

Accurate *in vivo* population sequencing uncovers drivers of within-host genetic diversity in viruses

Maoz Gelbart^{1,#}, Sheri Harari^{1,#}, Ya'ara Ben-Ari¹, Talia Kustin¹, Dana Wolf^{2,3},
Michal Mandelboim^{4,5}, Orna Mor^{4,6}, Pleuni Pennings⁷, Adi Stern^{1,*}

¹ School of Molecular Cell Biology and Biotechnology, George S. Wise Faculty of Life Sciences, Tel Aviv University, Tel Aviv, Israel

² Clinical Virology Unit, Hadassah Hebrew University Medical Center, Jerusalem, Israel

³ The Lautenberg Center for General and Tumor Immunology, IMRIC, the Faculty of Medicine, the Hebrew University, Jerusalem, Israel

⁴ Central Virology Laboratory, Ministry of Health, Sheba Medical Center, Ramat-Gan, Israel.

⁵ Department of Epidemiology and Preventive Medicine, School of Public Health, Sackler Faculty of Medicine, Tel Aviv University, Tel Aviv, Israel

⁶ Sackler School of Medicine, Tel Aviv University, Tel Aviv, Israel.

⁷ Department of Biology, San Francisco State University, San Francisco, CA, USA

Authors contributed equally

* To whom correspondence should be addressed. Email: sternadi@tau.ac.il

1 **ABSTRACT**

2 Mutations fuel evolution and facilitate adaptation to novel environments. However,
3 characterizing the spectrum of mutations in a population is obscured by high error rates of next
4 generation sequencing. Here, we present AccuNGS, a novel *in vivo* sequencing approach that
5 detects variants as rare as 1:10,000. Applying it to 46 clinical samples taken from early infections
6 of the human-infecting viruses HIV, RSV and CMV, revealed large differences in within-host
7 genetic diversity among virus populations. Haplotype reconstruction revealed that increased
8 diversity was mostly driven by multiple transmitted/founder viruses in HIV and CMV samples.
9 Conversely, we detected an abundance of defective virus genomes (DVGs) in RSV samples,
10 including hyper-edited genomes, nonsense mutations and single point deletions. Higher
11 proportions of DVGs correlated with increased viral loads, suggesting increased cellular co-
12 infection rates, which enable DVG persistence. AccuNGS establishes a general platform that
13 allows detecting DVGs, and in general, rare variants that drive evolution.

14 INTRODUCTION

15 Viruses are among the fastest evolving entities on earth. Thanks to short generation times, large
16 population sizes and high mutation rates, viruses and in particular RNA viruses rapidly
17 accumulate genetic diversity. This genetic diversity is key to successful adaptation of viruses to
18 novel challenges such as the immune system and drugs (Duffy, et al. 2008). The short time
19 window following virus transmission, termed *acute infection*, is extremely critical for virus
20 populations: they must develop their arsenal of genetic diversity to either escape from the
21 immune system, evade drugs or adapt to a new environment (possibly even a new host)
22 (Parrish, et al. 2008). While we know that many genetic variants are created in this critical time
23 window – the vast majority of these variants segregate at low frequencies that approach the
24 mutation rate of the virus, and are completely undetectable using current next generation
25 sequencing (NGS) approaches (Acevedo, et al. 2014). In fact, these variants will only be
26 discovered after they reach high frequencies, when they will already exert a deleterious effect
27 on the host. For example, today we are only able to detect an HIV drug resistance mutation
28 when it is at a relatively high frequency, after resistance phenotype is observed and already
29 unpreventable (Ram, et al. 2015; Boucher, et al. 2018; Döring, et al. 2018). The naturally
30 occurring frequency of drug resistance mutations could also allow us to infer the fitness cost of
31 these mutations in the absence of drugs (Theys, et al. 2018). An additional long-standing
32 question in the field of virus evolution is the measurement of mutation rates, which also
33 requires accurate sequencing (Acevedo, et al. 2014; Zanini, Puller, et al. 2017). Thus, an accurate
34 sequencing method that maintains high yield is required to thoroughly study evolving
35 populations.

36 Accurate population sequencing is critical in many different disciplines including genetics,
37 immunology and microbiology as well as tumor screening and prenatal diagnosis. However,
38 using the standard NGS protocols may result in significant background error rates. In fact,
39 following typical post-processing of NGS data, genetic variants that are measured at frequencies
40 lower than 1-5% are discarded (Meacham, et al. 2011; Casadella and Paredes 2017; Huber, et al.
41 2017; McCrone, et al. 2018). In the past few years, several innovative approaches were
42 suggested to reduce the background error rates of the NGS process: rolling-circle-based
43 redundant coding of the amplified fragments (Lou, et al. 2013; Acevedo, et al. 2014; Reid-Bayliss
44 and Loeb 2017; Wang, et al. 2017); consensus sequencing of barcoded genomic fragments
45 (Jabara, et al. 2011; Kennedy, et al. 2014; Zhou, et al. 2015; Jee, et al. 2016; Newman, et al.

46 2016; Wang, et al. 2017); error reduction by overlapping paired reads in paired-end sequencing
47 (Chen-Harris, et al. 2013; Schirmer, et al. 2015; Preston, et al. 2016); and usage of improved
48 polymerases (Imashimizu, et al. 2013). However, most experimental methods described above
49 are designed for samples with high biomass and are inapplicable for sequencing of clinical
50 samples, where the biomass may be extremely low. Furthermore, these experimental protocols
51 may introduce their own artifacts to the sequencing process (Lou, et al. 2013; Brodin, et al.
52 2015). On the computational side, it has been suggested that well-established variant callers do
53 not perform well on clinical virus samples (McCrone and Lauring 2016). Here, we sought to
54 develop a simple and rapid approach that can tackle the problem of accurate sequencing of
55 clinical samples, and applied it to study the early stages of virus infection.

56 We describe AccuNGS, a simple yet powerful approach for accurate population sequencing and
57 bioinformatics variant calling. We extensively optimize all stages of the method to ensure high
58 accuracy and maximal yield. We use AccuNGS to perform in-depth sequencing of 46 samples
59 from three different major human pathogenic viruses: human immunodeficiency virus (HIV),
60 respiratory syncytial virus (RSV), and cytomegalovirus (CMV), all sampled during the acute
61 infection stage. We compare the within-host genetic diversity among and within different virus
62 populations, and find patterns characteristic of each virus. We demonstrate the role of multiple
63 transmitted/founder viruses as major contributors to the genetic diversity in HIV and CMV.
64 Furthermore, we identify and quantify the impact of various host editing enzymes on the
65 mutational spectrum of viral genomes/populations in vivo. Intriguingly, we find that RSV
66 samples bear much higher levels of potentially defective virus genome (DVGs) than the two
67 other viruses analyzed herein. Finally, we propose a link between the level of DVGs in a
68 population and the levels of cellular co-infection.

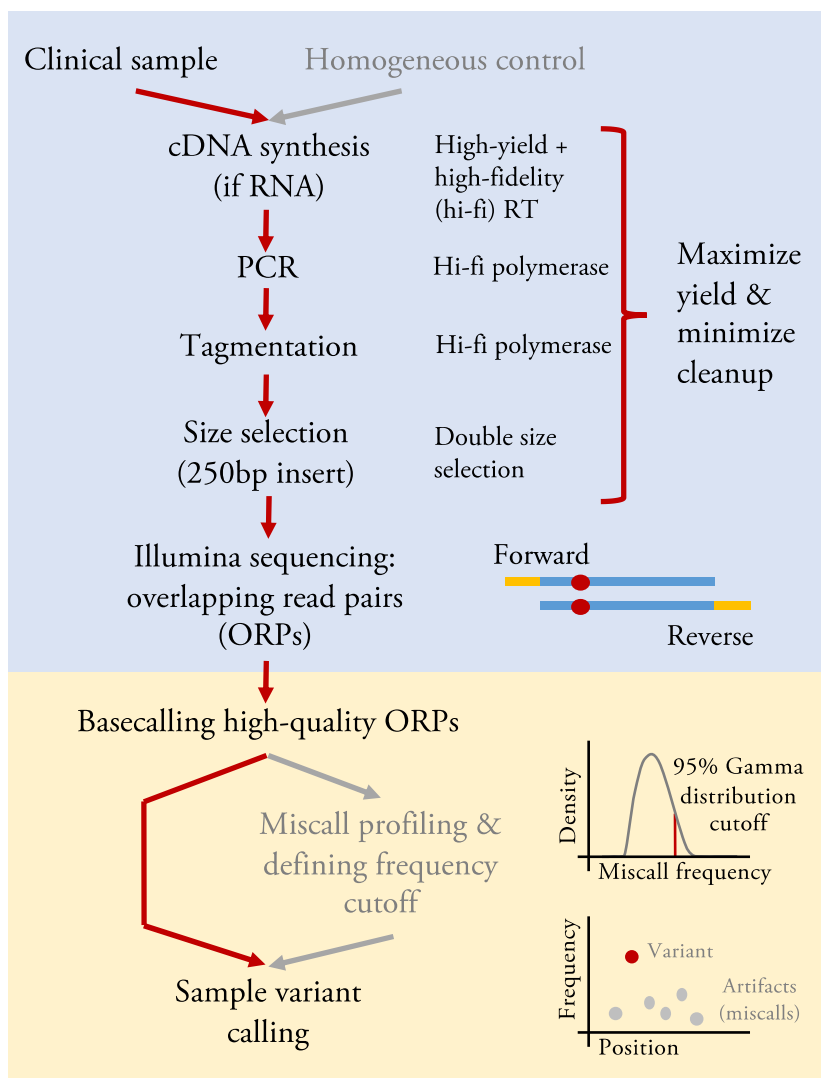
69 **RESULTS**

70 **An overview of the AccuNGS sequencing approach**

71 We have developed a new sequencing approach called AccuNGS, designed to maximize both the
72 accuracy of variant calling and template recovery. To achieve this, several key concepts were put
73 together: (i) usage of high-yield and high-fidelity polymerases to increase yield and reduce error
74 rates; (ii) sequencing error reduction through overlapping paired end reads; (iii) minimization of
75 template loss across different stages of the protocol, (iv) statistical modeling of error rates using
76 a control that forms the heart of a computational pipeline for variant calling and inference of

77 variant frequency (Fig. 1). We first extensively optimized the different experimental stages of
78 the protocol, by testing the accuracy of sequencing a clonally derived plasmid (Methods). This
79 allowed us to assess the contribution of the various stages to the error rate of the protocol
80 (Supplementary text, Fig. S1, Tables S1-S5). We were able to (a) mostly rule out the contribution
81 of PCR to error rates (Zanini, Brodin, et al. 2017), (b) note that the two reads generated by the
82 sequencer are dependent observations, and (c) reveal that AccuNGS is likely limited by the
83 accuracy of the sequencing machine itself (which we cannot affect) (Supplementary text, Fig.
84 S1B). We conclude that the mean error rate for transitions and transversions that is achieved by
85 AccuNGS is a little lower than 10^{-4} and 10^{-5} , respectively, and varies by types of mutation (Table
86 1). This allows us to statistically call variants that are as rare as 10^{-4} , which we further verified by
87 creating synthetic mixtures of two different HIV plasmids at different proportions and
88 sequencing them (Fig. S2).

89 Initially we used a barcoding approach (Jabara, et al. 2011) to estimate the number of viral
90 genomes we sequenced for one of the HIV clinical samples, allowing us to estimate that we had
91 sequenced ~15,000 separate genomes (Supplementary text). While it was reassuring that we
92 sequenced a large number of genomes, we also showed that the mere addition of a barcode led
93 to a reduction in the number of sequenced templates (Supplementary text, Fig. S3). In other
94 words, we found that without a barcode we sequence more genomes, but we cannot count how
95 many. We henceforth avoided the use of barcoding. Finally, we would like to emphasize that
96 while the *average* error rate of AccuNGS is very low, sampling and PCR biases may lead to a
97 given variant being detected at a lower or higher frequency than it should be (Illingworth, et al.
98 2017; Zhao and Illingworth 2019). This is especially true when the number of input templates is
99 low. We thus conclude that AccuNGS is useful for inferring average rates and average
100 frequencies (e.g., premature stop codons as elaborated below) rather than inferring the exact
101 frequency of one given variant in a given sample.



