

Research Article

Ultradeep Pyrosequencing of Hepatitis C Virus Hypervariable Region 1 in Quasispecies Analysis

Kamila Caraballo Cortés,^{1,2} Osvaldo Zagordi,³ Tomasz Laskus,¹ Rafał Płoski,⁴
Iwona Bukowska-Ośko,¹ Agnieszka Pawełczyk,¹ Hanna Berak,⁵ and Marek Radkowski¹

¹ Department of Immunopathology of Infectious and Parasitic Diseases, Medical University of Warsaw,
3c Pawińskiego Street, 02-106 Warsaw, Poland

² Postgraduate School of Molecular Medicine, Żwirki i Wigury 61 Street, 02-091 Warsaw, Poland

³ Institute of Medical Virology, University of Zurich, Winterthurerstrasse, 190 8057 Zurich, Switzerland

⁴ Department of Medical Genetics, Medical University of Warsaw, 3c Pawińskiego Street, 02-106 Warsaw, Poland

⁵ Hospital for Infectious Diseases, 37 Wolska Street, 01-201 Warsaw, Poland

Correspondence should be addressed to Kamila Caraballo Cortés; kcaraballo@wum.edu.pl

Received 22 October 2012; Accepted 12 February 2013

Academic Editor: Ozgur Cogulu

Copyright © 2013 Kamila Caraballo Cortés et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Genetic variability of hepatitis C virus (HCV) determines pathogenesis of infection, including viral persistence and resistance to treatment. The aim of the present study was to characterize HCV genetic heterogeneity within a hypervariable region 1 (HVR1) of a chronically infected patient by ultradeep 454 sequencing strategy. Three independent sequencing error correction methods were applied. First correction method (Method I) implemented cut-off for genetic variants present in less than 1%. In the second method (Method II), a condition to call a variant was bidirectional coverage of sequencing reads. Third method (Method III) used *Short Read Assembly into Haplotypes* (ShoRAH) program. After the application of these three different algorithms, HVR1 population consisted of 8, 40, and 186 genetic haplotypes. The most sensitive method was ShoRAH, allowing to reconstruct haplotypes constituting as little as 0.013% of the population. The most abundant genetic variant constituted only 10.5%. Seventeen haplotypes were present in a frequency above 1%, and there was wide dispersion of the population into very sparse haplotypes. Our results indicate that HCV HVR1 heterogeneity and *quasispecies* population structure may be reconstructed by ultradeep sequencing. However, credible analysis requires proper reconstruction methods, which would distinguish sequencing error from real variability *in vivo*.

1. Introduction

Genetic variability is a characteristic feature of hepatitis C virus (HCV), due to an absence of error correction mechanisms of the viral RNA-dependent RNA polymerase, fast replication, and recombination events [1–3]. As a consequence, HCV displays high intrahost population diversity, forming a pool of closely related but distinct genetic variants (*quasispecies*) [1]. The viral genetic variability is not evenly distributed through the entire genome; the highest variable regions include HVR1, HVR2, and HVR3 of the envelope E2 protein [4]. It is believed that HCV variability has significant clinical implications, since it may result in the generation of immune escape mutants, which may contribute to chronic infection and treatment resistance [5].

The detailed study of the minor variants within the *quasispecies* population is hampered by the absence of sensitive sequencing strategies which would allow for the detection of low-frequency genomes. The traditional method for studying viral *quasispecies* is based on Sanger sequencing of bacterially cloned viral sequences. However, this strategy requires extensive cloning to achieve the desired sensitivity for minor variants detection, a process, that is, costly and time consuming. Another limitation of the Sanger method is its difficulty in sequencing GC-rich regions.

Other studies employed single strand conformational polymorphism (SSCP), an electrophoretic method shown to detect variants constituting as little as 3% of the viral population [6]. However, SSCP it is not informative of the

nature of genetic changes or genetic distance between variants and therefore could not be used for some applications such as investigation of the drug resistance. In addition, in a mixture of heterogenous sequences, certain bands may overlap, underrating viral complexity.

With next-generation sequencing (NGS) platforms, it is now possible to investigate viral *quasispecies* at much greater detail. Their high throughput allows for generation of millions of reads in a single sequencing run, facilitating in-depth sequencing.

NGS can detect variants at low frequencies, which would go undetected by standard sequencing methods [7]. Nevertheless, in order to make reliable reconstruction of the viral *quasispecies* from the noisy, incomplete data obtained by NGS, a proper data analysis is required [8, 9].

In the present study we used ultradeep pyrosequencing (454/Roche) to characterize the complexity and heterogeneity of hypervariable region 1 (HVR1) in a patient persistently infected with HCV genotype 1b. This region was chosen as its protein product is under constant selection pressure of the host immune responses, especially of cytotoxic T cells and neutralizing antibodies [10, 11]. We sequenced this short region at very high coverage, aiming at detecting a large number of minority variants. We took into account sequencing errors in order to have a reliable reconstruction of the viral *quasispecies* on this region.

Reports taking advantage of deep sequencing to investigate HCV genetic diversity for clinical and epidemiological studies are currently available [12–15]. Likewise, there are also works reporting and comparing bioinformatic approaches to infer the viral population from clinical samples, mostly HIV [9, 16–19]. Our study contributes progress in the evaluation of reconstruction methods and extends it for HCV *quasispecies* phenomenon investigation.

2. Patient and Methods

2.1. Sample. A serum sample from a 66-year-old treatment-naive female patient with genotype 1b chronic HCV infection was used. The serum HCV viral load was 1.54×10^6 IU/mL. The patient provided informed consent and the study was approved by the Institutional Bioethical Committee.

2.2. HVR1 Amplification. Viral RNA was extracted from 250 μ L of serum by a modified guanidinium thiocyanate-phenol/chlorophorm method using a commercially available Trizol reagent (Invitrogen) and suspended in 20 μ L of water. Five μ L of the solution containing RNA was subjected to reverse transcription at 37°C for 30 minutes using AccuScript High Fidelity Reverse Transcriptase (Stratagene). HVR1 sequences were amplified in a two-step PCR using FastStart High Fidelity Taq DNA Polymerase (Roche) as described previously [20]. Primers used for reverse transcription (E2 AS) and first round HCV HVR1 amplification (E2 S) were as follows: 5'-CATTGCAGTTCAGGGCCGTGCTA-3' and 5'-GGTGCTCACTGGGGAGTCCT-3'. Primers for the second round PCR (E2 NS and E2 NAS) were as follows: 5'-CGT ATC GCC TCC CTC GCG CCA TCAG

TCC ATG GTG GGG AAC TGG GC-3' and 5'-CTA TGC GCC TTG CCA GCC CGC TCAG TGC CAA CTG CCA TTG GTG TT-3'. The latter contained tags recognized by GS Junior Sequencing System (underlined).

