

1 **Title**

- 2 • **Full Title:** Early Emergence and Long-Term Persistence of HIV-Infected T Cell Clones  
3 in Children  
4 • **Short Title:** Early and Persistent HIV-1 Cell Clones in Children

5 **Authors**

6 Michael J. Bale<sup>1†</sup>, Mary Grace Katusiime<sup>1†\*</sup>, Daria Wells<sup>2</sup>, Xiaolin Wu<sup>2</sup>, Jonathan Spindler<sup>1</sup>, Elias  
7 K. Halvas<sup>3</sup>, Joshua C. Cyktor<sup>3</sup>, Ann Wiegand<sup>1</sup>, Wei Shao<sup>3</sup>, Mark F. Cotton<sup>4</sup>, Stephen H.  
8 Hughes<sup>1</sup>, John W. Mellors<sup>3</sup>, John M. Coffin<sup>5</sup>, Gert U. Van Zyl<sup>6</sup>, Mary F. Kearney<sup>1</sup>

9 **Affiliations**

10 <sup>1</sup>HIV Dynamics and Replication Program, CCR, National Cancer Institute, Frederick, MD, US.  
11

12 <sup>2</sup>Leidos Biomedical Research, Inc., Frederick National Laboratory for Cancer Research,  
13 Frederick, MD, US.  
14

15 <sup>3</sup>Department of Medicine, University of Pittsburgh, Pittsburgh, PA, US.  
16

17 <sup>4</sup>Department of Pediatrics and Child Health, Tygerberg Children's Hospital and Family Center  
18 for Research with Ubuntu, Stellenbosch University, Cape Town, ZA.  
19

20 <sup>5</sup>Department of Molecular Biology and Microbiology, Tufts University, Boston, MA, US.  
21

22 <sup>6</sup>Division of Medical Virology, Stellenbosch University and National Health Laboratory Service  
23 Tygerberg, Cape Town, ZA.  
24  
25  
26  
27  
28

29 †These authors contributed equally to this study

30 \*Corresponding Author:

31 Mary Grace Katusiime, PhD  
32 HIV Dynamics and Replication Program  
33 National Cancer Institute at Frederick  
34 1050 Boyles Street, Building 535, Room 134  
35 Frederick, Maryland, US 21702  
36 1-301-846-6522  
37 [marygrace.katusiime@nih.gov](mailto:marygrace.katusiime@nih.gov)  
38

39 **Abstract**

40 Little is known about the emergence and persistence of HIV-infected T cell clones in  
41 perinatally-infected children. We analyzed peripheral blood mononuclear cells for clonal  
42 expansion in 11 children who initiated antiretroviral therapy (ART) between 1.8-17.4 months of  
43 age and with viremia suppressed for 6-9 years. We obtained 8,662 HIV-1 integration sites from  
44 pre-ART and 1,861 sites on ART. Expanded clones of infected cells were detected pre-ART in  
45 10/11 children. In 8 children, infected cell clones detected pre-ART persisted for 6-9 years on  
46 ART. A comparison of integration sites in the samples obtained on ART with healthy donor PBMC  
47 infected *ex-vivo* showed selection for cells with proviruses integrated in *BACH2* and *STAT5B*. Our  
48 analyses indicate that, despite marked differences in T cell composition and dynamics between  
49 children and adults, HIV-infected cell clones are established early in children, persist for up to 9  
50 years on ART, and can be driven by proviral integration in proto-oncogenes.

51

52

53

## 54 Introduction

55 Human Immunodeficiency Virus (HIV) remains a worldwide health crisis. Approximately  
56 37.9 million people are living with HIV globally and about 1 million die each year (1). Although  
57 current antiretroviral therapy (ART) is able to fully suppress HIV-1 replication in the blood (2-5),  
58 lymph nodes (6-9), and other tissues (10, 11), it does not cure the infection. If treatment is initiated  
59 before the immune system is heavily compromised and if there is lifelong adherence, ART can  
60 lead to a partial restoration of CD4+ T cell numbers (12, 13) and can prevent immunodeficiency  
61 in most individuals.

62 The main obstacle to a cure for HIV-1 is the persistence of replication-competent  
63 proviruses in long-lived and/or proliferating populations of infected T cells (14, 15). Most of the  
64 infected cells that persist on ART contain defective proviruses that are incapable of producing  
65 infectious virus (16, 17), although they may be complemented to generate infectious virus upon  
66 ART interruption (18, 19). These defective proviruses do not directly contribute to the HIV-1  
67 reservoir that persists on ART but complicate its measurement and may contribute to persistent  
68 immune activation. The fraction of infected cells that contains replication-competent  
69 (intact/infectious) proviruses has been estimated to be between 1 and 5% in individuals on long-  
70 term ART (16, 17, 20). Although the fraction of intact proviruses is small relative to the total  
71 number of infected cells, there are sufficient replication-competent proviruses, or defective  
72 proviruses that can be readily complemented, to fuel rapid viral rebound if ART is interrupted (21,  
73 22). In both adults and children, when ART is initiated soon after infection, the number of infected  
74 cells is reduced, sometimes to levels below the detection limit of current assays (20, 23, 24) and  
75 rebound viremia can be significantly delayed (25-28).

76 Studies of HIV-1 integration sites were initially performed in cell lines and showed that  
77 sites were widely distributed but favored highly expressed genes (29-31). Two studies in 2014  
78 were the first to demonstrate expansion of HIV-infected T cells *in vivo* (32, 33). These clones of  
79 infected T cells can be detected as early as Fiebig IV in acute infection (34), can persist in adults  
80 for at least 3 years on ART (35), and are distributed among different tissues (6). Studies of clones  
81 persisting in adults on ART revealed selection against proviruses in expressed genes with a  
82 stronger selection against those that are integrated in the same orientation as the host gene (32, 36)  
83 and selection for proviruses integrated into some proto-oncogenes—e.g. *BACH2*, *MKL2*, and  
84 *STAT5B* (32, 33). Although much is now known about the HIV-1 integration site landscape in  
85 adults prior to and on ART, there is little information on clonal expansion of infected cells in  
86 children who acquired HIV perinatally (PHIV). The largely anti-inflammatory and  
87 immunoregulatory environment of the immune systems in children (37, 38) could affect the  
88 behavior of infected T cells in ways that would alter the integration site landscape and selection of  
89 proviruses at specific sites in children, leading to differences compared to adults. Furthermore,  
90 infants have a high fraction of naïve T-cells and fewer clonally expanded T-cells than adults (39,  
91 40), a difference that could affect integration site selection and clonal expansion of infected cells.

92 To our knowledge, only two reports have investigated clonal expansion of infected cells in  
93 children, to date (33, 41). However, only a few children were studied and the integration site  
94 sampling in these studies was shallow because it is difficult to collect large numbers of PBMC  
95 from infants and children. One study reported clones of infected cells in 3 children initiating ART  
96 during chronic infection (33) and the other in 3 neonates on ART who were followed for 2 years  
97 (41). Here, we expand on these studies to perform a deep look at the integration site landscape in  
98 11 children initiating ART early and followed for 6-9 years of continual suppression of viremia.  
99 We performed a detailed analysis of the integration site landscape by comparing the findings to  
100 those in *ex vivo* infected adult PBMC and to those in infected adults on ART (6, 32). To study the  
101 emergence of infected CD4+ T cell clones before ART initiation, the dynamics of their long-term  
102 persistence, and their potential survival within select genes, we obtained 10,523 integration sites  
103 from CD4+ T cells in the perinatally-infected children using samples obtained prior to and during  
104 long-term ART. We compared longitudinal integration site datasets to look for evidence of long-  
105 term persistence of clones of infected T cells and to investigate the frequency and size of the  
106 infected cell clones in the children. Finally, to determine if there exists selective maintenance of  
107 infected cells within single genes, we analyzed the integration sites in children compared to sites  
108 obtained from *ex vivo*-infected, CD8-depleted PBMC (deposited at [rid.ncifcrf.gov](http://rid.ncifcrf.gov)) (42, 43).

