

1 **Long-range HIV genotyping using viral RNA and proviral DNA for analysis**
2 **of HIV drug-resistance and HIV clustering**

3

4 Vlad Novitsky,^a Melissa Zahraban-Steele,^a Mary Fran McLane,^a Sikhulile Moyo,^b

5 Erik Widenfelt,^b Simani Gaseitsiwe,^b Joe Makhema,^b and M. Essex^{a,b,#}

6

7 Harvard School of Public Health, Boston, Massachusetts, USA^a; Botswana

8 Harvard AIDS Institute Partnership, Gaborone, Botswana^b

9

10 Running Head: Long-range HIV genotyping

11

12 #Address correspondence to M. Essex, messex@hsph.harvard.edu.

13 **Abstract**

14

15 The goal of the study was to improve the methodology of HIV genotyping for
16 analysis of HIV drug resistance and HIV clustering. Using the protocol of Gall et
17 al. (J Clin Microbiol 50:3838–44, 2012), we developed a robust methodology for
18 amplification of two large fragments of viral genome covering about 80% of the
19 unique HIV-1 genome sequence. Importantly, this method can be applied to both
20 viral RNA and proviral DNA amplification templates, allowing genotyping in HIV-
21 infected subjects with suppressed viral load (e.g., subjects on ART). The two
22 amplicons cover critical regions across the HIV-1 genome (including *pol* and
23 *env*), allowing analysis of mutations associated with resistance to protease
24 inhibitors, reverse transcriptase inhibitors (NRTIs and NNRTIs), integrase strand
25 transfer inhibitors, and virus entry inhibitors. The two amplicons generated span
26 7,124 bp, providing substantial sequence length and number of informative sites
27 for comprehensive phylogenetic analysis and greater refinement of viral linkage
28 analyses in HIV prevention studies. The long-range HIV genotyping from proviral
29 DNA was successful in about 90% of 212 targeted blood specimens collected in
30 a cohort where the majority of patients had suppressed viral load, including 65%
31 of patients with undetectable levels of HIV-1 RNA load. The generated amplicons
32 could be sequenced by different methods, such as population Sanger
33 sequencing, single-genome sequencing, or next-generation ultra-deep
34 sequencing. The developed method is cost-effective – the cost of the long-range

- 35 HIV genotyping is under \$150 per subject (by Sanger sequencing), and has the
36 potential to enable the scale-up of public health HIV prevention interventions.

37 **Introduction**

38

39 HIV genotyping is a critical tool for antiviral drug-resistance testing that has
40 revolutionized HIV care and advanced HIV-related research. Routine
41 antiretroviral (ARV) drug-resistance testing is useful in choosing an optimal
42 treatment regimen and monitoring its efficiency in clinical practice (1-12). HIV
43 genotyping has been used successfully in research on HIV transmission clusters
44 and HIV transmission dynamics (13-35).

45

46 Initial broadly used ARV regimens included combinations of nucleoside reverse
47 transcriptase inhibitors (NRTIs) and non-nucleoside reverse transcriptase
48 inhibitors (NNRTIs). To monitor the emergence of drug-resistant mutations
49 associated with NRTIs and NNRTIs, HIV genotyping targeted viral sequences
50 spanning an approximately 1,000 to 1,300 bp region of the HIV-1 genome,
51 encoding viral protease and partial reverse transcriptase (RT) using viral RNA as
52 a template for amplification. While the RNA-based approach works well in ART-
53 naïve individuals, it is less successful if levels of viral replication are low, such as
54 in individuals on ART. The sequence length of traditional RNA-based HIV
55 genotyping for drug resistance is relatively short, and does not cover the HIV-1
56 region encoding viral integrase, or the viral envelope, hindering analysis of drug-
57 resistant mutations associated with integrase strand transfer inhibitors, or entry
58 inhibitors. The global scale-up of ARV treatment and successful introduction of
59 integrase strand transfer inhibitors and entry inhibitors into clinical trials and

60 clinical practice necessitates modification of traditional methods of HIV
61 genotyping.

62

63 Two commercial genotyping assays, ViroSeq HIV-1 from Abbott Molecular and
64 TruGene HIV-1 from Siemens Molecular Diagnostics, have been widely used for
65 analysis of HIV-1-associated drug resistance. Both genotyping kits were
66 extensively tested and validated (36-45). While the ViroSeq HIV-1 kit is still on
67 the market, Siemens discontinued selling and supporting the TruGene HIV-1 kit
68 in 2014. ViroSeq HIV-1 kit covers the entire protease coding region and the RT
69 region encoding the first 320 amino acids. The TruGene HIV-1 sequences span
70 the protease (amino acids 4 to 99) and RT (amino acids 40 to 240) coding
71 regions. The CDC supplies WHO-designated and CDC-supported PEPFAR
72 Genotyping Laboratories with the ATCC HIV-1 Drug Resistance Genotyping Kit
73 (46) for drug-resistance testing. Many experienced genotyping laboratories have
74 developed their own “in-house” amplification and sequencing protocols (11, 47-
75 56), including identification of minor viral variants that are normally missed by
76 commercial genotyping kits (57-61). All of these approaches generally include
77 smaller and more restricted regions for testing HIV-1 drug resistance.

78

79 Recently, the protocol developed by Gall et al. (62) has enabled high-throughput
80 near full-length HIV-1 genome genotyping in individuals infected with multiple
81 HIV-1 subtypes. This method has become the cornerstone of the PANGEA-HIV
82 Consortium (Phylogenetics and Networks for Generalized HIV Epidemics in

83 Africa; <http://www.pangea-hiv.org/>) aiming to establish worldwide scientific
84 collaborations across phylogenetics, public health and epidemiology. The Gall
85 protocol (62) targets viral RNA as a template for cDNA synthesis and
86 amplification, and is very robust and reproducible when HIV-1 RNA load is high
87 (e.g., above 10,000 cps/ml). However, specimens with levels of HIV-1 RNA
88 below 1,000 cps/ml, or lower thresholds, present a substantial challenge, and few
89 of those samples could be genotyped. This is consistent with the commercially
90 available assays for HIV drug-resistance genotyping, ViroSeq and TrueGene,
91 which are unable to genotype specimens with low or undetectable HIV-1 RNA
92 load.

93

94 In HIV infection, proviral DNA presents an alternative template for HIV
95 genotyping. Drug-resistance mutations detected in viral RNA from plasma and
96 proviral DNA from PBMCs or dried blood spots (DBS) show a substantial
97 correlation in treated patients, suggesting that either compartment is suitable for
98 the detection of mutations as a virological guide for clinical care (63-65).

99

100 It is known that amplified HIV sequences and sequences from proviral DNA could
101 have substantial numbers of G-to-A transitions. Such an inordinate number of
102 identical guanine-to-adenine transitions is a retroviral signature known as
103 hypermutation (66-69). G-to-A hypermutations produce multiple stop codons and
104 reduce HIV replication, leading to an evolutionary dead end. It is an innate host
105 intra-cellular defense mechanism. Host factors APOBEC3F and APOBEC3G

106 induce G-to-A substitutions in reverse transcribed nascent retroviral DNA (70). G-
107 to-A hypermutations play an important role in the evolution of antiretroviral drug
108 resistance (71, 72) and could be associated with ART failure (73). The extent of
109 G-to-A hypermutations is not associated with levels of HIV-1 RNA (74), although
110 hypermutations are frequent in virologic controllers (75). For sequence quality
111 control, it is important that G-to-A hypermutations are not products of PCR
112 amplification (76).

113

114 In this study, we present a technique for long-range HIV genotyping using
115 proviral DNA, as well as viral RNA, as templates for amplification and
116 sequencing. The outcome of the long-range HIV genotyping are two large
117 fragments that span about 80% of the unique full-length HIV-1 genome
118 sequence. The proposed technique is a modification of the method by Gall et al.
119 (62). The key modifications include using (1) a proviral DNA template, (2) extra-
120 round PCR, (3) selection of robust primers, and (4) modified running conditions.
121 To illustrate the potential utility of the long-range HIV genotyping, the technique
122 was applied to a set of specimens collected in Botswana.

123

124

125 **Materials and Methods**

126

127 *Study subjects*

128 The technique of long-range HIV genotyping was applied to specimens collected
129 within three Botswana-Harvard AIDS Institute Partnership (BHP) studies
130 performing viral genotyping: An HIV Prevention Program for Mochudi, Botswana
131 (Mochudi Prevention Project, or MPP; R01 AI083036; PI: M. Essex) (34, 35), the
132 GWAS on Determinants of HIV-1 Subtype C Infection study (RC4 AI092715; PI:
133 M. Essex), and the Botswana Combination Prevention Project (BCPP, or Ya
134 Tsie; U01 GH000447; PI: M. Essex) (77). All study subjects signed a consent
135 form and donated a blood sample for viral genotyping. The first large fragment of
136 HIV-1 genome, "Amplicon 1," was amplified and sequenced in 649 HIV-infected
137 subjects (single sequence per subject) originating from eight geographic localities
138 in Botswana: Digawana, Gaborone, Lobatse, Mochudi, Molapowabojang,
139 Molepolole, Otse, and Ranaka. The second large fragment of HIV-1 genome,
140 Amplicon 2, was amplified and sequenced in 90 subjects (work is still in
141 progress) originating from Mochudi, Molapowabojang, Otse, and Ranaka.

142

143 A total of 212 specimens from the BCPP study were used for analysis of
144 genotyping efficiency. These samples were collected consecutively from subjects
145 participating in the BCPP baseline household survey (20% of households) in the
146 first four communities, Ranaka, Digawana, Molapowabojang and Otse, from
147 November 2013 to June 2014. Specimens from two other studies, MPP and

148 GWAS, represented subsets successfully amplified in the past for a shorter
149 region of HIV-1 *env* gp120, V1C5. Due to potential selection bias, the MPP and
150 GWAS specimens were not used in analysis of genotyping efficiency.

151

152 For sequences used in this study, the accession numbers are KR860607–
153 KR861255 for 649 amplicon 1 sequences and KR861256–KR861345 for 90
154 amplicon 2 sequences.

155

156 *Analyzed regions of the HIV-1 genome*

157 The extent of HIV clustering was analyzed by using the following sub-genomic
158 regions across the HIV-1 genome: (1) Amplicon 1 spanning the 3'-end of *gag* and
159 almost the entire *pol* and corresponding to “amplicon 2” in Gall et al. (62), nt
160 positions 1,486–5,058; (2) Amplicon 2 spanning *vpu*, *env*, *nef* and TATA-box in
161 the U3 region of 3'-LTR and corresponding to “amplicon 4” in Gall et al. (62), nt
162 positions 5,967–9,517; (3) ViroSeq – a partial *pol* sequence spanning the region
163 encoding HIV-1 protease and the first 335 amino acids of reverse transcriptase,
164 and corresponding to the sequence produced by ViroSeq (39, 44, 45, 78), nt
165 positions 2,253–3,554; and (4) V1C5 – a partial *env* sequence spanning the
166 region encoding gp120 V1C5 (34, 79, 80), nt positions 6,570–7,757. In addition,
167 the following combinations of the sub-genomic regions included concatenated
168 Amplicon 1 + Amplicon 2 and Amplicon 1 + V1C5. All multiple sequence codon-
169 based alignments were generated using MUSCLE (81) in MEGA6 (82).

