

Inferring Viral Quasispecies Spectra from 454 Pyrosequencing Reads

Irina Astrovskaya*¹ , Bassam Tork¹ , Serghei Mangul¹ , Kelly Westbrook² , Ion Măndoiu³ , Peter Balfe⁴ , Alex Zelikovsky*¹

¹Department of Computer Science, Georgia State University, Atlanta, GA 30303, USA

²Life Technologies, Foster City, CA, USA

³Department of Computer Science & Engineering, University of Connecticut, Storrs, CT 06269, USA

⁴Institute of Biomedical Research, Birmingham University, Birmingham B15 2TT, UK

Email: Irina Astrovskaya* - iraa@cs.gsu.edu; Bassam Tork - btork@cs.gsu.edu; Serghei Mangul - serghei@cs.gsu.edu; Kelly Westbrook - kellywestbrooks@yahoo.com; Ion Măndoiu - ion@enr.uconn.edu; Peter Balfe - p.balfe@bham.ac.uk; Alex Zelikovsky* - alexz@cs.gsu.edu;

*Corresponding author

Abstract

Background: RNA viruses infecting a host usually exist as a set of closely related sequences, referred to as a quasispecies. The genomic diversity of viral quasispecies is a subject of great interest, particularly for chronic infections, since it can lead to resistance to existing therapies. High-throughput sequencing is a promising approach to characterizing viral diversity, but unfortunately standard assembly software was originally designed for single genome assembly and cannot be used to simultaneously assemble and estimate the abundance of multiple closely related quasispecies sequences.

Results: In this paper we introduce a new **Viral Spectrum Assembler (ViSpA)** method for quasispecies spectrum reconstruction and compare it with the state-of-the-art ShoRAH tool on both simulated and real 454 pyrosequencing shotgun reads from HCV and HIV quasispecies. Experimental results show that ViSpA outperforms ShoRAH on simulated error-free reads, correctly assembling 10 out of 10 quasispecies and 29 sequences out of 40 quasispecies. While ShoRAH has a significant advantage over ViSpA on reads simulated with sequencing errors due to its advanced error correction algorithm, ViSpA is better at assembling the simulated reads after they have been corrected by ShoRAH. ViSpA also outperforms ShoRAH on real 454 reads. Indeed, 7 most frequent sequences reconstructed by ViSpA from a real HCV

dataset are viable (do not contain internal stop codons), and the most frequent sequence was within 1% of the actual open reading frame obtained by cloning and Sanger sequencing. In contrast, only one of the sequences reconstructed by ShoRAH is viable. On a real HIV dataset, ShoRAH correctly inferred only 2 quasispecies sequences with at most 4 mismatches whereas ViSpA correctly reconstructed 5 quasispecies with at most 2 mismatches, and 2 out of 5 sequences were inferred without any mismatches. ViSpA source code is available at <http://alla.cs.gsu.edu/~software/VISPA/vispa.html>.

Conclusions: ViSpA enables accurate viral quasispecies spectrum reconstruction from 454 pyrosequencing reads. We are currently exploring extensions applicable to the analysis of high-throughput sequencing data from bacterial metagenomic samples and ecological samples of eukaryote populations.

Background

Viral Quasispecies. Many viruses (including SARS, influenza, HBV, HCV, and HIV) encode their genome in RNA rather than DNA. Unlike DNA viruses, RNA viruses lack the ability to detect and repair mistakes during replication [14] and, as a result, their mutation rate can be as high as 1 mutation per each 1,000-100,000 bases copied per replication cycle [13]. Many of the mutations are well tolerated and passed down to descendants producing a family of co-existing related variants of the original viral genome referred to as *quasispecies*, a concept that originally described a mutation-selection balance [10, 11, 21, 23, 29].

The diversity of viral sequences in an infected individual can cause the failure of vaccines and virus resistance to existing drug therapies [19]. Therefore, there is great interest in reconstructing genomic diversity of viral quasispecies. Knowing sequences of the most virulent variants can help to design effective drugs [5, 27] and vaccines [12, 17] targeting particular viral variants *in vivo*.

454 Pyrosequencing Technology. Briefly, the 454 pyrosequencing system shears the source genetic material into fragments of approximately 300-800 bases. Millions of single-stranded fragments are sequenced by synthesizing their complementary strands. Repeatedly, nucleotide reagents are flown over the fragments, one nucleotide (A, C, T, or G) at a time. Light is emitted at a fragment location when the nucleotide base flown complements the first unpaired base of the fragment [16, 22]. Multiple identical nucleotides may be incorporated in a single cycle, in which case the light intensity corresponds to the number of incorporated bases. However, since the number of incorporated bases (referred to as homopolymer length) cannot be estimated accurately for long homopolymers, this results in a

relatively high percentage of insertion and deletion sequencing errors (which represent 65%-75%, respectively 20%-30% of all sequencing errors [2,7]).

The software provided by instrument manufacturers was originally designed to assemble all reads into a single genome sequence, and cannot be used for reconstructing quasispecies sequences. Thus, in this paper we address the following problem:

Quasispecies Spectrum Reconstruction (QSR) Problem. *Given a collection of 454 pyrosequencing reads generated from a viral sample, reconstruct the quasispecies spectrum, i.e., the set of sequences and the relative frequency of each sequence in the sample population.*

A major challenge in solving the QSR problem is that the quasispecies sequences are only slightly different from each other. The amount and distribution along the genome of differences between quasispecies varies significantly between virus species, as different species have different mutation rates and genomic architectures. In particular, due to the lower mutation rate and longer conserved regions, HCV quasispecies are harder to reconstruct than quasispecies of HBV and HIV. Additionally, the QSR problem is made difficult by the limited read length and relatively high error rate of high throughput sequencing data generated by current technologies.

Related Work. The QSR problem is related to several well-studied problems: *de novo* genome assembly (see e.g., [8,24,31]), haplotype assembly [4,20], population phasing (see e.g., [6]) and metagenomics (see e.g., [32]). As noted above, *de novo* assembly methods are designed to reconstruct a single genome sequence, and are not well-suited for reconstructing a large number of closely related quasispecies sequences. Haplotype assembly does seek to reconstruct two closely related haplotype sequences, but existing methods do not extend easily to the reconstruction of a large (and *a priori* unknown) number of sequences. Computational methods developed for population phasing deal with large numbers of haplotypes, but rely on the availability of genotype data that conflates information about pairs of haplotypes. Metagenomic samples do consist of sequencing reads generated from the genomes of a large number of species. However, differences between the genomes of these species are considerably larger than those between viral quasispecies. Furthermore, existing tools for metagenomic data analysis focus on species identification, as reconstruction of complete genomic sequences would require much higher sequencing depth than that typically provided by current metagenomic datasets.

In contrast, achieving high sequencing depth for viral samples is very inexpensive, owing to the short length of viral genomes. Mapping based approaches to QSR are naturally preferred to *de novo* assembly since reference genomes are available (or easy to obtain) for viruses of interest, and viral genomes do not contain repeats. Thus, it is not surprising that such approaches were adopted in the two pioneering works on the QSR problem [15,33]. Eriksson et

al. [15] proposed a multi-step approach consisting of sequencing error correction via clustering, haplotype reconstruction via chain decomposition, and haplotype frequency estimation via expectation-maximization, with validation on HIV data. In Westbrook et al. [33], the focus is on haplotype reconstruction via transitive reduction, overlap probability estimation and network flows, with application to simulated error-free HCV data. Recently, the QSR software tool ShoRAH was developed [34] and applied to HIV data [35]. Another combinatorial method for QSR was also developed and applied to HIV and HBV data in [26], with results similar to those of ShoRAH.

Our contributions in this paper are as follows:

- A novel QSR tool called **Viral Spectrum Assembler (ViSpA)** taking into account sequencing errors at multiple steps,
- Comparison of ViSpA with ShoRAH on HCV synthetic data both with and without sequencing errors, and
- Statistical and experimental validation of the two methods on real 454 pyrosequencing reads from HCV and HIV samples.

Methods

Our method for inferring the quasispecies spectrum of a virus sample from 454 pyrosequencing reads consists of the following steps (see Fig. 1):

- Constructing the consensus virus genome sequence for the given sample and aligning the reads onto this consensus,
- Preprocessing aligned reads to correct sequencing errors,
- Constructing a transitively reduced read graph with vertices representing reads and edges representing overlaps between them,
- Selecting paths in the read graph that correspond to the most probable quasispecies sequences, and assembling candidate sequences for selected paths by weighted consensus of reads, and
- EM-based estimation of candidate sequence frequencies.

Below we describe each step separately.

Read Alignment and Consensus Genome Sequence Construction. We assume that a reference genome sequence of the particular virus strain is available (e.g., from NCBI [1]). Since viral genomes do not have sizable repeats and

the quasispecies sequences are usually close enough to the reference sequence, the majority of reads can typically be uniquely aligned onto the reference genome. However, a significant number of reads may remain unaligned due to differences between the reference genome and sequences in the viral sample. In order to recover as many of these reads as possible, we iteratively construct a consensus genome sequence from aligned reads.

In particular, we first align 454 pyrosequencing reads to the reference sequence using the SEGEMEHL software [30]. Then we extend the reference sequence with a placeholder I for each nucleotide inserted by at least one uniquely aligned read. Similarly, we add a placeholder D to the read sequence for each reference nucleotide missing from the aligned read. Then we perform sequential multiple alignment of the previously aligned reads against this extended reference sequence. Finally, the consensus genome sequence is obtained by (1) replacing each nucleotide in the extended reference with the nucleotide or placeholder in the majority of the aligned reads and (2) removing all I and D placeholders corresponding to rare insertions, respectively to deletions found in a majority of reads. Reads may contain a small portion of unidentified nucleotides denoted by N 's – we treat N as a special allele value matching any of nucleotides A, C, T, G , as well as placeholders I , and D .

Iteratively, we replace the reference with the consensus and try to align the reads for which we could not find any acceptable alignment previously. Our experiments on a dataset consisting of approximately 31,000 454 pyrosequencing reads generated from a 5.2Kbp-long HCV fragment (see data description in Results and Discussions) show that 85% of reads are uniquely aligned onto the reference sequence and an additional 9% of the reads are aligned onto the final consensus sequence. Reads that cannot be aligned onto the final consensus are removed from further consideration.

Preprocessing of Aligned Reads. Since aligned reads contain insertions and deletions, we use placeholders I and D to simplify position referencing among the reads. All placeholders are treated as additional allele values but they are removed from the final assembled sequences. First, we substitute each deletion in the aligned reads with placeholder D . Deletions supported by a single read are replaced either with the allele present in all the other reads overlapping this position if they agree with each other, or with N , signifying an unknown value, otherwise. Next, we fill with placeholder I each gap in a read corresponding to the insertions in the other reads. All insertions supported by a single read are removed from consideration.

