

## Establishment and stability of the latent HIV-1 DNA reservoir

Johanna Brodin<sup>1</sup>, Fabio Zanini<sup>2,3</sup>, Lina Thebo<sup>1</sup>, Christa Lanz<sup>2</sup>, Göran Bratt<sup>4</sup>, Richard A. Neher<sup>2</sup>, Jan Albert<sup>1,5</sup>

<sup>1</sup>*Department of Microbiology, Tumor and Cell Biology, Karolinska Institute, Stockholm, Sweden*

<sup>2</sup>*Max Planck Institute for Developmental Biology, 72076 Tübingen, Germany*

<sup>3</sup>*current address: Stanford University, Stanford, CA, USA*

<sup>4</sup>*Department of Clinical Science and Education, Venhälsan, Stockholm South General Hospital, Stockholm, Sweden*

<sup>5</sup>*Department of Clinical Microbiology, Karolinska University Hospital, Stockholm, Sweden*

(Dated: May 23, 2016)

HIV-1 infection currently cannot be cured because the virus persists as integrated proviral DNA in long-lived cells despite years of suppressive antiretroviral therapy (ART). To characterize establishment, turnover, and evolution of viral DNA reservoirs we deep-sequenced the p17gag region of the HIV-1 genome from samples obtained from 10 patients after 3-18 years of suppressive ART. For each of these patients, whole genome deep-sequencing data of HIV-1 RNA populations before onset of ART were available from 6-12 longitudinal plasma samples spanning 5-8 years of untreated infection. This enabled a detailed analysis of the dynamics and origin of proviral DNA during ART. A median of 14% (range 0-42%) of the p17gag DNA sequences were overtly defective due to G-to-A hypermutation. The remaining sequences were remarkably similar to previously observed RNA sequences and showed no evidence of evolution over many years of suppressive ART. Most sequences from the DNA reservoirs were very similar to viruses actively replicating in plasma (RNA sequences) shortly before start of ART. The results do not support persistent HIV-1 replication as a mechanism to maintain the HIV-1 reservoir during suppressive therapy. Rather, the data indicate that viral DNA variants are turning over as long as patients are untreated and that suppressive ART halts this turnover.

### Introduction

Combination antiretroviral therapy (ART) has had a dramatic effect on the morbidity and mortality of human immunodeficiency virus type 1 (HIV-1) infection. Even though ART is very effective in suppressing active virus replication, it cannot eradicate the infection because HIV-1 persists as integrated proviral DNA in long-lived cells that constitute a virus reservoir. Latently infected resting memory CD4+ T-lymphocytes (memory CD4 cells) represent the most solidly documented HIV-1 reservoir [1–3]. Thus, a small fraction of memory CD4 cells have fully functional integrated HIV-1 proviruses. These cells do not produce virus when they are in a resting state, but can be induced to produce virus upon activation *in vitro* and *in vivo* [1–4].

Because of their importance for HIV-1 cure efforts, many methods to quantify the HIV-1 reservoirs have been developed. The quantitative virus outgrowth assay (QVOA) represents the “gold standard” [4–6], but this assay underestimates the true size of the functional reservoir due to incomplete induction by PHA stimulation. Ho et al.[7] showed that the functional HIV-1 reservoir may be 60-fold larger than originally estimated. PCR based assays are also commonly used for quantifying the HIV-1 reservoir, but these assays overestimate the size of the functional reservoir because they cannot distinguish between replication-competent and defective viral genomes. Quantification of the HIV-1 reservoir by PCR-based methods typically give at least 100-fold higher numbers than the QVOA because of defective proviruses [4–6]. Many defective proviruses have large internal deletions [7, 8]. Defective proviruses are also the result of

APOBEC editing, which induces G-to-A hypermutation [9–11].

The HIV-1 reservoir is established early during primary infection and is remarkably stable in both quantitative and qualitative terms. Early ART reduces the size and the genetic complexity of the reservoir [12–15]. Siliciano et al.[16] documented a half-life of 44 months for latently infected cells capable of producing replication-competent virus in the QVOA. Similarly, HIV-1 DNA levels and genetic compositions are very stable in patients on long-term suppressive ART [5, 13, 17–21]. Most studies indicate that the HIV-1 reservoir is maintained by the physiological homeostasis of CD4 memory that in part involves occasional expansions and contractions of individual CD4 cell clones [5, 12, 22]. However, some studies have suggested that persistent virus replication may be an important contributor to the maintenance of the HIV-1 reservoir [23, 24]. In particular, Lorenzo-Redondo et al.[25] recently reported evidence of rapid HIV-1 evolution in lymphoid tissue reservoirs.

Despite their significance for HIV-1 cure efforts relatively little is known about the establishment and turnover of the HIV-1 reservoir before start of ART. In this study we have characterized how HIV-1 DNA reservoirs are established and maintained in 10 patients. Evolution of HIV-1 in these patients during 5-8 years prior to ART had been characterized in a recent study by Zanini et al.[26] by whole genome deep-sequencing. The patients were selected to later have gone on to many years of fully suppressive ART. We now sequenced HIV-1 DNA from peripheral blood mononuclear cells (PBMCs) and compared these reservoir sequences to replicating HIV populations prior to ART. The timing of all available samples

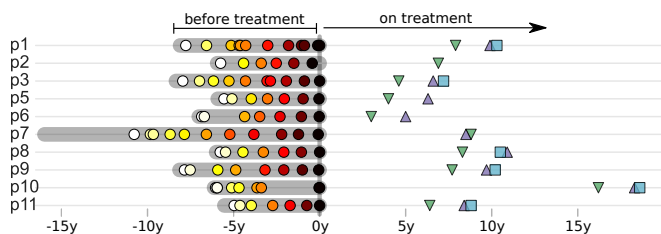


Figure 1: Sampling times before and after start of treatment. For each study participant, the thick grey bar indicates the period of untreated HIV-1 replication, while circles mark the collection times of RNA samples used for whole genome deep sequencing [26]. The collection dates of PBMC samples used for p17gag deep sequencing of integrated provirus are indicated by triangles and squares. All times are relative to start of treatment.

relative to start of treatment is summarized in Fig. 1.

We found that HIV-1 DNA populations remained genetically stable for up to 18 years after start of suppressive ART, which provides evidence against viral evolution and replication as a mechanisms to maintain HIV-1 reservoirs. Furthermore, we found that variants replicating shortly before start of therapy were overrepresented in the HIV-1 DNA reservoirs indicating that proviral HIV-1 variants were turning over as long as the patients were untreated.

## Results

### Patients and samples

The study included 10 HIV-1 infected patients who were diagnosed in Sweden between 1990 and 2003. The patients were selected on the following criteria: 1) A relatively well-defined time of infection; 2) Treatment-naive for a minimum of 5 years; and 3) Thereafter gone on to suppressive ART (plasma HIV-1 RNA levels continuously  $<50$  copies/ml) for a minimum of 2 years. In a recent study we performed whole-genome deep sequencing of replicating HIV-1 RNA populations in 9 of the 10 patients covering the time period before they started ART (6-12 longitudinal plasma samples per patient spanning 5-8 years) [26]. Here we included plasma RNA sequences from the tenth patient. Patient characteristics are summarized in Fig. 1 and Table 1.