170

171 To prevent sample contamination, basic lab rules were enforced and included
172 controlled flow of specimens, use of dedicated areas and equipment, proper
173 training, and routine implementation of QA/QC program.

174

175 *Analysis of drug resistance*

176 The WHO 2009 list of mutations for surveillance of transmitted drug-resistant HIV
177 strains was used for analysis of PI-, NRTI-, and NNRTI-associated mutations (2).

178 The list of PI-associated mutations included 40 mutations at 18 positions across
179 protease. The list of NRTI mutations included 34 mutations at 15 positions in RT.

180 The list of NNRTI mutations included 19 mutations at 10 positions across RT.

181 The IAS-USA list (2014 update) of drug-resistance mutations in HIV-1 was used
182 for analysis of integrase strand transfer inhibitors (20 mutations at 11 positions in
183 integrase) and entry inhibitors (10 mutations at 7 positions in gp41) (3).

184

185 *APOBEC-induced hypermutations*

186 The APOBEC-induced hypermutations were assessed by Hypermut (83) at the
187 LANL HIV Database (<http://www.hiv.lanl.gov/>). The HIV-1C consensus sequence
188 was used as a reference. Two parameters related to APOBEC-induced
189 hypermutations were analyzed, adjusted hypermutations and hypermutation ratio.

190 The adjusted hypermutations were expressed as a number of identified
191 hypermutations adjusted by sequence length. The hypermutation ratio was
192 computed as the ratio between weighted mutations (matched mutation out of
193 potential mutations) and weighted controls (control mutations out of potential

194 controls), and was derived as a statistical outcome of the Hypermute package
195 (83).

196

197 *Definition of HIV cluster*

198 The HIV cluster was defined as a viral lineage that gives rise to a monophyletic
199 sub-tree of the overall phylogeny with strong statistical support. The
200 bootstrapped Maximum Likelihood (ML) method (84-86) was used to determine
201 the statistical support of clusters. The four bootstrap thresholds for identification
202 of HIV clusters were ≥ 0.7 , ≥ 0.8 , ≥ 0.9 , and $= 1.0$. A viral lineage (group, sub-tree)
203 with at least two viral sequences and specified statistical support was considered
204 to be an HIV cluster. Clusters were identified using a depth-first algorithm (87,
205 88), a method for traversing or searching tree or graph data structures starting
206 from the root. This approach eliminated double-counting of viral sequences in
207 clusters when clusters had internal structure with strong support.

208

209 *Confidentiality*

210 The sharing of data, including generated HIV sequences, with the scientific
211 community for the purpose of research is of key importance in ensuring
212 continued progress in our understanding of how to contain the HIV epidemic. The
213 generated HIV sequences were deposited to GenBank. The confidentiality of
214 study subjects was protected by re-coding of deposited HIV sequences at the
215 country level (no community or village data).

216

217 *Phylogenetic inference*

218 The ML tree inference was implemented in RAxML (89, 90) under the GAMMA
219 model of rate heterogeneity. The statistical support for each node was assessed
220 by bootstrap analysis from 100 bootstrap replicates performed with the rapid
221 bootstrap algorithm implemented in RAxML (89). The RAxML runs were
222 performed using RAxML ver.8.0.20 at the high-performance computing cluster
223 Odyssey (<https://rc.fas.harvard.edu/resources/odyssey-architecture/>) at the
224 Faculty of Arts and Sciences, Harvard University (<https://rc.fas.harvard.edu/>).

225

226 *Proportion of HIV-1C sequences in clusters*

227 To test whether the extent of HIV clustering is associated with any sub-genomic
228 region, the proportion of clustered sequences was compared between long
229 (Amplicon 1, Amplicon 2, concatenated Amplicons 1 + 2, and concatenated
230 Amplicon 1 + V1C5) and short (ViroSeq and V1C5) HIV-1C sequences. The
231 proportion of HIV sequences in clusters was estimated at the bootstrap
232 thresholds for cluster definition from 0.7 to 1.0 under ML inference.

233

234 *Statistical analysis*

235 The HIV sequences in clusters were enumerated with PhyloPart v.2 (88) using
236 bootstrap thresholds 0.7, 0.8, 0.9 and 1.0. All confidence intervals of estimated
237 proportions are asymptotic 95% binomial confidence intervals (95% CI)
238 computed with the `prop.test()` function in R version 3.1.2 (91). Comparisons of
239 continuous outcomes between two groups were performed using the Wilcoxon

240 Rank Sum test. P-values less than 0.05 were considered statistically significant.
241 All reported p-values are 2-sided. Proportions of viral sequences in clusters
242 between targeted loci were compared by McNemar's test in R, and p-values less
243 than 1.0E-04 were considered statistically significant. All plots were produced in
244 R. All figures were finalized in Adobe Illustrator CS6.

245

246

247 **Results**

248

249 *Long-range HIV genotyping*

250 The original protocol for near full-length HIV-1 genome genotyping by
251 amplification of four large overlapping amplicons in a single round of RT-PCR
252 using viral RNA as a template was developed by Gall et al. (62). The protocol by
253 Gall et al. (62) is robust and highly reproducible for samples with relatively high
254 HIV-1 RNA. However, specimens with low, or undetectable levels of HIV-1 RNA
255 presented a substantial challenge for amplification from viral RNA. Attempts to
256 apply the original protocol to proviral DNA produced large number of non-specific
257 products evident from smeared "ladders" on the electrophoretic gel (data not
258 shown).

259

260 The modifications of the Gall et al. (62) protocol included the following steps: (1)
261 focus on 2 (amplicons #2 and #4 in the original protocols) instead of 4 amplicons;
262 (2) extra round of PCR: the amplified ~8.3 kb product was used as a template for

263 the second round of PCR; (3) highly-specific primers for the first round of PCR
264 and for cDNA synthesis (for viral RNA templates); and (4) modified PCR running
265 conditions.

266

267 The rationale for focusing on two instead of four amplicons was driven by a
268 balance between sequencing data and cost. Two amplicons have lengths of
269 3,574 bp and 3,550 bp (HXB2 nt length), which cumulatively covers about 80% of
270 the unique full length HIV-1 genome sequence (Figure 1). The first amplicon
271 (corresponds to amplicon #2 in Gall et al. paper (62)) spans partial *gag* at the 3'-
272 end and almost the entire *pol* (HXB2 nt positions 1,486–5,058). The second
273 amplicon (corresponds to amplicon #4 in Gall et al. paper (62)) spans *vpu*, *env*,
274 *nef* and 3'-LTR up to TATA-box in the U3 region (HXB2 nt positions 5,967–
275 9,517).

276

277 Amplification of a large fragment spanning almost the entire HIV-1 genome (Fig.
278 1; hatched bar) was introduced as the first round of PCR (RT-PCR for RNA
279 template). Primers OFM19 and SK145 (Supplementary Table S1) substantially
280 increased specificity of viral amplification. For proviral DNA template, the 1st
281 round of PCR was run with primers SK145 and OFM19. The PrimeSTAR GXL
282 DNA Polymerase (Takara; cat. #R050A) was used in 30 amplification cycles with
283 annealing temperature at 62 °C (98°C for 15 sec → 62°C for 30 sec → 68°C for 9
284 min cycling). For RNA template, cDNA synthesis with primer OFM19 was
285 followed by PCR with primers SK145 and OFM19 in a single tube RT-PCR. The

286 SuperScript III One-Step RT-PCR High Fidelity enzyme (Invitrogen; cat.
287 #12574035) was used with cDNA synthesis step of incubation at 50°C for 60 min
288 and 94°C for 2 min followed by 30 cycles of amplification in the 1st PCR round
289 (98°C for 15 sec → 62°C for 30 sec → 68°C for 9 min cycling). In cases of
290 specimens from subjects with low viral load, a lower annealing temperature
291 between 58 °C and 60 °C was used in the 1st round.

292

293 The 1st round product was used as template in two separate 2nd round PCRs with
294 specific primers (Table S1) to obtain amplicon 1 and amplicon 2 (Fig. 1; gray
295 bars). The PrimeSTAR GXL DNA Polymerase (Takara; cat. #R050A) was used
296 in 30 amplification cycles with annealing temperature at 62 °C (98°C for 15 sec →
297 62°C for 30 sec → 68°C for 4 min cycling). No additional extension step was
298 performed at the end of the run.

299

300 After standard purification with USB ExoSAP-IT (92) (Affymetrix, cat.
301 #782011ML), amplicon 1 was subject for direct Sanger sequencing on both
302 strands using a total of 12 sequencing primers (Supplementary Table S2). In
303 about 30% of cases direct sequencing of amplicon 1 failed apparently due to
304 heterogeneity of amplified product. These cases were cloned, and Sanger
305 sequenced on both strands. All amplicon 2 products were cloned before Sanger
306 sequencing on both strands with a total of 12 sequencing primers
307 (Supplementary Table S2). Direct Sanger sequencing was performed on the ABI
308 3730 DNA Analyzers using BigDye technology.

309

310 High diversity of HIV presents a challenge for direct Sanger sequencing.

311 Samples collected during the early stage of HIV infection are relatively

312 homogeneous (in case of transmission of a single HIV variant). In contrast,

313 samples obtained from chronically infected individuals are likely to include a

314 heterogeneous pool of viral quasispecies. High heterogeneity of viral

315 quasispecies combined with numerous insertions and deletions (indels) could

316 result in low quality of the directly sequenced specimens. In this case cloning

317 may be considered, as an alternative solution to direct sequencing. If time of HIV

318 infection is unknown, the diversity of the targeted region, or sub-region, could

319 guide the initial sequencing strategy. The amplicon 1 spans a relatively

320 conserved region of the HIV-1 genome. In contrast, amplicon 2 includes the most

321 variable regions of the HIV-1 genome, with multiple indels. Our preliminary

322 results suggest that applying cloning to about 30% of amplicon 1 sequences and

323 to 100% of amplicon 2 sequences is the most efficient sequencing strategy to

324 overcome the complexity of HIV quasispecies. The goal of this study was to

325 obtain a single HIV sequence per subject. Therefore, generation of a single

326 amplicon 1 and single amplicon 2 sequence was considered a success. If a study

327 aimed to address multiplicity of HIV infection, or diversity of viral quasispecies,

328 multiple sequences (e.g., 20 per targeted region per subject) could be generated

329 by appropriate amplification methods.

330

331 Cloning was performed by PCR Cloning Kit (NEB, cat. #E1202S) using Fast-
332 Media Amp XGal (Invivogen, cat. #fas-am-x). Ligation, transformation and plating
333 was performed according to manufacturer's instructions. Colonies were checked
334 for insert by EmeraldAmp GT PCR Master Mix (Takara, cat. #RR310A), and
335 submitted to GENEWIZ ([http://www.genewiz.com/public/DNA-sequencing-](http://www.genewiz.com/public/DNA-sequencing-services.aspx)
336 [services.aspx](http://www.genewiz.com/public/DNA-sequencing-services.aspx)) for colony sequencing. A list of sequencing primers used with
337 clones are presented in Supplementary Table S3.