Read Graph Construction. We begin with the definition of the read graph, introduced in [33] and independently in [15], and then describe the adjustments that need to be made to read graph construction and edge weights to account for sequencing errors as well as the high mutation rate between quasispecies.

The read graph $G = (V, E)$ is a directed graph with vertices corresponding to reads aligned with the consensus sequence. For a read u , we denote by $b(u)$, respectively $e(u)$, the genomic coordinate at which the first, respectively last, base of u gets aligned. A directed edge (u, v) connects read u to read v if a suffix of u overlaps with a prefix of v and they coincide across the overlap. Two auxiliary vertices – a source s and a sink t – are added such that s has edges into all reads with zero indegree and t has edges from all reads with zero outdegree. Then each $s - t$ -path corresponds to a possible candidate quasispecies sequence. The read graph is transitively reduced, i.e., each edge $e = (u, v)$ is removed if there is a $u - v$ -path not including edge e . Note that certain reads can be completely contained inside other reads. Let a *superread* refer to a read that is not contained in any other read and let the rest of the reads be called *subreads*. Subreads are not used in the construction of the read graph, but are taken into account in the final assembly of candidate sequences and frequency estimation.

Since the number of different $s - t$ -paths is exponential, we wish to generate a set of paths that have high probability to correspond to real quasispecies sequences. In order to estimate path probabilities, we independently estimate for each edge e the probability $p(e)$ that the two reads it connects come from the same quasispecies, and then multiply estimated probabilities for all edges on the path. Under the assumption of independence between edges, if we assign to each edge e a cost equal to $-\log(p(e)) = \log(1/p(e))$, then the minimum-cost $s - t$ -path will have the maximum probability to represent a quasispecies sequence.

For reads without errors, [33] estimated the probability that two reads u and v connected by edge (u, v) belong to the same quasispecies as

$$p_{\Delta} \approx \exp(-\Delta N/Lq) = \Theta(e^{-\Delta}) \quad (1)$$

where $\Delta = b(v) - b(u)$ is the *overhang* between reads u and v [33], $N = \#\text{reads}$, $q = \#\text{quasispecies}$, and $L = \#\text{starting positions}$. Thus, in this case the cost of an edge with overhang Δ can be approximated by $\Delta \propto \log(1/p_{\Delta})$.

To account for sequencing errors, we adjust the construction of the read graph to allow for mismatches. We use three parameters: (1) $n = \#\text{mismatches allowed between a read and a superread}$, (2) $m = \#\text{mismatches allowed in the overlap between two adjacent reads}$, and (3) $t = \#\text{mismatches expected between a read and a random quasispecies}$. The probability that two reads u and v with j mismatches within an overlap of length $o = e(u) - b(v)$ belong to the same quasispecies can be estimated as:

$$p_{\Delta_j} \approx e^{-\Delta} \binom{o}{j} (1 - \varepsilon)^{o-j} \varepsilon^j \quad (2)$$

where ε is the estimated 454 sequencing error rate. As in the case of error-free reads, defining the edge costs as $\Delta \log(\binom{o}{j}^{-1} (1 - \varepsilon)^{j-0} \varepsilon^{-j}) \propto \log(1/p_{\Delta_j})$ ensures that $s - t$ -paths with low cost correspond to most likely

quasispecies sequences.

Candidate Path Selection. To generate a set of high-probability (low-cost) paths rich enough to explain observed reads, we compute for each vertex in the read graph the minimum cost $s - t$ -path passing through it. Finding these paths is computationally fast. Indeed, we only need to compute in G two shortest-paths trees, one outgoing from s and one incoming into t ; the shortest $s - t$ -path passing through a vertex v is the concatenation of the shortest $s - v$ - and $v - t$ -paths.

Preliminary simulation experiments (see Additional File 1) show that better candidate sets are generated when edge costs c defined by (1) and (2) are replaced by e^c . In fact, if we use even faster dependency on c , then obtain even better candidate sets. The fastest growing cost effectively changes the shortest path into so called max-bandwidth paths, i.e., paths that minimize maximum edge cost for the entire path and for each subpath. So, ViSpA generates candidate paths using this strategy.

Candidate Sequence Assembly. When no mismatches are allowed in the construction of the read graph, finding the candidate sequence corresponding to an $s - t$ -path is trivial, since by definition adjacent superreads coincide across their overlap. When mismatches are allowed, we first assemble a consensus sequence using superreads used by the $s - t$ -path. This may be not the best choice, especially when the coverage with superreads is low. Hence, we replace each initial candidate sequence with a weighted consensus sequence obtained using both superreads and subreads of the path, as described below.

For each read r , we compute the probability that it belongs to a particular initial candidate sequence s as:

$$p(s, r) = \binom{l}{k} (1 - t/L)^{l-k} (t/L)^k \quad (3)$$

where l and L denote the lengths of the read and initial candidate sequence, respectively, k is the number of mismatches between the read and the initial candidate sequence s , and t/L is the estimated mutation rate. Then final candidate sequence is computed as the weighted consensus over all reads, where the weight of a read is the probability that it belongs to the sequence. Note that, unlike the case without mismatches, the same candidate sequence can be obtained from different candidate $s - t$ -paths, so at the end of this step we remove duplicates.

Estimation of Candidate Quasispecies Sequence Frequencies. We assume that reads R with observed frequencies $\{o_r\}_{r=1}^{|R|}$ where generated from a quasispecies population Q as follows. First, a quasispecies sequence $q \in Q$ is randomly chosen accordingly to its unknown frequency f_q . A read starting position is generated from the uniform distribution and then a read r is produced from quasispecies q with j sequencing errors. The probability of this event is calculated as $h_{q,r} = \binom{l}{j} (1 - \varepsilon)^{l-j} \varepsilon^j$, where l is the read length and ε is the sequencing error rate. Thus, the

probability of observing a read r under this model is $Pr(r) = \sum_{q \in Q} f_q h_{q,r}$.

Quasispecies frequencies $\{f_q\}_{q=1}^{|Q|}$ are estimated by maximizing the log-likelihood function:

$\ell(f_1, \dots, f_{|Q|}) = \sum_{r \in R} o_r \log Pr(r)$ using an EM algorithm [9] (see Additional File 1 for details). Currently, convergence of the EM algorithm is determined at the tolerance level 0.005.

Results and Discussions

In our simulation studies we use the following read data sets.

Reads simulated from known HCV quasispecies. In order to perform cross-validation on the assembly method, we simulate reads data from 1739-bp long fragment from the E1E2 region of 44 HCV sequences [18] when sequence frequencies are generated according to some specific distribution. In our simulation experiments, we use geometric distribution (i th sequence is constant factor more frequent than the $(i + 1)$ th sequence) to create sample quasispecies populations with different number of randomly selected above-mentioned quasispecies sequences.

We first simulate reads without sequencing errors: the length of a read follows normal distribution with a particular mean value and variance 400 and a starting position follows the uniform distribution. This simplified model of reads generation has two parameters: number of the reads that varies from 20K up to 100K and the averaged read length that varies from 200bp up to 600bp.

Additionally, we simulate 454 pyrosequencing reads for the 10 quasispecies sequences (following geometric distribution of frequencies) from the set of 44 HCV sequences [18] using FlowSim [3]. We generated 30K reads with average length 350bp.

454 Pyrosequencing Reads from HCV Samples. The data set Data1 has been received from HCV Research Group in Institute of Biomedical Research, at University of Birmingham. Data1 contains 30927 reads obtained from the 5.2Kb-long fragment of HCV-1a genome (which is more than a half of the entire HCV genome). The aligned read length average is 292bp but it significantly varies as well as the depth of position coverage (see Fig. 8 in Additional File 1). The depth of reads coverage variability is due to a strong bias in the sequence start points, reflecting the secondary structure of the template DNA or RNA used to generate the initial PCR products. As a result, shorter reads are produced by GC-rich sequences. Data is available upon request from the authors.

454 Pyrosequencing Reads from HIV Samples [35]. The HIV dataset contains 55611 reads from mixture of 10 different 1.5Kb-long region of HIV-1 quasispecies, including *pol* protease and part of the *pol* reverse transcriptase. The aligned reads length varies from 35 up to 584 with average about 345 bp (see Fig. 9 in Additional File 1). In contrast to [35], we do not filter out reads with low-quality scores.

Experimental Validation on Simulated Data. In all our experimental validations we compare the proposed algorithm ViSpA with state-of-the-art ShoRAH as well as with ViSpA on ShoRAH-corrected reads (ShoRAHreads + ViSpA). We say the quasispecies sequence is captured if one of the candidate sequences exactly matches it. We measure the quality of assembling by portion of the real quasispecies sequences being captured by candidate sequences (sensitivity = $\frac{TP}{TP+FN}$) and its portion among candidate sequences (positive predictive value ($PPV = \frac{TP}{TP+FP}$)) in cross-validation tests. Both sensitivity and PPV are analyzed as functions of the number of quasispecies in underlying sample population (see Fig. 2(left)). ViSpA can correctly assemble all sequences for 10 quasispecies and 29 sequences for 40 quasispecies if average read length is at least 300bp. If the average read length is smaller (for example, in range from 250bp till 299bp), the method can assemble at least 8 out of 10 sequences and 20 out of 40 sequences. Here, we see advantage of ViSpA over ShoRAH.

Following [15], we measure the prediction quality of frequency distribution with Kullback-Leibler divergence, or relative entropy. Given two probability distributions, relative entropy measures the "distance" between them, or, in other words, the quality of approximation of one probability distribution by the other distribution. Formally, the relative entropy between true distribution P and approximation distribution Q is given by the formula:

$$D_{KL}(P||Q) = \sum_{i \in I} P(i) \log \frac{P(i)}{Q(i)},$$

where summation is over all reconstructed original sequences $I = \{i | P(i) > 0, Q(i) > 0\}$, i. e., over all original sequences that have a match (exact or with at most k mismatches) among assembled sequences.

Relative entropy is decreasing with increasing of the average read length. It is expected since sensitivity is increasing with increasing of the average read length and EM predicts underlying distribution more accurate. ViSpA algorithm considerably outperforms ShoRAH (see Fig. 2(right)).

However, ShoRAH has a significant advantage over ViSpA on a read data simulated by FlowSim both in prediction power and in robustness of results (see Table 1). Indeed, ShoRAH correctly infers 3 out of 10 real quasispecies sequences whereas ViSpA reconstructs only 1 sequence. Additionally, 10 most frequent assemblies inferred by ShoRAH are more robust with repeating up to 45% of times on 10%-reduced data versus 1% of times for ViSpA's assemblies. This advantage can be explained by superior read correction in ShoRAH. If ViSpA is used on ShoRAH-corrected reads, the results drastically improves: 5 quasispecies sequences are inferred and exactly 95% of times are repeated on reduced data, confirming that ViSpA is better in assembling sequences (see Table 1).