2.3. SSCP Analysis of HVR1 Quasispecies. Second round PCR product was purified using Wizard SV Genomic DNA Purification System (Promega) and resuspended in 20 μ L of water. Next, 2–5 μ L of purified PCR product was subjected to SSCP assay as described previously [21]. Complexity of a population was reflected by the number of distinct bands.

2.4. Ultradeep Pyrosequencing. Pyrosequencing was carried out according to the manufacturer's protocol for amplicons using GS Junior System (454/Roche). In order to lower contamination with short sequences (i.e., primer residues), HVR1 product of the second round PCR was purified from agarose gel by QIAquick Gel Extraction Kit (Qiagen). The extracted product was measured fluorometrically using Quant-iT PicoGreen dsDNA Assay Kit (Molecular Probes), and the amount of DNA equivalent to 3×10^7 copies was subjected to emulsion PCR using GS Junior Titanium emPCR Kit (Lib-A). Pyrosequencing was performed according to the amplicon processing procedure for 100 cycles (recommended for amplicons up to 250 bp).

2.5. Data Analysis. Reads that did not match primer sequences or had undetermined bases (Ns) were excluded from further analysis. Retained sequences of 179 bp were visualized using GS Amplicon Variant Analyzer (Roche). Subsequently, primer sequences were trimmed from the target sequence and reads of 138 bp were aligned to the reference sequence for genotype 1b HCV (GenBank accession number AJ406073) and translated to amino acid sequences by (*Molecular Evolutionary Genetics Analysis*) MEGA, version 5.0 (<http://www.megasoftware.net/>) [22]. Phylogenetic analyses were conducted in MEGA5 using the Maximum Likelihood method based on the Tamura-Nei model [23] using MEGA 5.0 software. Genetic parameters such as genetic diversity and sequence polymorphisms within sequences were 5 assessed by DNA SP version (<http://www.ub.edu/dnasp/>). The program *diri_sampler* from the ShoRAH software was used to correct sequencing errors and infer haplotypes. Given the high number of reads obtained in the sequencing, the dataset was split equally in two, and the obtained sets were analyzed independently. Error correction included mismatches as well as insertions and deletions.

3. Results

3.1. Amplification and Sequencing Errors. As our experiment used RT-PCR-amplified material, we attempted to assess the error rate in the consecutive experimental steps taking into account error rates of employed enzymes. For reverse transcription, AccuScript High Fidelity Reverse Transcriptase (Stratagene) was used, which displays three times higher fidelity than commonly used MMLV reverse transcriptase [24]. The estimated AccuScript RT error rate is 2×10^{-5}

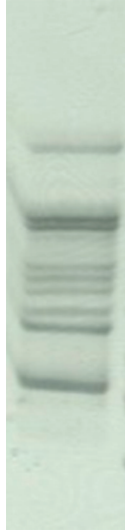


FIGURE 1: The SSCP image of HVR1 amplified from the serum of HCV-infected patient.

(manufacturers data). For PCR amplification, we used Fast-Start High Fidelity Taq DNA Polymerase (Roche), which has estimated error rate of 2×10^{-6} (three times lower than Taq DNA polymerase) [25]. Finally, the pyrosequencing error rate is estimated to be 1.07%, including mismatches (0.088%), insertions (0.541%), deletions (0.359%), and ambiguous base calls (0.085%) [26].

Studying clonal samples, or control samples where a set of clones are mixed in predetermined proportions are important to evaluate the error rate of the sequencing process and the performance of the haplotype reconstruction methods. Since these have already been reported elsewhere [16, 26], it seems not requisite to perform these experiments for every new study of the viral *quasispecies*.

3.2. Heterogeneity of HCV HVR1 Viral Variants Assessed by SSCP Analysis. Based on gel analysis, at least nine SSCP bands of HVR1 were observed (Figure 1). The frequency was not uniform across variants, as could be seen by the different intensities of the bands.

3.2.1. Heterogeneity of HCV HVR1 Assessed by Ultra-deep Sequencing. To check the applicability of ultra-deep pyrosequencing for HCV HVR1 heterogeneity analysis, the amplified product was sequenced by GS Junior System (454/Roche). Based on the data of GS Amplicon Variant Analyzer, the total number of sequenced nucleotides was 1.37×10^8 . The system read 76 332 individual sequences, among them 73 236 (95.9%) (28 098 forward and 45 138 reverse) were aligned to the reference sequence AJ406073 of genotype 1b HCV. The GS Amplicon Variant Analyzer software detected 15 917 haplotypes. Mean coverage of each variant (expressed by the number of identical reads) was 4.6. The most abundant haplotype coverage was 4540 reads. The rarest haplotypes comprised single sequence reads (74,6% of detected haplotypes). Our results are summarized in Table 1.

TABLE 1: HVR1 HCV characteristics obtained by pyrosequencing using GS Junior System (454/Roche).

Number of sequenced nucleotides	1.37×10^8
Number of individual sequences that passed the quality control*	76 332
Number of individual sequences aligned to reference genome	73 236
Mean coverage per sequence	4.6
Identified haplotypes	15 917

*No undetermined bases, 100% match with primer sequences.

3.2.2. Error Correction in Haplotype Reconstruction. In order to reflect the HVR1 HCV population *in vivo* as accurately as possible, we explored the effect of different strategies to take the sequencing error rate into account. In a very conservative approach (Method I), we only considered variants detected at a frequency higher than 1%. This amounts to discard most variants, even if, given the high coverage, they appear in hundreds of reads. With this strategy we only retained 8 haplotypes.

A second strategy (Method II) consisted in requiring bidirectional coverage, that is, in only retaining variants supported by at least one forward and one reverse read. This method identified 40 HVR1 variants.