For the present study we have obtained sequence data from HIV-1 DNA in viral reservoirs by deep sequencing of the p17gag region of the HIV-1 genome in DNA prepared from PBMC. Patient characteristics are summarized in Table 1. Longitudinal PBMC samples (1-3 samples per patient spanning up to 2.6 years) were obtained obtained 3 - 18 years after start of suppressive ART (Table 1). We define viral DNA reservoirs as HIV p17gag sequences that were still present in PBMC after a minimum of 2 years of suppressive ART. The HIV-1 DNA template numbers were quantified by limiting dilution. Identical p17gag

sequences were merged into haplotypes while preserving their abundance. Minor haplotypes were merged with major haplotypes if they differed by only one mutation (see Materials and Methods). Processed sequence data will be made available at [hiv.tuebingen.mpg.de](http://hiv.tuebingen.mpg.de). Raw sequencing reads from all HIV-1 DNA samples have been deposited in the European Nucleotide archive and will be available under study accession number PRJEB13841 (sample accession numbers ERS1138001-ERS1138025).

### Proviral DNA sequences reflect pretreatment RNA sequences

The HIV-1 DNA sequences recapitulate the diversity observed in RNA sequences before treatment, often with exact sequence matches, see Fig. 2 and Fig. S3. While we observed large variations in the abundance of haplotypes with sequence read frequencies varying between 0.1% and 50% (see Fig. S1), the close match between RNA and DNA sequences confirms that we characterized proviral diversity in a specific and sensitive manner. Variation in haplotype abundance likely reflects clonal expansions [5, 13], independent integrations of identical sequences, and resampling of the same original DNA templates during sequencing. The exact contribution by these distinct mechanisms is difficult dissect in our sequence data.

The estimated number of HIV DNA templates, the number of distinct haplotypes observed, and the fraction of haplotypes seen in multiple samples are given in Table S1. We typically recapture one third (median 0.29) of haplotypes observed at a frequency above 1% in another sample from the same patient.

### Hypermutated sequences are frequent in HIV-1 reservoirs

We found that a substantial proportion (median 14%; range 0-42%) of the p17gag DNA sequences from the viral reservoirs were hypermutated and therefore replication incompetent (see Fig. S2), which is consistent with other reports [5, 6, 10, 13]. A small proportion of sequences had stop codons not obviously due to G-to-A hypermutation (average 3%, range 0-12%). It is likely that a proportion of sequences without overt inactivating mutations were also replication incompetent due to mutations or deletions outside of p17gag.

Hypermutation in the HIV-1 DNA sequences complicates comparison with non-defective DNA and RNA sequences. For this reason we excluded hypermutated sequences from the main analyses, but we also performed complementary analyses that included hypermutated sequences.

Patient	Gender	Transmission	Subtype	Age <sup>a</sup>	HIV RNA from plasma			HIV DNA from PBMCs	
					# samples	first/last since EDI <sup>b</sup>	time on ART <sup>b</sup>	# templates	
p1	F	HET	01_AE	37	12	0.3	8.2	7.9/ 9.9/ 10.4	820/ 148/ 38
p2	M	MSM	B	32	6	0.2	5.5	6.9	75
p3	M	MSM	B	52	10	0.4	8.4	4.6/6.7/ 7.2	243/ 102/ 108
p5	M	MSM	B	38	7	0.4	5.9	4.0/ 6.3	180/ 72
p6	M	HET	C	31	7	0.2	7.0	3.0/ 5.0/ 5.5	115 /15/ nd
p7	M	MSM	B	31	11	6.3 <sup>c</sup>	16.1	6.3/ 8.4/ 8.8	88/ 279/ 108
p8	M	MSM	B	35	7	0.2	6.0	8.4/ 10.6/ 10.9	180/ 55/ 175
p9	M	MSM	B	32	8	0.3	8.1	7.7/ 9.7/ 10.2	60/ 72/ 72
p10	M	MSM	B	34	9	0.1	6.2	16.2/ 18.3/ 18.6	249/ 116/ 51
p11	M	MSM	B	53	7	0.6	5.6	6.4/ 8.4/ 8.8	124/ 120/ 123

Table I: **Summary of patient characteristics.** <sup>a</sup> at diagnosis; <sup>b</sup> EDI: estimated date of infection; all times are given in years; <sup>c</sup> sequencing failed in earlier samples due to low plasma HIV-1 RNA levels.

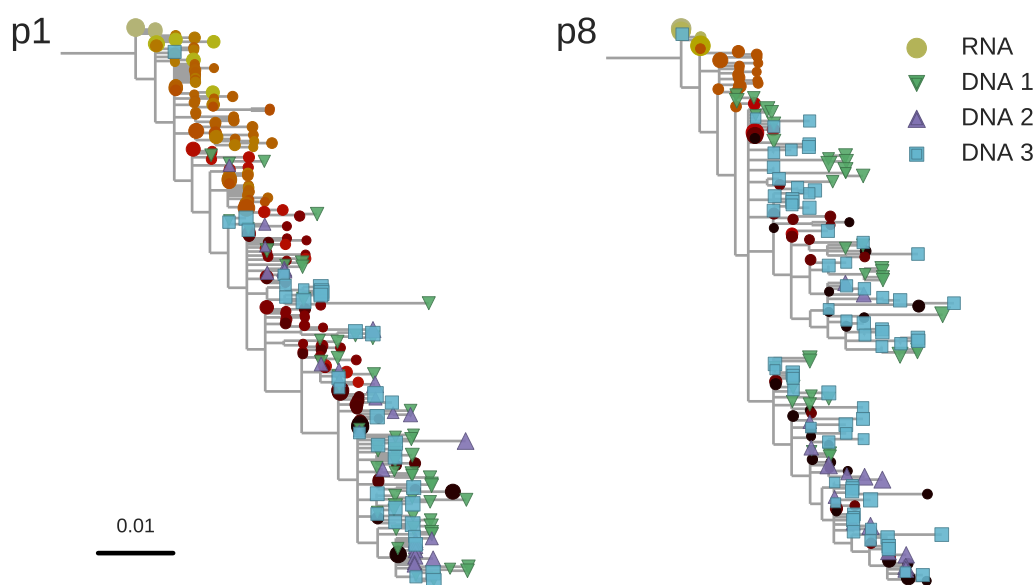


Figure 2: Reconstructed phylogenetic trees of RNA sequences (circles) and DNA (blue squares and triangles) from two patients. RNA sequences are colored by time since infection from yellow to red to black. Size of the symbols indicates the fraction of reads represented by the node. The trees were build using the software FastTree (see methods) [27].

### Lack of evidence of persistent replication in HIV-1 DNA reservoirs

It remains controversial whether or not HIV-1 reservoirs are maintained by persistent replication [5, 12, 20, 22–25]. We used the p17gag DNA sequences from viral reservoirs to search for evidence of sequence evolution, which would be expected to take place if the virus was replicating. The p17gag DNA sequences from viral reservoirs were obtained from 3.0 to 18.6 years after start of suppressive ART. HIV-1 RNA sequences from plasma samples obtained before start of therapy were used as reference materials.