Experimental Validation on 454 Pyrosequencing Reads from HCV Samples. We first discuss the choice of parameters of the read graph and candidate sequence assembly from $s - t$ -paths. Then we give statistical validation for obtained 10 most frequent quasispecies sequences. Finally, we show how to identify and fix the erroneous

homopolymer indels in the coding region of HCV.

We infer quasispecies spectrum based on the read graphs constructed with various numbers n and m (numbers of mismatches allowed for superreads and overlaps corresponding to edges). We sort the estimated frequencies in descending order and count the number of sequences which cumulative frequency is 85%, 90%, and 95%. Fig. 3 reports these numbers as a percent of the total number of candidate sequences. There is an obvious drop in percentage for all three categories if we allow up to $n = 6$ mismatches to cluster reads and up to $m = 15$ mismatches to create edges. In this case, the constructed read graph has no isolated vertices.

To refine assembled candidate sequences, we use all reads and parameter t varying from 80bp till 350 bp, or, in the other words, mutation rate varying from 1.75% up to 8% per sequence (which is in the range observed in [28]).

Out of 3207 max-bandwidth paths, we obtain as much as 938 distinct sequences ($t = 80$) and as low as 755 sequences ($t = 350$) for different value of $t \in [80, 350]$.

The neighbor-joining tree for the most frequent 10 candidate sequences (see Fig. 4) obtained by ViSpA and ShoRAH reminds a neighbor-joining tree for HCV quasispecies evolution. Additionally, the most frequent candidate sequence found by ViSpA is 99% identical to one of the actual ORFs obtained by cloning the quasispecies.

Viral sequences containing internal stop codons are not viable since the entire HCV genome consists of a single coding region for a large polyprotein. So the number of reconstructed viable sequences can serve as an accuracy measure for quasispecies assembly. Out of 10 most frequent sequences reconstructed by ViSpA, only 3 are not viable while ShoRAH is able to reconstruct only one viable sequence. This sequence has 99.94% similarity with the ViSpA's fourth most frequent assemblies. Both methods returned similar frequency estimations for this sequence: 0.017% (ShoRAH) and 0.019% (ViSpA).

Both ShoRAH and ViSpA ($n = 6, m = 15$) are run on 8 2.66GHz-CPU's with 8M cache. They take around 40 minutes to assemble sequences and estimate their frequencies. Smaller value of n increases ViSpA's runtime since its bottleneck (candidate sequences assembling) is proportional to the number of reads times number of paths. Indeed, smaller value of n results in larger number of superreads in built read graph, thus, in larger set of candidate paths. For example, ViSpA runs 90 minutes for $n = 2, m = 2$.

Statistical Validation of the Quasispecies Spectrum. The plot on Fig. 5 shows validation results for 10 most frequent quasispecies sequences with respect to EM estimations assembled on Data1 by ShoRAH and ViSpA ($n = 6, m = 15, \text{ and } t = 120$). Repeatedly, 100 times we have deleted randomly chosen 10% of reads and we run both methods on each reduced read instance to reconstruct quasispecies spectrum.

The plot reports the percentage of runs when each of 10 most frequent sequences assembled on Data1 are reproduced

among the 10 most frequent quasispecies inferred on the reduced instances with no mismatches ($k = 0$), or with $k = 1, 2, 5$ mismatches. For example, for $k = 0$ ShoRAH repeatedly (35% of times) reconstructs only the 3rd most frequent sequence, while ViSpA reconstructs 7 sequences in at least 15% times and the most frequent sequence is reconstructed 40% times. This plot shows that the found sequences are pretty much reproducible for ViSpa.

Experimental Validation on 454 Pyrosequencing Reads from HIV Samples. In order to compare ViSpA and ShoRAH, we run both of the methods on HIV dataset, used in the first experiment in [35]. As was said above, we do not preprocess reads with respect to its 454 quality score, and it can explain poorer performance of ShoRAH. Indeed, ShoRAH correctly infers only 2 quasispecies sequences with at most 4 mismatches: one assembly has 3 mismatches with real quasispecies sequence, and the other has 4 mismatches. These assemblies were 256-th and 388-th most frequent sequences. Among 100 most frequent assemblies, the smallest distance between assemblies and real quasispecies sequences is 28 bp.

ViSpA correctly reconstructs 5 quasispecies with at most 2 mismatches (3 of them among 10 most frequent assemblies): two sequences are inferred without any mismatches (one is among 10 most frequent assemblies), one assembly has 1 mismatch with real quasispecies sequence (and it is among 10 most frequent assemblies), and the rest sequences have 2 mismatches (one is among 10 most frequent assemblies). The assemblies correspond to a viable protein sequences.

Conclusions

In this paper, we have proposed and implemented ViSpA, a novel software tool for quasispecies spectrum reconstruction from high-throughput sequencing reads. The ViSpA assembler takes into account sequencing errors at multiple steps, including mapping-based read preprocessing, path selection based on maximum-bandwidth, and candidate sequence assembly using probability-weighted consensus techniques. Sequencing errors are also taken into account in ViSpA's EM-based estimation of quasispecies sequence frequencies.

We have validated our method on simulated error-free reads, FlowSim-simulated reads with sequencing errors, and real 454 pyrosequencing reads from HCV and HIV samples. We are currently exploring extensions of ViSpA to paired-reads; the main difficulty is selection of pair-aware candidate paths. We also foresee application of ViSpA's techniques to the analysis of high-throughput sequencing data from microbial communities [32] and ecological samples of eukaryote populations [25].

Availability

The ViSpA source code is available at <http://alla.cs.gsu.edu/~software/VISPA/vispa.html>.

Authors contributions

IA designed algorithms, developed software, performed analyzes and experiments, wrote the paper. BT performed analyses and experiments. SM contributed to developing software. KW designed algorithms and developed software. IM contributed to designing the algorithms and writing the paper. PB supplied the HCV data and contributed to performing the analysis. AZ designed the algorithms, wrote the paper and supervised the project. All authors have read and approved the final manuscript.

Competing interests

The authors declare that they have no competing interests.

Acknowledgements

This work has been partially supported by NSF awards IIS-0546457, IIS-0916401, IIS-0916948, and GSU Molecular Basis of Disease Fellowship. KW worked on the project while being a Ph.D. student in GSU Computer Science Department.

References

1. National center for biotechnology information, <http://www.ncbi.nlm.nih.gov>.
2. Michael P. Strömberg Gábor T. Marth Aaron R. Quinlan, Donald A. Stewart. PyroBayes: an improved base caller for SNP discovery in pyrosequences. *Nature Methods*, 5(2):179–181, February 2008.
3. S. Balsler, K. Malde, A. Lanzen, A. Sharma, and I. Jonassen. Characteristics of 454 pyrosequencing data—enabling realistic simulation with FlowSim. *Bioinformatics*, 26:i420–5, 2010.
4. V. Bansal and V. Bafna. HapCUT: an efficient and accurate algorithm for the haplotype assembly problem. *Bioinformatics*, 24(16):i153–i159, August 2008.
5. N. Beerenwinkel, T. Sing, T. Lengauer, J. Rahnenfuehrer, and K. Roomp et al. Computational methods for the design of effective therapies against drug resistant HIV strains. *Bioinformatics*, 21:3943–3950, 2005.
6. D. Brinza and A. Zelikovsky. 2SNP: Scalable phasing based on 2-SNP haplotypes. *Bioinformatics*, 22:371–373, 2006.
7. Alvarez P. Young S. Garber M. Giannoukos G. Lee W. L. Russ C. Lander-E.S. Nusbaum C.-Jaffe D. B. Brockman, W. Quality scores and SNP detection in sequencing-by-synthesis systems. *Genome Reasearch*, 18(5):763–770, 05 2008.
8. M.J. Chaisson and P.A. Pevzner. Short read fragment assembly of bacterial genomes. *Genome research*, 18:324–330, 2008.
9. Laird N. M. Rubin D. B. Dempster, A. P. Maximum likelihood from incomplete data via the EM algorithm (with discussions). *Journal of the Royal Statistical Society, B*, 39, pages 1–38, 1977.
10. E. Domingo and J.J. Holland. RNA virus mutations and fitness for survival. *Annu Rev Microbiol*, 51:151–178, 1997.
11. Martinez-Salas E. Sobrino F. de la Torre J.C. Portela A. Ortin J. Lopez-Galindez C. Perez-Brena P. Villanueva N. Najera R. Domingo, E. The quasispecies (extremely heterogeneous) nature of viral RNA genome populations: biological relevance – a review. *Gene*, 40, pages 1–8, 1985.
12. D.C. Douek, P.D. Kwong, and G.J. Nabel. The rational design of an AIDS vaccine. *Cell*, 124:677–681, 2006.
13. Holland-J.J. Drake, J.W. Mutation rates among RNA viruses. In *Proceedings of the National Academy of Sciences USA*, 96, pages 13910–13913, 1999.