338

339 All sequence contigs were assembled by SeqScape v.2.7.

340

341 *Troubleshooting*

342 Some amplification issues during long-range HIV genotyping such as lack or
343 insufficient amplification, over-amplified product, or presence of multiple bands
344 could be resolved by troubleshooting. The initial amplification results could guide
345 troubleshooting. A lack of visible bands (or weak bands) on the gel after second
346 round PCR could be resolved by decreasing annealing temperature in the first
347 round PCR to 58 °C, and/or increasing the number of cycles in the first round
348 PCR to 35, or increasing the amount of RNA template (e.g., up to 5 µl). The over-
349 amplified products could be overcome by reducing the number of cycles in the
350 first round PCR to 25, or in the second round PCR to 20-25, or by decreasing the
351 amount of input template. Multiple bands on the gel could be resolved by either
352 extracting the right size band from the gel using Wizard SV Gel and PCR Clean-

353 Up System (Promega, cat. #A9281), or re-running the first round PCR in
354 replicates and with serial dilutions.

355

356 *HIV genotyping results*

357 Amplicon 1 was amplified and sequenced in 649 HIV-infected subjects (single
358 sequence per subject), while amplicon 2 was amplified and sequenced in 90
359 subjects.

360

361 The long-range HIV genotyping from proviral DNA was applied to 212 specimens
362 collected from subjects participating in the BCPP baseline household survey in
363 the first four communities, Ranaka, Digawana, Molapowabojang and Otse, from
364 November 2013 till June 2014. The distribution of amplified and sequenced
365 samples from proviral DNA is presented in Table 1. Amplicon 1 was successfully
366 amplified in 89.6% (95% CI 84.5% to 93.2%) cases. Viral sequences were
367 obtained for all amplified samples. The majority of amplified Amplicon 1
368 sequences, 144 of 167 (86.2%; 95% CI 79.8% to 90.9%) were obtained by direct
369 Sanger sequencing. In 23 cases (12.1%; 95% CI 8.0% to 17.8%) amplicon 1
370 sequences obtained by direct Sanger sequencing had gaps that did not exceed
371 10% of the amplicon 1 length. Cloning followed by Sanger sequencing helped to
372 resolve gaps in all 23 cases.

373

374 Levels and distribution of HIV-1 RNA load in amplified and non-amplified
375 specimens from proviral DNA were of particular interest. HIV-1 RNA load data

376 were available for a subset of 202 HIV-positive subjects from BCPP study. The
377 proportion of successfully amplified cases was 89.6% (95% CI from 84.3% to
378 93.3%; Table 1). Sequences were obtained for all amplified products. Partial
379 sequences (less than 10% of missing data) were obtained in 11.6% (95% CI
380 7.5% to 17.4%) of amplified cases.

381

382 Distribution of HIV-1 RNA load among specimens amplified and sequenced from
383 proviral DNA is presented in Figure 2. The histogram shows the distribution of
384 HIV-1 RNA among 202 specimens with available viral load data (both amplified
385 and failed specimens). The distribution indicates that high proportion of subjects
386 participating in the baseline household survey in four BCPP communities had
387 suppressed levels of HIV-1 RNA, primarily due to high proportion of HIV-infected
388 individuals receiving ART. In fact, 71.3% (95% CI 64.4% to 77.3%) of HIV-
389 infected subjects had HIV-1 RNA below 1,000 cps/ml including 65.3% (95% CI
390 58.3% to 71.8%) with undetectable HIV-1 RNA below 40 cps/ml. Distributions of
391 HIV-1 RNA were similar among amplified (n=181) from proviral DNA and failed
392 (n=21) specimens (Figure 2, pie charts).

393

394 Amplification and sequencing of amplicon 2 was completed for 90 subjects.

395 Given that the first round (RT-) PCR product is used for amplification of both
396 amplicons 1 and 2, obtaining amplicon 1 suggests a successful amplification of
397 amplicon 2. Amplification of the overlapping product designated as “amplicon 3”

398 in Gall et al. paper (62) should be possible, as the first round PCR product
399 completely covers “amplicon 3”. This strategy has not been explored yet.

400

401 Overall, the long-range HIV genotyping from proviral DNA (for amplicon 1) was
402 successful in about 90% of targeted blood specimens collected in a cohort where
403 majority of patients had suppressed viral load including 65% patients with
404 undetectable HIV-1 RNA load.

405

406 *Amplification from viral RNA*

407 To assess the utility of long-range HIV genotyping for amplification and
408 sequencing from viral RNA template, we performed a small-scale genotyping
409 (n=32) from viral RNA in plasma (Table 2). HIV-1 RNA load was available for 31
410 of these samples, and was above 1,000 cps/ml in 29 cases. A subset of 23
411 specimens were successfully amplified and sequenced. Interestingly, two of nine
412 specimens that failed amplification from proviral DNA (HIV-1 RNA load 1,576
413 cps/ml and 5,620 cps/ml), were successfully amplified from viral RNA.

414

415 Nine failed cases included one sample with unknown and eight specimens with
416 available viral load. Among the later group, two samples had viral load below
417 1,000 cps/ml (181 and 497 cps/ml), 5 samples had viral load between 1,191 and
418 8,528 cps/ml, and one sample had viral load of 156,821 cps/ml. The later failed
419 sample with high viral load also failed amplification from proviral DNA, apparently
420 suggesting an intrinsic problem with mismatch of amplification primers.

421

422 *Analysis of mutations associated with antiretroviral drug-resistance*

423 Amplicon 1 covers almost the entire HIV-1 *pol* gene and allows analysis of

424 mutations associated with antiretroviral drug-resistance to Protease inhibitors

425 (PI), Nucleoside Reverse Transcriptase inhibitors (NRTI), non-Nucleoside

426 Reverse Transcriptase inhibitors (NNRTI), and Integrase strand transfer

427 inhibitors. Amplicon 2 covers the entire HIV-1 *env* gene and allows analysis of

428 mutations associated with drug-resistance to virus entry inhibitors.

429

430 To illustrate the validity of long-range HIV genotyping for analysis of mutations

431 associated with antiretroviral drug-resistance, we estimated drug-resistance

432 profiles within two groups of specimens originating from MPP and BCPP studies,

433 respectively. Amplicon 1 sets included 192 MPP sequences and 186 BCPP

434 sequences. Amplicon 2 sets included 35 MPP and 55 BCPP sequences.

435

436 Despite relatively rare use of Protease inhibitors in Botswana, mutations

437 associated with drug-resistance to PI were detected at five positions in Protease:

438 D30N (5% MPP and 6% BCPP), M46I (5% MPP and 10% BCPP), G73S (10%

439 MPP and 9% BCPP), I85V (1% MPP), and N88S (1% BCPP). The encoding

440 analysis revealed that all 22 D30N mutations were caused by GAT (Asp) to AAT

441 (Asn) substitution, 26 of 27 M46I mutations were due to ATG (Met) to ATA (Ile)

442 substitution, and 35 of 36 G73S mutations were found because of GGT (Gly) to

443 AGT (Ser) substitution. Thus, it is likely that the majority of identified mutations in
444 the Protease gene were caused by G-to-A hypermutations.

445

446 NRTI and NNRTI have been part of National antiretroviral program in Botswana
447 since 2002. Viral mutations associated with drug-resistance to NRTI were found
448 at the following positions across RT: M41L (1% BCPP), D67N (1% MPP and 2%
449 BCPP), K70R (1% BCPP), K70E (1% BCPP), V75M (1% BCPP), M184V (2%
450 BCPP), M184I (16% BCPP), and T215Y (1% BCPP). Almost all M184I (60 out of
451 61) mutations were caused by ATG (Met) to ATA (Ile) substitution. Mutations to
452 NNRTI were observed at multiple RT positions and demonstrated low frequency:
453 K101E (1% MPP and 1% BCPP), K103N (1% MPP and 3% BCPP), K103S (1%
454 MPP and 1% BCPP), Y181C (1% BCPP), Y188C (1% BCPP), G190A (1%
455 BCPP), G190S (1% BCPP), G190E (1% MPP and 1% BCPP), and P225H (1%
456 BCPP).

457

458 HIV mutations associated to Integrase strand transfer inhibitors were detected at
459 three positions in Integrase: L74M (1% MPP), T97A (3% MPP and 1% BCPP),
460 and E138K (3% MPP and 6% BCPP). Mutations to entry inhibitors were found at
461 the following positions in gp41: G36S (24% BCPP) and V38M (2% BCPP). All 13
462 out of 55 G36S mutations were caused by a switch from GGT (Gly) to AGT (Ser),
463 which is a likely effect of G-to-A hypermutation.

464

465 *G-to-A hypermutations*

466 Presence of G-to-A hypermutations in the products amplified from proviral DNA
467 is not surprising, as massive APOBEC-induced G-to-A transitions in retroviruses
468 are well recognized as a key innate defense by host. Distribution of identified
469 APOBEC-induced hypermutations in HIV-1C sequences amplified from proviral
470 DNA is presented in Figure 3.

471

472 The sequence length among the 649 cases analyzed differed from 3,190 bp to
473 3,625 bp. Therefore, the number of potentially G-to-A hypermutated sites
474 (compared to the HIV-1 subtype C consensus sequence) was adjusted for the
475 sequence length and expressed as a proportion. Two cutoff values, 0.02 and
476 0.05, were used to demonstrate the proportion of viral sequences in the analyzed
477 set with potential G-to-A hypermutations. Figures 3A and 3C demonstrate
478 distribution of G-to-A hypermutations adjusted by sequence length for amplicons
479 1 (n=649) and 2 (n=90), respectively. For example, 125 of 649 amplicon 1
480 sequences (19.3%; 95% CI 16.3% to 22.6%) and 37 of 90 amplicon 2 sequences
481 (41.1%; 95% CI 31.0% to 52.0%) exceeded the 0.02 level of adjusted
482 hypermutations. Or, 23 of 649 amplicon 1 sequences (3.5%; 95% CI 2.3% to
483 5.4%) and 11 of 90 amplicon 2 sequences (12.2%; 95% CI 6.6% to 21.2%) were
484 above the 0.05 level of adjusted hypermutations. Figures 3B and 3D show
485 distribution of the hypermutation ratio estimated by Hypermut (83). Both metrics
486 indicate presence of APOBEC-induced hypermutations among amplicon 1 and
487 amplicon 2 sequences amplified from proviral DNA.

488

489 The majority of viral sequences with antiretroviral mutations had high rates of
490 APOBEC-induced hypermutations, suggesting association between
491 hypermutations and drug-resistant mutations. Distribution of APOBEC-induced
492 hypermutations among 36 MPP sequences with drug-resistant mutations within
493 Protease, RT and Integrase is presented in Supplementary Table S4, while
494 hypermutations among 46 BCPP sequences with drug-resistant mutations are
495 presented in Supplementary Table S5. It is evident that many hypermutated
496 sequences have multiple drug-resistant mutations due to G-to-A transition.

497

498 HIV-1C sequences with identified drug-resistant mutations demonstrated high
499 rate of APOBEC-induced hypermutations. Horizontal boxplots in Figure 3
500 indicate distribution of adjusted number of hypermutations (Figs. 3A and 3C) and
501 hypermutation ratio (Figs. 3B and 3D) among viral sequences with drug-resistant
502 mutations in relation to the distribution of hypermutation parameters in the entire
503 set of sequences. Comparison of these distributions indicates association
504 between APOBEC-induced hypermutations and drug-resistant mutations.

