unrelated haplotypes; and 3) a known Transitions to Transversions ratio (Ti/Tv) ϵ , where transition is expected to occur more frequently than transversion under selective pressure. The default known rates for human genome are: $\theta = 0.001$, $\omega = 0.0001$, $\epsilon = 2.1$, where ϵ is set to the value for human and change between species.

Results

We benchmarked 16GT with GATK UnifiedGenotyper, GATK HaplotypeCaller [1], Freebayes [6], Fermikit [7] and ISAAC [8] using a set of very high-confidence variants developed by the Genome-in-a-bottle (GIAB) project for genome NA12878 [9] (version 2.19, Additional File 1: **Supplementary Note**). The results are shown in **Table 2**.

Table 2. Benchmark comparisons between 16GT and five other variant callers on a dataset from the Genome in a Bottle project consisting of 787M read pairs (53-fold) from genome NA12878. UG: GATK UnifiedGenotyper; HC: GATK HaplotypeCaller. FP: false positive, FN: false negative.

Caller	Time (minutes w/ 36 cores)	SNP						Indel				
		TP	FP				FN	TP	FP			FN
			Total	dbSNP 138	dbSNP 138 %	TP in Omni 2.5			Total	dbSNP 138	dbSNP 138 %	
16GT	121	2,663,179	5,346	4,220	79%	20/20	918	167,549	1,462	944	65%	3,180
UG	29	2,655,608	1,639	563	34%	15/15	8,489	163,839	624	546	88%	6,890
HC	539	2,653,684	419	143	34%	4/4	10,413	168,444	1,232	726	59%	2,285
Freebayes	52	2,655,513	724	353	49%	11/14	8,584	162,505	559	0	0%	8,224
Fermikit	45	2,567,672	2,036	509	25%	9/9	96,425	161,916	1,996	1,076	54%	8,813
ISAAC	63	2,659,438	1,115	586	53%	15/15	4,659	158,642	1,239	710	57%	12,087

140
141 For SNPs, 16GT produced the most true positive calls and the fewest false negative calls;
142 i.e. it has the highest sensitivity among all tools. 79% of 16GT's false positive calls were also
143 reported by dbSNP version 138, which is highest among other callers. However, we should
144 point out that the GIAB variant set is biased towards GATK because it was primarily derived
145 from GATK-based analyses, as reported previously [10]. As a less-biased test, we therefore
146 assessed the false positive calls against a set of unbiased calls made by the Illumina Omni
147 2.5 SNP array (Additional File 1: **Supplementary Note**). Among the 5,346 false positive
148 calls for 16GT, 20 were covered by the Omni array and all 20 (100%) had the correct
149 genotype. Although limited by the small number of measurable alleles in the Illumina Omni
150 2.5 SNP array, only allowing us to reassess 20 'false positive' calls as true positives, the
151 observation that all 20 genotypes out of the 20 covered alleles are correct suggests that a
152 number of the remaining "false positive" calls are actually correct.

153
154 For indels, 16GT produced fewer true positive calls and more false negative calls than
155 HaplotypeCaller, but less than half as many false negative calls as UnifiedGenotyper. 65%
156 of 16GT's false positive indels were covered by dbSNP version 138. Further investigation
157 into the 1,462 false positive indels shows that 981 (67%) of them meet all three of the
158 following criteria: 1) at least three reads supporting the variant; 2) at least one read

159 supporting both the positive and negative strands, and; 3) in over 80% of the reads that
160 support the variant, there exists no other variant in its flanking 10bp. This suggests that
161 some of these “false positives” might be correct, although further experimental validation
162 would be required to confirm this suggestion. Figure 1 shows three examples where the
163 putative false positive from 16GT is likely to be correct.

164

165 **Conclusions**

166 16GT is the firstly publicly available implementation using a 16-genotype probabilistic model
167 for variant calling. Compared with local assembly based variant callers, 16GT provides
168 better sensitivity in SNP calling and comparable sensitivity in indel calling. In the future, we
169 will improve 16GT to support somatic variant detection and extend the model to support
170 variant calling in species with more than two haplotypes.

