

# Long Single-Molecule Reads Can Resolve the Complexity of the Influenza Virus Composed of Rare, Closely Related Mutant Variants

ALEXANDER ARTYOMENKO<sup>1,\*</sup> NICHOLAS C WU<sup>2,\*</sup> SERGHEI MANGUL<sup>3,4,\*</sup>  
ELEAZAR ESKIN<sup>3</sup> REN SUN<sup>5</sup> and ALEX ZELIKOVSKY<sup>1</sup>

## ABSTRACT

As a result of a high rate of mutations and recombination events, an RNA-virus exists as a heterogeneous “swarm” of mutant variants. The long read length offered by single-molecule sequencing technologies allows each mutant variant to be sequenced in a single pass. However, high error rate limits the ability to reconstruct heterogeneous viral population composed of rare, related mutant variants. In this article, we present two single-nucleotide variants (2SNV), a method able to tolerate the high error rate of the single-molecule protocol and reconstruct mutant variants. 2SNV uses linkage between single-nucleotide variations to efficiently distinguish them from read errors. To benchmark the sensitivity of 2SNV, we performed a single-molecule sequencing experiment on a sample containing a titrated level of known viral mutant variants. Our method is able to accurately reconstruct clone with frequency of 0.2% and distinguish clones that differed in only two nucleotides distantly located on the genome. 2SNV outperforms existing methods for full-length viral mutant reconstruction.

**Keywords:** single-nucleotide variation, SMRT reads, RNA viral variants.

## 1. INTRODUCTION

MAJORITY OF THE EMERGING AND RE-EMERGING DISEASES (influenza, hantaviruses, Ebola virus, and Nipah virus), which represent a global threat to the public health, are caused by RNA viruses (Murphy and Kingsbury, 1996). RNA viruses can be featured by their robust adaptability and evolvability due to their high mutation rates and rapid replication cycles (Holland et al., 1982; Domingo, 1994). This enables a within-host RNA virus population to organize as a complex and dynamic mutant swarm of many highly similar viral genomes. This mutant spectrum, also known as quasispecies (Eigen, 1971), is continuously maintained and re-generated during viral infection (Domingo and Holland, 1997; Lauring and Andino, 2010). Deep sequencing has provided a new lens to monitor individual viral variants accelerating the understanding of escape and resistance

---

<sup>1</sup>Department of Computer Science, Georgia State University, Atlanta, Georgia.

<sup>2</sup>Department of Integrative Structural and Computational Biology, The Scripps Research Institute, La Jolla, California.

<sup>3</sup>Department of Computer Science, University of California, Los Angeles, Los Angeles, California.

<sup>4</sup>Institute for Quantitative and Computational Biosciences, University of California Los Angeles, Los Angeles, California.

<sup>5</sup>Molecular and Medical Pharmacology, University of California, Los Angeles, Los Angeles, California.

\*These authors contributed equally to this work.

mechanisms (Bushman et al., 2008; Margeridon-Thermet et al., 2009), in addition to providing insights about the viral evolutionary landscape and the genomic interactions (Palmer et al., 2006; Wang et al., 2007; Liu et al., 2011).

Short reads offered by commonly used fragmentation-based protocols are well suited to detect discrete genome components, such as the frequency of each single-nucleotide polymorphism. However, high similarity of the individual viral genomes imposes a huge challenge to assemble discrete components into a population of full-length viral genomes. In particular, mutations are often located on the distances unreachable by the short reads (Fig. 1). Therefore, even hybrid technologies based on error correction of PacBio reads with Illumina reads were not applied to sequencing of viral variants. Indeed, short reads cannot tell the allele—the same short read is equally well mapped to a variant with the major allele and a variant with the minor allele.

Single-molecule real-time (SMRT) sequencing is a parallelized single-molecule DNA sequencing method. PacBio SMRT sequencing reads are much longer than sequencing reads provided by Illumina, however, its throughput is much lower and the error rate is significantly higher. The read length offered by a single-molecule sequencing protocol (Eid et al., 2009) is comparable to the genome size of most RNA viruses. It allows each genome variant to be sequenced in a single pass, providing an accurate phasing of the distant mutations. The main drawbacks of the long single-molecule technologies are the high error rate and comparatively low throughput, limiting ability of those technologies to study the heterogeneous viral populations. Thus, a complete profiling of all viral genomes within a mutant spectrum is not yet possible.

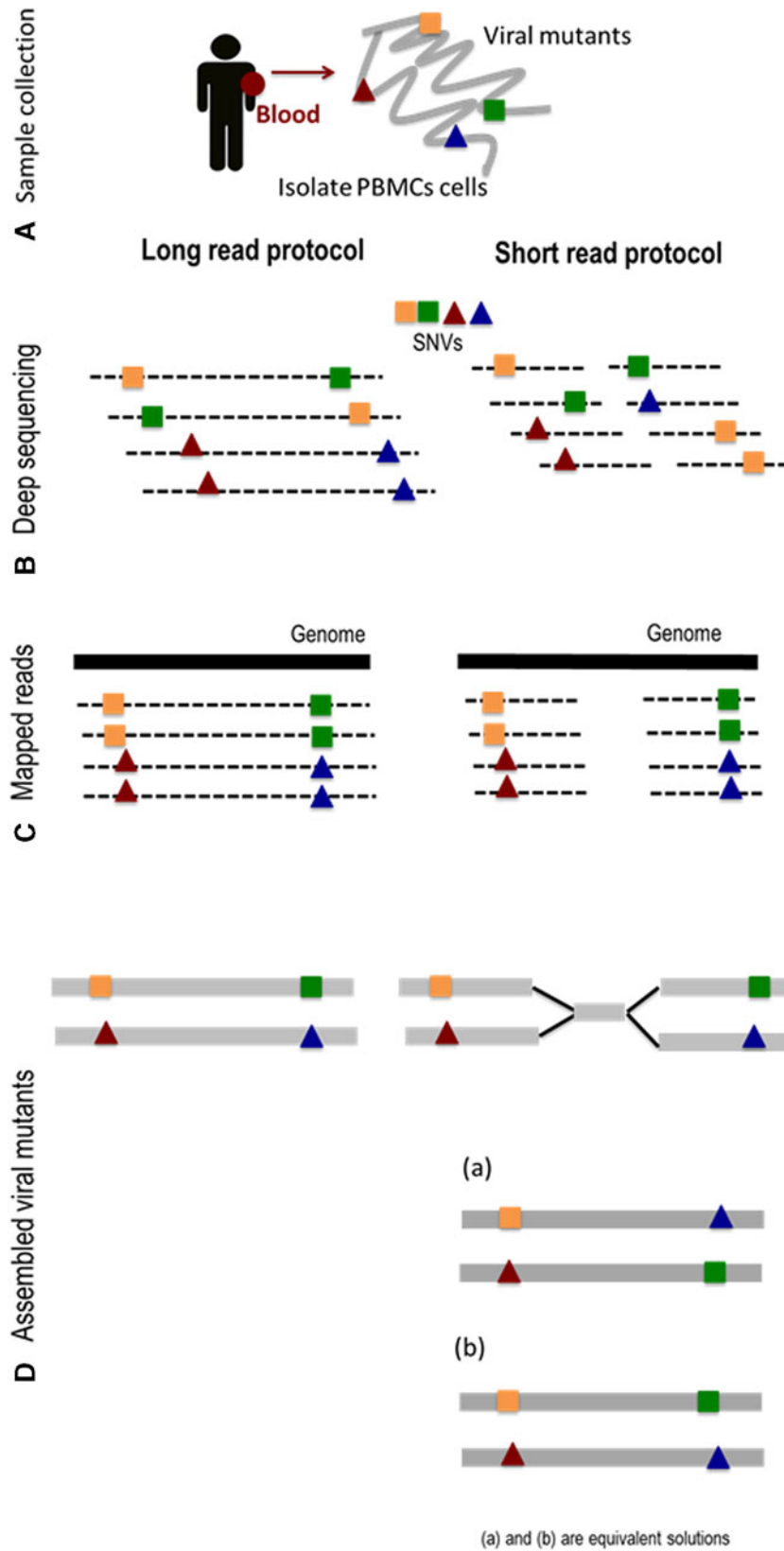
Recently, this problem has been addressed using various computational and statistical approaches implemented in QuasiRecomb (Töpfer et al., 2013), PredictHaplo (Prabhakaran et al., 2014), HaploClique (Töpfer et al., 2014), VGA (Mangul et al., 2014), and  $k$  GEM (Skums et al., 2015). These methods perform reasonably well on short reads with high coverage and low error rate, but our experimental validation shows far from satisfactory performance on the sequencing data provided by single-molecule technologies. Also, a workflow for reconstruction of closely related variants from raw reads generated during SMRT sequencing was proposed in Dilermia et al. (2015). Note that a recent method for haplotyping using PacBio reads proposed in Tilgner et al. (2014) is only applicable for diploid organisms and is not suitable for viral haplotyping with numerous variants.

