

1 HIV-1 evolutionary dynamics under non-suppressive, 2nd-line protease-inhibitor containing
2 antiretroviral therapy.

3

4 Steven A. Kemp^{1,3}, Oscar Charles¹, Anne Derache², Collins Iwuji^{2,5}, John Adamson², Katya Govender²,
5 Tulio de Oliveira^{2,6}, Nonhlanhla Okesola², Francois Dabis^{7,8}, Darren P. Martin⁹, on behalf of the French
6 National Agency for AIDS and Viral Hepatitis Research (ANRS) 12249 Treatment as Prevention (TasP)
7 Study Group, Deenan Pillay¹, Richard A. Goldstein¹ & Ravindra K. Gupta^{2,3}

8

- 9 1. Division of Infection & Immunity, University College London, London, UK
10 2. Africa Health Research Institute, Durban, South Africa
11 3. Cambridge Institute of Therapeutic Immunology & Infectious Disease (CITIID), Cambridge,
12 UK
13 4. Department of Medicine, University of Cambridge, Cambridge, UK
14 5. Research Department of Infection and Population Health, University College London, United
15 Kingdom.
16 6. KRISP - KwaZulu-Natal Research and Innovation Sequencing Platform, UKZN, Durban, South
17 Africa.
18 7. INSERM U1219-Centre Inserm Bordeaux Population Health, Université de Bordeaux, France.
19 8. Université de Bordeaux, ISPED, Centre INSERM U1219-Bordeaux Population Health, France.
20 9. Department of Integrative Biomedical Sciences, University of Cape Town, South Africa

21

22 Address for correspondence:

23 Steven A Kemp

24 Cambridge Institute for Therapeutic Immunology and Infectious Diseases

25 Jeffrey Cheah Biomedical Centre

26 Cambridge CB2 0AW, UK

27 sk2137@cam.ac.uk

28

29 or

30

31 Ravindra K Gupta

32 Cambridge Institute for Therapeutic Immunology and Infectious Diseases

33 Jeffrey Cheah Biomedical Centre

34 Cambridge CB2 0AW, UK

35 Rkg20@cam.ac.uk

36 **Abstract**

37 Viral population dynamics in long term viraemic antiretroviral therapy (ART) treated individuals have
38 not been well characterised. Prolonged virologic failure on 2nd-line protease inhibitor (PI) based ART
39 without emergence of major protease mutations is well recognised, providing an opportunity to
40 study within-host evolution.

41

42 Using next-generation Illumina short read sequencing and in silico haplotype reconstruction we
43 analysed whole genome sequences from longitudinal plasma samples of eight chronically infected
44 HIV-1 individuals failing 2nd-line regimens from the ANRS 12249 TasP trial, in the absence of high
45 frequency major PI resistance mutations. Plasma drug levels were measured by HPLC. Three
46 participants were selective for in-depth variant and haplotype analyses, each with five or more
47 timepoints spanning at least 16 months. Recombination and linkage disequilibrium between
48 haplotypes and genes was also explored.

49

50 During PI failure synonymous mutations were around twice as frequent as non-synonymous
51 mutations across participants. Prior to or during exposure to PI, we observed several polymorphic
52 amino acids in *gag* (e.g. T81A, T375N) which are have also been previously associated with exposure
53 to protease inhibitor exposure. Although overall SNP frequency at abundance above 2% appeared
54 stable across time in each individual, divergence from the consensus baseline sequence did increase
55 over time. Non-synonymous changes were enriched in known polymorphic regions such as *env*
56 whereas synonymous changes were more often observed to fluctuate in the conserved *pol* gene.
57 Phylogenetic analyses of whole genome viral haplotypes demonstrated two common features:
58 Firstly, evidence for selective sweeps following therapy switches or large changes in plasma drug
59 concentrations, with hitchhiking of synonymous and non-synonymous mutations. Secondly,
60 competition between multiple viral haplotypes that intermingled phylogenetically alongside soft
61 selective sweeps. The diversity of viral populations was maintained between successive timepoints
62 with ongoing viremia, particularly in *env*. Changes in haplotype dominance were often distinct from
63 the dynamics of drug resistance mutations in *reverse transcriptase* (RT), indicating the presence of
64 softer selective sweeps and/or recombination.

65

66 Large fluctuations in variant frequencies with diversification occur during apparently 'stable' viremia
67 on non-suppressive ART. Reconstructed haplotypes provided further evidence for sweeps during
68 periods of partial adherence, and competition between haplotypes during periods of low drug

69 exposure. Drug resistance mutations in RT can be used as markers of viral populations in the
70 reservoir and we found evidence for loss of linkage disequilibrium for drug resistance mutations,
71 indicative of recombination. These data imply that even years of exposure to PIs, within the context
72 of large stable populations displaying ongoing selective competition, may not precipitate emergence
73 of major PI resistance mutations, indicating significant fitness costs for such mutations. Ongoing viral
74 diversification within reservoirs may compromise the goal of sustained viral suppression.
75

76 Introduction

77 Even though HIV-1 infections are most commonly initiated with a single founder virus¹, acute and
78 chronic disease are characterised by extensive inter- and intra-participant genetic diversity^{2,3}. The
79 rate and degree of diversification is influenced by multiple factors, including selection pressures
80 imposed by the adaptive immune system, exposure, and penetration of the virus to drugs, and
81 tropism/fitness constraints relating to replication and cell-to-cell transmission in different tissue
82 compartments^{4,5}. During HIV-1 infection, high rates of reverse transcriptase- (RT) related mutation
83 and high viral turnover during replication result in swarms of genetically diverse variants⁶ which co-
84 exist as quasispecies. Existing literature on HIV-1 intrahost population dynamics is largely limited to
85 untreated infections, predominantly in subtype B infected individuals⁷⁻¹⁰. These works have shown
86 non-linear diversification of virus away from the founder strain during chronic untreated infection.

87

88 Viral population dynamics in long-term viraemic antiretroviral therapy (ART) treated individuals have
89 not been well characterised. HIV rapidly accumulates drug-resistance associated mutations (DRMs),
90 particularly during non-suppressive 1st-line ART^{5,11}. As a result, ART-experienced participants failing
91 1st-line regimens for prolonged periods of time are characterised by high frequencies of common
92 nucleoside reverse transcriptase (NRTI) and non-nucleoside reverse transcriptase (NNRTI) drug
93 resistance mutations (DRMs) such as M184V, K65R and K103N¹². Routinely, 2nd-line ART regimens
94 consist of two NRTIs and underpinned by a boosted protease inhibitor (PIs); PI DRMs are
95 uncommonly reported¹³ however, a situation that differs for less potent drugs used in the early PI
96 era⁵. A number of studies have indicated that less well characterised mutations accumulating in the
97 *gag* gene during PI failure might impact PI susceptibility¹⁴⁻²⁰, though common pathways have been
98 difficult to discern, likely reflecting plasticity to drug escape.

99

100 Prolonged virological failure on PI-based regimens, without emergence of PI DRMs provides an
101 opportunity to study evolution under partially-suppressive ART. The process of selective sweeps in
102 the context of HIV-1 infection has previously been described^{21,22} and it was reported that major PI
103 DRMs and other non-synonymous mutations in regulatory regions such as *pol*, significantly lower
104 fitness^{2,23,24}. However, this has been typically shown outside of the context of longitudinal sampling.
105 HIV has been shown to exhibit significant genetic diversity within infected hosts, with different
106 populations of virus accumulating beneficial mutations – these are referred to as ‘quasispecies’^{25,26}.
107 By sampling participants consistently over several years, we propose that ongoing evolution is driven
108 by the dynamic flux between selection, recombination, and genetic drift.

109

110 We have deployed next-generation sequencing of stored blood plasma specimens from participants
111 in the Treatment as Prevention (TasP) ANRS 12249 study ²⁷, conducted in Kwazulu-Natal, South
112 Africa. All participants were infected with HIV-1 subtype C and characterised as failing 2nd-line
113 regimens containing Lopinavir and Ritonavir (LPV/r), with prolonged virological failure in the absence
114 of major PI mutations ²⁸. In this manuscript, we report details of evolutionary dynamics during non-
115 suppressive 2nd-line ART, through investigation of individual individual quasispecies using a novel
116 computational haplotype reconstruction tool, Haplotype Reconstruction of Longitudinal Deep
117 sequencing data (HaROLD, ²⁹.

118

119 **Results**

120 **Participant Characteristics**

121 Eight south African participants with virological failure of 2nd-line PI based ART and at least two
122 timepoints, with viraemia above 1000 copies/ml were selected from the French ANRS TasP trial for
123 viral dynamic analysis. Participant metadata collected included viral loads, regimens and time since
124 ART initiation (**Table 1**). HIV RNA was isolated from venous blood samples and subject to whole-
125 genome sequencing (WGS) using Illumina technology; from this whole-genome haplotypes were
126 reconstructed. Prior to participation in the TasP trial, participants accessed 1st-line regimens for an
127 average of 5.6yrs (± 2.7 yrs). At baseline enrolment into TasP (whilst failing 1st-line regimens), median
128 viral load was 4.96×10^{10} (IQR: $4.17 \times 10^{10} - 5.15 \times 10^{10}$); 12 DRMs were found at a threshold of >2%; the
129 most common of which were RT mutations K103N, M184V and P225H, consistent with previous use
130 of d4T, NVP, EFV and FTC/3TC. Six of the eight participants had DRMs associated with PI failure at
131 minority frequencies (average 6.4%) and usually at single timepoints throughout the longitudinal
132 sampling. Observed mutations included L23I, I47V, M46I/L, G73S, V82A, N83D and I85V
133 (**Supplementary Tables 1a-3c**). Viral populations of four of the eight participants also carried major
134 integrase strand inhibitor (INSTI) mutations, at minority frequencies (average 5.0%) and usually at
135 single timepoints (T97A, E138K, Y143H, Q148K).

136

137

138 **SNP frequencies and measures of diversity/divergence over time**

139 WGS data was used to measure the changing frequencies of viral single nucleotide polymorphisms
140 (SNPs) relative to a dual-tropic subtype C reference sequence (AF411967) within individuals over
141 time (**Figure 1a-b**). The number of longitudinal synonymous SNPs approximately mirrored the
142 number of non-synonymous SNPs, but the former were two-to-three-fold more common. In most
143 participants, viral populations were homogenous. Diversification, by counting the number of SNPs

144 difference to the reference sequence was considered. There was largely idiosyncratic changes in the
145 number of SNPs over time, with both increases and decreases in number of SNPs, suggesting
146 different population competition, or selective sweeps occurring. From timepoint two onwards (all
147 participants now on 2nd-line, PI-containing regimens for >6 months), all participants (except 28545)
148 had increases in both synonymous and non-synonymous SNPs.

149

150 In previous literature, chronically infected, but untreated HIV-1 patients viral populations showed
151 some reversion to the founder or infecting virus⁷. We assessed this phenomenon in our chronically
152 infected, but treated HIV-1 population. HIV-1, subtype C and M group consensus sequences were
153 downloaded from the Los Alamos National Laboratory (LANL) database and an ‘ancestral’ consensus
154 sequence was constructed according to the Materials and Methods. Diversification towards or away
155 from the baseline (or infecting sequence) as well as both the subtype C and subtype M consensus
156 was measured.

157

158 When considering the whole genome, linear regression suggests that all sequences continued to
159 diversify away from the infecting/founder strain, though R^2 and p-values suggests the diversification
160 towards or away from the founder was not significant in most cases. Despite this, all participants
161 appeared to diversify away from the infecting strain (**Figure 1C**), all but two participants (15664 &
162 29447) diversified away from the ancestral subtype C (**Figure 1D**), and all but three participants
163 (15664, 22828 & 29447) diversified away from the ancestral Subtype M (**Figure 1E**). This is contrary
164 to current literature, and may be a novel feature for patients who fail 2nd-line regimens or
165 sporadically adhere to therapies. To determine whether specific genomic regions were responsible
166 for reversion or divergence, we re-examined the divergence from simulated ancestors across the
167 genome in a sliding window of 1000bp (**Supplementary Figures 1-3**). This revealed that whilst there
168 were heterogeneous patterns of divergence, when considering portions of each genome
169 independently there was a general trend towards divergence, although curiously the regions 1-
170 2000bp (covering the *gag* gene) appears to be show reversion (to the baseline strain). The highest
171 degree of diversification occurred in regions 2000-3000bp and 7000-8000bp, corresponding to the
172 *protease* and *reverse transcriptase* genes, and the *env* gene respectively. For any given patient their
173 HIV genomes are able to both revert in part, and diverge in others, likely enabled by recombination
174 which unlinks hyper-variable loci from strongly constrained neighboring sites.