505

506 To further address how G-to-A hypermutations can affect drug-resistance
507 mutations, we compared hypermutations between two groups, *with* and *without*
508 M184I mutations. Individuals with M184I mutations have higher adjusted
509 numbers of hypermutations (Figure 4A) and higher hypermutation ratios (Figure
510 4B). Summary statistics are presented at the bottom of Figure 4. The difference

511 was highly significant for both comparisons (p -value < 0.0001, Wilcoxon Sum
512 Rank test).

513

514 *HIV cluster analysis*

515 To exemplify utility of the long-range HIV genotyping for analysis of HIV
516 transmission dynamics and viral linkage, we compared extent of clustering within
517 viral sequences generated in this study. The concatenated amplicons 1 + 2 span
518 over about 80% of unique HIV-1 genome sequence. In this study the number of
519 matched amplicons 1 + 2 sequences was limited to 83. The extent of HIV
520 clustering within this small set was compared for three long loci, amplicon 1
521 (3,574 bp), amplicon 2 (3,550 bp), and concatenated amplicons 1 + 2 (7,124 bp),
522 and two short loci, ViroSeq (1,263 bp) and V1C5 (1,188 bp).

523

524 The proportion of clustered HIV sequences was compatible within long loci
525 (Table 3). For example, at bootstrap support of ≥ 0.80 , proportion of clustered HIV
526 sequences was 0.265, 0.289, and 0.337 for amplicon 1, amplicon 2 and
527 concatenated amplicons 1 + 2, respectively. For short loci ViroSeq and V1C5 at
528 the same bootstrap support of ≥ 0.80 , proportion of HIV sequences in clusters
529 was 0.157 and 0.145, respectively, The proportion of clustered sequences
530 seemed to be higher for long regions than for short loci, although the difference
531 reached significance of the 0.05 level in selected comparisons only.

532

533 A larger set of available HIV-1C sequences, n=547, included matched viral
534 sequences for amplicon 1 and V1C5 region of gp120 generated in our previous
535 studies (34, 35, 80). Clustering patterns were compared for two long loci,
536 amplicon 1 and concatenated amplicon 1 + V1C5, and for two short regions
537 across the HIV-1 genome, ViroSeq, and V1C5. Similarly to the small set of HIV
538 sequences (n=83), proportion of clustered sequences in the large set (n=547)
539 was higher for long loci than for short regions (Table 4).

540

541 To address whether longer loci are associated with higher extent of HIV
542 clustering, we analyzed congruent (++) and (--) and discordant (+- and -+)
543 clustering between different combinations of long and short HIV-1C sequences
544 (Figure 5). At all bootstrap thresholds from 0.70 to 1.0, amplicon 1 and
545 concatenated amplicon 1 + V1C5 demonstrated higher extent of HIV clustering
546 than ViroSeq and V1C5 sequences (significant difference is highlighted by gray
547 squares on background).

548

549 *Estimated cost of the long-range HIV genotyping*

550 The estimated cost for amplification and Sanger sequencing of both amplicons 1
551 and 2 in this study was \$137.50 for proviral DNA and \$139.75 for viral RNA. This
552 includes cost of reagents, materials and disposables for nucleic acid isolation,
553 amplification (RT-PCR and PCR), purification of amplicons, cloning up to 30% of
554 amplicon 1 products and 100% of amplicon 2 products, and Sanger sequencing.
555 The estimated cost does not include labor, training, supervision, or indirect costs.

556

557

558 **Discussion**

559

560 A technique for long-range HIV genotyping from both viral RNA and proviral DNA
561 has been presented. Using proviral DNA as a template, long-range HIV
562 genotyping was successfully performed in one of the BCPP cohorts with a high
563 proportion of virologically suppressed individuals with the success rate at about
564 90%.

565

566 Both clinical trials and clinical care could benefit from routine use of long-range
567 HIV genotyping. The proposed long-range HIV genotyping has potential to
568 improve methodology of drug-resistance testing, broaden spectrum of monitored
569 ARV's, and enable surveillance of transmitted drug-resistance. Mapping of HIV
570 transmission networks performed by long-range genotyping could help reveal
571 transmitting viral variants in Treatment-as-Prevention studies. Implementation of
572 long-range HIV genotyping could allow greater refinement of viral linkage
573 analyses in HIV prevention studies, and better coordination with evaluation of
574 prevention strategies based on such interventions as behavior change, male
575 circumcision, and treatment as prevention. The cost-effective long-range HIV
576 genotyping technique has the potential to enable the scale-up of public health
577 HIV prevention interventions across communities.

578

579 In this study we demonstrated that long-range HIV genotyping using proviral
580 DNA could be successfully applied to a population with high level of ART-
581 experienced individuals that normally presents a challenge for HIV genotyping
582 from viral RNA. In fact, more than 70% of individuals in the BCPP cohort
583 participating in the baseline household survey in the first four communities had
584 HIV-1 RNA below 1,000 cps/ml including 65% with undetectable levels of HIV-1
585 RNA below 40 cps/ml. The ongoing scale up of National ARV programs in Africa
586 leads to growing number of individuals with suppressed HIV-1 RNA load across
587 communities. The presented technique of long-range HIV genotyping from
588 proviral DNA should alleviate and enable analysis of HIV drug-resistance and
589 HIV transmission dynamics using samples collected from individuals on ART.

590

591 The comprehensive strategy of HIV genotyping could include two steps. First,
592 viral RNA template for amplification could be targeted, if HIV-1 RNA load is
593 relatively high (e.g., above 1,000 cps/ml). If amplification is successful, there is
594 no need for proviral DNA. However, if amplification from viral RNA does not work,
595 or HIV-1 RNA load is below 1,000 cps/ml, or undetectable, using proviral DNA is
596 a logical step toward successful HIV genotyping. A complimentary use of both
597 viral RNA and proviral DNA templates could be an efficient and cost-effective
598 approach for HIV genotyping.

599

600 The long-range HIV genotyping enables analysis of drug resistance (both
601 transmitted and acquired) for all major groups of ARVs, including protease

602 inhibitors, NRTI, NNRTI, integrase strand transfer inhibitors, and virus entry
603 inhibitors. A comprehensive analysis of HIV drug-resistance is feasible due to a
604 long sequence length of generated amplicons that span the HIV-1 *pol* and *env*
605 genes. While long-range HIV genotyping is able to identify drug-resistant
606 mutations, it cannot distinguish transmitted and acquired drug-resistant mutations
607 without additional information on sampling strategy and stage of HIV infection.
608 For example, drug-resistant mutations identified in individuals during early stages
609 of HIV infection (e.g., in seroconverters) are likely to represent transmitted drug-
610 resistant mutations. In contrast, specimens collected in chronic HIV infection, or
611 from individuals on ART, are more likely to be associated with acquired HIV drug
612 resistance.

613

614 The G-to-A hypermutations observed in sequences amplified from proviral DNA
615 and their relation to drug-resistant mutations should be interpreted cautiously in
616 the context of specific study. Our data suggest that G-to-A hypermutations are
617 likely to contribute to critical drug-resistant mutations, such as M184I. We
618 recommend controlling viral sequences generated from proviral DNA for the
619 adjusted number of hypermutations and/or hypermutation ratio using the online
620 package Hypermut (83) at the LANL HIV Database (<http://www.hiv.lanl.gov/>), and
621 the subtype consensus sequence, as a reference. Based on IQR boundaries in
622 our data, the adjusted number of hypermutations above 2.8% (1st IQR in
623 individuals with M184I) indicates a hypermutated sequence, and below 0.5% (3rd
624 IQR in individuals without M184I) suggests a non-hypermutated sequence.

625 Whether HIV-associated drug-resistant mutations should be interpreted
626 differentially depending on the extent of G-to-A hypermutations still needs to be
627 addressed in future studies.
628

629 The study showed utility of the long-range genotyping for analysis of HIV
630 transmission dynamics and HIV clustering. A higher extent of clustering for
631 longer HIV sequences in this study corroborates well with results of our recent
632 study (93) that used a set of near full-length HIV-1C sequences from LANL HIV
633 Database (<http://www.hiv.lanl.gov/>). Longer HIV sequences are more informative
634 for HIV cluster analysis due to a larger number of informative sites (93). The
635 technique of long-range HIV genotyping allows using amplicon 1 and amplicon 2
636 sequences either separately or in concatenation for a powerful cluster analysis.
637 The concatenated amplicons 1 and 2 span about 80% of the unique HIV-1
638 genome sequence, and could be considered as a cheaper alternative to near full-
639 length HIV-1 sequencing. [A combination of conserved \(amplicon 1\) and variable
640 \(amplicon 2\) regions could help to deal with different and/or unknown stages of
641 HIV infection in an analyzed set of viral sequences.](#) The choice of particular
642 bootstrap value and filtering by the threshold of pairwise distances and/or
643 internode certainty (94, 95) could depend on specific scientific question and take
644 into account specifics of analyzed set of sequences including sampling density
645 (35).
646

647 In summary, the presented technique of long-range HIV genotyping using viral
648 RNA and proviral DNA can help in analysis of HIV drug-resistance and HIV
649 clustering in cohorts and populations on ART when amplification from viral RNA
650 is unsuccessful due to low levels of HIV-1 RNA load.

651 **Acknowledgements**

652

653 We are very grateful to all participants of the BHP projects in Botswana. We
654 thank CDC and the Botswana Ministry of Health for their collaboration. The
655 Mochudi Prevention Project in Botswana was supported and funded by NIH grant
656 R01 AI083036, *An HIV Prevention Program for Mochudi, Botswana*. The GWAS
657 on Determinants of HIV-1 Subtype C Infection study was supported and funded
658 by NIH grant RC4 AI092715. The Botswana Combination Prevention Project,
659 BCPP, has been supported by the President's Emergency Plan for AIDS Relief
660 (PEPFAR) through the United States Centers for Disease Control and Prevention
661 under the terms of grant number U01 GH000447. We thank Lendsey Melton for
662 excellent editorial assistance.