In this article, we present two single-nucleotide variants (2SNV), a comprehensive method for the accurate reconstruction of the heterogeneous viral population from the long single-molecule reads. The 2SNV method hierarchically clusters together reads containing pairs of correlated (i.e., linked) SNVs until no cluster has correlated SNVs left and outputs consensus of each cluster. It allows to reduce error rate and differentiate true biological variants from sequencing artifacts, thus providing increased accuracy to study diversity and composition of the viral spectrum. To benchmark the sensitivity of 2SNV, we performed a single-molecule sequencing experiment on a sample containing a titrated level of known viral mutant variants. We were able to reconstruct a haplotype with a frequency of 0.2% and distinguish clones that differed in only two nucleotides. We also showed that 2SNV outperformed existing haplotype reconstruction tools. With a high sensitivity and accuracy, 2SNV is anticipated to facilitate not only viral quasispecies reconstruction but also other biological questions that require detection of rare haplotypes, such as genetic diversity in cancer cell population, and monitoring B cell and T cell receptor repertoire. The open source implementation of 2SNV is freely available for download (see the first Reference).

## 2. METHODS

Any method for reconstruction of viral variants from single-molecule reads should overcome low volume and high error rate of sequencing data combined with very high similarity and very low frequency of viral variants. This challenge is equivalent to extraction of an extremely weak signal from very noisy background with signal-to-noise ratio approaching zero. However impossible this task may seem, a satisfactory solution can be based on distinguishing randomness of the noise from systematic signal repetition. Previously, linkage between SNVs was used for distinguishing sequencing errors from SNVs (Macalalad et al., 2012), however, to the best of our knowledge, it was never applied for haplotyping.

Since all reads are from the same RNA region of very similar sequences, they can be reliably aligned to each other. In general, the errors in different positions are independent from each other and the further these positions are from each other the less likely any dependency can be caused by systematic errors. Therefore, even slightly more than expected co-occurrence of two rare alleles in nonadjacent positions may serve as a trustful signature of one or more rare variants having both rare alleles. Such single-nucleotide variations (SNVs) are called linked.



**FIG. 1.** Overview of the long single-molecule sequencing protocol. **(A)** Extract the viral genomic DNA from the whole blood sample. **(B)** DNA material from the viral mutants is cleaved into sequence fragments using any suitable restriction enzyme. Amplified fragments are sequenced. **(C)** Long single-molecule reads are mapped to the reference genome. **(D)** SNVs are detected and assembled into the viral mutant variants. The short read protocol produces equivalent solutions. SNV, single-nucleotide variant.

The proposed 2SNV method recursively clusters reads containing pairs of linked SNVs until no pair of SNVs exhibits statistically significant linkage in any cluster. Then, each cluster should contain just a single viral variant, which can be simply reconstructed as the consensus of all reads in the cluster.

In the remainder of the section, we derive statistical conditions of SNV linkage and then give detailed description of the 2SNV method that identifies rare variant-based SNV pairs satisfying these conditions.

### 2.1. Linkage of SNV pairs

In this section, we analyze statistical significance of the linkage between a pair of SNVs, which allows to distinguish reads emitted by a rare variant from background errors.

We assume that errors are random and a rare variant has at least two mismatches with other variants. Let us consider an arbitrary pair of two distinct positions  $I, J \in \{1, \dots, L\}, I \neq J$ , where  $L$  be the length of the amplicon (Fig. 2b). Let  $I_1$  and  $J_1$  be the alleles of the most frequent 2-haplotype  $(I_1J_1)$ . Note that  $(I_1J_1)$  should be a 2-haplotype from at least one true viral variant assuming that the error rates in the  $I$ -th and  $J$ -th positions are small and independent.

Let  $I_2 \neq I_1$  and  $J_2 \neq J_1$  be the alleles of another 2-haplotype. Let  $E_{kl}, k, l \in \{1, 2\}$ , be the expected number of reads with 2-haplotypes  $(I_kJ_l)$ . The following theorem can be used to decide if the haplotype  $I_2 \neq I_1$  exists.

**Theorem 1.** *Assume that the sequencing error is random, independent, and does not exceed 50%. If no viral variant with the haplotype  $(I_2J_2)$  exists, then the expected value of  $E_{22}$  is at most*

$$E_{22} \leq \frac{E_{21} \cdot E_{12}}{E_{11}}. \quad (1)$$

*The inequality (Equation 1) becomes an equality if at least one of 2-haplotypes  $(I_1J_2)$  or  $(I_2J_1)$  also does not exist.*

**Proof.** Let  $\varepsilon_I^{kl}$  and  $\varepsilon_J^{kl}, k, l \in \{1, 2\}$ , be the probabilities to observe the allele  $l$  instead of the true allele  $k$  in the positions  $I$  and  $J$ , respectively. We are not going to estimate the parameters  $\varepsilon_I^{kl}$ . The model only assumes that these parameters are random, independent, and do not exceed 50%.

Let  $T_{kl}, k, l \in \{1, 2\}$ , be the true count of 2-haplotypes  $(I_kJ_l)$ . Then error randomness and independence imply that

$$E_{kl} = \sum_{m, n=1, 2} \varepsilon_I^{mk} \varepsilon_J^{nl} T_{mn}.$$

To prove (Equation 1), it is sufficient to show that  $E_{11} \cdot E_{22} \leq E_{12} \cdot E_{21}$  assuming that  $T_{22}=0$ . Indeed,

