

1 Validation of multiplex PCR sequencing assay of SIV

2

3 Ryan V. Moriarty^a, Nico Fesser^a, Matthew S. Sutton^a, Vanessa Venturi^b, Miles P. Davenport^b,

4 Timothy Schlub^c, Shelby L O'Connor^a

5

6 ^aDepartment of Pathology and Laboratory Medicine, University of Wisconsin-Madison, WI 53711

7 ^bInfection Analytics Program, Kirby Institute for Infection and Immunity, UNSW Sydney, Sydney,
8 New South Wales 2052, Australia

9 ^cUniversity of Sydney, Faculty of Medicine and Health, Sydney School of Public Health, Sydney
10 2000, NSW, Australia

11

12 Corresponding Author: Dr. Shelby O'Connor, sfeinberg@wisc.edu, 555 Science Dr, Madison,
13 WI 53711

14

15 Abstract:

16 Background: The generation of accurate and reproducible viral sequence data is necessary to
17 understand the diversity present in populations of RNA viruses isolated from clinical samples.
18 While various sequencing methods are available, they often require high quality templates and
19 high viral titer to ensure reliable data.

20 Methods: We modified a multiplex PCR and sequencing approach to characterize populations of
21 simian immunodeficiency virus (SIV) isolated from nonhuman primates. We chose this approach
22 with the aim of reducing the number of required input templates while maintaining fidelity and
23 sensitivity. We conducted replicate sequencing experiments using different numbers of quantified
24 viral RNA (vRNA) or viral cDNA as input material. We performed assays with clonal SIVmac239
25 to detect false positives, and we mixed SIVmac239 and a variant with 24 point mutations
26 (SIVmac239-24X) to measure variant detection sensitivity.

27 Results: We found that utilizing a starting material of quantified viral cDNA templates had a lower
28 rate of false positives and increased reproducibility when compared to that of quantified vRNA
29 templates. This study identifies the importance of rigorously validating deep sequencing methods
30 and including replicate samples when using a new method to characterize low frequency variants
31 in a population with a small number of templates.

32 Conclusions: Because the need to generate reproducible and accurate sequencing data from
33 diverse viruses from low titer samples, we modified a multiplex PCR and sequencing approach to
34 characterize SIV from populations from non-human primates. We found that increasing starting
35 template numbers increased the reproducibility and decreased the number of false positives
36 identified, and this was further seen when cDNA was used as a starting material. Ultimately, we
37 highlight the importance of vigorously validating methods to prevent overinterpretation of low
38 frequency variants in a sample.

39 Keywords: SIV; multiplex PCR; deep sequencing; SNV detection

40

41 Introduction:

42 Characterizing the sequence diversity of RNA virus populations is an essential component
43 of studying viral pathogenesis and transmission in individuals (Vignuzzi et al. 2006, Poirier and
44 Vignuzzi 2017). This sequence data can be used to identify antiviral drug resistance mutations
45 (Dudley et al. 2014), understand how viruses evolve (Lessler et al. 2016, Henn et al. 2012), and
46 track virus transmission during epidemics (Hadfield et al 2018), such as the Ebola virus outbreak
47 in West Africa in 2014, the Zika virus outbreak in Brazil in 2015 (Quick et al. 2017, Carroll et al.
48 2015), and the current SARS-CoV-2 outbreak (Fauver et al 2020, CDC 2020).

49 The accumulation of mutations in RNA viruses can impact their pathogenesis (Sanjuán
50 and Domingo-Calap 2016). While many mutations can be deleterious or neutral, some are
51 beneficial for virus proliferation, survival, or transmission (Henn et al. 2012, Zanini et al. 2017).
52 Naturally elicited host immune responses that fail to eliminate replicating viruses select for
53 variants that avoid immune detection (Henn et al. 2012). Drug resistance mutations can also
54 accumulate when antiretroviral therapy does not fully suppress virus replication (Bangsberg et al
55 2007). Accurate detection of these variants in RNA virus populations can help determine whether
56 therapeutic interventions eliminate or exacerbate mutations from the replicating virus population.

57 Sequencing RNA viruses requires the generation of viral cDNA, followed by amplification
58 of either long (greater than 1000bp) or short (less than 400bp) DNA segments. Long amplicons
59 are used to study distantly linked nucleotides on the same virus templates using Pacific
60 Biosciences (Brese et al. 2018) or Oxford Nanopore instruments (Tyler et al. 2018,
61 Oikonomopoulos et al. 2016), In contrast, Illumina technology can generate sequence data from
62 shorter viral segments with higher throughput, better fidelity, and improved efficiency (Adey et al.
63 2010). While each approach has advantages and disadvantages, the desire to acquire sequence
64 data with newer assays often trumps taking the time to perform experiments required to validate
65 their sensitivity and reproducibility.

66 Our goal was to implement a multiplex PCR approach, similar to those used for Ebola
67 (Arias et al 2016), ZIKV (Quick et al 2017), and SARS-CoV-2 (Fauver et al 2020), to improve the
68 reproducibility and sensitivity of sequencing SIV derived from plasma with low virus titers or cell-
69 associated vRNA isolated from different tissues. SIV dynamics are frequently studied in
70 nonhuman primates (Harris et al. 2013, Gambhira et al. 2014, Hassounah et al. 2014), but the
71 samples collected from animals with interesting biological phenotypes often have a low virus titer.
72 With an ongoing emphasis on understanding the dynamics of SIV replication in nonhuman
73 primates (Kumar et al. 2016), we aimed to determine if the multiplex method could be applied to
74 SIV to improve the characterization of virus populations with improved sensitivity and
75 reproducibility.