In the third approach we used the program *diri_sampler* from the software suite ShoRAH [16]. In this analysis, inference of the viral *quasispecies* is done in probabilistic manner using a Bayesian approach. It does not rely on the input of an error rate, rather, it estimates it from the sequencing data. Reads are clustered together and the consensus sequence of each cluster represents the original haplotype. Together with the frequency of each variant, the program enables assessments of the posterior probability of each haplotype, a confidence value for their existence. The number of diverse reads sequenced was higher than what the program can handle on a desktop computer with 4 GB of RAM. In order to face this limitation, we split the reads equally in two subsets, and performed haplotype reconstruction independently. Only haplotypes with confidence value >95% were retained. As an additional measure of reliability, only haplotypes supported by at least 5 reads were included. Since we are dealing with a coding sequence, frameshift inducing insertions/deletions were resolved correcting to the most common nucleotide for that position. As a result of two independent computations on raw data halves, two populations (A and B), consisting of 333 and 315 haplotypes respectively, were obtained. Their frequencies varied from 10.54% and 10.44% (the most abundant variants in population A and B, resp.) down to 0.013% and 0.014% (the least abundant variants in population A and B, resp.). 186 haplotypes were common to both populations and their frequencies were all above 0.02%. Seventeen haplotypes were present with a frequency >1%, constituting in total 58.6% of the entire population.

3.2.3. Characteristics of Inferred HVR1 Populations. After application of error correction methods, such parameters as

TABLE 2: The impact of haplotype reconstruction method on the variability parameters of HVRI.

Correction method	I Cut-off 1%	II Bidirectional coverage	III ShoRAH
Number of haplotypes	8	40	186
Number of nucleotide substitutions within HVRI	51	59	70
Percentage of mutated amino acid positions within HVRI (%)	55.6	55.6	74.1
Genetic distance	3.874	0.065	0.110
Genetic diversity	0.923	0.998	0.984

percentage of mutated amino acid positions, genetic distance, genetic diversity as well as number of substitutions were calculated (Table 2). The highest genetic distance characterized population reconstructed by cut-off method (3.874) followed by ShoRAH method (0.110) and bi-directional coverage method (0.065), whereas genetic diversities were similar for all populations (0.923, 0.998 and 0.984 for method I, II and III, resp.). The highest number of nucleotide substitutions was detected in ShoRAH-reconstructed population (overall 70). 47 (67%) of them were present in genetic variants constituting more than 1% of the entire population.

HVRI populations were also compared on amino acid level (Figure 2). Within 27 amino acid stretch of HVRI, only 15 (55.5%) positions were polymorphic after application of methods I and II, and 20 (74.1%) after ShoRAH computations. Based on ShoRAH computation results, the most variable was the fourth HVRI position, where 11 amino acid substitutions were detected when compared to reference sequence (V/D, V/M, V/T, V/L, V/R, V/A, V/E, V/G, V/N, V/I, V/Q).

Viral populations were also analyzed phylogenetically. As shown in Figure 3, the general topology of three populations was similar. However, the tree topology based on ShoRAH computation was the most extensive.

4. Discussion

Pyrosequencing is a relatively novel technique which may help to decipher complex viral populations in terms of their diversity and structure. To date, it was successfully used in human immunodeficiency virus (HIV) research to identify minor drug resistant variants, analyze variable regions of heavy and light chains of neutralizing antibodies against HIV, as well as to determine HIV tropism, analyze superinfections and assess diversity of genital microbiota in HIV-infected women [27–31]. Ultradeep sequencing strategies also offers a new approach in HCV research. However, application of this method requires that several issues are taken into account. The foremost of these is the generation of mutations during reverse transcription and amplification reactions, due to enzyme errors [32]. Reverse transcriptase is the most error-prone, as it lacks a proofreading activity. For instance, error rate of common reverse transcriptases used *in vitro* to synthesize cDNA is at least 10^{-4} [24], and errors that occurred during this step are propagated during the subsequent PCR amplification. In the present study, in order to minimize errors, high fidelity enzymes were used in amplification reactions preceding sequencing (AccuScript High Fidelity

Reverse Transcriptase and FastStart High Fidelity Taq DNA Polymerase). Nevertheless, the resulting hypothetical error rate of amplification is estimated to be lower than the sequencing error rate itself. The sequencing step introduces various types of errors related to the pyrosequencing chemistry and detection technology. The major contributor to errors is the ambiguity of homopolymer length, which results from the difficulty to resolve intensity of luminescence when a homopolymer is encountered. Moreover, insufficient flushing may lead to single base insertions. Overall, it was estimated that the mean error rate of pyrosequencing (defined as the number of errors such as miscalled bases or inserted or deleted bases divided by the total number of sequenced bases) was 1.07% [26]. This value may be considered as the experimentally confirmed resolution of the method. For the above reasons, the raw data obtained from sequencing should be additionally processed in order to remove low-quality reads and reads containing errors.

Three different error correction methods were applied to the raw sequencing data, which resulted in three HVRI populations, differing in complexity and heterogeneity. The most sensitive was ShoRAH program reconstruction, which allowed to obtain the broadest spectrum of HVRI sequences. This method has already been shown to reliably detect variants down to about 0.1% [33]. In this study we detected variants down to 0.02%, confirmed in two independent computations. The cut-off method, in which variants present in less than 1% of the population were discarded, was the least sensitive, as it allowed to detect only 8 haplotypes. Similar cut-off was applied in analysis of pyrosequencing reads of *pol/gag* of HIV population as well as *PePHD E2* of HCV [14, 34]. It was reported that this method may result in inadequate haplotype reconstruction of low precision, low recall, or both, depending on the cut-off value. Too low cut-off value may result in low precision (fraction of true haplotypes among all called haplotypes) and conversely, high cut off may significantly lower recall (fraction of called haplotypes among all true haplotypes). For instance, based on the analysis of *gag/pol* HIV genes, it was shown that the cut-off of 50 read observations resulted in 80% precision but only 40% recall [33]. Application of the bi-directional coverage correction method II allowed us to determine the presence of 40 haplotypes, but such verification is laborious and raises concern regarding the acceptance of haplotypes characterized by high disproportion in forward and reverse strand counts. Among these forty sequences with bidirectional coverage, we could identify twenty that matched exactly one of the