$$\begin{aligned} E_{11} \cdot E_{22} &= \sum_{m, n=1, 2} \varepsilon_I^{m1} \varepsilon_J^{n1} T_{mn} \cdot \sum_{m, n=1, 2} \varepsilon_I^{m2} \varepsilon_J^{n2} T_{mn} \\ &= \varepsilon_I^{11} \varepsilon_J^{11} \varepsilon_I^{12} \varepsilon_J^{12} T_{11}^2 + \varepsilon_I^{21} \varepsilon_J^{11} \varepsilon_I^{22} \varepsilon_J^{12} T_{21}^2 + \varepsilon_I^{11} \varepsilon_J^{21} \varepsilon_I^{12} \varepsilon_J^{22} T_{12}^2 \\ &\quad + (\varepsilon_I^{11} \varepsilon_J^{11} \varepsilon_I^{12} \varepsilon_J^{22} + \varepsilon_I^{11} \varepsilon_J^{21} \varepsilon_I^{12} \varepsilon_J^{12}) T_{11} T_{12} \\ &\quad + (\varepsilon_I^{11} \varepsilon_J^{11} \varepsilon_I^{22} \varepsilon_J^{12} + \varepsilon_I^{21} \varepsilon_J^{11} \varepsilon_I^{12} \varepsilon_J^{12}) T_{11} T_{21} \\ &\quad + (\varepsilon_I^{11} \varepsilon_J^{21} \varepsilon_I^{22} \varepsilon_J^{12} + \varepsilon_I^{21} \varepsilon_J^{11} \varepsilon_I^{12} \varepsilon_J^{22}) T_{12} T_{21} \\ E_{12} \cdot E_{21} &= \sum_{m, n=1, 2} \varepsilon_I^{m1} \varepsilon_J^{n2} T_{mn} \cdot \sum_{m, n=1, 2} \varepsilon_I^{m2} \varepsilon_J^{n1} T_{mn} \\ &= \varepsilon_I^{11} \varepsilon_J^{12} \varepsilon_I^{12} \varepsilon_J^{11} T_{11}^2 + \varepsilon_I^{21} \varepsilon_J^{12} \varepsilon_I^{22} \varepsilon_J^{11} T_{21}^2 + \varepsilon_I^{11} \varepsilon_J^{22} \varepsilon_I^{12} \varepsilon_J^{21} T_{12}^2 \\ &\quad + (\varepsilon_I^{11} \varepsilon_J^{12} \varepsilon_I^{12} \varepsilon_J^{21} + \varepsilon_I^{11} \varepsilon_J^{22} \varepsilon_I^{12} \varepsilon_J^{11}) T_{11} T_{12} \\ &\quad + (\varepsilon_I^{11} \varepsilon_J^{12} \varepsilon_I^{22} \varepsilon_J^{11} + \varepsilon_I^{21} \varepsilon_J^{12} \varepsilon_I^{12} \varepsilon_J^{11}) T_{11} T_{21} \\ &\quad + (\varepsilon_I^{11} \varepsilon_J^{22} \varepsilon_I^{22} \varepsilon_J^{11} + \varepsilon_I^{21} \varepsilon_J^{12} \varepsilon_I^{12} \varepsilon_J^{21}) T_{12} T_{21} \end{aligned}$$

Note that only coefficients for  $T_{12}T_{21}$  are different for these products. Therefore, if either  $T_{12}=0$  or  $T_{21}=0$ , then  $E_{11} \cdot E_{22} = E_{12} \cdot E_{21}$ . Otherwise, let all three 2-haplotypes  $(I_1J_1)$ ,  $(I_1J_2)$ , and  $(I_2J_1)$  exist. Then,

$$\begin{aligned}
& E_{12}E_{21} - E_{11}E_{22} = \\
& = (\varepsilon_I^{11} \varepsilon_J^{22} \varepsilon_I^{22} \varepsilon_J^{11} + \varepsilon_I^{21} \varepsilon_J^{12} \varepsilon_I^{12} \varepsilon_J^{21} - \varepsilon_I^{11} \varepsilon_J^{21} \varepsilon_I^{22} \varepsilon_J^{12} - \varepsilon_I^{21} \varepsilon_J^{11} \varepsilon_I^{12} \varepsilon_J^{22}) T_{12} T_{21} \\
& = \left(1 - \frac{\varepsilon_I^{12} \varepsilon_J^{21}}{\varepsilon_I^{11} \varepsilon_J^{22}}\right) \left(1 - \frac{\varepsilon_J^{12} \varepsilon_I^{21}}{\varepsilon_I^{11} \varepsilon_J^{22}}\right) \varepsilon_I^{11} \varepsilon_J^{22} \varepsilon_I^{22} \varepsilon_J^{11} T_{12} T_{21} > 0
\end{aligned}$$

The last inequality holds since observing the true allele is more probable than observing the erroneous allele and, therefore,  $\varepsilon_I^{kl} < \varepsilon_I^{kk}$  and  $\varepsilon_J^{kl} < \varepsilon_J^{kk}$ ,  $k, l \in \{1, 2\}$ . QED

Note that Theorem 1 does not require linkage disequilibrium of haplotypes—the lack of linkage is explained by errors. The 2SNV method uses Theorem 1 to decide if the alleles  $I_2$  and  $J_2$  are linked as follows. Let  $O_{kl}$ ,  $k, l \in \{1, 2\}$ , be the observed number of reads with 2-haplotypes  $(I_k J_l)$ . Let  $n$  be the total number of reads covering both positions  $I$  and  $J$ , then

$$p = \frac{O_{21} \cdot O_{12}}{O_{11} \cdot n}, \quad (2)$$

is the largest probability of observing the 2-haplotype  $(I_2 J_2)$  among these  $n$  reads. The probability to observe at least  $O_{22}$  reads in the  $(n, p)$  binomial distribution equals

$$Pr(X \geq O_{22}) = 1 - \sum_{i=1}^{O_{22}-1} \binom{n}{i} p^i (1-p)^{n-i}. \quad (3)$$

Since we are looking for a pair of SNVs among  $\binom{L}{2}$  possible pairs, we also adjust to multiple testing using Bonferroni correction requiring

$$1 - \sum_{i=1}^{O_{22}-1} \binom{n}{i} p^i (1-p)^{n-i} \leq \frac{\mathcal{P}}{\binom{L}{2}}, \quad (4)$$

where  $p$  is defined in Equation (2) and  $\mathcal{P}$  is the user-defined  $p$ -value, by default  $\mathcal{P} = 0.01$ .

Finally, when the cluster is too small, the statistical test (Equation 4) may be not stringent enough to weed out spurious linkages. Therefore, we require the number of reads  $O_{22}$  to be at least an empirically defined value (by default equal 30), to decide whether there is an additional haplotype producing these reads.

Note that the binomial model used in Equation (4) may not be stringent enough to compensate for reducing PPV caused by overdispersion especially for higher coverage. In future releases of our tool, we plan to take into account additional variance modeling unknown experimental data processes contributing to variance, e.g., replacing the binomial distribution with the beta-binomial distribution.

#### Algorithm 1 2SNV Algorithm

##### Procedure 1: constructing the consensus haplotype for all reads:

Initialize the set of all clusters with a single cluster with all reads  $\mathcal{C} \leftarrow \{R\}$

For each position  $i$  find allele of highest frequency  $a_i$

$\text{Consensus}(\mathcal{C}) \leftarrow (a_1, \dots, a_L)$

##### Procedure 2: partitioning reads into simple clusters

**while** not all clusters are simple **do**

**for** each non-simple cluster  $C \in \mathcal{C}$  **do**

**if** no pair SNVs is linked according to (2–4) **then**

      Regard  $C$  as a simple cluster

**else**

      Find a pair of linked SNVs  $I_2$  and  $J_2$  minimizing (3)

      Find the set  $C_I$  of all reads with the 2-haplotype  $(I_2 J_2)$

      Find the consensus  $c_1 \leftarrow \text{Consensus}(C_I)$

$C_1 \leftarrow \text{Voronoi}(c_1)$

$C_2 \leftarrow C \setminus C_1$ ,  $c_2 \leftarrow \text{Consensus}(C_2)$

$\mathcal{C} \leftarrow \mathcal{C} \cup \{C_1\} \cup \{C_2\} \setminus \{C\}$

##### Procedure 3: estimating frequencies of the consensuses of simple clusters

Run  $k$ GEM algorithm for the set of haplotypes  $\{\text{Consensus}(C), C \in \mathcal{C}\}$ .