Root-to-tip distances for plasma RNA populations and PBMC DNA populations were calculated relative to the major RNA haplotype in the first plasma sample. Fig. 3 shows temporal changes of root-to-tip distances in HIV-1 RNA and DNA populations obtained before and after start of suppressive ART, respectively. As previously shown, plasma HIV-1 RNA populations obtained before start of ART evolved at a relatively constant rate [26], as evidenced by a steady increase of average root-to-tip distances over time. In sharp contrast, HIV-1 DNA populations obtained after 3 - 18 years of suppressive therapy showed stable root-to-tip distances. Hypermutated DNA sequences showed larger root-to-tip distances, but

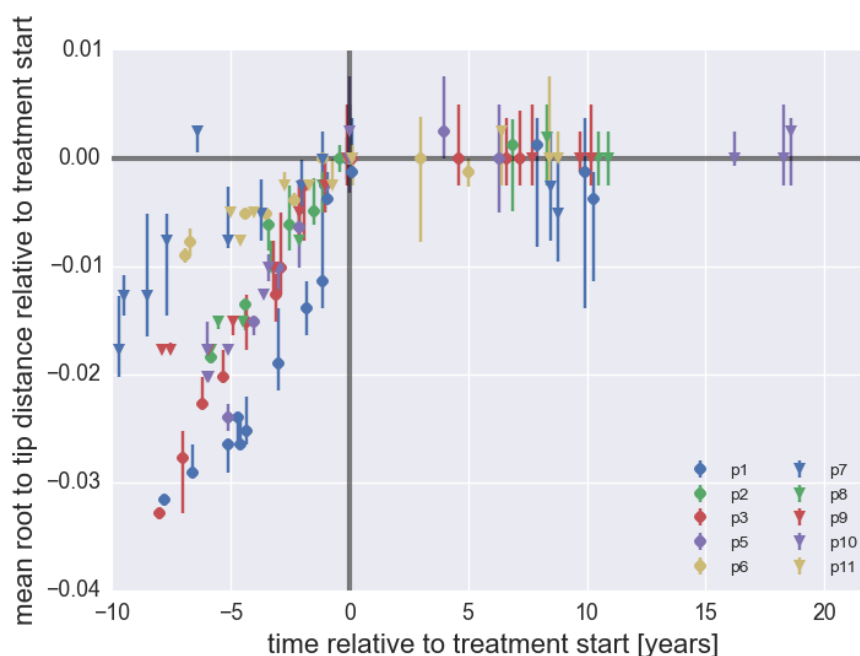


Figure 3: **Root-to-tip distances.** Plasma HIV-1 RNA sequences evolve steadily before start of ART, while no evolution is observed in PBMC HIV-1 DNA sequences obtained after start of ART. The figure contains data on DNA sequence not classified as hypermutants, the analogous figure for hypermutants is shown in Fig. S4. The error bars indicated the inter-quartile range of the root to tip distance.

also these distances were stable over time (Fig. S4)

Table 2 shows the rate of evolution before and after start of suppressive ART. Before start of therapy we observed statistically significant evolution of plasma RNA sequences with rates 1 to  $4 \times 10^{-3}$ /year in all 10 patients. In contrast, no statistically significant evolution was observed in DNA reservoirs during up to 18 years of suppressive ART.

Collectively, our results do not provide support for persistent HIV-1 replication as a mechanism to maintain the HIV-1 reservoir during suppressive therapy.

### Time for deposition of reservoir DNA sequences

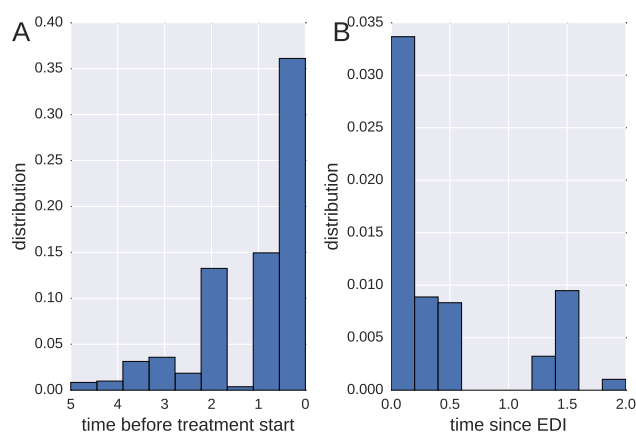
To investigate when PBMC HIV-1 DNA variants had been deposited the viral reservoirs we compared the on-treatment PBMC DNA sequences with the longitudinal pre-treatment plasma RNA sequences (Fig. 4 and Fig. S5). For each p17gag DNA sequence, we determined the RNA sample and haplotype that was the most likely source. By this procedure, most proviral sequences were assigned to the plasma samples closest to the start of treatment (Fig. 4, panel A). The representation of RNA haplotypes from earlier plasma samples dropped rapidly over 1-2 years. However, sequences representing earlier plasma sampling time points were also found as minor variants among the p17gag DNA sequences. Among these minor variants, DNA sequences matching the vari-

patient	RNA rate		DNA rate	
	[year <sup>-1</sup> ]	p-value	[year <sup>-1</sup> ]	p-value
p1	$4.2 \times 10^{-3}$	$< 10^{-6}$	$1 \times 10^{-4}$	0.866
p2	$3.3 \times 10^{-3}$	$< 10^{-3}$	$2 \times 10^{-4}$	–
p3	$4.2 \times 10^{-3}$	$< 10^{-6}$	0	1.0
p5	$4.6 \times 10^{-3}$	$< 10^{-3}$	$1 \times 10^{-4}$	0.93
p6	$1.3 \times 10^{-3}$	$< 10^{-3}$	$-6 \times 10^{-4}$	0.67
p7	$1.4 \times 10^{-3}$	$< 10^{-2}$	$-8 \times 10^{-4}$	0.3
p8	$3.0 \times 10^{-3}$	$< 10^{-4}$	$2 \times 10^{-5}$	0.94
p9	$2.4 \times 10^{-3}$	$< 10^{-4}$	$3 \times 10^{-7}$	0.87
p10	$3.5 \times 10^{-3}$	$< 10^{-5}$	$-1 \times 10^{-4}$	0.55
p11	$1.1 \times 10^{-3}$	$< 10^{-2}$	$6 \times 10^{-5}$	0.6

Table II: Rates of evolution in plasma RNA and PBMC DNA sequences obtained before start and after start of suppressive ART, respectively.

ant dominating the first plasma samples obtained within 6 month after infection were overrepresented in some patients (Fig. 4, panel B). Thus, the proportion of DNA sequences matching the initial plasma HIV-1 RNA haplotype was 14%, 2.4%, 42%, <1%, and 6.9% in patients 2, 3, 6, 8 and 11, respectively.

The overrepresentation in the HIV-1 DNA reservoir of viral variants that replicated shortly before start of suppressive ART indicates that cells carrying defective and non-defective proviral variants were turning over as



**Figure 4: Probable origin of sequences in the DNA reservoir.** Panel A shows the time difference between the start of treatment and the RNA sample that is most likely to have generated at sequence from the reservoir. A large proportion of the reservoir DNA sequences are most closely related to RNA sequences in the plasma samples obtained during the last year before start of therapy. Panel B focuses on DNA sequences resembling RNA sequences obtained during the first two years after estimated time of infection. These sequences represent a small fraction of all reads obtained and are observed in 50% of patients only.

long as the patients were untreated and that suppressive ART halted this turnover. The width of the peak in Fig. 4A suggest a half-life on the order of one year.

## Discussion

In this study we have investigated the composition and turnover of HIV-1 DNA sequences in viral reservoirs in patients on long-term suppressive therapy. Reservoir HIV-1 DNA populations were remarkably stable and showed no signs of ongoing replication. We also traced when during the course of HIV infection viruses in the the DNA reservoirs had been deposited and found that they mainly derived from the last year(s) before start of suppressive therapy.

