



**HAL**  
open science

# Impact of CCR5delta32 host genetic background and disease progression on HIV-1 intra-host evolutionary processes: efficient hypothesis testing through hierarchical phylogenetic models

Diana Edo-Matas, Philippe Lemey, Jennifer Tom, Cèlia Serna-Bolea, Agnes E. van den Blink, Angélique B. van 'T Wout, Hanneke Schuitemaker, Marc Suchard

## ► To cite this version:

Diana Edo-Matas, Philippe Lemey, Jennifer Tom, Cèlia Serna-Bolea, Agnes E. van den Blink, et al.. Impact of CCR5delta32 host genetic background and disease progression on HIV-1 intra-host evolutionary processes: efficient hypothesis testing through hierarchical phylogenetic models. *Molecular Biology and Evolution*, 2010, 10.1093/molbev/MSQ326 . hal-00648534

**HAL Id: hal-00648534**

**<https://hal.science/hal-00648534>**

Submitted on 6 Dec 2011

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 **Impact of CCR5delta32 host genetic background and disease progression on**  
2 **HIV-1 intra-host evolutionary processes: efficient hypothesis testing through**  
3 **hierarchical phylogenetic models**  
4

5 Research article

6  
7 Diana Edo-Matas<sup>1</sup>, Philippe Lemey<sup>2,†</sup>, Jennifer A. Tom<sup>3</sup>, Cèlia Serna-Bolea<sup>1,a</sup>, Agnes E. van den  
8 Blink<sup>1</sup>, Angélique B. van 't Wout<sup>1,b</sup>, Hanneke Schuitemaker<sup>1,b</sup>, Marc A. Suchard<sup>3,4</sup>  
9

10 <sup>1</sup> Department of Experimental Immunology, Sanquin Research, Landsteiner Laboratory, Center  
11 for Infection and Immunity Amsterdam (*CINIMA*) at the Academic Medical Center of the  
12 University of Amsterdam, Amsterdam, The Netherlands

13 <sup>2</sup> Rega Institute for Medical Research, Leuven, Belgium

14 <sup>3</sup> Department of Biostatistics, School of Public Health; University of California, Los Angeles, CA  
15 90095, USA

16 <sup>4</sup> Departments of Biomathematics and Human Genetics, David Geffen School of Medicine,  
17 University of California, Los Angeles, CA, 90095, USA

18 <sup>a</sup> Present address: Barcelona Centre for International Health Research, Hospital Clinic, Institut  
19 d'Investigacions Biomediques August Pi i Sunyer, Universtitat de Barcelona, Barcelona, Spain.

20 <sup>b</sup> Present address: Crucell Holland BV, Leiden, The Netherlands.

21

22

23 <sup>†</sup>Correspondence should be addressed to:

24 Philippe Lemey

25 Rega Institute, Minderbroedersstraat 10, 3000 Leuven, Belgium

26 Tel: +32 (0)16 332160

27 Fax: +32 (0)16 332131

28 Email: philippe.lemey@uz.kuleuven.be

29

30 Key words:

31 - CCR5

32 - Envelope

33 - HIV-1

34 - Hierarchical Phylogenetic models

35 - Disease progression

36 - Bayesian inference

37

38 Running head: HIV-1 Hierarchical Phylogenetic Hypothesis Testing

39

40

40 **ABSTRACT**

41 The interplay between CCR5 host genetic background, disease progression and intra-host HIV-1  
42 evolutionary dynamics remains unclear because differences in viral evolution between hosts limit  
43 the ability to draw conclusions across hosts stratified into clinically relevant populations. Similar  
44 inference problems are proliferating across many measurably evolving pathogens for which intra-  
45 host sequence samples are readily available. To this end, we propose novel hierarchical  
46 phylogenetic models (HPMs) that incorporate fixed-effects to test for differences in dynamics  
47 across host-populations in a formal statistical framework employing stochastic search variable  
48 selection and model averaging. To clarify the role of CCR5 host genetic background and disease  
49 progression on viral evolutionary patterns, we obtain *gp120* envelope sequences from clonal  
50 HIV-1 variants isolated at multiple time points in the course of infection from populations of  
51 HIV-1 infected individuals who only harbored CCR5-using HIV-1 variants at all time points.  
52 Presence or absence of a CCR5 wt/ $\Delta$ 32 genotype and progressive or long-term non-progressive  
53 course of infection stratify the clinical populations in a two-way design. As compared to the  
54 standard approach of analyzing sequences from each patient independently, the HPM provides  
55 more efficient estimation of evolutionary parameters such as nucleotide substitution rates and  
56  $d_N/d_S$  rate ratios, as shown by significant shrinkage of the estimator variance. The fixed-effects  
57 also corrects for non-independence of data between populations and results in even further  
58 shrinkage of individual patient estimates. Model selection suggests an association between  
59 nucleotide substitution rate and disease progression, but a role for CCR5 genotype remains  
60 elusive. Given the absence of clear  $d_N/d_S$  differences between patient groups, delayed onset of  
61 AIDS symptoms appears to be solely associated with lower viral replication rates rather than with  
62 differences in selection on amino acid fixation.

63

## 63 INTRODUCTION

64 The high mutation rate and rapid viral turnover that characterize HIV-1 infection (Ho et al. 1995;  
65 Wei et al. 1995) generate a highly diverse genetic viral population within an HIV-1 infected  
66 individual (Shankarappa et al. 1999). Continuous emergence of new HIV-1 variants facilitates  
67 rapid viral adaptation to humoral and cellular immune responses of the host (Borrow et al. 1997;  
68 Goulder et al. 1997; Wei et al. 2003; Jones et al. 2004), escape from antiretroviral drugs (Coffin  
69 1995) and the selection for optimal biological properties such as replication capacity and use of  
70 the entry complex (Koning et al. 2003; Kwa et al. 2003) (Sterjovski et al. 2007; Repits et al.  
71 2008).

72 Following primary infection, an asymptomatic phase with a gradual loss of CD4<sup>+</sup> T cells and T-  
73 cell function characterizes the clinical course of HIV-1 infection (Lane et al. 1985; Polk et al.  
74 1987; Miedema et al. 1988), resulting eventually in the development of AIDS. The duration of  
75 this asymptomatic phase in the absence of antiretroviral therapy varies among patients, from  
76 several months to more than two decades, and determines their rate of disease progression  
77 (Veugelers et al. 1994; Munoz, Sabin, and Phillips 1997). Many selective forces may play a role  
78 in intra-host viral evolution and disease progression such as neutralizing antibodies (nAbs) and  
79 cytotoxic T cell (CTL) response, immune activation, target cell availability, co-receptor  
80 expression levels and emergence of CXCR4-using viruses among others. The severity of HIV  
81 infection may be further complicated by co-infections and heritable viral genetic factors  
82 (Hollingsworth et al. 2010). Largely stimulated by a comprehensive longitudinal analysis  
83 demonstrating common patterns of sequence divergence, diversity and emergence of CXCR4-  
84 using variants in chronic HIV-1 infections (Shankarappa et al. 1999), phylogenetic analyses have  
85 been widely used as a means of elucidating how host factors impact HIV within-host dynamics.  
86 More specific evolutionary parameters such as evolutionary rate (Lemey et al. 2007; Lee et al.

87 2008), adaptation rates (Williamson 2003), positively selected sites (Ross and Rodrigo 2002),  
88 compartmentalization (Kemal et al. 2003) and recombination (Carvajal-Rodriguez et al. 2008)  
89 have been scrutinized, but consistent associations with disease progression have rarely been  
90 revealed.

91 Here, we focus on a polymorphism in the CCR5 gene, which is a host factor known to influence  
92 disease progression. The CCR5 gene encodes one of the main coreceptors required for HIV-1  
93 entry, and a heterozygous genotype for a 32 base pair deletion (CCR5 wt/ $\Delta$ 32) associates with a  
94 lower viral load set point, defined as the viral load between 18 and 24 months after  
95 seroconversion which is stable in most HIV-1 infected individuals and predictive for clinical  
96 course of infection (Mellors et al. 1996; de Wolf et al. 1997), and a slower HIV-1 disease  
97 progression (de Roda Husman et al. 1997; Ioannidis et al. 2001). Given the reported lower  
98 percentages of CCR5 expressing target cells and higher levels of RANTES production in HIV-1  
99 infected individuals with a CCR5 wt/ $\Delta$ 32 genotype (de Roda Husman et al. 1999a; Blaak et al.  
100 2000), it is likely that target cell and CCR5 availability influence HIV-1 intra-patient evolution  
101 and contributes to the progression to AIDS.

102 To investigate these influences, we compared the evolution of CCR5-using HIV-1 variants (R5)  
103 in individuals with either a CCR5 wt/wt or CCR5 wt/ $\Delta$ 32 genotype who only harbored CCR5-  
104 using HIV-1 variants in their progressive or long-term non-progressive course of infection. Such  
105 comparisons require asking questions across multiple populations of individuals about the  
106 evolutionary histories that occur within each individual. Traditional modelling of evolutionary  
107 histories across individuals generally assumes that within-individual processes vary  
108 independently and are fit separately from individual to individual (Shankarappa et al. 1999; Ross  
109 and Rodrigo 2002; Potter et al. 2006; Lemey et al. 2007; Carvajal-Rodriguez et al. 2008). Often,  
110 this approach results in poor estimates of the underlying evolutionary parameters, as the

111 informative content within a single intra-host dataset is sparse. Not surprisingly, Carvajal-  
112 Rodriguez et al. (2008) arrived at the conclusion that the statistical characterization of HIV  
113 within-host evolutionary processes in relationship to disease progression is a difficult task and  
114 suffers from a lack of power. To overcome the data sparsity, one may enforce strict equality  
115 between within-individual evolutionary parameters (Rodrigo et al. 2003). In both cases, however,  
116 the ability to formally assess similarities or differences between populations of individuals is lost.  
117 Hierarchical modelling (Laird and Ware 1982)(Gelman et al. 1995), and in particular hierarchical  
118 phylogenetic models [HPMs] (Suchard et al. 2003), furnish an advantageous statistical  
119 framework in which to consider drawing conclusions across populations of individuals about the  
120 evolutionary processes within individuals. In general, the Bayesian hierarchical framework  
121 allows different evolutionary histories of the intra-host variants and pressures driving their  
122 evolution from individual to individual, while providing overall or across-individual summaries  
123 of important evolutionary measures, such as the DNA sequence mutation rate or  
124 synonymous/non-synonymous rate ratio ( $d_N/d_S$ ) identifying positive selection. Critically, the  
125 HPM allows the within-individual-level parameters to vary about, for example, an unknown  
126 common mean for each population. This occurs through the employment of a hierarchical prior  
127 distribution on the parameters that are in turn characterized by unknown estimable  
128 hyperparameters. Then conveniently, hypothesis testing reduces to asking if these common mean  
129 parameters differ between populations. Fortuitously, the hierarchical prior embedded in the HPM  
130 also affords a borrowing of strength of information from one individual by another, providing  
131 more precise within-individual-level estimates (Suchard et al. 2003; Kitchen et al. 2004; Kitchen  
132 et al. 2006; Kitchen et al. 2009).

133 In this study, we extend the HPM across multiple populations of individuals through the  
134 introduction of population-specific, fixed effects. These effects allow the expected evolutionary

135 parameter estimated within a population to potentially vary across populations. We then exploit  
136 ideas from Bayesian model averaging (Hoeting et al. 1999) and selection (Suchard, Weiss, and  
137 Sinsheimer 2001) to formally ask if these effects statistically differ between populations. We use  
138 this approach to estimate viral evolutionary rates and selective pressures within hosts and to  
139 evaluate whether these quantities differ with respect to CCR5 wt/ $\Delta$ 32 host genetic background  
140 and disease progression.