## 2.2. 2SNV method for viral variant reconstruction

The input to 2SNV consists of a set of aligned PacBio reads (Fig. 2a). Alignment required to be in a form of multiple sequence alignment (MSA). The MSA algorithms are too slow to handle PacBio data sets; so instead, we use pairwise alignment by BWA (Li and Durbin, 2009) and b2w from Shorah (Zagordi et al., 2011) to transform pairwise alignment to MSA format.

The main novel step of the 2SNV algorithm identifies a pair of linked SNVs (Fig. 2b) with higher than expected portion of reads containing the 2-haplotype with both minor alleles according to Equations (2)–(4).

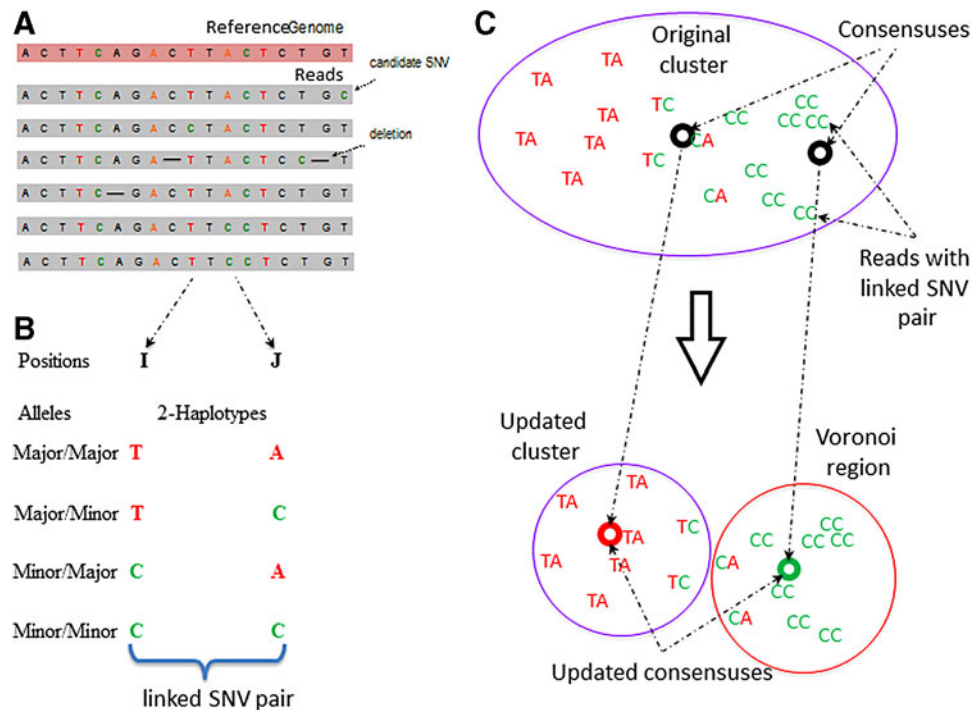
The 2SNV method maintains a partition of all reads into clusters. Each cluster is assumed to consist of the reads emitted by the single variant coinciding with the cluster consensus (Fig. 2c). Until no pair of SNVs in the cluster  $C$  is linked, we recursively partition  $C$  into two clusters  $C_1$  and  $C_2$ .  $C_1$  consists of reads with the linked pair of SNVs and  $C_2$  consists of the remaining reads of  $C$ . We further modify  $C_1$  and  $C_2$  by replacing them with the Voronoi regions of their consensus, where the Voronoi region of the consensus  $c_1$  of  $C_1$  consists of reads that are closer to  $c_1$  than to the consensus of  $C_2$ . Finally,  $k$  GEM finds maximum likelihood estimates of frequencies of haplotypes represented by cluster consensus using expectation-maximization algorithm (Skums et al., 2015).

Algorithm 1 describes the formal pseudocode of the 2SNV algorithm.

## 3. RESULTS

We were using three data sets: PacBio reads from a single IAV clone and 10 IAV clones, and simulated PacBio reads from 20 HCV clones.

Error-prone PCR was performed on the influenza A virus (A/WSN/33) PB2 segment using the GeneMorph II Random Mutagenesis Kits (Agilent Technologies, Westlake Village, CA) according to the manufacturer's instruction. The 2 kb region was amplified from the IAV viral population and subjected to PacBio RS II sequencing using two SMRT cells with P4-C2. The average read length was 1973 bp and ranges from 200 bp to



**FIG. 2.** Overview of the 2SNV method. (a) Multiple sequence alignment of reads from the same amplicon; (b) identification of a linked SNV pair in positions  $I$  and  $J$ ; (c) recursive cluster splitting: (i) finding consensus of reads with the linked SNV pair, (ii) finding Voronoi region of this consensus, and (iii) update the original cluster and the consensus for the two new clusters.

5 kb. Some reads are much longer than the amplified region due to long insertions that are sequencing errors. Raw sequencing data have been submitted to the NIH short read archive under accession number: BioProject PRJNA284802. The nucleotide sequences of the 10 clones are freely available (see the first Reference).

### 3.1. The data set with a single IAV clone

The total number of reads was 11,907 and the average Hamming distance between the true haplotype and reads is 14.4%.

### 3.2. The data set with 10 IAV clones

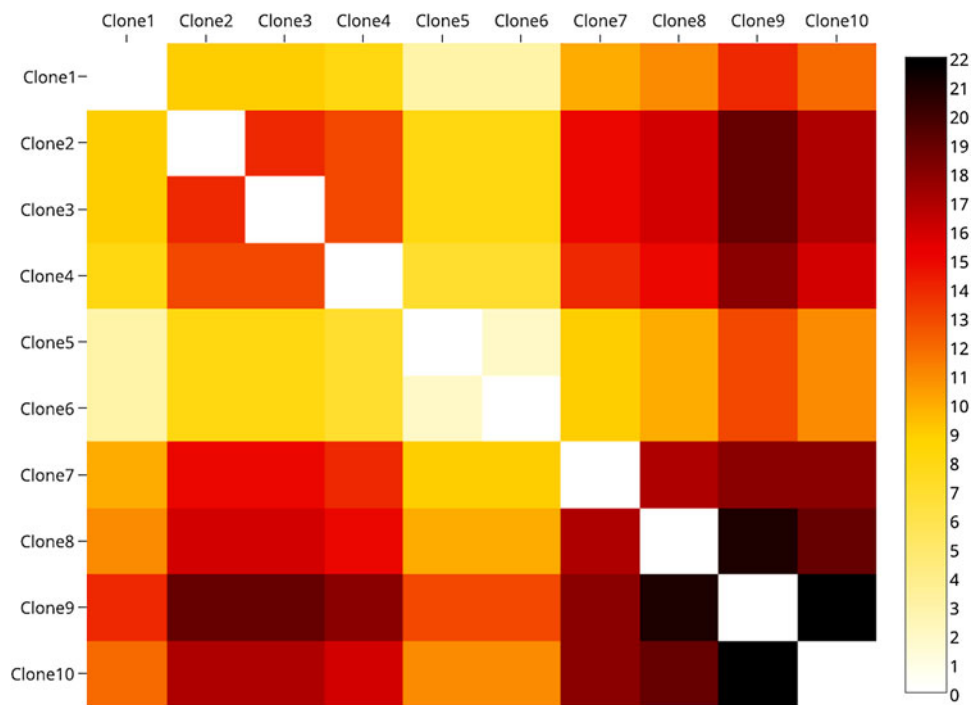
Ten independent clones, ranging from 1 to 13 mutations from the original single, were selected. These 10 clones were mixed at a geometric ratio with a twofold difference in occurrence frequency for consecutive clones starting with the maximum frequency of 50% and the minimum frequency of 0.1%. The pairwise edit distance between clones is given in the heat map on Figure 3. In total, there were 33,558 reads generated from 10 clones.