175

176 To assess the relationship of the observed divergent patterns, we examined nucleotide diversity by
177 considering all pairwise nucleotide distances of each consensus sequence, by timepoint and

178 participant using a multidimensional scaling approach³⁰. Intra-participant nucleotide diversity varied
179 considerably between participants (**Figure 2a**). Some participants showed little diversity between
180 timepoints (e.g. participant 16207), whereas others showed higher diversity between timepoints
181 (e.g. participant 22763). Some participants were tightly clustered, suggesting little change over time
182 (**Figure 3a**, participants 16207, 26892 & 47939), compared to others (participants 22828 & 28545).
183 To corroborate the MDS approach, we used a novel method of examining nucleotide diversity of
184 longitudinal timepoints using all positional information from BAM files (**Supplementary Figure 4**).

185

186 **Phylogenetic analysis of inferred haplotypes**

187 The preceding diversity assessments suggested the existence of distinct viral haplotypes within each
188 participant. We therefore used a recently reported computational tool²⁹ to infer 289 unique
189 haplotypes across all participants, with between 11 and 32 haplotypes (average 21) per participant.
190 The number haplotypes changed dynamically between successive timepoints indicative of
191 dynamically shifting populations (**Figure 2B**). To ensure that haplotypes were sensibly reconstructed,
192 a phylogeny of all consensus sequences was also inferred (**Supplementary Figure 5**). Furthermore, to
193 ensure that reconstructed haplotypes were sensible, a subsequent MDS plot of all viral haplotypes
194 was constructed (**Supplementary Figure 6**).

195

196 **Linkage Disequilibrium and Recombination**

197 LD between two pairwise loci is reduced by recombination, such that LD tends to be higher for loci
198 that are close and lower for more distant loci³¹. HIV is known to rapidly recombine such that
199 sequences are not generally in Linkage Disequilibrium (LD) beyond 400bp⁷. The significance of
200 recombination in intra-host, single infection setting is less well understood³². To assess if intra-patient
201 recombination was occurring between patient haplotypes for three most sampled participants, we
202 determined LD decay patterns, assuming that if there was random recombination, this would equate
203 to a smooth LD decay pattern. This was not observed, rather, each participant demonstrated a
204 complex decay pattern, consistent with non-random recombination along the genome (**Figure 3A**).
205 Given this, we characterized recombination patterns (**Figure 3b**). Perceived recombination
206 breakpoints were recurrent within participants and identifiable over successive timepoints. DRMs
207 were gained over successive timepoints for time for participant 22763 whereas in participant
208 15664, the reverse was true, whereby the converse is seen for patient 15664, whereby there was a
209 gradual loss of DRMs. This indicates that the not all DRMs were required to overcome sporadic drug
210 pressure and the original drug-resistant virus with lower fitness was selected preferentially by ART

211 drug pressure. Participant 16207 had recombinant breakpoints localised in the the *pol* gene in two
212 timepoints, though it retained its majority DRM (K103N) across all haplotype populations.

213

214

215 **Changing landscapes of non-synonymous and synonymous mutations**

216 In the absence of major PI mutations, we first examined non-synonymous mutations across the
217 whole genome (**Figures 4-6**), with a specific focus on *pol* (to observe first and second line NRTI-
218 associated mutations) and *gag* (given its involvement in PI susceptibility). We and others have
219 previously shown that *gag* mutations accumulate during non-suppressive PI therapy^{33,34}. There are
220 also data suggesting associations between *env* mutations and PI exposure^{35,36}. **Supplementary**
221 **Tables 1-3** summarise the changes in variant frequencies of *gag*, *pol* and *env* mutations in
222 participants over time. We found between two and four mutations at sites previously associated
223 with PI resistance in each participant, all at persistently high frequencies (>90%) even in the absence
224 of presumed drug pressure. This is explained by the fact that a significant proportion of sites
225 associated with PI exposure are also polymorphic across HIV-1 subtypes^{18,37}. To complement this
226 analysis, we examined underlying synonymous mutations across the genome. This revealed complex
227 changes in the frequencies of multiple nucleotide residues across all genes. These changes often
228 formed distinct ‘chevron-like’ patterns between timepoints (**Figures 4c & 5b**), indicative of linked
229 alleles dynamically shifting and suggestive of competition between viral haplotypes.

230

231 **Participant 15664** had consistently low drug plasma concentration of all drugs at each measured
232 timepoint, with detectable levels measured only at month 15 and beyond (**Figure 4a**). At baseline,
233 whilst on NNRTI-based 1st-line ART, known NRTI (M184V) and NNRTI (K103N and P225H) DRMs⁵
234 were at high prevalence in the virus populations which is as expected whilst adhering to 1st-line
235 treatments. Haplotype reconstruction and subsequent analysis inferred the presence of a majority
236 haplotype carrying all three of these mutations at baseline, as well as a minority haplotype with the
237 absence of P225H (**Figure 4d**, dark grey circles). Following the switch to a 2nd-line regimen, variant
238 frequencies of M184V and P225H dropped below detection limits (<2% of reads), whilst K103N
239 remained at high frequency (**Figure 4B**). Haplotype analysis was concordant, revealing that viruses
240 with K103N, M184V and P225H were replaced by haplotypes with only K103N (**Figure 4D**, light grey
241 circles). At timepoint two (month 8), there were also numerous synonymous mutations observed at
242 high frequency in both *gag* and *pol* genes, corresponding with the switch to a 2nd-line regimen. At
243 timepoint three (15 months post-switch to 2nd-line regimen) drug concentrations were highest,
244 though still low in absolute terms, indicating partial adherence. Between timepoints three and four

245 we observed a two-log reduction in viral load, with modest change in frequency of RT DRMs.
246 However, we observed synonymous variant frequency shifts predominantly in both *gag* and *pol*
247 genes, as indicated by multiple variants increasing and decreasing contemporaneously, creating
248 characteristic chevron patterning (**Figure 4b**). However many of the changes were between
249 intermediate frequencies, (e.g. between 20% and 60%), which differed from changes between time
250 points one and two where multiple variants changed more dramatically in frequency from <5% to
251 more than 80%, indicating harder selective sweep. These data are in keeping with a soft selective
252 sweep between time points three and five. Between timepoints five and six, the final two samples,
253 there was another population shift - M184V and P225H frequencies fell below the detection limit at
254 timepoint six, whereas the frequency of K103N dropped from almost 100% to around 75% (**Figure**
255 **4b**). This was consistent with haplotype reconstruction, which inferred a dominant viral haplotype at
256 timepoint six bearing only K103N, as well as a minor haplotype with no DRMs at all (**Figure 4d**, light
257 blue circles). Several inferred haplotypes without DRMs was nonetheless phylogenetically distinct
258 from the timepoint one minority haplotype (**Figure 4d**, compare small orange and pink circles in
259 lowest clade).

260
261 Upon examining the phylogenetic relationships of the inferred haplotype sequences, there were
262 several distinct clades with haplotypes from all timepoints interspersed throughout (except at
263 timepoint 4, which remained phylogenetically distinct). This is indicative of ongoing viral population
264 competition. DRMs showed some segregation by clade; viruses carrying a higher frequency of DRMs
265 were observed in upper cladea (Clade A, **Figure 4d**), and those with either K103N alone, or no DRMs
266 were preferentially located in the upper clade (Clade C, **Figure 4d**). However, this relationship was
267 not clear cut, and therefore consistent with competition between haplotypes during low drug
268 exposure. Soft sweeps were evident, given the increasing diversity (**Figure 1, Supplementary Figure**
269 **4**) of this participant, as well as constrained variant frequencies between 20-80% (**Figure 4b,c**).

270
271 **Participant 16207.** Viral load in this participant were consistently elevated >10,000 copies/ml (**Figure**
272 **5a**). As with participant 15664, drug concentrations in blood plasma remained extremely low or
273 absent at each measured timepoint, consistent with non-adherence to the prescribed regimen.
274 There was almost no change in the frequency of DRMs throughout the follow up period, even when
275 making the switch to the 2nd-line regimen. NNRTI resistance mutations such as K103N are known to
276 have minimal fitness costs²⁴ and can therefore persist in the absence of NNRTI pressure. Throughout
277 treatment the participant maintained K103N at a frequency of >95% but also carried several
278 integrase strand transfer inhibitor (INSTI) associated changes (E157Q) and PI-exposure associated

279 amino acid replacements (L23I and M46I) at low frequencies at timepoints two and three. Despite
280 little change in DRM site frequencies, very significant viral population shifts were observed at the
281 whole genome level, again indicative of selective sweeps (**Figures 5b-c**). Between timepoints one
282 and four, several linked mutations changed abundance contemporaneously, generating chevron-like
283 patterns of non-synonymous changes in *env* specifically (blue lines). A large number of alleles
284 increased in frequency from <20% to >80% at the same time as numerous others decreased in
285 frequency from above 80% to below 20%. Whereas large shifts in *gag* and *pol* alleles also occurred,
286 the mutations involved were almost exclusively synonymous (red and green lines). These putative
287 selective sweeps in *env* were evident in the phylogenetic analysis (**Figure 5d**, see long branch lengths
288 between timepoints one and four, and cladal structure) possibly driven by neutralising antibodies
289 and/or T-cell immune pressures.

290

291 Phylogenetic analysis of inferred whole genome haplotypes overall showed a distinct cladal
292 structure as observed in 15664 (**Figure 5d**), although the dominant haplotypes were equally
293 observed in the upper clade (A) and lower clade (C) (**Figure 5d**). K103N was the majority DRM at all
294 timepoints, except for a minority haplotype at timepoint three, also carrying E157Q. Haplotypes did
295 not cluster by time point. Significant diversity in haplotypes from this participant was confirmed by
296 MDS (**Supplementary Figure 6**).

297

298 **Participant 22763** was notable for a number of large shifts in variant frequencies across multiple
299 drug resistance associated residues and synonymous sites. Drug plasma concentration for different
300 drugs was variable yet detectable at most measured timepoints reflecting changing levels of
301 adherence across the treatment period (**Figure 6a**). Non-PI DRMs such as M184V, P225H and K103N
302 were present at baseline (time of switch from first to second line treatments). These mutations
303 persisted despite synonymous changes between time points one and two. Most of the highly
304 variable synonymous changes in this participant were found in the *gag* and *pol* genes (as in
305 participant 16207) (**Figure 6c**), but in this case *env* displayed large fluctuations in synonymous and
306 non-synonymous allelic frequencies over time. At timepoint three, therapeutic concentrations of
307 boosted lopinavir (LPV/r) and tenofovir (TDF) were measured in plasma and haplotypes clustered
308 separately from the first two timepoints (**Figure 6d**, light and dark grey circles). NGS confirmed that
309 the D67N, K219Q, K65R, L70R, M184V DRMs and NNRTI-resistance mutations were present at low
310 frequencies from timepoint three onwards. Of note, between timepoints three and six, therapeutic
311 concentration of TDF was detectable, and coincided with increased frequencies of the canonical TDF
312 DRM, K65R⁵. The viruses carrying K65R outcompeted those carrying the thymidine analogue mutants

313 (TAMs) D67N and K70R, whilst the lamivudine (3TC) associated resistance mutation, M184V,
314 persisted throughout. In the final three timepoints M46I emerged in *protease*, but never increased
315 in frequency above <6%. At timepoint seven, populations shifted again with some haplotypes
316 resembling those previously timepoint four, with D67N and K70R again being predominant over
317 K65R in *reverse transcriptase* (Figure 6d, green and blue circles). At the final timepoint (eight) the
318 frequency of K103N was approximately 85% and the TAM-bearing populations continued to
319 dominate over the K65R population, which at this timepoint had a low frequency.