102

103 **Fig. 1. AccuNGS principles.** The protocol requires side by side targeted sequencing of a
 104 biological sample together with a homogenous control (e.g., a plasmid encoding for the same
 105 genomic region sequenced in the sample). Stages of the protocol include: (i) High-fidelity and
 106 high-yield RT reaction, (ii) High-fidelity PCR reaction, (iii) Tagmentation and library construction
 107 with size selection for an insert the size of a single paired-end read, (iv) Paired-end sequencing
 108 where each base in the insert is sequenced twice, once in the forward read and again in the
 109 reverse read, (v) Alignment of reads, Q-score filtering on both reads and base-calling of both the
 110 sample and a homogeneous control, and (vi) Sample variant calling based on gamma
 111 distributions that are fitted to the process errors, defined as mutations found in the
 112 homogeneous control.

113

114 **Table 1. Error rates for AccuNGS.** Shown are parameters for fitted gamma distributions given a
 115 Q30 score cutoff. Transitions are shaded in darker gray than transversions.

Error type	Shape (κ)	Scale (θ)	95 th percentile of distribution	Mean error rate ^a	Median error rate ^a
A->G	7.292	8.51E-06	1.04E-04	6.21E-05	6.0E-05
C->T	3.349	1.27E-05	8.70E-05	4.27E-05	3.8E-05
G->A	4.026	1.22E-05	9.51E-05	4.91E-05	4.6E-05
T->C	6.704	9.80E-06	1.12E-04	6.57E-05	6.3E-05
A>C	3.295	2.27E-06	1.52E-05	7.48E-06	6.0E-06
A>T	2.709	1.63E-05	9.54E-05	4.41E-05	4.0E-05
C>A	2.149	8.70E-06	4.33E-05	1.86E-05	1.4E-05
C>G	3.747	2.74E-06	2.02E-05	1.02E-05	9.0E-06
G>C	3.810	2.95E-06	2.21E-05	1.12E-05	1.0E-05
G>T	3.095	5.70E-06	3.67E-05	1.76E-05	1.6E-05
T>A	2.874	1.21E-05	7.41E-05	3.48E-05	3.3E-05
T>G	2.342	3.51E-06	1.85E-05	8.22E-06	6.0E-06

116 ^a Calculated based on the raw frequencies of mutations found in the baseline control.

117 **Accurate sequencing of different virus populations during acute/early infections**

118 We obtained a total of 46 samples from patients recently infected by the RNA viruses HIV and
119 RSV, and the DNA virus CMV (Table 2, Table S6). We focused our sequencing efforts mostly on
120 conserved genes, precisely since we expect less diversity and we wanted to test our ability to
121 detect rare “hidden” variants that are otherwise unobservable. Hyper-variable genes such as the
122 HIV-1 envelope have been sequenced extensively using other sequencing approaches (e.g.,
123 Keele, et al. 2008; Salazar-Gonzalez, et al. 2008; Zhou, et al. 2016), and the presence of high
124 frequency variants is not surprising in such genes. We thus chose the Gag-Pol open reading
125 frame for HIV, the M2 and L open reading frame (encoding for the viral polymerase) for RSV,
126 and the UL54 (also encoding for the viral polymerase) for CMV. In order to allow for comparison,
127 we also sequenced the F and G envelope glycoproteins genes in RSV, and further compared our
128 results to previous sequencing results of the envelope gene in HIV.

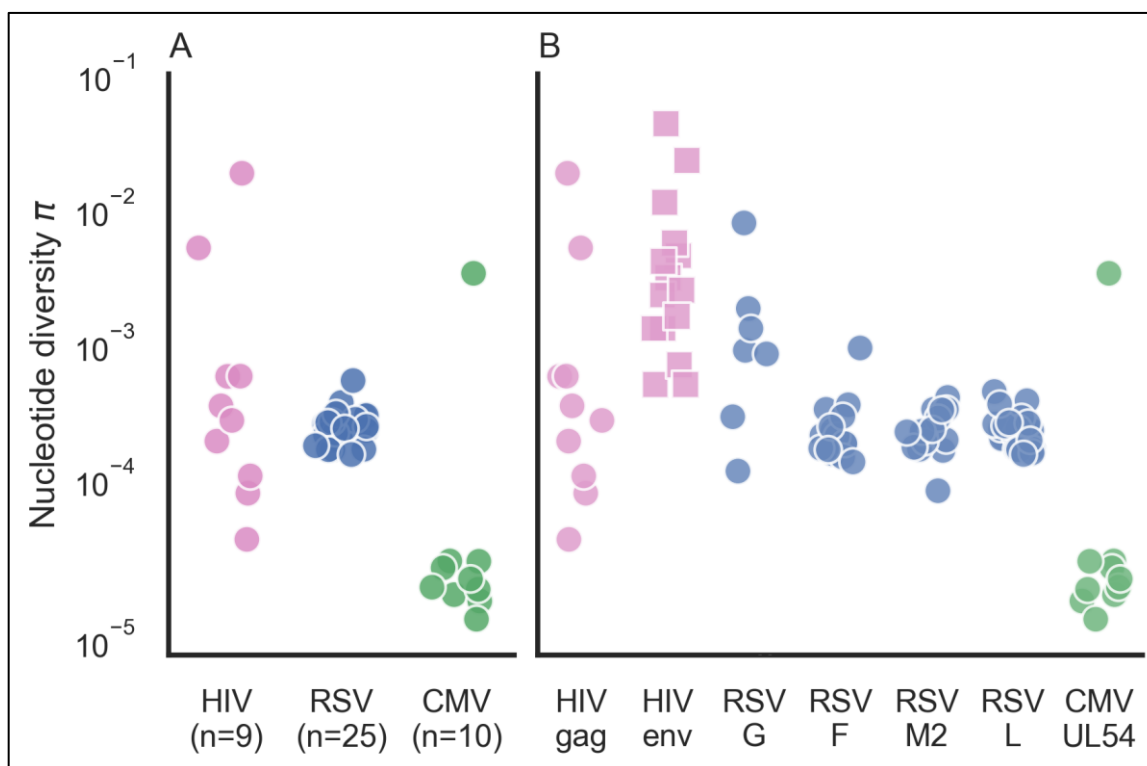
129 Each sample underwent population sequencing and variant calling using AccuNGS (see
130 Methods). We began by calculating the total nucleotide diversity π in each sample based on the
131 transition variants (Methods). This revealed different distributions of diversity within and
132 between viruses (Fig. 2A). In the HIV populations, diversity values spanned several orders of

133 magnitude. On the contrary, RSV populations exhibited very similar intermediate levels of
 134 diversity. Similarly, CMV populations usually displayed the lowest diversity with the exception of
 135 one sample. We set out to understand what factors drive the differences in diversity among the
 136 different samples.

137 **Table 2. Details of samples sequenced from clinical virus samples.**

Virus	#Samples	Samples' origin	Estimated weeks post infection	Region sequenced (NCBI ID)
HIV	9	Plasma	2-5	532:3280 (K03455)
RSV	25	Nasal/throat swabs	<1	4640-14350 (U39661)
CMV	12	Amniotic fluid / Urine / saliva	>=4	78200-81912 (NC_006273)

138



139

140 **Fig. 2. Nucleotide diversity π for acute infections across different virus samples and genes.** (A)
 141 Each point represents the π diversity of a single sample, across all genes sequenced. Diversity
 142 values were calculated using transition mutations only. (B) Gene by gene breakdown of
 143 nucleotide π diversity. Values for HIV envelope (squares) were taken from previously published
 144 data (Salazar-Gonzalez, et al. 2008).

145 **Mutation and selection**

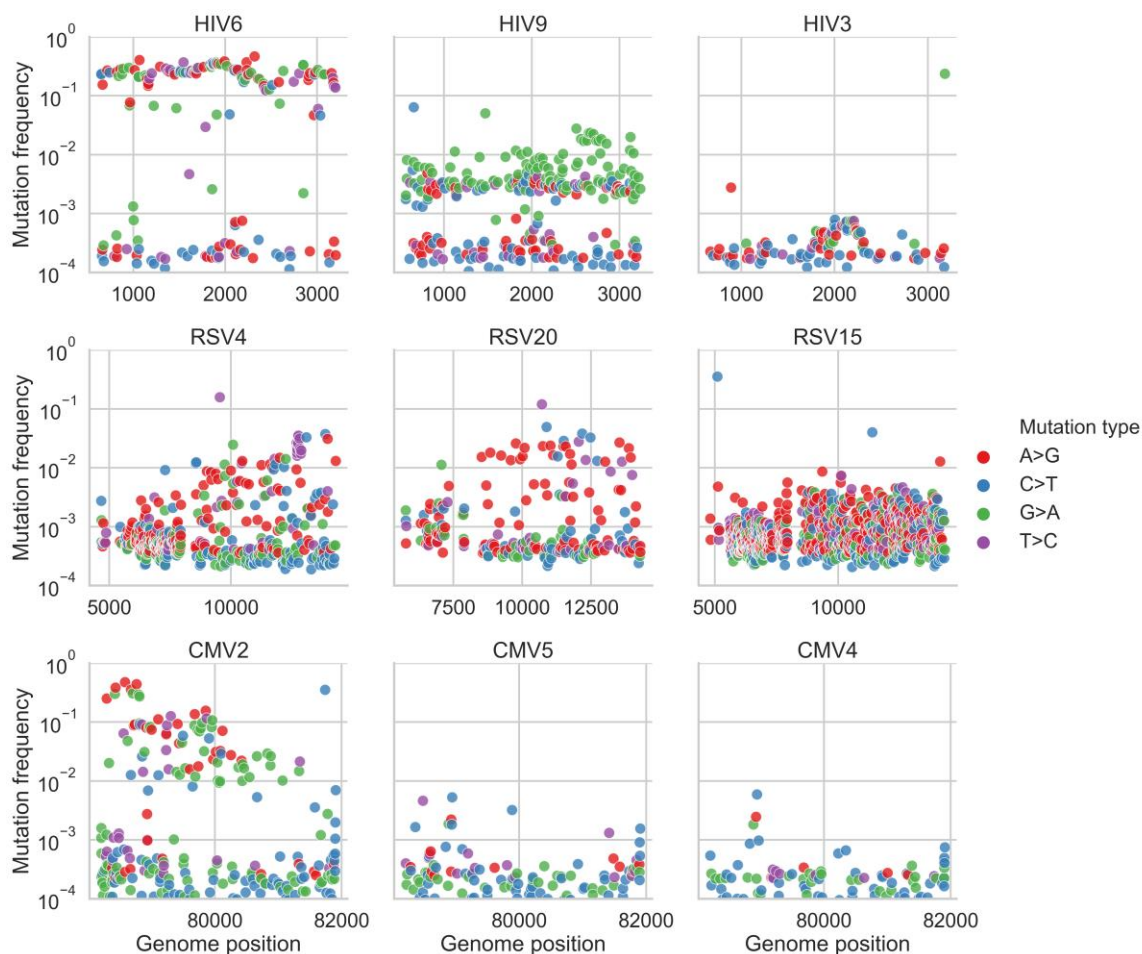
146 We first considered the two most evident evolutionary causes of differences in diversity:
147 mutation and selection. First, when considering the mutation rate of a virus, it is clear that the
148 only DNA virus in our data has much lower mutation rates than RNA viruses (Sanjuan, et al.
149 2010) and indeed displays lower diversity. The two RNA viruses display more diversity than the
150 DNA virus CMV, but the variation in diversity levels is much higher in HIV than in RSV.

151 We further considered whether differences in selection pressure cause the variation in diversity
152 we see among the RNA virus samples. This was unlikely to cause within-virus differences, since
153 we sequenced the same set of genes within each of the virus samples. We did note that the
154 immunogenic envelope proteins in this study (HIV Env and RSV G proteins), known to be under
155 positive selection (Seibert, et al. 1995; Nielsen and Yang 1998; Tan, et al. 2013), displayed on
156 average higher diversity than the conserved genes (Fig. 2B). This is not surprising given that
157 some mutations in envelope proteins will be under positive selection since they allow immune
158 evasion, and hence may reach higher frequencies. However, this could not explain why we saw
159 dramatic differences in diversity in different samples from the same virus when focusing on the
160 same gene (e.g., *gag* in HIV).

