

# Accurate assembly of minority viral haplotypes from next-generation sequencing through efficient noise reduction

Sergey Knyazev<sup>1,2,3,\*</sup>, Viachaslau Tsyvina<sup>1,†</sup>, Anupama Shankar<sup>2</sup>, Andrew Melnyk<sup>1</sup>, Alexander Artyomenko<sup>4</sup>, Tatiana Malygina<sup>5</sup>, Yuri B. Porozov<sup>6,7</sup>, Ellsworth M. Campbell<sup>2</sup>, William M. Switzer<sup>2</sup>, Pavel Skums<sup>1</sup>, Serghei Mangul<sup>8,‡</sup> and Alex Zelikovskiy<sup>1,6,\*</sup>

<sup>1</sup>Department of Computer Science, Georgia State University, Atlanta, GA 30302, USA, <sup>2</sup>Division of HIV Prevention, Centers for Disease Control and Prevention, Atlanta, GA 30333, USA, <sup>3</sup>Oak Ridge Institute for Science and Education, Oak Ridge, TN 37830, USA, <sup>4</sup>Guardant Health Inc., Redwood City, CA 94063, USA, <sup>5</sup>International Scientific and Research Institute of Bioengineering, ITMO University, St. Petersburg 197101, Russia, <sup>6</sup>World-Class Research Center “Digital biodesign and personalized healthcare”, I.M. Sechenov First Moscow State Medical University, Moscow 119991, Russia, <sup>7</sup>Department of Computational Biology, Sirius University of Science and Technology, Sochi 354340, Russia and <sup>8</sup>Department of Clinical Pharmacy, School of Pharmacy, University of Southern California, Los Angeles, CA 90089, USA

Received October 09, 2020; Revised May 25, 2021; Editorial Decision June 08, 2021; Accepted June 18, 2021

## ABSTRACT

Rapidly evolving RNA viruses continuously produce minority haplotypes that can become dominant if they are drug-resistant or can better evade the immune system. Therefore, early detection and identification of minority viral haplotypes may help to promptly adjust the patient’s treatment plan preventing potential disease complications. Minority haplotypes can be identified using next-generation sequencing, but sequencing noise hinders accurate identification. The elimination of sequencing noise is a non-trivial task that still remains open. Here we propose CliqueSNV based on extracting pairs of statistically linked mutations from noisy reads. This effectively reduces sequencing noise and enables identifying minority haplotypes with the frequency below the sequencing error rate. We comparatively assess the performance of CliqueSNV using an *in vitro* mixture of nine haplotypes that were derived from the mutation profile of an existing HIV patient. We show that CliqueSNV can accurately assemble viral haplotypes with frequencies as low as 0.1% and maintains consistent performance across short and long bases sequencing platforms.

## INTRODUCTION

Rapidly evolving RNA viruses, such as human immunodeficiency virus (HIV), hepatitis C virus (HCV), influenza A virus (IAV), SARS and SARS-CoV-2 form populations of closely related genomic variants inside infected hosts (1–10). The intra-host viral populations include minority viral variants that are frequently responsible for drug resistance, immune escape and disease transmission (11–24). Therefore, accurately predicting minority viral populations from extremely large and noisy viral genomic data is important for biomedical research, epidemiology and clinical applications. Although this problem has recently attracted significant interest from the biomedical research community (25–27), numerous obstacles still delay next-generation sequencing (NGS) integration into the viral studies. The last decade witnessed numerous attempts to employ NGS and bioinformatics methods for reconstructing intra-host viral populations. These methods are not accurate enough for clinical and epidemiological applications since they cannot reliably identify haplotypes accounting for a substantial portion of the population. Existing methods are ill-equipped to assemble closely related haplotypes and have elevated false-positive rates. Additionally, there is only one *in vitro* viral sequencing benchmark for validation of haplotyping tools (27), and to convincingly demonstrate that such tools are ready for clinical and epidemiological appli-

\*To whom correspondence should be addressed. Tel: +1 404 6631985; Fax: +1 404 4135717; Email: alexz@gsu.edu  
Correspondence may also be addressed to Sergey Knyazev. Tel: +1 470 2631752; Fax: +1 404 4135717; Email: sergey.n.knyazev@gmail.com

†The authors wish it to be known that, in their opinion, the first two authors should be regarded as Joint First Authors.

‡These authors jointly supervised this work.

**Disclaimer:** The findings and conclusions in this paper are those of the authors and do not necessarily represent the official position of the CDC.

cations, new comprehensive sequencing benchmarks are urgently required (28).

NGS technologies now provide versatile opportunities to study viral populations. In particular, the popular Illumina MiSeq/HiSeq platforms produce 25–320 million reads, which allow multiple coverage of highly variable viral genomic regions. This high coverage is essential for capturing rare variants. Ability of NGS technologies to efficiently identify minority variants have recently gained U.S. Food & Drug Administration (FDA) approval. However, *haplotyping* of heterogeneous viral populations (i.e. assembly of full-length genomic variants and estimation of their frequencies) is extremely complicated due to the vast number of sequencing reads, the need to assemble an unknown number of closely related viral sequences and to identify and preserve low-frequency variants. Single-molecule sequencing technologies, such as PacBio, provide an alternative to short-read sequencing by allowing full-length viral variants to be sequenced in a single pass. However, the high level of sequence noise due to background or platform-specific sequencing errors produced by all currently available platforms makes inference of low-frequency genetically close variants especially challenging, since it is required to distinguish between real and artificial genetic heterogeneity produced by sequencing errors.

Recently, a number of computational tools for inference of viral quasiespecies populations from NGS reads have been proposed (28), including Savage (26), PredictHaplo (29), aBayesQR (30), QuasiRecomb (31), HaploClique (32), VGA (33), VirA (34,35), SHORAH (36), ViSpA (37), QURE (38) and others (39–43). Even though these algorithms proved useful in many applications, accurate and scalable viral haplotyping remains a challenge. In particular, inference of low-frequency viral variants is still problematic, while many computational tools designed for the previous generation of sequencing platforms have severe scalability problems when applied to datasets produced by state-of-the-art technologies.

Previously, several tools, such as V-phaser (44), V-phaser2 (45) and CoVaMa (46) exploited linkage of mutations for single nucleotide variant (SNV) calling rather than haplotype assembly, but they do not accommodate sequencing errors when deciding whether two variants are linked. These tools are also unable to detect the frequency of mutations above sequencing error rates (47). The 2SNV algorithm (48) accommodates errors in links and was the first such tool to be able to correctly detect haplotypes with a frequency below the sequencing error rate.

We propose a novel method that can accurately identify minority haplotypes from NGS reads consisting of three steps. First, we extract pairs of statistically linked mutations. Second, we find maximal sets of pairwise linked mutations (cliques) where each clique corresponds to a set of mutations in a minority haplotype. Finally, we assign each read to the closest clique, and for each clique, we form a haplotype as a consensus of reads assigned to it.

All haplotyping tools require solid and convincing validation benchmarks (49,50). The true viral variants and their distribution are only known for simulated data (51), but sequencing errors, variation of coverage depth, polymerase chain reaction bias and systematic noise are diffi-

cult to simulate [see e.g. (52)]. Therefore experimental sequencing benchmarks that provide an adequate evaluation of haplotyping tools are necessary.

By now, there are only two experimental sequencing benchmarks—(i) Illumina sequencing reads consisting of a mixture of five HIV-1 strains (HIV5exp, see Table 1) (27) and (ii) PacBio sequencing reads from a sample consisting of ten IAV viral variants (IAV10exp, see Table 1) (48). In the HIV5exp, five different HIV-1 strains each having 20% frequency were prepared to mimic an intra-host viral population. Unfortunately, this benchmark is not realistic enough since the observed intra-host viral populations consist of variants that are much closer to each other than different strains and contain both frequent and rare variants (53). The IAV10exp benchmark significantly better mimics the intra-host viral population since its variants are very similar to each other and the variant frequencies are realistically non-uniform. Thus, similar to the IAV10exp benchmark, it would be beneficial to develop Illumina benchmarks which adequately imitate intra-host viral populations containing closely related minority variants.

To validate our method's performance, we have introduced two novel in vitro sequencing HIV-1 benchmarks, which consist of Illumina MiSeq experiments on haplotype mixtures based on the mutation profile from an existing patient.

Finally, there is an essential gap in existing quality measures of intra-host viral population assembly. Up-to-date, instead of *populations* (i.e. haplotypes with their frequencies), only *sets* of reconstructed and the ground truth haplotypes are compared (29). Here we propose to measure differences between haplotype populations using Matching Error and the Earth mover's distance (EMD) which account for both the distances between haplotypes and their frequencies.

## MATERIALS AND METHODS

### CliqueSNV algorithm idea

A schematic diagram of the CliqueSNV algorithm is shown in Figure 1. The algorithm takes aligned reads as input and infers haplotype sequences with their frequencies as output. The method consists of six steps:

- i Step 1 uses aligned reads to build the consensus sequence and identifies all SNVs. Then all pairs of SNVs are divided into three groups: *linked*, *forbidden* and *unclassified*. Each SNV is represented as a pair  $(p, n)$  of its position  $p$  and nucleotide value  $n$  in the aligned reads. If there are enough reads that have two SNVs  $(p, n)$  and  $(p', n')$  simultaneously, then we estimate probability that there are no haplotypes simultaneously containing both SNVs (see CliqueSNV algorithm details). If this probability is low, then the algorithm classifies these two SNVs as *linked*. Otherwise, we estimate probability that there is a sufficiently frequent haplotype simultaneously containing both SNVs (see CliqueSNV algorithm details). If this probability is low, then the algorithm classifies these two SNVs as a *forbidden* pair. If the both estimated probabilities are not small, then the pair of SNVs remains *unclassified*.