109 We report here that clones of infected cells are found in children as early as 1.8 months  
110 after birth and that some of the clones that arose early persisted for up to 9 years on ART.  
111 Strikingly, although there are noted differences between the immune environments in children  
112 compared to adults, our findings on the population of infected T cell clones are similar to what has  
113 been reported for adults, suggesting that clonal expansion is the main mechanism for persistence  
114 of HIV-1 in children whose viremia is suppressed by ART. We also found that the selection for  
115 proviruses integrated in certain genes is similar in adults and children and, importantly, that this  
116 selection occurs pre-ART. Integration events and selection for proviruses in these genes in children  
117 born with HIV-1 could have long-term effects in adulthood that have not been investigated and  
118 are not observed in adults who were not born with HIV infection.

119

## 120 **Results**

### 121 *Participants and Sampling*

122 PBMC were obtained from children enrolled in the Children with HIV and Early  
123 Antiretroviral therapy (CHER) randomized trial and post-CHER cohort (44) who were identified  
124 as plasma HIV-1 RNA positive by 7 weeks of age, initiated ART within 18 months of age (median:  
125 5.1 months; range: [1.8 to 17.4] months), and had long-term, sustained suppression on ART (3)  
126 (Table S1). Children were included based on the availability of pre-ART PBMC and PBMC  
127 obtained after at least 6 years of continuous suppression of viremia (median: 8.1 years; range 6.8  
128 to 9.1 years). The sex, pre-treatment plasma HIV-1 RNA, ART regimen, time to viral load  
129 suppression, and CD4 percentage after long-term ART are shown in Table S1. The pre-ART and  
130 on-ART enriched-CD4+ T cells were analyzed for the presence and persistence of clones of  
131 infected cells. We obtained between 197 and 1386 (median: 655) integration sites from each of the  
132 samples taken before ART was initiated and between 77 and 432 (median: 137) integration sites  
133 from those after at least 6 years on ART (Table 1). In total, we obtained 10,523 HIV-1 integration  
134 sites from the 11 children.

135

### 136 *Clones of HIV-1 infected cells are detected in children pre-ART and persist on long-term ART*

137 Clonal expansion of cells infected with replication-competent proviruses or defective  
138 proviruses that can be complemented during active replication (45, 46), is an important mechanism  
139 for HIV-1 persistence on ART (6, 14, 35, 47, 48). The detection of identical integration sites within  
140 a sample is the hallmark of clonal expansion of an infected cell, independent of the replication-  
141 competence of the integrated provirus. We defined an integration site as being from a clone using  
142 three separate criteria: 1) detection of the same integration site at least 3 times in pre-ART samples  
143 (to account for recently-infected cells that had duplicated their DNA but would die before  
144 establishing a clone), 2) detection of the same integration site at least twice in an on-ART sample  
145 (if a cell is dividing after long term ART, it is almost certainly part of a clone), and 3) detection of  
146 the same integration site in two different samples from the same donor. Additionally, the method  
147 we use to identify integration sites recovers the host-virus DNA junctions from both the 5' and 3'  
148 LTRs (32). Therefore, integration sites observed at both junctions were considered as a single  
149 integration site under the conservative assumption that they could have originated from the same  
150 provirus. In all but one of the 11 donors [Participant Identifier (PID) ZA009], we found at least  
151 one clone of infected cells in the pre-ART samples (range: [1 to 27]) (Table 1, column 4). We  
152 found at least 3 clones of infected cells in all on-ART samples (range: [3 to 32]) (Table 1, column  
153 9). Although we did not detect any clones of infected cells in the pre-ART samples from donor  
154 ZA009 by the stringent criteria described above, 16 of the integration sites were detected twice,  
155 suggesting that clones of infected T cells could have been present in this donor pre-ART (Table 1,  
156 column 4 parenthetical). We identified clones of infected cells in the pre-ART samples that  
157 persisted for up to 6-9 years on ART in 8 of the 11 children (range: [1, 7] clones) (Table 1, column

158 11). These data show that clonal expansion contributes to the persistence of total HIV-1 DNA in  
159 children, as was shown previously for adults (6, 15, 32, 35, 47).

160

161 *Size of infected cell clones is similar in children and adults*

162 We analyzed the size and frequency of infected cell clones using a modified Gini  
163 coefficient called the “oligoclonality index” (OCI) (49). Briefly, the OCI, which has a value  
164 between 0 and 1, is a measure of the non-uniformity of a given dataset; 0 indicates complete  
165 heterogeneity and 1 indicates complete homogeneity. In our analysis, 0 would mean each detected  
166 integration site was detected only once while a value of 1 would mean that all the integration sites  
167 would be from a single large clone. In the pre-ART samples, most integration sites were detected  
168 only once (Table 1, column 5). The pre-ART samples contained large numbers of recently infected  
169 cells that had not undergone clonal expansion. Thus, all pre-ART OCI values were less than 0.1  
170 (range: 0.006 to 0.085; median: 0.027). The pre-ART OCI positively correlated with the age at  
171 which ART was initiated – presumably because clones increase in size with time, which makes it  
172 easier for us to detect them (Adj.  $R^2=0.53$ ;  $p=0.011$ ) (Figure 1A). Stated differently, although  
173 clones can arise soon after infection (35, 50), they may require time to expand to a size that can be  
174 detected using the integration sites assay (35). As expected, the OCIs were significantly higher  
175 during long-term suppression on ART (range: 0.055 to 0.403; median: 0.161;  $p=0.002$ ) (Table 1,  
176 column 10, Figure 1B), suggesting that the short survival of most recently-infected T cells makes  
177 it easier to detect clones of infected cells after long-term ART (32, 35). It should be noted;  
178 however, that the on-ART OCI does not correlate with time on ART (Adj.  $R^2=-0.08$ ;  $p=0.63$ ),  
179 suggesting that clonal expansion during ART is not just a function of time, but rather a complex  
180 dependence on homeostatic, antigen-driven, and integration-driven proliferation (Figure 1B). We  
181 further compared the on-ART OCI in children to published datasets from 9 infected adults (6, 32)  
182 on long-term ART and found no statistical difference ( $p>0.99$ ; Figure 2; numerical data found in  
183 Jupyter Notebook – see methods) (51).

184

185 *Selection for cells with proviruses integrated in certain genes*

186 Recent reports show that HIV-1 proviruses integrated in one of a small number of genes  
187 (15, 32, 33, 35, 52) contribute to the growth, survival, and persistence of the infected cell clones  
188 *in vivo*. To look for evidence of similar selection in children born with HIV-1 and treated early  
189 with ART, we compared the distribution of integration sites from the children (pre-ART and on-  
190 ART) to integration sites obtained from *ex vivo* HIV-1 infected, CD8-depleted PBMC from healthy  
191 donors [deposited at [rid.ncifcrf.gov](http://rid.ncifcrf.gov); (42)]. We asked if there was evidence for enrichment of  
192 proviruses in specific genes *in vivo* (relative to *ex vivo*). We also analyzed the orientation of the  
193 proviruses relative to the host gene. Enrichment in the fraction of proviruses within, and oriented  
194 in the same direction as, the gene are evidence of post-integration selection. Enrichment of the

195 integration sites was determined by comparing the *ex vivo*-infected PBMC dataset against the *in*  
196 *vivo* datasets. For this analysis, clonally amplified sites were removed from the *in vivo* datasets by  
197 collapsing identical integration sites. Integration sites in intergenic regions (mapped to hg19) were  
198 not included in the analysis. The resulting datasets consisted of 335,614 integration sites from the  
199 *ex vivo* infected PBMC (87.2% of the initial data), 7039 sites from the pre-ART dataset from the  
200 children (83.9%), and 1202 (76.8%) sites from the on-ART dataset from the children. To detect  
201 enrichment in both the pre-ART and on-ART datasets relative to the *ex vivo* PBMC dataset,  
202 Fisher's Exact Tests were performed on genes in each library with post-hoc multiple tests  
203 correction. Adjusted p-values are reported with  $p_{\text{adj}} \leq 0.05$  being considered significant.