141

## 141 MATERIALS AND METHODS

### 142 Study subjects

143 18 men who have sex with men (MSM) participants in the Amsterdam Cohort Studies on HIV  
144 and AIDS, 11 with a CCR5 wt/wt genotype (patients P1 to P11) and 7 with a CCR5 wt/ $\Delta$ 32  
145 genotype (patients P12 to P18), who at all times tested during follow-up harbored only R5 HIV-1  
146 variants were selected. All patients were either seropositive at entry in the cohort studies  
147 (seroprevalent cases with an average imputed seroconversion (SC) date of 18 months before entry  
148 in the cohort (Geskus 2000)) or seroconverted during active follow-up in the cohort studies. Nine  
149 individuals were classified as long-term non-progressors (LTNP) (defined as HIV-1 infected  
150 patients that at the end of follow-up (April 1997) had an asymptomatic seropositive follow-up of  
151 at least 11 years with relatively stable CD4<sup>+</sup> T cell counts that were still above 400 cells/ml in the  
152 ninth year of follow-up in the absence of antiretroviral therapy). The remaining nine individuals  
153 progressed to AIDS during the study period (median time to AIDS = 8.2 (2.7-10.8) years after SC  
154 or imputed SC date) and were classified as Progressors (P). Individuals included in this study did  
155 not receive effective antiretroviral therapy during the study period. Clinical parameters and time  
156 points of virus isolation are shown per patient in Figure 1.

157 The Amsterdam Cohort Studies are conducted in accordance with the ethical principles set out in  
158 the declaration of Helsinki and written informed consent was obtained prior to data collection.

159 The study was approved by the Academic Medical Center institutional medical ethics committee.

160

### 161 Isolation of clonal HIV-1 variants

162 Clonal HIV-1 variants were isolated by co-cultivation of serial dilutions of patient Peripheral  
163 Blood Mononuclear Cells (PBMC) from two to eight time points in the course of their infection

164 and expanded to viral stocks for further study as described previously (Schuitemaker et al. 1992;  
165 van 't Wout, Schuitemaker, and Kootstra 2008). For each patient, time points of virus isolation  
166 and number of clonal HIV-1 variants per time point are summarized in Supplementary Table S1.  
167 The R5 phenotype of all clonal HIV-1 variants that were isolated was confirmed by inability to  
168 replicate in the MT2 cell-line, in PHA-PBMC from a donor with a CCR5 $\Delta$ 32 homozygous  
169 genotype and in astrogloma cells transfected with CD4 and CCR3 or CXCR4 (de Roda Husman  
170 et al. 1999b) and predicted co-receptor use based on the V3 amino acid sequence using the  
171 position specific scoring matrix (PSSM) (NSI/SI  
172 (<http://indra.mullins.microbiol.washington.edu/pssm/>))(Jensen et al. 2003).

173

#### 174 **DNA isolation, PCR and sequencing**

175 Total DNA was isolated from PBMCs infected with clonal HIV-1 variants using a modification  
176 of the L6 isolation method (Kootstra and Schuitemaker 1999). Precipitated DNA was dissolved  
177 in 100 $\mu$ l of distilled water and 5 $\mu$ l were used for PCR amplification of the gp120 (C1-C4) region  
178 corresponding to HXB2 nucleotide positions 6444 to 7595. Amplification was performed by PCR  
179 with primers TB3 forward (5'-GGCCTTATTAGGACACATAGTTAGCC-3') and OFM19  
180 reverse (5'-GCACTCAAGGCAAGCTTTATTGAGGCTTA-3') using the expand high fidelity  
181 Taq polymerase kit (Roche) and the following amplification cycles: 2 min 30s 94°C, 9 cycles of  
182 15s 94°C, 45s 50°C, 6 min 68°C, 30 cycles of 15s 94°C, 45s 53°C, 6 min 68°C, followed by a 10  
183 min extension at 68°C and subsequent cooling to 4°C. Nested PCR was performed with two  
184 different inner PCR primer combinations: Seq1 forward (5'-  
185 TACATAATGTTTGGGCCACACATGCC -3'), Seq4 reverse (5'-  
186 CTTGTATTGTTGTTGGGTCTTGAC -3'), Seq5 forward (5'-  
187 GTCAACTCAACTGCTGTAAATGGC -3') and Seq2 reverse (5'-

188 TCCTTCATATCTCCTCCTCCAGGTC -3'). Nested PCRs were performed using Promega Taq  
189 polymerase in the presence of 2mM MgCl<sub>2</sub> using the following amplification cycles: 5 min 94°C,  
190 40 cycles of 15s 95°C, 30s 59°C, 2 min 72°C, followed by a 10 min extension at 72°C and  
191 subsequent cooling to 4°C.

192 PCR products were purified using ExoSAP-IT (USB, Cleveland, Ohio, USA) according to  
193 manufacturer's protocol. Sequencing conditions consisted of 5' at 94°C, 30 cycles of 15'' at 94°C,  
194 10'' at 50°C, 2' at 60°C and a 10' extension at 60°C. Sequencing was performed using BigDye  
195 Terminator v1.1 Cycle Sequencing kit (ABI Prism, Applied Biosystems, Warrington, UK)  
196 according to the manufacturer's protocol using the nested PCR primers. Sequences were analyzed  
197 on the Applied Biosystems 3130 xl Genetic Analyzer. The nucleotide sequences are available  
198 from Genbank under the accession numbers EU743973.1-EU44009.1, EU744014.1-EU744046.1,  
199 EU744055.1-EU744093.1, EU744097.1-EU744129.1, EU744146.1-EU744175.1, GU455514-  
200 GU455525 and HQ644787-HQ645012.

201

## 202 **Bayesian inference of within-host HIV evolutionary rates and selection pressures**

203 Nucleotide sequences for all clonal HIV-1 gp120 (C1-C4) variants isolated from the individual  
204 patients were aligned using ClustalW (Thompson, Higgins, and Gibson 1994) and manually  
205 edited. Cross-contamination was excluded using phylogenetic analysis.

206 *(a) Independent estimates of within-host evolutionary rates.* Nucleotide substitution rates were  
207 estimated for each patient using strict and relaxed (uncorrelated lognormal) molecular clock  
208 models implemented in BEAST v.1.4.8 (Drummond et al. 2006; Drummond and Rambaut 2007).  
209 We used a general time-reversible (GTR) model of nucleotide substitution with discrete gamma-  
210 distributed rate variation among sites. Posterior distributions were obtained using Bayesian

211 Markov chain Monte Carlo (MCMC) analysis. MCMC chains were run sufficiently long to  
212 ensure stationarity and adequate effective sample sizes ( $ESS > 100$ ) as diagnosed using Tracer  
213 (<http://tree.bio.ed.ac.uk/software/tracer/>). The uncertainty of continuous parameter estimates is  
214 expressed as 95% highest posterior density (HPD) intervals.

215 *(b) Hierarchical estimates of evolutionary parameters.* To draw inference about different  
216 evolutionary patterns across populations of patients, we implement a novel HPM in BEAST  
217 (Suchard et al. 2003). HPMs analyze viral sequence data from multiple patients simultaneously  
218 and have found extensive use in uncovering common patterns of intra-host HIV evolution  
219 (Kitchen et al. 2004; Kitchen et al. 2006; Kitchen et al. 2009). At the heart of the HPM lies a  
220 Bayesian mixed effects model that pools information across patients. Pooling information  
221 through random effects affords more precise individual-patient parameter estimates when the data  
222 are sparse for a patient. Further, unique to the work here, the introduction of fixed effects (see  
223 below) offers a formal hypothesis testing framework from which to identify differences in  
224 evolutionary process between patient population groups.

225 Let  $\theta_i$  for  $i = 1, \dots, N$  patients represent the evolutionary process parameter of interest; this could  
226 be, for example, the overall rate of nucleotide substitution or the nonsynonymous/synonymous  
227 substitution rate ratio ( $d_N/d_S$ ) in a codon substitution process across the unknown genealogy  
228 relating the sequences from within patient  $i$ . In the analysis of four different patient groups:  
229 Progressors, Long-term non-progressors (LTNP), CCR5 wt/wt (WT) and CCR5 wt/ $\Delta 32$  ( $\Delta 32$ ),  
230 we assume that either  $\log \theta_i$  or  $\theta_i$  is drawn from an underlying normal distribution where the  
231 mean and variance of this underlying prior distribution are also unknown and simultaneously  
232 estimated along with all sequence data. The choice of a log transform is convenient for modeling  
233 strictly positive parameters. Importantly, fixing this mean and variance to known values does not

234 return a hierarchical model, but rather results in complete independence across individuals. On  
235 the other hand, estimating the mean or variance imparts both an approach to make comparisons  
236 across populations and the borrowing of strength for poorly informed within-individual model  
237 parameters.

238 For nucleotide analyses, we apply this hierarchical setup to the strict clock evolutionary rate (on  
239 the log-scale), the mean evolutionary rate parameter of the lognormal relaxed clock (log), the  
240 constant population size (log) of the demographic prior, the GTR substitution parameters (log)  
241 and the shape parameter (log) of the discrete gamma distribution modeling rate variation among  
242 sites. For codon model analyses, a hierarchical transition/transversion rate parameter and a  
243 hierarchical  $d_N/d_S$  rate ratio (Goldman and Yang 1994) replace the GTR model parameters.

244 *(c) Hierarchical estimation with population-specific, fixed effects.* For hypothesis testing  
245 purposes, we extend the HPM to include across-population fixed effects. Each patient belongs to  
246 one of four fixed population groups that we can designate using two indicator factors:  $LTNP_i = 0$   
247 (1) for short (long) term progressors and  $\Delta 32_i = 0$  (1) for deletion 32 absent (present) patients.  
248 Our HPM assumes

249

$$250 \quad \log \theta_i = \beta_0 + \delta_{LTNP} \beta_{LTNP} LTNP_i + \delta_{\Delta 32} \beta_{\Delta 32} \Delta 32_i + \varepsilon_i,$$

251

252 where  $\beta_0$  is an unknown grand-mean,  $\delta_{LTNP}$  and  $\delta_{\Delta 32}$  are binary indicator variables,  $\beta_{LTNP}$  and  $\beta_{\Delta 32}$   
253 are conditional effective sizes and  $\varepsilon_i$  are independent and normally distributed random variables  
254 with mean 0 and an estimable variance. The inclusion of the indicator variables follows from a  
255 Bayesian stochastic search variable selection approach (Kuo and Mallick 1998; Chipman, George,  
256 and McCulloch 2001) that simultaneously estimates the posterior probabilities of all possible

257 linear models that may or may not include LTNP or  $\Delta_{32}$  status effects. When an indicator equals  
258 1, this effect is included in the model, demonstrating that the evolutionary process parameter  
259 differs with high probability between patient population groups. Lemey et al. (2009) discuss  
260 Bayesian stochastic search variable selection in further detail (Lemey et al. 2009).

261 We complete this HPM model with variable selection through assigning independent Bernoulli  
262 prior probability distributions on  $\delta_{\text{LTNP}}$  and  $\delta_{\Delta_{32}}$ . These distributions place equal probability on  
263 each factor's inclusion and exclusion. We further assume diffuse priors on the unknown grand-  
264 mean and error variance and specify that *a priori*  $\beta_{\text{LTNP}}$  and  $\beta_{\Delta_{32}}$  are normally distributed with  
265 mean 0 and a variance of 1/2. We choose 1/2 as, before seeing the data, we believe that, if a  
266 factor does result in different evolutionary parameters across population groups, process  
267 parameters should differ by at most an order of magnitude on their original scale. The  
268 introduction of HPMs into BEAST necessitates the development of MCMC transition kernels to  
269 efficiently explore that space of the grand-mean and effect-size, model indicator, and random-  
270 effects variance parameters. Given our judicious prior choices, the full conditional distributions  
271 of these parameters are in standard-form: multivariate-normal, binomial and inverse-gamma,  
272 respectively. This enables us to build highly effective Gibbs samplers (Casella and George 1992;  
273 Suchard et al. 2003) over the joint space of these parameters. Suchard et al. (2003) provide  
274 detailed derivations of the full condition distributions and their Gibbs samplers (Suchard et al.  
275 2003). We implement these Gibbs samplers as regular BEAST “operators” that are now  
276 accessible to interested readers through BEAST’s XML model specification language.  
277 Supplementary material to this paper reports the transition kernels’ XML syntax and gives  
278 examples on their use to implement HPMs.