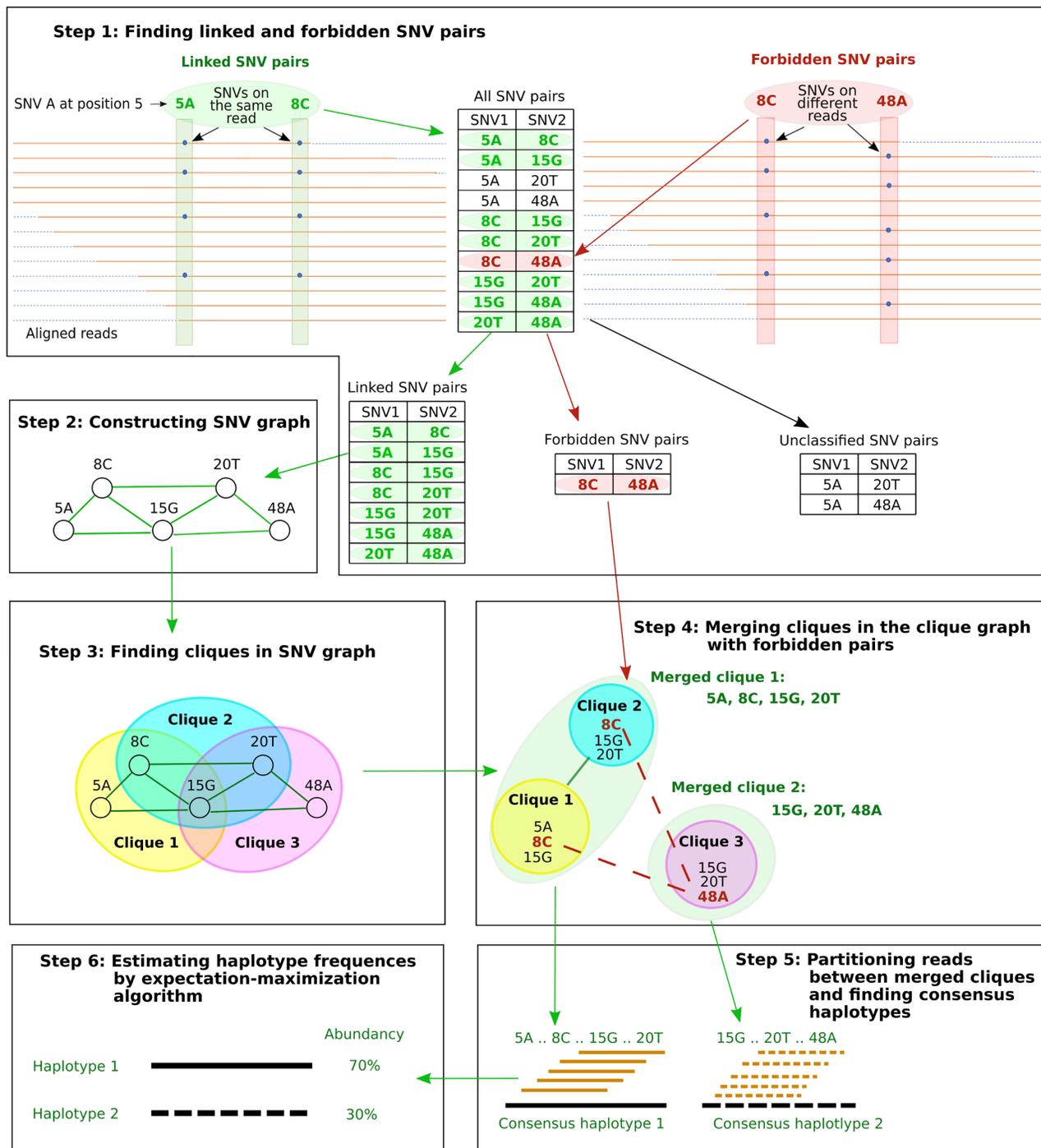


Figure 1. Schematic representation of the CliqueSNV algorithm. Here, SNV is a single nucleotide variation.

- ii In Step 2, we build a graph  $G = (V, E)$  with a set of nodes  $V$  representing SNVs, and a set of edges  $E$  connecting linked SNV pairs.
- iii Ideally, SNVs of each true minority haplotype form a clique in  $G$ . A maximal clique  $C \subseteq V$  is a set of nodes such that  $(u, v) \in E$  for any  $u, v \in C$  and for any  $x \notin C$  there is  $u \in C$  such that  $(x, u) \notin E$ . Step 3 finds all maximal cliques in  $G$ .
- iv For real sequencing data, the linkage between some SNV pairs may be undetected due to sequencing noise, uneven coverage, or the shortness of the NGS reads. As a result, a single clique corresponding to a haplotype will be split into several overlapping cliques. Step 4 merges such overlapping cliques. In order to avoid merging distinct haplotypes, two cliques are not merged if they contain a forbidden SNV pair.

**Table 1.** Four experimental and two simulated sequencing datasets of human immunodeficiency virus type 1 (HIV-1) and IAV

Name	Type	Virus	#haplotypes	Haplotype frequencies	Hamming distance
HIV9exp	experimental	HIV-1	9	0.2-50%	0.22-2.1%
HIV2exp	experimental	HIV-1	2	50-50%	1.2%
HIV5exp	experimental	HIV-1	5	20-20%	2-3.5%
IAV10exp	experimental	IAV	10	0.1-50%	0.1-1.1%
HIV7sim	simulated	HIV-1	7	14.3-14.3%	0.6-3%
IAV10sim	simulated	IAV	10	0.1-50%	0.1-1.1%

The datasets contain MiSeq and PacBio reads from intra-host viral populations consisting of two to ten variants each with frequencies in the range of 0.1–50%, and Hamming distances between variants in the range of 0.1–3.5%.

- v Step 5 assigns each read to a merged clique with which it shares the largest number of SNVs. Then CliqueSNV builds a consensus haplotype from all reads assigned to a single merged clique.
- vi Finally, haplotype frequencies are estimated via an expectation-maximization (EM) algorithm in Step 6.

Similarly to Predicthaplo, CliqueSNV splits the global problem into a sequel of several reconstruction tasks of increasing genome length. We start with a local reconstruction with the read-fragment long region of maximum coverage. Then the SNV graph  $G$  is enriched with links between SNVs from the same reconstructed haplotype in the processed region. We progressively increase the region currently analysed by the read-fragment length until it covers the entire haplotypes' length.

### Intra-host viral population sequencing benchmarks

We tested the ability of CliqueSNV to assemble haplotype sequences and estimate their frequencies from PacBio and MiSeq reads using four real (experimental) and two simulated datasets from HIV and IAV samples (Table 1). Each dataset contains between two to ten haplotypes with frequencies of 0.1 to 50%. The Hamming distances between pairs of variants for each dataset are shown in Supplementary Figure S1.

#### Experimental datasets.

- i. *HIV-1 subtype B plasmid mixtures and MiSeq reads (HIV2exp and HIV9exp)*. We designed nine *in silico* plasmid constructs comprising a 950-bp region of the HIV-1 subtype B polymerase (*pol*) gene that were then synthesized and cloned into pUCIDT-Amp (Integrated DNA Technologies, Skokie, IL, USA). Each clone was confirmed by Sanger sequencing. This 950-bp region at the beginning of *pol* contains known protease and reverse transcriptase genes that are monitored for drug-resistant mutations and is monitored with sequence analysis for patient care. Each of these plasmids contains a specific set of point mutations chosen using mutation profiles of patient p7 from a real clinical study (53) to create nine unique synthetic HIV-1 *pol* haplotypes. Different proportions of these plasmids were mixed and then sequenced using an Illumina MiSeq

protocol to obtain 2×300-bp reads (see Supplementary Methods). HIV2exp and HIV9exp are mixtures of two and nine variants, respectively.

- ii. *HIV-1 subtype B mixture and MiSeq reads (HIV5exp and HIV5full)*. This dataset consists of Illumina MiSeq 2×250-bp reads with an average read coverage of ~20 000× obtained from a mixture of five HIV-1 isolates: 89.6, HXB2, JRC5F, NL43 and YU2 available at (27). Isolates have pairwise Hamming distances in the range from 2 to 3.5% (27 to 46-bp differences). The original HIV-1 sequence length was 9.3 kb and the benchmark consisting of all reads forms the HIV5full benchmark. The biologically relevant beginning of *pol* with a length of 1.3 kb forms the HIV5exp benchmark.
- iii. *Influenza A mixture and PacBio reads (IAV10exp)*. This benchmark contains ten IAV virus clones that were mixed at a frequency of 0.1–50%. The Hamming distances between clones ranged from 0.1 to 1.1% (2–22-bp differences) (48). The 2 kb-amplicon was sequenced using the PacBio platform yielding a total of 33 558 reads with an average length of 1973 nucleotides.