320

321 Although the DRM profile suggested the possibility of a selective sweep, we observed the same
322 groups of other non-synonymous or synonymous alleles exhibiting dramatic frequency shifts, but to
323 a lesser degree than in the previous two participants i.e. 'chevron patterns' were less pronounced,
324 outside of the *env* gene (Figure 6b-c). Variable drug pressures placed on the viral populations
325 throughout the 2nd-line regimen appear to have played some role in limiting haplotype diversity.
326 Timepoints 1-4 all formed distinct clades, without intermingling, indicating that competition
327 between populations was not occurring to the same degree as in previous participants. Some
328 inferred haplotypes had K65R and others the TAMs D67N and K70R. K65R was not observed in
329 combination with D67N or K70R, consistent with previously reported antagonism between K65R and
330 TAMs whereby these mutations are not commonly found together within a single genome³⁸⁻⁴⁰. One
331 possible explanation for the disconnect between the trajectories of DRM frequencies over time and
332 haplotype phylogeny is competition between different viral populations. Alternatively, emergence of
333 haplotypes from previously unsampled reservoir with different DRM profile is possible, but one
334 might have expected other mutations to characterise such haplotypes that would manifest as
335 change in frequencies of large numbers of other mutations.

336

337 Discussion

338 The proportion of people living with HIV (PLWH) accessing ART has increased from 24% in 2010, to
339 68% in 2020^{41,42}. However, with the scale-up of ART, there has also been an increase in both pre-
340 treatment drug resistance (PDR)^{43,44} and acquired drug resistance^{12,45} to 1st-line ART regimens
341 containing NNRTIs. Integrase inhibitors (specifically dolutegravir) are now recommended for first-
342 line regimens by the WHO in regions where PDR exceeds 10%⁴⁶. Boosted PI-containing regimens
343 remain second line drugs following first 1st-line failure, though one unanswered question relates to
344 the nature of viral populations during failure on PI-based ART where major mutations in *protease*,
345 described largely for less potent PI, have not emerged. Here we have comprehensively analysed viral

346 populations present in longitudinally collected plasma samples of chronically-infected HIV-1
347 participants under non-suppressive 2nd-line ART.

348

349 With the vast majority of PLWH treated in the post-ART era, virus dynamics during non-suppressive
350 ART is important to understand, as there may be implications for future therapeutic success. For
351 example broadly neutralising antibodies (bNab) are being tested not only for prevention, but also as
352 part of remission strategies in combination with latency reversal agents. We know that HIV
353 sensitivity to broadly neutralising antibodies (bNab) is dependent on *env* diversity^{47,48}, and therefore
354 prolonged ART failure with viral diversification could compromise sensitivity to these agents.

355

356 Our understanding of virus dynamics largely stems from studies that were limited to untreated
357 individuals¹⁰, with largely subgenomic data analysed rather than whole genome¹⁰. Traditional
358 analysis of quasispecies distribution, for example as reported by Yu et al⁴⁹, suggests that the viral
359 diversity increases in longitudinal samples. However the findings of Yu et al were based entirely on
360 short-read NGS data without considering whole-genome haplotypes. The added benefit of
361 examining whole genome is that linked mutations can be identified statistically using an approach
362 that we recently developed²⁹. Indeed haplotype reconstruction has proved beneficial in the analysis
363 of compartmentalisation and diversification of several RNA viruses, including HIV-1, CMV and SARS-
364 CoV-2^{33,50,51}.

365

366 Key findings of this study were: firstly that diversity increased over time with variable trajectory
367 away from the consensus baseline sequence and also the reconstructed ancestral subtype C and M
368 consensus. Approximately half of the participants appeared to diversify away from the
369 reconstructed ancestral subtype C and M sequence, whereas three participants showed possible
370 reversion back towards the ancestral consensus C and M (albeit with insufficient statistical support).

371

372 Secondly, and in contrast to the fractions of synonymous and non-synonymous mutations reported
373 by Zanini et al in a longitudinal untreated dataset², we show that the fractions of synonymous
374 mutations are generally two-to-three fold higher than non-synonymous mutations during non-
375 suppressive ART in chronic infection. This finding may reflect early versus chronic infection and
376 differing selective pressures. Haplotype reconstruction revealed evidence for competing haplotypes,
377 with evidence for numerous soft selective sweeps in phylogenies, evidenced by intermingling of
378 haplotypes during periods where there was low drug concentration measured in participant's blood
379 plasma.

380

381 Individuals in the present study were treated with Ritonavir boosted Lopinavir along with two NRTIs
382 (typically Tenofovir + Emtricitabine). We observed significant change in the frequencies of NRTI
383 mutations in two of the three participants studied in-depth. These fluctuations likely reflected
384 adherence to the 2nd-line regimen though we saw evidence for possible archived virus populations
385 with DRMs emerging during follow-up because large changes in DRM frequency were not always
386 accompanied by changes at other sites. This is consistent with soft sweeps occurring and that non-
387 DRMs do not necessarily drift with other mutations to fixation²¹ and that the same mutations are
388 occurring on different backgrounds. As frequencies of RT DRMs did not always segregate with
389 haplotype frequencies, we suggest that a high number of recombination events, known to be
390 common in HIV infections, was responsible for the haplotypic diversity.

391

392 Although no participant developed major DRMS at consistently high frequencies to PIs
393 (<https://hivdb.stanford.edu/dr-summary/resistance-notes/PI/>), we did observe non-synonymous
394 mutations associated with PI exposure that are also known to be polymorphic; however, there was
395 no temporal evidence of specific changes being associated with selective sweeps. For example PI
396 exposure associated residues in matrix (positions 76 and 81) were observed in participant 16207
397 prior to PI initiation⁵². Furthermore, participant 16207 was one of few participants who achieved
398 two partial suppressions (<750 copies/ml). After both of these partial suppressions, the rebound
399 populations appeared to be less diverse, consistent with drug-resistant virus re-emerging.

400

401 Mutations in all genes that are further apart than 100bp are subject to shuffling via recombination⁵³.
402 Unlike the smooth LD decay curve as reported in the literature, we identified complex LD decay
403 patterns within patients, indicative of non-random recombination. Recombination appears as the
404 loss and gain of common genomic regions over successive timepoints between each participant's
405 haplotype populations (**Figure 3B**). Participant 15664 recombines between haplotype populations in
406 the *vif* and *vpr* genes in four of the six timepoints. In contrast, participant 22763 showed
407 recombination in the *gag-pol* genes in three of the eight timepoints. We explain these
408 recombination events in longitudinal sequences, as reflected in the previously discussed 'chevron'
409 patterns whereby variants increase and subsequently decrease between timepoints. HIV
410 quasispecies allows the virus to increase fitness through recombination when selectively
411 advantageous²⁶. The relationship between recombination and acquisition of DRMs is unclear with
412 each patient showing unique patterns; participant 16207 recombined in *pol* between haplotypes at
413 timepoint two and six and maintained the major DRM K103N. Participant 22763 recombined in *pol*

414 at three timepoints (two, four and six) resulting in no change of DRMs, gain of DRMs, and loss of
415 DRMs respectively. This occurring at the same time as antagonism between TAMs and DRMs (K65R
416 and D67N). Finally, participant 15664 steadily lost DRMs throughout longitudinal sampling, although
417 we did not see evidence of recombination driving this. This suggests that in the absence of strong
418 drug pressures, viral populations only maintained crucial DRMs which were useful to evading innate
419 immunity.

420

421 This study had some limitations – we examined in-detail only three participants with ongoing
422 viraemia and variable adherence to 2nd-line drug regimens. Despite the small sample size, this type
423 of longitudinal sampling of ART-experienced participants is unprecedented. We are confident that
424 the combination of computational analyses has provided a detailed understanding of viral dynamics
425 under non-suppressive ART may be applicable to wider datasets. The method used to reconstruct
426 viral haplotypes *in silico* is novel and has previously been validated in HIV-positive participants with
427 CMV⁵⁰. We are confident that the approach implemented by HaROLD has accurately, if
428 conservatively estimated haplotype frequencies and future studies should look to validate these
429 frequencies using an *in vitro* method such as single genome amplification. Despite there being high
430 viral loads present at each of the analysed timepoints, nuances of the sequencing method led in
431 some cases to suboptimal degrees of gene coverage, particularly in the *env* gene. To ensure that
432 uneven sequencing coverage did not bias our analyses, we ensured that variant analysis was only
433 performed where coverage was >10 reads.

434

435 In summary we have found compelling evidence of HIV-1 within-host viral diversification,
436 recombination and haplotype competition during non-suppressive ART. In future, participants
437 failing PI-based regimens are likely to be switched to INSTI-based ART (specifically Dolutegravir in
438 South Africa) prior to genotypic typing or resistance analysis. Although the prevalence of underlying
439 major INSTI resistance mutations is low in sub-Saharan Africa^{54,55}, this approach needs assessment
440 given data linking individuals with NNRTI resistance with poorer virological outcomes on
441 Dolutegravir⁵⁶, coupled with a history of intermittent adherence. Having shown that long-time intra-
442 host PI failure increases diversity of HIV viral populations, monitoring future drug-failure cases will
443 be of interest due to their capacity to maintain a reservoir of drug-resistant and transmissible virus.

444

445 **Methods**

446 **Study & Participant selection**

447 This cohort was nested within the French ANRS 12249 Treatment as Prevention (TasP) trial²⁷. TasP
448 was a cluster-randomised trial comparing an intervention arm who offered ART after HIV diagnosis
449 irrespective of participant CD4 + count, to a control arm which offered ART according to prevailing
450 South African guidelines. A subset of 44 longitudinal samples from eight chronically infected
451 participants. Participants were selected for examination if there were >3 timepoint samples
452 available. All samples were collected from blood plasma. The Illumina MiSeq platform was used and
453 an adapted protocol for sequencing⁵⁷. Adherence to 2nd-line regimens was measured by HPLC using
454 plasma concentration of drug levels as a proxy. Drug levels were measured at each timepoint with
455 detectable viral loads, post-PI initiation.

456

457 Ethical approval was originally grant by the Biomedical Research Ethics Committee (BFC 104/11) at
458 the University of KwaZulu-Natal, and the Medicines Control Council of South Africa for the TasP trial
459 (Clinicaltrials.gov: [NCT01509508](https://clinicaltrials.gov/ct2/show/study/NCT01509508); South African Trial Register: DOH-27-0512-3974). The study was
460 also authorized by the KwaZulu-Natal Department of Health in South Africa. Written informed
461 consent was obtained from all participants. Original ethical approval also included downstream
462 sequencing of blood plasma samples and analysis of those sequences to better understand drug
463 resistance. No additional ethical approval was required for this.

464

465 **Illumina Sequencing**

466 Sequencing of viral RNA was performed as previously described by Derache et al⁵⁸ using a modified
467 protocol previously described by Gall et al⁵⁹. Briefly, RNA was extracted from 1ml of plasma with
468 detectable viral load of >1000 copies/ml, using QIAamp Viral RNA mini kits (Qiagen, Hilden,
469 Germany), and eluted in 60µl of elution buffer. The near-full HIV genome was amplified with 4
470 subtype C primers pairs, generating 4 overlapping amplicons of between 2100 and 3900kb.

471

472 DNA concentrations of amplicons were quantified with the Qubit dsDNA HS Assay kit (Invitrogen,
473 Carlsbad, CA). Diluted amplicons were pooled equimolarly and prepared for library using the Nextera
474 XT DNA Library preparation and the Nextera XT DNA sample preparation index kits (Illumina, San
475 Diego, CA), following the manufacturer's protocol.