161 **Transmission bottleneck size as a contributor to genetic diversity during acute infections**

162 It has previously been noted that infections initiated by a few different divergent viruses are
163 characterized by higher genetic diversity (Keele, et al. 2008; Cudini, et al. 2019). Visual
164 inspection of our frequency plots (Fig. 3, Fig. S4-S6) suggested that often variant frequencies
165 were strongly imbalanced, also evident as “bands” of variants at similar frequencies. For
166 example, sample HIV6 (measured diversity 1.46×10^{-2}) contained many variants segregating at a
167 frequency of $\sim 2 \times 10^{-1}$ yet very few variants segregated at frequencies between 10^{-3} and 10^{-1} (Fig.
168 3A). We first considered how likely it is that such a sample would be initiated by only one
169 founder virus/genotype, where all variants begin at a defined frequency of zero. Given a large
170 enough population size and a mutation rate in the order of 10^{-5} mutations/site/day (Zanini,
171 Puller, et al. 2017), we expect neutral variants that are likely generated over and over almost
172 every day to roughly reach a frequency of 10^{-4} - 10^{-3} after a few weeks of infection, which is much
173 lower than 10^{-1} . Genetic drift or positive selection could drive a few variants to increase in
174 frequency over a short time; however, it seems extremely unlikely that there is such a large set
175 of sites under the exact same regime of positive selection, especially as we had sequenced a

176 gene where positive selection is all in all less prevalent, at least this early in the course of the
177 infection. Thus, it seems quite unlikely that very high diversity samples containing many high
178 frequency variants are founded by one virus genotype, and a more likely explanation is the
179 presence of multiple transmitted/founder viruses.



180

181 **Figure 3. Variants frequency plots in representative samples.** Shown are frequencies of
182 transition variants called by AccuNGS, for representative samples from each virus (HIV, top row,
183 RSV, middle row, CMV, bottom row). Samples exemplify multiple founder infections, mutation
184 biases, and relatively homogenous populations (see text for details).

185

186 **Inferring haplotypes and multiple founders**

187 To evaluate the number of founder viruses we require an estimation of the different haplotypes
188 present in a sample, and their abundances. However, reconstruction of virus haplotypes from
189 short reads and from one time-point is a longstanding problem (Schirmer, et al. 2014). This is

190 due to two conflicting features of viral population sequencing data: on the one hand, the data is
191 often too homogenous. In other words, most reads are identical or almost identical to the
192 consensus, and there may not be enough variants on one read that allow “linking” it with
193 another read. On the other hand, the mutation rates of viruses may scale with sequencing error
194 rates, throwing off most commonly used haplotype reconstruction methods.

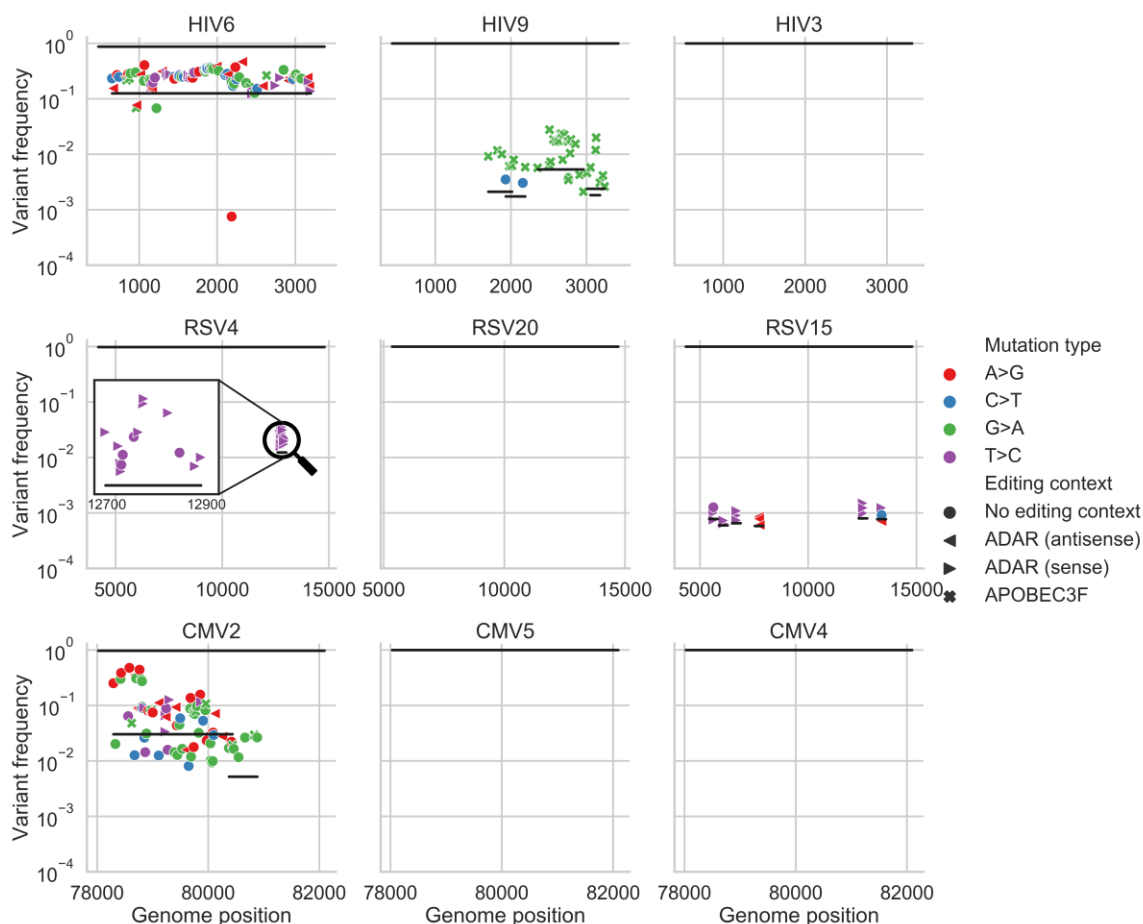
195 Two major advantages of our data are the much higher accuracy of AccuNGS, coupled with very
196 high sequencing depth. We thus set out to develop a new approach for inferring viral
197 haplotypes. Instead of attempting to reconstruct the entire haplotype, we mainly focused on
198 inferring if more than one haplotype (and hence more than one founder virus in our case) is
199 present in a sample. Our approach is based on looking for statistical enrichment for two variants
200 being present on the same read as opposed to each variant on its own, and then linking these
201 reads one with another based on shared variants (Methods) (Yang, et al. 2013). To validate our
202 approach, we used the samples where we had synthetically mixed two plasmids and discovered
203 that all and only true variant sites were identified as linked to each other in with an accurately
204 estimated frequency of 1:10,000 (Fig. S7). Importantly, this confirmed that AccuNGS does not
205 suffer from various artifacts such as PCR recombination that could break linkage between
206 adjacent sites, and that even a rare haplotype can be inferred.

207 Our haplotype reconstruction approach also led us to realize one of the combined strength and
208 pitfalls of accurate sequencing: we were able to initially detect minute contaminations (a few
209 hundred out of millions of reads) from one sample into another, which we were then able to
210 computationally filter out (Supplementary text). This contamination likely occurred during one
211 of the stages of the library preparation and sequencing, and emphasizes the sensitivity of
212 genomics approaches today, which may often exceed that of the molecular biology itself
213 (Kircher, et al. 2011; Gu, et al. 2019). On the other hand, AccuNGS also allows for the clear-cut
214 detection and evaluation of any contamination, which we believe are very important to capture.

215 We next applied our haplotype inference flow to all the filtered samples, and found that all of
216 the high diversity samples (total diversity $>10^{-3}$, Fig. 2A) exhibited strong evidence for containing
217 two or more divergent haplotypes (Figs. S4-S6). Two examples are shown in Figs. 3 and 4: HIV
218 sample 6 has a “band” of variant frequencies around 2×10^{-1} (Fig. 3), and indeed most of these
219 variants can be linked to each other in this sample (Fig. 4). CMV sample 2 has a wide “band” of
220 variant frequencies between 10^{-2} and 5×10^{-1} (Fig. 3), which were also mostly found to be linked,

221 and likely represent a founder haplotype and the associated variants that were created on the
222 background of this haplotype (Fig. 4). In general, we found no evidence for two or more
223 haplotypes in the less diverse samples, except for the most diverse RSV sample that also showed
224 limited evidence of a low frequency haplotype (Fig. S5) (see discussion).

225



226

227 **Figure 4. Haplotype reconstruction based on co-occurrence of variants on the same reads.**

228 Shown are inferred haplotypes (lines) based on consecutive significant associations of pairs of
229 variants (shapes) one to another on the same read. The uppermost line in each panel represents
230 the consensus sequence, which by definition is the major haplotype in each sample. Both HIV6
231 and CMV2 samples show strong evidence of an additional haplotype, which is likely a second
232 founder genotype. Sample HIV9 shows evidence of G>A hyper-mutation in the context of
233 APOBEC3 editing, samples RSV4 and RSV15 show evidence of T>C or A>G hyper-mutation in the
234 context of ADAR editing in regions spanning a few hundred bases. The hyper-mutated region in
235 RSV4 sample is magnified for clarity. "Empty" panels signify what are likely single haplotype
236 infections, with no evidence of hyper-mutation.

237 **Short hyper-mutated genomic stretches**

238 One well-known phenomenon of HIV infections is the potential of host APOBEC3 (A3) proteins
239 to induce hyper-editing on the negative strand of nascent HIV DNA during reverse transcription,
240 resulting in an excess of G>A mutations in regions of the RNA genome (Hache, et al. 2006;
241 Malim 2009). This hyper-mutation strategy is thought to lead to DVGs that are unable to
242 replicate. However, HIV encodes for a gene called *vif* that counteracts A3 proteins, and thus
243 most HIV viruses sequenced from blood samples show only minor evidence for A3 activity
244 (Cuevas, et al. 2015). Similarly, the family of human ADAR proteins have also been shown to
245 induce A>I mutations (read as A>G mutations) in a variety of viruses (Samuel 2012). We set out
246 to test if we detect signals of hyper-editing in our samples. In particular we sought to find
247 stretches of hyper-mutations using our haplotype reconstruction approach in order to evaluate
248 whether hyper-editing contributes to the observed genetic diversity, and to what extent.

249 Of all 46 samples, only one HIV sample (HIV9, Fig. 3) displayed strong evidence for G>A hyper-
250 editing. In this sample, editing seemed to be widespread, with multiple distinct and overlapping
251 hyper-mutated genomes (Fig. 4). Hyper G>A mutations were enriched in the context of GpA
252 which is the APOBEC3D/F/H favored editing context but not of the canonical APOBEC3G (Beale,
253 et al. 2004; Bishop, et al. 2004). Most variants on these hyper-mutated stretches were missense
254 variants; some of these stretches contained variants that lead to premature stop codons which
255 are presumably lethal for the virus (Fig. S4). The maximum frequency of such variants in the
256 sample was roughly 2×10^{-2} . To test whether this occurs due to an inactive *vif* gene, we
257 sequenced this gene in this sample using AccuNGS. We found no support for this hypothesis
258 since the consensus sequence of this gene was intact, but we once again noticed a relatively
259 high level of G>A mutations in the *vif* gene itself (Fig. S8). Notably, this sample was taken from
260 the patient with the highest viral load (VL) in our study, which was >10 million cp/ml.

261 Out of 25 independent RSV populations, 11 (44%) exhibited evidence of ADAR-mediated hyper-
262 edited genomes, manifested as at least three ADAR-associated mutations on the same
263 haplotype (Whitmer, et al. 2018). The presence of ADAR-like hyper-edited haplotypes was
264 positively correlated with viral loads ($R^2=0.59$; logistic regression; $p<0.0001$), for example all 9
265 samples with measured VL>650k cp/ml contained hyper-mutated inferred haplotypes (Fig. 5D).
266 When observed, ADAR-like linked variants were present at frequencies varying between $\sim 10^{-3}$
267 and $\sim 10^{-2}$, which by far exceed the mutation rate of any known virus. Out of 26 ADAR-like hyper-
268 edited haplotypes, 23 of them were on the negative strand and only 3 were on the positive

269 strand, in line with previous studies demonstrating that most ADAR-like mutations are acquired
270 on the negative strand of (-)ssRNA viruses (e.g., Whitmer, et al. 2018). Most of the ADAR-like
271 variants on these stretches were missense variants, suggesting they have a detrimental effect on
272 the virus (Fig. S5).