279 To assign statistical significance to differences between population groups, we employ Bayes  
280 factors (Jeffreys 1998; Suchard, Weiss, and Sinsheimer 2001) that report how much the data  
281 change our prior opinion (here, 1:1 odds) about the inclusion of each factor. These Bayes factors  
282 are straightforward to estimate through the variable selection procedure, as the Bayes factor  
283 equals the posterior odds that a factor indicator equals 1 divided by the corresponding prior odds.  
284 The posterior odds follow immediately from the marginal posterior probability that a factor  
285 indicator equals 1 that we estimate through the posterior expectation of the factor indicator. In  
286 cases where an estimate of this expectation approaches very closely to 0 or 1, an estimator based  
287 on a Rao-Blackwellization procedure is available (Casella and Robert 1996).  
288

## 288 **RESULTS**

### 289 **Independent versus hierarchical estimation of evolutionary parameters**

290 We first explored the nucleotide substitution rate as a hierarchical parameter estimated across  
291 patients in four separate patient groups: Progressors, LTNP, CCR5 wt/wt and CCR5 wt/ $\Delta$ 32.  
292 Using a strict clock model, a higher mean evolutionary rate was estimated in the Progressors  
293 group (mean =  $7.65 \times 10^{-4}$ , 95% HPD = [ $6.45 \times 10^{-4}$ ,  $8.84 \times 10^{-4}$ ]) compared to the LTNP group  
294 ( $5.87 \times 10^{-4}$  [ $4.30 \cdot 10^{-4}$ - $7.55 \times 10^{-4}$ ]) (Figure 2A). While these estimates demonstrate overlapping  
295 marginal posterior credible intervals (CIs), immediately concluding that their difference is not  
296 significant ignores the correlation between the rates; we return to a formal test later. A less  
297 pronounced difference in evolutionary rate was estimated between the CCR5 wt/wt ( $7.27 \times 10^{-4}$   
298 [ $5.74 \times 10^{-4}$ - $8.75 \times 10^{-4}$ ]) and CCR5 wt/ $\Delta$ 32 ( $6.00 \times 10^{-4}$  [ $4.21 \times 10^{-4}$ - $7.89 \times 10^{-4}$ ]) groups. Similar rate  
299 differences, with somewhat less overlapping CIs between Progressors ( $7.57 \times 10^{-4}$  [ $6.49 \times 10^{-4}$ -  
300  $8.67 \times 10^{-4}$ ]) and LTNPs ( $5.63 \times 10^{-4}$  [ $4.19 \times 10^{-4}$ - $7.06 \times 10^{-4}$ ]), were observed using a relaxed clock  
301 model (Figure 2B), in which the log of the mean evolutionary rate across all branches in a patient  
302 genealogy is drawn from an underlying normal distribution. For both strict and relaxed  
303 evolutionary rate estimates (Figure 3 A and B), as well as other substitution model and  
304 population genetic parameters (data not shown), we observed significant shrinkage in uncertainty  
305 under the standard hierarchical fit, which clearly demonstrates the HPM improvement. Moreover,  
306 separate fit of parameter-rich models such as the uncorrelated relaxed clock required informative  
307 priors to achieve efficient sampling. To demonstrate the impact of such priors on our posterior  
308 rate estimates obtained by separate model fitting, and compare these with the hierarchical  
309 estimates that did not require such priors, we plot the marginal posterior rate estimates for the  
310 three least informative (lowest number of time points and/or sequences per time point) and three

311 most informative patients within the LTNP group (P10, P16 and P17 versus P9, P11 and P13  
312 respectively) as violin plots in Figure 4. Violin plots are box plots overlaid with (rotated) kernel  
313 density estimates in order to show to the probability density at different parameter values. The  
314 patients for which only two or three time points were available resulted in rate estimates that only  
315 weakly diverged from their respective prior (uniform[0,0.004] or lognormal(-7.5,1); Figure 4A  
316 and C respectively), whereas many time points provide sufficient information to dominate over  
317 these priors (Figure 4 B and D). Under the hierarchical model, even weakly informative patient-  
318 specific data sets with extremely diffuse priors on the rate yield relative precise posteriors (Figure  
319 4E), and the individual patient estimates are only marginally higher than for the three most  
320 informative patients (Figure 4F). This demonstrates that comparing the mean rates for individual  
321 estimates would be inappropriate to assess differences among patient groups. Weakly informative  
322 patients result in relatively high mean rates, but their high variances ensure that the contribution  
323 to the population rate (LTNP group) in the hierarchical model remains low.

324 While the application of relaxed clock models to individual data sets with few time points or  
325 sequences may be questionable, analysis under a HPM, in which information is pooled between  
326 patients, enables us to side-step this limitation. Marginal likelihood estimates for the both strict  
327 and relaxed clock analyses of the different patient groups (Supplementary Table S2) indicate a  
328 better fit of the relaxed clock model, with log Bayes factors (BFs) of 7.8, 6.1, 4.4 and 4.2 in favor  
329 of the relaxed clock for Progressors, LTNP, CCR5 wt/wt and CCR5 wt/ $\Delta$ 32 respectively. The  
330 fact that a strict clock could often not be rejected for individual patient analysis also indicates the  
331 HPM draws on increased statistical power of HPMs to reject simpler models. Because of the  
332 increased model fit, we employ relaxed clocks in further codon model analyses and hypothesis  
333 tests incorporating fixed effects.

334 Analyses using a codon model revealed comparable codon substitution rate differences between  
335 Progressors/LTNP and between CCR5 wt/wt and CCR5 wt/ $\Delta$ 32 compared to the nucleotide  
336 analyses (Supplementary Figure 1A vs. Figure 2B). Hierarchical  $d_N/d_S$  estimates, however, were  
337 comparable for the four patient groups (Supplementary Figure 1B).

338

### 339 **Hypothesis testing using HPMs incorporating across-population fixed effects**

340 The four different groups considered previously are not comprised of independent patient sets;  
341 some patients fall in more than one group. Hence, direct comparison of the marginal parameter  
342 estimates fit to each group independently does not generate independent estimates. For more  
343 appropriate hypothesis testing of difference, the HPM for the evolutionary rate was extended to  
344 accommodate fixed effects (see methods), enabling estimation of hierarchical parameters across  
345 all patients. Successfully, hierarchical estimation with fixed effects across all patients resulted in  
346 even further shrinkage of individual patient estimates compared to hierarchal models applied to  
347 separate groups (Figure 3B). Bayes factor comparison of the fixed-effects HPM model to a model  
348 that assumes either completely linked or unlinked parameters (log BF of 51.7 and 57.2  
349 respectively) provides strong evidence that the shrinkage is accompanied by improved goodness-  
350 of-fit. The main results of the fixed effect HPM analyses are listed in Table 1. For the nucleotide  
351 analysis, the LTNP versus Progressor and CCR5 wt/wt versus CCR5 wt/ $\Delta$ 32 effects were  
352 employed to model the evolutionary rates. Through examining the posterior distribution of the  
353 rate indicators ( $\delta_{\text{effect}}$ ), we estimate the posterior probability for including the LTNP versus  
354 Progressor effect at 0.72 resulting in a moderate Bayes factor support of 2.6 in agreement with  
355 the group-by-group hierarchical rate estimates obtained above. Importantly, the rate decrease  
356 attributable to this fixed-effect returns a credible interval that does not include 0. This approach

357 appropriately controls for the non-independence missed in the group-by-group analyses and  
358 rejects the null hypothesis of no difference between LTNP and Progressor patients.  
359 There was no support in favor of a CCR5 wt/ $\Delta$ 32 effect. Even after conditioning on the effect-  
360 indicator equaling 1 to estimate the potential effect-size, the posterior CCR5 wt/ $\Delta$ 32 effect-size  
361 parameter distribution remained centered close to 0 with symmetric CIs. In the codon analysis,  
362 the same effects were tested on both the substitution rate and  $d_N/d_S$ . A very similar LTNP effect  
363 was observed for the codon substitution rate, although the CIs now included 0. Interestingly, the  
364 conditional effect size of LTNP versus Progressor on codon substitution rate remains very similar  
365 to the effect size on nucleotide substitution rate. Further, there was more support against than in  
366 favor of a CCR5 wt/ $\Delta$ 32 effect. Finally, no support for a LTNP effect or CCR5 wt/ $\Delta$ 32 effect was  
367 observed on the hierarchical  $d_N/d_S$  estimates.  
368

368 **DISCUSSION**

369 In this study, we adopted a HPM approach to estimate within-host HIV evolutionary parameters  
370 and test evolutionary hypotheses regarding host susceptibility and disease progression. We  
371 sought to investigate whether the CCR5 wt/ $\Delta$ 32 genotype, which is associated with a lower viral  
372 load set point and a slower HIV-1 disease progression (de Roda Husman et al. 1997; Ioannidis et  
373 al. 2001), also impacts the evolutionary rate of the virus by limiting target cell or CCR5  
374 availability. Furthermore, we wanted to evaluate the contribution of CCR5 availability and CCR5  
375 use on the selection pressure directed against the viral envelope protein by estimating  $d_N/d_S$ .

376 HPMs have been used for HIV evolutionary enquiry before, but this is the first study that  
377 develops HPMs to estimate evolutionary rate,  $d_N/d_S$  and demographic parameters. In a HPM  
378 framework, we assume that the patient-specific HIV-1 evolutionary parameters can be drawn  
379 from a population distribution. Estimations of the evolutionary process based on a limited sample  
380 from each patient are riddled with noise and the improvement of a HPM follows from the  
381 reduced uncertainty on individual patient estimates. Bayes factor comparison further confirms a  
382 considerable improvement in goodness-of-fit of the HPM with respect to a completely linked and  
383 unlinked model. This can be explained by the fact that the completely linked model  
384 inappropriately ignores any difference among patients on the one hand, and a completely linked  
385 model suffers from an unnecessarily high effective number of parameters (Spiegelhalter et al.  
386 2002) arising from the independent prior specifications on the other hand. The HPM sits in  
387 between these two extremes and reduces the effective number of parameters without sacrificing  
388 fit to the data. Furthermore, we demonstrate that the HPM is more powerful in rejecting simpler  
389 evolutionary models, like the constant rate assumption, which is frequently violated for HIV.

390 The hierarchical estimates for the Progressors, LTNP, CCR5 wt/wt and CCR5 wt/ $\Delta$ 32 groups  
391 indicated a pronounced strict and relaxed clock rate difference between the Progressors and

392 LTNP, whereas differences between CCR5 wt/wt and CCR5 wt/ $\Delta$ 32 rates were less pronounced.  
393 The same patterns were observed for relaxed codon substitution rates, but no real differences  
394 were noted in terms of  $d_N/d_S$ . These comparisons are based on non-independent data because  
395 patients will be part of two different groups. For more appropriate hypothesis testing, we  
396 incorporated fixed effects and employed Bayesian stochastic search variable selection to estimate  
397 the posterior probability that different patient group characteristics influence within-host  
398 evolutionary parameters. The advantage of a Bayesian model averaging approach that  
399 simultaneously explores the space of models and regression coefficients is the opportunity to  
400 distinguish between the relative size of an effect and its importance, which can be formalized in  
401 terms of standard Bayes factor support. The latter effectively becomes independent of the scale of  
402 the predictors, which otherwise may confound drawing conclusions on the effect sizes only.  
403 Because both predictors we considered only achieve 0 or 1, controlling for scale is not an issue in  
404 the current study, but it does contribute to a more general framework for evolutionary hypothesis  
405 testing. While the statistical support is not decisive, the fixed-effects HPM approach produces  
406 substantially more efficient parameter estimates and conditional effect sizes confirm rate  
407 differences among LTNP and Progressors. Despite the elevated power, more elaborate sampling  
408 in terms of numbers of patients, within-host time points or maybe even larger genome regions  
409 would be desirable.