476

477 **Genomics & Bioinformatics**

478 Poor quality reads (with Phred score <30) and adapter sequences were trimmed from FastQ files
479 with TrimGalore! v0.6.519⁶⁰ and mapped to a clade C South African reference genome (AF411967)
480 with BWA-MEM⁶¹. The reference genome was manually annotated in Geneious Prime v2020.3 with

481 DRMs according to the Stanford HivDB ⁶². Optical PCR duplicate reads were removed using Picard
482 tools (<http://broadinstitute.github.io/picard>). Finally, QualiMap2 ⁶³ was used to assess the mean
483 mapping quality scores and coverage in relation to the reference genome for the purpose of
484 excluding poorly mapped sequences from further analysis. Single nucleotides polymorphisms (SNPs)
485 were called using VarScan2 ⁶⁴ with a minimum average quality of 20, minimum variant frequency of
486 2% and in at least 10 reads. These were then annotated by gene, codon and amino acid alterations
487 using an in-house script ⁶⁵ modified to utilise HIV genomes.

488

489 All synonymous variants and DRMs were examined, and their frequency compared across successive
490 timepoints. Synonymous variants were excluded from analysis if their prevalence remained at $\leq 10\%$
491 or $\geq 90\%$ across all timepoints. DRMs were retained for analysis if they were present at over 2%
492 frequency and on at least two reads. A threshold of 2% is supported by a study evaluating different
493 analysis pipelines, which reported fewer discordances over this cut-off ⁶⁶.

494

495 **Haplotype Reconstruction & Phylogenetics**

496 Whole-genome viral haplotypes were constructed for each participant timepoint using Haplotype
497 Reconstruction for Longitudinal Samples (HaROLD) ⁶⁷. Briefly, SNPs were assigned to each haplotype
498 such that the frequency of variants was equal to the sum of the frequencies of haplotypes
499 containing a specific variant. Maximal log likelihood was used to optimise time-dependent
500 frequencies for longitudinal haplotypes which was calculated by summing over all possible
501 assignment of haplotype variants. Haplotypes were then constructed based on posterior
502 probabilities. After constructing haplotypes, a refinement process remapped reads from BAM files to
503 those constructed haplotypes. The number of haplotypes either increased or decreased as a result
504 of combination or division according to AIC scores, in order to present the most accurate
505 representation of viral populations at each timepoint.

506

507 Whole-genome nucleotide diversity was calculated from BAM files using an in-house script
508 (<https://github.com/ucl-pathgenomics/NucleotideDiversity>). Briefly diversity is calculated by fitting
509 all observed variant frequencies to either a beta distribution or four-dimensional Dirichlet
510 distribution plus delta function (representing invariant sites). These parameters were optimised by
511 maximum log likelihood.

512

513 Maximum-likelihood phylogenetic trees and ancestral reconstruction were performed using IQTree2
514 v2.1.3⁶⁸ and a GTR+F+I model with 1000 ultrafast bootstrap replicates⁶⁹. All trees were visualised

515 with Figtree v.1.4.4 (<http://tree.bio.ed.ac.uk/software/figtree/>), rooted on the AF411967.3
516 reference sequence, and nodes arranged in descending order. Phylogenies were manipulated and
517 annotated using ggtree v2.2.4.

518

519 **Multi Dimension Scaling (MDS) Plots**

520 Pairwise distances between these consensus sequences were calculated using the `dist.dna()`
521 package, with a TN93 nucleotide-nucleotide substitution matrix and with pairwise deletion
522 implemented in the R package Ape v.5.4. Non-metric Multi-dimensional scaling (MDS) was
523 implemented using the `metaMDS()` function in the R package, `vegan` v2.5.7. MDS is a method to
524 attempt to simplify high dimensional data into a simpler representation of reducing dimensionality
525 whilst retaining most of the variation relationships between points. We find that like network trees,
526 non-metric MDS better represents the true relative distances between sequences, whereas
527 eigenvector methods are less reliable in this sense. In a genomics context we can apply
528 dimensionality reduction on pairwise distance matrices, where each dimension is a sequence with
529 data points of $n-1$ sequences pairwise distance. The process was repeated with whole genome
530 haplotype sequences.

531

532 **Linkage Disequilibrium & Recombination**

533 Starting with a sequence alignment we determined the pairwise LD R^2 associations for all variable
534 sites using `WeightedLD`⁷⁰. This method allowed us to exclude sites with any insertions or ambiguous
535 characters easily where we used the option `--min-acgt 0.99` and `--min-variability 0.05`. The pairwise
536 R^2 values were then binned per 200bp comparison distance blocks along the genome and the mean
537 R^2 value were taken and represented graphically to assess LD decay. This analysis was run for the
538 three participants taken forward for in-detail analysis, and run using an alignment of all their
539 timepoint samples. Graphics were generated using `Rv4.04`.

540

541 We first performed an analysis for detecting individual recombination events in individual genome
542 sequences using `RDP`, `GENECONV`, `BOOTSCAN`, `MAXCHI`, `CHIMAERA`, `SISCAN`, and `3SEQ` methods
543 implemented in `RDP5`⁷¹ using default settings. Putative breakpoint hotspots were identified and
544 manually checked and adjusted if necessary using the `BURT` method with the `MAXCHI` matrix and
545 `LARD` two breakpoint scan methods. Final recombination hotspots were confirmed if at least three
546 or more methods supported the breakpoint.

547

548 **Funding**

549 SAK is supported by the Bill and Melinda Gates Foundation: OPP1175094. RKG is supported by
550 Wellcome Trust Senior Fellowship in Clinical Science: WT108082AIA. OC is supported by a PhD
551 studentship/UKRI MRC grant : MR/N013867/1.

552

553 **Competing Interests**

554 RKG has received ad hoc consulting fees from Gilead, ViiV and UMOVIS Lab.

555

556 **Author Contributions**

557 Conceptualization of study: S.A.K, R.K.G, A.D, Bioinformatic processes: A.D, S.A.K, O.C, D.P.M,

558 Writing and revising manuscript: S.A.K, O.C, A.D, D.P.M, D.P, R.A.G, R.K.G

559

560 **Data Availability Statement**

561 All fasta files have been deposited on Genbank with the following accession numbers [pending
562 approval].

563

564 **Code Availability Statement**

565 Custom code used to produce figures and graphs can be found at: [https://github.com/Steven-](https://github.com/Steven-Kemp/21-2_hiv_tasp/tree/main/scripts)

566 [Kemp/21-2_hiv_tasp/tree/main/scripts](https://github.com/Steven-Kemp/21-2_hiv_tasp/tree/main/scripts) or within the references manuscripts.

567 **References**

568 1 Abrahams, M. R. *et al.* Quantitating the multiplicity of infection with human
569 immunodeficiency virus type 1 subtype C reveals a non-poisson distribution of transmitted
570 variants. *J Virol* **83**, 3556-3567, doi:10.1128/JVI.02132-08 (2009).

571 2 Zanini, F., Puller, V., Brodin, J., Albert, J. & Neher, R. A. In vivo mutation rates and the
572 landscape of fitness costs of HIV-1. *Virus Evol* **3**, vex003, doi:10.1093/ve/vex003 (2017).

573 3 Salemi, M. The intra-host evolutionary and population dynamics of human
574 immunodeficiency virus type 1: a phylogenetic perspective. *Infect Dis Rep* **5**, e3,
575 doi:10.4081/idr.2013.s1.e3 (2013).

576 4 Lemey, P., Rambaut, A. & Pybus, O. G. HIV evolutionary dynamics within and among hosts.
577 *Aids Reviews* **8**, 125-140 (2006).

578 5 Collier, D. A., Monit, C. & Gupta, R. K. The Impact of HIV-1 Drug Escape on the Global
579 Treatment Landscape. *Cell host & microbe* **26**, 48-60, doi:10.1016/j.chom.2019.06.010
580 (2019).

- 581 6 Biebricher, C. K. & Eigen, M. What is a quasispecies? *Curr Top Microbiol Immunol* **299**, 1-31,
582 doi:10.1007/3-540-26397-7_1 (2006).
- 583 7 Zanini, F. *et al.* Population genomics of inpatient HIV-1 evolution. *Elife* **4**, e11282,
584 doi:10.7554/eLife.11282 (2015).
- 585 8 Lythgoe, K. A. & Fraser, C. New insights into the evolutionary rate of HIV-1 at the within-host
586 and epidemiological levels. *Proceedings of the Royal Society B-Biological Sciences* **279**, 3367-
587 3375, doi:10.1098/rspb.2012.0595 (2012).
- 588 9 Hedskog, C. *et al.* Dynamics of HIV-1 Quasispecies during Antiviral Treatment Dissected
589 Using Ultra-Deep Pyrosequencing. *PLoS one* **5**, e11345, doi:ARTN e11345
590 10.1371/journal.pone.0011345 (2010).
- 591 10 Shankarappa, R. *et al.* Consistent viral evolutionary changes associated with the progression
592 of human immunodeficiency virus type 1 infection. *J Virol* **73**, 10489-10502,
593 doi:10.1128/JVI.73.12.10489-10502.1999 (1999).
- 594 11 Masikini, P. & Mpondo, B. C. HIV drug resistance mutations following poor adherence in HIV-
595 infected patient: a case report. *Clin Case Rep* **3**, 353-356, doi:10.1002/ccr3.254 (2015).
- 596 12 TenoRes Study, G. Global epidemiology of drug resistance after failure of WHO
597 recommended first-line regimens for adult HIV-1 infection: a multicentre retrospective
598 cohort study. *Lancet Infect Dis* **16**, 565-575, doi:10.1016/S1473-3099(15)00536-8 (2016).
- 599 13 Collier, D. *et al.* Virological Outcomes of Second-line Protease Inhibitor-Based Treatment for
600 Human Immunodeficiency Virus Type 1 in a High-Prevalence Rural South African Setting: A
601 Competing-Risks Prospective Cohort Analysis. *Clinical infectious diseases : an official
602 publication of the Infectious Diseases Society of America* **64**, 1006-1016,
603 doi:10.1093/cid/cix015 (2017).
- 604 14 Giandhari, J. *et al.* Genetic Changes in HIV-1 Gag-Protease Associated with Protease
605 Inhibitor-Based Therapy Failure in Pediatric Patients. *AIDS Res Hum Retroviruses* **31**, 776-
606 782, doi:10.1089/AID.2014.0349 (2015).
- 607 15 Kelly Pillay, S., Singh, U., Singh, A., Gordon, M. & Ndungu, T. Gag drug resistance mutations
608 in HIV-1 subtype C patients, failing a protease inhibitor inclusive treatment regimen, with
609 detectable lopinavir levels. *Journal of the International AIDS Society* **17**, 19784 (2014).

- 610 16 Sutherland, K. A. *et al.* Evidence for Reduced Drug Susceptibility without Emergence of
611 Major Protease Mutations following Protease Inhibitor Monotherapy Failure in the SARA
612 Trial. *PLoS one* **10**, e0137834, doi:10.1371/journal.pone.0137834 (2015).
- 613 17 Sutherland, K. A. *et al.* Phenotypic characterization of virological failure following
614 lopinavir/ritonavir monotherapy using full-length Gag-protease genes. *The Journal of*
615 *antimicrobial chemotherapy* **69**, 3340-3348, doi:10.1093/jac/dku296 (2014).
- 616 18 Sutherland, K. A. *et al.* Gag-Protease Sequence Evolution Following Protease Inhibitor
617 Monotherapy Treatment Failure in HIV-1 Viruses Circulating in East Africa. *AIDS research and*
618 *human retroviruses* **31**, 1032-1037, doi:10.1089/aid.2015.0138 (2015).
- 619 19 Day, C. L. *et al.* Proliferative capacity of epitope-specific CD8 T-cell responses is inversely
620 related to viral load in chronic human immunodeficiency virus type 1 infection. *Journal of*
621 *virology* **81**, 434-438, doi:10.1128/JVI.01754-06 (2007).
- 622 20 Blanch-Lombarte, O. *et al.* HIV-1 Gag mutations alone are sufficient to reduce darunavir
623 susceptibility during virological failure to boosted PI therapy. *The Journal of antimicrobial*
624 *chemotherapy* **75**, 2535-2546, doi:10.1093/jac/dkaa228 (2020).
- 625 21 Feder, A. F. *et al.* More effective drugs lead to harder selective sweeps in the evolution of
626 drug resistance in HIV-1. *Elife* **5**, e10670, doi:10.7554/eLife.10670 (2016).
- 627 22 Harris, R. B., Sackman, A. & Jensen, J. D. On the unfounded enthusiasm for soft selective
628 sweeps II: Examining recent evidence from humans, flies, and viruses. *PLoS genetics* **14**,
629 e1007859, doi:10.1371/journal.pgen.1007859 (2018).
- 630 23 Dam, E. *et al.* Gag mutations strongly contribute to HIV-1 resistance to protease inhibitors in
631 highly drug-experienced patients besides compensating for fitness loss. *PLoS pathogens* **5**,
632 e1000345 (2009).
- 633 24 Cong, M. E., Heneine, W. & Garcia-Lerma, J. G. The fitness cost of mutations associated with
634 human immunodeficiency virus type 1 drug resistance is modulated by mutational
635 interactions. *Journal of Virology* **81**, 3037-3041, doi:10.1128/Jvi.02712-06 (2007).
- 636 25 Wilke, C. O. Quasispecies theory in the context of population genetics. *BMC Evol Biol* **5**, 44,
637 doi:10.1186/1471-2148-5-44 (2005).