204 Consistent with what has been observed in virally suppressed adults (32), we found a strong  
205 enrichment for proviruses during ART integrated into both *BACH2* ( $p_{\text{adj}}=2.7*10^{-15}$ ) and *STAT5B*  
206 ( $p_{\text{adj}}=4.0*10^{-29}$ ) (Table 2), but not *MKL2* ( $p_{\text{adj}}>0.05$ ) during ART. The question of enrichment in  
207 samples prior to ART initiation in either adults or children has not previously been addressed.  
208 Strikingly, we observed a signal for enrichment of integrations into *BACH2* in children even prior  
209 to ART initiation ( $p_{\text{adj}}=8.9*10^{-17}$ ) showing that selection can occur early in PHIV infection.  
210 Although not statistically significant, we also observed a trend toward selection for integration  
211 events in *STAT5B* ( $p_{\text{adj}}=0.14$ ) (Table 2) prior to ART initiation.

212 Previous studies in adults have shown that, if there is post-integration selection for an HIV  
213 provirus in a gene, like *STAT5B* and *BACH2*, the proviruses are highly enriched for the same  
214 orientation as the gene (32, 33). We analyzed the genes for which there were at least 15 unique  
215 integrations in the *ex vivo* dataset and at least 1 integration in the *in vivo* dataset so that there would  
216 be a signal sufficient to detect selection. Although 18 genes were retained for analysis in the pre-  
217 ART dataset, only 2 met these criteria in the on-ART dataset (Table S2, S3). Despite the global  
218 preference for proviruses detected on ART to be integrated against the gene (*ex vivo* PBMC: 50.0%  
219 vs. children on-ART: 54.7%;  $p=0.0011$ ), there was no evidence for such global selection prior to  
220 initiation of ART (*ex vivo* PBMC: 50.0%; children pre-ART: 50.7%;  $p=0.26$  for the difference)  
221 (Figure 3A). However, of the 18 genes in which there were sufficient numbers of integrations in  
222 pre-ART samples, we found selection for with-the-gene integration in both *BACH2* ( $p_{\text{adj}}=2.0*10^{-3}$ )  
223 and *STAT5B* ( $p_{\text{adj}}=7.8*10^{-3}$ ) (Table S2, Figure 3B) and an against-the-gene bias in an ankyrin  
224 repeat protein, *ANKRD11* ( $p_{\text{adj}}=0.028$ ) (Table S3). Although these data provide evidence for strong  
225 selection for both *BACH2* and *STAT5B* pre-ART, we do not consider the against-gene bias for  
226 *ANKRD11* to be evidence of selection specific to that gene because of the global bias for against-  
227 gene integrations and the lack of an enrichment signal in this and previous datasets.

228 Likewise, integration sites recovered from children on ART in *BACH2* and *STAT5B* were  
229 significantly selected for with-the-gene orientation ( $p_{\text{BACH2}}=0.034$ ;  $p_{\text{STAT5B}}=6.2*10^{-5}$ ) (Table S3,  
230 Figure 3B). Taken together with the enrichment analyses, we conclude that cells containing  
231 proviruses integrated in *BACH2* and *STAT5B* in the same orientation as the genes were selected in  
232 children both prior to and on-ART.

233 We also compared the within-gene distribution of proviruses in *BACH2* and *STAT5B* in the  
234 children vs. the *ex vivo*-infected PBMC using an in-house mapping application (36) (Figure 4). In  
235 both genes, clearly visible clusters of integration sites in the same orientation as the gene (shown  
236 in blue) in a single intron upstream of the start of translation were observed in the children both  
237 before and during ART (Figure 4B, C, E, F). The *ex vivo*-infected PBMC have a broader, randomly  
238 oriented (equal red to blue) distribution (Figure 4A, D) in comparison. The different distributions  
239 highlight the selection for directional and clustered integration events into *BACH2* and *STAT5B* in  
240 children both prior to and on-ART.

241

242 *Sub-genomic sequencing datasets do not accurately characterize clonality within individuals*

243 Proviruses in a subset of the children in this study were previously characterized using  
244 single-genome sequencing (SGS) of the *gag-pol* genes (encoding P6, protease, and the first 900  
245 nucleotides of reverse transcriptase) (3). We assessed the clonality of the infected cells using  
246 integration site analysis compared to the identical sequences found in the SGS analysis. We found  
247 that proviruses with identical sub-genomic sequences were more common and constituted larger  
248 fractions of the data than the clones detected by sequencing integration site analyses (Figure 5,  
249 Figure S1). We also calculated the OCI for each set of data and found that the OCIs were  
250 significantly higher (average fold difference: 3.4x) for the sub-genomic single-genome sequences  
251 than for the integration site datasets ( $p = 0.0078$ ) (Figure 5). These data suggest that either  
252 proviruses with identical sub-genomic sequences have different sites of integration, as has been  
253 shown for adults (47), or that many of the integration sites that were detected contained proviruses  
254 for which the *gag-pol* regions could not be amplified and sequenced due to deletions, PCR primer  
255 mismatches, or both.



## 256 Discussion

257 Despite effective therapies, which have reduced the rate of mother-to-child HIV-1  
258 transmission (53-55), approximately 180,000 infants were infected worldwide in 2018 (53). These  
259 children must be included in the larger quest for effective HIV-1 curative interventions and such  
260 interventions may need to be tailored to their developing immune systems. Although the  
261 contribution of clonal expansion to HIV-1 persistence is well-studied in adults (6, 32, 35), this  
262 mechanism has not been well-described in children. Additionally, no analysis has been done in  
263 children on the clonal expansion of infected cells prior to the initiation of ART. To compare the  
264 mechanisms that underlie the persistence of HIV-1 infected cells during ART in adults and  
265 vertically-infected children, we performed HIV-1 integration site analysis on samples obtained  
266 from perinatally infected infants (prior to ART initiation) and from the same children during long-  
267 term suppression of viremia on ART (6-9 years of full suppression on ART). Despite inherent  
268 differences in T cell composition between children and adults (40) the clones of HIV-1 infected  
269 cells obtained from the blood of PHIV children in our study were not statistically different from  
270 adults (6, 32).

271 A study by Coffin et al. showed that infected cell clones can arise in adults in the first few  
272 weeks post-infection (35). In this study, we found that infected cell clones were detectable, using  
273 the integration sites assay, in 4 of the 5 samples collected from infants <3 months of age, consistent  
274 with early detection of clones in adults (35). In 2 of the 5 infants first sampled at <3 months old,  
275 we detected multiple proviruses with identical integration sites in both the pre-ART sample and  
276 the 6-9 years on-ART sample, demonstrating that clones of cells arose prior to ART initiation and  
277 persisted for years on ART. The other 7 donors, who were >3 months of age when initiating ART,  
278 also had detectable infected cell clones that persisted for at least 6 years of treatment. The  
279 frequency of clonal detection in the pre-ART populations tracked linearly with the estimated  
280 duration of infection prior to ART – using age as a surrogate – suggesting that the number of  
281 infected cell clones that expanded to detectable levels increased with the time of untreated  
282 infection, at least during the relatively short periods our donors were infected pre-ART. Our  
283 finding that infected cell clones had expanded and become large enough to be detected before two  
284 months of age supports the idea that the HIV-1 reservoir is generated rapidly, in actively dividing  
285 cells, in both adults and children (35, 56).

286 These results, in conjunction with previous studies showing that ongoing HIV-1 replication  
287 does not occur in children when viremia is fully suppressed on ART (3, 57, 58) and the fact that  
288 intact proviruses persist for years both in adults treated early (16) and in children treated early (20),  
289 supports the conclusion that the HIV-1 reservoir is maintained in vertically-infected children  
290 through the proliferation of cells infected prior to ART initiation, as it is in adults (6, 32, 35, 59).  
291 However, the available data are limited by the rarity of infected cells and the very small subset of  
292 HIV-infected cells that harbor intact, replication-competent proviruses in children (20). Although  
293 further studies are required to increase our understanding of the clonal expansion of intact  
294 proviruses as a mechanism by which the reservoir persists in both children and adults, it is possible

295 that defective proviruses can undergo complementation upon ART interruption and contribute to  
296 viral rebound (60).