410 The HPM estimates suggest an association between evolutionary rate and disease progression,  
411 but the CCR5 genotype does not account for the rate differences. Given the absence of clear  
412  $d_N/d_S$  differences – if anything, they are slightly higher in LTNP – we cannot attribute the rate  
413 nuances to differences in selection on amino acid fixation. Therefore, we conclude that these  
414 differences are due to variations in the product of mutation rate and generation time. In particular,  
415 lower replication rates may be associated with delayed onset of AIDS symptoms. In agreement

416 with this, a codon-model extension of the Bayesian relaxed-clock analysis of more extensively  
417 sampled patients has shown that absolute synonymous substitutions are correlated with disease  
418 progression (Lemey et al. 2007). These authors argued that synonymous substitutions were a  
419 marker of replication rate and most probably reflect the action of immune activation, which in  
420 itself is a marker of disease progression. In the current study, we employed standard codon model  
421 implementation in the Bayesian framework, rather than evaluating genealogies under nucleotide  
422 models as a proxy. This approach comes at a computational expense, and further extensions -  
423 such as codon models to estimate absolute rates of synonymous and nonsynonymous  
424 substitutions (Seo, Kishino, and Thorne 2004) - may prove even more computationally intensive.  
425 Fortunately, recent advances in GPU computation provide significant increases in computation  
426 speed for high state space models (Suchard and Rambaut 2009). These advances promise to  
427 stimulate further development of various codon models in the Bayesian framework, the  
428 parameters of which could be efficiently estimated in hierarchical models.

429 CCR5 genotype has a measurable impact on disease progression (de Roda Husman et al. 1997;  
430 Ioannidis et al. 2001) but there appears to be no absolute relationship (not all CCR5 wt/ $\Delta$ 32  
431 infected individuals are LTNP). This implies a more complex scenario, in which the combination  
432 of CCR5 availability with other host genetic factors, in particular cellular and humoral immune  
433 pressures, and immune activation, will determine the viral replication rate and progression of the  
434 disease in a patient. While lower CCR5 availability does not appear to exert selection pressure on  
435 the viral envelope during the chronic phase of infection, it cannot be excluded that in HIV-1  
436 infected individuals with CCR5 WT/ $\Delta$ 32 genotype, in whom CCR5<sup>+</sup> target cells and CCR5  
437 expression are already limiting in the acute phase, selection for viruses with optimal CCR5 use  
438 occurs in a very early stage. Moreover, we performed analyses on sequences in which  
439 ambiguously aligned hypervariable regions were deleted, which may play an important role in

440 both humoral immune responses (Cao et al. 1997; Chackerian, Rudensey, and Overbaugh 1997;  
441 Stamatatos and Cheng-Mayer 1998; Pinter et al. 2004; Sagar et al. 2006; Gray et al. 2007) and  
442 selection for optimal CCR5 use (Hubert and Arabie 1985; Stamatatos, Wiskerchen, and Cheng-  
443 Mayer 1998; Wang et al. 1999; Sagar et al. 2006; Repits et al. 2008).

444 Studying evolutionary dynamics within hosts has become an integral part of HIV research, but  
445 one that still faces the challenge of fully unraveling the relationship between evolutionary  
446 parameters and clinical outcome. There may be several reasons for the difficulty in establishing  
447 the role of evolutionary processes in disease progression. Within-host dynamics appear to be  
448 highly complex, with many host-specific and environmental (co-infections) factors interacting  
449 with various evolutionary processes such as hypermutation, diversifying and directional selection,  
450 recombination and compartmentalization. Untangling this complex interplay requires accurate  
451 measurement of all host factors involved and evolutionary models that explicitly accommodate  
452 the relevant evolutionary forces. Without the latter, many simplifying assumptions are at risk of  
453 being violated when considering HIV evolution. Parameter-rich models may be limited by  
454 current sampling as they require highly informative data. To our knowledge, the most elaborate  
455 sampling dates back to over a decade ago (Shankarappa et al. 1999), which, differently from this  
456 study, included patients with HIV populations harboring CXCR4-using variants. Next generation  
457 sequencing may offer new opportunities for within host HIV genetic analyses, but produces data  
458 with particular challenges for comparative analyses (Vrancken et al. 2010). Here, we have  
459 adopted a modeling approach that efficiently pools the information from multiple individuals and  
460 we demonstrate how this can be employed for rigorous testing across patient populations. We  
461 hope that this stimulates further model-based inference of evolutionary processes, which  
462 ultimately may lead to more profound insights into persistent viral infections.

463

463 **FUNDING**

464 This work was supported by Netherlands AIDS fund (grant nr 6006) and The European  
465 Community's Seventh Framework Programme NGIN (FP7/2007-2013) under grant agreement n°  
466 201433. PL was supported by a postdoctoral fellowship from the Fund for Scientific Research  
467 (FWO) Flanders. MAS and JAT are supported by the National Institutes of Health R01 grant  
468 GM86887. The research leading to these results has received funding from the European  
469 Research Council under the European Community's Seventh Framework Programme (FP7/2007-  
470 2013) / ERC Grant agreement n° 260864.

471

472

473 **ACKNOWLEDGEMENTS**

474

475 The Amsterdam Cohort Studies on HIV infection and AIDS, a collaboration between the Public  
476 Health Service of Amsterdam, the Academic Medical Center of the University of Amsterdam, the  
477 Sanquin Blood Supply Foundation, the University Medical Center Utrecht, and the Jan van  
478 Goyen Medical Center, are part of the Netherlands HIV Monitoring Foundation and financially  
479 supported by the Center for Infectious Disease Control of the Netherlands National Institute for  
480 Public Health and the Environment.

481

482 **TABLES**

483 **Table 1. Estimates of the long-term non-progressor (LTNP) and  $\Delta 32$  effects on nucleotide**  
 484 **substitution rates, codon substitution rates and  $d_N/d_S$ .**

485

Evolutionary parameter	Effect support/size	LTNP effect	$\Delta 32$ effect
Nucleotide substitution rate	Posterior probability $\delta_{\text{effect}} = 1$	0.72	0.27
	Bayes factor <sub>effect</sub>	2.6	0.4
	$\beta_{\text{effect}}   \delta_{\text{effect}} = 1^*$	-0.275 (-0.524,-0.016)	-0.007 (-0.940,0.920)
Codon substitution rate:	Posterior probability $\delta_{\text{effect}} = 1$	0.726	0.324
	Bayes factor <sub>effect</sub>	2.6	0.5
	$\beta_{\text{effect}}   \delta_{\text{effect}} = 1^*$	-0.265 (-0.523,0.019)	-0.012 (-0.700,0.692)
$d_N/d_S$	Posterior probability $\delta_{\text{effect}} = 1$	0.502	0.393
	Bayes factor <sub>effect</sub>	1.0	0.6
	$\beta_{\text{effect}}   \delta_{\text{effect}} = 1^*$	0.083 (-0.101,0.25)	-0.005 (-0.228,0.242)

486

487 \*these are effective sizes conditional on the effect being included (the binary effect indicator

488  $\delta_{\text{effect}}$  being 1). For the rates these effective sizes are in log space.

489 **FIGURES**

490 **Figure 1. CD4<sup>+</sup> T cell numbers, viral loads, and antiretroviral treatments of 18 participants**  
491 **from the Amsterdam Cohort Studies who were selected for this study.** Time points of clinical  
492 AIDS diagnosis are indicated with open downward triangles. Arrows indicate time points of  
493 clonal virus isolation. The length and type of antiretroviral therapy are indicated in the top part of  
494 the panels.

495  
496 **Figure 2. Evolutionary rate estimates using a hierarchical phylogenetic model applied**  
497 **separately to four patient groups (Progressors, LTNP, CCR5 wt/wt and CCR5 wt/ $\Delta$ 32).**  
498 Evolutionary rate estimated under strict clock model (A). Mean evolutionary rate estimated under  
499 relaxed clock model (B). CCR5 wt/wt (WT); CCR5 wt/ $\Delta$ 32 ( $\Delta$ 32).

500  
501 **Figure 3. Improved statistical efficiency (shrinkage effect) of the hierarchical phylogenetic**  
502 **model.** Strict clock (A). Relaxed clock (B). Posterior variance of estimated evolutionary rate  
503 from the independent analyses of each patient (white); evolutionary rate variance from the  
504 hierarchical analysis of LTNPs and Progressors (black); evolutionary rate variance from the  
505 hierarchical analysis of LTNPs and Progressors incorporating fixed effects (grey).

506  
507 **Figure 4. Marginal posterior rate distributions for LTNP patients with different numbers**  
508 **of sampling time points.** Least informative patients (lowest number of time points or sequences  
509 per time point): P10, P16 and P17. Most informative patients: P9, P11 and P13. A & B:  
510 Assuming a uniform[0,0.004] rate prior. C & D: lognormal(-7.5,1) rate prior. E & F: hierarchical  
511 phylogenetic model with unknown mean and variance and diffuse priors.

512

512 **REFERENCES**

- 513
- 514
- 515 Blaak, H., L. J. Ran, R. Rientsma, and H. Schuitemaker. 2000. Susceptibility of in vitro  
516 stimulated PBMC to infection with NSI HIV-1 is associated with levels of CCR5  
517 expression and beta-chemokine production. *Virology* **267**:237-246.
- 518 Borrow, P., H. Lewicki, X. Wei, M. S. Horwitz, N. Pfeffer, H. Meyers, J. A. Nelson, J. E. Gairin,  
519 B. H. Hahn, M. B. Oldstone, and G. M. Shaw. 1997. Antiviral pressure exerted by HIV-1-  
520 specific cytotoxic T lymphocytes (CTLs) during primary infection demonstrated by rapid  
521 selection of CTL escape virus. *Nat Med* **3**:205-211.
- 522 Cao, J., N. Sullivan, E. Desjardin, C. Parolin, J. Robinson, R. Wyatt, and J. Sodroski. 1997.  
523 Replication and neutralization of human immunodeficiency virus type 1 lacking the V1  
524 and V2 variable loops of the gp120 envelope glycoprotein. *J Virol* **71**:9808-9812.
- 525 Carvajal-Rodriguez, A., D. Posada, M. Perez-Losada, E. Keller, E. J. Abrams, R. P. Viscidi, and  
526 K. A. Crandall. 2008. Disease progression and evolution of the HIV-1 env gene in 24  
527 infected infants. *Infect Genet Evol* **8**:110-120.
- 528 Casella, G., and E. I. George. 1992. Explaining the Gibbs sampler. *The American Statistician*  
529 **46**:167-174.
- 530 Casella, G., and C. Robert. 1996. Rao-Blackwellisation of sampling schemes. *Biometrika* **83**:81-  
531 94.
- 532 Chackerian, B., L. M. Rudensey, and J. Overbaugh. 1997. Specific N-linked and O-linked  
533 glycosylation modifications in the envelope V1 domain of simian immunodeficiency  
534 virus variants that evolve in the host alter recognition by neutralizing antibodies. *J Virol*  
535 **71**:7719-7727.
- 536 Chipman, H., E. George, and R. McCulloch. 2001. The practical implementation of Bayesian  
537 model selection. *IMS Lecture Notes - Monograph Series* **38**:67-134.
- 538 Coffin, J. M. 1995. HIV population dynamics in vivo: implications for genetic variation,  
539 pathogenesis, and therapy. *Science* **267**:483-489.
- 540 de Roda Husman, A. M., H. Blaak, M. Brouwer, and H. Schuitemaker. 1999a. CC chemokine  
541 receptor 5 cell-surface expression in relation to CC chemokine receptor 5 genotype and  
542 the clinical course of HIV-1 infection. *J Immunol* **163**:4597-4603.
- 543 de Roda Husman, A. M., M. Koot, M. Cornelissen, I. P. Keet, M. Brouwer, S. M. Broersen, M.  
544 Bakker, M. T. Roos, M. Prins, F. de Wolf, R. A. Coutinho, F. Miedema, J. Goudsmit, and  
545 H. Schuitemaker. 1997. Association between CCR5 genotype and the clinical course of  
546 HIV-1 infection. *Ann Intern Med* **127**:882-890.
- 547 de Roda Husman, A. M., R. P. van Rij, H. Blaak, S. Broersen, and H. Schuitemaker. 1999b.  
548 Adaptation to promiscuous usage of chemokine receptors is not a prerequisite for human  
549 immunodeficiency virus type 1 disease progression. *J Infect Dis* **180**:1106-1115.
- 550 de Wolf, F., I. Spijkerman, P. T. Schellekens, M. Langendam, C. Kuiken, M. Bakker, M. Roos, R.  
551 Coutinho, F. Miedema, and J. Goudsmit. 1997. AIDS prognosis based on HIV-1 RNA,  
552 CD4+ T-cell count and function: markers with reciprocal predictive value over time after  
553 seroconversion. *Aids* **11**:1799-1806.
- 554 Drummond, A. J., S. Y. W. Ho, M. J. Phillips, and A. Rambaut. 2006. Relaxed phylogenetics and  
555 dating with confidence. *PLoS Biol* **4**.
- 556 Drummond, A. J., and A. Rambaut. 2007. BEAST: Bayesian evolutionary analysis by sampling  
557 trees. *BMC Evol Biol* **7**:214.