#seq_34	frequency_0.566364E..TI	..V.RTTS	LSG..RA..H	..I...
#seq_35	frequency_0.539622DP.DR	..V.RTTS	LSG..RA..Q	..I...
#seq_36	frequency_0.537195E..TI	..V.RTTS	..SG..RA..H	..I...
#seq_37	frequency_0.521421E..TI	..V.RTTS	LSG..RA..Q	..I...
#seq_38	frequency_0.511224DS.MI	..SV..G.R	LSS..TA..Q	..I...
#seq_39	frequency_0.506563DS.MI	..E.R..S	LSG..TR..Y	..I...
#seq_40	frequency_0.489579DS.DR	..E.R..S	LSG..TR..Y	..I...
#seq_41	frequency_0.477541E..TI	..V.RTTS	LSG..RA..H	..I...
#seq_42	frequency_0.475100E..TI	..SV..G.R	LSS..TA..Q	..I...
#seq_43	frequency_0.416991DP.DR	..SV..G.R	LSS..TA..Q	..I...
#seq_44	frequency_0.403574DS.DR	..E.R..G	LSA..TR..Y	..I...
#seq_45	frequency_0.398722DS.DR	..E.R..S	LSG..TR..Y	..I...
#seq_46	frequency_0.386373E..TI	..SV..G.R	LSS..TA..Q	..I...
#seq_47	frequency_0.375498DP.DR	..V.RTTS	LSG..RA..H	..I...
#seq_48	frequency_0.368952E..TI	..V.RTTS	LSG..TA..Q	..I...
#seq_49	frequency_0.353843E..TI	..V.RTTS	..SS..RA..H	..I...
#seq_50	frequency_0.350268E..PN.DR	..SV..G.R	LSS..TA..Q	..I...
#seq_51	frequency_0.346676E..PN.DR	..SV..G.R	LSS..TA..Q	..I...
#seq_52	frequency_0.336838DP.DR	..V.RTTS	LSG..RA..H	..I...
#seq_53	frequency_0.332207E..TI	..V.RTTS	LSG..RA..H	..I...
#seq_54	frequency_0.330468DP.DR	..V.RTTS	..SG..RA..Q	..I...
#seq_55	frequency_0.325899DS.DR	..E.R..S	LSG..TA..Q	..I...
#seq_56	frequency_0.324684DP.DR	..V.RTTS	..SS..RA..H	..I...
#seq_57	frequency_0.323406E..TI	..SV..G.R	LSS..TA..Q	..I...
#seq_58	frequency_0.322620DS.MS	..E.R..S	LSG..TR..Y	..I...
#seq_59	frequency_0.291277DS.MI	..E.R..S	LSG..TA..Q	..I...
#seq_60	frequency_0.290212DS.DR	..A-E.R..S	LSG..RA..H	..I...
#seq_61	frequency_0.265723DS.MI	..E.R..G	LSA..TR..Y	..I...
#seq_62	frequency_0.256780DP.DR	..V.RTTS	LSG..RA..H	..I...
#seq_63	frequency_0.243733DS.MI	..E.R..S	LSG..RA..H	..I...
#seq_64	frequency_0.232532DP.DR	..V.RTTS	..SG..TA..Q	..I...
#seq_65	frequency_0.232330E..TI	..V.RTTS	..SS..TA..Q	..I...
#seq_66	frequency_0.231031E..PN.DR	..SV..G.R	LSS..TA..Q	..I...
#seq_67	frequency_0.226299P..TI	..SV..G.R	LSS..TA..Q	..I...
#seq_68	frequency_0.226212E..TI	..V.RTTS	LSG..RA..H	..I...
#seq_69	frequency_0.226001E..TI	..V.RTTS	LSS..TA..Q	..I...
#seq_70	frequency_0.214432P..TI	..S...TR	..SS..TL..Q	..I...
#seq_71	frequency_0.198007E..TI	..V.RTTS	..SG..TA..Q	..I...
#seq_72	frequency_0.189106E..PN.DR	..SV..G.R	LSS..TA..Q	..I...
#seq_73	frequency_0.171091DP.DR	..V.RTTS	..SS..TA..Q	..I...
#seq_74	frequency_0.161615DP.DR	..V.RTTS	LSG..RA..H	..I...
#seq_75	frequency_0.158776E.RAI	..E.R..S	LSG..TR..Y	..I...
#seq_76	frequency_0.156317K..TI	..SV..G.R	LSS..TA..Q	..I...
#seq_77	frequency_0.148586DS.DR	..RE.R..S	LSG..TR..Y	..I...
#seq_78	frequency_0.147589E..LTLI	..SV..G.R	LSS..TA..Q	..I...
#seq_79	frequency_0.144746DP.DR	..V.RTTS	LSS..TA..Q	..I...
#seq_80	frequency_0.143410K..TI	..SV..G.R	LSS..TA..Q	..I...
#seq_81	frequency_0.141619DS.DR	..V.RTTS	LSG..RA..H	..I...
#seq_82	frequency_0.140555R..TI	..SV..G.R	LSS..TA..Q	..I...
#seq_83	frequency_0.138271DP.DR	..V..G.R	LSS..TA..Q	..I...
#seq_84	frequency_0.132354E..G..TI	..SEP.G.R	LSS..TA..Q	..I...
#seq_85	frequency_0.127754DP.DR	..V.RTTS	LSG..TA..Q	..I...
#seq_86	frequency_0.123660E..TI	..V.RTH	LSG..RA..H	..I...
#seq_87	frequency_0.121928G..TI	..SV..G.R	LSS..TA..Q	..I...
#seq_88	frequency_0.120898DS.MI	..SV..G.R	LSS..TA..Q	..I...
#seq_89	frequency_0.118086DS.EV	..E.R..S	LSG..TR..Y	..I...
#seq_90	frequency_0.117558E..TI	..V.RTTS	..SG..RA..H	..I...
#seq_91	frequency_0.117161R..TI	..SV..G.R	LSS..TA..Q	..I...
#seq_92	frequency_0.112927P..TI	..S...TR	..SS..TA..Q	..I...
#seq_93	frequency_0.112631E..TI	..V.RTTS	LSG..RA..H	..I...
#seq_94	frequency_0.112472DS.MI	..V.RTTS	LSG..RA..H	..I...
#seq_95	frequency_0.112397Q..RTI	..E.R..S	LSG..TR..Y	..I...
#seq_96	frequency_0.111663DP.DR	..RV.RTTS	..SG..RA..H	..I...
#seq_97	frequency_0.107630E..TI	..V.RTTS	LSG..RA..H	..I...
#seq_98	frequency_0.106809E..TI	..V.RTTS	..S.G..RA..H	..I...
#seq_99	frequency_0.103442E..PN.DR	..V.RTTS	LSG..RA..H	..I...
#seq_100	frequency_0.100934E..TI	..V.RTTS	..SG..RA..Q	..I...
#seq_101	frequency_0.100065E..PTRV	..SV..G.R	LSS..TA..Q	..I...
#seq_102	frequency_0.099957E..PDTI	..SV..G.R	LSS..TA..Q	..I...
#seq_103	frequency_0.099954E..TI	..V.RTTS	LSG..TR..Y	..I...
#seq_104	frequency_0.097549E..PN.DR	..V.RTTS	LSG..RA..H	..I...
#seq_105	frequency_0.095057EW.DS.DR	..E.R..S	LSG..TR..Y	..I...
#seq_106	frequency_0.095048G..TI	..SV..G.R	LSS..TA..Q	..I...
#seq_107	frequency_0.094927DS.DR	..E..G.R	LSS..TA..Q	..I...
#seq_108	frequency_0.093505DP.DR	..SV..G.R	LSS..TA..Q	..I...
#seq_109	frequency_0.090986DP.DR	..V.RTTS	LSG..TR..Y	..I...
#seq_110	frequency_0.090194E..TI	..E.R..S	LSG..TR..Y	..I...
#seq_111	frequency_0.089734E..G..TI	..SV..G.R	LSS..TR..Y	..I...
#seq_112	frequency_0.089089DP.DR	..E.R..S	LSG..TR..Y	..I...
#seq_113	frequency_0.088302DS.MI	..E.R.G*	LSG..TR..Y	..I...
#seq_114	frequency_0.081532K..TI	..SV..G.R	LSS..TA..Q	..I...
#seq_115	frequency_0.081455E..PTQI	..SV..G.R	LSS..TA..Q	..I...
#seq_116	frequency_0.079521GR.P..TI	..V.RTTS	..SG..RA..H	..I...
#seq_117	frequency_0.079391E..TI	..SV..G.R	LSS..TA..Q	..I...
#seq_118	frequency_0.078571E..TI	..R*RTTS	LSG..RA..H	..I...
#seq_119	frequency_0.076269DS.G*	..PE.R..S	LSG..TR..Y	..I...
#seq_120	frequency_0.074612DP.DR	..V.RTTS	..S.G..RA..H	..I...
#seq_121	frequency_0.074134DP.EV	..SV..G.R	LSS..TA..Q	..I...
#seq_122	frequency_0.074079E..TI	..R*PRTTS	LSG..RA..H	..I...
#seq_123	frequency_0.073492E..G..TI	..SE..G.R	LSS..TA..Q	..I...
#seq_124	frequency_0.072251DS.G*	..E.R..S	LSG..TR..Y	..I...
#seq_125	frequency_0.071245E..G..TI	..SE..G.R	LSS..TA..Q	..I...
#seq_126	frequency_0.070688E..G..TI	..SE..G.R	LSS..TA..Q	..I...
#seq_127	frequency_0.070584DP.DR	..R.TS	..SG...V..H	..I...
#seq_128	frequency_0.070491DS.DR	..E.R..G	LSA..TR..Y	..I...
#seq_129	frequency_0.069469E..TI	..R	LTG..TL..Q	..I...
#seq_130	frequency_0.065577E..G..TI	..SV..G.R	LSS..TR..Y	..I...
#seq_131	frequency_0.065258E..TI	..RV.RTTS	LSG..RA..H	..I...
#seq_132	frequency_0.064382DS.MI	..E..G.R	LSS..TA..Q	..I...
#seq_133	frequency_0.064154E..TI	..V.RTTS	..SG..RA..H	..I...
#seq_134	frequency_0.063206E..TI	..EPRTTS	LSG..RA..H	..I...
#seq_135	frequency_0.061966R..TI	..SV..G.R	LSS..TA..Q	..I...