76 We developed a multiplex PCR approach to amplify and sequence SIV. To validate this
77 method, we sequenced different numbers of vRNA and viral cDNA templates of clonal
78 SIVmac239, as well as variable ratios of two clonal SIV strains differing at 24 nucleotide positions.
79 We found improved sensitivity and reproducibility of variant calling when normalizing to the
80 number of viral cDNA templates added to the reaction when compared to the number of vRNA
81 templates added to the reaction. By validating the SIV multiplex sequencing method here, we
82 identify the strengths and limitations of this method, which are essential for defining the usability
83 of any new technique.

84

85 Results:

86 **Design of a multiplex PCR assay for SIV.**

87 Candidate multiplex primers for SIV were designed in Primal Scheme, a tool developed
88 by Quick et al. (2017). Each primer set was tested individually and then pooled such that the
89 amplicon products would not overlap with each other (Figure 1A, Table 1). Primer pools were
90 tested to verify that individual primer pairs would generate amplicons spanning the entire viral

91 genome when combined. Final primer pair concentrations and corresponding sequences can be
92 found in Table 1.

93 We first isolated vRNA from a stock of clonal SIVmac239. For Method 1, we quantified the
94 vRNA stock and then diluted it to 10^6 vRNA templates per reaction. Serial dilutions of quantified
95 vRNA were converted to viral cDNA by reverse-transcription. The multiplex PCR was performed
96 on the viral cDNA. For Method 2, we prepared total cDNA from 10^7 vRNA templates of each stock.
97 We then quantified the viral cDNA with a qPCR reaction specific for *gag*. The quantified viral
98 cDNA was diluted to 10^6 cDNA *gag* copies per reaction and the multiplex PCR was then
99 performed. After multiplex PCR for either Method 1 or 2, 75ng of each pool of PCR products were
100 combined into a single tube to generate a 150ng DNA pool containing all the generated PCR
101 amplicons. This amplicon library was then tagged using an Illumina TruSeq kit, and sequenced
102 on an Illumina MiSeq.

103

104 **Detection of false positives in clonal SIVmac239**

105 We first sequenced clonal SIVmac239 to determine the frequency of false positives when
106 using either Method 1 or 2. We used serially diluted 100% SIVmac239 vRNA or viral cDNA for
107 this part of the project. For each replicate using Method 1, new cDNA was prepared and then
108 multiplex PCR and sequencing were performed. These experiments were performed in triplicate.
109 For each replicate using Method 2, the same prepared cDNA was used for all of the multiplex
110 PCR reactions. These experiments were performed in duplicate.

111 FASTQ sequences were examined using a modified version of a custom pipeline
112 previously used to analyze multiplex PCR ZIKV sequences (Dudley et al. 2017,
113 https://github.com/SLO-Lab/SIV_MultiplexPCR). Using this tool, we randomly subsampled up to
114 2000 reads per amplicon across each data set and mapped them to SIVmac239 (Accession:
115 M33262), as described in the Materials and Methods. Amplification of each PCR product does
116 not occur equally, so by subsampling up to 2000 reads, we could attempt to informatically

117 normalize the depth of coverage, while not oversampling any one single amplicon. VarScan
118 (<https://sourceforge.net/projects/varscan/>) was then used to identify nucleotides present in the
119 virus population that were different from the reference at a frequency of 1% or greater and had a
120 depth of coverage of at least 1800 nucleotides, or 90% of our maximum subsampled depth.
121 SNPeff (Cingolani et al. 2012) was used to annotate variants and their effect on each coding
122 sequence. Any single nucleotide variant (SNV) present at a frequency of 1% or greater and with
123 a depth of coverage of at least 1800 nucleotides was categorized as a false positive for our
124 analysis. These thresholds are more conservative than the 3% cutoff and 400x coverage required
125 by Grubaugh et. al (2019).

126 We began by assessing false positives present in sequences generated by Method 1. The
127 average rate of false positives in a single replicate was related to the number of input templates,
128 with samples containing 10^3 input copies having a higher average rate of false positives at $1.13 \times$
129 10^{-2} false positives per nucleotide, and samples containing 10^6 input copies having a lower
130 average rate of 2.6×10^{-3} false positives per nucleotide (Figure 2A, left panel, closed circles, $p <$
131 0.0001 , Tukey's multiple comparisons test). We then determined whether the rate of false
132 positives declined when considering two or more replicates. We found there was not a significant
133 copy-dependent decrease in false positives when we used two replicates compared to a single
134 replicate (Figure 2A, left panel, open circles, $p = 0.83$, Tukey's multiple comparisons test).

135 We investigated the individual nucleotide positions where we detected false positives in
136 multiple replicates when using Method 1 (Figure 2B). We found 11 positions with false positives
137 at all input copy numbers, with 4 being insertions at nucleotide positions 1254, 1480, 5428, and
138 7396, and 9 being substitutions at nucleotide positions 6181, 6186, 6188, 6190, 6192, 6201, 6205,
139 6207, and 6713. We found the median false positive frequency did not depend on the number of
140 input copies ($p = 0.92$, Kruskal-Wallis) (Figure 2B). Each of the insertions occurred in a poly-A
141 region containing 6 consecutive adenines. Although the chemistry of Illumina sequencing does
142 not lead to the same errors in homopolymers that are notorious in other sequencing platforms