558 Gelman, A., J. B. Carlin, H. S. Stern, and D. B. Rubin. 1995. Bayesian Data Analysis. Chapman  
559 & Hall/CRC, New York.

560 Geskus, R. B. 2000. On the inclusion of prevalent cases in HIV/AIDS natural history studies  
561 through a marker-based estimate of time since seroconversion. *Stat Med* **19**:1753-1769.

562 Goldman, N., and Z. Yang. 1994. A codon-based model of nucleotide substitution for protein-  
563 coding DNA sequences. *Mol. Biol. Evol.* **11**:725-736.

564 Goulder, P. J., R. E. Phillips, R. A. Colbert, S. McAdam, G. Ogg, M. A. Nowak, P. Giangrande,  
565 G. Luzzi, B. Morgan, A. Edwards, A. J. McMichael, and S. Rowland-Jones. 1997. Late  
566 escape from an immunodominant cytotoxic T-lymphocyte response associated with  
567 progression to AIDS. *Nat Med* **3**:212-217.

568 Gray, E. S., P. L. Moore, I. A. Choge, J. M. Decker, F. Bibollet-Ruche, H. Li, N. Leseke, F.  
569 Treurnicht, K. Mlisana, G. M. Shaw, S. S. Karim, C. Williamson, and L. Morris. 2007.  
570 Neutralizing antibody responses in acute human immunodeficiency virus type 1 subtype  
571 C infection. *J Virol* **81**:6187-6196.

572 Ho, D. D., A. U. Neumann, A. S. Perelson, W. Chen, J. M. Leonard, and M. Markowitz. 1995.  
573 Rapid turnover of plasma virions and CD4 lymphocytes in HIV-1 infection. *Nature*  
574 **373**:123-126.

575 Hoeting, J., D. Madigan, A. Raftery, and C. Volinsky. 1999. Bayesian Model Averaging.  
576 *Statistical Science* **14**:382-401.

577 Hollingsworth, T. D., O. Laeyendecker, G. Shirreff, C. A. Donnelly, D. Serwadda, M. J. Wawer,  
578 N. Kiwanuka, F. Nalugoda, A. Collinson-Streng, V. Ssempijja, W. P. Hanage, T. C.  
579 Quinn, R. H. Gray, and C. Fraser. 2010. HIV-1 transmitting couples have similar viral  
580 load set-points in Rakai, Uganda. *PLoS Pathog* **6**:e1000876.

581 Hubert, L., and P. Arabie. 1985. Comparing Partitions. *Journal of Classification* **2**:193.

582 Ioannidis, J. P., P. S. Rosenberg, J. J. Goedert, L. J. Ashton, T. L. Benfield, S. P. Buchbinder, R.  
583 A. Coutinho, J. Eugen-Olsen, T. Gallart, T. L. Katzenstein, L. G. Kostrikis, H. Kuipers, L.  
584 G. Louie, S. A. Mallal, J. B. Margolick, O. P. Martinez, L. Meyer, N. L. Michael, E.  
585 Operskalski, G. Pantaleo, G. P. Rizzard, H. Schuitemaker, H. W. Sheppard, G. J. Stewart,  
586 I. D. Theodorou, H. Ullum, E. Vicenzi, D. Vlahov, D. Wilkinson, C. Workman, J. F.  
587 Zagury, and T. R. O'Brien. 2001. Effects of CCR5-Delta32, CCR2-64I, and SDF-1 3'A  
588 alleles on HIV-1 disease progression: An international meta-analysis of individual-patient  
589 data. *Ann Intern Med* **135**:782-795.

590 Jeffreys. 1998. Theory of Probability. Oxford University Press, New York.

591 Jensen, M. A., F. S. Li, A. B. van 't Wout, D. C. Nickle, D. Shriner, H. X. He, S. McLaughlin, R.  
592 Shankarappa, J. B. Margolick, and J. I. Mullins. 2003. Improved coreceptor usage  
593 prediction and genotypic monitoring of R5-to-X4 transition by motif analysis of human  
594 immunodeficiency virus type 1 env V3 loop sequences. *J. Virol.* **77**:13376-13388.

595 Jones, N. A., X. Wei, D. R. Flower, M. Wong, F. Michor, M. S. Saag, B. H. Hahn, M. A. Nowak,  
596 G. M. Shaw, and P. Borrow. 2004. Determinants of human immunodeficiency virus type  
597 1 escape from the primary CD8+ cytotoxic T lymphocyte response. *J Exp Med* **200**:1243-  
598 1256.

599 Kemal, K. S., B. Foley, H. Burger, K. Anastos, H. Minkoff, C. Kitchen, S. M. Philpott, W. Gao,  
600 E. Robison, S. Holman, C. Dehner, S. Beck, W. A. Meyer, 3rd, A. Landay, A. Kovacs, J.  
601 Bremer, and B. Weiser. 2003. HIV-1 in genital tract and plasma of women:  
602 compartmentalization of viral sequences, coreceptor usage, and glycosylation. *Proc Natl*  
603 *Acad Sci U S A* **100**:12972-12977.

604 Kitchen, C. M., J. Lu, M. A. Suchard, R. Hoh, J. N. Martin, D. R. Kuritzkes, and S. G. Deeks.  
605 2006. Continued evolution in gp41 after interruption of enfuvirtide in subjects with  
606 advanced HIV type 1 disease. *AIDS Res Hum Retroviruses* **22**:1260-1266.

607 Kitchen, C. M., S. Philpott, H. Burger, B. Weiser, K. Anastos, and M. A. Suchard. 2004.  
608 Evolution of human immunodeficiency virus type 1 coreceptor usage during antiretroviral  
609 Therapy: a Bayesian approach. *J Virol* **78**:11296-11302.

610 Kitchen, C. M. R., V. Marconi, D. R. Kuritzkes, E. W. Bloomquist, S. G. Deeks, and M. A.  
611 Suchard. 2009. Two-way Bayesian hierarchical phylogenetic models: an application to  
612 the co-evolution of gp120 and gp41 during partial treatment interruptions of enfuvirtide.  
613 *Computational Statistics and Data Analysis* **53**:766-775.

614 Koning, F. A., D. Kwa, B. Boeser-Nunnink, J. Dekker, J. Vingerhoed, H. Hiemstra, and H.  
615 Schuitemaker. 2003. Decreasing sensitivity to RANTES (regulated on activation,  
616 normally T cell-expressed and -secreted) neutralization of CC chemokine receptor 5-using,  
617 non-syncytium-inducing virus variants in the course of human immunodeficiency virus  
618 type 1 infection. *J Infect Dis* **188**:864-872.

619 Kootstra, N. A., and H. Schuitemaker. 1999. Phenotype of HIV-1 lacking a functional nuclear  
620 localization signal in matrix protein of gag and Vpr is comparable to wild-type HIV-1 in  
621 primary macrophages. *Virology* **253**:170-180.

622 Kuo, L., and B. Mallick. 1998. Variable selection for regression models. *Sankhya B* **60**:65-81.

623 Kwa, D., J. Vingerhoed, B. Boeser, and H. Schuitemaker. 2003. Increased in vitro cytopathicity  
624 of CC chemokine receptor 5-restricted human immunodeficiency virus type 1 primary  
625 isolates correlates with a progressive clinical course of infection. *J Infect Dis* **187**:1397-  
626 1403.

627 Laird, N. M., and J. H. Ware. 1982. Random-effects models for longitudinal data. *Biometrics*  
628 **38**:963-974.

629 Lane, H. C., J. M. Depper, W. C. Greene, G. Whalen, T. A. Waldmann, and A. S. Fauci. 1985.  
630 Qualitative analysis of immune function in patients with the acquired immunodeficiency  
631 syndrome. Evidence for a selective defect in soluble antigen recognition. *N Engl J Med*  
632 **313**:79-84.

633 Lee, H. Y., A. S. Perelson, S. C. Park, and T. Leitner. 2008. Dynamic correlation between  
634 intrahost HIV-1 quasispecies evolution and disease progression. *PLoS Comput Biol*  
635 **4**:e1000240.

636 Lemey, P., S. L. Kosakovsky Pond, A. J. Drummond, O. G. Pybus, B. Shapiro, H. Barroso, N.  
637 Taveira, and A. Rambaut. 2007. Synonymous substitution rates predict HIV disease  
638 progression as a result of underlying replication dynamics. *PLoS Comput Biol* **3**:e29.

639 Lemey, P., A. Rambaut, A. J. Drummond, and M. A. Suchard. 2009. Bayesian phylogeography  
640 finds its roots. *PLoS Comput Biol* **5**:e1000520.

641 Mellors, J. W., C. R. Rinaldo, Jr., P. Gupta, R. M. White, J. A. Todd, and L. A. Kingsley. 1996.  
642 Prognosis in HIV-1 infection predicted by the quantity of virus in plasma. *Science*  
643 **272**:1167-1170.

644 Miedema, F., A. J. Petit, F. G. Terpstra, J. K. Schattenkerk, F. de Wolf, B. J. Al, M. Roos, J. M.  
645 Lange, S. A. Danner, J. Goudsmit, and et al. 1988. Immunological abnormalities in  
646 human immunodeficiency virus (HIV)-infected asymptomatic homosexual men. HIV  
647 affects the immune system before CD4+ T helper cell depletion occurs. *J Clin Invest*  
648 **82**:1908-1914.

649 Munoz, A., C. A. Sabin, and A. N. Phillips. 1997. The incubation period of AIDS. *AIDS* **11**  
650 **Suppl A**:S69-76.