(c)

FIGURE 2: Continued.

```

#seq_136_frequency_0.061551 .....DP,DR..-V,RTTS..SG..TA..Q..I...
#seq_137_frequency_0.059904 .....E..G..TI..SEP,G,R..LSS..TA..Q..I...
#seq_138_frequency_0.059792 .....DP,DR..-V,RTTS..SG..RA..Q..I...
#seq_139_frequency_0.059624 .....DP,DR..-V,RTH..LSG..RA..H..I...
#seq_140_frequency_0.059526 .....DPDTI..SV..G,R..LSS..TA..Q..I...
#seq_141_frequency_0.059048 .....E..G..TI..E,R..S..LSG..TR..Y..I...
#seq_142_frequency_0.057471 .....DP,DR.....R..LTG..TA..Q..I...
#seq_143_frequency_0.057002 .....L,LIG..-T..RTTS..SN..KL..Q..I...
#seq_144_frequency_0.054072 .....DS,DR..A-E,R..S..LSG..RA..H..I...
#seq_145_frequency_0.051865 .....E..G..TI..V,RTTS..LSG..RA..H..I...
#seq_146_frequency_0.051459 .....DS,MI..V,RTTS..SG..RA..Q..I...
#seq_147_frequency_0.050898 .....DP,DR.....R..LTG..TL..Q..I...
#seq_148_frequency_0.049974 .....E..PN,DR..-E,R..S..LSG..TR..Y..I...
#seq_149_frequency_0.049031 .....E..TI.....R..LTG..TA..Q..I...
#seq_150_frequency_0.048870 .....E..G..TI..V,RTTS..SG..RA..Q..I...
#seq_151_frequency_0.047691 .....GW,DSDMI..E,R..S..LSG..TR..Y..I...
#seq_152_frequency_0.046979 .....DP,DR..-E,R..S..LSG..TR..Y..I...
#seq_153_frequency_0.046820 .....E..TI.....R,T,S..SG..V..H..I...
#seq_154_frequency_0.046768 .....E..-PTRV..SV..G,R..LSS..TA..Q..I...
#seq_155_frequency_0.046693 .....E..TI..V,RTTS..SS..TA..Q..I...
#seq_156_frequency_0.045850 .....E..G..TI..SV..G,R..LSS..RA..H..I...
#seq_157_frequency_0.044736 .....DS,MI..V,RTTS..LSG..RA..H..I...
#seq_158_frequency_0.044263 .....DS,MI..E,R..S..LSG..TA..Q..I...
#seq_159_frequency_0.044067 .....-R..TI..SV..G,R..LSS..TA..Q..I...
#seq_160_frequency_0.043337 .....-LALI..E,R..S..LSG..TR..Y..I...
#seq_161_frequency_0.043222 .....E..G..TI..SEPTG,R..LSS..TA..Q..I...
#seq_162_frequency_0.042263 .....-R..TI..SV..G,R..LSS..TA..Q..I...
#seq_163_frequency_0.042250 .....E..-PTRV..SV..G,R..LSS..TA..Q..I...
#seq_164_frequency_0.041116 .....E..TI.....R..LTG..TR..Y..I...
#seq_165_frequency_0.040595 .....E..TI.....EPGTT,S..SG..RA..H..I...
#seq_166_frequency_0.040021 .....DS,DR..-E,R..S..LSG..TA..Q..I...
#seq_167_frequency_0.038500 .....DS,EV..E,R..S..LSG..TR..Y..I...
#seq_168_frequency_0.038116 .....E..G..TI..SD..G,R..LSS..TA..Q..I...
#seq_169_frequency_0.037373 .....EW,DS,DR..A-E,R..S..LSG..RA..H..I...
#seq_170_frequency_0.036679 .....DS,DR..V,RTTS..LSG..RA..H..I...
#seq_171_frequency_0.036009 .....D,E..TI..V,RTTS..SG..RA..H..I...
#seq_172_frequency_0.035691 .....DP,NR..V,RTTS..LSG..RA..H..I...
#seq_173_frequency_0.035587 .....DP,DR..V,RTTS..SG..TR..Y..I...
#seq_174_frequency_0.034870 .....E..G..TI..V,RTTS..SG..RA..Q..I...
#seq_175_frequency_0.033978 .....DS,MI..E,R..G..LSA..TR..Y..I...
#seq_176_frequency_0.031252 .....E..TI..E,R..S..LSG..RA..H..I...
#seq_177_frequency_0.026769 .....-Q,RTI..E,R..S..LSG..TR..Y..I...
#seq_178_frequency_0.026603 .....E..TI..E,R..S..LSG..TA..Q..I...
#seq_179_frequency_0.025699 .....DP,DR.....R,T,S..SG..TA..Q..I...
#seq_180_frequency_0.025357 .....E..TI..S...TR..SS..TL..Q..I...
#seq_181_frequency_0.024865 .....DP,DR.....G,R..LSS..TA..Q..I...
#seq_182_frequency_0.024580 .....DS,G*..E,R..S..LSG..RA..H..I...
#seq_183_frequency_0.024399 .....DS,DR..-T..RTTS..SN..KL..Q..I...
#seq_184_frequency_0.023851 .....P..TI..V,RTTS..LSG..RA..H..I...
#seq_185_frequency_0.021642 .....DP,DR..-E,R..S..LSG..TA..Q..I...
#seq_186_frequency_0.019821 .....DP,DR..-E,R..S..LSG..TR..Y..I...