### 3.3. The simulated data set with 20 HCV clones

21K simulated PacBio reads were generated from 1739-bp-long fragment from the E1E2 region of 20 HCV sequences (Von Hahn et al., 2007) using simulator pbsim (Ono et al., 2013). The reads were simulated with mean accuracy 98% and minimum accuracy 95% reflecting advancements in PacBio technology. We have generated reads 10 times for two distributions of the clone frequencies—uniform (all frequencies are 5%) and skewed (a single clone has 90.5% and every other clone has frequency of 0.5%).

### 3.4. Reconstruction of viral variants

2SNV was compared with two tools originally tuned to handle HIV variants [PredictHaplo; Prabhakaran et al. (2014), QuasiRecomb; Töpfer et al. (2013), and  $k$  GEM; Skums et al. (2015)] tuned for a short HCV amplicon. We could not compare with HaploClique (Töpfer et al., 2014) since it is no longer maintained by the authors. A workflow (Dilernia et al., 2015) is not currently available and we were not able to run it on our data. Also, the experimental data in Dilernia et al. (2015) are also not fully available and we were not able to run 2SNV on these data.



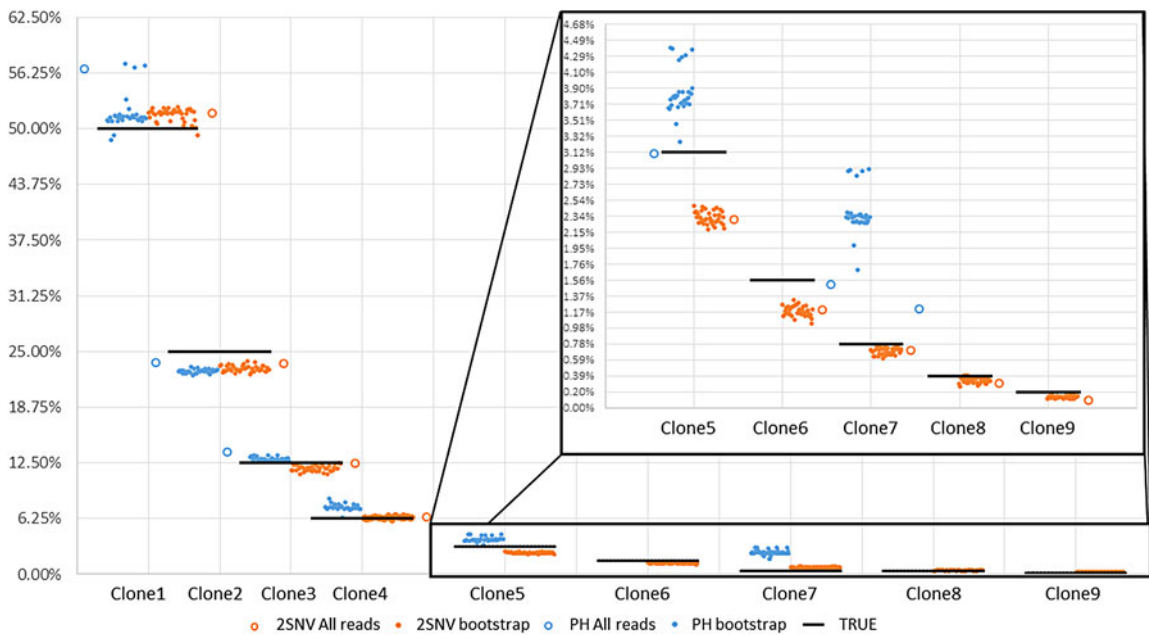
**FIG. 3.** The heatmap representing pairwise edit distance between the 10 IAV clones.

For the data set with a single IAV clone 2SNV,  $k$  GEM and PredictHaplo were able to reconstruct no more than a single variant, which perfectly matches the original clone. QuasiRecomb reported multiple variants, none fully matching the original clone.

For the data set with 10 IAV clones, 2SNV reported 10 haplotypes: the 9 most frequent haplotypes exactly matching 9 most frequent clones and the least frequent haplotype (1%) not matching any clone. The correlation between the estimated and true frequencies of the nine correctly reconstructed haplotypes is 99.4%. PredictHaplo was able to reconstruct only six true variants missing four variants with total frequency of 8% while not having any false positives. To reliably compare the reconstruction rate of two methods, we have applied them to 40 subsamples of the original data (each subsample consists of 33,558 reads randomly selected with repetition from the original data). The results are presented in Figure 4 and Table 1.  $k$  GEM was able to reconstruct only two most frequent clones, and QuasiRecomb failed to reconstruct even a single clone.

To estimate how accuracy of reconstruction methods depends on the coverage, we have randomly subsampled  $N$  reads ( $N=500, 1000, 2000, 4000, 8000, 16,000$ ) from the original 33,558 reads and run 2SNV and PredictHaplo. The results are shown in Figure 5. For each coverage and each clone (except Clone5), 2SNV more accurately estimates the frequency. Clone6 and Clone8 for all subsamples, Clone4 for  $N \leq 8000$ , and Clone 3 for  $N \leq 1000$  are missed by PredictHaplo but reconstructed by 2SNV. Clone6, which is only two mutations away from the more frequent Clone5, was successfully reconstructed for  $N \geq 4000$ , while PredictHaplo was never able to reconstruct Clone6. Note that since these two SNVs between Clone5 and Clone6 are far apart, only long reads can reconstruct this rare variant. From the last plot one can see that the false-positive rate for PredictHaplo is also higher than for 2SNV, for example, 2SNV does not report false positives for  $N \leq 8000$ . The averages of all runs are given in Table 2.

For the simulated data set with 20 HCV variants, we have compared 2SNV only with PredictHaplo. For the uniform frequency distribution, the average sensitivity and PPV for 2SNV are 85% and 100%, respectively, while for PredictHaplo, the corresponding values are 72% and 53%, respectively. For the skewed frequency distribution, the average sensitivity and PPV for 2SNV are 99% and 69%, respectively, while for PredictHaplo, the corresponding values are 36% and 46%, respectively.



**FIG. 4.** The results of running 2SNV and PredictHaplo (PH) on the original sample with all 33,558 reads and on 40 bootstrapped samples (only 35 runs of PH were successful), y axis labels are clone frequencies and x axis labels are clone ids. Horizontal black bars are representing true clone frequency and colored dots are representing frequency reported by corresponding method in each of 40 runs. Clones 6 and 9 were never reconstructed by PH and clone 8 was reconstructed only on full data.