14. Novella-I.S. Weaver S.C. Domingo E. Wain-Hobson S. Clarke D.K. Moya A. Elena-S.F. de la Torre J.C. Holland J.J. Duarte, E.A. RNA virus quasispecies:significance for viral disease and epidemiology. *Infectious Agents and Disease*, 3, pages 201–214, 1994.
15. N. Eriksson, L. Pachter, Y. Mitsuya, S.Y. Rhee, and C. Wang et al. Viral population estimation using pyrosequencing. *PLoS Comput Biol*, 4:e1000074, 2008.
16. H. Fakhrai-Rad, N. Pourmand, and M. Ronaghi. Pyrosequencing: An accurate detection platform for single nucleotide polymorphisms. *Hum Mutat*, 19:479–485, 2002.
17. B. Gaschen, J. Taylor, K. Yusim, B. Foley, and F. Gao et al. Diversity considerations in HIV-1 vaccine selection. *Science*, 296:2354–2360, 2002.
18. T. Von Hahn, J.C. Yoon, H. Alter, C.M. Rice, B. Rehmann, P. Balfe, and J.A. Mckeating. Hepatitis C virus continuously escapes from neutralizing antibody and T-cell responses during chronic infection in vivo. *Gastroenterology*, 132:667–678, 2007.
19. de la Torre J.C. Steinhauer D.A. Holland, J.J. RNA virus populations as quasispecies. *Current Topics in Microbiology and Immunology*, 176, pages 1–20, 1992.
20. R. Lippert, R. Schwartz, G. Lancia, and S. Istrail. Algorithmic strategies for the single nucleotide polymorphism haplotype assembly problem. *Briefings in Bioinformatics*, 3:23–31, 2002.
21. M. Eigen M, J. McCaskill, and P. Schuster. The molecular quasi-species. *Adv Chem Phys*, 75:149–263, 1989.
22. M. Margulies, M. Egholm, W.E. Altman, S. Attiya, and J.S. Bader et al. Genome sequencing in microfabricated high-density picolitre reactors. *Nature*, 437:376–380, 2005.
23. Esteban J.I. Quer J. Genesca J. Weiner-A. Esteban R. Guardia J. Gomez J. Martell, M. Hepatitis C virus (HCV) circulates as a population of different but closely related genomes: quasispecies nature of HCV genome distribution. *Journal of Virology*, 66, pages 3225–3229, 1992.
24. G. Myers. Building fragment assembly string graphs. In *Proc. ECCB*, pages 79–85, 2005.
25. S.T. O’Neil and S. Emrich. Robust haplotype reconstruction of eukaryotic read data with Hapler In *Proc. ICCABS*, pages 141–146, 2011.
26. Mattia Prosperi, Luciano Prosperi, Alessandro Bruselles, Isabella Abbate, Gabriella Rozera, Donatella Vincenti, Maria Solmone, Maria Capobianchi, and Giovanni Ulivi. Combinatorial analysis and algorithms for quasispecies reconstruction using next-generation sequencing. *BMC Bioinformatics*, 12(1):5+, 2011.
27. S-Y. Rhee, T.F. Liu, S.P. Holmes, and R.W. Shafer. HIV-1 subtype B protease and reverse transcriptase amino acid covariation. *PLoS Comput Biol*, 3:e87, 2007.
28. Andrea D. Branch Sarah L. Fishman. The quasispecies nature and biological implications of the hepatitis C virus. *Infection, Genetics and Evolution*, 9:1158–1167, 2009.
29. Holland J.J. Steinhauer, D.A. Rapid evolution of RNA viruses. *Annual Review of Microbiology*, 41, pages 409–433, 1987.
30. Stefan Kurtz Cynthia M. Sharma-Philipp Khaitovich Jörg Vogel Peter F. Stadler Steve Hoffmann, Christian Otto and Jörg Hackermüller. Fast mapping of short sequences with mismatches, insertions and deletions using index structures. *PLoS Comput Biol*, 5(9):e1000502, 09 2009.
31. A. Sundquist, M. Ronaghi, H. Tang, P. Pevzner, and S. Batzoglou. Whole-genome sequencing and assembly with high-throughput, short-read technologies. *PLoS ONE*, 2:e484, 2007.
32. J.C. Venter et al. Environmental genome shotgun sequencing of the Sargasso sea. *Science*, 304:66–74, 2004.
33. K. Westbrooks, I. Astrovskaya, D. Campo, Y. Khudyakov, P. Berman, and A. Zelikovsky. HCV quasispecies assembly using network flows. In *Proc. ISBRA*, pages 159–170, 2008.
34. Osvaldo Zagordi, Lukas Geyrhofer, Volker Roth, and Niko Beerenwinkel. Deep sequencing of a genetically heterogeneous sample: local haplotype reconstruction and read error correction. *Journal of computational biology : a journal of computational molecular cell biology*, 17(3):417–428, March 2010.
35. Osvaldo Zagordi, Rolf Klein, Martin Daumer, and Niko Beerenwinkel. Error correction of next-generation sequencing data and reliable estimation of HIV quasispecies. *Nucleic Acids Research*, 38(21):7400–7409, 2010.

Figures

Figure 1 - ViSpA's flowchart.

Figure 2 - Statistical validation on error-free reads from known HCV quasispecies.

Left: PPV and sensitivity as a function of the number of quasispecies in the original population (40K reads with average read length 300). Right: the relative entropy as a function of the average read length (40K reads from 10 quasispecies).

Figure 3 - Percentage of candidate sequences which cumulative frequencies is 85%, 90%, and 95%.

The values on x-axis corresponds to the number of allowed mismatches during read graph construction. n_m means that up to n mismatches are allowed in superreads and up to m mismatches are allowed in edges.

Figure 4 - The neighbor-joining phylogenetic tree for 10 most frequent HCV quasispecies variants on a 5.205bp long fragment obtained by ViSpA and ShoRAH.

Sequences are labeled with software name and its rank among 10 most frequent assembled sequences.

Figure 5 - Percentage of runs when the i -th most frequent sequences is reproduced among 10 most frequent quasispecies assembled on the 10%-reduced set of reads.

The i -th point at x-axis corresponds to the i -th most frequent sequence assembled on the 100% of reads. No data are shown for the sequences that are reproduced less than 5% of runs.

Tables

Table 1 - Comparison of three methods – ViSpA, ShoRAH, and ShoRAHreads+ViSpA – on the read data simulated by FlowSim.

The quasispecies sequence is considered found if one of candidate sequences matches it exactly ($k = 0$) or with at most k (1 or 9) mismatches. All methods are run 100 times on 10% - reduced data. For the i -th ($i = 1, \dots, 10$) most frequent sequence assembled on the whole data, we record its reproducibility, i.e., percentage of runs when there is a match (exact or with at most k mismatches) among 10 most frequent sequences found on reduced data.

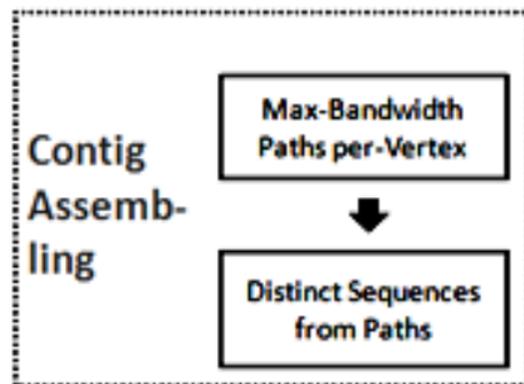
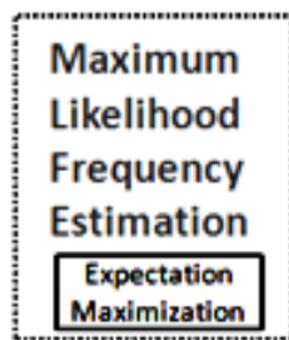
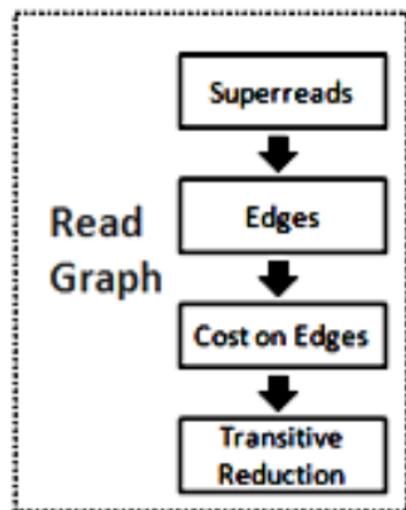
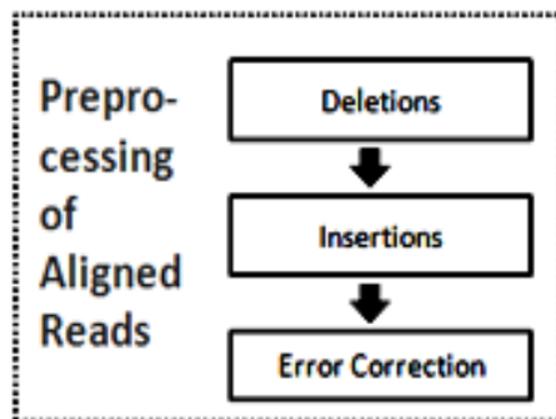
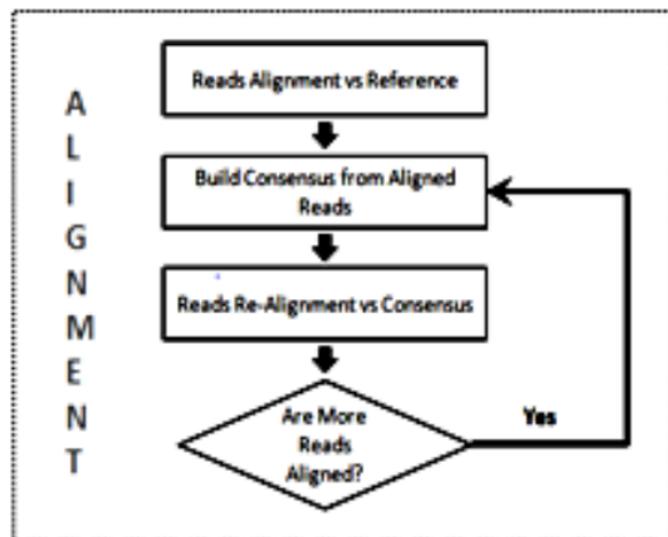
”Reproducibility: Max” and ”Reproducibility: Average” report respectively maximum and average of those percentages.”

	ShoRAH				ViSpA				ShoRAHreads+ViSpA			
	PPV	Sensitivity	Reproducibility		PPV	Sensitivity	Reproducibility		PPV	Sensitivity	Reproducibility	
			Max	Average			Max	Average			Max	Average
k=0	0.0097	0.3	0.45	0.11	0.0008	0.1	0.1	0.1	0.5	0.5	0.95	0.95
k=1	0.0129	0.4	0.6	0.32	0.0008	0.1	0.1	0.1	0.5	0.5	0.95	0.95
k=9	0.0162	0.5	0.95	0.64	0.0015	0.2	0.1	0.1	0.5	1	0.95	0.95

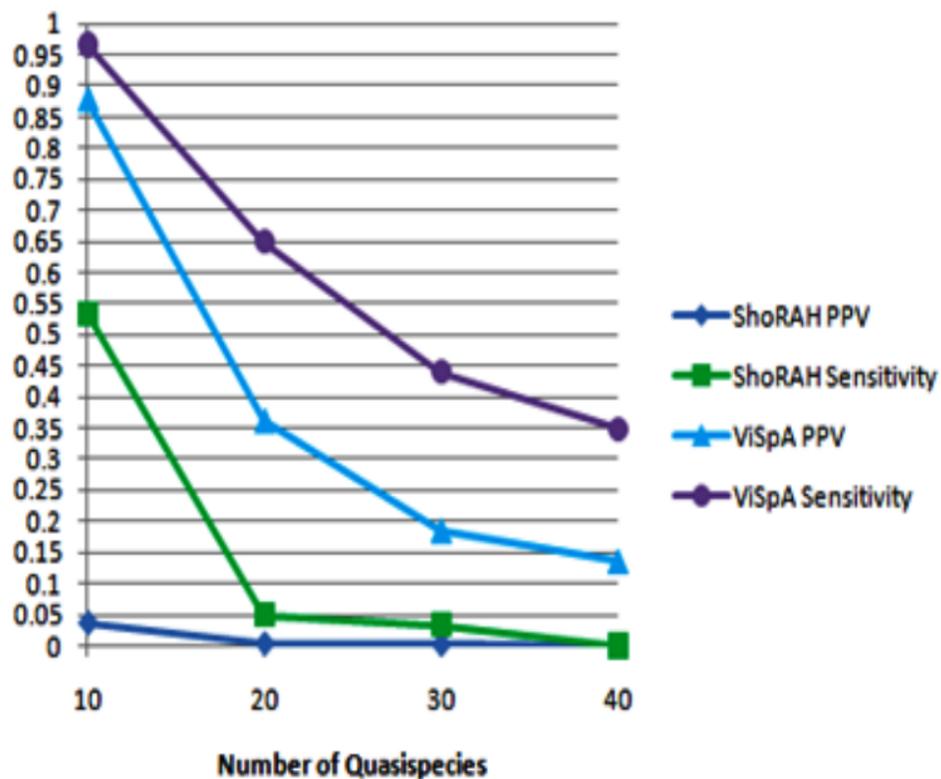
Additional Files

Additional file 1 — Supplementary Materials.