```

(c)

FIGURE 2: Amino acid sequences of HVRI populations inferred after the application of three different error correction methods. (a) Cut-off method >1% (I), (b) bidirectional coverage (II), and (c) ShoRAH computation (III). Top sequence corresponds to reference sequence AJ406073 for genotype 1b HCV. Dots indicate consensus positions. Dashes indicate positions not present in the sequence. Asterisks indicate stop codons.

sequences obtained with ShoRAH. This is probably a very precise dataset, although, for example, some sequences might have a very biased forward/reverse read ratio and yet be included in it. Undoubtedly, using the strand information while reconstructing haplotypes is a strategy worth pursuing. A promising approach seems to be a proper statistical treatment of the strand bias, implemented together with the error correction of ShoRAH (McElroy, unpublished data).

In our study, using two independent ShoRAH computations, 186 haplotypes were reconstructed. In contrast, in the study of Bull et al. [13], 100 E2 variants were detected by ShoRAH. However, patients in that study were in an early phase of HCV infection and the study was performed along the whole genome, at much lower coverage, and thus HVRI complexity could have been lower [13]. Based on ShoRAH reconstruction, we found that the most frequent HVRI variants constitute a relatively small percentage of the entire population. Thus, the most abundant variant constituted 10.5%, and only 17 haplotypes were present in proportions higher than 1%. These data suggest that during the chronic phase of infection the *quasispecies* population is highly dispersed into minor variants, with no predominant

sequences present. Similar haplotype frequency distribution was observed in foot-and-mouth virus population reconstructed by next generation genome sequencing [35]. In the only other HCV study investigating this issue, two to five variants were detected in frequency higher than 2.5%, whereas we detected eight such sequences. However, as already mentioned, viral samples in that study were drawn in the acute phase of infection [13]. The highest number of substitutions (70) within HVRI was detected in population reconstructed by ShoRAH. Importantly, 47 (67%) of these were detected in variants constituting more than 1%. In contrast, during the acute phase of infection, less than 50% of substitutions were detected in variants present in more than 1% [13]. This suggests that, during the acute phase of infection, rare variants contribute more to the population diversity than during the chronic infection.

SSCP analysis, which has the sensitivity limit of 3%, revealed the presence of 9 bands. If the same cut-off value were applied in ShoRAH-based reconstruction, 7 haplotypes would have been identified. This fact suggests that the SSCP sensitivity might be higher than expected.

However, it must be stressed that the absolute SSCP band number may not reflect the haplotype number as it represents