#### Simulated datasets.

- i. *HIV-1 subtype B mixture and MiSeq reads (HIV7sim)*. This benchmark contains simulated Illumina MiSeq reads with a 10k-coverage of 1-kb *pol* sequences. The reads were simulated from seven equally distributed HIV-1 variants chosen from the NCBI database: AY835778, AY835770, AY835771, AY835777, AY835763, AY835762, and AY835757. The Hamming distances between clones are in the range from 0.6-3.0% (6 to 30-bp differences). We used SimSeq (54) for generating reads.
- ii. *Influenza A mixture and MiSeq reads (IAV10sim)*. This benchmark contains simulated IAV Illumina MiSeq reads with the same IAV haplotypes and their frequencies as for the IAV10exp benchmark. The sequencing of a 2kb-amplicon with 40k coverage with paired Illumina MiSeq reads was simulated by SimSeq (54) with the default sequencing error profile in SimSeq.
- iii. *10-strain HCV mixture C (HCV10sim)*. This is a mixture of 10 strains of HCV, Subtype 1a, with a total sequencing depth of 20000× (i.e. 400000 reads) (55). The haplotypes were obtained from true HCV genomes in the NCBI nucleotide database and have a pairwise divergence varying from 6% to 9%. Paired-end reads were simulated at relative frequencies between 5% and 13% per haplotype, i.e. a sequencing depth varies from 1000× to 4600× per haplotype.
- iv. *3- and 15-strain ZIKV mixture (ZIKV3sim and ZIKV15sim)*. ZIKV3sim consists of three master strains extracted from the NCBI nucleotide database. ZIKV15sim consists of the same three master strains and four mutants for each master strain (55). The pairwise divergence varies between 1% and 12% and the reads were simulated at relative frequencies varying from 2% to 13.3%. The total sequencing depth for this dataset is again 20000×.

### Validation metrics for viral population inference

**Precision and recall.** Inference quality is typically measured by precision and recall.

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

where  $TP$  is the number of true predicted haplotypes,  $FP$  is the number of false predicted haplotypes and  $FN$  is the number of undiscovered haplotypes.

It is also important to count precision and recall when a certain number of mismatches are allowed. Following (29) we introduced an acceptance threshold, which is the number of mismatches permitted for a predicted haplotype to count as a  $TP$ . We report the numbers  $TP$  and  $FP$  for acceptance allowing from 0 to 30 mismatches (see Figure 2).

**Matching errors between populations.** However, precision and recall do not take into account (i) distances between true and inferred viral variants as well as (ii) the frequencies of the true and inferred viral variants. Instead, we chose to use analogues of precision and recall defined for populations as follows.

Let  $T = \{(t, f_t)\}$ , be the true haplotype population, where  $f_t$  is the frequency of the true haplotype  $t$ ,  $\sum_{t \in T} f_t = 1$ . Similarly, let  $P = \{(p, f_p)\}$ , be the reconstructed haplotype population, where  $f_p$  is the frequency of the reconstructed haplotype  $p$ ,  $\sum_{p \in P} f_p = 1$ . Let  $d_{pt}$  be the distance between haplotypes  $p$  and  $t$ . Thus, instead of precision, we used the *matching error*  $E_{T \rightarrow P}$  which measures how well each reconstructed haplotype  $p \in P$  weighted by its frequency is matched by the closest true haplotype.

$$E_{T \rightarrow P} = \sum_{p \in P} f_p \min_{t \in T} d_{pt}$$

Indeed, precision increases while  $E_{T \rightarrow P}$  decreases and reaches 100% when  $E_{T \rightarrow P} = 0$ . Similarly, instead of recall, we propose to use the *matching error*  $E_{T \leftarrow P}$  which measures how well each true haplotype  $t \in T$  weighted by its frequency is matched by the closest reconstructed haplotype (56).

$$E_{T \leftarrow P} = \sum_{t \in T} f_t \min_{p \in P} d_{pt}$$

Note that recall increases while  $E_{T \leftarrow P}$  decreases and reaches 100% when  $E_{T \leftarrow P} = 0$ .

**Earth mover's distance (EMD) between populations.** The matching errors described above match haplotypes of true and reconstructed populations but do not match their frequencies. In order to simultaneously match haplotype sequences and their frequencies, we allowed for a fractional matching when portions of a single haplotype  $p$  of population  $P$  are matched to portions of possibly several haplotypes of  $T$  and *vice versa*. Thus, we separated  $f_p$  into  $f_{pt}$ 's each denoting portion of  $p$  matched to  $t$  such that  $f_p = \sum_{t \in T} f_{pt}$ ,  $f_{pt} \geq 0$ . Symmetrically,  $f_t$ 's are also separated into  $f_{pt}$ 's, i.e.

$\sum_{p \in P} f_{pt} = f_t$ . Finally, we chose  $f_{pt}$ 's minimizing the total error of matching  $T$  to  $P$  which is also known as Wasserstein metric or the EMD between  $T$  and  $P$  (57,58).

$$\text{EMD}(T, P) = \min_{f_{pt} > 0} \sum_{t \in T} \sum_{p \in P} f_{pt} d_{pt}$$

$$\text{s.t. } \sum_{t \in T} f_{pt} = f_p, \text{ and } \sum_{p \in P} f_{pt} = f_t$$

EMD is efficiently computed as an instance of the transportation problem using network flows.

EMDs can vary a lot over different benchmarks since they may have different complexities, which depends on the number of true variants, the frequency distribution, the similarity between haplotypes, sequencing depth, sequencing error rate, and many other parameters. Hence, we measured the complexity of a benchmark as the EMD between the true population and a population consisting of a single consensus haplotype (59).

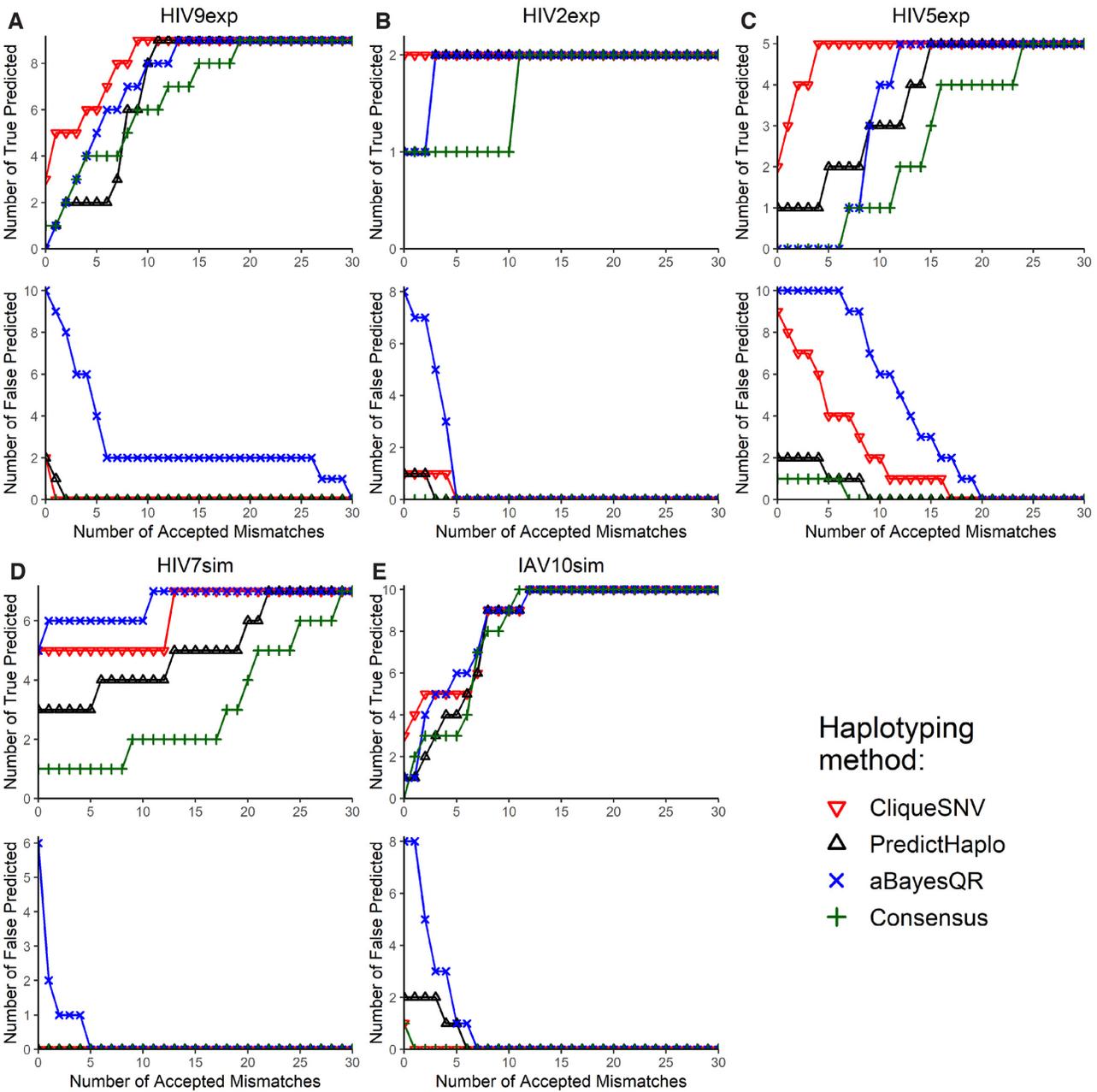
### CliqueSNV algorithm details

Data input for CliqueSNV consists of PacBio or Illumina reads from an intra-host viral population aligned to a reference genome. A deletion is treated as a special nucleotide, i.e. it is just one of five possible alleles. Insertions are identified after alignment and should be used to extend the reference. Output is the set of inferred viral variant RNA sequences with their frequencies. The formal high-level pseudocode of the CliqueSNV algorithm is described in the supplementary materials. Below we describe in detail the six major steps of CliqueSNV that are schematically presented in Figure 1.