297         Although the number of infected cells in children on ART is small, we were able to detect  
298 an enrichment in the number and the orientation of proviruses in both *BACH2* and *STAT5B* in the  
299 pre-ART and on-ART samples, suggesting that proviruses in a specific intron and oriented with  
300 these genes can promote the survival of these clones *in vivo*, as in adults (32). While the selection  
301 for the survival of cells harboring *BACH2* and *STAT5B* proviruses has been previously described  
302 in adults on-ART (32, 33), no data had been presented to show that such selection exists prior to  
303 ART initiation. In both pre- and on-ART, we saw clear evidence for selection of cells containing  
304 proviruses in the exon immediately upstream of the start site of translation in *BACH2* and *STAT5B*.  
305 Although the selection of *BACH2* integrants pre-ART was largely driven by a single child  
306 (ZA002) who did not initiate ART until 17 months of age, this single example nonetheless shows  
307 that clonal selection due to integration in specific genes is not strictly an on-ART phenomenon.  
308 The duration of untreated infection in this child may have allowed enough time for the selection  
309 of the cells with the *BACH2* proviruses to become detectable. Similar conclusions can be drawn  
310 for selection for proviruses integrated in the first intron of *STAT5B*, where there was clear  
311 evidence of selection for cells containing proviruses in the first intron, despite it's being a very  
312 strong target for integration *ex vivo*. The trend towards enrichment of *STAT5B* integrants in pre-  
313 ART samples was due to the high level of sampling required to overcome the background of  
314 integration events in this gene compared to the *ex vivo*-PBMC infected dataset; however, the  
315 statistically significant orientation bias prior to ART demonstrates that pre-ART selection exists  
316 for *STAT5B*.

317         Samples from a subset of the children studied here were previously characterized in  
318 experiments that showed that ART is effective in suppressing on-going cycles of viral replication  
319 in children (3). Thus, proviral SGS data were available at the same on-ART timepoint. The OCIs  
320 obtained using the P6-PR-RT SGS results were significantly higher than the OCIs obtained from  
321 the on-ART integration site data. The observation that a higher OCI was obtained from the SGS  
322 data than the ISA data adds to the growing number of studies (15, 47, 50) suggesting that viruses  
323 with identical sub-genomic sequences may not all come from a clonal population of infected cells.  
324 These data strongly suggest that sub-genomic sequencing does not always accurately identify  
325 clones of infected cells or sufficiently characterize the genetic diversity of the intra-patient HIV-  
326 1 populations that persist on ART (47). Although the results here are consistent with previous  
327 studies showing that sub-genomic sequences are not sufficient to define clonality, it should be  
328 noted that calculating an OCI for small-N datasets can result in artificially high OCI values. Studies  
329 that are based on integration site analysis, rather than SGS, are more appropriate to study the clonal  
330 expansion of infected cells.

331         It is important to note that because these children were diagnosed within a few weeks of  
332 birth it is not known whether the transmission of HIV-1 occurred at birth or in utero. Because of  
333 this ambiguity, the age of the participant may not accurately reflect the duration of infection,

334 although we found evidence of clonal expansion as early as 1.8 months after birth. Furthermore,  
335 the integration site libraries only represent a small fraction of the total number of infected cells in  
336 the blood. It is therefore likely that many of the integration sites that were recovered only once  
337 belong to clones of infected cells.

338         Despite these caveats, we have presented here the largest dataset yet of integration sites  
339 from pediatric HIV-1 infections both prior to ART and after durable suppression on ART. Because  
340 children primarily have naïve T cells, which do not have the HIV coreceptor CCR5 as a surface  
341 marker (40), as well as an immune environment that promotes quiescence (37, 38), and a more  
342 diverse T cell receptor repertoire (39), it is important to determine if there are differences between  
343 the observed frequency of clones and patterns of integration and post-insertional selection in  
344 children and adults. However, despite the differences in the immune systems of adults and  
345 children, our data suggest that these differences do not influence the infection and clonal expansion  
346 of T cells to a degree that is detectable by our integration site analysis. It is possible that by 6 to 9  
347 years of age the immune system may be similar enough to that of an adult to account for the striking  
348 similarities in the on-ART libraries of these children and the published data from adults. Although  
349 these data suggest that the role of clonal expansion as the mechanism for HIV-1 persistence during  
350 ART is similar in children and adults, further studies are warranted to better understand how the  
351 developing immune system affects clonal expansion and what effects proposed curative  
352 interventions might have in both children and adults.

353

354

355

## 356 **Materials and Methods**

### 357 *Study Approval and Ethics Statement*

358 The CHER trial is registered with ClinicalTrials.gov (NCT00102960). Guardians of all  
359 donors provided written informed consent and the study was approved by the Stellenbosch  
360 University Internal Review Board.

361

### 362 *Total HIV-1 DNA quantification*

363 HIV-1 DNA levels were determined using the integrase cell-associated DNA (iCAD) assay  
364 as previously described (61) with the following primers for use with HIV-1 Subtype C:

365 Forward primer      HIV\_Int\_FP    CCCTACAATCCCCAAAGTCA    4653 → 4672

366 Reverse primer      HIV\_Int\_RP    CACAATCATCACCTGCCATC    5051 → 5070

367

### 368 *Integration Sites Assay*

369 ISA was performed and analyzed as previously described (32, 62) using patient-specific  
370 primers to the 5' and 3' LTRs. Importantly, our protocol includes a shearing step (63) that  
371 effectively tags each DNA molecule, allowing determination of the relative numbers of cells in the  
372 initial pool with identical sites of integration (i.e., clonality). The full set of integration sites  
373 obtained has been submitted to the Retroviral Integration Sites Database (<https://rid.ncifcrf.gov/>)  
374 (43) and the primer sequences are available in Supplemental Table 5.

375 A comparison integration site dataset was prepared from CD8-depleted PBMC isolated  
376 from two HIV negative human donors infected *in vitro* with replication-competent HIV-1, subtype  
377 B (BAL) (64). After 2 days the cells were harvested and DNA was prepared and integration sites  
378 analyzed as previously described (36). The global distribution of the integration sites from the two  
379 donors, was indistinguishable; therefore all comparisons were performed with combined data from  
380 the two donors.

381

### 382 *Oligoclonality Index*

383 The oligoclonality Index (OCI) was calculated using a python script available at  
384 [https://github.com/michaelbale/python\\_stuff/](https://github.com/michaelbale/python_stuff/). Full details of the calculation are described in the  
385 supplemental text of Gillet, *et al* (49). Briefly, the LTR-corrected counts of all unique integration  
386 sites are sorted into descending order and the cumulative abundance of the clones are summed as  
387 a fraction of the total number of unique integration sites and normalized to have a maximal value  
388 of 1. Mathematically, the OCI is calculated as below:

389  $s_i$  – LTR-corrected count of integration site  $i$

390  $S$  – Number of unique integration sites in library

391  $N = \sum_{i=1}^S s_i$  – Total number of integration sites in library

392  $p_i = \frac{s_i}{N}$  – Relative abundance of integration site  $i$

393  $X_i = \sum_{k=1}^i p_k$  – cumulative abundance of all integration sites of size  $\{s_i\}$  or greater

394  $OCI = 2 * (\sum_{k=1}^S \left(\frac{X_k}{S}\right) - 0.5)$  – Oligoclonality index of library

395

### 396 *Statistical Analysis*

397 Clonality was assessed by grouping sequenced integration sites with identical pseudo-  
398 3’LTR genomic coordinates and different shear points into count data in R. These count data were  
399 used to generate the OCI. Independent integration sites into genes were pooled and assessed for  
400 selection by Fisher’s Exact test by either the pre-ART or on-ART library vs. the ex vivo infected  
401 PBMC library as null set. P-values for this gene-enrichment analysis were corrected post-hoc by  
402 the Benjamini-Hochberg method. Orientation biases were assessed in a similar manner with post-  
403 hoc corrections only in the pre-ART comparison. All adjusted p-values are presented as  $p_{adj}$  where  
404 appropriate. All other statistical analyses are noted where appropriate and performed in R v3.5.2.  
405 A Jupyter notebook (51) with the R commands and visualizations for the unedited figures available  
406 at [github.com/michaelbale/cher\\_bale/](https://github.com/michaelbale/cher_bale/).

407

### 408 *Phylogenetic Analyses*

409 HIV-1 P6-PR-RT sequences were aligned to HIV Consensus C using MUSCLE and  
410 neighbor joining phylogenetic p-distance trees were built using MEGA 7  
411 (<https://www.megasoftware.net/>) (65) and outgroup rooted to Consensus C. Distance matrix  
412 generation for calculation of the sequence-based OCI was performed using Hamming distance.