143 (Bently et al. 2008), there can still be PCR-based errors in homopolymeric regions (McInerney et
144 al. 2014, Potapov and Ong 2017). All substitutions, aside from position 6713, were present within
145 a stretch of 27 nucleotides that are adjacent to a primer binding site. These SNVs are contained
146 within an overlap region between Amplicons 20 and 21. Notably, these variants were present in
147 Amplicon 21, but not Amplicon 20, suggesting that Amplicon 21 may be more prone to the
148 incorporation of PCR-based substitutions than Amplicon 20. While unfortunate, inaccuracies in
149 variant reporting is not an uncommon phenomenon at the ends of amplicons and has been
150 reported previously (Schirmer et al 2016, McCoy et al 2014). We also observed that when using
151 a different analysis pipeline that does not normalize coverage across the genome through
152 subsampling, these variants were not reported in the vcf file, highlighting the importance of
153 validating the analysis methods prior to calling variants as true variants. However, we felt that the
154 benefit of standardizing variant calling with normalized coverage across the genome outweighed
155 the complexity associated with variabilities related to relative oversampling of individual
156 amplicons.

157 We then used the same metrics to identify false positives using Method 2 (Figure 2C).
158 Similar to Method 1, the average number of false positives per nucleotide in at least one replicate
159 was related to the number of input templates, with 10^3 input cDNA templates having an average
160 of 1.05×10^{-2} false positives per nucleotide and 10^6 input cDNA templates having an average of
161 1.52×10^{-3} false positives per nucleotide (Figure 2A, right panel, closed circles, $p < 0.001$, Tukey's
162 multiple comparisons test). When only including false positives detected in at least two replicates,
163 there was no difference in the rate of false positives between 10^3 and 10^6 cDNA templates (Figure
164 2A, right panel).

165 We identified 9 nucleotide positions with false positives in at least one replicate of all input
166 template levels using Method 2 (Figure 2C). All of the false positives detected by Method 2 were
167 also detected by Method 1. Since these false positives are present in nearly every sample and
168 this is a clonal virus stock (Figures 2B and 2C), it is likely an artifact of the method rather than

169 true variants, highlighting the importance of validating novel methods with virus stocks of known
170 composition. Additionally, it is important to understand the effects of nucleotide sequence and
171 primer binding sites on false positive detection, as primer slippage may be a confounding factor.

172 To help determine if the rate of false positives was related to coverage depth, we
173 calculated the frequency of nucleotide sites that had sufficient coverage (a nucleotide depth of at
174 least 1800) for our cDNA and vRNA data sets. There was no significant difference between the
175 percentage of bases with at least 1800x coverage using Method 1 or 2 (Method 1 mean = 76.27%
176 nucleotides over 1800, Method 2 mean = 75.03% nucleotides over 1800; $p = 0.95$, Mann-Whitney,
177 data not shown), indicating that the differences in false positive frequency are more likely a result
178 of starting template than coverage alone.

179

180 **Detection of genome-wide variants using multiplex SIV sequencing**

181 We then examined the sensitivity and reproducibility of detecting individual SNVs in SIV
182 by Methods 1 and 2. We used two stock viruses, SIVmac239 and SIVmac239-24x, that differed
183 at 24 nucleotides throughout the entire viral coding sequence (Figure 1b). Viral RNA was isolated
184 from these two stocks and quantified with a *gag* qPCR assay. We proceeded with Method 1 by
185 mixing the two stocks of vRNA to a total number of 10^6 copies at the following SIVmac239 to
186 SIVmac239-24x ratios: 100:0, 95:5, 90:10, 75:25, 50:50, and 0:100 (Figure 3A). Each mixture of
187 vRNA was serially diluted to 10^5 , 10^4 , and 10^3 templates per 11 μ l. We also tested Method 2 by
188 first preparing viral cDNA from 10^7 vRNA templates of each of the two stocks, quantifying viral
189 cDNA, and then mixing the cDNA templates to a total of 10^6 templates in the same ratios as the
190 vRNA templates were mixed (Figure 3B). The same quantified vRNA or viral cDNA mixtures were
191 used for the entire experiment.

192 The remaining multiplex PCR procedures were performed for the different numbers of
193 input templates and for each of the individual ratios. PCR products were tagged, and sequencing
194 was performed on the Illumina MiSeq. FASTQ reads were mapped to SIVmac239 and the

195 frequencies of each individual SNV relative to SIVmac239 were determined as described for the
196 clonal SIVmac239 data.

197 We compared the observed to the expected variant frequencies for all 24 positions in the
198 genome for both Methods 1 and 2. We generated a linear regression for each number of input
199 templates (Figure 4A) to determine if the relationship between the expected and observed SNV
200 frequency was the same. We did not find a significant difference when we compared the slopes
201 for all four linear regression lines with either Method 1 ($p = 0.069$, Figure 4A) or Method 2 ($p =$
202 0.185 , Figure 4B). Notably, all of these data sets had an SNV present at position 9110 (Figure 4A
203 and 4B, open circles) that was consistently detected inaccurately. While there did appear to be a
204 slight increase in observed variant frequency when compared to expected variant frequency, site
205 9110 was a clear outlier in the data sets (Figure 4C).

206 To further understand how the number of templates and the type of quantified starting
207 material affects the reproducibility of the detected SNV frequencies, we compared the observed
208 frequencies of each of the 24 individual SNVs across all the data sets. We found that when using
209 Method 2, there was less variability in variant frequencies across the number of input templates
210 when compared to using Method 1. (Figure 5A-F). This observation is consistent with data
211 indicating that the process of reverse transcription is inefficient and variable (Bustin et al. 2015),
212 such that when 10^3 vRNA input templates are used in the assay, it is unlikely that there are actually
213 10^3 viral cDNA templates available for subsequent PCR. For both input types, it was not surprising
214 that as the number of templates increased, the SNV frequencies tended to be more consistent
215 and reproducible across the genome.