**Step 1: Finding linked and forbidden SNV pairs.** At a given genomic position  $I$ , the most frequent nucleotide is referred to as a *major variant* and is denoted 1. Let us fix one of the less frequent nucleotide (referred to as a *minor variant*) and denote it 2. A pair of variants at two distinct genomic positions  $I$  and  $J$  is referred to as a 2-haplotype. Let  $O_{22}$  be the observed count of the 2-haplotype (22) in the reads covering positions  $I$  and  $J$ . In this step, CliqueSNV tries to decide whether the observed  $O_{22}$  reads are sequencing errors or they are produced by an existing haplotype containing the 2-haplotype (22).

The pairs of minor variants (referred to as SNV pairs) are classified into three categories: linked, forbidden, and unclassified. An SNV pair is *linked* if it is extremely unlikely that there is no sufficiently frequent haplotype containing both minor variants is very low. On the other side, an SNV pair is *forbidden* if it is extremely unlikely that the corresponding minor variants belong to the same haplotype of sufficient frequency. All other SNV pairs are referred to as *unclassified*.

Below we estimate the probability of observing at least  $x \geq O_{22}$  reads given that the true frequency  $T_{22}$  of the 2-haplotype (22) is at most  $t$  (by default  $t = 0.001$ ). This probability should be low enough so that false positive linked pairs would be virtually impossible, i.e. we require that the expected number of false positive linked pairs be  $< 0.05$ . Therefore, this probability should be less than  $0.05/\binom{L}{2}$ ,



**Figure 2.** The number of true and false predicted haplotypes depending on the number of accepted mismatches for five benchmarks: (A) HIV9exp; (B) HIV2exp; (C) HIV5exp; (D) HIV7sim; (E) IAV10sim. Two haplotypes are regarded identical if the Hamming distance between them is at most the number of accepted mismatches.

where  $L$  is the haplotype length.

$$\begin{aligned}
 Pr(x \geq O_{22} | T_{22} \leq t) &= 1 - Pr(x < O_{22} | T_{22} \leq t) \\
 &\leq 1 - \sum_{i=0}^{O_{22}-1} \binom{n}{i} t^i (1-t)^{n-i} \\
 &\leq \frac{0.05}{\binom{L}{2}} \tag{1}
 \end{aligned}$$

Pairs of SNVs passing this linkage test (1) are classified as a *linked SNV pairs*.

For every other pair of SNVs, we check whether they can be classified as a *forbidden SNV pair*, i.e. whether the probability of observing at most  $x \leq O_{22}$  reads is low enough ( $< 0.05$ ) given that the 2-haplotype (22) has frequency  $T_{22} \geq t$  (by default  $t = 0.001$ ). Similarly, we require that the expected number of false positive forbidden pairs be  $< 0.05$ .

$$\begin{aligned}
 Pr(x \leq O_{22} | T_{22} \geq t) &\leq \sum_{i=0}^{O_{22}} \binom{n}{i} t^i (1-t)^{n-i} \\
 &\leq \frac{0.05}{\binom{L}{2}} \tag{2}
 \end{aligned}$$

Pairs of SNVs passing this linkage test (2) are classified as a *forbidden* SNV pairs.

**Step 2: Constructing the SNV graph.** The SNV graph  $G = (V, E)$  consists of vertices corresponding to minor variants and edges corresponding to linked pairs of minor variants from different positions. If the intra-host population consists of very similar haplotypes, then the number of true SNVs corresponding to non-isolated vertices in the graph  $G$  is very small which makes the number of edges very small as well. Indeed, the PacBio dataset for IAV encompassing 2500 positions is split into 10 000 vertices, while the SNV graph contains only 700 edges, and, similarly, the simulated Illumina read dataset for the same haplotypes contains only 368 edges.

Note that the isolated minor variants correspond to genotyping errors unless they have a significant frequency. This fact allows us to estimate the number of errors per read, assuming that all isolated SNVs are errors. As expected, the distribution of the PacBio reads has a heavy tail (see Supplementary Figure S4), which implies that most reads are (almost) error free, while a small number of heavy-tail reads accumulate most of the errors. The heavy-tail reads are difficult to correct, and they are uniformly distributed across haplotypes. Therefore, their removal drastically reduces number of errors but does not affect frequency of rare haplotypes. Our analysis allows the identification of such reads, which can then be filtered out. By default, we filter out  $\approx 10\%$  of PacBio reads, but we do not filter out any Illumina reads. The SNV graph is then constructed for the reduced set of reads. Such filtering allows the reduction of systematic errors and refines the SNV graph significantly.

**Step 3: Finding cliques in the SNV graph  $G$ .** Although the MAX CLIQUE is a well-known NP-complete problem and there may be an exponential number of maximal cliques in  $G$ , a standard Bron–Kerbosch algorithm requires little computational time since  $G$  is very sparse (60).

**Step 4: Merging cliques in the clique graph  $C_G$ .** The clique graph  $C_G = (C, F, L)$  consists of vertices corresponding to cliques in the SNV graph  $G$  and two sets of edges  $F$  and  $L$ . A *forbidding edge*  $(p, q) \in F$  connects two cliques  $p$  and  $q$  with at least one forbidden pair of minor variants from  $p$  and  $q$  respectively. A *linking edge*  $(p, q) \in L$  connects two cliques  $p$  and  $q$ ,  $(p, q) \notin F$ , with at least one linked pair of minor variants from  $p$  and  $q$ , respectively. Any true haplotype corresponds to a maximal  $(L \setminus F)$ -connected subgraph  $H$  of  $C_G$  which is connected with edges from  $L$  and does not contain any edge from  $F$  [see Figure 1 (4)].

Unfortunately, even deciding whether there is a  $L$ -path between  $p$  and  $q$  avoiding forbidding edges is known to be NP-hard (61). We find all subgraphs  $H$  as follows (see Supplementary Figure S5): (i) connect all pairs of vertices except connected with forbidding edges, (ii) find all maximal super-cliques in the resulted graph  $C'_G = (C, C^{(2)} - F)$  using (60), (iii) split each super-clique into  $L$ -connected components and (iv) output maximal  $L$ -connected components.

**Step 5: Partitioning reads between merged cliques and finding consensus haplotypes.** Let  $S$  be the set of all positions containing at least one minor variant in  $V$ . Let  $q_S$  be a *major clique* corresponding to a haplotype with all major variants in  $S$ . The distance between a read  $r$  and a clique  $q$  equals the number of variants in  $q$  that are different from

the corresponding nucleotides in  $r$ . Each read  $r$  is assigned to the closest clique  $q$  (which can possibly be  $q_S$ ). In case of a tie, the read  $r$  contributes only  $1/n$  frequency to consensus, where  $n$  is the number of closest cliques. Usually, the number of closest cliques is 1, although in a case of IAV, when most cliques share the same position, a significant portion of reads is assigned to multiple cliques (see Supplementary materials S6–8).

Finally, for each clique  $q$ , CliqueSNV finds the consensus  $v(q)$  of all reads assigned to  $q$ . Then  $v(q)$  is extended from  $S$  to a full-length haplotype by setting all non- $S$  positions to major SNVs.

**Step 6: Estimating haplotype frequencies by using the EM algorithm.** CliqueSNV estimates the frequencies of the assembled intra-host haplotypes via an EM algorithm similar to the one used in IsoEM (62). Let  $K$  be the number of assembled viral variants, and let  $\alpha$  be the probability of sequencing error. EM algorithm works as follows:

- i. Initialize frequencies of viral variants  $f_j^{(0)} \leftarrow \frac{1}{K}$ ,  
Compute the probability of  $l_i$ -long read  $r_i$   $i = 1, \overline{N}$ , being emitted by viral variant  $j = 1, \overline{K}$ ,  
$$h_{ji} = \prod_{l=1}^{l_i} ((1 - \alpha)M_{ji,l} + \frac{\alpha}{3}(1 - M_{ji,l})),$$
where  $M_{ji,l}$  - indicator if  $i$ -th read coincides with  $j$ -th viral variant in the position  $l$
- ii. (Expectation) Update the amount of read  $r_i$  emitted by the  $j$ th viral variant  $p_{ij} \leftarrow \frac{f_j^{(n-1)}h_{ji}}{\sum_{u=1}^k f_u^{(n-1)}h_{ui}}$
- iii. (Maximization) Update the frequency of the  $j$ th viral variant  $f_j^{(n)} \leftarrow \frac{\sum_{i=1}^N p_{ij}}{\sum_{u=1}^k \sum_{i=1}^N p_{iu}}$
- iv. if  $\|f_j^{(n-1)} - f_j^{(n)}\| > \varepsilon$ , then  $n \leftarrow n + 1$  and go to step 2
- v. Output estimated frequencies  $f_j^{(n)}$

## RESULTS

### Performance of haplotyping methods

We compared CliqueSNV to the 2SNV, PredictHaplo and aBayesQR haplotyping methods. Since CliqueSNV, PredictHaplo and aBayesQR use Illumina reads, we compared them using the HIV9exp, HIV2exp, HIV5exp, HIV7sim and IAV10sim datasets. Since CliqueSNV, 2SNV and PredictHaplo can also use PacBio reads, we compared them using the IAV10exp dataset. We also used consensus sequences in the comparisons (59) because of its simplicity and to evaluate sequences most similar to those generated by the Sanger sequencing method (63). Finally, we validated scalability of CliqueSNV, PredictHaplo and aBayesQR with respect to the length of the genomic region using benchmarks from Table 2.