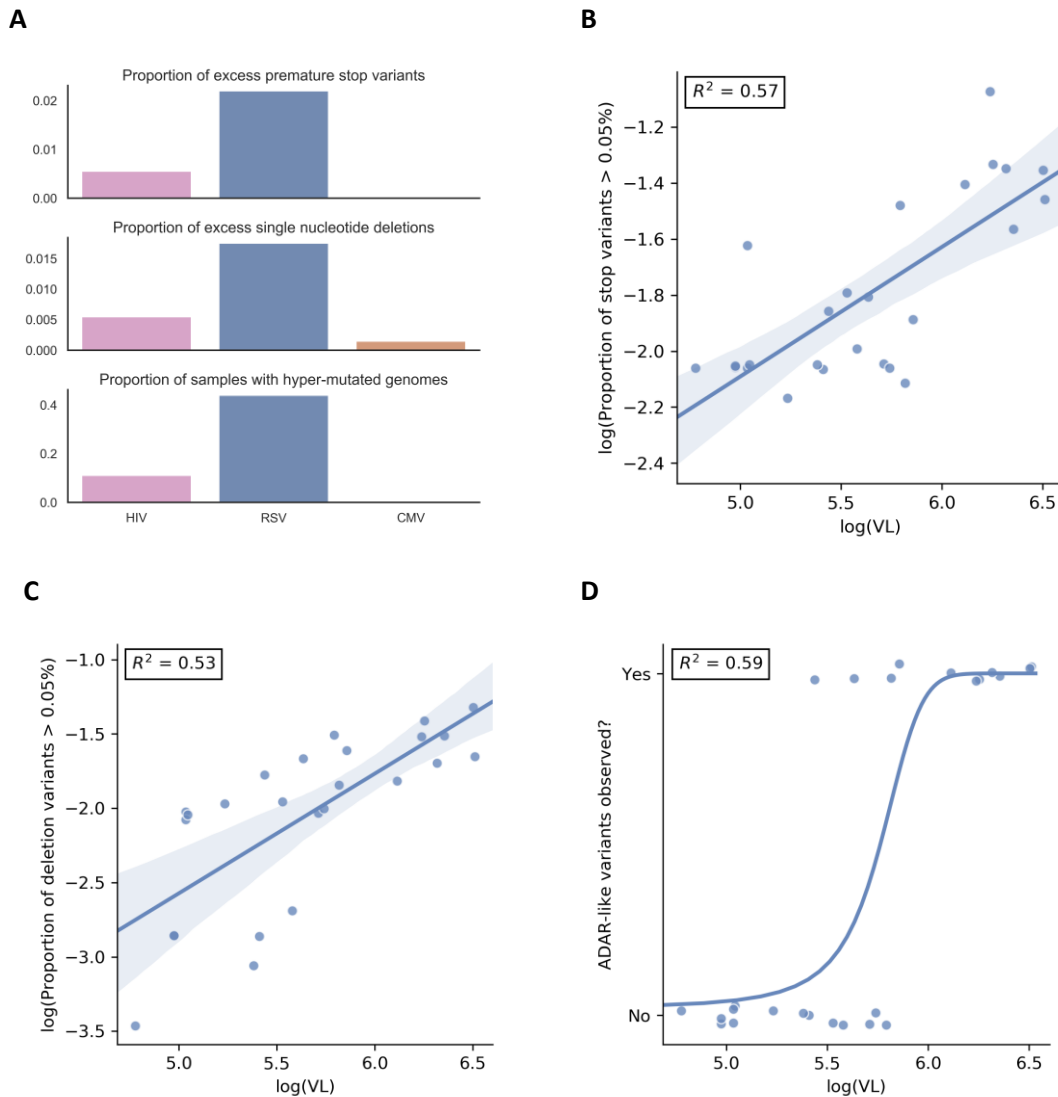
273 None of the CMV populations exhibited any A3, ADAR, or other pattern of hyper-mutation,
274 suggesting that these hyper-mutating enzymes do not act on CMV populations, at least not for
275 the gene sequenced here, or at the level of detection of AccuNGS (but see Weisblum, et al.
276 2017).

277 **Defective virus genomes**

278 Under mutation-selection balance, lethal variants, which are the ultimate DVGs, are expected to
279 segregate at the frequency at which they are generated, which is the mutation rate (Acevedo, et
280 al. 2014; Cuevas, et al. 2015). Based on the expected mutation rates in the viruses we
281 sequenced ($\leq 10^{-4}$ for HIV and RSV, and $\sim 10^{-7}$ for CMV (Sanjuan, et al. 2010)), our initial
282 expectation was that lethal variants should rarely be observed. We were therefore surprised to
283 note the existence of premature stop codons, one of the most obvious forms of a lethal
284 mutation, at a frequency of 10^{-3} - 10^{-2} in some of the RSV samples (Fig. S5, Fig. S8). All in all we
285 found that RSV substantially differed from the two other viruses with respect to three markers
286 of excess lethal variants: high frequency premature stop codons, high frequency single
287 nucleotide deletions, and presence of hyper-edited haplotypes. RSV had higher proportions of
288 all three types of lethal mutations (Fig. 5A). It thus appears that mutation-selection balance does
289 not necessarily hold for RSV (at least, for some RSV infections). Similar to the presence of hyper-
290 edited haplotypes for RSV, an excess of stop variants and deletion variants in RSV positively
291 correlated with viral load (VL) (Fig. 5B-D). This raises the possibility that extensive
292 complementation occurs in some RSV samples and in the HIV sample described above, allowing
293 the “rescue” of DVGs through high rates of co-infection.

294

295



296 **Figure 5. RSV contains high levels of DVGs.** (A) Proportion of sites across all samples identified
297 as containing excess lethal variants (top and middle panels), proportion of samples containing at
298 least one hyper-mutated haplotype with a frequency above 5×10^{-4} (bottom panel). For RSV, viral
299 load (VL) positively correlated with the (B) excess premature stop codons, (C) excess single
300 nucleotide deletions, and (D) presence of ADAR-like hyper-edited stretches. An excessive lethal
301 variant is defined as one exceeding a frequency of 5×10^{-4} . (B) and (C) display linear regression
302 lines, whereas (D) displays a logistic regression line.

303

304

305 **DISCUSSION**

306 Application of next generation sequencing to clinical samples is still limited by the ability to
307 reliably capture minor variants (Clutter, et al. 2016; McCrone and Luring 2016; Boucher, et al.
308 2018). Here we describe AccuNGS, a simple, rapid and accurate experimental protocol and
309 associated computational pipeline for detecting ultra-rare variants from low-biomass clinical
310 RNA and DNA samples. AccuNGS aims to accurately detect minor variants present in a
311 population of genomes at frequencies of 1:10,000 or lower, close to the mutation rate of RNA
312 viruses (Sanjuan, et al. 2010).

313 We used AccuNGS to characterize HIV, RSV, and CMV diversity. The HIV-1 samples were from
314 early stages of infection, typically 2-5 weeks post infection, based on serology testing. The RSV
315 samples were taken from children hospitalized due to respiratory problems, about 3-5 days post
316 infection (Lessler, et al. 2009), while the CMV samples taken from amniotic fluid or newborn
317 urine and saliva are several weeks post infection. Most previous studies have reported a nearly
318 completely homogenous virus population during early stages of infection (Henn, et al. 2012;
319 Zanini, et al. 2015; Kijak, et al. 2017). Nevertheless, in all samples sequenced, AccuNGS captured
320 dozens to thousands of minor transition and transversion variants, mostly segregating at
321 frequencies between 1:100 and 1:10,000 (Figs. S4-S6). These variants are most likely the genetic
322 reservoir that serves to fuel the future evolution of these virus populations.

323 Our results suggest that a prominent factor in determining the intra-host genetic diversity of a
324 sample during acute infection is the number of founder viruses, and in fact we find that the
325 samples with the highest levels of diversity always show evidence for multiple founder infection.
326 A debate has arisen regarding the diversity in CMV samples, where one study has claimed that
327 diversity in this DNA virus is comparable to diversity in RNA viruses (Renzette, et al. 2011), and
328 others suggest that diversity is low in single founder infections and is elevated only when
329 multiple founders/genotypes initiate the infection (Cudini, et al. 2019). Our results strongly
330 support the latter hypothesis, and in fact the high resolution of AccuNGS suggests even lower
331 diversity for single founder new CMV infections than what has been previously reported (Cudini,
332 et al. 2019).

333 For RSV samples, we found no evidence for multiple haplotypes in the samples, which is
334 somewhat surprising given that this is an airborne virus that reaches high titers. Previous work
335 has suggested that the number of founders in RSV infections is 25 ± 35 (Lau, et al. 2017). Notably,

336 these values were obtained for adults experimentally inoculated with RSV, whereas our study
337 represent natural infection of infants. However, another explanation for this discrepancy is that
338 our data does not allow us to detect infection with multiple founders when they share very
339 similar genotypes, since our haplotype reconstruction method relies on detecting short reads
340 that share two or more mutations. Given the very short duration of RSV infection, it is possible
341 that relatively little genetic diversity is created do novo, and hence very little genetic diversity is
342 transmitted. In other words, an infection may be initiated by several very genetically similar
343 founder genotypes, but we would not detect it. On the other hand, CMV and HIV create longer
344 infections, and the potential to generate and transmit more diverse genotypes within a single
345 carrier is higher.

346 Our results enabled pinpointing the activity of viral hyper-editing by host enzymes, namely
347 APOBEC3 enzymes and ADAR. The latter was particularly prominent in the RSV infections, where
348 we found distinct clusters of mutations matching ADAR context. Surprisingly, the frequency of
349 these clustered mutations was often relatively quite high, as discussed above. There is a debate
350 today surrounding the role of ADAR in viral infections: in some case it was found to be pro-viral
351 whereas in other cases it has been shown to be anti-viral. Pro-viral activity may be plausible
352 when considering that it has been found that ADAR protects cellular transcripts from being
353 detected by intracellular innate immune response (Liddicoat, et al. 2015; Pfaller, et al. 2018).

354 Is it possible that the ADAR signatures we find represent edited viral genomes that escape
355 innate immunity? If so, this would mean these are not DVGs but rather haplotypes with a
356 selective advantage. We consider that this is unlikely: many of the ADAR-like mutations we find
357 are non-synonymous, with often 5-10 such mutations found in a short region. It is highly
358 improbable that so many mutations would yield a “viable” genome, and we hence conclude that
359 ADAR-like hyper-editing yields DVGs. This comes together with evidence for other DVGs
360 (premature stop mutants and point deletions) that also segregate at similar frequencies in the
361 high viral load RSV samples. We find that the most likely explanation for this phenomenon is
362 that the rate of cellular co-infection is very high in some RSV samples, which may be promoted
363 by the syncytia that RSV creates, allowing for complementation of these DVGs. In fact, given
364 rough estimates of RSV mutation rates around 10^{-4} - 10^{-5} , we are able to infer that the rate of co-
365 infection should be between 0.9 and 0.99 in some of our RSV samples (Wilke and Novella 2003),
366 which is clearly very high. Moreover, we observe that samples with more DVGs are samples with

367 higher viral loads. Putting this together, we suggest that RSV infections probably occur in a
368 relatively dense site, which allows for so many co-infections, and for the propagation of DVGs.

369 The relationship between viral loads and the levels of DVGs in a sample is quite intuitive. This
370 suggests that in many cases, DVGs may be incorrectly neglected in downstream analyses. For
371 example, the assumption that such mutations occur at mutation selection balance is used to
372 estimate mutation rates and/or fitness costs (Acevedo, et al. 2014; Cuevas, et al. 2015; Theys, et
373 al. 2018), yet high levels of co-infection may lead to deviations from this balance. We suggest
374 that the use of AccuNGS can detect the rate at which DVGs occur, allowing the inference of the
375 co-infection rate. More generally, AccuNGS has the power to uncover the presence and rate of
376 rare genomic events, such as mutations and editing, directly from clinical samples, allowing a
377 better and more detailed understanding of the processes that govern evolution.

378

379

380 **MATERIALS AND METHODS**

381 **Ethics declaration**

382 The study was approved by the local institutional review boards of Tel-Aviv University, Sheba
383 Medical Center (approval number SMC 4631-17 for HIV and SMC 5653-18 for RSV), and
384 Haddasah Medical Center (approval number HMO-063911 for CMV).