663 **References:**

664

- 665 1. **Johnson VA, Calvez V, Gunthard HF, Paredes R, Pillay D, Shafer RW,**
666 **Wensing AM, Richman DD.** 2013. Update of the drug resistance mutations in
667 HIV-1: March 2013. *Top Antivir Med* **21**:6-14.
- 668 2. **Bennett DE, Camacho RJ, Otelea D, Kuritzkes DR, Fleury H, Kiuchi M,**
669 **Heneine W, Kantor R, Jordan MR, Schapiro JM, Vandamme AM,**
670 **Sandstrom P, Boucher CA, van de Vijver D, Rhee SY, Liu TF, Pillay D,**
671 **Shafer RW.** 2009. Drug resistance mutations for surveillance of transmitted HIV-
672 1 drug-resistance: 2009 update. *PLoS One* **4**:e4724.
- 673 3. **Wensing AM, Calvez V, Gunthard HF, Johnson VA, Paredes R, Pillay D,**
674 **Shafer RW, Richman DD.** 2014. 2014 Update of the drug resistance mutations
675 in HIV-1. *Top Antivir Med* **22**:642-650.
- 676 4. **Rojas Sanchez P, Holguin A.** 2014. Drug resistance in the HIV-1-infected
677 paediatric population worldwide: a systematic review. *J Antimicrob Chemother*
678 **69**:2032-2042.
- 679 5. **Ssemwanga D, Lihana RW, Ugoji C, Abimiku A, Nkengasong J, Dakum P,**
680 **Ndembi N.** 2014. Update on HIV-1 Acquired and Transmitted Drug Resistance in
681 Africa. *AIDS Rev* **17**.
- 682 6. **Smit E.** 2014. Antiviral resistance testing. *Curr Opin Infect Dis* **27**:566-572.
- 683 7. **Bhargava M, Cajas JM, Wainberg MA, Klein MB, Pant Pai N.** 2014. Do
684 HIV-1 non-B subtypes differentially impact resistance mutations and clinical

- 685 disease progression in treated populations? Evidence from a systematic review. *J*
686 *Int AIDS Soc* **17**:18944.
- 687 8. **Snedecor SJ, Sudharshan L, Nedrow K, Bhanegaonkar A, Simpson KN,**
688 **Haider S, Chambers R, Craig C, Stephens J.** 2014. Burden of nonnucleoside
689 reverse transcriptase inhibitor resistance in HIV-1-infected patients: a systematic
690 review and meta-analysis. *AIDS Res Hum Retroviruses* **30**:753-768.
- 691 9. **Ambrosioni J, Nicolas D, Sued O, Agüero F, Manzardo C, Miro JM.** 2014.
692 Update on antiretroviral treatment during primary HIV infection. *Expert Rev Anti*
693 *Infect Ther* **12**:793-807.
- 694 10. **Iwuji CC, Orne-Gliemann J, Tanser F, Boyer S, Lessells RJ, Lert F, Imrie J,**
695 **Barnighausen T, Rekacewicz C, Bazin B, Newell ML, Dabis F, Group ATS.**
696 2013. Evaluation of the impact of immediate versus WHO recommendations-
697 guided antiretroviral therapy initiation on HIV incidence: the ANRS 12249 TasP
698 (Treatment as Prevention) trial in Hlabisa sub-district, KwaZulu-Natal, South
699 Africa: study protocol for a cluster randomised controlled trial. *Trials* **14**:230.
- 700 11. **Manasa J, Danaviah S, Pillay S, Padayachee P, Mthiyane H, Mkhize C,**
701 **Lessells RJ, Seebregts C, de Wit TF, Viljoen J, Katzenstein D, De Oliveira T.**
702 2014. An affordable HIV-1 drug resistance monitoring method for resource
703 limited settings. *J Vis Exp* doi:10.3791/51242.
- 704 12. **Siliciano JD, Siliciano RF.** 2013. Recent trends in HIV-1 drug resistance. *Curr*
705 *Opin Virol* **3**:487-494.

- 706 13. **Brenner BG, Roger M, Moisi DD, Oliveira M, Hardy I, Turgel R, Charest H,**
707 **Routy JP, Wainberg MA.** 2008. Transmission networks of drug resistance
708 acquired in primary/early stage HIV infection. *AIDS* **22**:2509-2515.
- 709 14. **Brenner BG, Roger M, Routy JP, Moisi D, Ntemgwa M, Matte C, Baril JG,**
710 **Thomas R, Rouleau D, Bruneau J, Leblanc R, Legault M, Tremblay C,**
711 **Charest H, Wainberg MA.** 2007. High rates of forward transmission events after
712 acute/early HIV-1 infection. *J Infect Dis* **195**:951-959.
- 713 15. **Brenner BG, Roger M, Stephens D, Moisi D, Hardy I, Weinberg J, Turgel R,**
714 **Charest H, Koopman J, Wainberg MA.** 2011. Transmission Clustering Drives
715 the Onward Spread of the HIV Epidemic Among Men Who Have Sex With Men
716 in Quebec. *J Infect Dis* **204**:1115-1119.
- 717 16. **Brenner BG, Wainberg MA.** 2013. Future of phylogeny in HIV prevention. *J*
718 *Acquir Immune Defic Syndr* **63 Suppl 2**:S248-254.
- 719 17. **Leventhal GE, Gunthard HF, Bonhoeffer S, Stadler T.** 2014. Using an
720 epidemiological model for phylogenetic inference reveals density dependence in
721 HIV transmission. *Mol Biol Evol* **31**:6-17.
- 722 18. **Leventhal GE, Kouyos R, Stadler T, Wyl V, Yerly S, Boni J, Cellerai C,**
723 **Klimkait T, Gunthard HF, Bonhoeffer S.** 2012. Inferring epidemic contact
724 structure from phylogenetic trees. *PLoS Comput Biol* **8**:e1002413.
- 725 19. **Stadler T, Bonhoeffer S.** 2013. Uncovering epidemiological dynamics in
726 heterogeneous host populations using phylogenetic methods. *Philos Trans R Soc*
727 *Lond B Biol Sci* **368**:20120198.

- 728 20. **Stadler T, Kouyos R, von Wyl V, Yerly S, Boni J, Burgisser P, Klimkait T,**
729 **Joos B, Rieder P, Xie D, Gunthard HF, Drummond AJ, Bonhoeffer S.** 2012.
730 Estimating the basic reproductive number from viral sequence data. *Mol Biol Evol*
731 **29**:347-357.
- 732 21. **Stadler T, Kuhnert D, Bonhoeffer S, Drummond AJ.** 2013. Birth-death skyline
733 plot reveals temporal changes of epidemic spread in HIV and hepatitis C virus
734 (HCV). *Proc Natl Acad Sci U S A* **110**:228-233.
- 735 22. **Frost SD, Volz EM.** 2010. Viral phylodynamics and the search for an 'effective
736 number of infections'. *Philos Trans R Soc Lond B Biol Sci* **365**:1879-1890.
- 737 23. **Volz EM, Ionides E, Romero-Severson EO, Brandt MG, Mokotoff E,**
738 **Koopman JS.** 2013. HIV-1 Transmission during Early Infection in Men Who
739 Have Sex with Men: A Phylodynamic Analysis. *PLoS Med* **10**:e1001568.
- 740 24. **Volz EM, Koelle K, Bedford T.** 2013. Viral phylodynamics. *PLoS Comput Biol*
741 **9**:e1002947.
- 742 25. **Volz EM, Koopman JS, Ward MJ, Brown AL, Frost SD.** 2012. Simple
743 Epidemiological Dynamics Explain Phylogenetic Clustering of HIV from Patients
744 with Recent Infection. *PLoS Comput Biol* **8**:e1002552.
- 745 26. **Volz EM, Kosakovsky Pond SL, Ward MJ, Leigh Brown AJ, Frost SD.** 2009.
746 Phylodynamics of infectious disease epidemics. *Genetics* **183**:1421-1430.
- 747 27. **Leigh Brown AJ, Lycett SJ, Weinert L, Hughes GJ, Fearnhill E, Dunn DT,**
748 **Collaboration UHDR.** 2011. Transmission network parameters estimated from
749 HIV sequences for a nationwide epidemic. *J Infect Dis* **204**:1463-1469.

- 750 28. **Wertheim JO, Leigh Brown AJ, Hepler NL, Mehta SR, Richman DD, Smith**
751 **DM, Kosakovsky Pond SL.** 2014. The global transmission network of HIV-1. *J*
752 *Infect Dis* **209**:304-313.
- 753 29. **Kuhnert D, Stadler T, Vaughan TG, Drummond AJ.** 2014. Simultaneous
754 reconstruction of evolutionary history and epidemiological dynamics from viral
755 sequences with the birth-death SIR model. *J R Soc Interface* **11**:20131106.
- 756 30. **Bezemer D, Faria NR, Hassan AS, Hamers RL, Mutua G, Anzala O,**
757 **Mandaliya KN, Cane PA, Berkley JA, Rinke de Wit TF, Wallis CL, Graham**
758 **SM, Price MA, Coutinho R, Sanders EJ.** 2013. HIV-1 transmission networks
759 amongst men having sex with men and heterosexuals in Kenya. *AIDS Res Hum*
760 *Retroviruses* doi:10.1089/AID.2013.0171.
- 761 31. **Bezemer D, Ratmann O, van Sighem A, Dutilh BE, Faria N, van den Hengel**
762 **R, Gras L, Reiss P, de Wolf F, Fraser C, ATHENA observational cohort.**
763 2014. Ongoing HIV-1 Subtype B Transmission Networks in the Netherlands,
764 abstr *CROI 2014*, Boston, MA,
- 765 32. **Bezemer D, van Sighem A, Lukashov VV, van der Hoek L, Back N,**
766 **Schuurman R, Boucher CA, Claas EC, Boerlijst MC, Coutinho RA, de Wolf**
767 **F, cohort Ao.** 2010. Transmission networks of HIV-1 among men having sex
768 with men in the Netherlands. *AIDS* **24**:271-282.
- 769 33. **Carnegie NB, Wang R, Novitsky V, De Gruttola V.** 2014. Linkage of Viral
770 Sequences among HIV-Infected Village Residents in Botswana: Estimation of
771 Linkage Rates in the Presence of Missing Data. *PLoS Comput Biol* **10**:e1003430.

- 772 34. **Novitsky V, Bussmann H, Logan A, Moyo S, van Widenfelt E, Okui L,**
773 **Mmalane M, Baca J, Buck L, Phillips E, Tim D, McLane MF, Lei Q, Wang**
774 **R, Makhema J, Lockman S, DeGruttola V, Essex M.** 2013. Phylogenetic
775 Relatedness of Circulating HIV-1C Variants in Mochudi, Botswana. *PLoS One*
776 **8:e80589.**
- 777 35. **Novitsky V, Moyo S, Lei Q, DeGruttola V, Essex M.** 2014. Impact of Sampling
778 Density on the Extent of HIV Clustering. *AIDS Res Hum Retroviruses* **30**:1226-
779 1235.
- 780 36. **Maes B, Schrooten Y, Snoeck J, Derdelinckx I, Van Ranst M, Vandamme**
781 **AM, Van Laethem K.** 2004. Performance of ViroSeq HIV-1 Genotyping System
782 in routine practice at a Belgian clinical laboratory. *J Virol Methods* **119**:45-49.
- 783 37. **Ribas SG, Heyndrickx L, Ondoa P, Franssen K.** 2006. Performance evaluation
784 of the two protease sequencing primers of the Trugene HIV-1 genotyping kit. *J*
785 *Virol Methods* **135**:137-142.
- 786 38. **Church JD, Mwatha A, Bagenda D, Omer SB, Donnell D, Musoke P,**
787 **Nakabiito C, Eure C, Bakaki P, Matovu F, Thigpen MC, Guay LA,**
788 **McConnell M, Fowler MG, Jackson JB, Eshleman SH.** 2009. In utero HIV
789 infection is associated with an increased risk of nevirapine resistance in ugandan
790 infants who were exposed to perinatal single dose nevirapine. *AIDS Res Hum*
791 *Retroviruses* **25**:673-677.
- 792 39. **Cunningham S, Ank B, Lewis D, Lu W, Wantman M, Dileanis JA, Jackson**
793 **JB, Palumbo P, Krogstad P, Eshleman SH.** 2001. Performance of the applied
794 biosystems ViroSeq human immunodeficiency virus type 1 (HIV-1) genotyping