Our study provides evidence against persistent HIV-1 replication as a mechanism for maintenance of HIV-1 reservoirs during suppressive therapy. This is at variance with a recent report by Lorenzo-Redondo et al. [25] and a few earlier reports [23, 24], but agrees with several other earlier studies [5, 13, 17–21]. Lorenzo-Redondo et al. [25] compared genetic diversity in samples HIV-1 RNA in plasma at start of therapy with HIV-1 DNA sequences obtained from blood and tissues at baseline, three and six months after the start of treatment. They report a signal of evolution between the different time points at an extraordinary high rate ( $7.4 - 12 \times 10^{-3}$  changes per site per year) – about 5-fold higher than typically observed in gag and pol of replicating RNA populations. This observa-

tion is incompatible with the lack of observable changes in reservoir sequences over 20 times longer time intervals reported here. Without longitudinal data on the evolution of the HIV-1 population prior to treatment, the nature of the change reported by Lorenzo-Redondo et al. [25] is difficult to discern. One possible explanation for the apparent conflict between the two studies could be that after onset of therapy, short-lived cells that sampled the most recent circulating virus populations start to disappear, leaving longer-lived cells that sample deeper into the history of the infection. This scenario would not correspond to evolution, but quite oppositely a sampling of earlier variants. Another difference is that Lorenzo-Redondo et al. [25] investigated HIV-1 DNA sequences in tissue as well as PBMC samples whereas we only studied PBMC samples. However, tissue and blood HIV-1 DNA variants should be well-mixed over the time frame that we investigate [5, 13, 25]. Both Lorenzo-Redondo et al. [25] and we studied DNA sequences in HIV-1 reservoirs, which are known to contain a high proportion of defective virus. These proviruses serve as markers of T-cell clones, rather than replication-competent virus. Hence the absence of evolution or turnover of provirus that we found does not exclude the possibility that there is replication and evolution of replication-competent virus in reservoirs. However, if such replication exists, it happens very low levels that do not contribute substantially to the pool of proviral DNA in PBMCs. To enrich for replication competent and putatively evolving virus, QVOA followed by sequencing of virus released into to supernatants should performed, rather than sequencing of total HIV-1 DNA as done by Lorenzo-Redondo et al. [25], us and others [10, 13, 20]. In agreement with our finding of genetic stability in the DNA reservoirs Josefsson et al. [13] and Stockenstrom et al. [5] have reported that defective HIV-1 DNA integrants present during long-term effective ART appear to be maintained by proliferation and longevity of infected cells rather than by ongoing viral replication.

Because we had access to detailed longitudinal data on the evolution of the plasma HIV-1 RNA population from time of infection to start of suppressive ART, we could trace when during the course of untreated HIV-1 infection the viruses in the DNA reservoirs had been deposited. We found that a majority of variants in the HIV-1 DNA reservoirs were derived from HIV-1 RNA variants that had actively replicated during the last year(s) before start of suppressive ART, with no evidence for evolution after treatment start. Frenkel et al.[28], in contrast to us, reported persistence of a greater number of early compared to recent viruses in a few children on suppressive ART; more research is warranted to assess the origin of this difference.

Defective HIV-1 proviruses can be regarded as unique in vivo labels of individual memory CD4 cell clones which can be used to track their fate similar to sequencing of T-cell receptors [29]. This strategy was used by Imamichi et al.[30] to demonstrate that a T-cell clone persisted

more than 17 years. Similarly, prenatally formed T-cell receptors shared by twins have been reported to have lifetimes > 30 years [31]. During suppressive ART the turnover of infected memory CD4 cell clones is likely to follow the same dynamics as in uninfected people. In contrast, we observe a strong overrepresentation in the reservoirs of “late” HIV-1 RNA variants, which indicates that HIV-1 target cells, primarily CD4+ T-lymphocytes, were turning over with a half-life of about one year in absence of treatment. This turnover was dramatically slowed by suppressive ART. Earlier studies, based on different types of labelling of CD4 cells, have indicated a 3 - 4 fold increased rate of CD4 cell death in untreated HIV-1-infected patients as compared with uninfected persons and patients on suppressive ART [32–34]. The more dramatic difference we observe is likely explained by different methodologies. Earlier studies estimated the lifespan of individual cells whereas we primarily have estimated the lifespan of CD4 cell clones carrying defective proviruses (i.e. infected cells as well as their daughter cells).

Our study has several limitations. We have not sorted cells and therefore cannot investigate if there are differences in HIV-1 turnover between different types and subsets of cells, such as memory CD4 cells and their subsets. However, it is reasonable to assume that a majority of our HIV-1 DNA sequences came from memory CD4 cells because others have shown that these cells constitute the main HIV-1 reservoir [1–3]. We sequenced a relatively short region of the HIV-1 genome and therefore cannot reliably distinguish between replication-competent and defective viruses. While we observe no evolution in these proviral DNA sequences, we cannot rule out the possibility that a small subset of viruses indeed was replicating but remained undetected among the many replication-incompetent viruses. We observed large variations in the abundance of sequence haplotypes that likely reflect both clonal expansions [5, 13], independent integrations of identical sequences, and resampling of the same original DNA templates during sequencing. With our sequencing method we could not exactly determine the relative contribution by these distinct mechanisms. We are attempting Primer ID sequencing [35] to even better understand the in vivo dynamics of different viral haplotypes.

In summary, we provide compelling evidence against persistent viral replication as a mechanism to maintain the latent HIV-1 DNA reservoir during suppressive therapy. Furthermore, we show that most latently infected cells during long-term suppressive ART are infected shortly before ART start and that the rate of T-cell turnover is reduced upon starting suppressive ART.

## Materials and methods

*Ethical statement* The study was conducted according to the Declaration of Helsinki. Ethical approval was granted by the Regional Ethical Review board in Stock-

holm, Sweden (Dnr 2012/505 and 2014/646). Patients participating in the study gave written and oral informed consent to participate.

*Patients* The study included 10 HIV-1-infected patients who were diagnosed in Sweden between 1990 and 2003. Prior to the present study the patients were included in a recent study on the population genomics of inpatient HIV-1 evolution [26]. The patients were selected based on the following inclusion criteria: 1) A relatively well-defined time of infection based on a negative HIV antibody test less than two years before a first positive test or a laboratory documented primary HIV infection; 2) No ART during a minimum of approximately five years following diagnosis; 3) Availability of biobank plasma samples covering this time period; and 4) Later have started successful ART (plasma viral levels < 50 copies/ $\mu$ l) for a minimum of two years. As previously described 6 - 12 plasma samples per patient were retrieved from biobanks and used for full-genome HIV-1 RNA sequencing [26]. The same patient nomenclature is used in both studies. For the present study the same patient were asked to donate 70 ml of fresh EDTA-treated blood on up to three occasions over a time period of 2.5 years. These blood samples were obtained 3 - 18 years after start of successful ART. Estimated time of infection (ETI) was calculated as previously described using clinical and laboratory findings including Fiebig staging and BED testing [26]. Information about the patients and the samples are summarized in Table 1.

*HIV-1 RNA sequencing from plasma* Whole-genome deep-sequencing of virus RNA populations in plasma samples obtained before start of therapy was performed as previously described [26]. In short, total RNA in plasma was extracted using RNeasy® Lipid Tissue Mini Kit (Qiagen Cat No. 74804) and amplified using a one-step RT-PCR with outer primers for six overlapping regions and Superscript® III One-Step RT-PCR with Platinum® Taq High Fidelity High Enzyme Mix (Invitrogen, Carlsbad, California, US). An optimized Illumina Nextera XT library preparation protocol was used together with a kit from the same supplier to build DNA libraries, which were sequenced on the Illumina MiSeq instrument with 2 x 250bp or 2x 300bp sequencing kits (MS-102-2003/MS-10-3003). For the present study a part of the p17gag region of the HIV-1 genome (see below) was extracted from the entire full-genome RNA data set. The median number of high quality reads covering the entire p17 sequence was 146 (inter-quartile range 56 - 400) and the cDNA template numbers are available in Zanini et al.[26].