171

172 **Declarations**

173 ***Acknowledgements***

174 We thank United Electronics Co. Limited for providing code samples for the bam2snapshot
175 function.

176

177 ***Funding***

178 This work has been supported by the U.S. National Institutes of Health under grants R01-
179 HL129239 and R01-HG006677.

180

181 ***Availability of data and materials***

182 Project name: 16GT

183 Project homepage: <https://github.com/aquaskyline/16GT>

184 Archived version: <https://github.com/aquaskyline/16GT/releases/tag/1.0>

185 Operating system: Platform independent

186 Programming language: C++ and Perl

187 Other requirements: See GitHub page

188 License: GPLv3

189 Any restrictions to use by non-academics: None

190

191 ***Authors' contribution***

192 RL, MCS and SLS conceived the study. RL developed and implemented the 16GT algorithm
193 and benchmarked 16GT with other variant callers. RL, MCS and SLS wrote the paper. All
194 authors have read and approved the final version of the manuscript.

195

196 ***Competing interests***

197 The authors declare that they have no competing interests.

198

199 **References**

- 200 1. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella
201 K, Altshuler D, Gabriel S, Daly M, DePristo MA: **The Genome Analysis Toolkit: a
202 MapReduce framework for analyzing next-generation DNA sequencing data.**
203 *Genome Res* 2010, **20**:1297-1303.
- 204 2. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,
205 Durbin R, Genome Project Data Processing S: **The Sequence Alignment/Map
206 format and SAMtools.** *Bioinformatics* 2009, **25**:2078-2079.
- 207 3. Albers CA, Lunter G, MacArthur DG, McVean G, Ouwehand WH, Durbin R: **Dindel:
208 accurate indel calls from short-read data.** *Genome Res* 2011, **21**:961-973.

1 209 4. Li H: **Improving SNP discovery by base alignment quality.** *Bioinformatics* 2011,
2
3 210 **27:1157-1158.**

4
5 211 5. Luo R, Wong YL, Law WC, Lee LK, Cheung J, Liu CM, Lam TW: **BALSA: integrated**
6
7 212 **secondary analysis for whole-genome and whole-exome sequencing,**
8
9 213 **accelerated by GPU.** *PeerJ* 2014, **2:e421.**

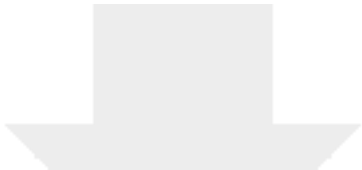
10
11 214 6. Garrison E, Marth G: **Haplotype-based variant detection from short-read**
12
13 215 **sequencing.** *arXiv preprint arXiv:12073907* 2012.

14
15 216 7. Li H: **FermitKit: assembly-based variant calling for Illumina resequencing data.**
16
17 217 *Bioinformatics* 2015:btv440.

18
19 218 8. Racz C, Petrovski R, Saunders CT, Chorny I, Kruglyak S, Margulies EH, Chuang H-
20
21 219 Y, Källberg M, Kumar SA, Liao A: **Isaac: ultra-fast whole-genome secondary**
22
23 220 **analysis on Illumina sequencing platforms.** *Bioinformatics* 2013:btt314.

24
25 221 9. Zook JM, Chapman B, Wang J, Mittelman D, Hofmann O, Hide W, Salit M: **Integrating**
26
27 222 **human sequence data sets provides a resource of benchmark SNP and indel**
28
29 223 **genotype calls.** *Nat Biotechnol* 2014, **32:246-251.**

30
31 224 10. Chiang C, Layer RM, Faust GG, Lindberg MR, Rose DB, Garrison EP, Marth GT,
32
33 225 Quinlan AR, Hall IM: **SpeedSeq: ultra-fast personal genome analysis and**
34
35 226 **interpretation.** *Nat Methods* 2015, **12:966-968.**
36
37 227
38 228
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65



Click here to access/download
Supplementary Material
Additional File 1.docx