- 795 system for sequence-based analysis of HIV-1 in pediatric plasma samples. *J Clin*
796 *Microbiol* **39**:1254-1257.
- 797 40. **Eshleman SH, Crutcher G, Petrauskene O, Kunstman K, Cunningham SP,**
798 **Trevino C, Davis C, Kennedy J, Fairman J, Foley B, Kop J.** 2005. Sensitivity
799 and specificity of the ViroSeq human immunodeficiency virus type 1 (HIV-1)
800 genotyping system for detection of HIV-1 drug resistance mutations by use of an
801 ABI PRISM 3100 genetic analyzer. *J Clin Microbiol* **43**:813-817.
- 802 41. **Eshleman SH, Guay LA, Mwatha A, Brown ER, Cunningham SP, Musoke P,**
803 **Mmiro F, Jackson JB.** 2004. Characterization of nevirapine resistance mutations
804 in women with subtype A vs. D HIV-1 6-8 weeks after single-dose nevirapine
805 (HIVNET 012). *J Acquir Immune Defic Syndr* **35**:126-130.
- 806 42. **Eshleman SH, Hackett J, Jr., Swanson P, Cunningham SP, Drews B,**
807 **Brennan C, Devare SG, Zekeng L, Kaptue L, Marlowe N.** 2004. Performance
808 of the Celera Diagnostics ViroSeq HIV-1 Genotyping System for sequence-based
809 analysis of diverse human immunodeficiency virus type 1 strains. *J Clin*
810 *Microbiol* **42**:2711-2717.
- 811 43. **Eshleman SH, Hoover DR, Chen S, Hudelson SE, Guay LA, Mwatha A,**
812 **Fiscus SA, Mmiro F, Musoke P, Jackson JB, Kumwenda N, Taha T.** 2005.
813 Nevirapine (NVP) resistance in women with HIV-1 subtype C, compared with
814 subtypes A and D, after the administration of single-dose NVP. *J Infect Dis*
815 **192**:30-36.
- 816 44. **Mracna M, Becker-Pergola G, Dileanis J, Guay LA, Cunningham S, Jackson**
817 **JB, Eshleman SH.** 2001. Performance of Applied Biosystems ViroSeq HIV-1

- 818 Genotyping System for sequence-based analysis of non-subtype B human
819 immunodeficiency virus type 1 from Uganda. *J Clin Microbiol* **39**:4323-4327.
- 820 45. **Sturmer M, Berger A, Doerr HW.** 2003. Modifications and substitutions of the
821 RNA extraction module in the ViroSeq HIV-1 genotyping system version 2:
822 effects on sensitivity and complexity of the assay. *J Med Virol* **71**:475-479.
- 823 46. **CDC.** 2015. ATCC HIV-1 Drug Resistance Genotyping Kit.
824 [http://www.atcc.org/products/cells_and_microorganisms/hiv-](http://www.atcc.org/products/cells_and_microorganisms/hiv-1_drug_resistance_genotyping_kit.aspx)
825 [1_drug_resistance_genotyping_kit.aspx](http://www.atcc.org/products/cells_and_microorganisms/hiv-1_drug_resistance_genotyping_kit.aspx). Accessed 1/6/2015.
- 826 47. **Wallis CL, Papathanasopoulos MA, Lakhi S, Karita E, Kamali A, Kaleebu P,**
827 **Sanders E, Anzala O, Bekker LG, Stevens G, de Wit TF, Stevens W.** 2010.
828 Affordable in-house antiretroviral drug resistance assay with good performance in
829 non-subtype B HIV-1. *J Virol Methods* **163**:505-508.
- 830 48. **Youngpairoj AS, Masciotra S, Garrido C, Zahonero N, de Mendoza C,**
831 **Garcia-Lerma JG.** 2008. HIV-1 drug resistance genotyping from dried blood
832 spots stored for 1 year at 4 degrees C. *J Antimicrob Chemother* **61**:1217-1220.
- 833 49. **Devereux HL, Youle M, Johnson MA, Loveday C.** 1999. Rapid decline in
834 detectability of HIV-1 drug resistance mutations after stopping therapy. *Aids*
835 **13**:F123-127.
- 836 50. **Neogi U, Sahoo PN, De Costa A, Shet A.** 2012. High viremia and low level of
837 transmitted drug resistance in anti-retroviral therapy-naive perinatally-infected
838 children and adolescents with HIV-1 subtype C infection. *BMC Infect Dis* **12**:317.
- 839 51. **Zhou Z, Wagar N, DeVos JR, Rottinghaus E, Diallo K, Nguyen DB, Bassey**
840 **O, Ugbena R, Wadonda-Kabondo N, McConnell MS, Zulu I, Chilima B,**

- 841 **Nkengasong J, Yang C.** 2011. Optimization of a low cost and broadly sensitive
842 genotyping assay for HIV-1 drug resistance surveillance and monitoring in
843 resource-limited settings. *PLoS One* **6**:e28184.
- 844 52. **Inzaule S, Yang C, Kasembeli A, Nafisa L, Okonji J, Oyaro B, Lando R,**
845 **Mills LA, Laserson K, Thomas T, Nkengasong J, Zeh C.** 2013. Field
846 evaluation of a broadly sensitive HIV-1 in-house genotyping assay for use with
847 both plasma and dried blood spot specimens in a resource-limited country. *J Clin*
848 *Microbiol* **51**:529-539.
- 849 53. **Zhang G, Cai F, Zhou Z, DeVos J, Wagar N, Diallo K, Zulu I, Wadonda-**
850 **Kabondo N, Stringer JS, Weidle PJ, Ndongmo CB, Sikazwe I, Sarr A, Kagoli**
851 **M, Nkengasong J, Gao F, Yang C.** 2013. Simultaneous detection of major drug
852 resistance mutations in the protease and reverse transcriptase genes for HIV-1
853 subtype C by use of a multiplex allele-specific assay. *J Clin Microbiol* **51**:3666-
854 3674.
- 855 54. **Acharya A, Vaniawala S, Shah P, Misra RN, Wani M, Mukhopadhyaya PN.**
856 2014. Development, validation and clinical evaluation of a low cost in-house
857 HIV-1 drug resistance genotyping assay for Indian patients. *PLoS One* **9**:e105790.
- 858 55. **Chen JH, Wong KH, Chan K, Lam HY, Lee SS, Li P, Lee MP, Tsang DN,**
859 **Zheng BJ, Yuen KY, Yam WC.** 2007. Evaluation of an in-house genotyping
860 resistance test for HIV-1 drug resistance interpretation and genotyping. *J Clin*
861 *Virology* **39**:125-131.
- 862 56. **Steege K, Demecheleer E, De Cabooter N, Nges D, Temmerman M,**
863 **Ndumbe P, Mandaliya K, Plum J, Verhofstede C.** 2006. A sensitive in-house

- 864 RT-PCR genotyping system for combined detection of plasma HIV-1 and
865 assessment of drug resistance. *J Virol Methods* **133**:137-145.
- 866 57. **MacLeod IJ, Rowley CF, Thior I, Wester C, Makhema J, Essex M, Lockman**
867 **S.** 2010. Minor resistant variants in nevirapine-exposed infants may predict
868 virologic failure on nevirapine-containing ART. *J Clin Virol* **48**:162-167.
- 869 58. **Rowley CF, Boutwell CL, Lee EJ, MacLeod IJ, Ribaud HJ, Essex M,**
870 **Lockman S.** 2010. Ultrasensitive detection of minor drug-resistant variants for
871 HIV after nevirapine exposure using allele-specific PCR: clinical significance.
872 *AIDS Res Hum Retroviruses* **26**:293-300.
- 873 59. **Rowley CF, Boutwell CL, Lockman S, Essex M.** 2008. Improvement in allele-
874 specific PCR assay with the use of polymorphism-specific primers for the
875 analysis of minor variant drug resistance in HIV-1 subtype C. *J Virol Methods*
876 **149**:69-75.
- 877 60. **Rozera G, Abbate I, Bruselles A, Vlassi C, D'Offizi G, Narciso P, Chillemi G,**
878 **Prosperi M, Ippolito G, Capobianchi MR.** 2009. Massively parallel
879 pyrosequencing highlights minority variants in the HIV-1 env quasispecies
880 deriving from lymphomonocyte sub-populations. *Retrovirology* **6**:15.
- 881 61. **Wang C, Mitsuya Y, Gharizadeh B, Ronaghi M, Shafer RW.** 2007.
882 Characterization of mutation spectra with ultra-deep pyrosequencing: application
883 to HIV-1 drug resistance. *Genome Res* **17**:1195-1201.
- 884 62. **Gall A, Ferns B, Morris C, Watson S, Cotten M, Robinson M, Berry N, Pillay**
885 **D, Kellam P.** 2012. Universal amplification, next-generation sequencing, and
886 assembly of HIV-1 genomes. *J Clin Microbiol* **50**:3838-3844.

- 887 63. **Devereux HL, Loveday C, Youle M, Sabin CA, Burke A, Johnson M.** 2000.
888 Substantial correlation between HIV type 1 drug-associated resistance mutations
889 in plasma and peripheral blood mononuclear cells in treatment-experienced
890 patients. *AIDS Res Hum Retroviruses* **16**:1025-1030.
- 891 64. **Steege K, Luchters S, Demecheleer E, Dauwe K, Mandaliya K, Jaoko W,**
892 **Plum J, Temmerman M, Verhofstede C.** 2007. Feasibility of detecting human
893 immunodeficiency virus type 1 drug resistance in DNA extracted from whole
894 blood or dried blood spots. *J Clin Microbiol* **45**:3342-3351.
- 895 65. **Diallo K, Murillo WE, de Rivera IL, Albert J, Zhou Z, Nkengasong J, Zhang**
896 **G, Sabatier JF, Yang C.** 2012. Comparison of HIV-1 resistance profiles in
897 plasma RNA versus PBMC DNA in heavily treated patients in Honduras, a
898 resource-limited country. *Int J Mol Epidemiol Genet* **3**:56-65.
- 899 66. **Vartanian JP, Henry M, Wain-Hobson S.** 2002. Sustained G-->A
900 hypermutation during reverse transcription of an entire human immunodeficiency
901 virus type 1 strain Vau group O genome. *J Gen Virol* **83**:801-805.
- 902 67. **Vartanian JP, Meyerhans A, Asjo B, Wain-Hobson S.** 1991. Selection,
903 recombination, and G----A hypermutation of human immunodeficiency virus type
904 1 genomes. *J Virol* **65**:1779-1788.
- 905 68. **Vartanian JP, Meyerhans A, Sala M, Wain-Hobson S.** 1994. G-->A
906 hypermutation of the human immunodeficiency virus type 1 genome: evidence for
907 dCTP pool imbalance during reverse transcription. *Proc Natl Acad Sci U S A*
908 **91**:3092-3096.