*HIV-1 DNA sequencing from PBMCs* Approximately 70 ml of fresh whole blood was obtained in 7 Vacutainer tubes with EDTA as anticoagulant. PMBC were isolated by Ficoll-Paque PLUS (GE Healthcare BioSciences AB, Uppsala, Sweden) centrifugation according to the instructions by the manufacturer. Total DNA was extracted from PBMC using the OMEGA E.Z.N.A.® Blood DNA Mini Kit (Omega bio-tek, Norcross, Geor-

gia) or the QIAamp DNA Blood Mini Kit (Qiagen GmbH, Hilden, Germany) according to the instructions by the manufacturer. The amount of DNA was measured with Qubit® dsDNA HS Assay Kit (Invitrogen™ Eugene, Oregon, USA). Patient-specific nested primers (Integrated DNA Technologies) were used to amplify a 387-bp long portion of the p17gag gene corresponding to positions 787 to 1173 in the HxB2 reference sequence. The primers were designed based on the plasma RNA sequences from each patient (Tab. S2). Outer primers were used together with Platinum® Taq DNA Polymerase High Fidelity (Invitrogen™ Carlsbad, California, US) for the first PCR. The program started with a denaturation step at 94°C for 2 min followed by 15 PCR cycles of denaturation at 94°C for 20 s, annealing at 50°C for 20 s and extension at 72°C for 30 s and a final extension step at 72°C for 6 min. For the second PCR, 2.5 µl of the product from the first PCR was amplified with inner primers and the cycle profile and enzyme as for the first PCR. Amplified DNA was purified using Agencourt AMPure XP (Beckman Coulter Beverly, Massachusetts) and quantified using Qubit. For each sample the number of HIV-1 DNA templates used for sequencing was roughly quantified in triplicate by limiting dilution using the same PCR conditions, three dilutions (usually 0.5, 0.1, 0.02 µg of DNA) and Poisson statistics. Control experiments were performed to evaluate PCR-induced recombination using the plasmids NL4-3 and SF162, which were spiked in equal proportion into human DNA and amplified using the same PCR conditions as above. The results showed that there was minimal PCR-induced recombination in this short amplicon.

*Sequencing and read processing* The HIV specific primers were flanked by NexteraXT adapters. To construct sequencing libraries, indices and sequencing primers were added in 12-15 cycles of additional PCR. Amplicons were sequenced on an Illumina MiSeq machine with 2x250 cycle kits. Between 6,500 and 190,000 (median 35,000) paired-end reads were generated per sample. The overlapping paired-end sequencing reads were merged to create synthetic reads spanning the entire p17 amplicon. In case of disagreement between paired reads, the nucleotide on the read with the higher quality score was used. We counted the number of times a particular p17 sequence was observed and did subsequent analysis with read-abundance pairs. To reduce the influence of

sequencing and PCR errors, we combined rare sequences (below frequency 0.002) with common sequences if they differed by no more than one position. Specifically, starting with the rarest sequences, we merged rare sequences with the most common sequence that was one base away. The cutoff 0.002 is the typical error frequency of the pipeline as determined earlier [26]. All analysis is done in Python using the libraries numpy, biopython, and matplotlib [36–38].

*Hypermutation detection* To classify sequences into obvious hypermutants and sequences representative of circulating virus, we counted mutations at positions that are not variable in the RNA samples obtained prior to therapy. If more than 4 mutations were observed and at least half of them were G→A, the sequence was considered a hypermutant. The distribution of the different transition mutations relative to the closest genome found in RNA samples are shown in Fig. S2 for reads classified as hypermutants or not. Results we obtained for sequences classified as non-hypermutants are very similar to results obtained when using only sequences without stop codons.

*Phylogenetic analysis* We reconstructed phylogenetic trees using the approximate maximum likelihood method implemented by FastTree [27]. Tips were annotated with frequency, source and sample date using custom python scripts.

*Statistical analysis* Root-to-tip distances were calculated as the average distance between a sample and the founder sequence approximated by the consensus sequence of the first RNA sample. To determine the rate of evolution in absence of treatment, this root-to-tip distance was regressed against time. To determine the rate of evolution on treatment, the root-to-tip sequence of the last RNA sample and the DNA samples was regressed against time. To determine the most likely seeding time for a p17gag DNA sequence obtained from PBMCs, we calculated the likelihood of sampling this sequence given the SNP frequencies in each RNA sample and assigned the sequence to the sample where this likelihood was highest.

*Acknowledgements* This work was supported by the European Research Council through grant Stg. 260686 and the Swedish Research Council through grant K2014-57X-09935. We would also like to express our gratitude to the study participants.

- 
- [1] S. Eriksson, E. H. Graf, V. Dahl, M. C. Strain, S. A. Yukl, E. S. Lysenko, R. J. Bosch, J. Lai, S. Chioma, F. Emad, et al., *PLoS Pathog* **9**, 1 (2013).
  - [2] T.-W. Chun, L. Carruth, D. Finzi, X. Shen, J. A. DiGiuseppe, H. Taylor, M. Hermankova, K. Chadwick, J. Margolick, T. C. Quinn, et al., *Nature* **387**, 183 (1997).
  - [3] T.-W. Chun, D. Finzi, J. Margolick, K. Chadwick, D. Schwartz, and R. F. Siliciano, *Nat Med* **1**, 1284 (1995), ISSN 1078-8956.
  - [4] M. Massanella and D. D. Richman, *Journal of Clinical Investigation* **126**, 464 (2016), ISSN 00219738.
  - [5] S. von Stockenstrom, L. Odevall, E. Lee, E. Sinclair, P. Bacchetti, M. Killian, L. Epling, W. Shao, R. Hoh, T. Ho, et al., *Journal of Infectious Diseases* **212**, 596 (2015).
  - [6] K. M. Bruner, N. N. Hosmane, and R. F. Siliciano, *Trends in microbiology* **23**, 192 (2015).
  - [7] Y.-C. Ho, L. Shan, N. N. Hosmane, J. Wang, S. B.