216

217 Discussion:

218 The goal of this study was to adapt a multiplex PCR and sequencing approach (Quick et
219 al. 2017) to sequence SIV from low quality starting material. This would include occasions where
220 SIV is present at low titer or as partially degraded vRNA. Recognizing that different sequencing

221 methods have their limitations, we set out to validate this approach in a series of assays described
222 in this study.

223 We first sequenced clonal SIVmac239 to determine the false positive rate using both
224 Methods 1 and 2. We found 4 nucleotide sites (1254, 2480, 5428, and 7396) in homopolymeric
225 regions with consecutive adenines where false positive indels were detected. We predict that
226 these insertions were introduced during the PCR step. We also found 8 individual false positives
227 in the stock that were attributed to substitutions consistently present in the same 27-nucleotide
228 region of Amplicon 21, but not in the adjacent Amplicon 20. We hypothesize that these
229 substitutions are specific to the generation of Amplicon 21 and the analysis pipeline, rather than
230 actually being real substitutions.

231 We also found that sequencing replicates reduced the detection of false positives,
232 particularly when there are low numbers of input templates, consistent with previous results
233 (Grubaugh et al 2019). While we realize that there are not always enough resources available to
234 sequence a sample in duplicate, our data highlights that caution should be taken when interpreting
235 data from a single assay of a sample with low virus titer. Importantly, the process of validating a
236 method with a known clonal virus stock is key to distinguishing between false positives,
237 sequencing error, and true variants. Without doing the validation assays in this study, it would be
238 impossible to know the benefits and technical limitations of using the multiplex PCR approach to
239 sequence virus isolated from animals infected with SIVmac239. Detecting these method-
240 dependent systematic errors by characterizing false positives in a clonal stock is important so that
241 investigators using this method can perform the assay with knowledge of which variants are real
242 and which are technical artifacts.

243 By mixing SIVmac239 and SIVmac239-24x, we detected variants at a frequency of 5%
244 with as few as 1000 input copies. We opted for this conservative threshold because we already
245 knew that there were some false positives detected when a threshold of 1% was used (Figure 2),
246 and the most relevant variants accumulate over time to a higher frequency. Thus, detection of

247 variants at a frequency of <5% was less critical for broad analyses of SIV population diversity.
248 Future studies that require more sensitive variant detection could address whether variants
249 present between 1% and 5% can be accurately detected.

250 For these mixing studies, we chose two viruses with variants scattered throughout the
251 genome, with at least one variant present in each gene. This let us determine whether we could
252 effectively detect variants throughout the genome and across a large number of PCR amplicons
253 generated by either Method 1 or 2. We were surprised to find it difficult to interpret the SNV
254 frequency at position 9110. This site lies in a region dense with adenines and guanines which
255 may contribute to some inconsistencies as a result of PCR slippage or misincorporation of
256 nucleotides during PCR amplification (Pfeiffer et al 2018). In addition, the forward primer for
257 Amplicon 33 is one nucleotide different from its complementary sequence in SIVmac239-24x due
258 to the modified nucleotide 9110 present in the SIVmac239-24x sequence. While we did trim
259 primers computationally, this would not prevent PCR error from occurring. As a result, some
260 SIVmac239-24x templates may not be amplified as efficiently because of a single nucleotide
261 difference, which may also lead to amplicon dropout and skewed results.

262 Throughout our study, we compared the results obtained using Methods 1 and 2, which
263 used quantified vRNA and quantified viral cDNA, respectively. Reverse transcription is inefficient
264 (Bustin et al 2015), so we wanted to determine if there were fewer false positives and more
265 consistent detection of SNVs when quantified cDNA was used as the starting material rather than
266 vRNA. We found the observed SNV frequencies were more similar to expected frequencies when
267 quantified cDNA was used as a starting template (Figure 4b) and, not surprisingly, when
268 increased numbers of vRNA or cDNA templates were used. Even though our quantification of
269 viral cDNA was based only on the copies of *gag*, we found that using quantified viral cDNA as the
270 input improved the reproducibility of variant detection when we mixed two clonal virus inocula at
271 predefined ratios, even when using only 10^3 quantified templates. This observation further raises
272 concerns that using quantified vRNA as starting material gives an overestimation of the number

273 of vRNA templates that are actually converted to cDNA and amplified to yield the reported
274 sequence data.

275 Overall, we found that the multiplex PCR approach could be successfully used to generate
276 genome wide sequences of SIV, but our results strongly imply that any new sequencing and
277 analysis methods be validated before using them widely to characterize variant frequency in a
278 virus population. While it was possible to generate sequence data from 10^3 vRNA templates, the
279 use of quantified cDNA was more consistent. Further, although this method could be used to
280 successfully detect SNVs across the genome, we found there were key features in the viral
281 genome that affected the accuracy of the multiplex PCR approach. Thus, while the multiplex PCR
282 method has many advantages for deep sequencing virus populations, validation experiments and
283 visualization of the output alignments are essential for correct data reporting, as expected for any
284 sequencing approach.