385 **Reagents and Kits**

386 Unless stated otherwise, all the described reactions in this paper were carrier with the described
387 products according to the manufacturer's instructions: gel purifications were performed using
388 Wizard® SV Gel and PCR Clean-Up System (Promega, Madison, WI, USA); beads purifications
389 were performed using AMPure XP beads (Beckman Coulter, Brea, CA, USA); concentrations were
390 determined using Qubit fluorometer (Thermo Fisher Scientific, Waltham, MA, USA); reverse
391 transcription (RT) reactions were performed using SuperScript III or IV Reverse Transcriptase
392 (Thermo Fisher Scientific); polymerase chain reactions (PCR) were made using Platinum™
393 SuperFi™ high-fidelity DNA Polymerase (Thermo Fisher Scientific) or Q5 high-fidelity DNA
394 Polymerase (New England Biolabs (NEB), Ipswich, MA, USA).

395 **Generation of amplicons from HIV-1 clinical samples**

396 **Clinical HIV-1 samples.** Plasma samples from nine recently diagnosed HIV-1 patients with viral
397 loads of 5×10^5 - 1×10^7 cp/ml were provided by the National HIV Reference Laboratory, Chaim
398 Sheba Medical Center, Ramat-Gan, Israel (Table S6). HIV-1 viral loads were determined and RNA
399 extracted from 0.5 mL as described above. From each sample a maximum of 300,000 HIV-1
400 copies were reverse transcribed using random hexamer priming.

401 **HIV-1 control sample.** A pLAI.2 plasmid was linearized using the Sall restriction enzyme (NEB)
402 according to manufacturer's instruction. In order to create RNA a PCR reaction was set up using
403 primer containing the sequence of the T7 RNA polymerase promotor 5'TAA TAC GAC TCA CTA
404 TAG CTG GGA GCT CTC TGG CTA AC and the primer pLAI 5761-5782 5'GAG ACT CCC TGA CCC
405 AGA TGC C using Q5 DNA polymerase and the following PCR program: initial denaturation for
406 3min at 98°C, followed by 40 cycles of denaturation for 10sec at 98°C, annealing for 30sec at
407 65°C and extension for 3min at 72°C, and final extension for 5min at 72°C. Eight microliters from
408 PCR reaction were carried to in-vitro transcription reaction using HiScribe™ T7 High Yield RNA
409 Synthesis Kit (NEB) according to the manufacturer's instruction. Finally, the RNA was beads
410 purified (1X ratio). RT was accomplished using random hexamers priming.

411 **Generation of Gag-Pol amplicons.** The cDNA of the 9 HIV-1 clinical samples and the HIV-1
412 control sample were used to generate amplicons. To remove excess primers, the resulting cDNA
413 was beads purified (0.5X ratio) and eluted with 30µl nuclease-free water. Fifteen microliters of
414 each sample were then used for PCR amplification using SuperFi DNA polymerase. To amplify
415 ~2500 bp spanning entire Gag and part of Pol HIV-1 regions (HXB2 coordinates 524-3249, the
416 following primers were used: GAG FW 5'CTC AAT AAA GCT TGC CTT GAG TGC and RT gene RV
417 5'ACT GTC CAT TTA TCA GGA TGG AG, and the following PCR program: initial denaturation for
418 3min at 98°C, followed by 40 cycles of denaturation for 20sec at 98°C, annealing for 30sec at
419 62°C and extension for 2.5min at 72°C, and final extension for 5min at 72°C. The amplicons were
420 gel purified and their concentration was determined. The purified products were further used
421 for library construction.

422 **Generation of Gag amplicon with primer-ID from HIV9 sample.** A primer specific to the entire
423 Gag gene of HIV-1 (HXB2 position 2347) was designed with a 15 N-bases unique barcode
424 followed by a linker sequence for subsequent PCR, Gag ID RT 5'TAC CCA TAC GAT GTT CCA GAT
425 TAC GNN NNN NNN NNN NNN NAC TGT ATC ATC TGC TCC TG TRT CT. Based on the measured
426 viral load and sample concentration, 4 µl (containing roughly 300,000 HIV-1 copies) were taken
427 for reverse transcription reaction. Reverse transcription was performed using SuperScript IV RT
428 with the following adjustments: (1) In order to maximize the primer annealing to the viral RNA,
429 the sample was allowed to cool down gradually from 65°C to room temperature for 10 minutes
430 before it was transferred to ice for 2min; And (2) The reaction was incubated for 30min at 55°C
431 to increase the overall reaction yield. To remove excess primers, the resulting cDNA was beads
432 purified (0.5X ratio) and eluted with 35µl nuclease-free water. To avoid loss of barcoded primers
433 ("primer-ID"s) due to coverage drop at the ends of a read as a result of the NexteraXT
434 tagmentation process (see "Library construction for Miseq"), the PCR forward primer was
435 designed with a 60bp overhang so the barcode is far from the end of the read. The primers used
436 for amplification were Gag ID FW 5'CTC AAT AAA GCT TGC CTT GAG TGC and Gag ID RV 5'AAG
437 CGA GGA GCT GTT CAC TGC CAT CCT GGT CGA GCT ACC CAT ACG ATG TTC CAG ATT ACG. PCR
438 amplification was accomplished using SuperFi DNA polymerase in a 50µl reaction with 33.5µl of
439 the purified cDNA as input using the following conditions: initial denaturation for 3min at 98°C,
440 followed by 40 cycles of denaturation for 20sec at 98°C, annealing for 30sec at 60°C and
441 extension for 1min at 72°C, and final extension for 2min at 72°C. The Gag amplicon was gel

442 purified and the concentration was determined. The purified product was further used for
443 library construction.

444 **Generation of Vif amplicon from HIV9 sample.** One and a half microliters from clinical sample
445 HIV9 were reverse transcribed using SuperScript IV RT and random hexamer priming. Five
446 microliters of the purified RT reaction were used to set-up a PCR reaction with SuperFi DNA
447 polymerase to amplify ~600 bp region spanning HIV-1 Vif gene using primers vif FW 5'AGG GAT
448 TAT GGA AAA CAG ATG GCA GGT and vif RV 5'CTT AAG CTC CTC TAA AAG CTC TAG TG, and the
449 following program: initial denaturation for 3min at 98°C, followed by 40 cycles of denaturation
450 for 20sec at 98°C, annealing for 30sec at 60°C and extension for 30min at 72°C, and final
451 extension for 5min at 72°C. The amplicon was gel purified and the concentration was
452 determined. The purified product was further used for library construction.

453 **Generation of amplicons from RSV clinical samples**

454 **Clinical RSV samples.** Nasopharyngeal samples of 25 patients hospitalized at Chaim Sheba
455 Medical Center (Table S6) were collected into Virocult liquid viral transport medium (LVTM)
456 (Medical Wire & Equipment Co, Wiltshire, United Kingdom) and stored at -70°C. Five hundred
457 microliters of each sample were extracted and purified using easyMAG according to the
458 manufacturer's instructions. A primer specific to the glycoprotein protein G was designed with a
459 15 N-bases unique barcode followed by a linker sequence for subsequent PCR, RSV G RT 5'TAC
460 CCA TAC GAT GTT CCA GAT TAC GNN NNN NNN NNN NNN NGC AAA TGC AAM CAT GTC CAA AA.
461 Eight microliters of each sample were reverse transcribed as described in "Generation of Gag
462 amplicon with primer-ID from HIV-1" section.

463 **RSV control sample.** In the absence of an RSV plasmid, we used human rhinovirus (RV) plasmid
464 (a kind gift by Ann Palmenberg (University of Wisconsin-Madison, WI, USA) to generate a
465 homogeneous control that was run on the same Nextseq run as the RSV samples, similar to the
466 described above. In vitro transcribed RNA underwent RT using SuperScript IV with the following
467 primer: RV14 5' TAC GCA TAC GAT GTT CCA GAN NNN NNN NNN NNN NNN NAT AAA CTC
468 CTA CTT CTA CTC AAA TTA AGT GTC. PCR amplification using Q5 DNA polymerase with the
469 following primers was performed: p3.26 FW 5' TTA AAA CAG CGG ATG GGT ATC CCA C and p3.26
470 RV 5'ATG GTG AGC AAG GGC GAG GAG CTG TTC ACC GGG GTG GTG CTA CGC ATA CGA TGT TCC
471 AGA.

472 **Generation of a glycoprotein-fusion protein amplicon and polymerase amplicons.** PCR
473 amplification was accomplished using Q5 DNA polymerase in 50µl reactions with 15µl of the
474 purified cDNA as input. The following conditions were used for the glycoprotein-fusion protein
475 amplicon: initial denaturation for 3min at 98°C, followed by 40 cycles of denaturation for 20sec
476 at 98°C, annealing for 30sec at 58°C and extension for 3.5min at 72°C, and final extension for
477 5min at 72°C, using the following primers: Extension FW 5'AAG CGA GGA GCT GTT CAC TGC CAT
478 CCT GGT CGA GCT ACC CAT ACG ATG TTC CAG ATT ACG and RSV G and F RV 5'TGA CAG TAT TGT
479 ACA CTC TTA. For the polymerase amplicon, the following conditions were used: initial
480 denaturation for 3min at 98°C, followed by 40 cycles of denaturation for 20sec at 98°C,
481 annealing for 30sec at 60°C and extension for 8min at 72°C, and final extension for 5min at 72°C,
482 using the following primers: RSV L FW 5'GGA CAA AAT GGA TCC CAT TAT T and RSV L RV 5'GAA
483 CAG TAC TTG CAY TTT CTT AC. The amplicons were beads purified and joint together at equal
484 amounts. Concentration was determined, and the product was further used for NextSeq library
485 construction.

486 **Generation of a UL54 amplicon from CMV clinical samples**

487 Clinical DNA samples of recently infected patients (see Table S6) were obtained and purified as
488 described previously (Weisblum, et al. 2017). Since CMV is a DNA virus, no reverse transcription
489 step was needed. To generate a homogeneous control sample, the UL54 gene from TB40/E
490 strain was cloned onto a pGEM-t plasmid as described previously (Weisblum, et al. 2017). The
491 samples were diluted to 30,000 copies per PCR amplification reaction, which was set-up using
492 the Q5 DNA polymerase. The primers used to amplify the UL54 gene were UL54 FW 5'TCA ACA
493 GCA TTC GTG CGC CTT and UL54 RV 5'ATG TTT TTC AAC CCG TAT CTG AGC GGC, and the
494 following PCR protocol was executed: initial denaturation for 3min at 98C, followed by 38 cycles
495 of denaturation for 20sec at 98C, annealing for 20sec at 65C and extension for 3min at 72C, and
496 final extension for 5min at 72C. The amplicons were beads purified and their concentrations
497 were determined. The purified products were further used for MiSeq library construction with
498 the following change, 0.875ng of DNA were used as input for tagmentation instead of 0.85ng.

499 **MiSeq/Nextseq Libraries construction**

500 PCR fragmentation and indexing of samples for sequencing was performed using the Nextera XT
501 DNA Library Prep Kit (Illumina, San Diego, CA, USA) with the following adjustments to the
502 manufacturer instructions; (1) In order to get a short insert size of ~250bp, 0.85 ng of input DNA

503 was used for tagmentation; (2) No neutralization of the tagmentation buffer was done, as
504 described previously (Baym, et al. 2015); (3) For library amplification of the tagmented DNA, the
505 Nextera XT DNA library prep PCR reagents were replaced with high-fidelity DNA polymerase
506 reagents (the same DNA polymerase that was used for the amplicon generation). The PCR
507 reaction (50 μ l total) was set as depicted. Directly to the tagmented DNA, index 1 (i5, illumina,
508 5 μ l), index 2 (i7, illumina, 5 μ l), buffer (10 μ l), high-fidelity DNA polymerase (0.5 μ l), dNTPs
509 (10mM, 1 μ l) and nuclease-free water (8.5 μ l) were added; (4) Amplification was performed with
510 annealing temperature set to 63°C instead of 55°C, as introduced previously (Baym, et al. 2015)
511 and final extension for 2min; (5) Post-amplification clean-up was achieved using AMPure XP
512 beads in a double size-selection manner (Bronner, et al. 2014), to remove both too large and too
513 small fragments in order to maximize the fraction of fully overlapping read pairs. For the first
514 size-selection, 32.5 μ l of beads (0.65X ratio) were added to bind the large fragments. These
515 beads were separated and discarded. For the second-size selection, 10 μ l of beads (0.2X ratio)
516 were added to the supernatant to allow binding of intermediate fragments, and the supernatant
517 containing the small fragments was discarded. The intermediate fragments were eluted and
518 their size was determined using a high-sensitivity DNA tape in TapeStation 4200 (Agilent, Santa
519 Clara, CA, USA). A mean size of \sim 370bp, corresponding to the desired insert size of \sim 250bp, was
520 achieved; And (6) Normalization and pooling was performed manually.