- 638 26 Lauring, A. S. & Andino, R. Quasispecies theory and the behavior of RNA viruses. *PLoS*
639 *Pathog* **6**, e1001005, doi:10.1371/journal.ppat.1001005 (2010).
- 640 27 Iwuji, C. C. *et al.* Evaluation of the impact of immediate versus WHO recommendations-
641 guided antiretroviral therapy initiation on HIV incidence: the ANRS 12249 TasP (Treatment
642 as Prevention) trial in Hlabisa sub-district, KwaZulu-Natal, South Africa: study protocol for a
643 cluster randomised controlled trial. *Trials* **14**, 230, doi:10.1186/1745-6215-14-230 (2013).
- 644 28 World Health Organization. *Consolidated guidelines on the use of antiretroviral drugs for*
645 *treating and preventing HIV infection: recommendations for a public health approach.*
646 (World Health Organization, 2016).
- 647 29 Pang, J. *et al.* Haplotype assignment of longitudinal viral deep-sequencing data using co-
648 variation of variant frequencies. *bioRxiv*, 444877, doi:10.1101/444877 (2020).
- 649 30 Cox, M. A. & Cox, T. F. in *Handbook of data visualization* 315-347 (Springer, 2008).
- 650 31 Stephens, M. & Scheet, P. Accounting for Decay of Linkage Disequilibrium in Haplotype
651 Inference and Missing-Data Imputation. *The American Journal of Human Genetics* **76**, 449-
652 462, doi:<https://doi.org/10.1086/428594> (2005).
- 653 32 Song, H. *et al.* Tracking HIV-1 recombination to resolve its contribution to HIV-1 evolution in
654 natural infection. *Nature Communications* **9**, 1928, doi:10.1038/s41467-018-04217-5 (2018).
- 655 33 Datir, R. *et al.* In Vivo Emergence of a Novel Protease Inhibitor Resistance Signature in HIV-1
656 Matrix. *mBio* **11**, e02036-02020, doi:10.1128/mBio.02036-20 (2020).
- 657 34 Kletenkov, K. *et al.* Role of Gag mutations in PI resistance in the Swiss HIV cohort study:
658 bystanders or contributors? *J Antimicrob Chemother* **72**, 866-875, doi:10.1093/jac/dkw493
659 (2017).
- 660 35 Rabi, S. A. *et al.* Multi-step inhibition explains HIV-1 protease inhibitor pharmacodynamics
661 and resistance. *The Journal of clinical investigation* **123**, 3848-3860, doi:10.1172/JCI67399
662 (2013).
- 663 36 Manasa, J. *et al.* Evolution of gag and gp41 in Patients Receiving Ritonavir-Boosted Protease
664 Inhibitors. *Sci Rep* **7**, 11559, doi:10.1038/s41598-017-11893-8 (2017).

- 665 37 Datir, R., El Bouzidi, K., Dakum, P., Ndembi, N. & Gupta, R. K. Baseline PI susceptibility by
666 HIV-1 Gag-protease phenotyping and subsequent virological suppression with PI-based
667 second-line ART in Nigeria. *The Journal of antimicrobial chemotherapy* **74**, 1402-1407,
668 doi:10.1093/jac/dkz005 (2019).
- 669 38 Parikh, U. M., Zelina, S., Sluis-Cremer, N. & Mellors, J. W. Molecular mechanisms of
670 bidirectional antagonism between K65R and thymidine analog mutations in HIV-1 reverse
671 transcriptase. *Aids* **21**, 1405-1414 (2007).
- 672 39 Parikh, U. M., Bachelier, L., Koontz, D. & Mellors, J. W. The K65R mutation in human
673 immunodeficiency virus type 1 reverse transcriptase exhibits bidirectional phenotypic
674 antagonism with thymidine analog mutations. *Journal of virology* **80**, 4971-4977 (2006).
- 675 40 Parikh, U. M., Barnas, D. C., Faruki, H. & Mellors, J. W. Antagonism between the HIV-1
676 reverse-transcriptase mutation K65R and thymidine-analogue mutations at the genomic
677 level. *The Journal of infectious diseases* **194**, 651-660 (2006).
- 678 41 Department of Health. 2019 ART Clinical Guidelines for the Management of HIV in Adults,
679 Pregnancy, Adolescents, Children, Infants and Neonates. (Republic of South Africa National
680 Department of Health, 2019).
- 681 42 UNAIDS. *Global HIV & AIDS statistics — 2020 fact sheet*,
682 <<https://www.unaids.org/en/resources/fact-sheet>> (2020), Accessed 3rd March 2021.
- 683 43 Gupta, R. K. *et al.* HIV-1 drug resistance before initiation or re-initiation of first-line
684 antiretroviral therapy in low-income and middle-income countries: a systematic review and
685 meta-regression analysis. *Lancet Infect Dis* **18**, 346-355, doi:10.1016/S1473-3099(17)30702-
686 8 (2018).
- 687 44 Gupta, R. K. *et al.* Global trends in antiretroviral resistance in treatment-naive individuals
688 with HIV after rollout of antiretroviral treatment in resource-limited settings: a global
689 collaborative study and meta-regression analysis. *Lancet* **380**, 1250-1258,
690 doi:10.1016/S0140-6736(12)61038-1 (2012).
- 691 45 Gregson, J. *et al.* Occult HIV-1 drug resistance to thymidine analogues following failure of
692 first-line tenofovir combined with a cytosine analogue and nevirapine or efavirenz in sub
693 Saharan Africa: a retrospective multi-centre cohort study. *Lancet Infect Dis*,
694 doi:10.1016/S1473-3099(16)30469-8 (2017).

- 695 46 WHO, C. Global Fund. HIV drug resistance report. 2017. *World Health Organisation* (2017).
- 696 47 Stefic, K., Bouvin-Pley, M., Braibant, M. & Barin, F. Impact of HIV-1 Diversity on Its Sensitivity
697 to Neutralization. *Vaccines (Basel)* **7**, 74, doi:10.3390/vaccines7030074 (2019).
- 698 48 Pancera, M. *et al.* Structure and immune recognition of trimeric pre-fusion HIV-1 Env.
699 *Nature* **514**, 455-461, doi:10.1038/nature13808 (2014).
- 700 49 Yu, F. *et al.* The Transmission and Evolution of HIV-1 Quasispecies within One Couple: a
701 Follow-up Study based on Next-Generation Sequencing. *Scientific reports* **8**, 1404,
702 doi:10.1038/s41598-018-19783-3 (2018).
- 703 50 Pang, J. *et al.* Mixed cytomegalovirus genotypes in HIV-positive mothers show
704 compartmentalization and distinct patterns of transmission to infants. *Elife* **9**, e63199,
705 doi:10.7554/eLife.63199 (2020).
- 706 51 Boshier, F. A. T. *et al.* Remdesivir induced viral RNA and subgenomic RNA suppression, and
707 evolution of viral variants in SARS-CoV-2 infected patients. *medRxiv*,
708 2020.2011.2018.20230599, doi:10.1101/2020.11.18.20230599 (2020).
- 709 52 Parry, C. M. *et al.* Three residues in HIV-1 matrix contribute to protease inhibitor
710 susceptibility and replication capacity. *Antimicrobial agents and chemotherapy* **55**, 1106-
711 1113, doi:10.1128/AAC.01228-10 (2011).
- 712 53 Neher, R. A. & Leitner, T. Recombination rate and selection strength in HIV intra-patient
713 evolution. *PLoS Comput Biol* **6**, e1000660, doi:10.1371/journal.pcbi.1000660 (2010).
- 714 54 El Bouzidi, K. *et al.* High prevalence of integrase mutation L74I in West African HIV-1
715 subtypes prior to integrase inhibitor treatment. *J Antimicrob Chemother* **75**, 1575-1579,
716 doi:10.1093/jac/dkaa033 (2020).
- 717 55 Derache, A. *et al.* Predicted antiviral activity of tenofovir versus abacavir in combination with
718 a cytosine analogue and the integrase inhibitor dolutegravir in HIV-1-infected South African
719 patients initiating or failing first-line ART. *The Journal of antimicrobial chemotherapy*,
720 doi:10.1093/jac/dky428 (2018).
- 721 56 Siedner, M. J. *et al.* Reduced efficacy of HIV-1 integrase inhibitors in patients with drug
722 resistance mutations in reverse transcriptase. *Nat Commun* **11**, 5922, doi:10.1038/s41467-
723 020-19801-x (2020).

- 724 57 Iwuji, C. *et al.* Universal test and treat is not associated with sub-optimal antiretroviral
725 therapy adherence in rural South Africa: the ANRS 12249 TasP trial. *J Int AIDS Soc* **21**,
726 e25112, doi:10.1002/jia2.25112 (2018).
- 727 58 Derache, A. *et al.* Impact of Next-generation Sequencing Defined Human Immunodeficiency
728 Virus Pretreatment Drug Resistance on Virological Outcomes in the ANRS 12249 Treatment-
729 as-Prevention Trial. *Clinical infectious diseases : an official publication of the Infectious*
730 *Diseases Society of America* **69**, 207-214, doi:10.1093/cid/ciy881 (2019).
- 731 59 Gall, A. *et al.* Universal amplification, next-generation sequencing, and assembly of HIV-1
732 genomes. *Journal of clinical microbiology* **50**, 3838-3844, doi:10.1128/JCM.01516-12 (2012).
- 733 60 Martin, M. J. E. j. Cutadapt removes adapter sequences from high-throughput sequencing
734 reads. **17**, pp. 10-12 (2011).
- 735 61 Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv*
736 *preprint arXiv:1303.3997* (2013).
- 737 62 Shafer, R. W. Rationale and uses of a public HIV drug-resistance database. *The Journal of*
738 *infectious diseases* **194 Suppl 1**, S51-58, doi:10.1086/505356 (2006).
- 739 63 Okonechnikov, K., Conesa, A. & Garcia-Alcalde, F. Qualimap 2: advanced multi-sample
740 quality control for high-throughput sequencing data. *Bioinformatics (Oxford, England)* **32**,
741 292-294, doi:10.1093/bioinformatics/btv566 (2016).
- 742 64 Koboldt, D. C. *et al.* VarScan 2: somatic mutation and copy number alteration discovery in
743 cancer by exome sequencing. *Genome Res* **22**, 568-576, doi:10.1101/gr.129684.111 (2012).
- 744 65 Charles, O. J., Venturini, C. & Breuer, J. cmvdr - An R package for Human Cytomegalovirus
745 antiviral Drug Resistance Genotyping. *bioRxiv*, 2020.2005.2015.097907,
746 doi:10.1101/2020.05.15.097907 (2020).
- 747 66 Perrier, M. *et al.* Evaluation of different analysis pipelines for the detection of HIV-1 minority
748 resistant variants. *PloS one* **13**, e0198334, doi:10.1371/journal.pone.0198334 (2018).
- 749 67 Goldstein, R. A., Tamuri, A. U., Roy, S. & Breuer, J. Haplotype assignment of virus NGS data
750 using co-variation of variant frequencies. *bioRxiv*, 444877 (2018).