- 651 Pinter, A., W. J. Honnen, Y. He, M. K. Gorny, S. Zolla-Pazner, and S. C. Kayman. 2004. The  
652 V1/V2 domain of gp120 is a global regulator of the sensitivity of primary human  
653 immunodeficiency virus type 1 isolates to neutralization by antibodies commonly induced  
654 upon infection. *J Virol* **78**:5205-5215.
- 655 Polk, B. F., R. Fox, R. Brookmeyer, S. Kanchanaraksa, R. Kaslow, B. Visscher, C. Rinaldo, and J.  
656 Phair. 1987. Predictors of the acquired immunodeficiency syndrome developing in a  
657 cohort of seropositive homosexual men. *N Engl J Med* **316**:61-66.
- 658 Potter, S. J., P. Lemey, W. B. Dyer, J. S. Sullivan, C. B. Chew, A. M. Vandamme, D. E. Dwyer,  
659 and N. K. Saksena. 2006. Genetic analyses reveal structured HIV-1 populations in serially  
660 sampled T lymphocytes of patients receiving HAART. *Virology* **348**:35-46.
- 661 Repits, J., J. Sterjovski, D. Badia-Martinez, M. Mild, L. Gray, M. J. Churchill, D. F. Purcell, A.  
662 Karlsson, J. Albert, E. M. Fenyo, A. Achour, P. R. Gorry, and M. Jansson. 2008. Primary  
663 HIV-1 R5 isolates from end-stage disease display enhanced viral fitness in parallel with  
664 increased gp120 net charge. *Virology* **379**:125-134.
- 665 Rodrigo, A. G., M. Goode, R. Forsberg, H. A. Ross, and A. Drummond. 2003. Inferring  
666 evolutionary rates using serially sampled sequences from several populations. *Mol Biol*  
667 *Evol* **20**:2010-2018.
- 668 Ross, H. A., and A. G. Rodrigo. 2002. Immune-mediated positive selection drives human  
669 immunodeficiency virus type 1 molecular variation and predicts disease duration. *J Virol*  
670 **76**:11715-11720.
- 671 Sagar, M., X. Wu, S. Lee, and J. Overbaugh. 2006. Human immunodeficiency virus type 1 V1-  
672 V2 envelope loop sequences expand and add glycosylation sites over the course of  
673 infection, and these modifications affect antibody neutralization sensitivity. *J Virol*  
674 **80**:9586-9598.
- 675 Schuitemaker, H., M. Koot, N. A. Kootstra, M. W. Dercksen, R. E. de Goede, R. P. van  
676 Steenwijk, J. M. Lange, J. K. Schattenkerk, F. Miedema, and M. Tersmette. 1992.  
677 Biological phenotype of human immunodeficiency virus type 1 clones at different stages  
678 of infection: progression of disease is associated with a shift from monocytotropic to T-  
679 cell-tropic virus population. *J Virol* **66**:1354-1360.
- 680 Seo, T. K., H. Kishino, and J. L. Thorne. 2004. Estimating absolute rates of synonymous and  
681 nonsynonymous nucleotide substitution in order to characterize natural selection and date  
682 species divergences. *Mol Biol Evol* **21**:1201-1213.
- 683 Shankarappa, R., J. B. Margolick, S. J. Gange, A. G. Rodrigo, D. Upchurch, H. Farzadegan, P.  
684 Gupta, C. R. Rinaldo, G. H. Learn, X. He, X. L. Huang, and J. I. Mullins. 1999.  
685 Consistent viral evolutionary changes associated with the progression of human  
686 immunodeficiency virus type 1 infection. *J. Virol.* **73**:10489-10502.
- 687 Spiegelhalter, D. J., N. G. Best, B. P. Carlin, and A. Van Der Linde. 2002. Bayesian measures of  
688 model complexity and fit. *Journal of the Royal Statistical Society: Series B (Statistical*  
689 *Methodology)* **64**:583-639.
- 690 Stamatatos, L., and C. Cheng-Mayer. 1998. An envelope modification that renders a primary,  
691 neutralization-resistant clade B human immunodeficiency virus type 1 isolate highly  
692 susceptible to neutralization by sera from other clades. *J Virol* **72**:7840-7845.
- 693 Stamatatos, L., M. Wiskerchen, and C. Cheng-Mayer. 1998. Effect of major deletions in the V1  
694 and V2 loops of a macrophage-tropic HIV type 1 isolate on viral envelope structure, cell  
695 entry, and replication. *AIDS Res Hum Retroviruses* **14**:1129-1139.
- 696 Sterjovski, J., M. J. Churchill, A. Ellett, L. R. Gray, M. J. Roche, R. L. Dunfee, D. F. Purcell, N.  
697 Saksena, B. Wang, S. Sonza, S. L. Wesselingh, I. Karlsson, E. M. Fenyo, D. Gabuzda, A.

698 L. Cunningham, and P. R. Gorry. 2007. Asn 362 in gp120 contributes to enhanced  
699 fusogenicity by CCR5-restricted HIV-1 envelope glycoprotein variants from patients with  
700 AIDS. *Retrovirology* **4**:89.

701 Suchard, M. A., C. M. Kitchen, J. S. Sinsheimer, and R. E. Weiss. 2003. Hierarchical  
702 phylogenetic models for analyzing multipartite sequence data. *Syst. Biol.* **52**:649-664.

703 Suchard, M. A., and A. Rambaut. 2009. Many-core algorithms for statistical phylogenetics.  
704 *Bioinformatics* **25**:1370-1376.

705 Suchard, M. A., R. E. Weiss, and J. S. Sinsheimer. 2001. Bayesian selection of continuous-time  
706 Markov chain evolutionary models. *Mol Biol Evol* **18**:1001-1013.

707 Thompson, J. D., D. G. Higgins, and T. J. Gibson. 1994. CLUSTAL W: improving the sensitivity  
708 of progressive multiple sequence alignment through sequence weighting, position-specific  
709 gap penalties and weight matrix choice. *Nucleic Acids Res.* **22**:4673-4680.

710 van 't Wout, A. B., H. Schuitemaker, and N. A. Kootstra. 2008. Isolation and propagation of  
711 HIV-1 on peripheral blood mononuclear cells. *Nat Protoc* **3**:363-370.

712 Veugelers, P. J., K. A. Page, B. Tindall, M. T. Schechter, A. R. Moss, W. W. Winkelstein, Jr., D.  
713 A. Cooper, K. J. Craib, E. Charlebois, R. A. Coutinho, and et al. 1994. Determinants of  
714 HIV disease progression among homosexual men registered in the Tricontinental  
715 Seroconverter Study. *Am J Epidemiol* **140**:747-758.

716 Vrancken, B., S. Lequime, K. Theys, and P. Lemey. 2010. Covering all bases in HIV research:  
717 unveiling a hidden world of viral evolution. *AIDS Rev* **12**:89-102.

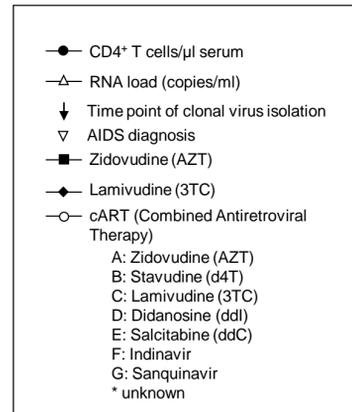
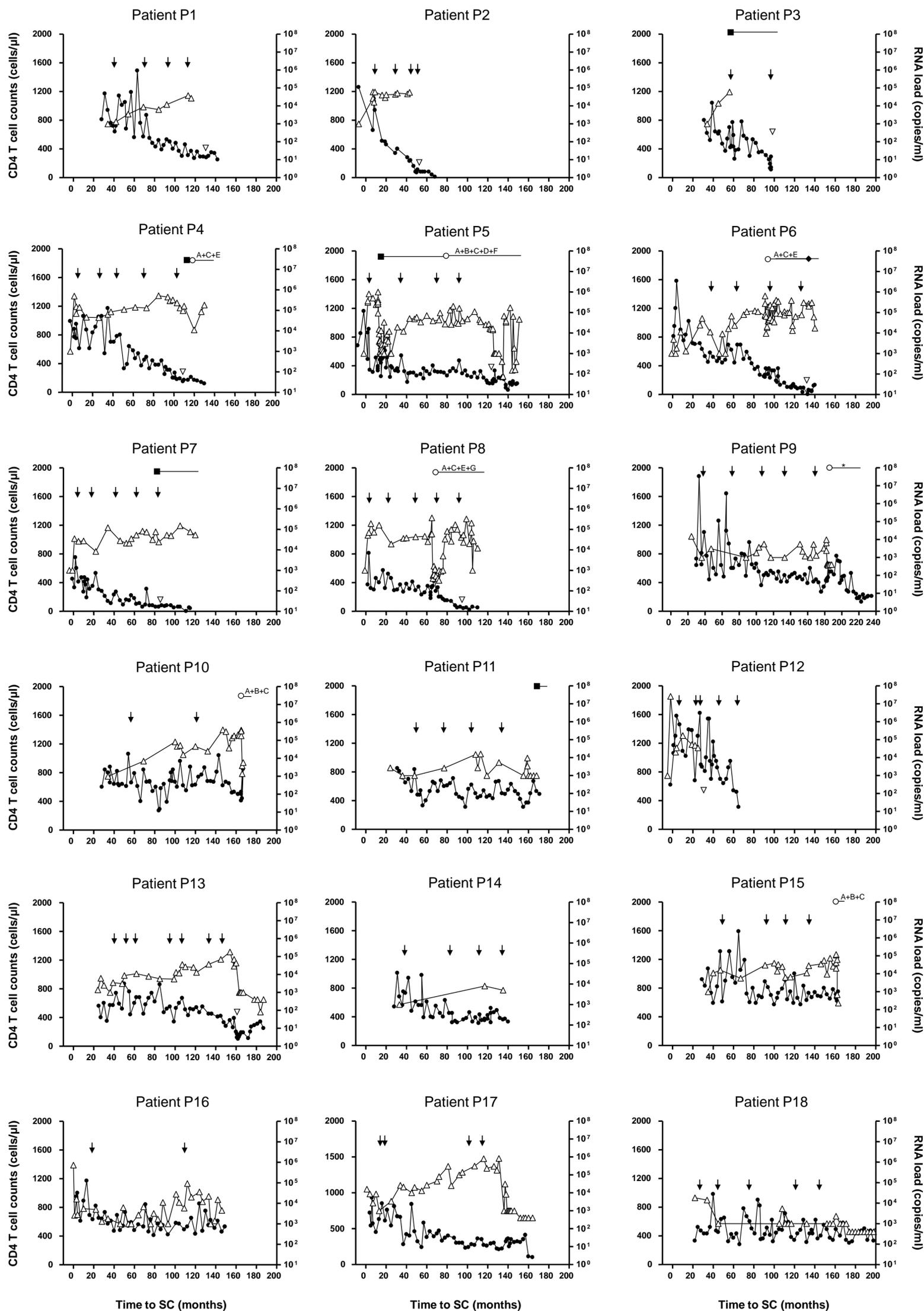
718 Wang, W. K., T. Dudek, M. Essex, and T. H. Lee. 1999. Hypervariable region 3 residues of HIV  
719 type 1 gp120 involved in CCR5 coreceptor utilization: therapeutic and prophylactic  
720 implications. *Proc Natl Acad Sci U S A* **96**:4558-4562.