TABLE 1. COMPARISON OF 2SNV AND PREDICTHAPLO ON FULL DATA WITH BOOTSTRAPPING

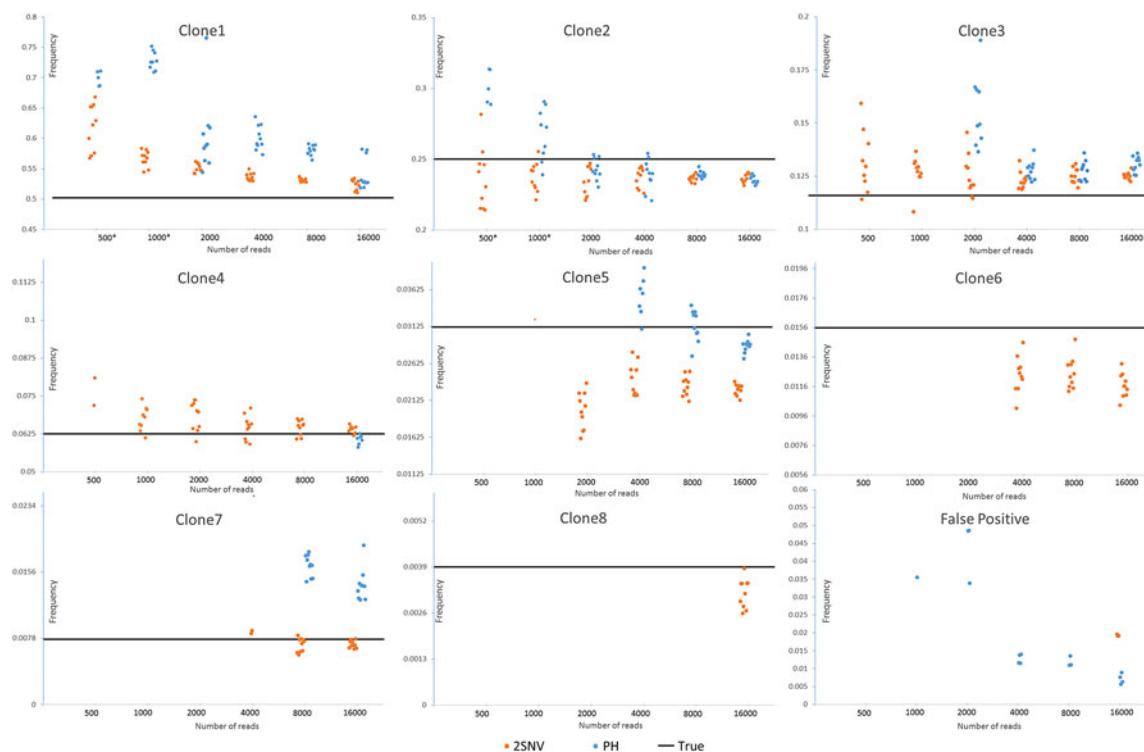
Method	True frequency (%)	Clone 1	Clone 2	Clone 3	Clone 4	Clone 5	Clone 6	Clone 7	Clone 8	Clone 9	Clone 10	FP
			50	25	12.5	6.25	3.125	1.56	0.78	0.39	0.19	0.097
2SNV	Match	✓	✓	✓	✓	✓	✓	✓	✓	✓	×	1
	Frequency (%)	51.8	23.7	12.5	6.4	2.3	1.2	0.7	0.3	0.1	0	1.9
	95% CI											
	Low	50.0	22.5	11.2	6.2	2.2	1.1	0.6	0.3	0.09	—	1.7
	Upper	52.2	23.8	12.4	6.6	2.4	1.3	0.8	0.4	0.14	—	4.0
PH	Match	✓	✓	✓	×	✓	×	✓	✓	×	×	0
	Frequency (%)	56.7	23.8	13.7	0	3.1	0	1.5	1.2	0	0	0
	95% CI											
	Low	49.2	22.3	12.6	0	3.5	—	1.7	—	—	—	0
	Upper	57.0	23.0	13.3	8.0	4.3	—	2.9	—	—	—	7.4

For all 33.5K reads, the sign “✓” (respectively, “×”) denotes fully matched (respectively, unmatched) true variant and the column FP reports the number of incorrectly predicted variants (false positives) and their total frequency. 95 CI showing results on 40 bootstraps, “—” means variant was never reconstructed.

2SNV, two single-nucleotide variants; CI, confidence interval.

### 3.5. Runtime

The runtime of 2SNV is linear with respect to the number of reads, however, implementation is  $O(n \log n)$  due to parallelization (Fig. 6) and quadratic with respect to the length of the amplicon region. For all experiments, we used the same PC (Intel(R) Xeon(R) CPU X5550 2.67 GHz x2 8 cores per CPU, DIMM DDR3 1333 MHz RAM 4 Gb × 12) with operating system CentOS 6.4.



**FIG. 5.** Dependency of accuracy on coverage represented by the number of reads.  $N$  reads ( $N=500, 1000, 2000, 4000, 8000, 16000$ ) were randomly selected 10 times from the original data and both methods 2SNV and PredictHaplo were applied. For  $N=500$  and  $N=1000$ , PredictHaplo gave results only in 5 and 9 runs, respectively. Each dot represents the reconstructed frequency of the clone in the respective runs.

TABLE 2. COMPARISON OF 2SNV AND PREDICTHAPLO ON FULL AND SUBSAMPLED DATA

No. of reads	Method	Clones True frequency, %	Clones	Clones	Clones	Clones	Clones	Clones	Clones	Clones	Clones	Clones	FP	
			1	2	3	4	5	6	7	8	9	10		
33.5K (all)	2SNV	Match	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	×	1
		Frequency (%)	51.8	23.7	12.5	6.4	2.3	1.2	0.7	0.3	0.1	0	0	1.0
	PH	Match	✓	✓	✓	×	✓	×	✓	✓	×	×	×	0
		Frequency (%)	56.7	23.8	13.7	0	3.1	0	1.5	1.2	0	0	0	0
16K	2SNV	Match (%)	100	100	100	100	100	100	100	100	0	0	0	0.2
		Frequency (%)	52.4	23.7	12.5	6.4	2.3	1.1	0.7	0.3	0	0	0	0.6
	PH	Match	100	100	100	70	100	0	100	40	0	0	0	0.3
		Frequency (%)	54.2	23.5	13.1	6.0	2.9	0	1.4	1.0	0	0	0	0.5
8K	2SNV	Match (%)	100	100	100	100	100	100	100	0	0	0	0	0
		Frequency (%)	53.1	23.7	12.5	6.5	2.3	1.25	0.7	0	0	0	0	0
	PH	Match (%)	100	100	100	0	100	0	100	20	0	0	0	0.2
		Frequency (%)	58.1	24.0	12.7	0	3.1	0	1.6	1.3	0	0	0	0.5
4K	2SNV	Match (%)	100	100	100	100	100	100	20	0	0	0	0	0
		Frequency (%)	53.7	23.7	12.3	6.5	2.4	1.2	0.9	0	0	0	0	0
	PH	Match (%)	100	100	100	0	70	0	10	0	0	0	0	0.3
		Frequency (%)	60.1	23.9	12.8	0	3.5	0	2.5	0	0	0	0	0.5
2K	2SNV	Match (%)	100	100	100	100	100	0	0	0	0	0	0	0
		Frequency (%)	55.2	23.4	12.5	6.9	2.0	0	0	0	0	0	0	0
	PH	Match (%)	100	100	90	0	0	0	0	0	0	0	0	0.3
		Frequency (%)	60.4	24.3	15.6	0	0	0	0	0	0	0	0	0.4
1K	2SNV	Match (%)	100	100	100	100	10	0	0	0	0	0	0	0
		Frequency (%)	56.7	23.7	12.7	6.6	3.2	0	0	0	0	0	0	0
	PH	Match (%)	90	90	0	0	0	0	0	0	0	0	0	0.1
		Frequency (%)	72.8	26.8	0	0	0	0	0	0	0	0	0	0.4
0.5K	2SNV	Match (%)	100	100	100	20	0	0	0	0	0	0	0	0
		Frequency (%)	62.0	23.7	12.8	7.6	0	0	0	0	0	0	0	0
	PH	Match (%)	50*	50*	0	0	0	0	0	0	0	0	0	0
		Frequency (%)	69.9	30.1	0	0	0	0	0	0	0	0	0	0

For all 33.5K reads, the sign “✓” (respectively, “×”) denotes fully matched (respectively, unmatched) true variant and the column FP reports the number of incorrectly predicted variants (false positives) and their total frequency. For each subsample size (16K, ..., 0.5K), the table reports the percent of runs when a variant is completely matched and its average frequency. Similarly, the column FP reports the average number of false-positive variants and their average total frequency.