413

414

415

416

417

418

## 419 **References**

- 420 1. Tavhi F, Carter A, Jahagirdar D, Biehl M, Douwes-Schultz D, Larson S, Arora M, Dwyer-  
421 Lindgren L, Steuben K, Abbastabar H, Abu-Raddad L, Abyu D, Adabi M, Adebayo O,  
422 Adekanmbi V, Adetokunboh O, Ahmadi A, Ahmadian E, Ahmadpour E, Ahmed M, Akal  
423 C, Alahdab F. 2019. Global, regional, and national incidence, prevalence, and mortality of  
424 HIV, 1980-2017, and forecasts to 2030, for 195 countries and territories: a systematic  
425 analysis for the Global Burden of Disease, Injuries, and Risk Factors Study 2017. *THE*  
426 *LANCET* 6:e831 - e859.
- 427 2. Kearney MF, Spindler J, Shao W, Yu S, Anderson EM, O'Shea A, Rehm C, Poethke C,  
428 Kovacs N, Mellors JW, Coffin JM, Maldarelli F. 2014. Lack of Detectable HIV-1  
429 Molecular Evolution during Suppressive Antiretroviral Therapy. *PLoS Pathog*  
430 10:e1004010.
- 431 3. Van Zyl GU, Katusiime MG, Wiegand A, McManus WR, Bale MJ, Halvas EK, Luke B,  
432 Boltz VF, Spindler J, Laughton B, Engelbrecht S, Coffin JM, Cotton MF, Shao W, Mellors  
433 JW, Kearney MF. 2017. No evidence of HIV replication in children on antiretroviral  
434 therapy. *J Clin Invest* doi:10.1172/JCI94582.
- 435 4. Josefsson L, von Stockenstrom S, Faria NR, Sinclair E, Bacchetti P, Killian M, Epling L,  
436 Tan A, Ho T, Lemey P, Shao W, Hunt PW, Somsouk M, Wylie W, Douek DC, Loeb L,  
437 Custer J, Hoh R, Poole L, Deeks SG, Hecht F, Palmer S. 2013. The HIV-1 reservoir in  
438 eight patients on long-term suppressive antiretroviral therapy is stable with few genetic  
439 changes over time. *Proc Natl Acad Sci U S A* 110:E4987-96.
- 440 5. Vancoillie L, Hebberecht L, Dauwe K, Demecheleer E, Dinakis S, Vaneechoutte D,  
441 Mortier V, Verhofstede C. 2017. Longitudinal sequencing of HIV-1 infected patients with  
442 low-level viremia for years while on ART shows no indications for genetic evolution of  
443 the virus. *Virology* 510:185-193.
- 444 6. McManus WR, Bale MJ, Spindler J, Wiegand A, Musick A, Patro SC, Sobolewski MD,  
445 Musick VK, Anderson EM, Cyktor JC, Halvas EK, Shao W, Wells D, Wu X, Keele BF,  
446 Milush JM, Hoh R, Mellors JW, Hughes SH, Deeks SG, Coffin JM, Kearney MF. 2019.  
447 HIV-1 in lymph nodes is maintained by cellular proliferation during antiretroviral therapy.  
448 *J Clin Invest* 130.
- 449 7. Bozzi G, Simonetti FR, Watters SA, Anderson EM, Gouzoulis M, Kearney MF, Rote P,  
450 Lange C, Shao W, Gorelick R, Fullmer B, Kumar S, Wank S, Hewitt S, Kleiner DE, Hattori  
451 J, Bale MJ, Hill S, Bell J, Rehm C, Grossman Z, Yarchoan R, Uldrick T, Maldarelli F.  
452 2019. No evidence of ongoing HIV replication or compartmentalization in tissues during  
453 combination antiretroviral therapy: Implications for HIV eradication. *Sci Adv* 5:eaav2045.
- 454 8. Josefsson L, Palmer S, Faria NR, Lemey P, Casazza J, Ambrozak D, Kearney M, Shao W,  
455 Kottlilil S, Sneller M, Mellors J, Coffin JM, Maldarelli F. 2013. Single cell analysis of  
456 lymph node tissue from HIV-1 infected patients reveals that the majority of CD4+ T-cells  
457 contain one HIV-1 DNA molecule. *PLoS Pathog* 9:e1003432.

- 458 9. Lee E, von Stockenstrom S, Morcilla V, Odevall L, Hiener B, Shao W, Hartogensis W,  
459 Bacchetti P, Milush J, Liegler T, Sinclair E, Hatano H, Hoh R, Somsouk M, Hunt P, Boritz  
460 E, Douek D, Fromentin R, Chomont N, Deeks SG, Hecht FM, Palmer S. 2019. The impact  
461 of antiretroviral therapy duration on the HIV-1 infection of T-cells within anatomic sites.  
462 *J Virol* doi:10.1128/JVI.01270-19.
- 463 10. Imamichi H, Degray G, Dewar RL, Mannon P, Yao M, Chairez C, Sereti I, Kovacs JA.  
464 2011. Lack of compartmentalization of HIV-1 quasispecies between the gut and peripheral  
465 blood compartments. *J Infect Dis* 204:309-14.
- 466 11. Evering TH, Mehandru S, Racz P, Tenner-Racz K, Poles MA, Figueroa A, Mohri H,  
467 Markowitz M. 2012. Absence of HIV-1 evolution in the gut-associated lymphoid tissue  
468 from patients on combination antiviral therapy initiated during primary infection. *PLoS*  
469 *Pathog* 8:e1002506.
- 470 12. Maartens G, Celum C, Lewin SR. 2014. HIV infection: epidemiology, pathogenesis,  
471 treatment, and prevention. *Lancet* 384:258-71.
- 472 13. Li TS, Tubiana R, Katlama C, Calvez V, Ait Mohand H, Autran B. 1998. Long-lasting  
473 recovery in CD4 T-cell function and viral-load reduction after highly active antiretroviral  
474 therapy in advanced HIV-1 disease. *Lancet* 351:1682-6.
- 475 14. Simonetti FR, Sobolewski MD, Fyne E, Shao W, Spindler J, Hattori J, Anderson EM,  
476 Watters SA, Hill S, Wu X, Wells D, Su L, Luke BT, Halvas EK, Besson G, Penrose KJ,  
477 Yang Z, Kwan RW, Van Waes C, Uldrick T, Citrin DE, Kovacs J, Polis MA, Rehm CA,  
478 Gorelick R, Piatak M, Keele BF, Kearney MF, Coffin JM, Hughes SH, Mellors JW,  
479 Maldarelli F. 2016. Clonally expanded CD4+ T cells can produce infectious HIV-1 in vivo.  
480 *Proc Natl Acad Sci U S A* 113:1883-8.
- 481 15. Einkauf KB, Lee GQ, Gao C, Sharaf R, Sun X, Hua S, Chen SM, Jiang C, Lian X,  
482 Chowdhury FZ, Rosenberg ES, Chun TW, Li JZ, Yu XG, Lichterfeld M. 2019. Intact HIV-  
483 1 proviruses accumulate at distinct chromosomal positions during prolonged antiretroviral  
484 therapy. *J Clin Invest* 129:988-998.
- 485 16. Bruner KM, Murray AJ, Pollack RA, Soliman MG, Laskey SB, Capoferri AA, Lai J, Strain  
486 MC, Lada SM, Hoh R, Ho YC, Richman DD, Deeks SG, Siliciano JD, Siliciano RF. 2016.  
487 Defective proviruses rapidly accumulate during acute HIV-1 infection. *Nat Med* 22:1043-  
488 9.
- 489 17. Ho YC, Shan L, Hosmane NN, Wang J, Laskey SB, Rosenbloom DI, Lai J, Blankson JN,  
490 Siliciano JD, Siliciano RF. 2013. Replication-competent noninduced proviruses in the  
491 latent reservoir increase barrier to HIV-1 cure. *Cell* 155:540-51.
- 492 18. Lori F, Hall L, Lusso P, Popovic M, Markham P, Franchini G, Reitz MS. 1992. Effect of  
493 reciprocal complementation of two defective human immunodeficiency virus type 1 (HIV-  
494 1) molecular clones on HIV-1 cell tropism and virulence. *Journal of Virology* 66:5553-  
495 5560.