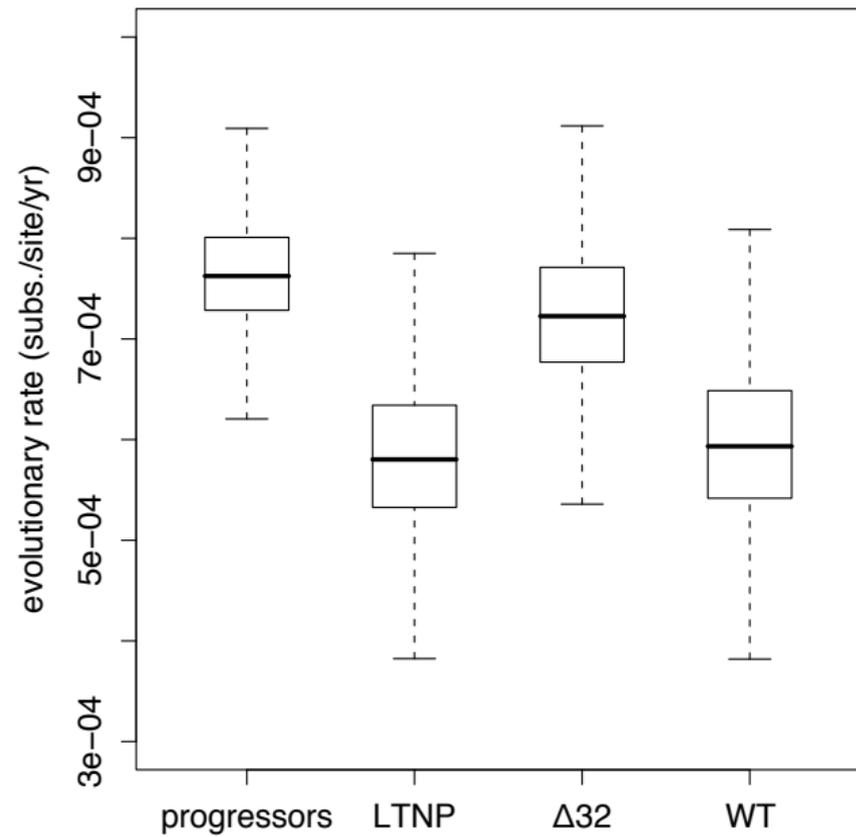
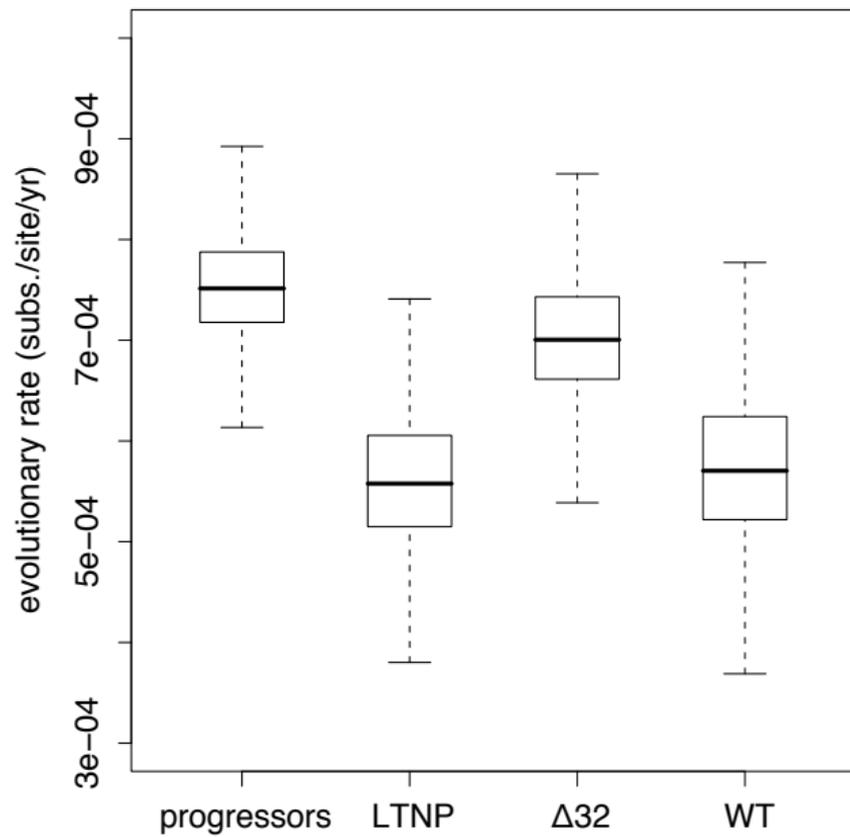
721 Wei, X., J. M. Decker, S. Wang, H. Hui, J. C. Kappes, X. Wu, J. F. Salazar-Gonzalez, M. G.  
722 Salazar, J. M. Kilby, M. S. Saag, N. L. Komarova, M. A. Nowak, B. H. Hahn, P. D.  
723 Kwong, and G. M. Shaw. 2003. Antibody neutralization and escape by HIV-1. *Nature*  
724 **422**:307-312.

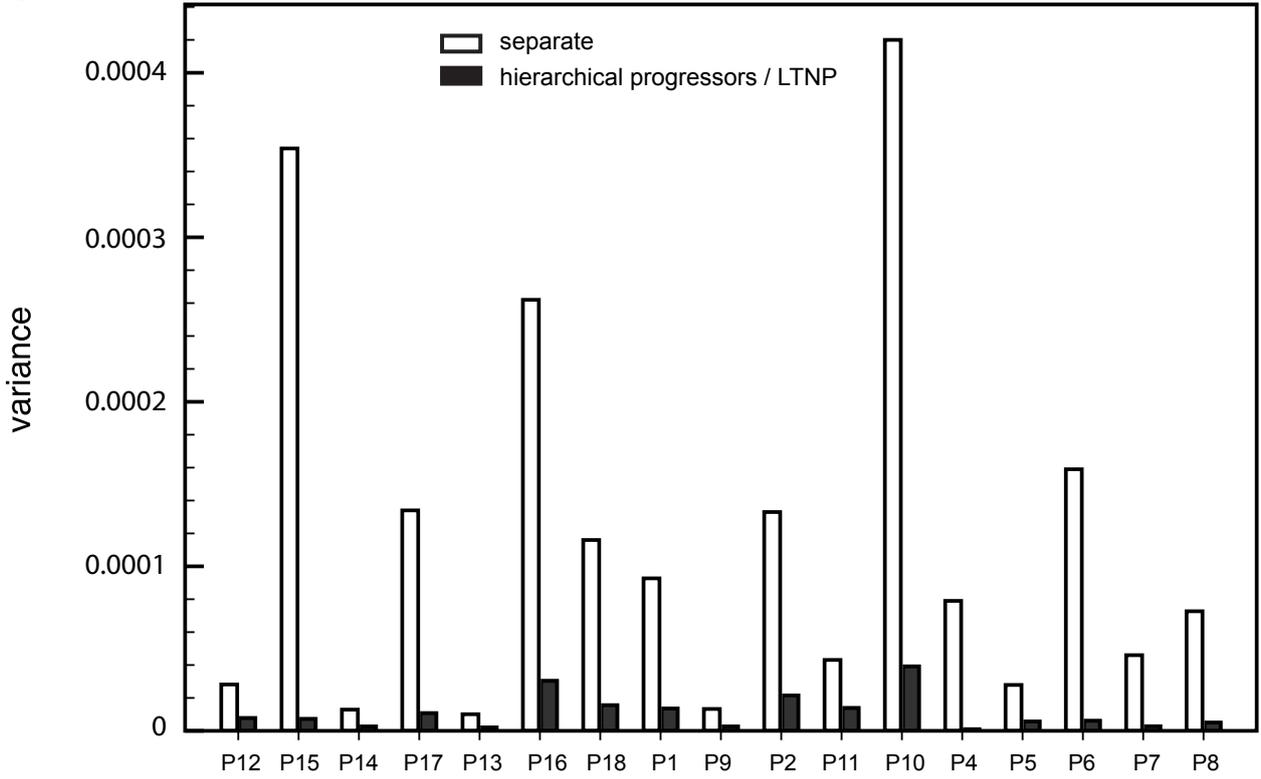
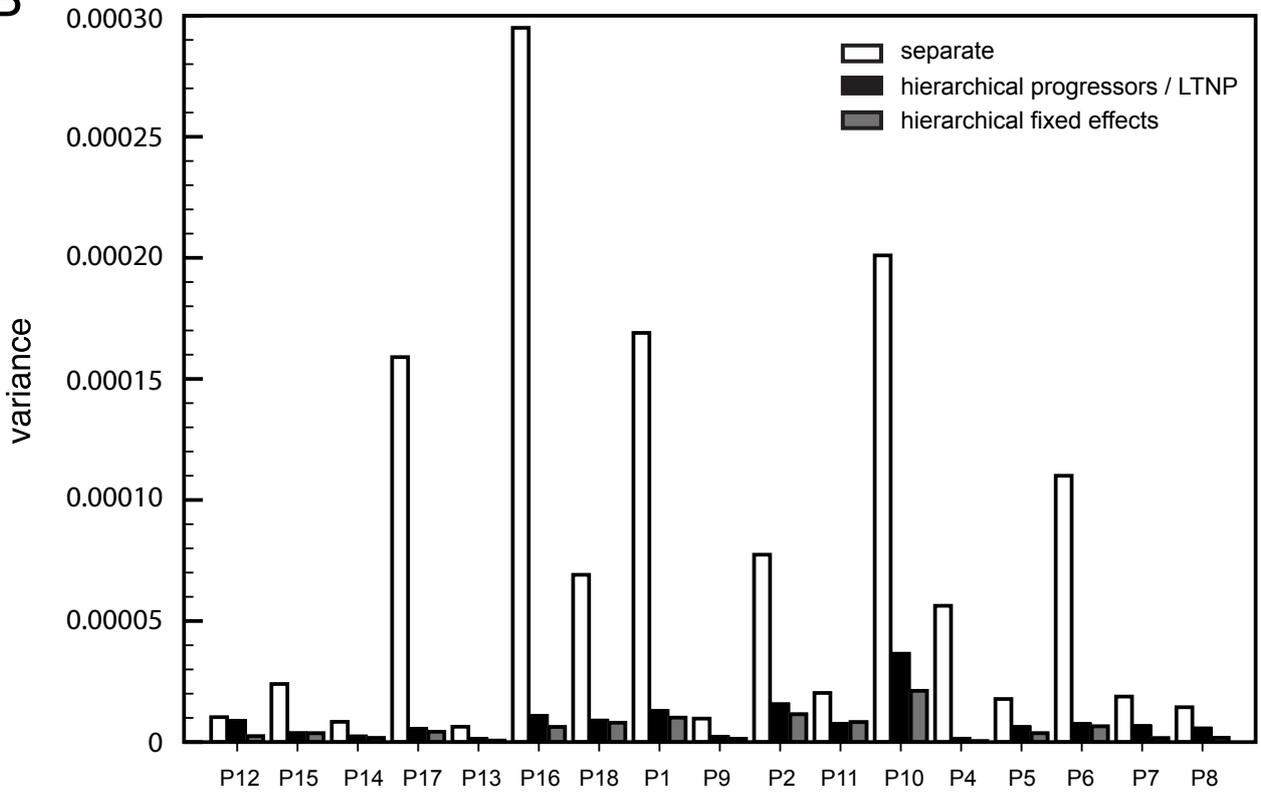
725 Wei, X., S. K. Ghosh, M. E. Taylor, V. A. Johnson, E. A. Emini, P. Deutsch, J. D. Lifson, S.  
726 Bonhoeffer, M. A. Nowak, B. H. Hahn, M. S. Saag, and G. M. Shaw. 1995. Viral  
727 dynamics in human immunodeficiency virus type 1 infection. *Nature* **373**:117-122.