285

286 Conclusions:

287 Our initial goal of this study was to generate a sequencing approach that was able to
288 characterize viral population diversity with low input templates. Multiplex PCR has been used to
289 accurately sequence other viruses, including Zika (Quick et al. 2017), Dengue, and Chikungunya
290 (Kafetzopoulou et al. 2018), at titers between 10^3 and 10^6 vRNA copies per mL, and most recently
291 with SARS-CoV-2 (CDC et al 2020, Fauver et al 2020) and we were hoping this would extend to
292 SIV. However, many publications fail to state the viral input titer when describing their sequencing
293 methods. We learned that increasing numbers of input SIV templates and utilization of quantified
294 cDNA as a starting material improved reproducibility of variant calling. Further, our data suggests
295 that the multiplex PCR and sequencing approach may not be as sensitive at low numbers of input
296 templates for SIV, when compared to other using low numbers of templates for other viruses.
297 Most importantly, our study demonstrates the need to validate new sequencing approaches
298 because the same method may not be viable for sequencing all viruses with the same sensitivity

299 and reproducibility. We now understand the limitations of the assay so that experiments can be
300 designed to maximize the likelihood of success and minimize the overinterpretation of data.

301

302 Methods:

303 **Primer design:** Primers were designed using Primal Scheme, as previously described by Quick
304 et al, 2017. FASTA files of SIVmac239 and three consensus sequences of virus populations
305 isolated from animals infected with SIVmac239 were used as the foundation for the Primal
306 Scheme tool. 37 primer pairs (Table 1) were generated to span the entire SIV genome. The
307 lengths of the resulting amplicons ranged from 285bp to 397bp, with an average length of 351bp.
308 The number of overlapping nucleotides for each amplicon ranged from 40bp to 149bp, with an
309 average length of 100bp. Primer pairs were split into two pools to ensure that the amplicons
310 generated within each pool would not overlap. Primer sequences, pools, and concentrations can
311 be found in Table 1. Final concentration of Pool 1 was 35uM and Pool 2 was 24uM.

312

313 **Isolation of vRNA for sequencing:** SIVmac239 and SIVmac239-24x vRNA were isolated from
314 clonal virus stocks. Briefly, 1ml of each virus stock was centrifuged at 13,000rpm for 30 seconds
315 to pellet any cells that were present. The supernatant was transferred to a 1.5mL Eppendorf tube
316 and spun at 13,000rpm for 1 hour at 4C to concentrate virus particles. After spinning the sample,
317 we removed all the supernatant, except 200ul of liquid, so as not to disturb the viral pellet. The
318 vRNA was then extracted using the Qiagen MinElute vRNA extraction kit, according to
319 manufacturer's protocols (Qiagen). Prior to elution, 25uL of Buffer AVE was added directly to the
320 MinElute Column membrane and incubated for 5 minutes.

321

322 **Preparation of viral cDNAs:** The vRNA isolated from the SIVmac239 and SIVmac239-24x virus
323 stocks were each diluted to 10^6 copies/11ul in nuclease-free water. They were mixed at
324 SIVmac239:SIVmac239-24x ratios of 100:0, 95:5, 90:10, 75:25, 50:50, and 0:100. These mixtures

325 were diluted 1:10 in nuclease-free water to generate vRNA template concentration dilution series
326 of 10^6 , 10^5 , 10^4 , and 10^3 templates per 11uL. From each mixture, we used 11ul of vRNA and
327 performed cDNA synthesis using SuperScript IV Reverse Transcriptase (Invitrogen), according
328 to manufacturer's protocol. For experiments where quantified viral cDNA was used as starting
329 material, approximately 10^7 viral templates were used for cDNA synthesis using SuperScript IV
330 Reverse Transcriptase (Invitrogen), according to manufacturer's protocol. Viral cDNA and vRNA
331 was then quantified using a *gag* qPCR assay as previously described (Cline et al. 2005).

332

333 ***Multiplex PCR reactions:*** Each tube of viral cDNA generated from the virus stocks or biological
334 samples was split equally, such that 10uL of viral cDNA was PCR amplified with the two separate
335 primer pools. Amplification was performed with the Q5 polymerase and the following reaction
336 conditions: 98°C for 30 seconds, 35 cycles of 95°C for 15 seconds and 65°C for 5 minutes, and
337 then cooled to 4°C. Products were verified using a 1% agarose gel and were quantified using the
338 Qiagen High Sensitivity DNA kit (Thermo Fisher).

339

340 ***Library Preparation and Sequencing:*** After the two amplicon pools were generated, 75ng of
341 each pool was mixed to generate a total of 150ng DNA. This pool of PCR products was tagged
342 with the Illumina TruSeq Nano HT kit, according to the manufacturer's protocol (Illumina).
343 Following tagging and purifying, the libraries were quantified using the Qiagen High Sensitivity
344 DNA kit. The quality of each library was characterized with a High Sensitivity DNA kit (Agilent) on
345 an Agilent Bioanalyzer. If unligated adapter dimers were detected at 140bp, an additional bead
346 clean up step was performed. The average tagged library size was approximately 503bp (range
347 491 to 512). Tagged libraries were pooled at equimolar concentrations and diluted so that the
348 final concentration of DNA molecules per run was 10pM. This diluted pool and 10pM PhiX were
349 denatured with 0.2N sodium hydroxide for 5 minutes at room temperature. Denatured PhiX was

350 then added to the pool at a final frequency of 10 percent. Each pool was loaded at 10pM
351 concentration onto a 500-cycle v2 MiSeq cartridge and sequenced.