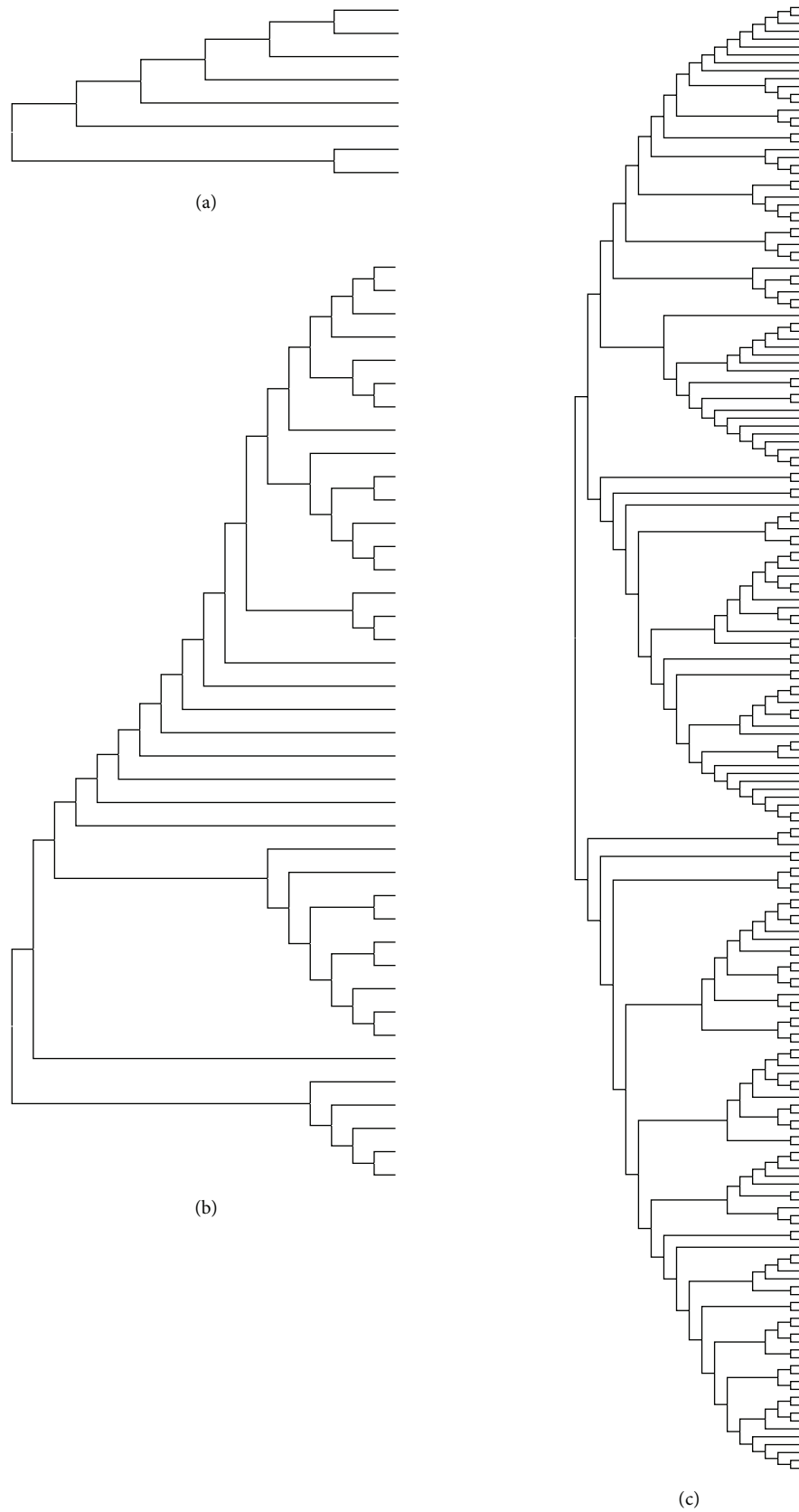


FIGURE 3: Molecular phylogenetic analysis of HVR1 populations inferred after the application of three different error correction methods. (a) Cut-off method >1% (I), (b) bidirectional coverage method (II), and (c) ShoRAH algorithm (III). The evolutionary history was inferred by using the Maximum Likelihood method based on the Tamura-Nei model [23]. Evolutionary analyses were conducted using MEGA 5.0 [22].

resolved single DNA strands. Importantly, only sequencing provides information about the nucleotide sequence.

5. Conclusions

The newly available pyrosequencing technique opens a new approach to the analysis of complex viral genomes as it allows for detection of rare molecular variants. Better understanding of population genetics of complex viral populations seems crucial for understanding *quasispecies* phenomenon, viral evolution, and drug resistance.

In the evaluation presented here, we used ShoRAH to obtain the broadest spectrum of HVR1 variants while trying to preserve their reliability. The use of different sequencing platforms, the optimization of library preparation, and data analysis will further improve the reconstruction of viral *quasispecies*.

Conflict of Interests

The authors declare no conflict of interests.

Acknowledgment

This work was supported by Projects from Polish National Science Centre number NN 401 64 67 40 and 1 M24/PM12/12.

References

- [1] M. Martell, J. I. Esteban, J. Quer et al., "Hepatitis C virus (HCV) circulates as a population of different but closely related genomes: quasispecies nature of HCV genome distribution," *Journal of Virology*, vol. 66, no. 5, pp. 3225–3229, 1992.
- [2] F. Kurbanov, Y. Tanaka, D. Avazova et al., "Detection of hepatitis C virus natural recombinant RFL2k/1b strain among intravenous drug users in Uzbekistan," *Hepatology Research*, vol. 38, no. 5, pp. 457–464, 2008.
- [3] P. Moreno, M. Alvarez, L. Lápez et al., "Evidence of recombination in Hepatitis C Virus populations infecting a hemophiliac patient," *Virology Journal*, vol. 6, article 203, 2009.
- [4] J. M. Cuevas, M. Torres-Puente, N. Jiménez-Hernández et al., "Refined analysis of genetic variability parameters in hepatitis C virus and the ability to predict antiviral treatment response," *Journal of Viral Hepatitis*, vol. 15, no. 8, pp. 578–590, 2008.
- [5] E. A. Duarte, I. S. Novella, S. C. Weaver et al., "RNA virus quasispecies: significance for viral disease and epidemiology," *Infectious Agents and Disease*, vol. 3, no. 4, pp. 201–214, 1994.
- [6] T. Laskus, J. Wilkinson, J. F. Gallegos-Orozco et al., "Analysis of hepatitis C virus quasispecies transmission and evolution in patients infected through blood transfusion," *Gastroenterology*, vol. 127, no. 3, pp. 764–776, 2004.
- [7] L. Barzon, E. Lavezzo, V. Militello, S. Toppo, and G. Palù, "Applications of next-generation sequencing technologies to diagnostic virology," *International Journal of Molecular Sciences*, vol. 12, no. 11, pp. 7861–7884, 2011.
- [8] N. Beerenwinkel, H. F. Gunthard, V. Roth, and K. J. Metzner, "Challenges and opportunities in estimating viral genetic diversity from next-generation sequencing data," *Frontiers in Microbiology*, vol. 3, article 329, 2012.
- [9] N. Beerenwinkel, "Ultra-deep sequencing for the analysis of viral populations," *Current Opinion in Virology*, vol. 1, no. 5, pp. 413–418, 2011.
- [10] S. Guglietta, A. R. Garbuglia, V. Pacciani et al., "Positive selection of cytotoxic T lymphocytes escape variants during acute hepatitis C virus infection," *European Journal of Immunology*, vol. 35, no. 9, pp. 2627–2637, 2005.
- [11] C. Di Lorenzo, A. G. Angus, and A. H. Patel, "Hepatitis C virus evasion mechanisms from neutralizing antibodies," *Viruses*, vol. 3, no. 11, pp. 2280–2300, 2011.
- [12] A. Escobar-Gutierrez, M. Vazquez-Pichardo, M. Cruz-Rivera et al., "Identification of hepatitis C virus transmission using a next-generation sequencing approach," *Journal of Clinical Microbiology*, vol. 50, no. 4, pp. 1461–1463, 2012.
- [13] R. A. Bull, F. Luciani, K. McElroy et al., "Sequential bottlenecks drive viral evolution in early acute hepatitis c virus infection," *PLoS Pathogens*, vol. 7, no. 9, Article ID e1002243, 2011.
- [14] F. Bolcic, M. Sede, F. Moretti et al., "Analysis of the PKR-eIF2alpha phosphorylation homology domain (PePHD) of hepatitis C virus genotype 1 in HIV-coinfected patients by ultra-deep pyrosequencing and its relationship to responses to pegylated interferon-ribavirin treatment," *Archives of Virology*, vol. 157, no. 4, pp. 703–711, 2012.
- [15] S. Fonseca-Coronado, A. Escobar-Gutierrez, K. Ruiz-Tovar et al., "Specific detection of naturally occurring hepatitis C virus mutants with resistance to telaprevir and boceprevir (protease inhibitors) among treatment-naive infected individuals," *Journal of Clinical Microbiology*, vol. 50, no. 2, pp. 281–287, 2012.
- [16] O. Zagordi, A. Bhattacharya, N. Eriksson, and N. Beerenwinkel, "ShoRAH: estimating the genetic diversity of a mixed sample from next-generation sequencing data," *BMC Bioinformatics*, vol. 12, article 119, 2011.
- [17] I. Astrovskaia, B. Tork, S. Mangul et al., "Inferring viral quasispecies spectra from 454 pyrosequencing reads," *BMC Bioinformatics*, vol. 12, supplement 6, 2011.
- [18] C. Quince, A. Lanzen, R. J. Davenport, and P. J. Turnbaugh, "Removing noise from pyrosequenced amplicons," *BMC Bioinformatics*, vol. 12, article 38, 2011.
- [19] P. Skums, Z. Dimitrova, D. S. Campo et al., "Efficient error correction for next-generation sequencing of viral amplicons," *BMC Bioinformatics*, vol. 13, Supplement 10, p. S6, 2012.
- [20] M. Gerotto, F. Dal Pero, S. Loffreda et al., "A 385 insertion in the hypervariable region 1 of hepatitis C virus E2 envelope protein is found in some patients with mixed cryoglobulinemia type 2," *Blood*, vol. 98, no. 9, pp. 2657–2663, 2001.
- [21] T. Laskus, M. Radkowski, L. F. Wang, M. Nowicki, and J. Rakela, "Uneven distribution of hepatitis C virus quasispecies in tissues from subjects with end-stage liver disease: confounding effect of viral adsorption and mounting evidence for the presence of low-level extrahepatic replication," *Journal of Virology*, vol. 74, no. 2, pp. 1014–1017, 2000.
- [22] K. Tamura, D. Peterson, N. Peterson, G. Stecher, M. Nei, and S. Kumar, "MEGA5: Molecular evolutionary genetics analysis using maximum likelihood, evolutionary distance, and maximum parsimony methods," *Molecular Biology and Evolution*, vol. 28, no. 10, pp. 2731–2739, 2011.
- [23] K. Tamura and M. Nei, "Estimation of the number of nucleotide substitutions in the control region of mitochondrial DNA in humans and chimpanzees," *Molecular Biology and Evolution*, vol. 10, no. 3, pp. 512–526, 1993.
- [24] I. Malet, M. Belnard, H. Agut, and A. Cahour, "From RNA to quasispecies: a DNA polymerase with proofreading activity is