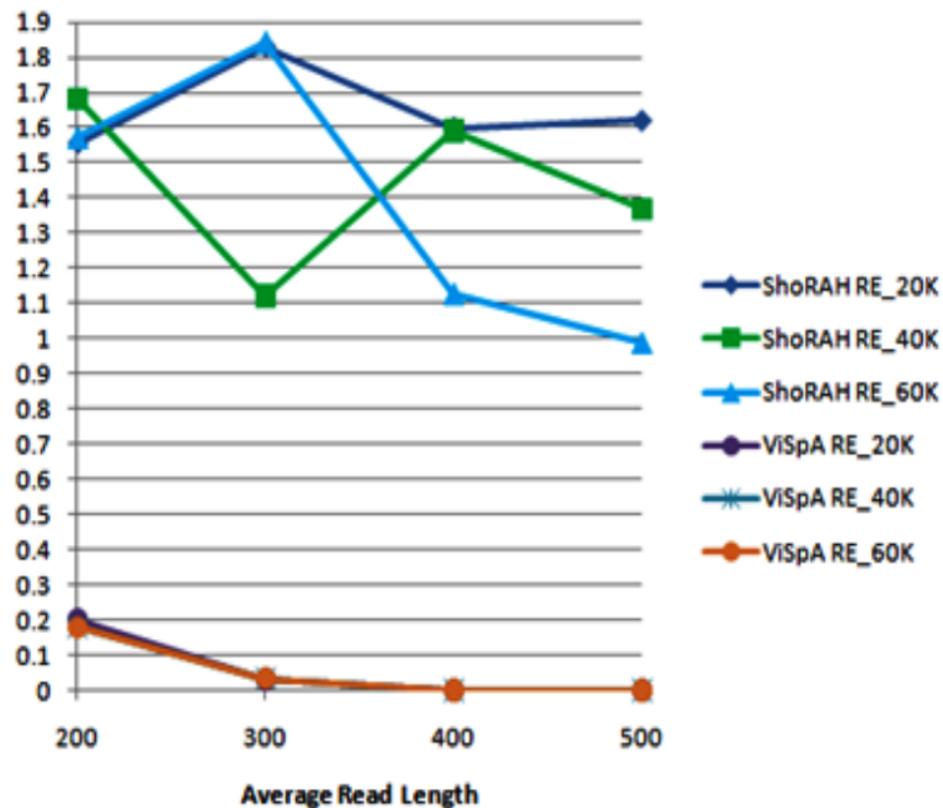
The file contains derivation of edge cost formula (2) and EM algorithm, example of read graph construction and analysis of 454 pyrosequencing data.