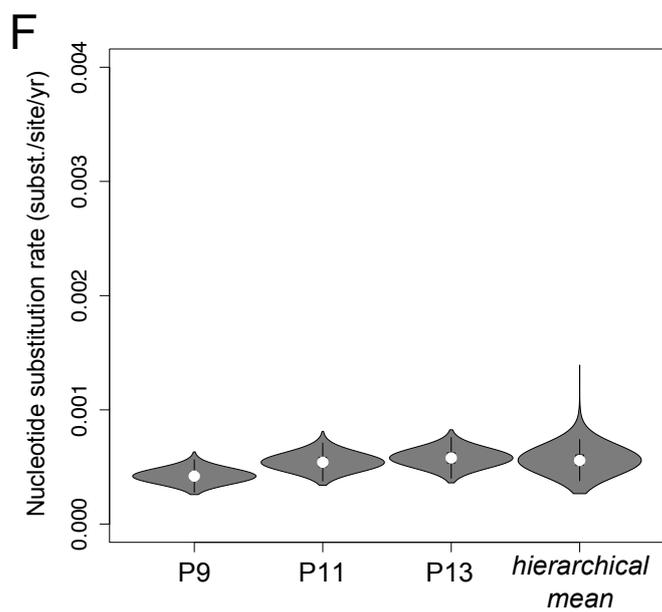
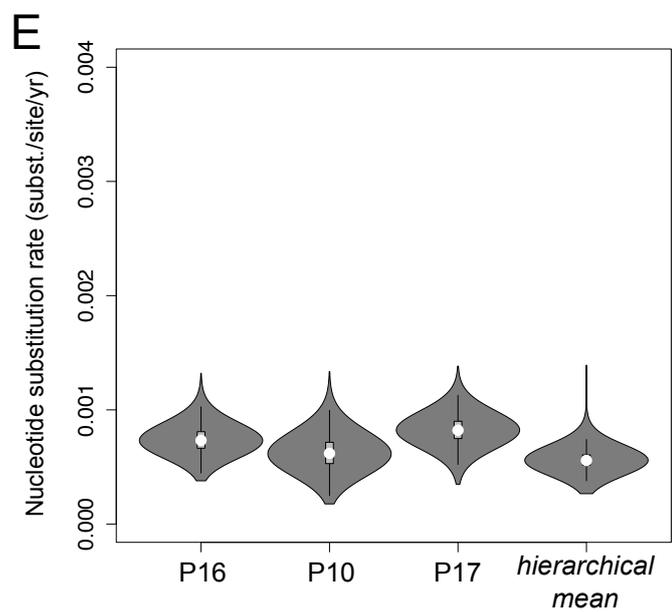
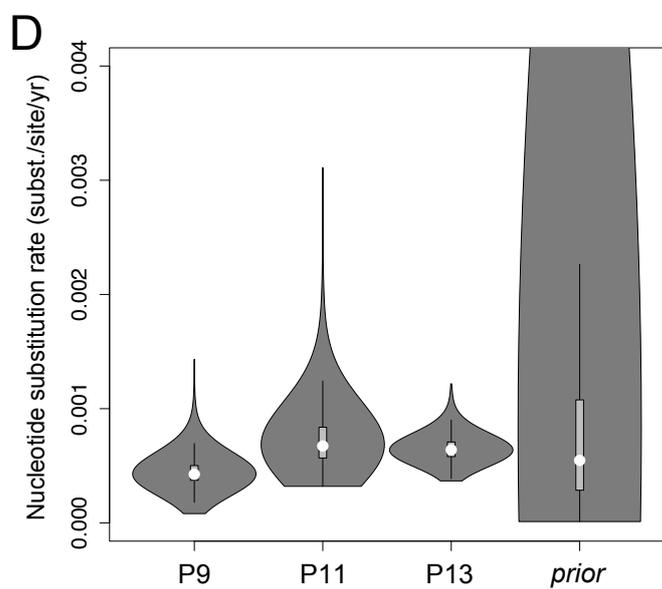
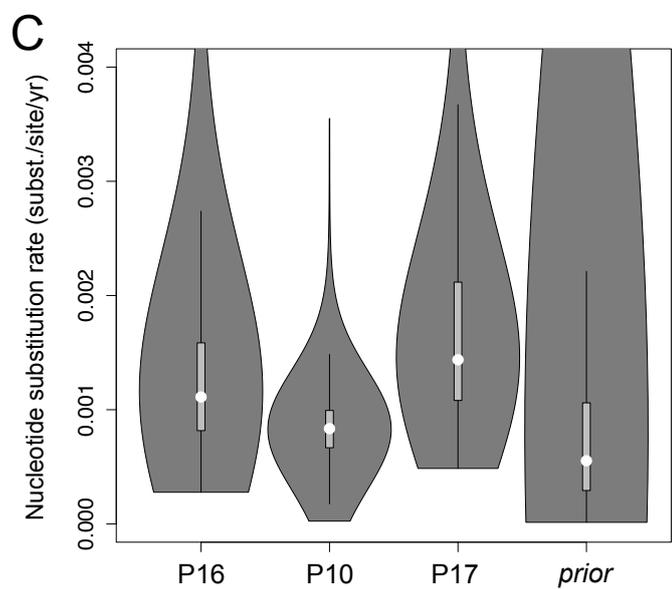
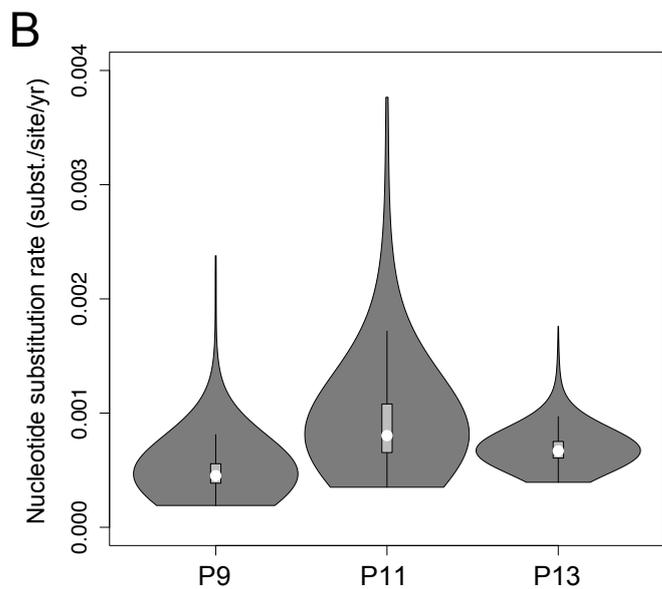
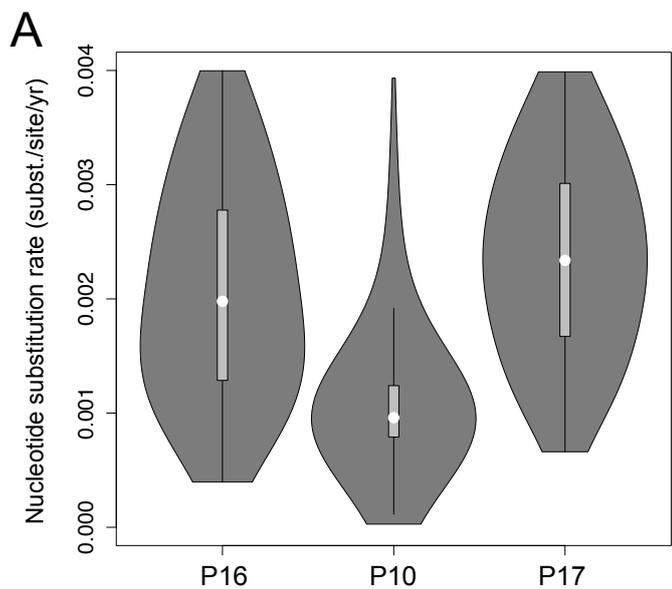
728 Williamson, S. 2003. Adaptation in the env gene of HIV-1 and evolutionary theories of disease  
729 progression. *Mol Biol Evol* **20**:1318-1325.

730  
731



**A****B**

**A****B**



## SUPPLEMENTARY INFORMATION

### SUPPLEMENTARY TABLES

**Table S1. Patients, time points and number of sequences analyzed.**

Subject	Patient number	CCR5 genotype	Disease progression	Sampling time after SC (months)	gp120 env (nr of clones)
P1	19858	WT/WT	P	42*	8
				69*	5
				92*	4
				113*	6
P2	19576	WT/WT	P	7	2
				29	4
				43	5
				51	5
P3	19947	WT/WT	P	56*	3
				98*	3
P4	19999	WT/WT	P	4	14
				26	16
				42	5
				74	7
				107	20
P5	19768	WT/WT	P	2	21
				36	15
				67	17
				93	12
P6	19659	WT/WT	P	2	1
				30	7
				62	22
				95	20
				128	5
P7	19542	WT/WT	P	2	4
				20	5
				43	7
				63	17
				86	15
P8	18969	WT/WT	P	2	25
				22	21
				47	10
				68	7
				91	15
P9	19559	WT/WT	LTNP	39*	3
				71*	5
				106*	1
				133*	5
				170*	3
P10	19932	WT/WT	LTNP	54*	3
				120*	5
P11	19417	WT/WT	LTNP	48*	3
				77*	6
				101*	5
				131*	5
P12	19828	Δ32/WT	P	4	5
				22	2
				25	4

				47	11
				63	4
P13	19383	$\Delta 32$ /WT	LTNP	39*	2
				50*	2
				62*	4
				71*	4
				95*	3
				107*	6
				133*	7
				148*	2
P14	19922	$\Delta 32$ /WT	LTNP	39*	5
				82*	6
				111*	5
				135*	5
P15	19663	$\Delta 32$ /WT	LTNP	47*	5
				91*	6
				111*	6
				140*	5
P16	19984	$\Delta 32$ /WT	LTNP	19	6
				109	4
P17	19566	$\Delta 32$ /WT	LTNP	13	2
				19	7
				101	3
				116	4
P18	19956	$\Delta 32$ /WT	LTNP	28*	5
				51*	3
				78*	2
				123*	1
				146*	1

P: Progressor; LTNP: Long-term non-progressor; SC: seroconversion; \*Sampling time after imputed SC date.

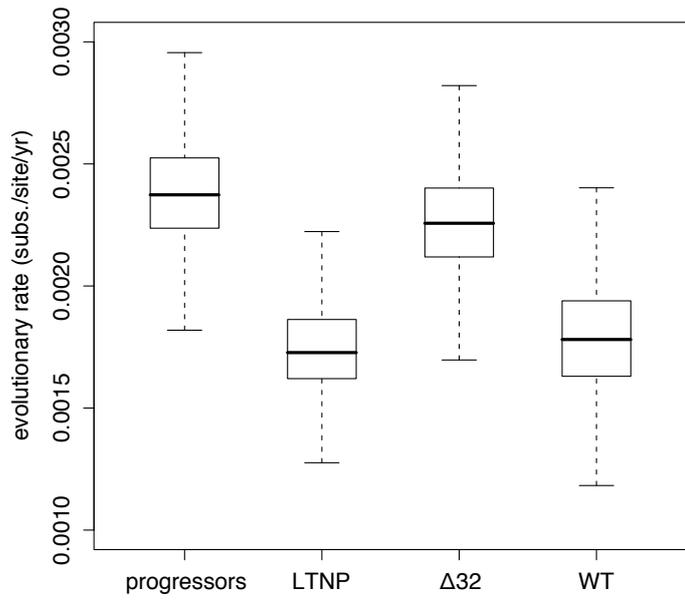
**Table S2. Log marginal likelihood estimates for strict and relaxed clock analyses of four patient groups.**

	Progressors	LTNP	WT	$\Delta 32$
Strict clock	-28534.3	-26762.8	-33501.0	-21793.5
Relaxed clock	-28526.5	-26756.7	-33496.6	-21789.3

WT: CCR5 wt/wt;  $\Delta 32$ : CCR5 wt/ $\Delta 32$ .

## SUPPLEMENTARY FIGURES

**Supplementary Figure 1. Codon model estimates.** Mean evolutionary rate estimated under relaxed clock model with codon model for four patient groups: Progressors, LTNP, CCR5 wt/wt and CCR5 wt/ $\Delta$ 32 (A).  $dN/dS$  rate ratios estimated for the same for patient groups (B). CCR5 wt/wt (WT); CCR5 wt/ $\Delta$ 32 ( $\Delta$ 32).

**A****B**