751 68 Minh, B. Q. *et al.* IQ-TREE 2: New models and efficient methods for phylogenetic inference in
752 the genomic era. *bioRxiv*, 849372, doi:10.1101/849372 (2019).

753 69 Minh, B. Q., Nguyen, M. A. & von Haeseler, A. Ultrafast approximation for phylogenetic
754 bootstrap. *Mol Biol Evol* **30**, 1188-1195, doi:10.1093/molbev/mst024 (2013).

755 70 Charles, O. J., Roberts, J., Breuer, J. & Goldstein, R. A. WeightedLD: The Application of
756 Sequence Weights to Linkage Disequilibrium. *bioRxiv*, 2021.2006.2004.447093,
757 doi:10.1101/2021.06.04.447093 (2021).

758 71 Martin, D. P. *et al.* RDP5: a computer program for analyzing recombination in, and removing
759 signals of recombination from, nucleotide sequence datasets. *Virus Evol* **7**, veaa087,
760 doi:10.1093/ve/veaa087 (2021).

761

762

763

764

765

766

767

768

769

770

771

772

773

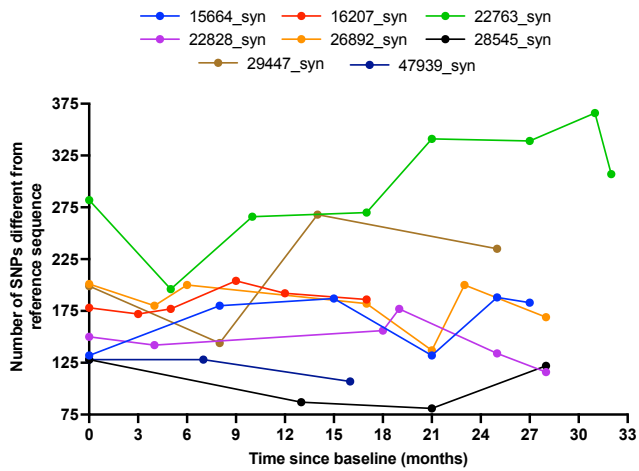
774 **Table 1.** Regimens and viral load at final timepoint for all participants. Participants initiated and
775 maintained 1st-line regimens for between 1-10 years before being switched to 2nd-line regimens as
776 part of the TasP trial. Eight of the nine participants were failing 2nd-line regimens at the final
777 timepoint.

778

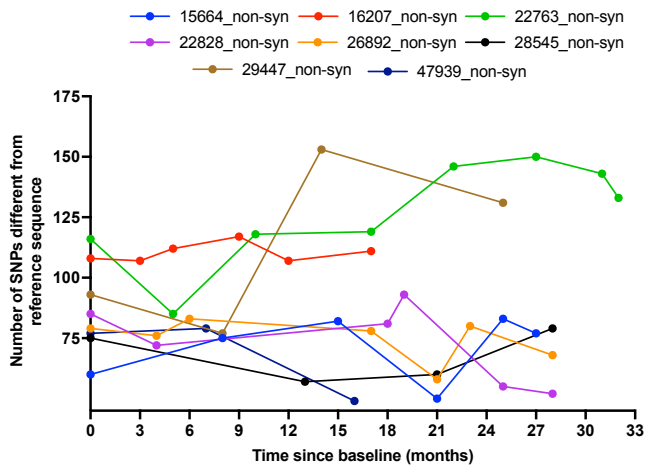
Participant	No. of timepoints	1st-line regimen	Time since initiation of 1 st -line treatment (yrs.)	2 nd -line regimen	Viral Load at final timepoint (copies/ml)
15664	6	d4T, 3TC, FTC	6.2	LPV/r, TDF, FTC	28655
16207	5	d4T, 3TC,	5.9	LPV/r, TDF, FTC	56660

		NVP			
22763	8	d4T, 3TC, EFV	6.2	LPV/r, TDF, 3TC	15017
22828	6	d4T, 3TC, NVP	6.4	LPV/r, TDF, 3TC/FTC	947
26892	7	d4T, 3TC, EFV	6	LPV/r, TDF, FTC	12221
28545	5	TDF, FTC, EFV	1.3	LPV/r, AZT, 3TC	12964
28841	2	-	-	-	199011
29447	4	TDF, FTC, EFV	2.8	LPV/r, TDF, FTC	64362
47939	3	d4T, 3TC, EFV	10.1	LPV/r, AZT, 3TC/FTC	6328
779	NRTI: Stavudine, d4T; Lamivudine, 3TC; Tenofovir, TDF; Emtricitabine, FTC; Zidovudine, AZT. NNRTI:				
780	Efavirenz, EFV; Nevirapine, NVP. PI: Lopinavir/ritonavir, LPV/r. “-”, data missing.				
781					
782					

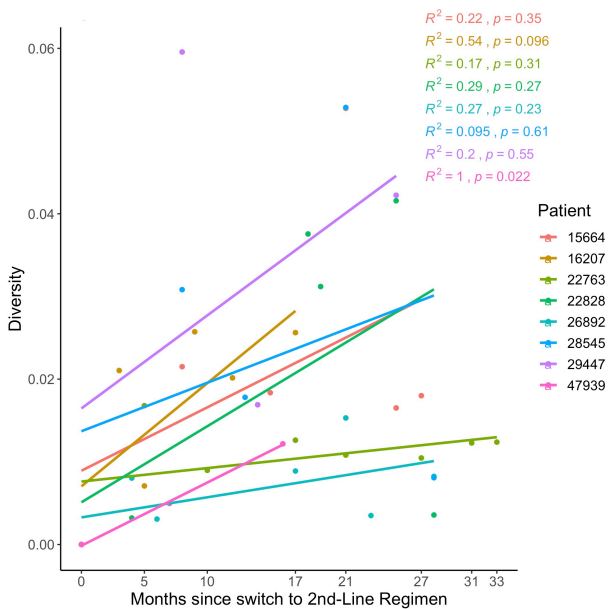
a)



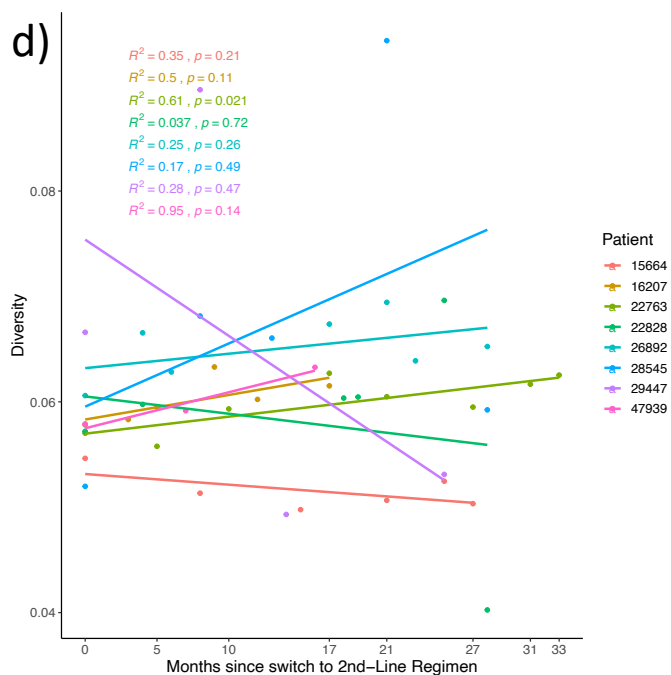
b)



c)



d)



e)

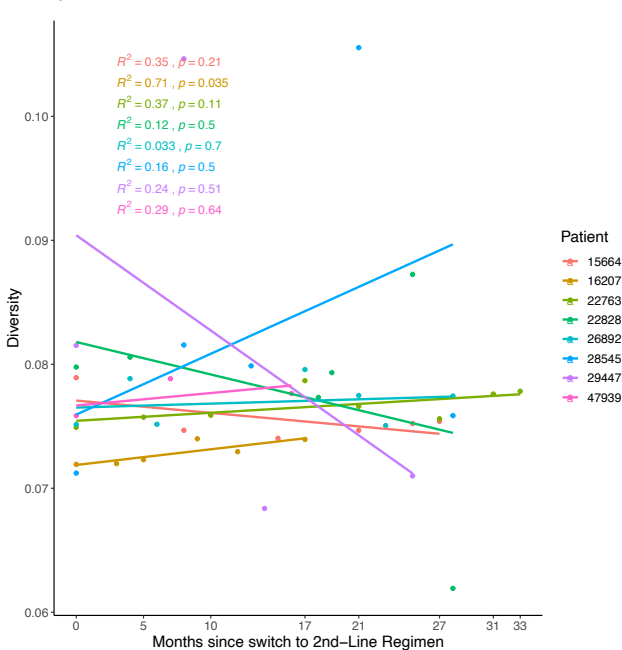
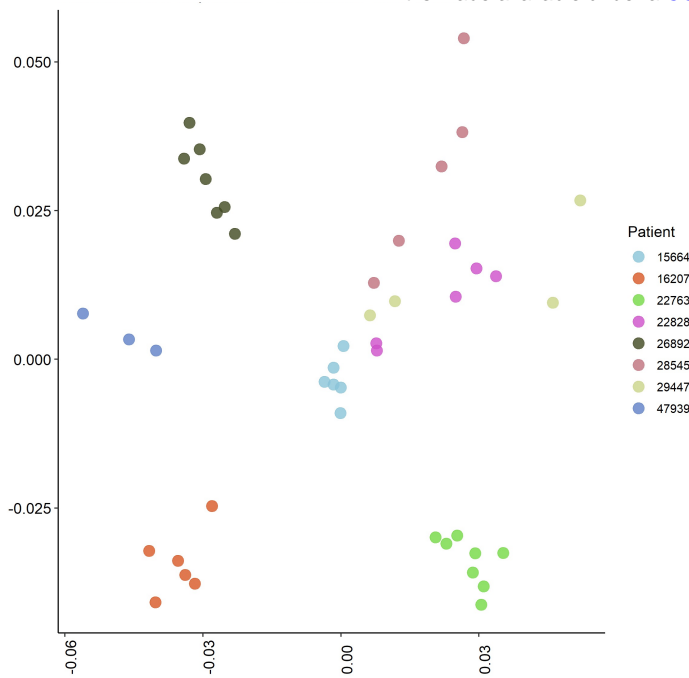


Figure 1. Measure of sequence divergence for eight participants under non-suppressive ART, relative to the subtype C reference strain at successive timepoints. These data were for SNPs detected by Illumina NGS at <2% abundance. Sites had coverage of at least 10 reads. In both a) synonymous and b) non-synonymous mutations, there was idiosyncratic change in number of SNPs relative to the reference strain over time. **1c-e) Linear regression of average pairwise distance relative to C) The baseline timepoint, D) a reconstructed subtype C consensus and E) a reconstructed subtype M consensus.** Average pairwise distances were estimated under a TN93 substitution model and reveal divergence from the initial samples. The Ancestral/Consensus HIV-1 subtype C and M was downloaded from the Los Alamos National Laboratory (<https://www.hiv.lanl.gov/content/sequence/NEWALIGN/align.html>). All ancestral HIV-1 subtype were downloaded from the same alignment and a consensus was created, as a proxy for an ancestral HIV-1 group M sequence. R^2 and p -values for logistic regression fits are indicated.

a)



b)