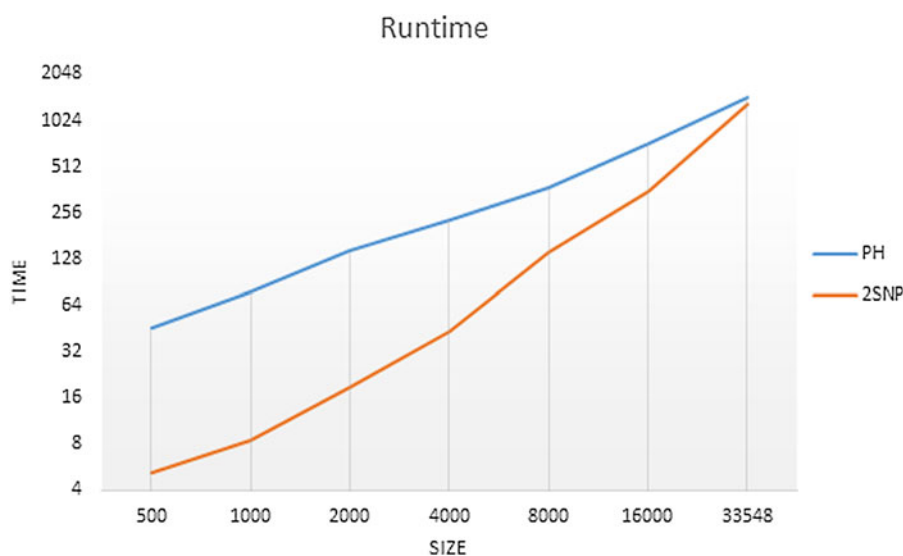


FIG. 6. Runtime of PredictHaplo (PH) and 2SNV on data sets with different sizes. The runtime of 2SNV includes processing of alignment with b2w.

#### 4. DISCUSSION

Haplotype phasing represents one of the biggest challenges in next-generation sequencing due to the short read length. The recent development of single-molecule sequencing platform produces reads that are sufficiently long to span the entire gene or small viral genome. It not only benefits the assembly of genomic regions with tandem repeat (Doi et al., 2014; Ummat and Bashir, 2014; Krsticevic et al., 2015) but also offers the opportunity to examine the genetic linkage between mutations. In fact, it is shown that the long read in single-molecule sequencing aids haplotype phasing in diploid genome (Pendleton et al., 2015), and in polyploid genome (Aguilar and Istrail, 2013). Nonetheless, the sequencing error rate of single-molecule sequencing platform is extremely high ( $\approx 14\%$  as estimated by this study), which hampers its ability to reconstruct rare haplotypes. This drawback prohibits single-molecule sequencing platform from applications, in which a high sensitivity of haplotypes is needed, such as quasispecies reconstruction.

In this study, we have developed 2SNV, which allows quasispecies reconstruction using single-molecule sequencing despite the high sequencing error rate. The high sensitivity of 2SNV permits the detection of extremely rare haplotypes and distinguishes between closely related haplotypes. Based on titrated levels of known haplotypes, we demonstrate that 2SNV is able to detect a haplotype that has a frequency as low as 0.2%. This sensitivity is comparable to many deep sequencing-based point mutation detection methods (Harismendy et al., 2011; Flaherty et al., 2012; Forshew et al., 2012; Li and Stoneking, 2012). In addition, 2SNV successfully distinguishes between Clone5 and Clone6 in this study, which are only two nucleotides away from each other. It highlights the sensitivity of 2SNV to distinguish closely related haplotypes. Our results also show that the sensitivity is coverage dependent, implying that the sensitivity of 2SNV may further improve when sequencing depth increases. Therefore, the constant increase of sequencing throughput offered by single-molecule sequencing technology provides the unprecedented resolution promising to increase number of discovered rare haplotypes.

The ability to accurately determine the genomic composition of the viral populations and identify closely related viral genomes makes our tool applicable for dissecting evolutionary trajectories and examining mutation interactions in RNA viruses. Evolutionary trajectories and mutation interactions have been shown to play an important role in viral evolution, such as drug resistance (Beerenwinkel et al., 2002; Wang et al., 2007; Bushman et al., 2008; Margeridon-Thermet et al., 2009), immune escape (Goepfert et al., 2008), and cross-species adaptation (Herfst et al., 2012; Imai et al., 2012). An unbiased and accurate understanding of the genomic composition of the RNA viruses opens a new avenue to study the underlying mechanism of adaptation, persistence, and virulence factors of the pathogen, which are yet to be comprehended.

While viral quasispecies reconstruction is used as a proof-of-concept in this study, the application of 2SNV can be extended to detect haplotype variants in any sample with high genetic heterogeneity and diversity, such as B cell and T cell receptor repertoire, cancer cell populations, and metagenomes. It is shown that monitoring B cell and T cell receptor repertoire helps investigate virus–host interaction dynamics (Wu et al., 2011; Klarenbeek et al., 2012; Miconnet, 2012; Zhu et al., 2013a,b). Furthermore, examining the genetic composition of the cancer cell populations in high sensitivity can facilitate diagnosis and treatment (Mardis and Wilson, 2009). Therefore, we anticipate that 2SNV will benefit different subfields of biomedical research in the genomic era. We also propose that 2SNV can be applied to increase the resolution of metagenomics profiling from species level to strain level. In summary, 2SNV is a widely applicable tool as a single-molecule sequencing technology being popularized.

#### ACKNOWLEDGMENTS

We thank H. Hao for performing the PacBio sequencing at Johns Hopkins Deep Sequencing & Microarray Core Facility. A.A. was supported by GSU Molecular Basis of Disease Fellowship. A.A. and A.Z. were supported by NSF grants 1619110 and 1564899. S.M. and E.E. were supported by NSF grants 0513612, 0731455, 0729049, 0916676, 1065276, 1302448, 1320589, and 1331176, and NIH grants K25-HL080079, U01-DA024417, P01-HL30568, P01-HL28481, R01-GM083198, R01-ES021801, R01-MH101782, and R01-ES022282. E.E. is supported, in part, by the NIH BD2K award, U54EB020403. S.M. was supported, in part, by the Institute for Quantitative & Computational Biosciences Fellowship, UCLA.

## AUTHOR DISCLOSURE STATEMENT

No competing financial interests exist.