The precision and recall of haplotype discovery for each method is provided in Table 3. CliqueSNV had the best precision and recall for five of the six datasets. For the HIV5exp dataset, PredictHaplo was more conservative and predicted less false positive variants (better precision) than CliqueSNV.

Following study (29), we also showed how precision and recall grew with the reduction of restriction on mismatches (Figure 2). The number of true predicted haplotypes for

**Table 2.** Full-length viral genome datasets

Name	Type	Virus	#haplotypes	Haplotype frequencies	Hamming distance	Region length
HCV10sim	simulated	HCV	10	5–13%	6–9%	8992
ZIKV3sim	simulated	ZIKV	3	16–60%	3–10%	9929
ZIKV15sim	simulated	ZIKV	15	2–13%	1–12%	9929
HIV5full	experimental	HIV-1	5	20–20%	1–6%	9275

Three simulated sequencing datasets HCV10sim, ZIKV3sim, ZIKV15sim and one experimental sequencing dataset HIV5full.

**Table 3.** Prediction statistics of haplotype reconstruction methods using experimental and simulated (a) MiSeq and (b) PacBio datasets

(a)						
Benchmark	CliquesNV		aBayesQR		PredictHaplo	
	Precision	Recall	Precision	Recall	Precision	Recall
HIV9exp	<b>0.60</b>	<b>0.33</b>	0.00	0.00	0.00	0.00
HIV2exp	<b>0.66</b>	<b>1.00</b>	0.11	0.50	0.50	0.50
HIV5exp	0.18	<b>0.40</b>	0.00	0.00	<b>0.33</b>	0.20
HIV7sim	<b>1.00</b>	<b>0.71</b>	<b>1.00</b>	0.42	0.45	<b>0.71</b>
IAV10sim	<b>0.75</b>	<b>0.30</b>	0.11	0.10	0.33	0.10
(b)						
Benchmark	CliquesNV		2SNV		PredictHaplo	
	Precision	Recall	Precision	Recall	Precision	Recall
IAV10exp	<b>1.00</b>	<b>1.00</b>	0.82	0.90	0.70	0.70

The precision and recall was evaluated stringently such that if a predicted haplotype has at least one mismatch to its closest answer, then that haplotype is scored as a false positive.

CliquesNV was always greater than that of the other methods on real experimental sequencing benchmarks indicating that CliquesNV more accurately identified the true haplotypes. The number of falsely predicted haplotypes for CliquesNV was always lower than those for aBayesQR, but similar to those predicted by PredictHaplo on four out of five datasets indicating that both CliquesNV and PredictHaplo had the best precision with MiSeq datasets.

Matching distance analysis showed that matching distances  $E_{T \leftarrow P}$  and  $E_{T \rightarrow P}$  are better for CliquesNV than for both PredictHaplo and aBayesQR on four out of five MiSeq datasets (Figure 3). For HIV7sim,  $E_{T \leftarrow P}$  for aBayesQR was slightly better than for CliquesNV. Using HIV9exp, HIV2exp, HIV7sim and IAV10sim datasets, the  $E_{T \leftarrow P}$  and  $E_{T \rightarrow P}$  for CliquesNV were very close to zero indicating that the predictions were almost perfect. Since  $E_{T \leftarrow P}$  and  $E_{T \rightarrow P}$  correlate with precision and recall, matching distance analysis indicates that CliquesNV had a better precision, and significantly outperformed both PredictHaplo and aBayesQR. Since aBayesQR had a higher  $E_{T \rightarrow P}$  on MiSeq datasets, it is more likely to make more false predictions. Notably, on the HIV7sim dataset, aBayesQR outperformed both CliquesNV and PredictHaplo by  $E_{T \leftarrow P}$ .

The EMD between the predicted and true haplotype populations for all five MiSeq datasets are shown in Figure 4. The exact EMD values are provided in Table 4. CliquesNV provided the lowest (the best) EMD across all tools on four out of five MiSeq benchmarks. For the simulated and PacBio datasets, CliquesNV had almost a zero EMD indicating a low error in predictions. PredictHaplo had a lower EMD than aBayesQR on four out of five MiSeq datasets. aBayesQR has almost a zero EMD with the HIV7sim dataset and outperformed CliquesNV, while us-

ing the HIV5exp dataset, aBayesQR performed poorer than other methods.

Next, CliquesNV, 2SNV and PredictHaplo were compared using the IAV10exp benchmark dataset (see Supplementary Table S1). CliquesNV correctly recovered all ten true variants, including the haplotype with frequencies significantly below the sequencing error rate. 2SNV recovered nine true variants but found one false positive. PredictHaplo recovered only seven true variants and falsely predicted three variants. To further explore the precision of these three methods with the IAV10exp data, we simulated low-coverage datasets by randomly subsampling  $n = 16K, 8K, 4K$  reads from the original data. For each dataset, CliquesNV found at least one true variant more than both 2SNV and PredictHaplo.

Finally, Table 5 reports the performance of three methods on full-length genomes. We normalize EMD over the genomic length so that the resulted EMD are in the same range and can be compared for different genomic regions. On average, CliquesNV for all lower bounds on frequency (2, 5 and 10%) outperforms PredictHaplo, but for two out of four full-length benchmarks PredictHaplo is more accurate than CliquesNV.

### Runtime comparison

To compare the computational run time of each method, we used the same PC (Intel(R) Xeon(R) CPU X5550 2.67GHz  $\times 2$  8 cores per CPU, DIMM DDR3 1333 MHz RAM 4Gb  $\times 12$ ) with the CentOS 6.4 operating system. The runtime of CliquesNV is sublinear with respect to the number of reads while the runtime of PredictHaplo and 2SNV exhibit super-linear growth. For the 33k IAV10sim reads

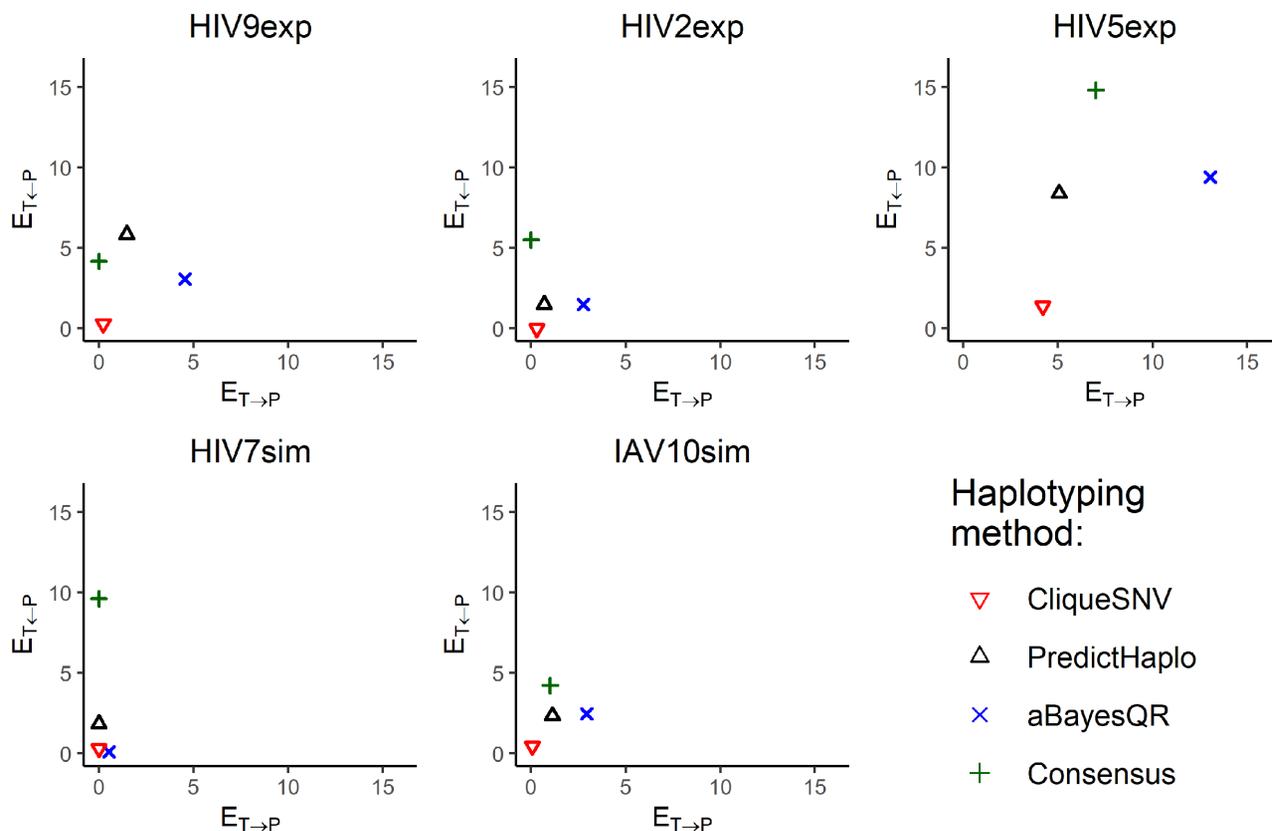


Figure 3. Matching distances  $E_{T \leftarrow P}$  and  $E_{T \rightarrow P}$  between the true haplotype population  $T$  and the reconstructed haplotype population  $P$  for five benchmarks.