352

353 **Data analysis:** FASTQ reads were demultiplexed and then processed using a modified pipeline
354 from our lab, called the Zequencer. All scripts used can be found on our github repository,
355 https://github.com/SLO-Lab/SIV_MultiplexPCR. Briefly, reads were trimmed, merged, and
356 normalized using bbtools (<https://jgi.doe.gov/data-and-tools/bbtools/>) and Seqtk
357 (<https://github.com/lh3/seqtk>). A FASTA file was generated that contained the nucleotide
358 reference sequences for all 37 amplicons, as they would exist in SIVmac239 (Accession:
359 M33262). Up to 2000 merged reads that mapped at low sensitivity to each of the 37 reference
360 amplicons were extracted from the data set. These reads were then aligned to SIVmac239 using
361 NovoAlign (<http://www.novocraft.com/products/novoalign/>). A pileup file was generated from the
362 BAM alignment. Variants with a frequency of 1% or higher were called by VarScan
363 (<https://sourceforge.net/projects/varscan/>) and annotated by SNPeff (Cingolani et al. 2012). VCF
364 files were processed and analyzed in R(v3.6.1). Variants with a sample depth less than 1800
365 were discarded to reduce bias as a result of poor sample depth. Position 9609 codes for a stop
366 codon in the *nef* protein in the M33262 Genbank reference for SIVmac239, but our stock virus is
367 SIVmac239-nef-open, which has a T to G transversion at this position, converting the stop codon
368 (TAA) to a glutamate (GAA) amino acid.

369

370 List of abbreviations:

371 HIV – Human immunodeficiency virus; SIV – Simian immunodeficiency virus; ZIKV – Zika virus;
372 vRNA – Viral RNA; SNV – single nucleotide variant.

373

374 Declarations and Acknowledgements:

375 Ethics approval: Not applicable.

376 Consent to publication: not applicable.

377 Availability of data and materials: code used to generate data can be found on the lab's GitHub
378 page (see Data analysis section).

379 Competing interests: All authors read and approved the manuscript and declare no competing
380 interests.

381 Funding: Research reported in this publication was supported by the Office Of The Director,
382 National Institutes of Health under Award Number P51OD011106 to the Wisconsin National
383 Primate Research Center, University of Wisconsin-Madison. This work was also supported by
384 an Australian National Health and Medical Research Council (NHMRC) grants 1080001 and
385 1052979 (to MPD) and NHMRC Career Development Fellowship 1067590 (to VV).

386 Contributions: The content is solely the responsibility of the authors and does not necessarily
387 represent the official views of the National Institutes of Health. RVM and SLO wrote the
388 manuscript. Primer design and validation by NF and MSS. Data analysis was aided by MPD,
389 TS, and VV, and conducted by RVM and SLO.

390 Acknowledgements: We would like to thank the Virology Services unit at the Wisconsin National
391 Primate Research Center for quantifying SIV vRNA and cDNA, which allowed us to dilute our
392 samples to predetermined levels. We would also like to thank Josh Quick and Nick Loman for
393 helping us with primer design.

394

395 References:

396

- 397 1. Adey, A., Morrison, H. G., Asan, Xun, X., Kitzman, J. O., Turner, E. H. et al. (2010). Rapid,
398 low-input, low-bias construction of shotgun fragment libraries by high-density in vitro
399 transposition. *Genome Biol*, 11(12), R119.
- 400 2. Bustin, S., Dhillon, H. S., Kirvell, S., Greenwood, C., Parker, M., Shipley, G. L. et al. (2015).
401 Variability of the reverse transcription step: practical implications. *Clin Chem*, 61(1), 202-
402 212.
- 403 3. Carroll, M. W., Matthews, D. A., Hiscox, J. A., Elmore, M. J., Pollakis, G., Rambaut, A. et
404 al. (2015). Temporal and spatial analysis of the 2014-2015 Ebola virus outbreak in West
405 Africa. *Nature*, 524(7563), 97-101.