- 496 19. Salzwedel K, Berger EA. 2000. Cooperative subunit interactions within the oligomeric  
497 envelope glycoprotein of HIV-1: Functional complementation of specific defects in gp120  
498 and gp41. *Proceedings of the National Academy of Sciences* 97:12794-12799.
- 499 20. Katusiime MG, Halvas EK, Wright I, Joseph K, Bale MJ, Kirby-McCullough B,  
500 Engelbrecht S, Shao W, Hu WS, Cotton MF, Mellors JW, Kearney MF, van Zyl GU. 2019.  
501 Intact HIV proviruses persist in children 7-9 years after initiation of ART in the first year  
502 of life. *J Virol* doi:10.1128/JVI.01519-19.
- 503 21. Chun TW, Justement JS, Murray D, Hallahan CW, Maenza J, Collier AC, Sheth PM, Kaul  
504 R, Ostrowski M, Moir S, Kovacs C, Fauci AS. 2010. Rebound of plasma viremia following  
505 cessation of antiretroviral therapy despite profoundly low levels of HIV reservoir:  
506 implications for eradication. *AIDS* 24:2803-8.
- 507 22. Hamlyn E, Ewings FM, Porter K, Cooper DA, Tambussi G, Schechter M, Pedersen C,  
508 Okulicz JF, McClure M, Babiker A, Weber J, Fidler S, Insight S, Investigators S. 2012.  
509 Plasma HIV viral rebound following protocol-indicated cessation of ART commenced in  
510 primary and chronic HIV infection. *PLoS One* 7:e43754.
- 511 23. de Souza MS, Pinyakorn S, Akapirat S, Pattanachaiwit S, Fletcher JL, Chomchey N, Kroon  
512 ED, Ubolyam S, Michael NL, Robb ML, Phanuphak P, Kim JH, Phanuphak N,  
513 Ananworanich J, Group RSS. 2016. Initiation of Antiretroviral Therapy During Acute  
514 HIV-1 Infection Leads to a High Rate of Nonreactive HIV Serology. *Clin Infect Dis*  
515 63:555-61.
- 516 24. Ananworanich J, Chomont N, Eller LA, Kroon E, Tovanabutra S, Bose M, Nau M, Fletcher  
517 JLK, Tipsuk S, Vandergeeten C, O'Connell RJ, Pinyakorn S, Michael N, Phanuphak N,  
518 Robb ML, Rv, groups RSs. 2016. HIV DNA Set Point is Rapidly Established in Acute HIV  
519 Infection and Dramatically Reduced by Early ART. *EBioMedicine* 11:68-72.
- 520 25. Henrich TJ, Hanhauser E, Marty FM, Sirignano MN, Keating S, Lee TH, Robles YP, Davis  
521 BT, Li JZ, Heisey A, Hill AL, Busch MP, Armand P, Soiffer RJ, Altfeld M, Kuritzkes DR.  
522 2014. Antiretroviral-free HIV-1 remission and viral rebound after allogeneic stem cell  
523 transplantation: report of 2 cases. *Ann Intern Med* 161:319-27.
- 524 26. Luzuriaga K, Gay H, Ziemniak C, Sanborn KB, Somasundaran M, Rainwater-Lovett K,  
525 Mellors JW, Rosenbloom D, Persaud D. 2015. Viremic relapse after HIV-1 remission in a  
526 perinatally infected child. *N Engl J Med* 372:786-8.
- 527 27. Persaud D, Gay H, Ziemniak C, Chen YH, Piatak M, Jr., Chun TW, Strain M, Richman D,  
528 Luzuriaga K. 2013. Absence of detectable HIV-1 viremia after treatment cessation in an  
529 infant. *N Engl J Med* 369:1828-35.
- 530 28. Henrich TJ, Hatano H, Bacon O, Hogan LE, Rutishauser R, Hill A, Kearney MF, Anderson  
531 EM, Buchbinder SP, Cohen SE, Abdel-Mohsen M, Pohlmeier CW, Fromentin R, Hoh R,  
532 Liu AY, McCune JM, Spindler J, Metcalf-Pate K, Hobbs KS, Thanh C, Gibson EA,  
533 Kuritzkes DR, Siliciano RF, Price RW, Richman DD, Chomont N, Siliciano JD, Mellors  
534 JW, Yukl SA, Blankson JN, Liegler T, Deeks SG. 2017. HIV-1 persistence following



- 535 extremely early initiation of antiretroviral therapy (ART) during acute HIV-1 infection: An  
536 observational study. *PLoS Med* 14:e1002417.
- 537 29. Craigie R, Bushman FD. 2012. HIV DNA integration. *Cold Spring Harb Perspect Med*  
538 2:a006890.
- 539 30. Schröder AR, Shinn P, Chen H, Berry C, Ecker JR, Bushman F. 2002. HIV-1 integration  
540 in the human genome favors active genes and local hotspots. *Cell* 110:521-9.
- 541 31. Singh PK, Plumb MR, Ferris AL, Iben JR, Wu X, Fadel HJ, Luke BT, Esnault C, Poeschla  
542 EM, Hughes SH, Kvaratskhelia M, Levin HL. 2015. LEDGF/p75 interacts with mRNA  
543 splicing factors and targets HIV-1 integration to highly spliced genes. *Genes Dev* 29:2287-  
544 97.
- 545 32. Maldarelli F, Wu X, Su L, Simonetti FR, Shao W, Hill S, Spindler J, Ferris AL, Mellors  
546 JW, Kearney MF, Coffin JM, Hughes SH. 2014. HIV latency. Specific HIV integration  
547 sites are linked to clonal expansion and persistence of infected cells. *Science* 345:179-83.
- 548 33. Wagner TA, McLaughlin S, Garg K, Cheung CY, Larsen BB, Styrchak S, Huang HC,  
549 Edlefsen PT, Mullins JI, Frenkel LM. 2014. HIV latency. Proliferation of cells with HIV  
550 integrated into cancer genes contributes to persistent infection. *Science* 345:570-3.
- 551 34. Fiebig EW, Wright DJ, Rawal BD, Garrett PE, Schumacher RT, Peddada L, Heldebrant C,  
552 Smith R, Conrad A, Kleinman SH, Busch MP. 2003. Dynamics of HIV viremia and  
553 antibody seroconversion in plasma donors: implications for diagnosis and staging of  
554 primary HIV infection. *AIDS* 17:1871-1879.
- 555 35. Coffin JM, Wells DW, Zerbato JM, Kuruc JD, Guo S, Luke BT, Eron JJ, Bale M, Spindler  
556 J, Simonetti FR, Hill S, Kearney MF, Maldarelli F, Wu X, Mellors JW, Hughes SH. 2019.  
557 Clones of infected cells arise early in HIV-infected individuals. *JCI Insight* 4.
- 558 36. Coffin JM, Bale MJ, Wells DW, Guo S, Luke B, Zerbato JM, Sobolewski M, Sia T, Shao  
559 W, Xiaolin W, Maldarelli F, Kearney MF, Mellors JW, Hughes SH. 2020. What Shapes  
560 the in Vivo Distribution of HIV Integration Sites? Submitted.
- 561 37. Surendran N, Simmons A, Pichichero ME. 2018. TLR agonist combinations that stimulate  
562 Th type I polarizing responses from human neonates. *Innate Immun* 24:240-251.
- 563 38. Kollmann TR, Crabtree J, Rein-Weston A, Blimkie D, Thommai F, Wang XY, Lavoie PM,  
564 Furlong J, Fortuno ES, 3rd, Hajjar AM, Hawkins NR, Self SG, Wilson CB. 2009. Neonatal  
565 innate TLR-mediated responses are distinct from those of adults. *J Immunol* 183:7150-60.
- 566 39. Wedderburn LR, Patel A, Varsani H, Woo P. 2001. The developing human immune  
567 system: T-cell receptor repertoire of children and young adults shows a wide discrepancy  
568 in the frequency of persistent oligoclonal T-cell expansions. *Immunology* 102:301-9.
- 569 40. Shearer WT, Rosenblatt HM, Gelman RS, Oyomopito R, Plaeger S, Stiehm ER, Wara DW,  
570 Douglas SD, Luzuriaga K, McFarland EJ, Yogev R, Rathore MH, Levy W, Graham BL,  
571 Spector SA, Pediatric ACTG. 2003. Lymphocyte subsets in healthy children from birth  
572 through 18 years of age: the Pediatric AIDS Clinical Trials Group P1009 study. *J Allergy*  
573 *Clin Immunol* 112:973-80.