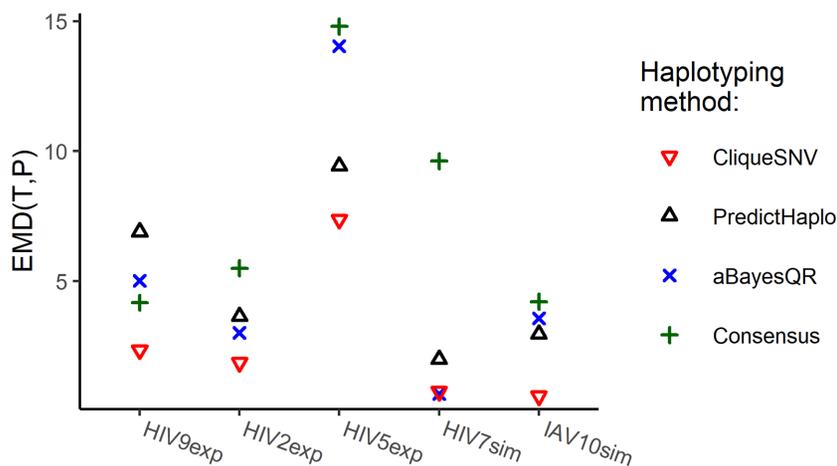


Figure 4. EMD between true and reconstructed haplotype populations for five benchmarks.

the CliqueSNV analysis took 21 s, while PredictHaplo and 2SNV took around 30 min. The runtime of CliqueSNV is quadratic with respect to the number of SNVs rather than by the length of the sequencing region (Supplementary Figure S2).

We also generated five HIV-1 variants within 1% Hamming distance from each other, which is the estimated genetic distance between related HIV variants from the same person (64). Then we simulated 1M Illumina reads for se-

quence regions of length 566, 1132, 2263 and 9181 nucleotides for which CliqueSNV required 37, 144, 227 and 614 seconds, respectively, for analyzing these datasets (Supplementary Figure S3). For the HIV2exp benchmark, aBayesQR, PredictHaplo and CliqueSNV required over 10 h, 24 min and only 79 s, respectively. Finally, for four full-length benchmarks, the CliqueSNV is 2–10× faster than PredictHaplo depending on lower bound on the haplotype frequency, while ABayesQR is more than 100× slower and

**Table 4.** EMD from predicted haplotypes to the true haplotype population and haplotyping method improvement

(a)							
Benchmark	Consensus	CliquesNV		aBayesQR		PredictHaplo	
	EMD	EMD	Impr.	EMD	Impr.	EMD	Impr.
HIV9exp	4.18	<b>2.35</b>	<b>1.78</b>	5.02	0.83	6.90	0.61
HIV2exp	5.50	<b>1.87</b>	<b>2.94</b>	3.02	1.82	3.65	1.51
HIV5exp	14.80	<b>7.37</b>	<b>2.01</b>	14.05	1.05	9.43	1.57
HIV7sim	9.63	0.76	12.72	<b>0.67</b>	<b>14.4</b>	2.00	4.80
IAV10sim	4.22	<b>0.59</b>	<b>7.2</b>	3.57	1.18	2.97	1.42
(b)							
Benchmark	Consensus	CliquesNV		2SNV		PredictHaplo	
	EMD	EMD	Impr.	EMD	Impr.	EMD	Impr.
IAV10exp	4.22	<b>0.22</b>	<b>19.18</b>	0.23	18.35	0.38	11.12

Four haplotyping methods (aBayesQR, CliquesNV, Consensus, PredictHaplo) are benchmarked using five MiSeq (A) and one PacBio datasets (B). The column Impr. (improvement) shows how much better is prediction of haplotyping method over inferred consensus, and it is calculated as  $\frac{EMD_m}{EMD_c}$  where  $EMD_c$  is an EMD for consensus, and  $EMD_m$  is an EMD for method.

**Table 5.** EMD from predicted haplotypes to the true haplotype population for 1K-, 2K-, 5K-long and full-length genomic regions for HIV, HCV and ZIKV benchmarks

Benchmark	Length	Consensus	CliquesNV			PredictHaplo	aBayesQR
			2%	5%	10%		
HCV10sim	1K	1.109	0.080	<b>0.037</b>	0.143	1.479	did not finish
	2K	3.397	<b>0.714</b>	0.854	1.923	0.981	did not finish
	5K	3.719	3.729	3.729	1.431	<b>0.450</b>	did not finish
	full-length	3.335	3.335	3.335	1.217	<b>0.140</b>	did not finish
ZIKV3sim	1K	2.408	<b>0.015</b>	<b>0.015</b>	<b>0.015</b>	0.100	0.086
	2K	2.388	<b>0.020</b>	<b>0.020</b>	<b>0.020</b>	0.092	did not finish
	5K	2.182	<b>0.010</b>	<b>0.010</b>	<b>0.010</b>	0.031	did not finish
	full-length	2.168	<b>0.007</b>	<b>0.007</b>	<b>0.007</b>	0.016	did not finish
ZIKV15sim	1K	3.368	<b>1.123</b>	1.335	1.840	3.640	1.345
	2K	3.571	1.469	<b>1.339</b>	3.128	7.037	1.807
	5K	3.506	<b>1.379</b>	1.778	1.698	6.701	did not finish
	full-length	3.803	<b>1.254</b>	1.643	2.801	6.014	did not finish
HIV5full	1K	1.480	0.696	<b>0.350</b>	0.590	0.943	1.405
	2K	3.460	<b>2.081</b>	2.100	2.626	2.825	3.100
	5K	2.384	2.219	<b>1.500</b>	1.997	2.267	did not finish
	full-length	3.189	2.968	2.557	2.273	<b>1.811</b>	did not finish
Average over all benchmarks		2.957	1.401	<b>1.371</b>	1.438	2.203	1.548*

The EMD is normalized by the genomic region length, i.e. we report the EMD per 100 genomic positions. Comparison is performed for CliquesNV with three different frequency thresholds: 2, 5 (default) and 10%, aBayesQR and PredictHaplo. The aBayesQR did not finish in majority of samples. The best EMD for each benchmark is in bold font.

manged to handle only short genomic regions (see Supplementary Table S2).

## DISCUSSION

Assembly of haplotype populations from noisy NGS data is one of the most challenging problems of computational genomics. High-throughput sequencing technologies, such as Illumina MiSeq and HiSeq, provide deep sequence coverage that allows discovery of rare, clinically relevant haplotypes. However, the short reads generated by the Illumina technology require assembly that is complicated by sequencing errors, an unknown number of haplotypes in a sample, and the genetic similarity of haplotypes within a sample. Furthermore, the frequency of sequencing errors in Illumina reads is comparable to the frequencies of true minor mutations (41). The recent development of single-molecule sequencing platforms such as PacBio produce reads that are

sufficiently long to span entire genes or small viral genomes. Nonetheless, the error rate of single-molecule sequencing is exceptionally high reaching 13 – 14% (65), which hampers PacBio sequencing to detect and assemble rare viral variants.

We developed CliquesNV, a new reference-based assembly method for reconstruction of rare genetically related viral variants such as those observed during infection with rapidly evolving RNA viruses like HIV, HCV and IAV. We demonstrated that CliquesNV infers accurate haplotyping in the presence of high sequencing error rates and is also suitable for both single-molecule and short-read sequencing. In contrast to other haplotyping methods, CliquesNV infers viral haplotypes by detection of clusters of statistically linked SNVs rather than through assembly of overlapping reads used with methods such as Savage (26).

Applied to the novel in vitro sequencing HIV-1 benchmark, CliquesNV correctly reconstructed 87% of the intra-

host haplotype population. At the same time, other state-of-the-art tools were not able to recover even a single haplotype without errors. Additionally, we have used the only previously known and commonly used *in vitro* benchmark (27) and simulated datasets to evaluate the accuracy of existing haplotyping methods. In contrast to the existing methods, CliqueSNV was able to detect minority haplotypes at a low 0.1% frequency and distinguish minority haplotypes differently in only two base pairs. We have also validated CliqueSNV on Illumina reads from full-length genomes where it is faster and more accurate than PredictHaplo on average.

We also believe that CliqueSNV can be applied to Nonopore sequencing data. We ran CliqueSNV on 4k-long spike gene region for 9 samples of SARS-CoV-2 with 300–900 bp read length. In each sample, CliqueSNV managed to reconstruct two to six haplotypes which is within the number of haplotypes identified for the same region from Illumina reads for similar SARS-CoV-2 samples.

Although very accurate and fast, CliqueSNV has some limitations. Unlike Savage (26), CliqueSNV is not a *de novo* assembly tool and requires a reference viral genome. This obstacle could easily be addressed by using Vicuna (59) or other analogous tools to first assemble a consensus sequence from the NGS reads, which can then be used as a reference. Another limitation is for variants that differ only by isolated SNVs separated by long conserved genomic regions longer than the read length which may not be accurately inferred by CliqueSNV. While such situations usually do not occur for viruses, where mutations are typically densely concentrated in different genomic regions, we plan to address this limitation in the next version of CliqueSNV.