- Laskey, D. I. S. Rosenbloom, J. Lai, J. N. Blankson, J. D. Siliciano, and R. F. Siliciano, *Cell* **155**, 540 (2013), ISSN 1097-4172.
- [8] G. Sanchez, X. Xu, J. C. Chermann, and I. Hirsch, *J. Virol.* **71**, 2233 (1997), ISSN 0022-538X.
- [9] Q. Yu, R. König, S. Pillai, K. Chiles, M. Kearney, S. Palmer, D. Richman, J. M. Coffin, and N. R. Landau, *Nat. Struct. Mol. Biol.* **11**, 435 (2004), ISSN 1545-9993.
- [10] T. L. Kieffer, P. Kwon, R. E. Nettles, Y. Han, S. C. Ray, and R. F. Siliciano, *J. Virol.* **79**, 1975 (2005), ISSN 0022-538X.
- [11] K. Stopak, C. de Noronha, W. Yonemoto, and W. C. Greene, *Mol. Cell* **12**, 591 (2003), ISSN 1097-2765.
- [12] N. Chomont, M. El-Far, P. Ancuta, L. Trautmann, F. A. Procopio, B. Yassine-Diab, G. Boucher, M.-R. Boulassel, G. Ghattas, J. M. Brenchley, et al., *Nat Med* **15**, 893 (2009), ISSN 1078-8956.
- [13] L. Josefsson, S. v. Stockenstrom, N. R. Faria, E. Sinclair, P. Bacchetti, M. Killian, L. Epling, A. Tan, T. Ho, P. Lemey, et al., *PNAS* **110**, E4987 (2013), ISSN 0027-8424, 1091-6490.
- [14] F. Lori, H. Jessen, J. Lieberman, D. Finzi, E. Rosenberg, C. Tinelli, B. Walker, R. F. Siliciano, and J. Lisziewicz, *J Infect Dis.* **180**, 1827 (1999), ISSN 0022-1899, 1537-6613.
- [15] M. C. Strain, S. J. Little, E. S. Daar, D. V. Havlir, H. F. Günthard, R. Y. Lam, O. A. Daly, J. Nguyen, C. C. Ignacio, C. A. Spina, et al., *J Infect Dis.* **191**, 1410 (2005), ISSN 0022-1899, 1537-6613.
- [16] J. D. Siliciano, J. Kajdas, D. Finzi, T. C. Quinn, K. Chadwick, J. B. Margolick, C. Kovacs, S. J. Gange, and R. F. Siliciano, *Nat Med* **9**, 727 (2003), ISSN 1078-8956.
- [17] G. J. Besson, C. M. Lalama, R. J. Bosch, R. T. Gandhi, M. A. Bedison, E. Aga, S. A. Riddler, D. K. McMahon, F. Hong, and J. W. Mellors, *Clin. Infect. Dis.* **59**, 1312 (2014), ISSN 1537-6591.
- [18] M. F. Kearney, J. Spindler, W. Shao, S. Yu, E. M. Anderson, A. O'Shea, C. Rehm, C. Poethke, N. Kovacs, J. W. Mellors, et al., *PLoS Pathog* **10**, e1004010 (2014).
- [19] H. F. Günthard, S. D. W. Frost, A. J. Leigh-Brown, C. C. Ignacio, K. Kee, A. S. Perelson, C. A. Spina, D. V. Havlir, M. Hezareh, D. J. Looney, et al., *J. Virol.* **73**, 9404 (1999), ISSN 0022-538X, 1098-5514.
- [20] T. H. Evering, S. Mehandru, P. Racz, K. Tenner-Racz, M. A. Poles, A. Figueroa, H. Mohri, and M. Markowitz, *PLoS Pathog.* **8**, e1002506 (2012), ISSN 1553-7374.
- [21] T. L. Kieffer, M. M. Finucane, R. E. Nettles, T. C. Quinn, K. W. Broman, S. C. Ray, D. Persaud, and R. F. Siliciano, *J Infect Dis* **189**, 1452 (2004).
- [22] N. Chomont, S. DaFonseca, C. Vandergaeten, P. Ancuta, and R.-P. Sékaly, *Curr Opin HIV AIDS* **6**, 30 (2011), ISSN 1746-6318.
- [23] M. J. Buzón, M. Massanella, J. M. Llibre, A. Esteve, V. Dahl, M. C. Puertas, J. M. Gatell, P. Domingo, R. Paredes, M. Sharkey, et al., *Nat Med* **16**, 460 (2010), ISSN 1078-8956.
- [24] S. A. Yukl, A. Shergill, K. McQuaid, S. Gianella, H. Lampiris, C. B. Hare, M. Pandori, E. Sinclair, H. F. Günthard, M. Fischer, et al., *AIDS* **24**, 2451 (2010), ISSN 0269-9370.
- [25] R. Lorenzo-Redondo, H. R. Fryer, T. Bedford, E.-Y. Kim, J. Archer, S. L. Kosakovsky Pond, Y.-S. Chung, S. Penugonda, J. G. Chipman, C. V. Fletcher, et al., *Nature* **530**, 51 (2016), ISSN 0028-0836.
- [26] F. Zanini, J. Brodin, L. Thebo, C. Lanz, G. Bratt, J. Albert, and R. A. Neher, *eLife Sciences* **4**, e11282 (2016), ISSN 2050-084X.
- [27] M. N. Price, P. S. Dehal, and A. P. Arkin, *PLoS ONE* **5**, e9490 (2010), ISSN 1932-6203.
- [28] L. M. Frenkel, Y. Wang, G. H. Learn, J. L. McKernan, G. M. Ellis, K. M. Mohan, S. E. Holte, S. M. D. Vange, D. M. Pawluk, A. J. Melvin, et al., *J. Virol.* **77**, 5721 (2003), ISSN 0022-538X, 1098-5514.
- [29] H. Robins, *Current Opinion in Immunology* **25**, 646 (2013), ISSN 0952-7915.
- [30] H. Imamichi, V. Natarajan, J. W. Adelsberger, C. A. Rehm, R. A. Lempicki, B. Das, A. Hazen, T. Imamichi, and H. C. Lane, *AIDS* **28**, 1091 (2014), ISSN 0269-9370.
- [31] M. V. Pogorelyy, Y. Elhanati, Q. Marcou, A. L. Sycheva, E. A. Komech, V. I. Nazarov, O. V. Britanova, D. M. Chudakov, I. Z. Mamedov, Y. B. Lebedev, et al., arXiv:1602.03063 [q-bio] (2016), arXiv: 1602.03063.
- [32] M. Hellerstein, M. B. Hanley, D. Cesar, S. Siler, C. Pappageorgopoulos, E. Wieder, D. Schmidt, R. Hoh, R. Neese, D. Macallan, et al., *Nat. Med.* **5**, 83 (1999), ISSN 1078-8956.
- [33] J. M. McCune, M. B. Hanley, D. Cesar, R. Halvorsen, R. Hoh, D. Schmidt, E. Wieder, S. Deeks, S. Siler, R. Neese, et al., *J. Clin. Invest.* **105**, R1 (2000), ISSN 0021-9738.
- [34] R. M. Ribeiro, H. Mohri, D. D. Ho, and A. S. Perelson, *PNAS* **99**, 15572 (2002), ISSN 0027-8424, 1091-6490.
- [35] C. B. Jabara, C. D. Jones, J. Roach, J. A. Anderson, and R. Swanstrom, *PNAS* **108**, 20166 (2011), ISSN 0027-8424, 1091-6490.
- [36] P. J. A. Cock, T. Antao, J. T. Chang, B. A. Chapman, C. J. Cox, A. Dalke, I. Friedberg, T. Hamelryck, F. Kauff, B. Wilczynski, et al., *Bioinformatics* **25**, 1422 (2009).
- [37] S. van der Walt, S. Colbert, and G. Varoquaux, *Computing in Science Engineering* **13**, 22 (2011).
- [38] J. D. Hunter, *Computing In Science & Engineering* **9**, 90 (2007).



## Supplementary Figures and Tables

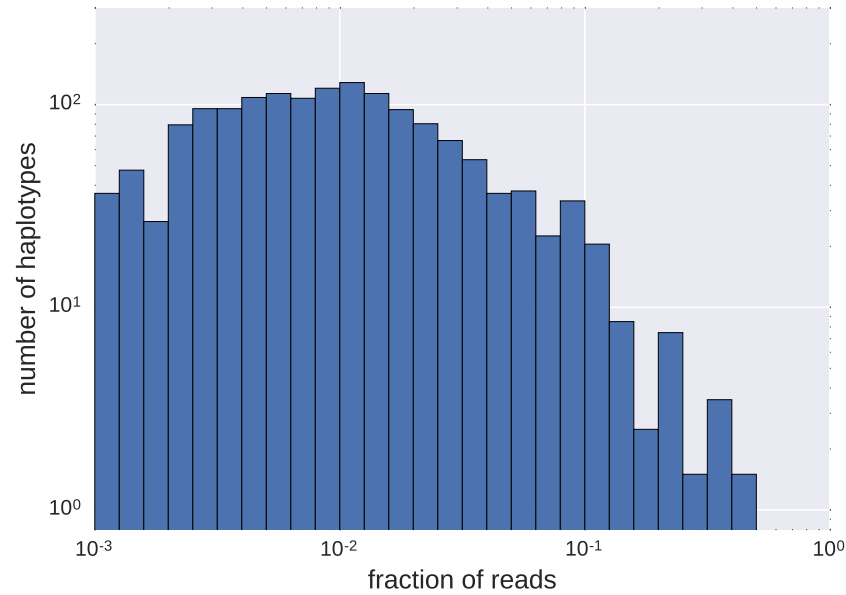


Figure S1: Distributions of frequencies of haplotypes. We observe a wide variation of haplotype frequencies – measured by fraction of reads – from below 0.001 to above 0.5. The majority of haplotypes are seen at frequencies around 0.01. Note that both scales are logarithmic.