- 574 41. Garcia-Broncano P, Maddali S, Einkauf K, Jiang C, Gao C, Chevalier J, Chowdhury F,  
575 Maswabi K, Ajibola G, Moyo S, Mohammed T, Ncube T, Makhema J, Jean-Philippe P,  
576 Yu X, Powis K, Lockman S, Kuritzkes D, Shapiro R, Lichterfeld M. 2019. Early  
577 antiretroviral therapy in neonates with HIV-1 infection restricts viral reservoir size and  
578 induces a distinct innate immune profile. *Science Translational Medicine* 11.
- 579 42. Ferris AL, Wells DW, Guo S, Del Prete GQ, Swanstrom AE, Coffin JM, Wu X, Lifson JD,  
580 Hughes SH. 2019. Clonal expansion of SIV-infected cells in macaques on antiretroviral  
581 therapy is similar to that of HIV-infected cells in humans. *PLoS Pathog* 15:e1007869.
- 582 43. Shao W, Shan J, Kearney MF, Wu X, Maldarelli F, Mellors JW, Luke B, Coffin JM,  
583 Hughes SH. 2016. Retrovirus Integration Database (RID): a public database for retroviral  
584 insertion sites into host genomes. *Retrovirology* 13:47.
- 585 44. Cotton MF, Violari A, Otway K, Panchia R, Dobbels E, Rabie H, Josipovic D, Liberty  
586 A, Lazarus E, Innes S, van Rensburg AJ, Pelser W, Truter H, Madhi SA, Handelsman E,  
587 Jean-Philippe P, McIntyre JA, Gibb DM, Babiker AG, Team CS. 2013. Early time-limited  
588 antiretroviral therapy versus deferred therapy in South African infants infected with HIV:  
589 results from the children with HIV early antiretroviral (CHER) randomised trial. *Lancet*  
590 382:1555-63.
- 591 45. Burke DS. 1997. Recombination in HIV: an important viral evolutionary strategy. *Emerg*  
592 *Infect Dis* 3:253-9.
- 593 46. Robertson DL, Sharp PM, McCutchan FE, Hahn BH. 1995. Recombination in HIV-1.  
594 *Nature* 374:124-6.
- 595 47. Patro SC, Brandt LD, Bale MJ, Halvas EK, Joseph KW, Shao W, Wu X, Guo S, Murrell  
596 B, Wiegand A, Spindler J, Raley C, Hautman C, Sobolewski M, Fennessey CM, Hu WS,  
597 Luke B, Hasson JM, Niyongabo A, Capoferri AA, Keele BF, Milush J, Hoh R, Deeks SG,  
598 Maldarelli F, Hughes SH, Coffin JM, Rausch JW, Mellors JW, Kearney MF. 2019.  
599 Combined HIV-1 sequence and integration site analysis informs viral dynamics and allows  
600 reconstruction of replicating viral ancestors. *Proc Natl Acad Sci U S A*  
601 doi:10.1073/pnas.1910334116.
- 602 48. Bui JK, Sobolewski MD, Keele BF, Spindler J, Musick A, Wiegand A, Luke BT, Shao W,  
603 Hughes SH, Coffin JM, Kearney MF, Mellors JW. 2017. Proviruses with identical  
604 sequences comprise a large fraction of the replication-competent HIV reservoir. *PLoS*  
605 *Pathog* 13:e1006283.
- 606 49. Gillet NA, Malani N, Melamed A, Gormley N, Carter R, Bentley D, Berry C, Bushman  
607 FD, Taylor GP, Bangham CR. 2011. The host genomic environment of the provirus  
608 determines the abundance of HTLV-1-infected T-cell clones. *Blood* 117:3113-22.
- 609 50. Garcia-Broncano P, Maddali S, Einkauf KB, Jiang C, Gao C, Chevalier J, Chowdhury FZ,  
610 Maswabi K, Ajibola G, Moyo S, Mohammed T, Ncube T, Makhema J, Jean-Philippe P,  
611 Yu XG, Powis KM, Lockman S, Kuritzkes DR, Shapiro R, Lichterfeld M. 2019. Early  
612 antiretroviral therapy in neonates with HIV-1 infection restricts viral reservoir size and  
613 induces a distinct innate immune profile. *Sci Transl Med* 11.

- 614 51. Kluyver T, Ragan-Kelley B, Perez F, Granger B, Bussonnier M, Frederic J, Kelley K,  
615 Hamrick J, Grout J, Corlay S, Ivanov P, Avila D, Abdalla S, Willing C. 2016. Jupyter  
616 Notebooks – a publishing format for reproducible computational workflows.
- 617 52. Ikeda T, Shibata J, Yoshimura K, Koito A, Matsushita S. 2007. Recurrent HIV-1  
618 integration at the BACH2 locus in resting CD4<sup>+</sup> T cell populations during effective highly  
619 active antiretroviral therapy. *J Infect Dis* 195:716-25.
- 620 53. Kempton J, Hill A, Levi JA, Heath K, Pozniak A. 2019. Most new HIV infections, vertical  
621 transmissions and AIDS-related deaths occur in lower-prevalence countries. *J Virus Erad*  
622 5:92-101.
- 623 54. Ghoma Linguissi LS, Sagna T, Soubeiga ST, Gwom LC, Nkenfou CN, Obiri-Yeboah D,  
624 Ouattara AK, Pietra V, Simpore J. 2019. Prevention of mother-to-child transmission  
625 (PMTCT) of HIV: a review of the achievements and challenges in Burkina-Faso. *HIV*  
626 *AIDS (Auckl)* 11:165-177.
- 627 55. Nduati EW, Hassan AS, Knight MG, Muema DM, Jahangir MN, Mwaringa SL, Etyang  
628 TJ, Rowland-Jones S, Urban BC, Berkley JA. 2015. Outcomes of prevention of mother to  
629 child transmission of the human immunodeficiency virus-1 in rural Kenya--a cohort study.  
630 *BMC Public Health* 15:1008.
- 631 56. Jones BR, Kinloch NN, Horacek J, Ganase B, Harris M, Harrigan PR, Jones RB,  
632 Brockman MA, Joy JB, Poon AFY, Brumme ZL. 2018. Phylogenetic approach to recover  
633 integration dates of latent HIV sequences within-host. *Proc Natl Acad Sci U S A*  
634 115:E8958-E8967.
- 635 57. Persaud D, Ray SC, Kajdas J, Ahonkhai A, Siberry GK, Ferguson K, Ziemniak C, Quinn  
636 TC, Casazza JP, Zeichner S, Gange SJ, Watson DC. 2007. Slow human immunodeficiency  
637 virus type 1 evolution in viral reservoirs in infants treated with effective antiretroviral  
638 therapy. *AIDS Res Hum Retroviruses* 23:381-90.
- 639 58. Ruff CT, Ray SC, Kwon P, Zinn R, Pendleton A, Hutton N, Ashworth R, Gange S, Quinn  
640 TC, Siliciano RF, Persaud D. 2002. Persistence of wild-type virus and lack of temporal  
641 structure in the latent reservoir for human immunodeficiency virus type 1 in pediatric  
642 patients with extensive antiretroviral exposure. *J Virol* 76:9481-92.
- 643 59. Reeves DB, Duke ER, Wagner TA, Palmer SE, Spivak AM, Schiffer JT. 2018. A majority  
644 of HIV persistence during antiretroviral therapy is due to infected cell proliferation. *Nat*  
645 *Commun* 9:4811.
- 646 60. Lu C-L, Pai JA, Nogueira L, Mendoza P, Gruell H, Oliveira TY, Barton J, Lorenzi JCC,  
647 Cohen YZ, Cohn LB, Klein F, Caskey M, Nussenzweig MC, Jankovic M. 2018.  
648 Relationship between intact HIV-1 proviruses in circulating CD4<sup>+</sup> T cells and  
649 rebound viruses emerging during treatment interruption. *Proceedings of the National*  
650 *Academy of Sciences* 115:E11341-E11348.
- 651 61. Hong F, Aga E, Cillo AR, Yates AL, Besson G, Fyne E, Koontz DL, Jennings C, Zheng  
652 L, Mellors JW. 2016. Novel Assays for Measurement of Total Cell-Associated HIV-1  
653 DNA and RNA. *J Clin Microbiol* 54:902-11.