- 909 69. **Wain-Hobson S, Sonigo P, Guyader M, Gazit A, Henry M.** 1995. Erratic G--
910 >A hypermutation within a complete caprine arthritis-encephalitis virus (CAEV)
911 provirus. *Virology* **209**:297-303.
- 912 70. **Harris RS, Liddament MT.** 2004. Retroviral restriction by APOBEC proteins.
913 *Nat Rev Immunol* **4**:868-877.
- 914 71. **McCallum M, Oliveira M, Ibanescu RI, Kramer VG, Moisi D, Asahchop EL,**
915 **Brenner BG, Harrigan PR, Xu H, Wainberg MA.** 2013. Basis for early and
916 preferential selection of the E138K mutation in HIV-1 reverse transcriptase.
917 *Antimicrob Agents Chemother* **57**:4681-4688.
- 918 72. **Fourati S, Malet I, Lambert S, Soulie C, Wirden M, Flandre P, Fofana DB,**
919 **Sayon S, Simon A, Katlama C, Calvez V, Marcelin AG.** 2012. E138K and
920 M184I mutations in HIV-1 reverse transcriptase coemerge as a result of
921 APOBEC3 editing in the absence of drug exposure. *AIDS* **26**:1619-1624.
- 922 73. **Neogi U, Shet A, Sahoo PN, Bontell I, Ekstrand ML, Banerjea AC,**
923 **Sonnerborg A.** 2013. Human APOBEC3G-mediated hypermutation is associated
924 with antiretroviral therapy failure in HIV-1 subtype C-infected individuals. *J Int*
925 *AIDS Soc* **16**:18472.
- 926 74. **Ulena NK, Sarr AD, Hamel D, Sankale JL, Mboup S, Kanki PJ.** 2008. The
927 level of APOBEC3G (hA3G)-related G-to-A mutations does not correlate with
928 viral load in HIV type 1-infected individuals. *AIDS Res Hum Retroviruses*
929 **24**:1285-1290.

- 930 75. **Eyzaguirre LM, Charurat M, Redfield RR, Blattner WA, Carr JK, Sajadi**
931 **MM.** 2013. Elevated hypermutation levels in HIV-1 natural viral suppressors.
932 *Virology* **443**:306-312.
- 933 76. **Cheyrier R, Gratton S, Vartanian JP, Meyerhans A, Wain-Hobson S.** 1997.
934 G --> A hypermutation does not result from polymerase chain reaction. *AIDS Res*
935 *Hum Retroviruses* **13**:985-986.
- 936 77. **Wang R, Goyal R, Lei Q, Essex M, De Gruttola V.** 2014. Sample size
937 considerations in the design of cluster randomized trials of combination HIV
938 prevention. *Clin Trials* **11**:309-318.
- 939 78. **Eshleman SH, Jones D, Flys T, Petrauskene O, Jackson JB.** 2003. Analysis of
940 HIV-1 variants by cloning DNA generated with the ViroSeq HIV-1 Genotyping
941 System. *Biotechniques* **35**:614-618, 620, 622.
- 942 79. **Novitsky V, Lagakos S, Herzig M, Bonney C, Kebaabetswe L, Rossenkhan R,**
943 **Nkwe D, Margolin L, Musonda R, Moyo S, Woldegabriel E, van Widenfelt E,**
944 **Makhema J, Essex M.** 2009. Evolution of proviral gp120 over the first year of
945 HIV-1 subtype C infection. NIHMSID # 79286. *Virology* **383**:47-59. PMID:
946 PMC2642736.
- 947 80. **Novitsky V, Wang R, Rossenkhan R, Moyo S, Essex M.** 2013. Intra-host
948 evolutionary rates in HIV-1C env and gag during primary infection. *Infect Genet*
949 *Evol* **19C**:361-368.
- 950 81. **Edgar RC.** 2004. MUSCLE: multiple sequence alignment with high accuracy and
951 high throughput. *Nucleic Acids Res* **32**:1792-1797.

- 952 82. **Tamura K, Stecher G, Peterson D, Filipski A, Kumar S.** 2013. MEGA6:
953 Molecular Evolutionary Genetics Analysis version 6.0. *Mol Biol Evol* **30**:2725-
954 2729.
- 955 83. **Rose PP, Korber BT.** 2000. Detecting hypermutations in viral sequences with an
956 emphasis on G --> A hypermutation. *Bioinformatics* **16**:400-401.
- 957 84. **Felsenstein J.** 1985. Confidence limits on phylogenies: an approach using a
958 bootstrap. *Evolution* **39**:783-791.
- 959 85. **Felsenstein J.** 2004. *Inferring phylogenies*. Sinauer Associates, Inc.
- 960 86. **Nei M, Kumar S.** 2000. *Molecular evolution and phylogenetics*. Oxford
961 University Press, New York, NY.
- 962 87. **Even S.** 2011. Graphic Algorithms, 2nd ed., p 202. *In* Even G (ed). Cambridge
963 University Press.
- 964 88. **Prosperi MC, Ciccozzi M, Fanti I, Saladini F, Pecorari M, Borghi V, Di**
965 **Giambenedetto S, Bruzzone B, Capetti A, Vivarelli A, Rusconi S, Re MC,**
966 **Gismondo MR, Sighinolfi L, Gray RR, Salemi M, Zazzi M, De Luca A, group**
967 **Ac.** 2011. A novel methodology for large-scale phylogeny partition. *Nat Commun*
968 **2**:321.
- 969 89. **Stamatakis A.** 2006. RAxML-VI-HPC: maximum likelihood-based phylogenetic
970 analyses with thousands of taxa and mixed models. *Bioinformatics* **22**:2688-2690.
- 971 90. **Stamatakis A.** 2014. RAxML version 8: a tool for phylogenetic analysis and
972 post-analysis of large phylogenies. *Bioinformatics* **30**:1312-1313.

- 973 91. **R Core Team.** 2014. *R: A language and environment for statistical computing.*
974 <http://www.R-project.org/>. R Foundation for Statistical Computing, Vienna,
975 Austria. .
- 976 92. **Dugan KA, Lawrence HS, Hares DR, Fisher CL, Budowle B.** 2002. An
977 improved method for post-PCR purification for mtDNA sequence analysis. *J*
978 *Forensic Sci* **47**:811-818.
- 979 93. **Novitsky V, Moyo S, Lei Q, DeGruttola V, Essex M.** 2015. Importance of Viral
980 Sequence Length and Number of Variable and Informative Sites in Analysis of
981 HIV Clustering. *AIDS Res Hum Retroviruses* doi:10.1089/AID.2014.0211.
- 982 94. **Salichos L, Rokas A.** 2013. Inferring ancient divergences requires genes with
983 strong phylogenetic signals. *Nature* **497**:327-331.
- 984 95. **Salichos L, Stamatakis A, Rokas A.** 2014. Novel information theory-based
985 measures for quantifying incongruence among phylogenetic trees. *Mol Biol Evol*
986 **31**:1261-1271.
- 987 96. **Kampstra P.** 2008. Beanplot: A Boxplot Alternative for Visual Comparison of
988 Distributions. *Journal of Statistical Software* **28**:1-9.
- 989
990

991 **Figure Legends**

992

993 **Figure 1.** Overview of the long-range HIV genotyping: 1st and 2nd round products
994 are mapped against the HIV-1 genome structure. First round (RT-) PCR product
995 is shown at the bottom as a hatched bar. Second round PCR products, amplicon
996 1 and amplicon 2 are shown as gray bars.

997

998 **Figure 2.** Distribution of HIV-1 RNA load in BCPP samples (n=202) that were
999 subjects for long-range HIV genotyping using proviral DNA as a template for
1000 amplification. Histogram depicts distribution of HIV-1 RNA in all samples, n=202.
1001 The x axis shows HIV-1 RNA on a log₁₀ scale. Two pie charts illustrate
1002 distribution of HIV-1 RNA among successfully amplified (n=181) and failed
1003 (n=21) samples. Legend at the top right outlines break down intervals of HIV-1
1004 RNA presented on pie charts.

1005

1006 **Figure 3.** Distribution of APOBEC-induced hypermutations in sequences
1007 amplified from proviral DNA (histograms). Horizontal boxplots outline distribution
1008 of APOBEC-induced hypermutations in subsets of sequences with identified
1009 drug-resistance mutations. A: Amplicon 1, n=649, distribution of hypermutations
1010 adjusted by sequence length. B: Amplicon 1, n=649, distribution of hypermutation
1011 ratio data (see Materials and Methods). C: Amplicon 2, n=90, distribution of

1012 hypermutations adjusted by sequence length. D: Amplicon 2, n=90, distribution of
1013 hypermutation ratio data (see Materials and Methods).

1014

1015 **Figure 4.** G-to-A hypermutations in HIV-1C sequences with and without M184I
1016 mutation. Beanplots – a combination of a box plot, density plot, and a rug with
1017 ticks for each value in the middle – are shown (96). Comparison between groups
1018 was performed by Wilcoxon signed rank test. **A:** Hypermutations adjusted by
1019 sequence length. **B:** Hypermutation ratio.

1020

1021 **Figure 5.** Clustering of HIV-1C sequences by loci, n=547. Proportion of HIV-1C
1022 sequences in clusters was estimated by bootstrapped ML inference. The extent
1023 of HIV clustering was analyzed at bootstrap thresholds for cluster definition
1024 ≥ 0.70 , ≥ 0.80 , ≥ 0.90 , and 1.0. The number of viral sequences found in clusters for
1025 specified locus and at specified bootstrap support was compared between loci.
1026 Four loci were used: Amplicon 1 concatenated with V1C5 region of gp120 shown
1027 as “Amp 1 + v1c5”, Amplicon 1 alone as “Amp 1”, ViroSeq sequence as
1028 “ViroSeq”, and V1C5 region of gp120 as “V1C5”. Pie charts show concordant (++)
1029 and (--) and discordant (+- and -+) clustering between specified sequence loci
1030 (the first sign corresponds to the first sequence locus listed). Cases of
1031 significantly different clustering between loci with p-values less than $1.0E-04$ in
1032 McNemar’s test are highlighted by gray squares on the background.

Table 1. Summary of HIV genotyping from proviral DNA, amplicon 1 (BCPP)

Proviral DNA specimens	Total		Subset with available HIV-1 RNA load data	
	n	proportion (95% CI)	n	proportion (95% CI)
Attempted cases	212		202	
Amplified	190	0.896 (0.845–0.932)	181	0.896 (0.843–0.933)
Amplification failed*	22	0.104 (0.068–0.155)	21	0.104 (0.067–0.157)
Sequenced by direct Sanger sequencing**	167	0.879 (0.822–0.920)	160	0.884 (0.826–0.925)
Partial sequences by direct Sanger sequencing: gaps at ≤10% of sequence length; resolved by cloning**	23	0.121 (0.080–0.178)	21	0.116 (0.075–0.174)

* - proportion of failed cases is calculated from the number of attempted cases;

** - proportion of sequenced cases is calculated from the number of amplified cases.

Table 2. Summary of HIV genotyping from viral RNA, amplicon 1 (BCPP)

Viral RNA specimens	Total		Subset with available HIV-1 RNA load data	
	n	proportion (95% CI)	n	proportion (95% CI)
Attempted cases	32		31	
Amplified	23	0.719 (0.530–0.856)	23	0.742 (0.551–0.875)
Amplification failed*	9	0.281 (0.144–0.470)	8	0.258 (0.125–0.449)
Sequenced**	23	1.0 (0.822–1.0)	23	1.0 (0.822–1.0)

* - proportion of failed cases is calculated from the number of attempted cases;

** - proportion of sequenced cases is calculated from the number of amplified cases.

Table 3. Observed proportion of HIV-1C sequences in clusters, small set of sequences (n=83)

Loci	Number (proportion) of HIV-1C sequences in clusters at specified bootstrap support of splits			
	≥0.70	≥0.8	≥0.9	=1.0
Amplicon 1	32 (0.386) ^{♦§}	22 (0.265)	19 (0.229)	11 (0.133)
Amplicon 2	32 (0.386) ^{♦§}	24 (0.289) [§]	24 (0.289) ^{♦§}	12 (0.145)
Amplicon 1 + 2	33 (0.398) ^{♦§}	28 (0.337) ^{♦§}	24 (0.289) ^{♦§}	11 (0.133)
ViroSeq	11 (0.133)	13 (0.157)	11 (0.133)	7 (0.084)
V1C5	16 (0.193)	12 (0.145)	9 (0.108)	4 (0.048)