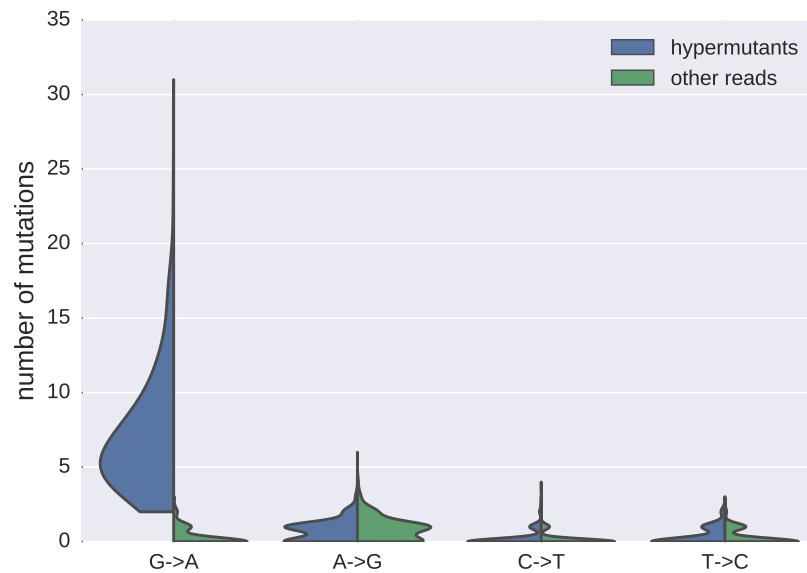


Figure S2: Distributions of mutations classified as hypermutants or regular reads.

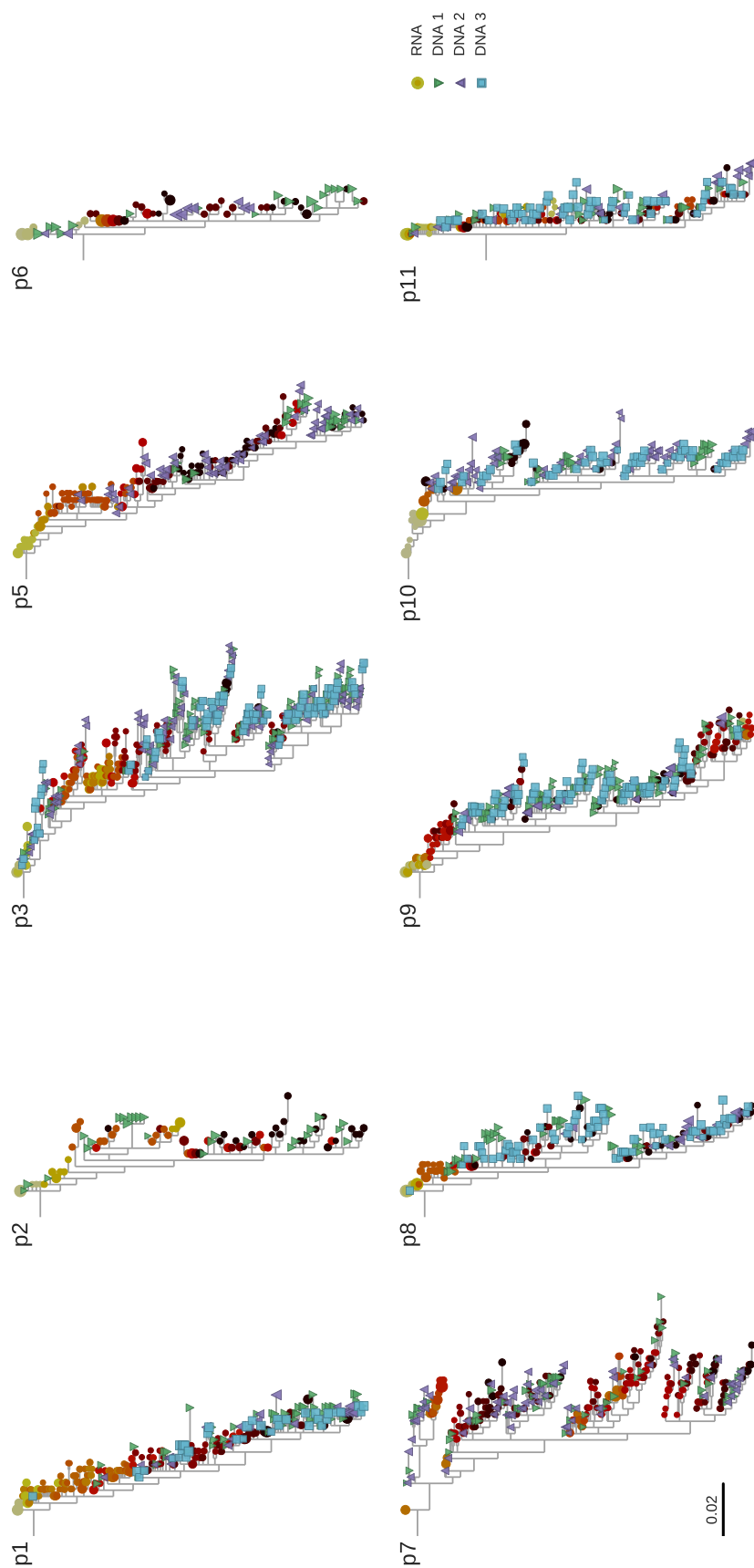


Figure S3: Phylogenetic trees of RNA and DNA samples.

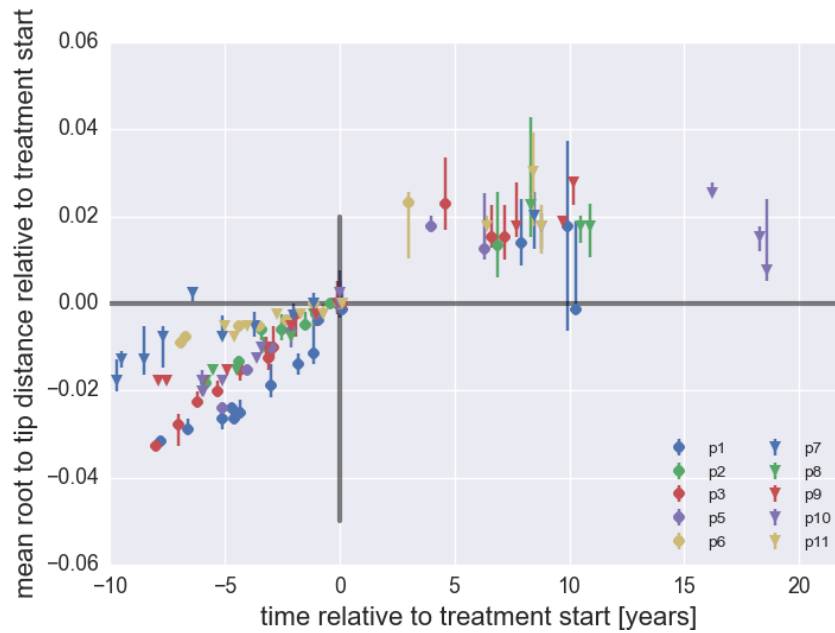


Figure S4: Average root-to-tip distances for plasma HIV-1 RNA sequences obtained before start of ART and PBMC HIV-1 DNA sequences obtained after start of ART. This figure is analogous to Fig. 2 in the main text, but shows root-to-tip distance of DNA sequences classified as hypermutants. While hypermutant sequences are between 2 and 4% more distance from the approximate founder sequence, the distances do not change over time.