- 406 4. Cingolani, P., Platts, A., Wang, L. L., Coon, M., Nguyen, T., Wang, L. et al. (2012). A
407 program for annotating and predicting the effects of single nucleotide polymorphisms,
408 SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly*
409 (*Austin*), 6(2), 80-92.
- 410 5. Cline, A. N., Bess, J. W., Piatak, M., & Lifson, J. D. (2005). Highly sensitive SIV plasma
411 viral load assay: practical considerations, realistic performance expectations, and
412 application to reverse engineering of vaccines for AIDS. *J Med Primatol*, 34(5-6), 303-312.
- 413 6. Dudley, D. M., Bailey, A. L., Mehta, S. H., Hughes, A. L., Kirk, G. D., Westergaard, R. P.
414 et al. (2014). Cross-clade simultaneous HIV drug resistance genotyping for reverse
415 transcriptase, protease, and integrase inhibitor mutations by Illumina MiSeq.
416 *Retrovirology*, 11(1), 122.
- 417 7. Dudley, D. M., Newman, C. M., Lalli, J., Stewart, L. M., Koenig, M. R., Weiler, A. M. et al.
418 (2017). Infection via mosquito bite alters Zika virus tissue tropism and replication kinetics
419 in rhesus macaques. *Nat Commun*, 8(1), 2096.
- 420 8. Gambhira, R., Keele, B. F., Schell, J. B., Hunter, M. J., Dufour, J. P., Montefiori, D. C. et
421 al. (2014). Transmitted/founder simian immunodeficiency virus envelope sequences in
422 vesicular stomatitis and Semliki forest virus vector immunized rhesus macaques. *PLoS*
423 *One*, 9(10), e109678.
- 424 9. Grubaugh, N. D., Gangavarapu, K., Quick, J., Matteson, N. L., De Jesus, J. G., Main, B.
425 J. et al. (2019). An amplicon-based sequencing framework for accurately measuring
426 intrahost virus diversity using PrimalSeq and iVar. *Genome Biol*, 20(1), 8.
- 427 10. Harris, M., Burns, C. M., Becker, E. A., Braasch, A. T., Gostick, E., Johnson, R. C. et al.
428 (2013). Acute-phase CD8 T cell responses that select for escape variants are needed to
429 control live attenuated simian immunodeficiency virus. *J Virol*, 87(16), 9353-9364.
- 430 11. Hassounah, S. A., Mesplède, T., Quashie, P. K., Oliveira, M., Sandstrom, P. A., &
431 Wainberg, M. A. (2014). Effect of HIV-1 integrase resistance mutations when introduced
432 into SIVmac239 on susceptibility to integrase strand transfer inhibitors. *J Virol*, 88(17),
433 9683-9692.
- 434 12. Henn, M. R., Boutwell, C. L., Charlebois, P., Lennon, N. J., Power, K. A., Macalalad, A. R.
435 et al. (2012). Whole genome deep sequencing of HIV-1 reveals the impact of early minor
436 variants upon immune recognition during acute infection. *PLoS Pathog*, 8(3), e1002529.
- 437 13. Kafetzopoulou, L. E., Efthymiadis, K., Lewandowski, K., Crook, A., Carter, D., Osborne, J.
438 et al. (2018). Assessment of metagenomic Nanopore and Illumina sequencing for
439 recovering whole genome sequences of chikungunya and dengue viruses directly from
440 clinical samples. *Euro Surveill*, 23(50).
- 441 14. Kumar, N., Chahroudi, A., & Silvestri, G. (2016). Animal models to achieve an HIV cure.
442 *Curr Opin HIV AIDS*, 11(4), 432-441.
- 443 15. Lessler, J., Chaisson, L. H., Kucirka, L. M., Bi, Q., Grantz, K., Salje, H. et al. (2016).
444 Assessing the global threat from Zika virus. *Science*, 353(6300), aaf8160.
- 445 16. McInerney, P., Adams, P., & Hadi, M. Z. (2014). Error Rate Comparison during
446 Polymerase Chain Reaction by DNA Polymerase. *Mol Biol Int*, 2014, 287430.
- 447 17. Oikonomopoulos, S., Wang, Y. C., Djambazian, H., Badescu, D., & Ragoussis, J. (2016).
448 Benchmarking of the Oxford Nanopore MinION sequencing for quantitative and qualitative
449 assessment of cDNA populations. *Sci Rep*, 6, 31602.

- 450 18. Poirier, E. Z., & Vignuzzi, M. (2017). Virus population dynamics during infection. *Current*
451 *opinion in virology*, 23, 82-87.
- 452 19. Potapov, V., & Ong, J. L. (2017). Examining Sources of Error in PCR by Single-Molecule
453 Sequencing. *PLoS One*, 12(1), e0169774.
- 454 20. Quick, J., Grubaugh, N. D., Pullan, S. T., Claro, I. M., Smith, A. D., Gangavarapu, K. et al.
455 (2017). Multiplex PCR method for MinION and Illumina sequencing of Zika and other virus
456 genomes directly from clinical samples. *Nat Protoc*, 12(6), 1261-1276.
- 457 21. Quick, J., Loman, N. J., Duraffour, S., Simpson, J. T., Severi, E., Cowley, L. et al. (2016).
458 Real-time, portable genome sequencing for Ebola surveillance. *Nature*, 530(7589), 228-
459 232.
- 460 22. Sanjuán, R., & Domingo-Calap, P. (2016). Mechanisms of viral mutation. *Cell Mol Life Sci*,
461 73(23), 4433-4448.
- 462 23. Tyler, A. D., Mataseje, L., Urfano, C. J., Schmidt, L., Antonation, K. S., Mulvey, M. R. et
463 al. (2018). Evaluation of Oxford Nanopore's MinION Sequencing Device for Microbial
464 Whole Genome Sequencing Applications. *Sci Rep*, 8(1), 10931.
- 465 24. Vignuzzi, M., Stone, J. K., Arnold, J. J., Cameron, C. E., & Andino, R. (2006).
466 Quasispecies diversity determines pathogenesis through cooperative interactions in a
467 viral population. *Nature*, 439(7074), 344-348.
- 468 25. Zanini, F., Puller, V., Brodin, J., Albert, J., & Neher, R. A. (2017). *In vivo* mutation
469 rates and the landscape of fitness costs of HIV-1. *Virus Evol*, 3(1), vex003.
- 470 26. Schirmer, M., D'Amore, R., Ijaz, U.Z. et al. Illumina error profiles: resolving fine-scale
471 variation in metagenomic sequencing data. *BMC Bioinformatics* 17, 125 (2016).
472 <https://doi.org/10.1186/s12859-016-0976-y>
- 473 27. McCoy RC, Taylor RW, Blauwkamp TA, Kelley JL, Kertesz M, et al. (2014) Illumina
474 TruSeq Synthetic Long-Reads Empower *De Novo* Assembly and Resolve Complex,
475 Highly-Repetitive Transposable Elements. PLOS ONE 9(9):
476 e106689. <https://doi.org/10.1371/journal.pone.0106689>
- 477 28. Bangsberg, D. R., Kroetz, D. L., & Deeks, S. G. (2007). Adherence-resistance
478 relationships to combination HIV antiretroviral therapy. *Current HIV/AIDS Reports*, 4(2),
479 65.
- 480 29. CDC, C. O. V. I. D.-1. R. T., Jorden, M. A., Rudman, S. L., Villarino, E., Hoferka, S.,
481 Patel, M. T. et al. (2020). Evidence for Limited Early Spread of COVID-19 Within the
482 United States, January-February 2020. *MMWR Morb Mortal Wkly Rep*, 69(22), 680-684
- 483 30. Fauver, J. R., Petrone, M. E., Hodcroft, E. B., Shioda, K., Ehrlich, H. Y., Watts, A. G. et
484 al. (2020). Coast-to-coast spread of SARS-CoV-2 during the early epidemic in the United
485 States. *Cell*.
- 486 31. Hadfield, J., Megill, C., Bell, S. M., Huddleston, J., Potter, B., Callender, C. et al. (2018).
487 Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics*, 34(23), 4121-4123.
- 488 32. Arias, A., Watson, S. J., Asogun, D., Tobin, E. A., Lu, J., Phan, M. V. T. et al. (2016).
489 Rapid outbreak sequencing of Ebola virus in Sierra Leone identifies transmission chains
490 linked to sporadic cases. *Virus Evol*, 2(1), vew016.
- 491 33. Pfeiffer, F., Gröber, C., Blank, M., Händler, K., Beyer, M., Schultze, J. L. et al. (2018).
492 Systematic evaluation of error rates and causes in short samples in next-generation
493 sequencing. *Sci Rep*, 8(1), 10950.
- 494