521 NextSeq: The longest NextSeq read length is 150bp, we hence selected for a shorter insert size
522 of 270bp, compared to the desired 370bp insert size for the MiSeq platform. The first size
523 selection of the post-NexteraXT amplification cleanup was performed using 42.5 μ l of AMPure XP
524 beads (0.85X ratio) (Bronner, et al. 2014).

525 **AccuNGS Development**

526 The AccuNGS protocol was evaluated using HIV-1 DNA plasmid (Peden, et al. 1991)). Our
527 underlying assumption was that this DNA starting material is homogenous with respect to the
528 theoretical error rate we calculated. This assumption was based on the fact that we used low-
529 copy plasmids that were grown in *Escherichia coli*, and only a single colony was subsequently
530 sequenced. The mutation rate of *E. coli* is in the order of 1×10^{-10} errors/base/replication (Jee, et
531 al. 2016), and accordingly, error rates in the purified plasmids are expected to be much lower
532 than the original expected protocol mean error of $\sim 10^{-5}$ (Table S7). Table S1 summarizes the
533 different differential sequencing that was performed during the Methods development, in an

534 attempt to test the contribution of each one of the stages of the protocol to the error rate of
535 the method (see Supplementary Text).

536 **Preparation of plasmids.** In order to maintain the plasmid stock as homogenous as possible,
537 plasmids were transformed to a chemically competent bacteria cells [DH5alpha (BioLab, Israel)
538 or TG1 [A kind gift by Itai Benhar (Tel Aviv University, Tel Aviv, Israel)]] using a standard heat-
539 shock protocol. Based on the fact that *E. coli* doubling time is 20 minutes in average using rich
540 growing medium (Sezonov, et al. 2007), a single colony was selected and grown to a maximum
541 of 100 generations. Plasmids were column purified using HiYield™ Plasmid Mini Kit (RBC
542 Bioscience, New Taipei City, Taiwan) and stored at -20°C until use.

543 **Construction of baseline control amplicon.** A baseline control amplicon (Table S1) was based on
544 clonal amplification and sequencing of the pLAI.2 plasmid, which contains a full-length HIV-1_{LAI}
545 proviral clone (Peden, et al. 1991) (obtained through the NIH AIDS Reagent Program, Division of
546 AIDS, NIAID, NIH: pLAI.2 from Dr. Keith Peden, courtesy of the MRC AIDS Directed Program). The
547 Integrase region of pLAI.2 was amplified using primers: KLV70 - 5'TTC RGG ATY AGA AGT AAA
548 YAT AGT AAC AG and KLV84 - 5'TCC TGT ATG CAR ACC CCA ATA TG (Moscona, et al. 2017). PCR
549 amplification was conducted using SuperFi DNA Polymerase in a 50µl reaction using 20-40 ng of
550 the plasmid as input. Amplification in a thermal cycler was performed as follows: initial
551 denaturation for 3min at 98°C, followed by 40 cycles of denaturation for 20sec at 98°C,
552 annealing for 30sec at 60°C and extension for 1min at 72°C, and final extension for 2min at 72°C.
553 In parallel, an alternative PCR reaction was up using Q5 DNA Polymerase. The Integrase
554 amplicon was gel purified and concentration was determined. The purified product was further
555 used for library construction.

556 **Construction of AmpR and RpoB control amplicons.** For generating the AmpR amplicon, the
557 conserved *AmpR* gene was amplified from pLAI.2 plasmid using primers: AmpR FW - 5'AAA GTT
558 CTG CTA TGT GGC GC and AmpR RV - 5'GGT CTG ACA GTT ACC AAT GC. PCR amplification was
559 carried out as described above, except for extension duration of 30sec instead of 1min.
560 Similarly, the conserved *RpoB* gene was amplified from the bacteria genome using the following
561 primers: RpoB FW 5'ATG GTT TAC TCC TAT ACC GA and RpoB RV 5'GTG ATC CAG ATC GTT GGT G
562 and the following PCR program: initial denaturation for 3min at 98°C, followed by 40 cycles of
563 denaturation for 10sec at 98°C, annealing for 10sec at 60°C and extension for 4sec at 72°C, and

564 final extension for 2min at 72°C. The AmpR and RpoB amplicons were gel purified and their
565 concentration was determined. The purified product was further used for library construction.

566 **Construction of alternative purification amplicons.** The agarose gel purification step of the
567 amplified integrase gene was replaced with other purification methods; (1) For the gel-free
568 sample, the amplified integrase gene was purified using 25µl of AMPure XP beads (0.5X ratio),
569 and (2) For the ExoSap sample, 10µl of the amplified integrase gene were mixed with 4µl of
570 ExoSAP-IT™ PCR Product Cleanup Reagent (Thermo Fisher Scientific) and incubated according to
571 the manufacturer's instructions. No other changes in the generation of amplicon protocol were
572 made.

573 **Construction of a PCR-free control amplicon.** For the PCR-free sample, 10µg of pLAI.2 plasmid
574 was digested using the restriction enzymes: NheI, StuI and XcmI (NEB) according to the
575 manufacturer's instructions. A ~1500bp fragment containing the integrase region was gel
576 purified and concentration was determined by Qubit. The purified product was further used for
577 library construction.

578 **Construction of RNA control amplicons.** We used a plasmid containing the full cDNA of
579 Coxsackie virus B3 (CVB3) under a T7 promoter that was a kind gift from Marco Vignuzzi (Institut
580 Pasteur, Paris, France). Ten micrograms of this plasmid were linearized using Sall (NEB), beads
581 purified (0.5X ratio) and then *in-vitro* transcribed using T7 RNA polymerase (NEB) according to
582 the manufacturer's instructions. The transcribed RNA was bead purified (0.5X ratio) and reverse
583 transcribed with random hexamers using SuperScript III RT. Four microliters of the reverse
584 transcription reaction were used as template for a PCR reaction using primers: CVB FW 5'GGA
585 GAG AAG GTC AAC TCT ATG GAA GC and CVB RV 5'TAC CAC CCT GTA GTT CCC CA, which amplify
586 a ~1500bp fragment within the CVB genome. PCR reaction (50µl total) was set and amplified
587 using SuperFi DNA polymerase as follows: initial denaturation for 3min at 98°C, followed by 40
588 cycles of denaturation for 20sec at 98°C, annealing for 30sec at 60°C and extension for 15sec at
589 72°C, and final extension for 2min at 72°C. The CVB amplicon was gel purified and the
590 concentration measured. The purified product was further used for library construction.

591 **Alternative library purification methods.** For the AMPure XP beads-free sample, post-
592 amplification clean-up by double size-selection was replaced with an agarose gel purification of
593 a ~370bp fragment, with no other changes in the library construction protocol.

594 **Alternative tagmentation sample.** For the alternative tagmentation sample, a 250bp amplicon
595 within the integrase region was designed, using specific primers with an overhang
596 corresponding to the sequence inserted during the tagmentation step of the NexteraXT DNA
597 library prep kit, NexteraXT free FW 5'TCG TCG GCA GCG TCA GAT GTG TAT AAG AGA CAG ACT
598 TGT CCA TGC ATG GCT TCT C and NexteraXT free RV 5'GTC TCG TGG GCT CGG AGA TGT GTA TAA
599 GAG ACA GTC TAT CTG GCA TGG GTA CCA GCA. PCR reaction was set up using SuperFi DNA
600 polymerase and carried out as follows: initial denaturation for 3min at 98°C, followed by 40
601 cycles of denaturation for 20sec at 98°C, annealing for 30sec at 62°C and extension for 15sec at
602 72°C, and final denaturation for 2min at 72°C. The PCR product was gel purified and the
603 concentration was measured by Qubit. The purified product was indexed by a succeeding PCR
604 amplification using primers corresponding to i5 and i7 NexteraXT primers (IDT) as mentioned
605 previously (Baym, et al. 2015) at a final concentration of 1uM. The PCR reaction was set up using
606 SuperFi DNA polymerase and amplified as detailed: initial denaturation for 3min at 98°C,
607 followed by 12 cycles of denaturation for 20sec at 98°C, annealing for 30sec at 63°C and
608 extension for 30sec at 72°C, and final extension for 2min at 72C. Size selection was achieved by
609 gel purification of ~370bp fragments.

610 **Sequencing.** Sequencing of all synthetic samples, the HIV-1, and CMV samples was performed
611 on the Illumina MiSeq platform using MiSeq Reagent Kit v2 (500-cycles, equal to 250x2 paired-
612 end reads) (Illumina). Sequencing of the RSV-1 samples and a dedicated synthetic sample was
613 performed on the Illumina NextSeq 500 platform using NextSeq 500/550 High Output Kit (300-
614 cycles, equal to 150x2 paired-end reads) (Illumina).

615 **Construction of synthetically mixed populations.** A pLAI.2 plasmid was mixed with a pNL4.3
616 plasmid [a kind gift from Eran Bacharach (Tel Aviv University, Tel Aviv, Israel)] using serial
617 dilutions to achieve the following ratios: 1:100, 1:500, 1:1,000, 1:5,000 and 1:10,000 with pLAI.2
618 held as the major strain. Concentrations were measured and each dilution was generated
619 independently to obtain three technical replicas for each ratio. The plasmids mixture was
620 further used for library construction.

621 **Barcode serial dilution test**

622 The pLAI.2 plasmid was used to generate an RNA pool. Five micrograms of this plasmid were
623 linearized using Sall (NEB) and beads purified (0.5X ratio). T7 polymerase promotor was added
624 to the linearized plasmid using T7 extension FW 5'TAA TAC GAC TCA CTA TAG CTG GGA GCT CTC

625 TGG CTA AC and the RV 5'GAG ACT CCC TGA CCC AGA TGC C in a PCR reaction using Q5 DNA
626 polymerase with the following program: initial denaturation for 3min at 98°C, followed by 40
627 cycles of denaturation for 10sec at 98°C, annealing for 10sec at 65°C and extension for 3min at
628 72°C, and final extension for 5min at 72°C. Four microliters of the reaction was in-vitro
629 transcribed using T7 RNA polymerase according to the manufacturer's instructions. The
630 transcribed RNA was beads purified (0.5X ratio). The purified RNA was serially diluted and for
631 each dilution two reactions were set-up: a primer-ID reaction (as described in the section
632 "Generation of Gag amplicon with primer-ID from HIV-1") and a random hexamer based RT
633 reaction (as described in the section "Construction of RNA control amplicons"). In order to
634 compare these reactions, for the PCR amplification of the random hexamer based RT reaction,
635 we used the following primers: GAG FW 5'CTC AAT AAA GCT TGC CTT GAG TGC and RTgene RV
636 5'ACT GTA TCA TCT GCT CCT GTA TCT corresponding to the primer-ID reaction primers without a
637 barcode. The same PCR program was used for both reactions. The PCR reactions were gel
638 purified and concentration was measured.

639 **Reads processing and base calling**

640 The paired-end reads from each control library were aligned against the reference sequence of
641 that control using an in-house script that relies on BLAST command-line tool (Altschul, et al.
642 1990; Altschul, et al. 1997; Camacho, et al. 2009). The paired-end reads from the clinical
643 samples were aligned against: HIV-1 subtype B HXB2 reference sequence (GenBank accession
644 number K03455.1), RSV reference sample (GenBank accession number U39661), CMV reference
645 sample Merlin (GenBank accession number NC_006273), and then realigned against the
646 consensus sequence obtained for each sample. Bases were called using an in-house script only if
647 the forward and reverse reads agreed and their average Q-score was above an input threshold
648 (30 or 38). At each position, for each alternative base, we calculate mutation frequencies by
649 dividing the number of reads bearing the mutation by loci coverage. Positions were retained for
650 analysis only if sequenced to a depth of at least 100,000 reads. In order to analyze the errors in
651 the sequencing process we used Python 3.7.3 (Anaconda distribution) with the following
652 packages: pandas 0.25.1 (McKinney 2010), matplotlib 3.1.0 (Caswell, et al. 2019), seaborn 0.9.0
653 (Waskom, et al. 2018), numpy 1.16.3 (Oliphant 2006; Walt, et al. 2011) and scipy 1.2.1 (Jones, et
654 al. 2016). Distributions of errors on control plasmids were compared using two-tailed t-test or
655 two-tailed Mann-Whitney U test.