[♦] - p-value < 0.05 for comparison to ViroSeq (Fisher exact test);

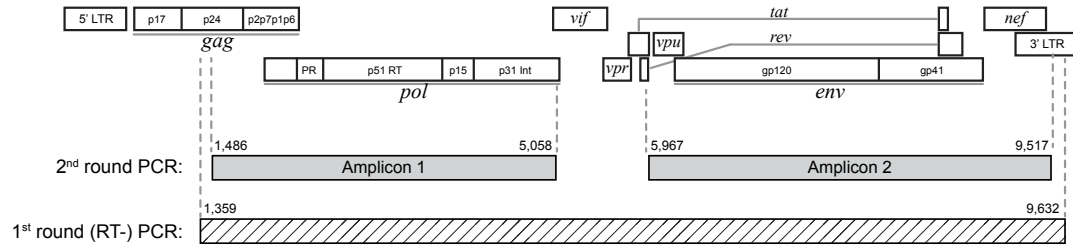
[§] - p-value < 0.05 for comparison to V1C5 (Fisher exact test)

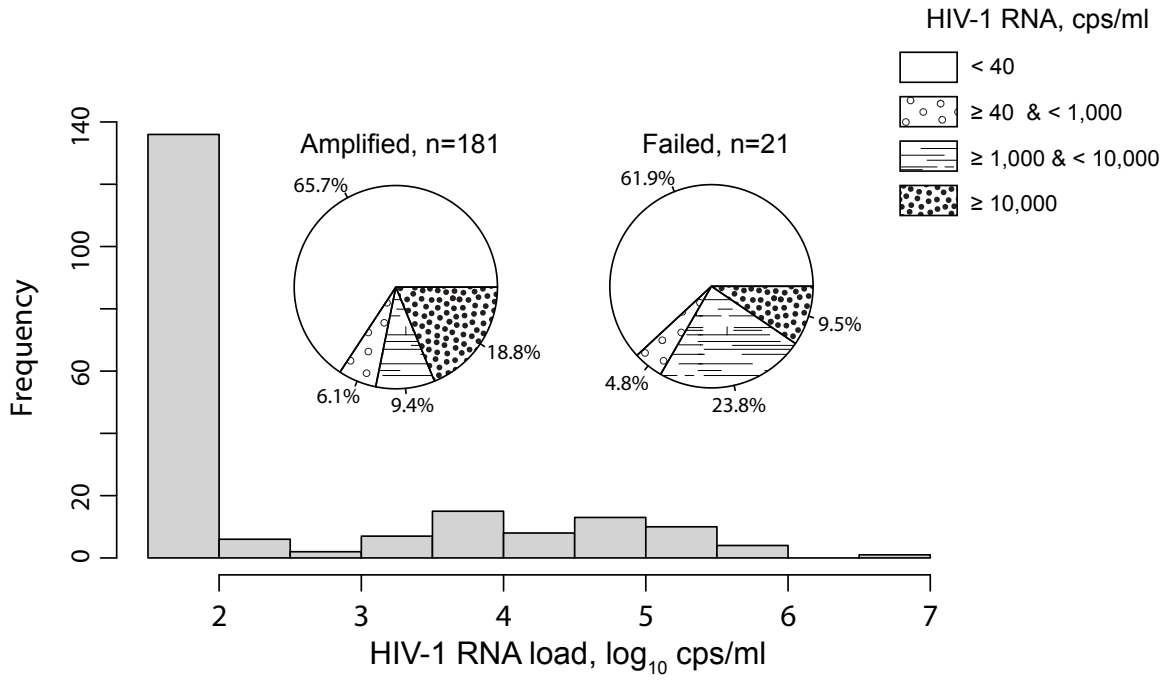
Table 4. Observed proportion of HIV-1C sequences in clusters, large set of sequences (n=547)

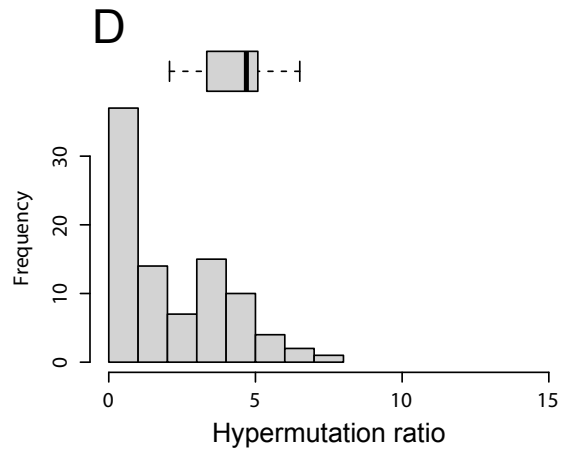
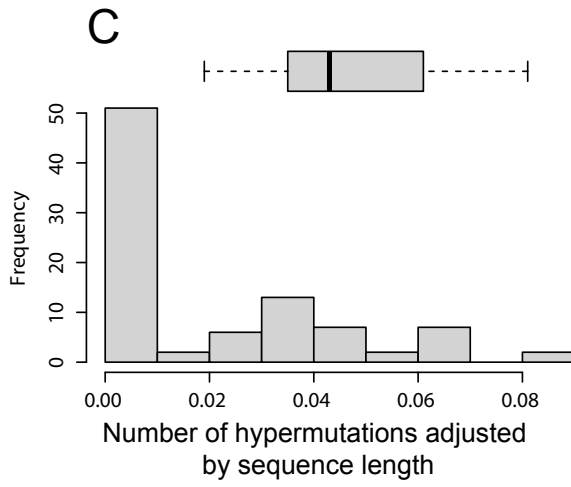
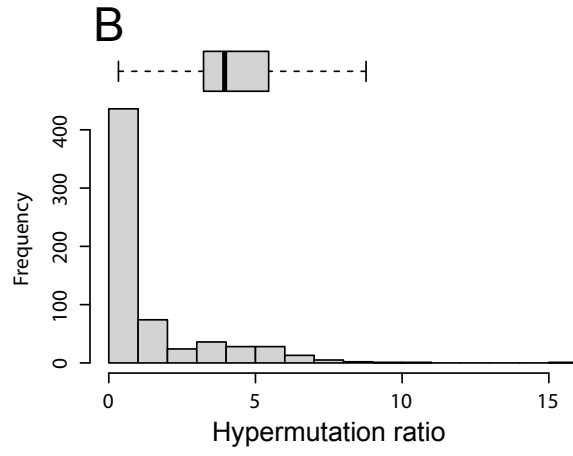
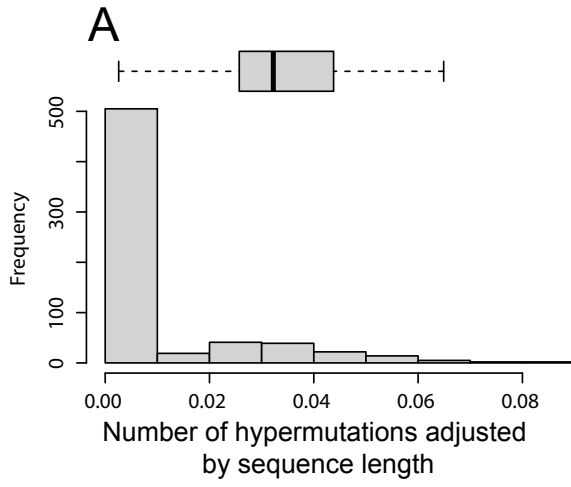
Loci	Number (proportion) of HIV-1C sequences in clusters at specified bootstrap support of splits			
	≥0.70	≥0.8	≥0.9	=1.0
Amplicon 1	251 (0.459) ^{◆§}	215 (0.393) ^{◆§}	159 (0.291) ^{◆§}	88 (0.161) ^{◆§}
Amplicon 1 + V1C5	267 (0.488) ^{◆§}	220 (0.402) ^{◆§}	181 (0.331) ^{◆§}	122 (0.223) ^{◆§}
ViroSeq	120 (0.219)	90 (0.165)	73 (0.133)	34 (0.062)
V1C5	135 (0.247)	114 (0.208)	88 (0.161)	44 (0.080)

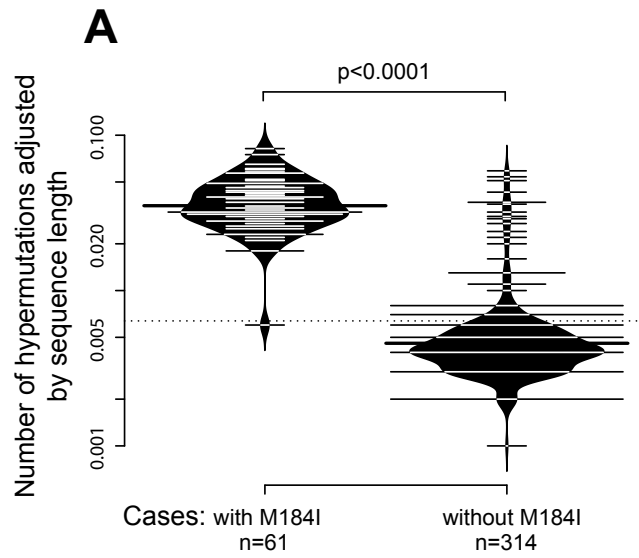
[◆] - p-value < 0.001 for comparison to ViroSeq (Fisher exact test);

[§] - p-value < 0.001 for comparison to V1C5 (Fisher exact test)

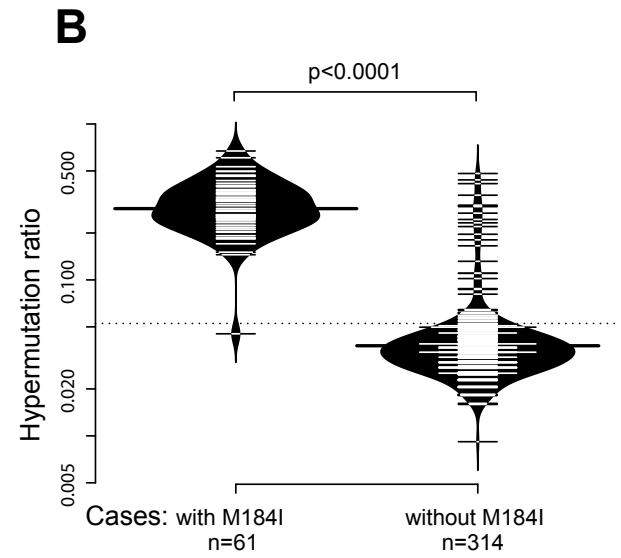








Minimum:	0.0060	0.0010
1 st Quartile:	0.0280	0.0030
Median:	0.0340	0.0040
Mean:	0.0379	0.0058
3 rd Quartile:	0.0460	0.0050
Maximum:	0.0820	0.0590



Minimum:	0.0451	0.0092
1 st Quartile:	0.2330	0.0276
Median:	0.2803	0.0346
Mean:	0.3089	0.0480
3 rd Quartile:	0.3724	0.0431
Maximum:	0.6711	0.4814