The ability to accurately infer the structure of intra-host viral populations makes CliqueSNV applicable for studying viral evolution, transmission and examining the genomic compositions of RNA viruses. In addition, we envision that the application of our method could be extended to other highly heterogeneous genomic populations, such as metagenomes, immune repertoires and cancer cell genes.

## DATA AVAILABILITY

The datasets HIV2exp and HIV9exp have been deposited in the Sequence Read Archive under accession number SRR12042289 and SRR12042290, respectively.

The links to the datasets and the consensus sequences of the individual strains are available at [https://github.com/Sergey-Knyazev/CliqueSNV-validation/blob/master/relevant\\_haplotypes/HIV9exp.fasta](https://github.com/Sergey-Knyazev/CliqueSNV-validation/blob/master/relevant_haplotypes/HIV9exp.fasta)

## Software

CliqueSNV is available at <https://github.com/vtsyvina/CliqueSNV>

## Validation scripts

Validation notebook is available through docker container that contains all instructions, related dependencies and input datasets at [https://hub.docker.com/r/amelnyk3/clique\\_snv\\_validation/](https://hub.docker.com/r/amelnyk3/clique_snv_validation/).

## SUPPLEMENTARY DATA

Supplementary Data are available at NAR Online.

## ACKNOWLEDGEMENTS

Use of trade names is for identification only and does not imply endorsement by the U.S. Department of Health and Human Services, the Public Health Service or the Centers for Disease Control and Prevention (CDC).

## FUNDING

NIH [1R01EB025022-01, in part to A.Z., P.S.]; NSF [CCF-1619110, in part to A.Z., CCF-2047828 to P.S., DBI-2041984, in part to S.M.]; Ministry of Science and Higher Education of the Russian Federation Grant [075-15-2020-926 to Y.P.]; Molecular Basis of Disease at Georgia State University (in part) (to S.K., V.T., A.M.). Funding for open access charge: NSF [DBI-2041984].

*Conflict of interest statement.* None declared.

## REFERENCES

- Kilmarx, P.H. (2009) Global epidemiology of HIV. *Curr. Opin. HIV AIDS*, **4**, 240–246.
- Hajarizadeh, B., Grebely, J. and Dore, G.J. (2013) Epidemiology and natural history of HCV infection. *Nat. Rev. Gastro. Hepat.*, **10**, 553–562.
- Lozano, R., Naghavi, M., Foreman, K., Lim, S., Shibuya, K., Aboyans, V., Abraham, J., Adair, T., Aggarwal, R., Ahn, S.Y. *et al.* (2012) Global and regional mortality from 235 causes of death for 20 age groups in 1990 and 2010: a systematic analysis for the Global Burden of Disease Study 2010. *Lancet*, **380**, 2095–2128.
- Eigen, M., McCaskill, J. and Schuster, P. (1989) The molecular quasi-species. *Adv. Chem. Phys.*, **75**, 149–263.
- Martell, M., Esteban, J., Quer, J., Genesca, J., Weiner, A., Esteban, R., Guardia, J. and Gomez, J. (1992) Hepatitis C virus (HCV) circulates as a population of different but closely related genomes: quasispecies nature of HCV genome distribution. *J. Virol.*, **66**, 3225–3229.
- Steinhauer, D. and Holland, J. (1987) Rapid evolution of RNA viruses. *Annu. Rev. Microbiol.*, **41**, 409–433.
- Domingo, E., Sheldon, J. and Perales, C. (2012) Viral quasispecies evolution. *Microbiol. Mol. Biol. R.*, **76**, 159–216.
- Rodriguez-Frias, F., Buti, M., Taberner, D. and Homs, M. (2013) Quasispecies structure, cornerstone of hepatitis B virus infection: mass sequencing approach. *World J. Gastroenterol.*, **19**, 6995–7023.
- Xu, D., Zhang, Z. and Wang, F.-S. (2004) SARS-associated coronavirus quasispecies in individual patients. *N. Engl. J. Med.*, **350**, 1366–1367.
- Shen, Z., Xiao, Y., Kang, L., Ma, W., Shi, L., Zhang, L., Zhou, Z., Yang, J., Zhong, J., Yang, D. *et al.* (2020) Genomic diversity of severe acute respiratory syndrome–coronavirus 2 in patients with coronavirus disease 2019. *Clin. Infect. Dis.*, **71**, 713–720.
- Beerenwinkel, N., Sing, T., Lengauer, T., Rahnenfuehrer, J. and Roomp, K. (2005) Computational methods for the design of effective therapies against drug resistant HIV strains. *Bioinformatics*, **21**, 3943–3950.
- Douek, D.C. and Kwong, P.D. (2006) The rational design of an AIDS vaccine. *Cell*, **124**, 677–681.
- Gaschen, B., Taylor, J., Yusim, K., Foley, B. and Gao, F. (2002) Diversity considerations in HIV-1 vaccine selection. *Science*, **296**, 2354–2360.
- Holland, J., De La Torre, J. and Steinhauer, D. (1992) RNA virus populations as quasispecies. *Curr. Top. Microbiol. Immunol.*, **176**, 1–20.
- Rhee, S.-Y., Liu, T., Holmes, S. and Shafer, R. (2007) HIV-1 subtype B protease and reverse transcriptase amino acid covariation. *PLoS Comput. Biol.*, **3**, e87.
- Campo, D.S., Skums, P., Dimitrova, Z., Vaughan, G., Forbi, J.C., Teo, C.-G., Khudyakov, Y. and Lau, D.T. (2014) Drug resistance of a