495

496 Figure Legends

497

498 Figure 1a) SIV multiplex primer scheme. Two non-overlapping pools corresponding to even (red)
499 and odd (blue) primer sets were designed using Primal Scheme to generate small amplicons
500 spanning the entire SIVmac239 genome. Primer pairs were pooled at varying concentrations
501 described in Table 1. 1b) Locations of SNPs in SIVmac239-24x. SNPs denoted with a red X.
502 SNPs are present across the entire SIVmac239 genome and are present in all genes.

503

504 Figure 2a) Number of false positives detected per nucleotide with coverage of at least 1800 and
505 a variant frequency of at least 1% in at least one replicate per input copy (closed circle) or at least
506 two replicates per input copy (open circle) for our Method 1 (vRNA) (left) and Method 2 (cDNA)
507 (right) data sets. Lines represent median +/- 95% confidence interval. No significant differences
508 were identified between data sets by Kruskal-Wallis tests. 2b) False positive variant frequency of
509 variants identified in all input templates for Method 1 data sets. 2c) False positive variant
510 frequency of variants identified in all input templates for Method 2 data sets. cDNA input template
511 numbers denoted by colors. Lines represent mean and standard deviation for each variant's
512 replicate. All variants shown are present at a frequency of 1% or greater, have a nucleotide depth
513 of at least 1800, and are detected in at least two samples.

514

515 Figure 3) Schematic of experimental design. (a) Method 1: viral RNA was isolated from original
516 stock and quantified via qRT-PCR. SIVmac239 and SIVmac239-24x were diluted to 10^6
517 copies/reaction and mixed at the following SIVmac239:SIVmac239-24x ratios: 100:0, 95:5, 90:10,
518 75:25, 40:60, and 0:100. Serial dilutions were preformed to 10^5 , 10^4 , and 10^3 copies per reaction.
519 Viral cDNA was generated from viral RNA mixes, with one cDNA reaction per vRNA mix. (b)
520 Method 2: viral RNA was isolated and approximately 10^7 viral RNA copies were added to each
521 cDNA synthesis reaction. Viral cDNA copies were quantified using qRT-PCR and each was
522 diluted to 10^6 copies per reaction. SIVmac239:SIVmac239-24x mixes were generated at the
523 following ratios: 100:0, 95:5, 90:10, 75:25, 50:50, and 0:100. cDNA mixes were then serially
524 diluted to 10^5 , 10^4 , and 10^3 copies per reaction. (c) cDNA was used for multiplex PCR. PCR
525 products were then combined at equimolar ratios and library prepped according to TruSeq Library
526 Preparation documentation (Illumina). Libraries were quantified, pooled, and sequenced using a
527 2x250 v2 MiSeq cartridge.

528

529 Figure 4a) Observed versus expected variant frequencies identified in the Method 1 (vRNA) mixed
530 data sets. Observed variant frequency indicates percent SIVmac239-24x identified. Error bars
531 indicate standard deviation for each replicate. Linear regressions colored by vRNA templates.
532 Open circles indicate SNV 9110. 4b) Observed versus expected variant frequencies identified in
533 Method 2 (cDNA) mixed data sets. Observed variant frequency indicates percent SIVmac239-24x
534 identified. Error bars indicate standard deviation for each replicate. Linear regressions colored by
535 cDNA templates. Open circles indicate SNV 9110. No significant difference is observed between
536 input templates and slope. 4c) Observed versus expected variant frequency for SNV at 9110 for
537 vRNA data sets. Lines represent medians and error bars indicate 95% confidence interval.
538 Asterisk indicate p-value less than 0.05 as determined by Kruskal-Wallis.

539

540 Figure 5 A-F) Observed variant frequency of SIVmac239-24x SNPs in varying
541 SIVmac239:SIVmac239-24x ratios by vRNA (left) or cDNA (right) input templates. Dotted line
542 indicates expected variant frequency. Significance determined by Kruskal-Wallis and designated
543 by asterisks. $p < 0.05$ (*), $p < 0.01$ (**), $p < 0.001$ (***), $p < 0.0001$ (****).