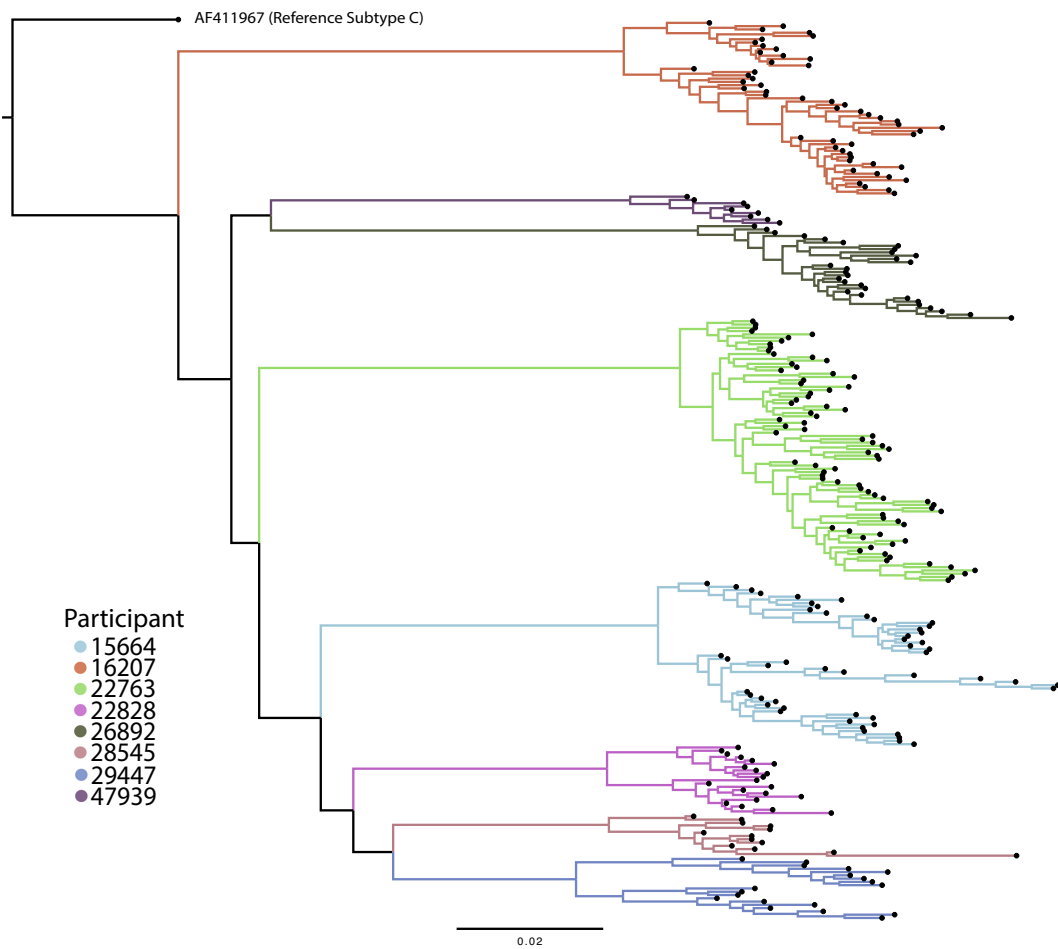
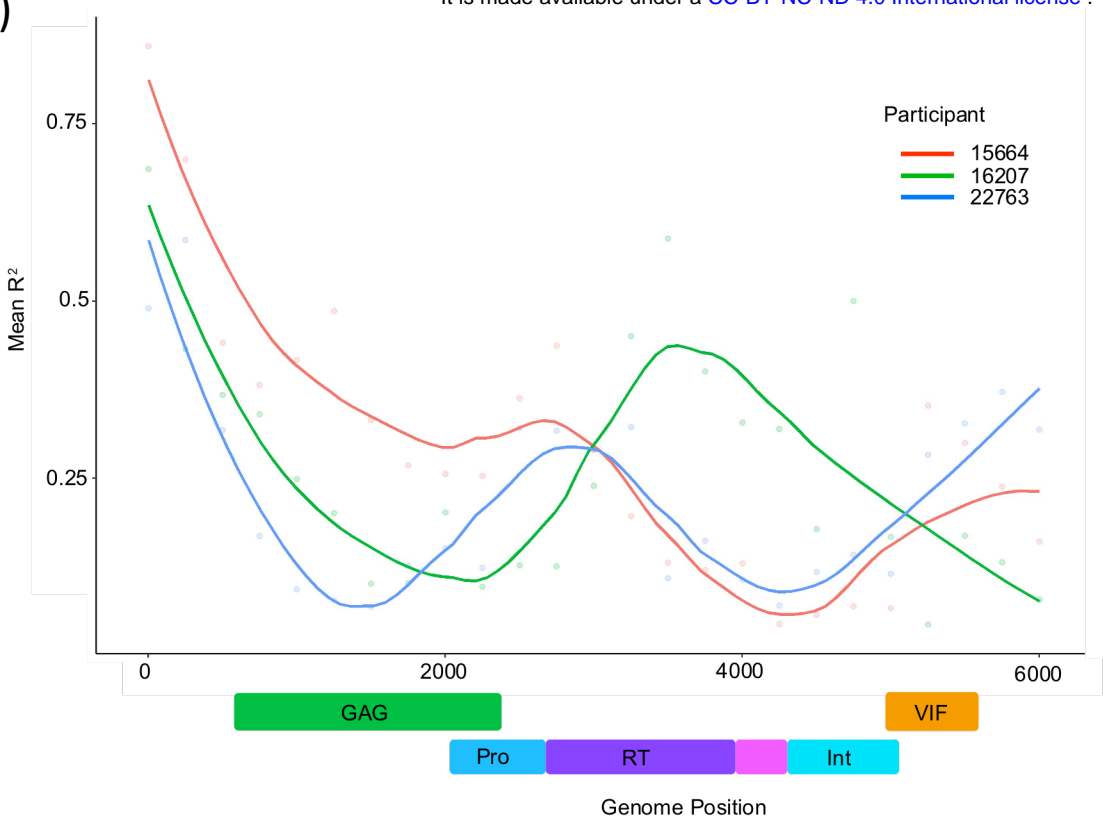


Figure 2. Multi-dimensional scaling showing A) clustering of HIV whole genomes from consensus sequences with high intra-participant diversity. Multi-dimensional scaling (MDS) were created by determining all pairwise distance comparisons under a TN93 substitution model, coloured by participant. Axis are MDS-1 and MDS-2. **B) Maximum likelihood phylogeny of constructed viral haplotypes for all participants.** The phylogeny was rooted on the AF411967 clade C reference genome. Reconstructed haplotypes were genetically diverse and did not typically cluster by timepoint.

A)



B)

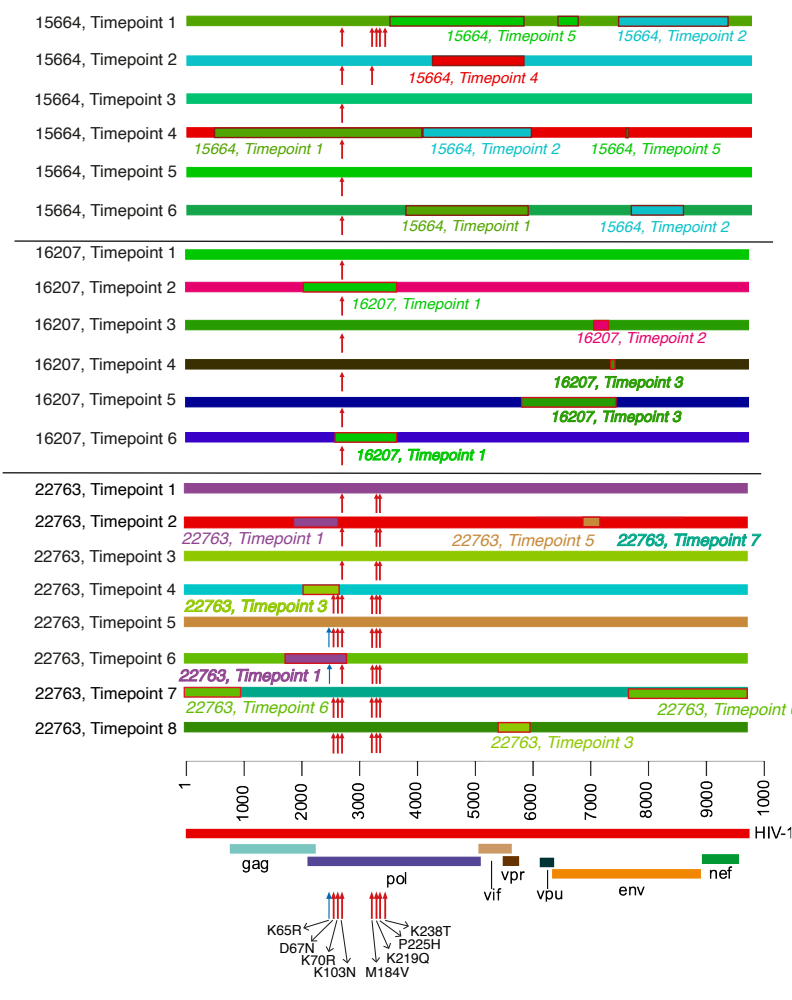


Figure 3A) Pairwise linkage disequilibrium decays rapidly with increasing distance between SNPs. The line indicates the average LD of all eight patients. There was a constant decrease in linkage disequilibrium over the first 800bp. **B) Perceived recombination breakpoints and drug-resistance associated mutations of all longitudinal consensus sequences belonging to three participants: 15664, 16207 and 22763.** All sequences were coloured uniquely uniquely; perceived recombination events supported by 4 or more methods implemented in RDP5 are highlighted with a red border and italic text to show the major parent and recombinant portion of the sequence. Drug-resistance associated mutations are indicated with a red arrow, relative to the key at the bottom of the image. For ease of distinguishment, the K65R mutations is indicated with a blue arrow.