- viral population and its individual intrahost variants during the first 48 hours of therapy. *Clin. Pharmacol. Ther.*, **95**, 627–635.
17. Skums,P., Bunimovich,L. and Khudyakov,Y. (2015) Antigenic cooperation among intrahost HCV variants organized into a complex network of cross-immunoreactivity. *Proc. Natl Acad. Sci. U.S.A.*, **112**, 6653–6658.
  18. Campo,D.S., Xia,G.-L., Dimitrova,Z., Lin,Y., Forbi,J.C., Ganova-Raeva,L., Punkova,L., Ramachandran,S., Thai,H., Skums,P. et al. (2016) Accurate genetic detection of hepatitis C virus transmissions in outbreak settings. *J. Infect. Dis.*, **213**, 957–965.
  19. Glebova,O., Knyazev,S., Melnyk,A., Artyomenko,A., Khudyakov,Y., Zelikovsky,A. and Skums,P. (2017) Inference of genetic relatedness between viral quasispecies from sequencing data. *BMC Genomics*, **18**(Suppl.10), 918.
  20. Skums,P., Zelikovsky,A., Singh,R., Gussler,W., Dimitrova,Z., Knyazev,S., Mandric,I., Ramachandran,S., Campo,D., Jha,D. et al. (2017) QUENTIN: reconstruction of disease transmissions from viral quasispecies genomic data. *Bioinformatics*, **34**, 163–170.
  21. Wymant,C., Hall,M., Ratmann,O., Bonsall,D., Golubchik,T., de Cesare,M., Gall,A., Cornelissen,M. and Fraser,C. (2017) PHYLOSCANNER: inferring transmission from within- and between-host pathogen genetic diversity. *Mol. Biol. Evol.*, **35**, 719–733.
  22. Melnyk,A., Knyazev,S., Vannberg,F., Bunimovich,L., Skums,P. and Zelikovsky,A. (2020) Using Earth mover's distance for viral outbreak investigations. *BMC Genomics*, **21**(Suppl. 5), 582.
  23. Boskova,V. and Stadler,T. (2020) PIQMEE: Bayesian phylodynamic method for analysis of large datasets with duplicate sequences. *Mol. Biol. Evol.*, **37**, 3061–3075.
  24. Icer Baykal,P.B., Lara,J., Khudyakov,Y., Zelikovsky,A. and Skums,P. (2020) Quantitative differences between intra-host HCV populations from persons with recently established and persistent infections. *Virus Evol.*, **7**, veaa103.
  25. Döring,M., Büch,J., Friedrich,G., Pironti,A., Kalaghatgi,P., Knops,E., Heger,E., Obermeier,M., Däumer,M., Thielen,A. et al. (2018) geno2pheno[ngs-freq]: a genotypic interpretation system for identifying viral drug resistance using next-generation sequencing data. *Nucleic Acids Res.*, **46**, W271–W277.
  26. Baaijens,J.A., El Aabidine,A.Z., Rivals,E. and Schönhuth,A. (2017) *De novo* assembly of viral quasispecies using overlap graphs. *Genome Res.*, **27**, 835–848.
  27. Giallonardo,F.D., Töpfer,A., Rey,M., Prabhakaran,S., Duport,Y., Leemann,C., Schmutz,S., Campbell,N.K., Joos,B., Lecca,M.R. et al. (2014) Full-length haplotype reconstruction to infer the structure of heterogeneous virus populations. *Nucleic Acids Res.*, **42**, e115.
  28. Knyazev,S., Hughes,L., Skums,P. and Zelikovsky,A. (2021) Epidemiological data analysis of viral quasispecies in the next-generation sequencing era. *Brief. Bioinform.*, **22**, 96–108.
  29. Prabhakaran,S., Rey,M., Zagordi,O., Beerenwinkel,N. and Roth,V. (2014) HIV haplotype inference using a propagating Dirichlet process mixture model. *IEEE/ACM Trans. Comput. Biol. Bioinform.*, **11**, 182–191.
  30. Ahn,S. and Vikalo,H. (2018) aBayesQR: a Bayesian method for reconstruction of viral populations characterized by low diversity. *J. Comput. Biol.*, **25**, 637–648.
  31. Töpfer,A., Zagordi,O., Prabhakaran,S., Roth,V., Halperin,E. and Beerenwinkel,N. (2013) Probabilistic inference of viral quasispecies subject to recombination. *J. Comput. Biol.*, **20**, 113–123.
  32. Töpfer,A., Marschall,T., Bull,R.A., Luciani,F., Schönhuth,A. and Beerenwinkel,N. (2014) Viral quasispecies assembly via maximal clique enumeration. *PLoS Comput. Biol.*, **10**, e1003515.
  33. Mangul,S., Wu,N.C., Mancuso,N., Zelikovsky,A., Sun,R. and Eskin,E. (2014) Accurate viral population assembly from ultra-deep sequencing data. *Bioinformatics*, **30**, i329–i337.
  34. Skums,P., Mancuso,N., Artyomenko,A., Tork,B., Mandoiu,I., Khudyakov,Y. and Zelikovsky,A. (2013) Reconstruction of viral population structure from next-generation sequencing data using multicommodity flows. *BMC Bioinformatics*, **14**, S2.
  35. Mancuso,N., Tork,B., Skums,P., Ganova-Raeva,L., Mandoiu,I. and Zelikovsky,A. (2011) Reconstructing viral quasispecies from NGS amplicon reads. *In Silico Biol.*, **11**, 237–249.
  36. Zagordi,O., Bhattacharya,A., Eriksson,N. and Beerenwinkel,N. (2011) ShoRAH: estimating the genetic diversity of a mixed sample from next-generation sequencing data. *BMC Bioinformatics*, **12**, 119.
  37. Astrovskaya,I., Tork,B., Mangul,S., Westbrooks,K., Mandoiu,I., Balfe,P. and Zelikovsky,A. (2011) Inferring viral quasispecies spectra from 454 pyrosequencing reads. *BMC Bioinformatics*, **12**(Suppl.6), S1.
  38. Prosperi,M.C. and Salemi,M. (2012) QuRe: software for viral quasispecies reconstruction from next-generation sequencing data. *Bioinformatics*, **28**, 132–133.
  39. Töpfer,A., Zagordi,O., Prabhakaran,S., Roth,V., Halperin,E. and Beerenwinkel,N. (2013) Probabilistic inference of viral quasispecies subject to recombination. *J. Comput. Biol.*, **20**, 113–123.
  40. Skums,P., Dimitrova,Z., Campo,D.S., Vaughan,G., Rossi,L., Forbi,J.C., Yokosawa,J., Zelikovsky,A. and Khudyakov,Y. (2012) Efficient error correction for next-generation sequencing of viral amplicons. *BMC Bioinformatics*, **13**(Suppl.10), S6.
  41. Skums,P., Artyomenko,A., Glebova,O., Campo,D.S., Dimitrova,Z., Zelikovsky,A. and Khudyakov,Y. (2016) Error correction of NGS reads from viral populations. In: Mandoiu,I.I. and Zelikovsky,A. (eds). *Computational Methods for Next Generation Sequencing Data Analysis*. Wiley, 331–354.
  42. Barik,S., Das,S. and Vikalo,H. (2018) Viral quasispecies reconstruction via correlation clustering. *Genomics*, **110**, 375–381.
  43. Westbrooks,K., Astrovskaya,I., Rendon,D.C., Khudyakov,Y., Berman,P. and Zelikovsky,A. (2008) HCV quasispecies assembly using network flows. In: *Proc. of International Symposium on Bioinformatics Research & Applications*. Vol. **4983**, pp. 159–170.
  44. Macalalad,A.R., Zody,M.C., Charlebois,P., Lennon,N.J., Newman,R.M., Malboeuf,C.M., Ryan,E.M., Boutwell,C.L., Power,K.A., Brackney,D.E. 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.
  45. Yang,X., Charlebois,P., Macalalad,A., Henn,M.R. and Zody,M.C. (2013) V-Phaser 2: variant inference for viral populations. *BMC Genomics*, **14**, 674.
  46. Routh,A., Chang,M.W., Okulicz,J.F., Johnson,J.E. and Torbett,B.E. (2015) CoVaMa: co-variation mapper for disequilibrium analysis of mutant loci in viral populations using next-generation sequence data. *Methods*, **91**, 40–47.
  47. Verbist,B.M., Thys,K., Reumers,J., Wetzels,Y., Van der Borgh,K., Talloen,W., Aerssens,J., Clement,L. and Thas,O. (2014) VirVarSeq: a low-frequency virus variant detection pipeline for Illumina sequencing using adaptive base-calling accuracy filtering. *Bioinformatics*, **31**, 94–101.
  48. Artyomenko,A., Wu,N.C., Mangul,S., Eskin,E., Sun,R. and Zelikovsky,A. (2017) Long single-molecule reads can resolve the complexity of the influenza virus composed of rare, closely related mutant variants. *J. Comput. Biol.*, **24**, 558–570.
  49. Mangul,S., Martin,L.S., Hill,B.L., Lam,A. K.-M., Distler,M.G., Zelikovsky,A., Eskin,E. and Flint,J. (2019) Systematic benchmarking of omics computational tools. *Nat. Commun.*, **10**, 1393.
  50. Mitchell,K., Brito,J.J., Mandric,I., Wu,Q., Knyazev,S., Chang,S., Martin,L.S., Karlsberg,A., Gerasimov,E., Littman,R. et al. (2020) Benchmarking of computational error-correction methods for next-generation sequencing data. *Genome Biol.*, **21**, 71.
  51. Eliseev,A., Gibson,K.M., Avdeyev,P., Novik,D., Bendall,M.L., Pérez-Losada,M., Alexeev,N. and Crandall,K.A. (2020) Evaluation of haplotype callers for next-generation sequencing of viruses. *Infect. Genet. Evol.*, **82**, 104277.
  52. Griebel,T., Zacher,B., Ribeca,P., Raineri,E., Lacroix,V., Guigó,R. and Sammeth,M. (2012) Modelling and simulating generic RNA-Seq experiments with the flux simulator. *Nucleic Acids Res.*, **40**, 10073–10083.
  53. Zanini,F., Brodin,J., Thebo,L., Lanz,C., Bratt,G., Albert,J. and Neher,R.A. (2015) Population genomics of inpatient HIV-1 evolution. *eLife*, **4**, e11282.
  54. Benidt,S. and Nettleton,D. (2015) SimSeq: a nonparametric approach to simulation of RNA-sequence datasets. *Bioinformatics*, **31**, 2131–2140.
  55. Baaijens,J.A., der Roest,B.V., Köster,J., Stougie,L. and Schönhuth,A. (2019) Full-length *de novo* viral quasispecies assembly through variation graph construction. *Bioinformatics*, **35**, 5086–5094.
  56. Tork,B., Nenastjeva,E., Artyomenko,A., Mancuso,N., Khan,M.I., O'Neill,R., Mandoiu,I.I. and Zelikovsky,A. (2016) Reconstruction of Infectious Bronchitis Virus Quasispecies from NGS Data. In:

- Mandoiu, I.I. and Zelikovsky, A.Z. (eds). *Computational Methods for Next Generation Sequencing Data Analysis*. Wiley, pp. 383–400.
57. Levina, E. and Bickel, P. (2001) The Earth mover's distance is the mallows distance: some insights from statistics. *Proc. ICCV 2001*, **2**, 251–256.
58. Mallows, C.L. (1972) A note on asymptotic joint normality. *Ann. Math. Stat.*, **43**, 508–515.
59. Yang, X., Charlebois, P., Gnerre, S., Coole, M.G., Lennon, N.J., Levin, J.Z., Qu, J., Ryan, E.M., Zody, M.C. and Henn, M.R. (2012) *De novo* assembly of highly diverse viral populations. *BMC Genomics*, **13**, 475.
60. Bron, C. and Kerbosch, J. (1973) Algorithm 457: finding all cliques of an undirected graph. *Commun. ACM*, **16**, 575–577.
61. Kováč, J. (2013) Complexity of the path avoiding forbidden pairs problem revisited. *Discrete Appl. Math.*, **161**, 1506–1512.
62. Nicolae, M., Mangul, S., Mandoiu, I. and Zelikovsky, A. (2011) Estimation of alternative splicing isoform frequencies from RNA-Seq data. *Algorithm. Mol. Biol.*, **6**, 9.
63. Kireev, D.E., Lopatukhin, A.E., Murzakova, A.V., Pimkina, E.V., Speranskaya, A.S., Neverov, A.D., Fedonin, G.G., Fantin, Y.S. and Shipulin, G.A. (2018) Evaluating the accuracy and sensitivity of detecting minority HIV-1 populations by Illumina next-generation sequencing. *J. Virol. Methods*, **261**, 40–45.
64. Wertheim, J.O., Leigh Brown, A.J., Hepler, N.L., Mehta, S.R., Richman, D.D., Smith, D.M. and Kosakovsky Pond, S.L. (2014) The global transmission network of HIV-1. *J. Infect. Dis.*, **209**, 304–313.
65. Quail, M.A., Smith, M., Coupland, P., Otto, T.D., Harris, S.R., Connor, T.R., Bertoni, A., Swerdlow, H.P. and Gu, Y. (2012) A tale of three next generation sequencing platforms: comparison of ion torrent, pacific biosciences and illumina MiSeq sequencers. *BMC Genomics*, **13**, 341.