656 **Variant calling**

657 In order to facilitate discrimination of true variants from AccuNGS process artifacts, we created
658 a variant caller based on two principles: (i) positions that exhibit relatively high level of error on
659 a control sample are error-prone for the clinical sample as well; and (ii) process errors on a
660 control sample follow a gamma distribution. A gamma distribution was fitted for each mutation
661 type in the control sequence. In order to detect and remove outliers from the fitting process we
662 used the “three-sigma-rule”, and positions that showed error higher than three standard
663 deviations from the mean of the fitted distribution were removed. For these rare loci a base was
664 called only if the mutation was more prevalent in the sample by an order of magnitude. For G>A
665 transition mutations, four distinct gamma distributions were fitted, corresponding for all four
666 G>A combinations with preceding nucleotide. Accordingly, for C>T transition mutations four
667 gamma distributions were fitted as well, on the four C>T reverse complement mutations of the
668 G>A mutations. For establishing Figures 2-5, variants were called on the input sample only if a
669 mutation was in the extreme 1% of the relevant gamma distribution fitted using the
670 corresponding control.

671 **Serial dilutions analysis**

672 All dilutions were mapped against the pLAI.2 reference sequence using the default parameters
673 and Q30 as the minimal base quality threshold. Positions with insufficient coverage, defined as
674 having less than 5 times the inverse of the frequency of the minor strain were filtered out, as
675 well as positions with minor variant frequency of above 0.5% in the homogeneous control
676 sample.

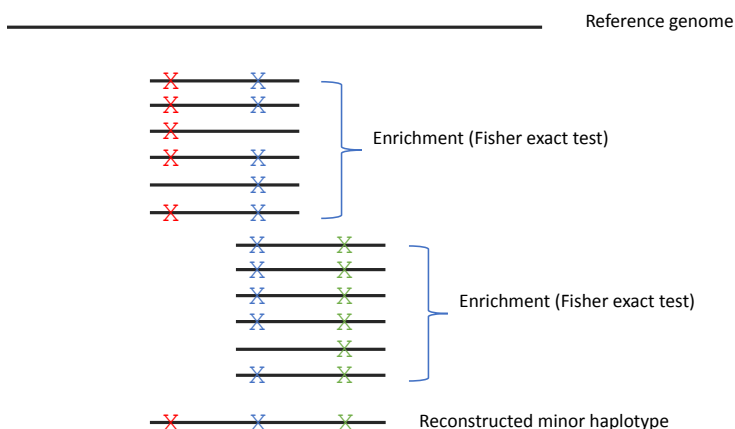
677 **Diversity calculation**

678 Transition nucleotide diversity π was calculated per sample using positions with at least 5,000x
679 coverage, using the formulas described in (Zhao and Illingworth 2019), but excluding
680 transversion variants. Sites whose variant frequencies weren't statistically different than the
681 background control were considered as variants of count = 0 for the calculation.

682 **Haplotypes inference**

683 To infer potential haplotypes, we used a two-step process, illustrated in Figure 6. First we
684 identified all pairs of non-consensus variants (the most common minor variant at each site) that
685 were statistically enriched when present on the same reads. Next we attempted to
686 "concatenate" multiple pairs into a longer stretch based on a shared mutation present in two

687 different pairs of variants. In order to find statistically enriched pairs, we considered all sites that
688 may be linked on the same reads (up to 250 bases, which is the maximal length of an Illumina
689 read). For each pair of loci we created a contingency table for the appearance of each variant
690 alone, the two variants together and no variant at all. We then used a one-tailed Fisher exact
691 test to obtain a p-value for the pair, and considered only p-values lower than 10^{-15} , to account of
692 multiple testing. From this contingency table we also extracted the frequency at which the two
693 variants co-occur. We repeated the process for all possible pairs of loci. This resulted in many
694 short haplotype stretches of 250 bases spanning two loci each. We then performed
695 "concatenation" of pairs of loci that had (1) at least one shared position and (2) a similar
696 frequency of co-occurrence, defined here as up to an order of magnitude in difference. Such
697 concatenated loci formed a longer stretch and its frequency was calculated as the mean
698 frequency of its components, i.e., the average frequency of all individual pairs added to this
699 stretch so far. For each sample, we iteratively attempted to concatenate all pairs of loci, starting
700 from the highest frequency pair to the least common pair, until no pairs could further merge.



701

702 **Figure 6. Illustration of method for haplotype reconstruction.** The method searches for
703 enrichment of pairs of mutations on the same read, and concatenation of enriched reads that
704 share a mutation into a reconstructed minor haplotype. Notably, the concatenation approach is
705 suitable for populations with limited diversity, as is the case in acute infections; in highly diverse
706 populations, many haplotypes may share the "blue" mutation illustrated in the figure.

707 **CODE AVAILABILITY**

708 We have developed the following computational resources that complement the AccuNGS
709 sequencing protocol:

- 710 (a) Base coverage calculator. AccuNGS relies on overlapping read pairs and high Q-scores
711 for both reads of a pair. The calculator receives as input the length of the target regions
712 and the desired coverage, and outputs the recommended number of reads required for
713 sequencing each sample.
- 714 (b) Computational pipeline for computing the number of unique RNA molecules sequenced,
715 based on primer-ID barcodes (see Supplementary Text).
- 716 (c) Computational pipeline for base-calling and inferring site by site base frequencies.
- 717 (d) Computational pipeline for inferring haplotypes.

718 All resources are freely available at <https://github.com/SternLabTAU/AccuNGS>.

719 **ACCESSION NUMBERS**

720 The datasets generated and reported in this study were deposited in the Sequencing Read
721 Archive (SRA, available at <https://www.ncbi.nlm.nih.gov/sra>), under BioProject PRJNA476431.
722 Frequencies of mutations following base calling will be available in Zenodo (<https://zenodo.org/>)
723 upon acceptance.

REFERENCES

- 724 Acevedo A, Brodsky L, Andino R. 2014. Mutational and fitness landscapes of an RNA virus
725 revealed through population sequencing. *Nature* 505:686-690.
- 726 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *J*
727 *Mol Biol* 215:403-410.
- 728 Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. 1997. Gapped BLAST
729 and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res*
730 25:3389-3402.
- 731 Baym M, Kryazhimskiy S, Lieberman TD, Chung H, Desai MM, Kishony R. 2015. Inexpensive
732 multiplexed library preparation for megabase-sized genomes. *PLoS one* 10:e0128036.
- 733 Beale RC, Petersen-Mahrt SK, Watt IN, Harris RS, Rada C, Neuberger MS. 2004. Comparison of
734 the differential context-dependence of DNA deamination by APOBEC enzymes: correlation with
735 mutation spectra in vivo. *J Mol Biol* 337:585-596.
- 736 Bishop KN, Holmes RK, Sheehy AM, Davidson NO, Cho SJ, Malim MH. 2004. Cytidine
737 deamination of retroviral DNA by diverse APOBEC proteins. *Curr Biol* 14:1392-1396.
- 738 Boucher CA, Bobkova MR, Geretti AM, Hung CC, Kaiser R, Marcelin AG, Streinu-Cercel A, van
739 Wyk J, Dorr P, Vandamme AM. 2018. State of the Art in HIV Drug Resistance: Science and
740 Technology Knowledge Gap. *AIDS reviews* 20:27-42.
- 741 Brodin J, Hedskog C, Heddini A, Benard E, Neher RA, Mild M, Albert J. 2015. Challenges with
742 using primer IDs to improve accuracy of next generation sequencing. *PLoS one* 10:e0119123.
- 743 Bronner IF, Quail MA, Turner DJ, Swerdlow H. 2014. Improved Protocols for Illumina Sequencing.
744 *Curr Protoc Hum Genet* 80:18 11-42.
- 745 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009. BLAST+:
746 architecture and applications. *BMC bioinformatics* 10:421.

747 Casadella M, Paredes R. 2017. Deep sequencing for HIV-1 clinical management. *Virus Res*
748 239:69-81.

749 Caswell T, Droettboom M, Hunter J. 2019. matplotlib/matplotlib v3. 1.0, doi: 10.5281/zenodo.
750 2893252. In.

751 Chen-Harris H, Borucki MK, Torres C, Slezak TR, Allen JE. 2013. Ultra-deep mutant spectrum
752 profiling: improving sequencing accuracy using overlapping read pairs. *BMC genomics* 14:96.

753 Clutter DS, Jordan MR, Bertagnolio S, Shafer RW. 2016. HIV-1 drug resistance and resistance
754 testing. *Infect Genet Evol* 46:292-307.

755 Cudini J, Roy S, Houldcroft CJ, Bryant JM, Depledge DP, Tutill H, Veys P, Williams R, Worth AJJ,
756 Tamuri AU, et al. 2019. Human cytomegalovirus haplotype reconstruction reveals high diversity
757 due to superinfection and evidence of within-host recombination. *Proc Natl Acad Sci U S A*
758 116:5693-5698.

759 Cuevas JM, Geller R, Garijo R, Lopez-Aldeguer J, Sanjuan R. 2015. Extremely High Mutation Rate
760 of HIV-1 In Vivo. *PLoS Biol* 13:e1002251.

761 Döring M, Büch J, Friedrich G, Pironti A, Kalaghatgi P, Knops E, Heger E, Obermeier M, Däumer
762 M, Thielen A, et al. 2018. geno2pheno[ngs-freq]: a genotypic interpretation system for
763 identifying viral drug resistance using next-generation sequencing data. *Nucleic Acids Research*
764 46:W271-W277.

765 Duffy S, Shackelton LA, Holmes EC. 2008. Rates of evolutionary change in viruses: patterns and
766 determinants. *Nat Rev Genet* 9:267-276.

767 Gu W, Miller S, Chiu CY. 2019. Clinical metagenomic next-generation sequencing for pathogen
768 detection. *Annual Review of Pathology: Mechanisms of Disease* 14:319-338.

769 Hache G, Mansky LM, Harris RS. 2006. Human APOBEC3 proteins, retrovirus restriction, and HIV
770 drug resistance. *AIDS reviews* 8:148-157.

771 Henn MR, Boutwell CL, Charlebois P, Lennon NJ, Power KA, Macalalad AR, Berlin AM, Malboeuf
772 CM, Ryan EM, Gnerre S, et al. 2012. Whole genome deep sequencing of HIV-1 reveals the impact
773 of early minor variants upon immune recognition during acute infection. *Plos Pathog*
774 8:e1002529.

775 Huber M, Metzner KJ, Geissberger FD, Shah C, Leemann C, Klimkait T, Boni J, Trkola A, Zagordi O.
776 2017. MinVar: A rapid and versatile tool for HIV-1 drug resistance genotyping by deep
777 sequencing. *J Virol Methods* 240:7-13.

778 Illingworth CJR, Roy S, Beale MA, Tutill H, Williams R, Breuer J. 2017. On the effective depth of
779 viral sequence data. *Virus Evol* 3:vex030.

780 Imashimizu M, Oshima T, Lubkowska L, Kashlev M. 2013. Direct assessment of transcription
781 fidelity by high-resolution RNA sequencing. *Nucleic Acids Res* 41:9090-9104.

782 Jabara CB, Jones CD, Roach J, Anderson JA, Swanstrom R. 2011. Accurate sampling and deep
783 sequencing of the HIV-1 protease gene using a Primer ID. *Proc Natl Acad Sci U S A* 108:20166-
784 20171.

785 Jee J, Rasouly A, Shamovsky I, Akivis Y, Steinman SR, Mishra B, Nudler E. 2016. Rates and
786 mechanisms of bacterial mutagenesis from maximum-depth sequencing. *Nature* 534:693-696.