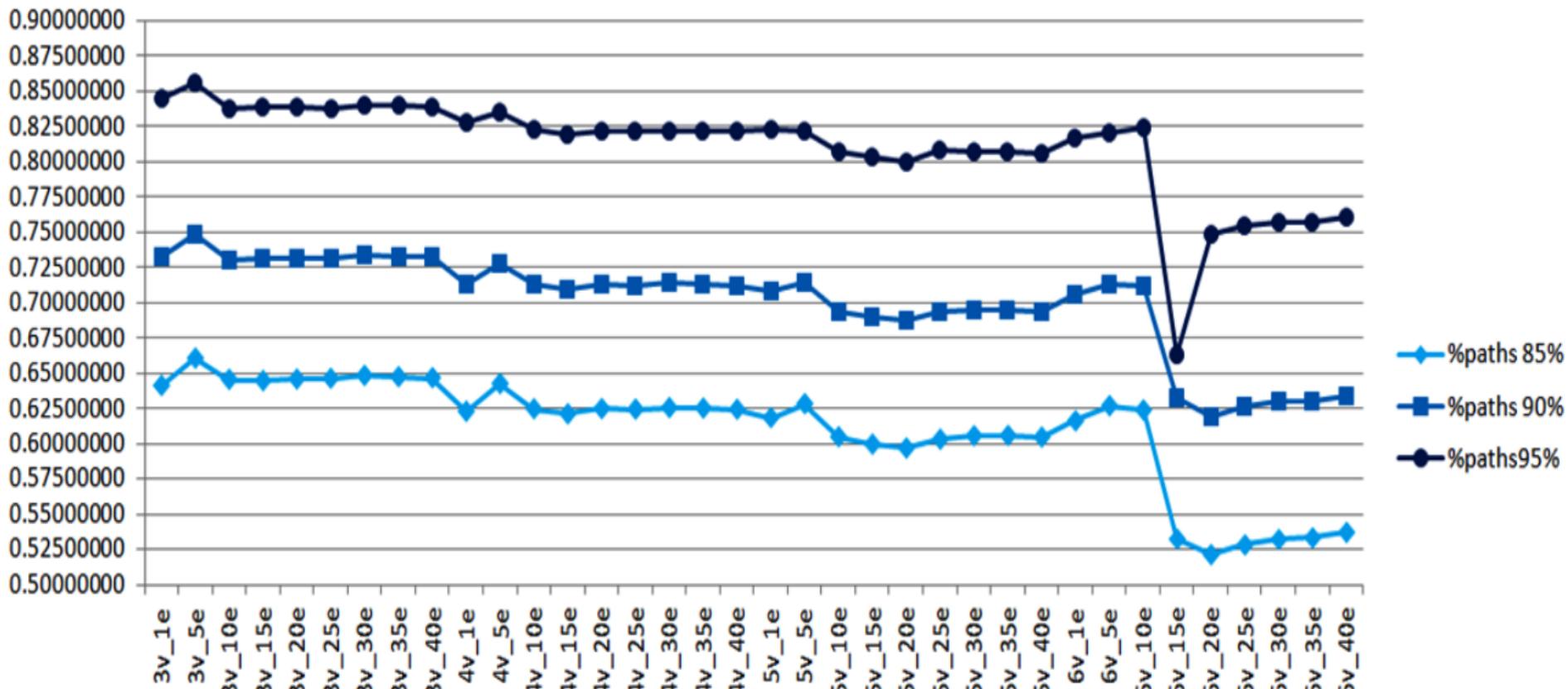
PPV and Sensitivity

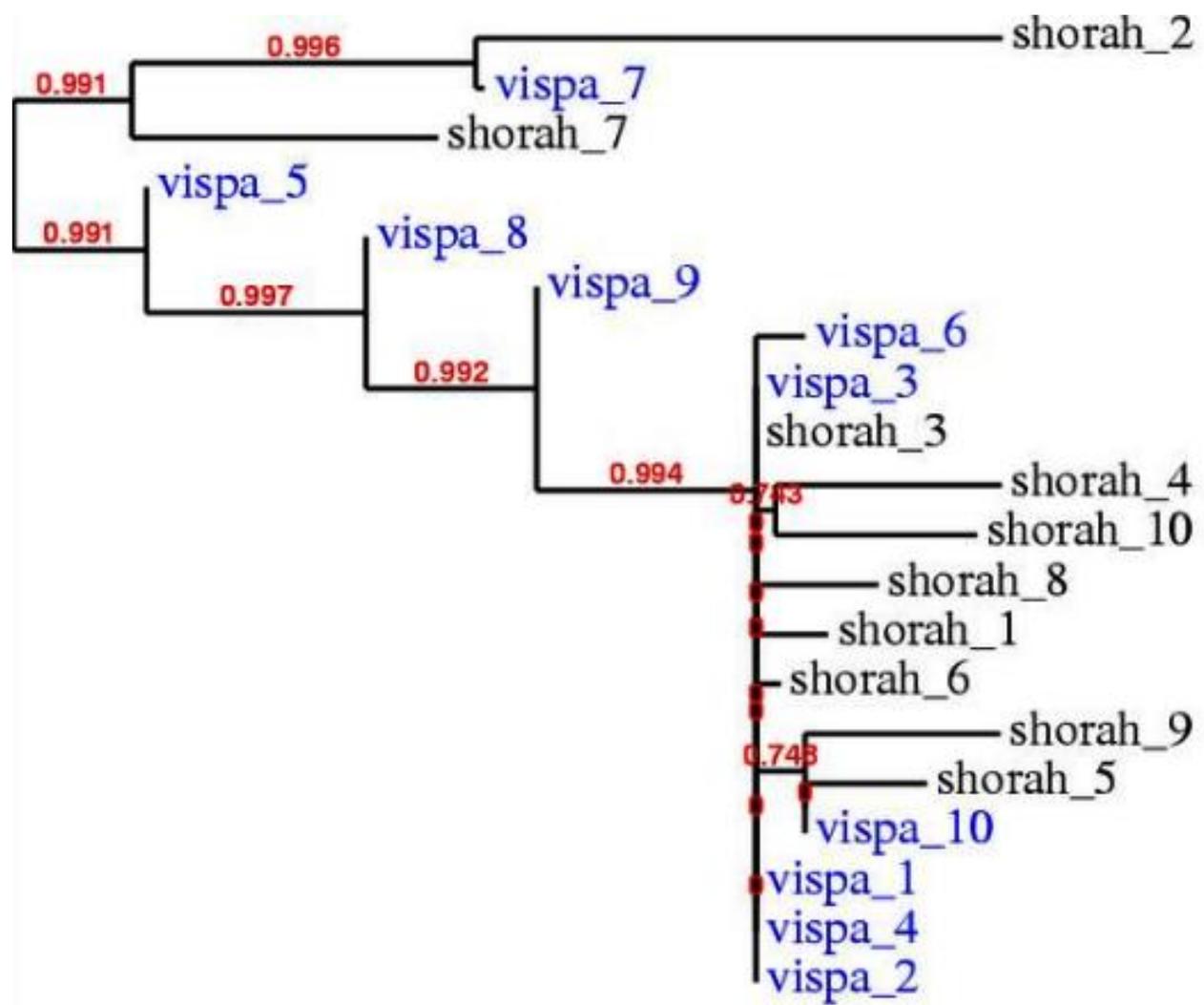


Relative Entropy



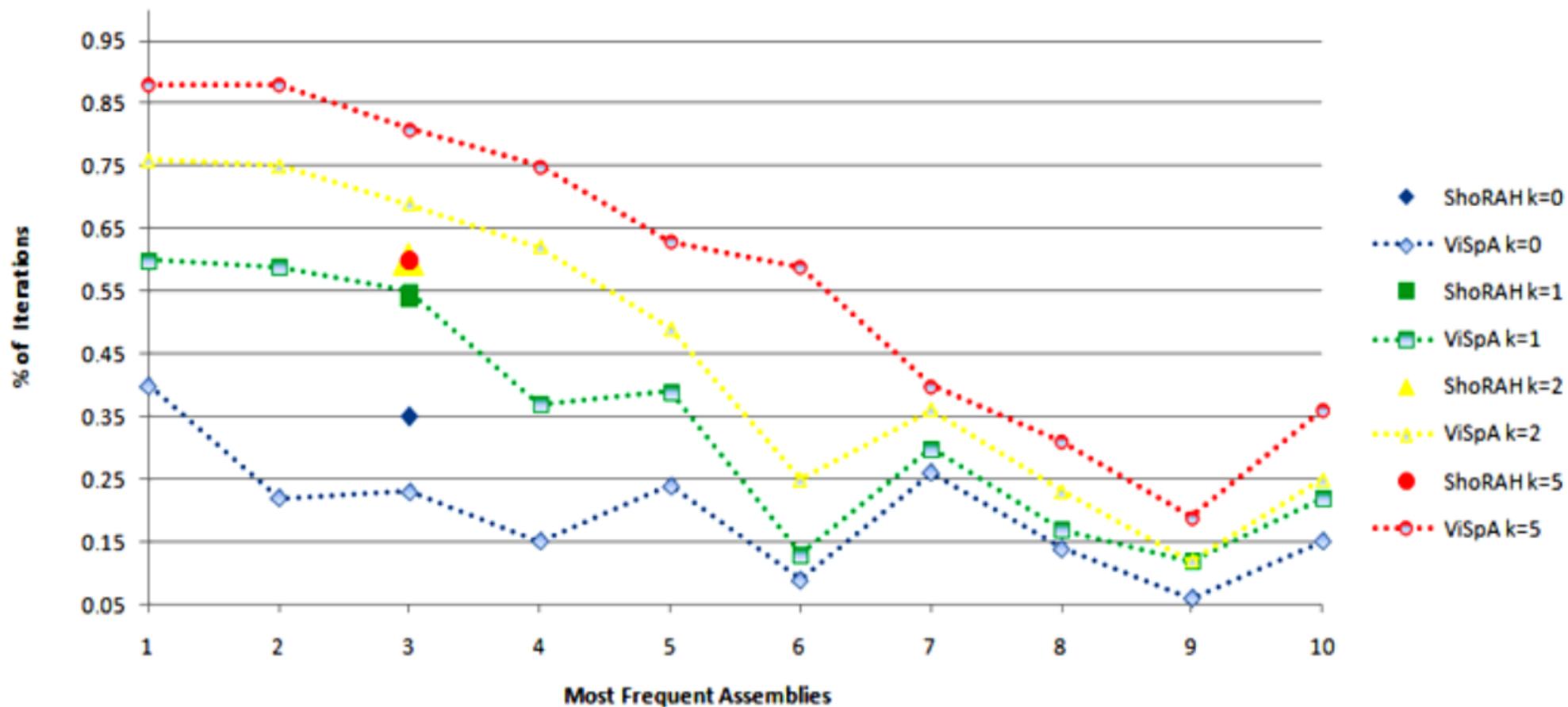
Percentage of Max-Bandwidth paths for R8_MID2_ID10 without N's





0.002

ShoRAH vs ViSpA: % of Iterations



Supplementary Materials for Inferring Viral Spectrum from 454 Pyrosequencing Reads

Irina Astrovskaya, Bassam Tork, Serghei Mangul, Kelly Westbrook, Ion Mandoiu, Peter Balfe, and
Alex Zelikovsky

1 Methods

Cost Formula Derivation. First, let us consider a simplified model where all quasispecies are uniformly distributed and reads' beginning positions follow uniform distribution as well. Let $b(u)$ and $e(u)$ be the beginning and the end positions of the read u , respectively. And let A be an event that two reads u, v from the same quasispecies Q are connected with an edge (u, v) in the transitively reduced graph and j differences among u and v occur in the overlap due to sequencing errors. In the other words, the event A is a product of the following independent events: (1) read u exists, (2) read v exists, starting at position $b(v) = b(u) + \Delta$ (overhang Δ is some shift) (3) no read w from the same quasispecies Q satisfies $b(u) < b(w) < b(v)$, and (4) there are j sequencing errors in overlap of reads u and v .

Given N reads, originated from q quasispecies of length L , the probability that a read u starts at a position $b(u)$ is $\frac{N}{Lq}$. The probability of the event $\Delta > k$ is the probability that there is no read from quasispecies Q starting at any position in $\{b(u) + 1, b(u) + 2, \dots, b(u) + k\}$, that is

$$p_k = \left(1 - \frac{N}{Lq}\right)^k \approx \exp(-kN/Lq).$$

Then

$$Pr(\Delta = k) = \left(\frac{N}{Lq}\right)^2 p_{k-1}.$$

Next we calculate the probability of having j genotyping errors among overlapped positions as follows:

$$p(j) = \binom{e(u) - b(v)}{j} (1 - \epsilon)^{e(u) - b(v) - j} \epsilon^j,$$

where ϵ is a genotyping error rate. Finally, the probability of event A is:

$$Pr(A) = p(\Delta_j) = Pr(\Delta = k)p(j) \approx \left(\frac{N}{Lq}\right)^2 \exp(-\Delta N/Lq) \binom{e(u) - b(v)}{j} (1 - \epsilon)^{e(u) - b(v) - j} \epsilon^j$$

where $\Delta = b(v) - b(u)$ is a shift (overhang) between starting positions of reads u and v .

If $\Delta \gg Lq/N$ then v is more likely to be a read from another quasispecies Q' and differences from quasispecies Q and Q' in interval between $b(u)$ and $b(v)$ are the cause of large Δ .

Therefore,

$$1/p_{\Delta_j} \approx \frac{\exp(\Delta N/Lq)}{\binom{e(u) - b(v)}{j} (1 - \epsilon)^{e(u) - b(v) - j} \epsilon^j} \quad (1)$$

measures our uncertainty that (u, v) is a true edge.

In practice, quasispecies frequencies are under some unknown distribution F . Since it is impossible to approximate, we assume that two identical copies of the same quasispecies genome correspond to two different entities, which follow uniform distribution.

In experimental results, reads starting positions tend to follow uniform distribution with small amount of outliers. If it is not a case, the cost formula can be adjusted by taking into account how many reads start in interval between $b(u)$ and $b(v)$. However, when a candidate path is constructed, at each step, algorithm chooses between edges spanning almost the same positions, thus, non-adjusting cost formula to the coverage does not influence the construction of candidate path.

Frequency Estimation via EM-algorithm.

Once a set of candidate sequences is obtained, their maximum-likelihood frequencies are calculated by the Expectation Maximization-based algorithm. Iteratively, we estimate missing probability $p_{q,r}$ that read r comes from a candidate sequence q with j sequencing errors and maximize likelihood of an approximated model.

First, we create a bipartite graph $G = \{Q \cup R, E\}$ such that each candidate sequence is represented as a vertex $q \in Q$, and each read is represented as a vertex $r \in R$. With each vertex $q \in Q$, we associate unknown frequency f_q of the candidate sequence. And with each vertex $r \in R$, we associate read's observed frequency o_r . Then for each pair q, r , we add an edge (q, r) weighted by probability of the read r being produced by the candidate sequence q with j genotyping errors:

$$h_{q,r} = \binom{l}{j} (1 - \epsilon)^{l-j} \epsilon^j,$$

where l is length of read sequence, and ϵ is the genotyping error rate.

After initializing frequencies $f_{q \in Q}$ at random, the algorithm repeatedly performs the next two steps until convergence:

E-step: For each pair q, r , compute the expected value $p_{q,r}$ that read r comes from candidate sequence q under the assumption that frequencies $f_{q \in Q}$ are correct by the following formula:

$$p_{q,r} = \frac{f_q \cdot h_{q,r}}{\sum_{q': (q', r) \in E} f_{q'} \cdot h_{q', r}}.$$

M-step: For each $q \in Q$, update value of f_q to the portion of reads being originated by the candidate sequence q among all observed reads in the sample, i.e.:

$$f_q = \frac{\sum_{r: (q, r) \in E} p_{q,r} \cdot o_r}{\sum_{r \in R} o_r}.$$

Currently, convergence of EM algorithm is determined at the tolerance level 0.005.

Rationale for Max-Bandwidth Path. Previously, we show that cost of edge should be correlated to overhang (shift) Δ . If we view the costs on edges as edge lengths, then the most probably paths for quasispecies are the shorter paths in the graph. So we can choose shortest path for each vertex to build the set of candidate paths.

In experiments on error-free reads, we consider the family of edge cost functions $cost_k(u, v) = e^{\frac{\Delta(u, v)}{k}}$. Figure 1 shows the number of the shortest paths in candidate set as a function of k . Smaller values of k yield fewer paths, and surprisingly, no correct candidate quasispecies are lost with decreasing of k . In the limiting case $k = 0$, the resulted paths are maximum bandwidth paths.