- highly recommended for accurate assessment of viral diversity,” *Journal of Virological Methods*, vol. 109, no. 2, pp. 161–170, 2003.
- [25] J. Cline, J. C. Braman, and H. H. Hogrefe, “PCR fidelity of Pfu DNA polymerase and other thermostable DNA polymerases,” *Nucleic Acids Research*, vol. 24, no. 18, pp. 3546–3551, 1996.
- [26] A. Gilles, E. Megléc, N. Pech, S. Ferreira, T. Malausa, and J. F. Martin, “Accuracy and quality assessment of 454 GS-FLX Titanium pyrosequencing,” *BMC Genomics*, vol. 12, article 245, 2011.
- [27] B. B. Simen, J. F. Simons, K. H. Hullsiek et al., “Low-abundance drug-resistant viral variants in chronically HIV-infected, antiretroviral treatment-naive patients significantly impact treatment outcomes,” *Journal of Infectious Diseases*, vol. 199, no. 5, pp. 693–701, 2009.
- [28] X. Wu, T. Zhou, J. Zhu et al., “Focused evolution of HIV-1 neutralizing antibodies revealed by structures and deep sequencing,” *Science*, vol. 333, no. 6049, pp. 1593–1602, 2011.
- [29] A. Gonzalez-Serna, R. A. McGovern, P. R. Harrigan et al., “Correlation of the virological response to short-term maraviroc monotherapy with standard and deep-sequencing-based genotypic tropism prediction methods,” *Antimicrobial Agents and Chemotherapy*, vol. 56, no. 3, pp. 1202–1207, 2012.
- [30] A. D. Redd, A. Collinson-Streng, C. Martens et al., “Identification of HIV superinfection in seroconcordant couples in Rakai, Uganda, by use of next-generation deep sequencing,” *Journal of Clinical Microbiology*, vol. 49, no. 8, pp. 2859–2867, 2011.
- [31] G. T. Spear, M. Sikaroodi, M. R. Zariffard, A. L. Landay, A. L. French, and P. M. Gillevet, “Comparison of the diversity of the vaginal microbiota in HIV-infected and HIV-uninfected women with or without bacterial vaginosis,” *Journal of Infectious Diseases*, vol. 198, no. 8, pp. 1131–1140, 2008.
- [32] I. Vandenbroucke, H. Van Marck, P. Verhasselt et al., “Minor variant detection in amplicons using 454 massive parallel pyrosequencing: experiences and considerations for successful applications,” *Biotechniques*, vol. 51, no. 3, pp. 167–177, 2011.
- [33] O. Zagordi, R. Klein, M. Däumer, and N. Beerenwinkel, “Error correction of next-generation sequencing data and reliable estimation of HIV quasispecies,” *Nucleic Acids Research*, vol. 38, no. 21, pp. 7400–7409, 2010.
- [34] N. Eriksson, L. Pachter, Y. Mitsuya et al., “Viral population estimation using pyrosequencing,” *PLoS Computational Biology*, vol. 4, no. 4, Article ID e1000074, 2008.
- [35] C. F. Wright, M. J. Morelli, G. Thébaud et al., “Beyond the consensus: dissecting within-host viral population diversity of foot-and-mouth disease virus by using next-generation genome sequencing,” *Journal of Virology*, vol. 85, no. 5, pp. 2266–2275, 2011.