## REFERENCES

- The open source implementation of 2SNV is freely available for download at <http://alan.cs.gsu.edu/NGS/?q=content/2snv>
- Aguiar, D., and Istrail, S. 2013. Haplotype assembly in polyploid genomes and identical by descent shared tracts. *Bioinformatics* 29, i352–i360.
- Beerenwinkel, N., Schmidt, B., and Walter, H., et al. 2002. Diversity and complexity of hiv-1 drug resistance: A bioinformatics approach to predicting phenotype from genotype. *Proc. Natl Acad. Sci.* 99, 8271–8276.
- Bushman, F.D., Hoffmann, C., and Ronen, K., et al. 2008. Massively parallel pyrosequencing in hiv research. *AIDS* 22, 1411–1415.
- Dilernia, D.A., Chien, J.-T., and Monaco, D.C., et al. 2015. Multiplexed highly-accurate DNA sequencing of closely-related HIV-1 variants using continuous long reads from single molecule, real-time sequencing. *Nucl. Acids Res.* 43(20), e129.
- Doi, K., Monjo, T., Hoang, et al. 2014. Rapid detection of expanded short tandem repeats in personal genomics using hybrid sequencing. *Bioinformatics* 30, 815–822.
- Domingo, E. 1994. Mutation rates and rapid evolution of RNA viruses, 161–184. *In The Evolutionary Biology of Viruses*. Raven Press.
- Domingo, E.J.J.H., and Holland, J.J. 1997. RNA virus mutations and fitness for survival. *Annu. Rev. Microbiol.* 51, 151–178.
- Eid, J., Fehr, A., and Gray, J., et al. 2009. Real-time dna sequencing from single polymerase molecules. *Science* 323, 133–138.
- Eigen, M. 1971. Selforganization of matter and the evolution of biological macromolecules. *Naturwissenschaften* 58, 465–523.
- Flaherty, P., Natsoulis, G., Muralidharan, O., et al. 2012. Ultrasensitive detection of rare mutations using next-generation targeted resequencing. *Nucleic Acids Res.* 40, e2.
- Forsheew, T., Murtaza, M., Parkinson, C., et al. 2012. Noninvasive identification and monitoring of cancer mutations by targeted deep sequencing of plasma DNA. *Sci. Transl. Med.* 4, 136ra68.
- Goepfert, P.A., Lumm, W., and Farmer, P., et al. 2008. Transmission of hiv-1 gag immune escape mutations is associated with reduced viral load in linked recipients. *J. Exp. Med.* 205, 1009–1017.
- Harismendy, O., Schwab, R.B., Bao, L., et al. 2011. Detection of low prevalence somatic mutations in solid tumors with ultra-deep targeted sequencing. *Genome Biol.* 12, R124.
- Herfst, S., Schrauwen, E.J.A., Linster, M., et al. 2012. Airborne transmission of influenza a/h5n1 virus between ferrets. *Science* 336, 1534–1541.
- Holland, J., Spindler, K., Horodyski, F., et al. 1982. Rapid evolution of RNA genomes. *Science* 215, 1577–1585.
- Imai, M., Watanabe, T., Hatta, M., et al. 2012. Experimental adaptation of an influenza h5 ha confers respiratory droplet transmission to a reassortant h5 ha/h1n1 virus in ferrets. *Nature* 486, 420–428.
- Klarenbeek, P.L., Remmerswaal, E.B., ten Berge, I.J., et al. 2012. Deep sequencing of antiviral T-cell responses to HCMV and EBV in humans reveals a stable repertoire that is maintained for many years. *PLoS Pathog.* 8, e1002889.
- Krsticevic, F.J., Schrago, C.G., and Carvalho, A.B. 2015. Long-read single molecule sequencing to resolve tandem gene copies: The Mst77Y region on the *Drosophila melanogaster* Y chromosome. *G3 (Bethesda)* 5, 1145–1150.
- Lauring, A.S. and Andino, R. 2010. Quasispecies theory and the behavior of RNA viruses. *PLoS Pathog.* 6, e1001005.
- Li, H. and Durbin, R. 2009. Fast and accurate short read alignment with burrows–wheeler transform. *Bioinformatics* 25, 1754–1760.
- Li, M., and Stoneking, M. 2012. A new approach for detecting low-level mutations in next-generation sequence data. *Genome Biol.* 13, R34.
- Liu, J., Miller, Danovich, R.M., et al. 2011. Analysis of low-frequency mutations associated with drug resistance to raltegravir before antiretroviral treatment. *Antimicrob. Agents Chemother.* 55, 1114–1119.
- Macalalad, A.R., Zody, M.C., Charlebois, P., et al. 2012. Highly sensitive and specific detection of rare variants in mixed viral populations from massively parallel sequence data. *PLoS Comput. Biol.* 8, e1002417.
- Mangul, S., Wu, N.C., Mancuso, N., et al. 2014. Accurate viral population assembly from ultra-deep sequencing data. *Bioinformatics* 30, i329–i337.
- Mardis, E.R., and Wilson, R.K. 2009. Cancer genome sequencing: A review. *Hum. Mol. Genet.* 18(R2), R163–168.
- Margeridon-Thermet, S., Shulman, N.S., Ahmed, A., et al. 2009. Ultra-deep pyrosequencing of hepatitis b virus quasispecies from nucleoside and nucleotide reverse-transcriptase inhibitor (NRTI)-treated patients and NRTI-naive patients. *J. Infect. Dis.* 199, 1275–1285.

- Miconnet, I. 2012. Probing the T-cell receptor repertoire with deep sequencing. *Curr. Opin. HIV AIDS* 7, 64–70.
- Murphy, F.A. and Kingsbury, D.W. 1996. Virus taxonomy. *Fields Virol.* 2, 15–57.
- Ono, Y., Asai, K., and Hamada, M. 2013. PBSIM: Pacbio reads simulator toward accurate genome assembly. *Bioinformatics* 29, 119–121.
- Palmer, S., Boltz, V., Maldarelli, F., et al. 2006. Selection and persistence of non-nucleoside reverse transcriptase inhibitor-resistant HIV-1 in patients starting and stopping non-nucleoside therapy. *AIDS* 20, 701–710.
- Pendleton, M., Sebra, R., Pang, A.W., et al. 2015. Assembly and diploid architecture of an individual human genome via single-molecule technologies. *Nat. Methods.* 12, 780–786
- Prabhakaran, S., Rey, M., Zagordi, O., et al. 2014. HIV haplotype inference using a propagating Dirichlet process mixture model. *IEEE/ACM Trans. Comput. Biol. Bioinf.* 11, 182–191.
- Skums, P., Artyomenko, A., Glebova, O., et al. 2015. Computational framework for next-generation sequencing of heterogeneous viral populations using combinatorial pooling. *Bioinformatics* 31, 682–690.
- Tilgner, H., Grubert, F., Sharon, D., et al. 2014. Defining a personal, allele-specific, and single-molecule long-read transcriptome. *Proc. Natl Acad. Sci.* 111, 9869–9874.
- Töpfer, A., Marschall, T., Bull, R.A., et al. 2014. Viral quasispecies assembly via maximal clique enumeration. *RECOMB* 309–310. Proceedings. Vol. 8394. Springer, 2014.
- Töpfer, A., Zagordi, O., Prabhakaran, S., et al. 2013. Probabilistic inference of viral quasispecies subject to recombination. *J. Comput. Biol.* 20, 113–123.
- Ummat, A., and Bashir, A. 2014. Resolving complex tandem repeats with long reads. *Bioinformatics* 30, 3491–3498.
- Von Hahn, T., Yoon, J.C., Alter, H., et al. 2007. Hepatitis c virus continuously escapes from neutralizing antibody and t-cell responses during chronic infection in vivo. *Gastroenterology* 132, 667–678.
- Wang, C., Mitsuya, Y., Gharizadeh, B., et al. 2007. Characterization of mutation spectra with ultra-deep pyrosequencing: Application to HIV-1 drug resistance. *Genome Res.* 17, 1195–1201.
- Wu, X., Zhou, T., Zhu, J., et al. 2011. Focused evolution of HIV-1 neutralizing antibodies revealed by structures and deep sequencing. *Science* 333, 1593–1602.
- Zagordi, O., Bhattacharya, A., Eriksson, N., et al. 2011. Shorah: Estimating the genetic diversity of a mixed sample from next-generation sequencing data. *BMC Bioinf.* 12, 119.
- Zhu, J., Ofek, G., Yang, Y., et al. 2013a. Mining the antibodyome for HIV-1-neutralizing antibodies with next-generation sequencing and phylogenetic pairing of heavy/light chains. *Proc. Natl. Acad. Sci. U.S.A.* 110, 6470–6475.
- Zhu, J., Wu, X., Zhang, B., et al. 2013b. De novo identification of VRC01 class HIV-1-neutralizing antibodies by next-generation sequencing of B-cell transcripts. *Proc. Natl. Acad. Sci. U.S.A.* 110, E4088–E4097.

Address correspondence to:

*Prof. Alex Zelikovsky*  
*Department of Computer Science*  
*Georgia State University*  
*25 Park Place*  
*Room 751*  
*Atlanta, GA 30303*

*E-mail: alexz@cs.gsu.edu*