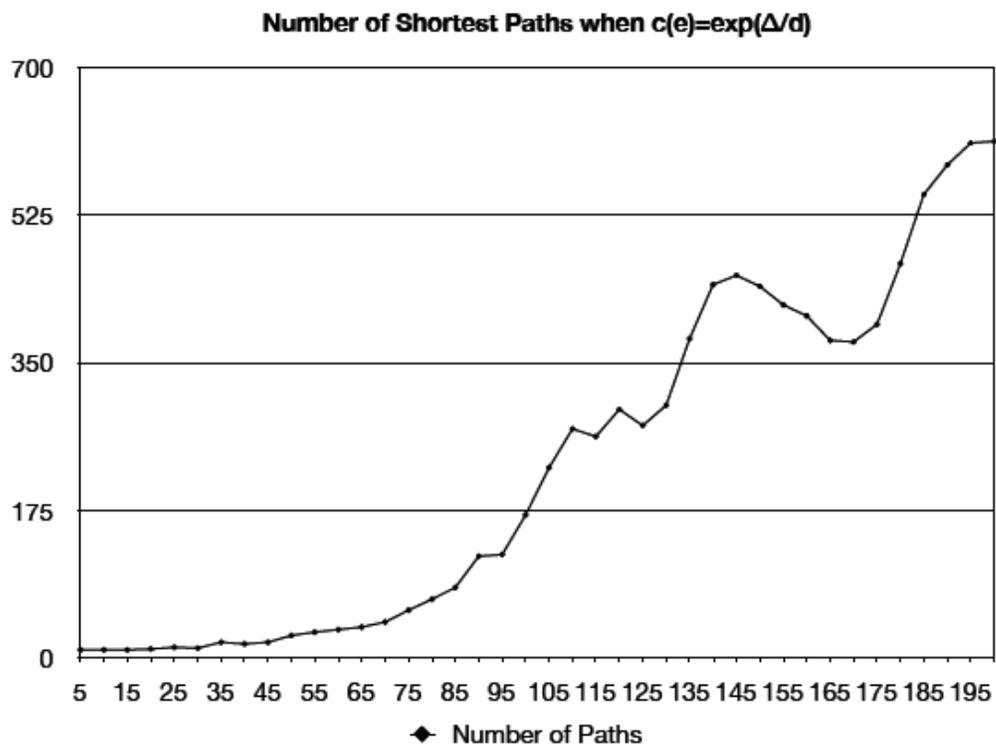


Fig. 1. The number of shortest paths in the read graph as a function of k when the edge cost function is $e^{\frac{\Delta}{k}}$.

Example of Read Graph Construction.

```
Reference: ...AGCGT---GAAGCT--T...
Read1: ...AG-GTG--GAAGCT
Read2: ...AGA-TGGAGAAA-T--T...
Read3: ...AGCGA---GTCG-TAAT...
Read4: ...AGCGA---GTCACT--T...
```

```
Reference: ...AGCGTIIIIGAAGCTIIIT...
Read1: ...AGDGTGIIIGAAGCT
Read2: ...AGADTGGAGAAADTIIIT...
Read3: ...AGCGAIIIGTCGDTAAT...
Read4: ...AGCGAIIIGTCACTIIIT...
```

```
Reference: ...AGCGTIIIIGAAGCTT...
Read1: ...AGNGTGIIIGAAGCT
Read2: ...AGAGTGGAGAAADTT...
Read3: ...AGCGAIIIGTCGDTT...
Read4: ...AGCGAIIIGTCACTT...
```

Fig. 2. Example of indels preprocessing and simple error correction. At the top, multiple reads alignment is given. At the middle, *I* and *D* placeholders are added. At the bottom, deletions, supported by a single read, are corrected, insertions, confirmed by a single read, are removed.

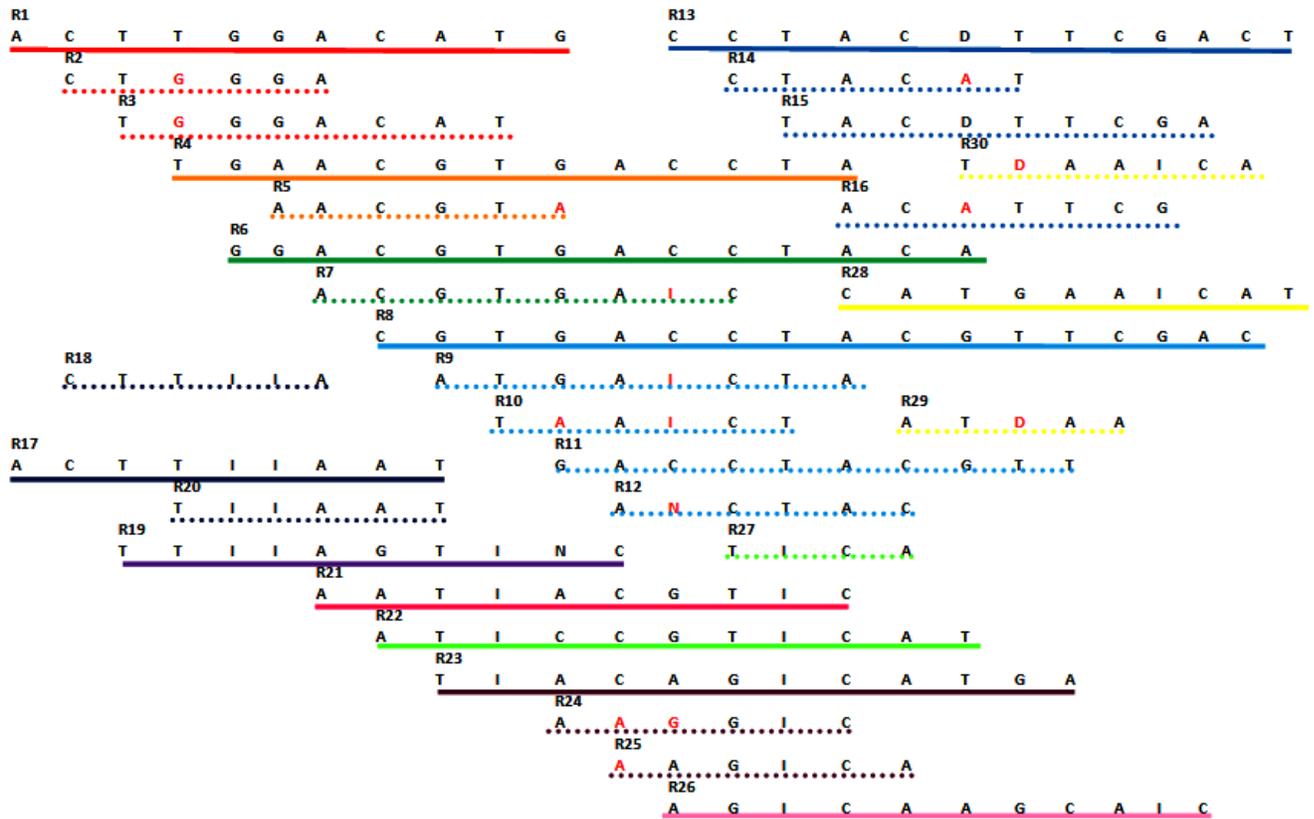


Fig. 3. Example of aligned reads. Reads are divided into superreads and subreads ($n = 2$). Reads underlined by bold lines are superreads. Reads underlined by dashed lines are subreads. Superread and its subreads are underlined by line of the same color. If subread has mismatch with its superread, the mismatch is colored in red. For example, $R2$ and $R3$ are subreads of $R1$ superread, $R2$ and $R3$ have only one difference with $R1$.

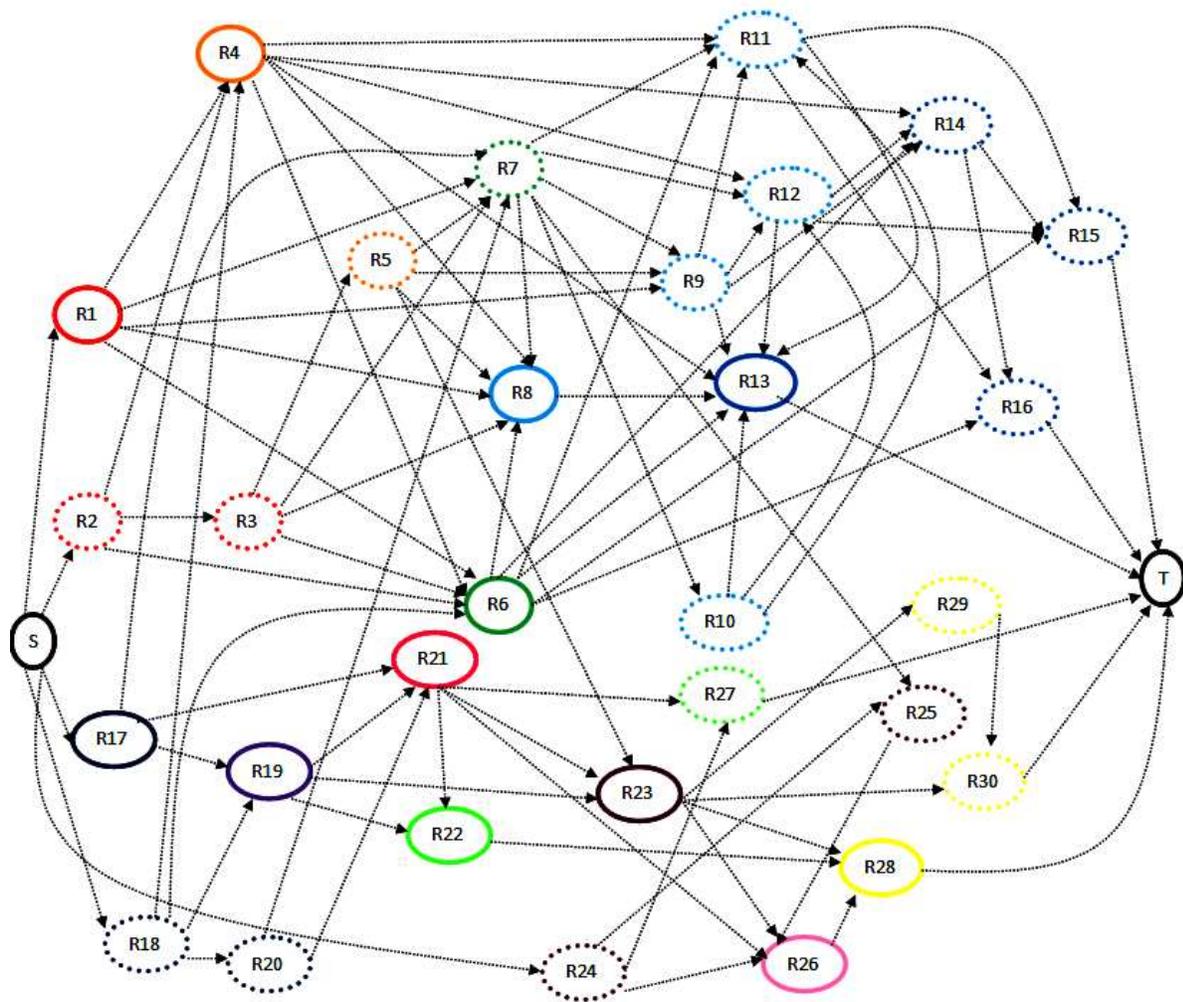


Fig. 4. Example of read graph ($n = 2, m = 2$) if both superreads and subreads are represented by vertices in the read graph. Vertices circled by bold lines correspond to superreads. Vertices circled by dashed lines correspond to subreads. Vertices that correspond to superread and its subreads are circled by the same color.