787 Jones E, Oliphant T, Peterson P. 2016. SciPy: Open source scientific tools for Python, 2001. In.
788 Keele BF, Giorgi EE, Salazar-Gonzalez JF, Decker JM, Pham KT, Salazar MG, Sun C, Grayson T,
789 Wang S, Li H, et al. 2008. Identification and characterization of transmitted and early founder
790 virus envelopes in primary HIV-1 infection. *Proc Natl Acad Sci U S A* 105:7552-7557.

791 Kennedy SR, Schmitt MW, Fox EJ, Kohn BF, Salk JJ, Ahn EH, Prindle MJ, Kuong KJ, Shen JC,
792 Risques RA, et al. 2014. Detecting ultralow-frequency mutations by Duplex Sequencing. *Nat*
793 *Protoc* 9:2586-2606.

794 Kijak GH, Sanders-Buell E, Chenine AL, Eller MA, Goonetilleke N, Thomas R, Leviyang S, Harbolick
795 EA, Bose M, Pham P, et al. 2017. Rare HIV-1 transmitted/founder lineages identified by deep
796 viral sequencing contribute to rapid shifts in dominant quasispecies during acute and early
797 infection. *Plos Pathog* 13:e1006510.

798 Kircher M, Sawyer S, Meyer M. 2011. Double indexing overcomes inaccuracies in multiplex
799 sequencing on the Illumina platform. *Nucleic Acids Research* 40:e3-e3.

800 Lau JW, Kim YI, Murphy R, Newman R, Yang X, Zody M, DeVincenzo J, Grad YH. 2017. Deep
801 sequencing of RSV from an adult challenge study and from naturally infected infants reveals
802 heterogeneous diversification dynamics. *Virology* 510:289-296.

803 Lessler J, Reich NG, Brookmeyer R, Perl TM, Nelson KE, Cummings DA. 2009. Incubation periods
804 of acute respiratory viral infections: a systematic review. *Lancet Infect Dis* 9:291-300.

805 Liddicoat BJ, Piskol R, Chalk AM, Ramaswami G, Higuchi M, Hartner JC, Li JB, Seeburg PH,
806 Walkley CR. 2015. RNA editing by ADAR1 prevents MDA5 sensing of endogenous dsRNA as
807 nonself. *Science* 349:1115-1120.

808 Lou DI, Hussmann JA, McBee RM, Acevedo A, Andino R, Press WH, Sawyer SL. 2013. High-
809 throughput DNA sequencing errors are reduced by orders of magnitude using circle sequencing.
810 *Proc Natl Acad Sci U S A* 110:19872-19877.

811 Malim MH. 2009. APOBEC proteins and intrinsic resistance to HIV-1 infection. *Philos Trans R Soc*
812 *Lond B Biol Sci* 364:675-687.

813 McCrone JT, Lauring AS. 2016. Measurements of Intrahost Viral Diversity Are Extremely Sensitive
814 to Systematic Errors in Variant Calling. *J Virol* 90:6884-6895.

815 McCrone JT, Woods RJ, Martin ET, Malosh RE, Monto AS, Lauring AS. 2018. Stochastic processes
816 constrain the within and between host evolution of influenza virus. *Elife* 7:e35962.

817 McKinney W editor. *Proceedings of the 9th Python in Science Conference*. 2010.

818 Meacham F, Boffelli D, Dhahbi J, Martin DI, Singer M, Pachter L. 2011. Identification and
819 correction of systematic error in high-throughput sequence data. *BMC bioinformatics* 12:451.

820 Moscona R, Ram D, Wax M, Bucris E, Levy I, Mendelson E, Mor O. 2017. Comparison between
821 next-generation and Sanger-based sequencing for the detection of transmitted drug-resistance
822 mutations among recently infected HIV-1 patients in Israel, 2000–2014. *J INT AIDS SOC* 20.

823 Newman AM, Lovejoy AF, Klass DM, Kurtz DM, Chabon JJ, Scherer F, Stehr H, Liu CL, Bratman SV,
824 Say C, et al. 2016. Integrated digital error suppression for improved detection of circulating
825 tumor DNA. *Nat Biotechnol* 34:547-555.

826 Nielsen R, Yang Z. 1998. Likelihood models for detecting positively selected amino acid sites and
827 applications to the HIV-1 envelope gene. *Genetics* 148:929-936.

828 Oliphant TE. 2006. *A guide to NumPy*: Trelgol Publishing USA.

829 Parrish CR, Holmes EC, Morens DM, Park E-C, Burke DS, Calisher CH, Laughlin CA, Saif LJ, Daszak
830 P. 2008. Cross-species virus transmission and the emergence of new epidemic diseases.
831 *Microbiol. Mol. Biol. Rev.* 72:457-470.

832 Peden K, Emerman M, Montagnier L. 1991. Changes in growth properties on passage in tissue
833 culture of viruses derived from infectious molecular clones of HIV-1LAI, HIV-1MAL, and HIV-1ELI.
834 *Virology* 185:661-672.

835 Pfaller CK, Donohue RC, Nersisyan S, Brodsky L, Cattaneo R. 2018. Extensive editing of cellular
836 and viral double-stranded RNA structures accounts for innate immunity suppression and the
837 proviral activity of ADAR1p150. *PLoS Biol* 16:e2006577.

838 Preston JL, Royall AE, Randel MA, Sikkink KL, Phillips PC, Johnson EA. 2016. High-specificity
839 detection of rare alleles with Paired-End Low Error Sequencing (PELE-Seq). *BMC genomics*
840 17:464.

841 Ram D, Leshkowitz D, Gonzalez D, Forer R, Levy I, Chowers M, Lorber M, Hindiyeh M, Mendelson
842 E, Mor O. 2015. Evaluation of GS Junior and MiSeq next-generation sequencing technologies as
843 an alternative to Trugene population sequencing in the clinical HIV laboratory. *Journal of*
844 *virological methods* 212:12-16.

845 Reid-Bayliss KS, Loeb LA. 2017. Accurate RNA consensus sequencing for high-fidelity detection of
846 transcriptional mutagenesis-induced epimutations. *Proc Natl Acad Sci U S A* 114:9415-9420.

847 Renzette N, Bhattacharjee B, Jensen JD, Gibson L, Kowalik TF. 2011. Extensive genome-wide
848 variability of human cytomegalovirus in congenitally infected infants. *Plos Pathog* 7:e1001344.

849 Salazar-Gonzalez JF, Bailes E, Pham KT, Salazar MG, Guffey MB, Keele BF, Derdeyn CA, Farmer P,
850 Hunter E, Allen S. 2008. Deciphering human immunodeficiency virus type 1 transmission and
851 early envelope diversification by single-genome amplification and sequencing. *Journal of*
852 *Virology* 82:3952-3970.

853 Samuel CE. 2012. ADARs: viruses and innate immunity. *Curr Top Microbiol Immunol* 353:163-
854 195.

855 Sanjuan R, Nebot MR, Chirico N, Mansky LM, Belshaw R. 2010. Viral mutation rates. *J Virol*
856 84:9733-9748.

857 Schirmer M, Ijaz UZ, D'Amore R, Hall N, Sloan WT, Quince C. 2015. Insight into biases and
858 sequencing errors for amplicon sequencing with the Illumina MiSeq platform. *Nucleic Acids Res*
859 43:e37.

860 Schirmer M, Sloan WT, Quince C. 2014. Benchmarking of viral haplotype reconstruction
861 programmes: an overview of the capacities and limitations of currently available programmes.
862 *Brief Bioinform* 15:431-442.

863 Seibert SA, Howell CY, Hughes MK, Hughes AL. 1995. Natural selection on the gag, pol, and env
864 genes of human immunodeficiency virus 1 (HIV-1). *Mol Biol Evol* 12:803-813.

865 Sezonov G, Joseleau-Petit D, D'Ari R. 2007. Escherichia coli physiology in Luria-Bertani broth. *J*
866 *Bacteriol* 189:8746-8749.

867 Tan L, Coenjaerts FE, Houspie L, Viveen MC, van Bleek GM, Wiertz EJ, Martin DP, Lemey P. 2013.
868 The comparative genomics of human respiratory syncytial virus subgroups A and B: genetic
869 variability and molecular evolutionary dynamics. *J Virol* 87:8213-8226.

870 Theys K, Feder AF, Gelbart M, Hartl M, Stern A, Pennings PS. 2018. Within-patient mutation
871 frequencies reveal fitness costs of CpG dinucleotides and drastic amino acid changes in HIV.
872 *PLoS Genet* 14:e1007420.

873 Walt Svd, Colbert SC, Varoquaux G. 2011. The NumPy Array: A Structure for Efficient Numerical
874 Computation. *Computing in Science & Engineering* 13:22-30.

875 Wang K, Lai S, Yang X, Zhu T, Lu X, Wu Cl, Ruan J. 2017. Ultrasensitive and high-efficiency screen
876 of de novo low-frequency mutations by o2n-seq. *Nat Commun* 8:15335.

877 Waskom M, Botvinnik O, O'Kane D, Hobson P, Ostblom J, Lukauskas S, Qalieh A. 2018.
878 mwaskom/seaborn: v0. 9.0 (July 2018). DOI: <https://doi.org/10.5281/zenodo.1313201>.

879 Weisblum Y, Oiknine-Djian E, Zakay-Rones Z, Vorontsov O, Haimov-Kochman R, Nevo Y,
880 Stockheim D, Yagel S, Panet A, Wolf DG. 2017. APOBEC3A Is Upregulated by Human
881 Cytomegalovirus (HCMV) in the Maternal-Fetal Interface, Acting as an Innate Anti-HCMV
882 Effector. *J Virol* 91.

883 Whitmer SLM, Ladner JT, Wiley MR, Patel K, Dudas G, Rambaut A, Sahr F, Prieto K, Shepard SS,
884 Carmody E, et al. 2018. Active Ebola Virus Replication and Heterogeneous Evolutionary Rates in
885 EVD Survivors. *Cell Rep* 22:1159-1168.

886 Wilke CO, Novella IS. 2003. Phenotypic mixing and hiding may contribute to memory in viral
887 quasispecies. *BMC Microbiol* 3:11.

888 Yang X, Charlebois P, Macalalad A, Henn MR, Zody MC. 2013. V-Phaser 2: variant inference for
889 viral populations. *BMC genomics* 14:674.
890 Zanini F, Brodin J, Albert J, Neher RA. 2017. Error rates, PCR recombination, and sampling depth
891 in HIV-1 whole genome deep sequencing. *Virus Res* 239:106-114.
892 Zanini F, Brodin J, Thebo L, Lanz C, Bratt G, Albert J, Neher RA. 2015. Population genomics of
893 inpatient HIV-1 evolution. *Elife* 4:e11282.
894 Zanini F, Puller V, Brodin J, Albert J, Neher RA. 2017. In vivo mutation rates and the landscape of
895 fitness costs of HIV-1. *Virus Evol* 3:vex003.
896 Zhao L, Illingworth CJR. 2019. Measurements of intrahost viral diversity require an unbiased
897 diversity metric. *Virus Evol* 5:vey041.
898 Zhou S, Bednar MM, Sturdevant CB, Hauser BM, Swanstrom R. 2016. Deep Sequencing of the
899 HIV-1 env Gene Reveals Discrete X4 Lineages and Linkage Disequilibrium between X4 and R5
900 Viruses in the V1/V2 and V3 Variable Regions. *J Virol* 90:7142-7158.
901 Zhou S, Jones C, Mieczkowski P, Swanstrom R. 2015. Primer ID Validates Template Sampling
902 Depth and Greatly Reduces the Error Rate of Next-Generation Sequencing of HIV-1 Genomic
903 RNA Populations. *J Virol* 89:8540-8555.

904

905 **ACKNOWLEDGEMENTS**

906 The authors would like to thank Oded Kushnir, Danielle Miller and Yiska Weisblum for valuable
907 support, and for Drs. Neta Zuckerman, Tzachi Hagai and Shaul Pollak for critical reading of the
908 manuscript and helpful discussions.

909 **FUNDING**

910 This work was supported by the SAIA foundation; by the Israeli Science Foundation [1333/16 to
911 AS]; by the German Israeli Foundation [I-1096-411.8-2015 to AS]; by the United-States-Israel
912 Binational Science Foundation [2016555 to AS]; by the Edmond J. Safra center for bioinformatics
913 in Tel Aviv University [to MG, SH, TK].