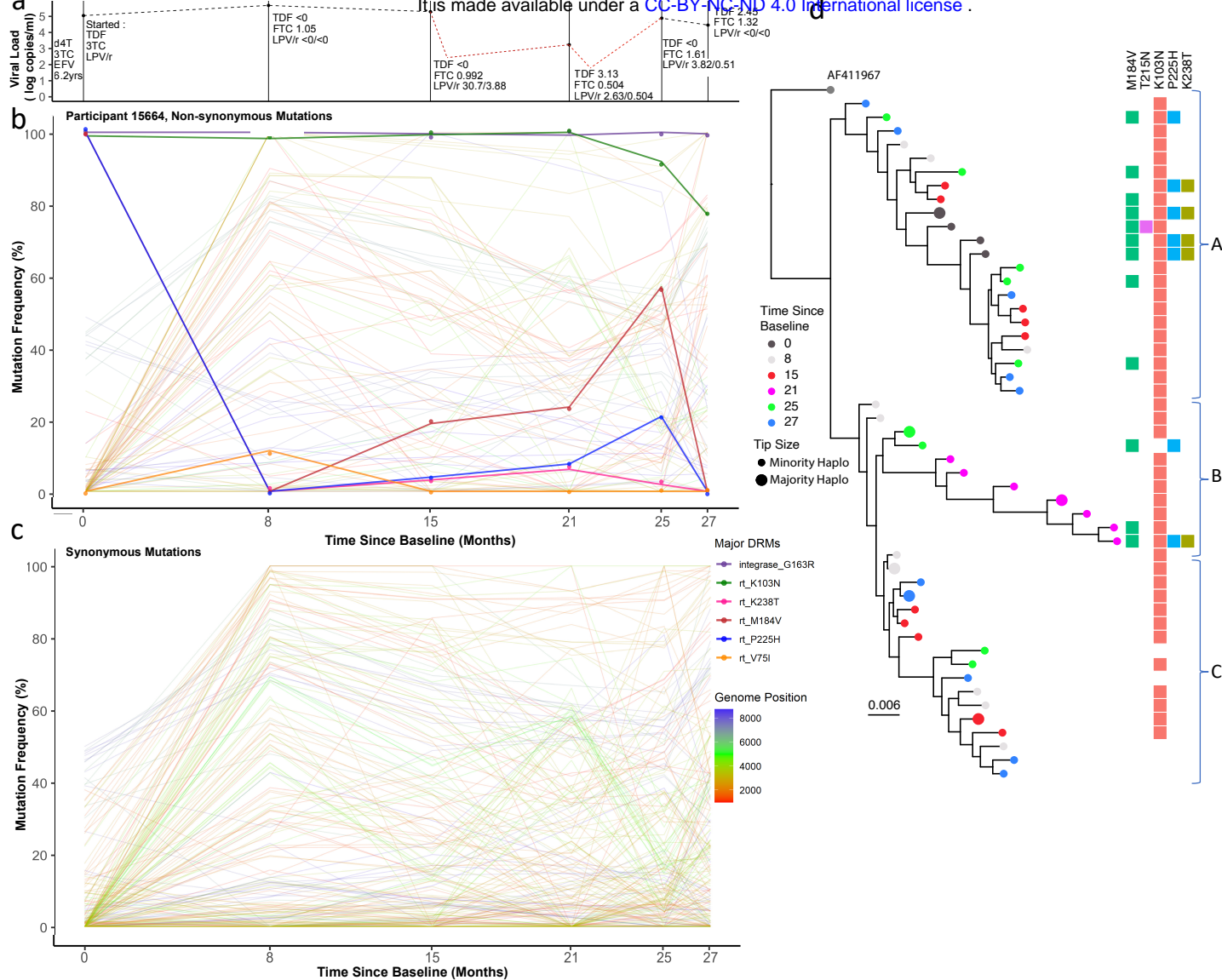


Figure 4. Drug regimen, adherence and viral dynamics within participant 15664. a) Viral load and drug levels. At successive timepoints drug regimen was noted and plasma drug concentration measured by HPLC (nmol/l). The participant was characterised by multiple partial suppression (<750 copies/ml, 16 months; <250 copies/ml, 22 months) and rebound events (red dotted line) and poor adherence to the drug regimen. **b) Drug resistance and non-drug resistance associated non-synonymous mutation frequencies by Illumina NGS.** The participant had large population shifts between timepoints 1-2, consistent with a hard selective sweep, coincident with the shift from 1st-line regimen to 2nd-line. **c) Synonymous mutation frequencies.** All mutations with a frequency of <10% or >90% at two or more timepoints were tracked over successive timepoints. Most changes were restricted to *gag* and *pol* regions and had limited shifts in frequency i.e. between 20-60%. **d) Maximum-likelihood phylogeny of reconstructed haplotypes.** Haplotypes largely segregated into three major clades (labelled A-C). Majority and minority haplotypes, some carrying lamivudine resistance mutation M184V. Clades referred to in the text body are shown to the right of the heatmap.

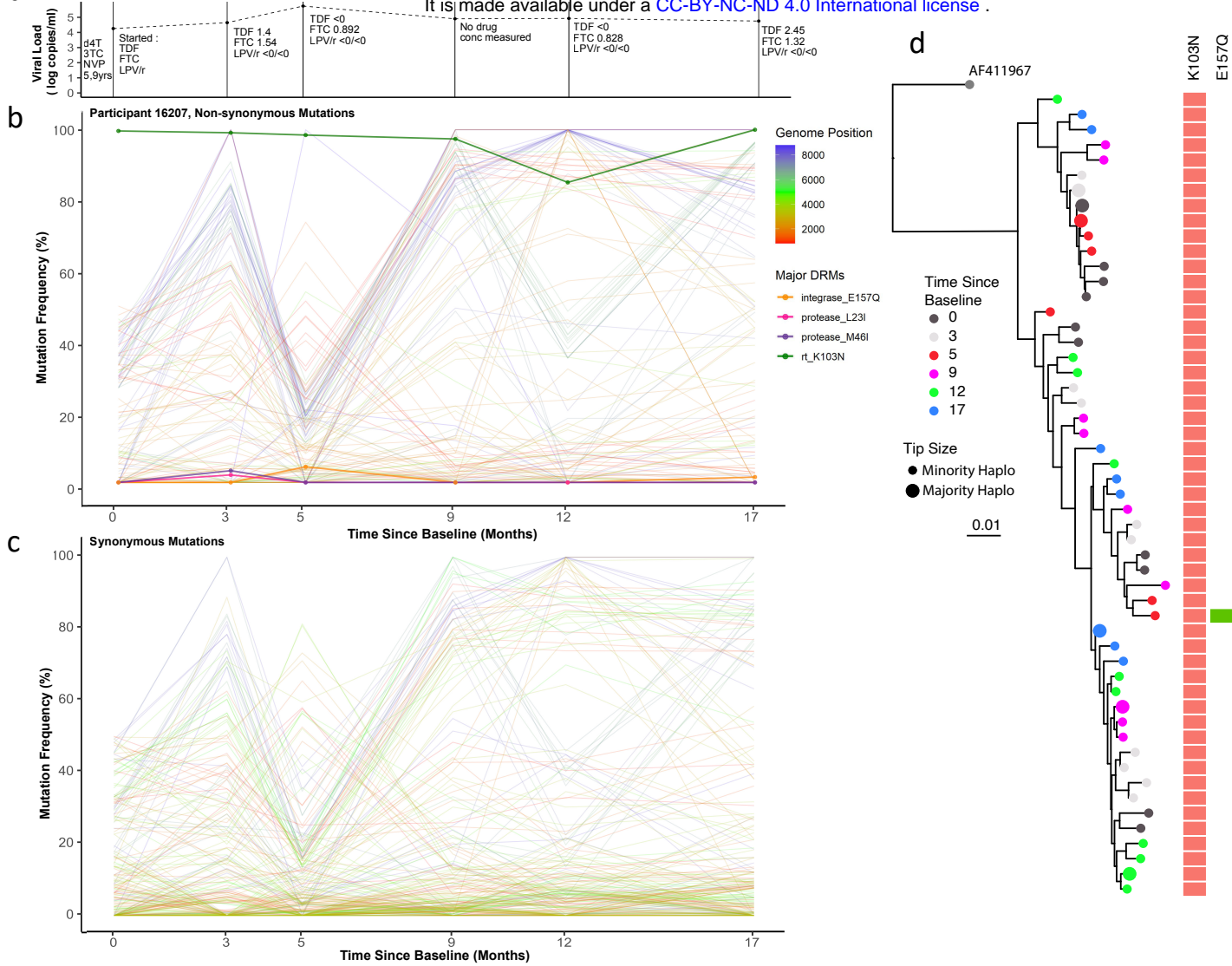


Figure 5. Drug regimen, adherence and viral dynamics within participant 16207. A) Viral load and drug levels. At successive timepoints regimen was noted and plasma drug concentration measured by HPLC (nmol/l). The participant displayed ongoing viraemia and poor adherence to the prescribed drug regimen. **B) Drug resistance and non-drug resistance associated non-synonymous mutations frequencies.** The participant had only one major RT mutation - K103N for the duration of the treatment period. Several antagonistic non-synonymous switches in predominantly *env* were observed between timepoints 1-4. **C) Synonymous mutation frequencies.** All mutations with a frequency of <10% or >90% at two or more timepoints were followed over successive timepoints. In contrast to non-synonymous mutations, most synonymous changes were in *pol*, indicative of linkage to the *env* coding changes. **D) Maximum-likelihood phylogeny of reconstructed haplotypes.** Haplotypes were again clearly divided into three distinct clades; each clade contained haplotypes from all timepoints, suggesting lack of hard selective sweeps and intermingling of viral haplotypes with softer sweeps. that most viral competition occurred outside of drug pressure.

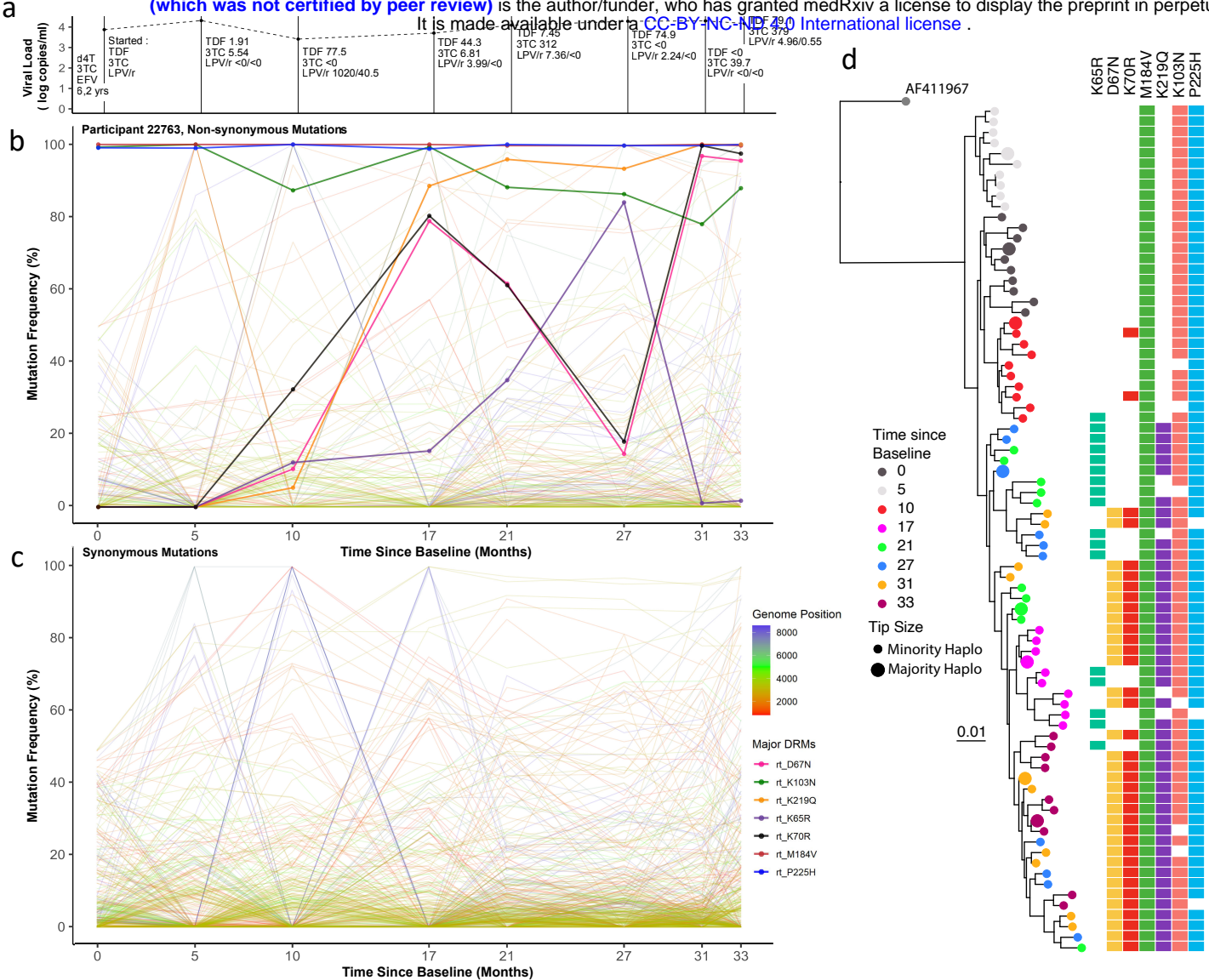


Figure 6. Drug regimen, adherence and viral dynamics of participant 22763. A) Viral load and regimen adherence. At successive timepoints the regimen was noted, and plasma drug concentration measured by HPLC (nmol/l). The participant had therapeutic levels of drug at several timepoints (3, 5 and 8), indicating variable adherence to the prescribed drug regimen. **B) Drug resistance and non-drug-resistance-associated non-synonymous mutation frequencies.** The participant had numerous drug resistance mutations in dynamic flux. Between timepoints 4-7, there was a complete population shift, indicated by reciprocal competition between the RT mutations K65R and the TAMs K67N and K70R. **C) Synonymous mutations frequencies.** All mutations with a frequency of <10% or >90% at two or more timepoints were followed over successive timepoints. Several *env* mutations mimicked the non-synonymous shifts observed between timepoints 2-4, suggestive of linkage. **D) Maximum-likelihood phylogeny of reconstructed haplotypes.** timepoints 1-4 were found in distinct lineages. In later timepoints, from 5-8, haplotypes became more intermingled, whilst maintaining antagonism between K65R and K67N bearing viruses.