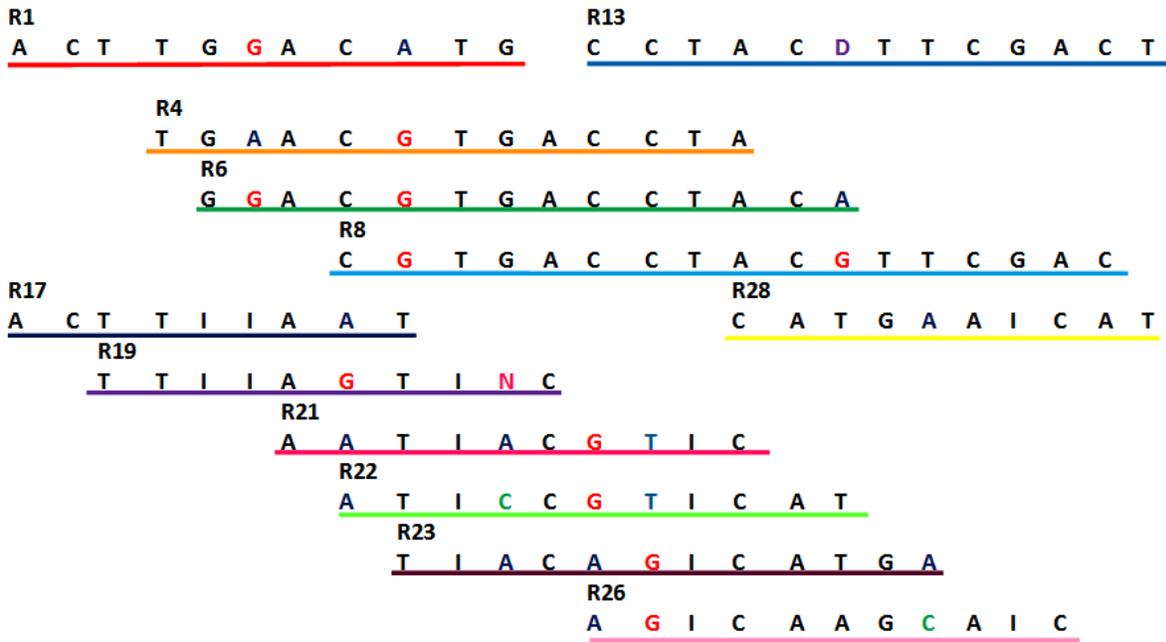


Fig. 5. Example of aligned superreads ($n = 2$). If there are several allelic values across reads in a position, they are marked by different colors to illustrate differences in overlap. Differences are marked across reads $R1 - R13$ and $R17 - R28$, separately. Colors are as follows: red is for "G", dark blue is for "A", green is for "C", blue is for "T", pink is for "N" and purple is for "D". For example, $R1$ and $R4$ have 2 differences in the overlap, $R1$ and $R6$ as well as $R1$ and $R6$ have 1 difference.

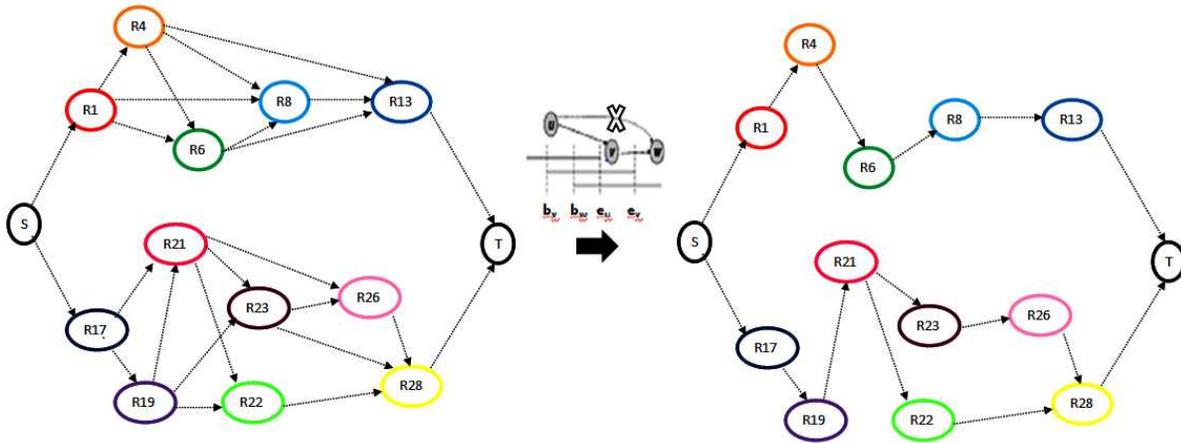


Fig. 6. Example of read graph ($n = 2, m = 2$) if only superreads are represented by vertices in the read graph. Read graph is shown before transitive reduction (at the left) and after transitive reduction was applied (at the right).

2 Data Sets

Reads generated by FlowSim from known HCV quaspecies.

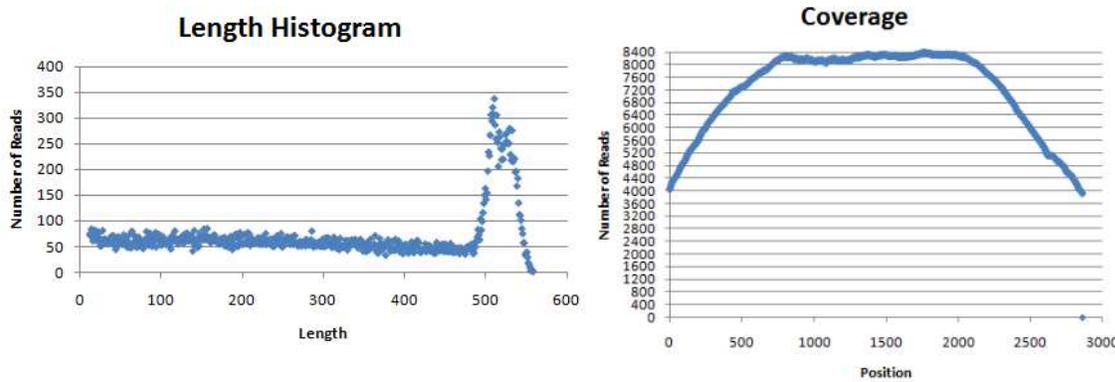


Fig. 7. Length histogram and position coverage for reads generated by FlowSim. Left: number of aligned reads with given length. Right: number of aligned reads covering extended reference positions. Each position is covered by at least 4000 reads, except the position at the very end.

Additional Read Statistics. 99.96% of aligned reads has at least one indel with respect to the reference: 99.97% of deletions and 99.6% of insertions are 1bp long. Only 1.1% of aligned reads has unknown value(s).

454 Pyrosequencing Reads from HCV Samples.

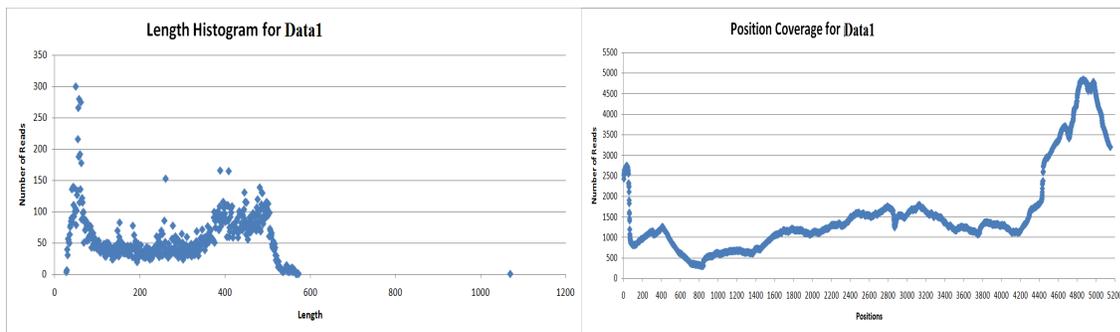


Fig. 8. Left: number of aligned reads with given length. There is a single read 1050bp long, major peak at 450bp long and 50bp long. Right: number of aligned reads covering extended reference positions. The global minimum is 800th nucleotide covered by 310 reads.

Additional Read Statistics. 72% of aligned reads has at least one deletion with respect to the reference: 98% of deletions are 1bp long, 1.5% has length 2, and the rest 0.5% has length 3. 77% of aligned reads has at least one insertion: 86% of insertions have length equalled to 1, and 9.8% have length equalled to 3. Only 7% of aligned reads has at least one unknown value. Assuming that only once encountered insertions caused by typing errors, we found that the insertion error rate is at least 0.025%.

454 Pyrosequencing Reads from HIV Samples.

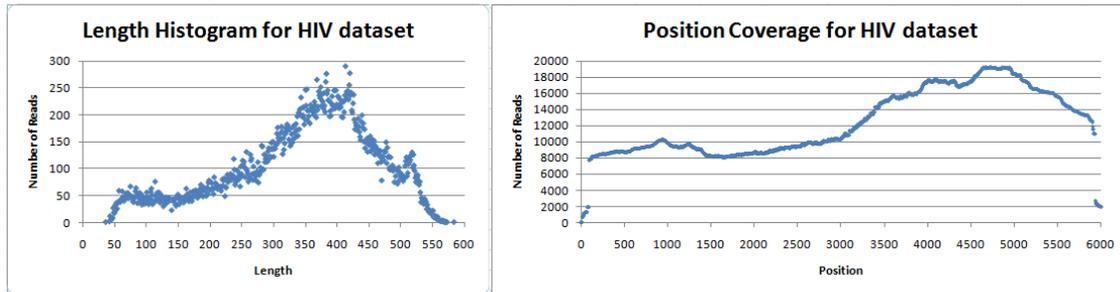


Fig. 9. Left: number of aligned reads with given length. Right: number of aligned reads covering extended reference positions.

Additional Read Statistics. 87% of aligned reads has at least one deletion with respect to the reference: 99.97% of deletions are 1bp long. 99% of aligned reads has at least one insertion: 85% of insertions have length equalled to 1, 10% have length equalled to 2, and 3.5% have length equalled to 3. 11.6% of aligned reads has at least one unknown value.