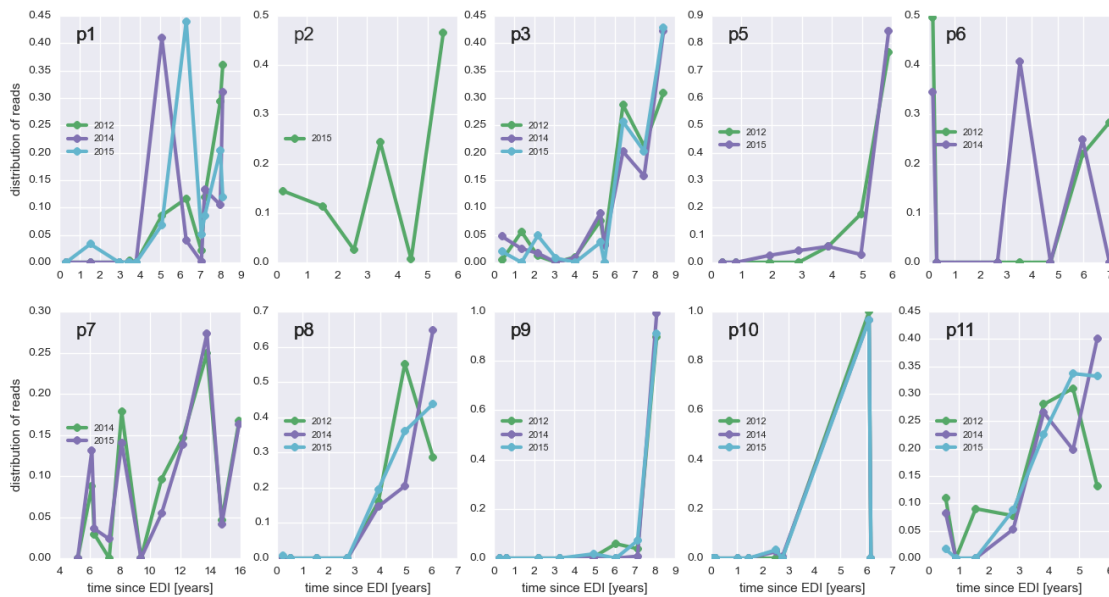


Figure S5: The distribution of plausible seeding times of reservoir sequences. For each read obtained from DNA samples, the RNA sample from which this read is most likely drawn is determined. The plots show the distribution of these most likely origin samples across all available RNA samples.

patient	sample	#reads	#templates	#seq> 0.002	#seq> 0.01	% recaptured	% hyper	#stop good	#stop hyper	% stop good	%stop hyper
p1	2012-10-02	19363	820	58	27	33	13	8	1.85	1.4	83
p1	2014-10-07	36454	148	23	13	31	14	9	0.81	1.3	36
p1	2015-02-12	44434	38	36	22	36	17	7	1.43	0.8	85
p2	2015-06-09	16586	75	26	12	NA	29	9	1.65	3.9	99
p3	2012-09-18	23347	243	116	26	50	19	8	1.23	4.6	75
p3	2014-10-13	28450	102	117	31	45	30	12	1.78	4.9	82
p3	2015-04-24	44850	108	98	34	38	24	10	1.21	5.0	61
p5	2012-10-26	11190	180	23	14	14	21	7	1.99	0.9	94
p5	2015-03-16	39762	72	56	29	7	6	8	1.36	3.8	99
p6	2012-10-24	191468	115	17	11	9	10	9	0.94	1.0	60
p6	2014-11-03	33887	15	9	3	33	0	9	nan	1.0	nan
p7	2014-12-01	35565	28	71	43	0	13	5	1.43	2.5	88
p7	2015-04-17	38323	108	67	41	0	11	5	1.34	0.5	94
p8	2012-09-21	6553	180	32	12	17	42	6	1.75	1.5	87
p8	2014-11-18	75473	279	16	11	0	30	6	2.00	0.7	99
p8	2015-04-07	18750	175	68	35	6	33	6	1.44	3.6	88
p9	2012-10-05	43306	60	64	43	19	13	5	0.52	1.1	44
p9	2014-10-02	11620	72	27	16	38	0	3	0.0	0.5	100
p9	2015-03-18	65003	72	71	38	16	12	6	1.86	0.7	98
p10	2012-10-09	7513	249	32	16	25	4	02	2.0	0.3	100
p10	2014-10-24	50537	116	74	29	28	14	16	1.22	7.1	99
p10	2015-02-27	59041	51	55	29	41	00	6	0.55	0.6	27
p11	2012-10-10	47701	124	36	23	30	16	20	1.00	6.0	86
p11	2014-10-22	32126	120	29	17	53	28	21	2.08	8.4	84
p11	2015-02-25	6834	123	57	37	35	33	11	0.99	3.7	75

Sequencing and hypermutation statistic of all samples. “good” refers to proviral sequences that are not obviously defective, while “hyper” refers to those with an excess of G→A mutations.

Patient/Plasmid	Primer name	Sequence
pAll	all <sup>a</sup> _fw <sup>b</sup> _1 <sup>c</sup> _689 <sup>d</sup> _705 <sup>e</sup>	AGG CAG GAC TCG GCT TGC
pAll	all_fw2(nex) <sup>f</sup> _770_787	( <sup>g</sup> TGC TCG GCA CGG TCA GAT GTG TAT AAG AGA CAG) <sup>h</sup> CGG AGG CTA GAA GGA GAG
p1-p5, p7, p8-p11	p1_p2_p3_p4_p5_p7_p8_p9_p10_p11_rev1_1339_1321	AAT CTT GTG GGG TGG CTC C
p6	p6_rev1_1339_1321	AAT CTG CTG GRG TGG CTC C
p1, p7, p8, p11	p1_p7_p8_p11_rev2(nex)_1101_1176	(GTC TCG TGG GCT CGG AGA TGT GTA TAA GAG ACA G)GT ATA GGG TAA TTT TGG CTG
p2	p2_rev2(nex)_1191-1176	(GTC TCG TGG GCT CGG AGA TGT GTA TAA GAG ACA G)GT ATA GGG TAA TTT TGG CTG
p3	p3_rev2(nex)_1191-1176	(GTC TCG TGG GCT CGG AGA TGT GTA TAA GAG ACA G)GT ATA GGG TAA TTT TGG CTG
p5	p5_rev2(nex)_1191-1176	(GTC TCG TGG GCT CGG AGA TGT GTA TAA GAG ACA G)GT ATA GGG TAA TTT TGG CTG
p6	p6_rev2(nex)_1191-1173	(GTC TCG TGG GCT CGG AGA TGT GTA TAA GAG ACA G)GT ATA GGA TAA TTT TGG CTG
p9, p10	p9_p10_rev2(nex)_1191-1176	(GTC TCG TGG GCT CGG AGA TGT GTA TAA GAG ACA G)GT ATA GGG TAA TTT TGR CTG
NL4-3	NL4-3_fw1_689_705	ACG CAG GAC TCG GCT TGC
NL4-3	NL4-3_rev1_1339_1321	AAT CTT GTG GGG TGG CTC C
NL4-3	NL4-3_fw2(nex)_770_787	(TGC TCG GCA GGG TCA GAT GTG TAT AAG AGA CAG) CGG AGG CTA GAA GGA GAG
NL4-3	NL4-3_rev2(nex)_1191_1173	(GTC TCG TGG GCT CGG AGA TGT GTA TAA GAG ACA G) TAT AGG GTA ATT TTG GCT G

**Primer table.** <sup>a</sup> Patient number. <sup>b</sup> Primer direction. <sup>c</sup> Primer position. <sup>d</sup> Inner primer. <sup>e</sup> Start position of the primer according to HXB2 coordinates. <sup>f</sup> Stop position of the primer according to HXB2 coordinates. <sup>g</sup> Nextera adapter. <sup>h</sup> Start position of the NexteraXT adapter. <sup>i</sup> Stop position of the NexteraXT adapter.