- 654 62. Wells DW, Guo S, Shao W, Bale MJ, Coffin JM, Hughes SH, Wu X. 2020. An analytical  
655 pipeline for identifying and mapping the integration sites of HIV and other retroviruses.  
656 BMC Genomics 21:216.
- 657 63. Berry CC, Gillet NA, Melamed A, Gormley N, Bangham CR, Bushman FD. 2012.  
658 Estimating abundances of retroviral insertion sites from DNA fragment length data.  
659 Bioinformatics 28:755-62.
- 660 64. Gartner S, Markovits P, Markovitz DM, Kaplan MH, Gallo RC, Popovic M. 1986. The  
661 role of mononuclear phagocytes in HTLV-III/LAV infection. Science 233:215-9.
- 662 65. Kumar S, Stecher G, Tamura K. 2016. MEGA7: Molecular Evolutionary Genetics Analysis  
663 Version 7.0 for Bigger Datasets. Mol Biol Evol 33:1870-4.

664

## 665 **Acknowledgements**

666

667 **General:** We thank the participants and their caregivers for their willingness to take part in this  
668 study as well as the clinic staff for their assistance.

669 **Funding:** This study was supported by funding from the National Institutes of Allergy and  
670 Infectious Diseases (Comprehensive International Program for Research on AIDS (CIPRA)–South  
671 Africa (U19 AI153217), U.S.-South Africa Program for Collaborative Biomedical Research,  
672 National Cancer Institute: U01CA200441, the Office of AIDS Research (to M.F.K.), Departments  
673 of Health of the Western Cape and Gauteng, South Africa, ViiV Healthcare, and by National  
674 Institute of Mental Health under Award R01MH105134. JMC was a Research Professor of the  
675 American Cancer Society and was supported by the NCI through a Leidos subcontract, I3XS110,  
676 and research grant R35 CA200421. JWM was supported by the NCI through a Leidos subcontract,  
677 12XS547.

## 678 **Author Contributions**

679 MJB – Analyzed data, performed statistical analyses, wrote the paper

680 MGK – Processed samples, performed SGS, analyzed data, wrote the paper

681 DW – Performed ISA, analyzed data

682 XW – Performed ISA, analyzed data

683 JS – Processed samples, performed SGS

684 EKH – Processed samples, performed DNA quantitation

685 JCC – Performed DNA quantitation

686 AW – Performed SGS

687 WS – Analyzed Data

688 MFC – Conceived of idea, reviewed manuscript

689 SHH – Conceived of idea, analyzed data, wrote the paper

690 JWM – Conceived of idea, wrote the paper

691 JMC–Analyzed data, wrote the paper  
692 GUVZ – Conceived of idea, wrote the paper  
693 MFK – Conceived of idea, analyzed data, wrote the paper

694

695 **Competing interests:** JWM is consultant to Gilead Sciences, Accelevir Diagnostics, Merck, and  
696 Infectious Disease Connect, has received grants from Gilead Sciences to the University of  
697 Pittsburgh, and owns share options in Co-crystal Pharmaceuticals, Inc. and Abound Bio, Inc.  
698 unrelated to the current work. The remaining authors have declared that no conflict of interest  
699 exist.

700

701 **Data and materials availability:** Previously published sequences from Van Zyl et al (3) were  
702 accessed from GenBank with accession numbers (KY820119-KY820376) retaining only the  
703 sequences associated with the on-ART timepoint of PIDs ZA003 – ZA010. The integration site  
704 libraries for the donors in this study are available in the Retrovirus Integration Database (RID)  
705 (43) at [rid.ncifcrf.gov](http://rid.ncifcrf.gov). The ex vivo infected PBMC integration sites library is also available in the  
706 RID under the Pubmed ID 31291371 (deposited at [rid.ncifcrf.gov](http://rid.ncifcrf.gov)) (42).

707

708

709

710

711

712

713

714

715

716

717

718

719

720

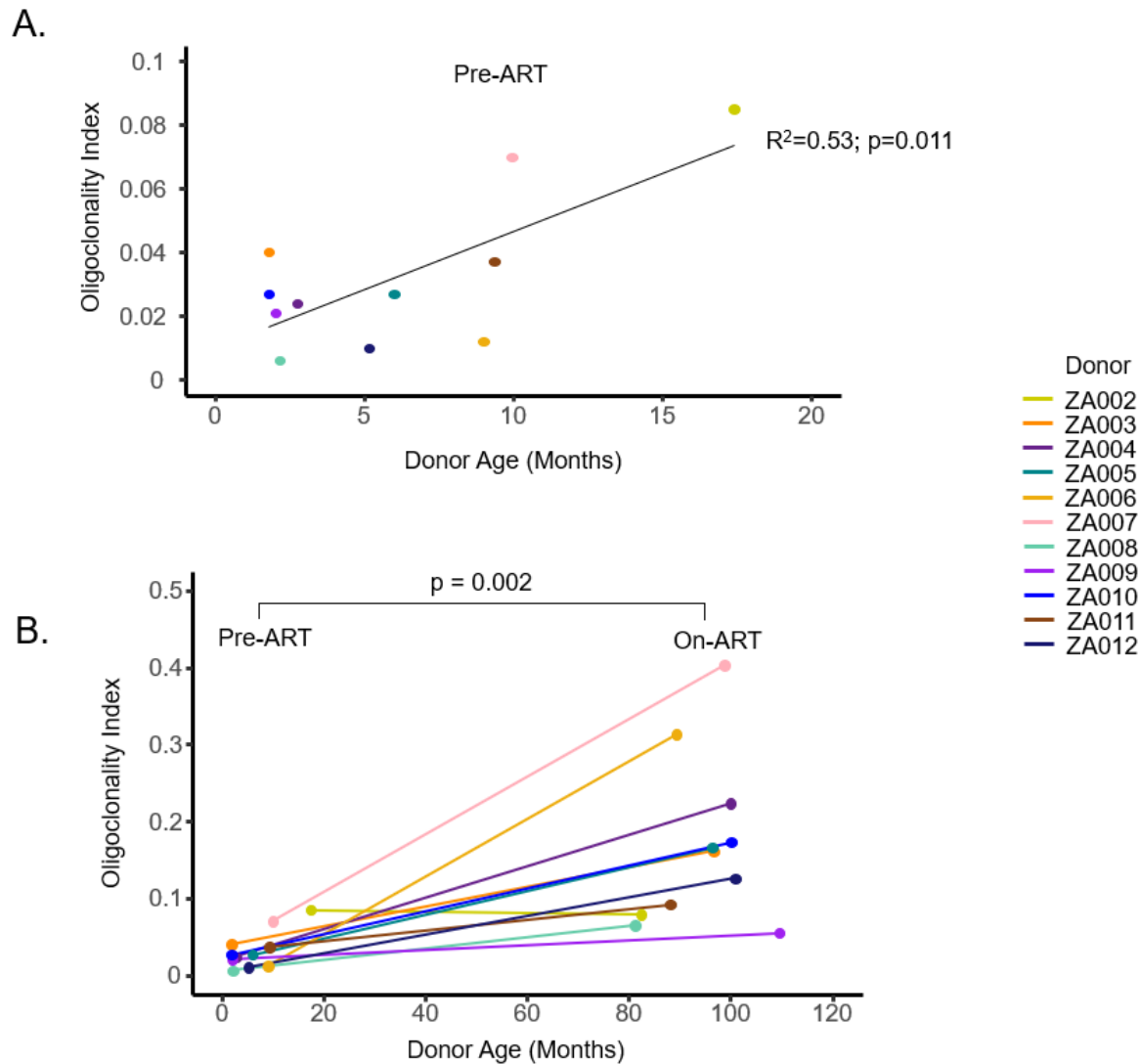
721

722

723 **Figures and Tables**

724

725 **Figure 1. Oligoclonality indexes correlate both with time on ART and duration of infection**  
726 **prior to ART.**



727

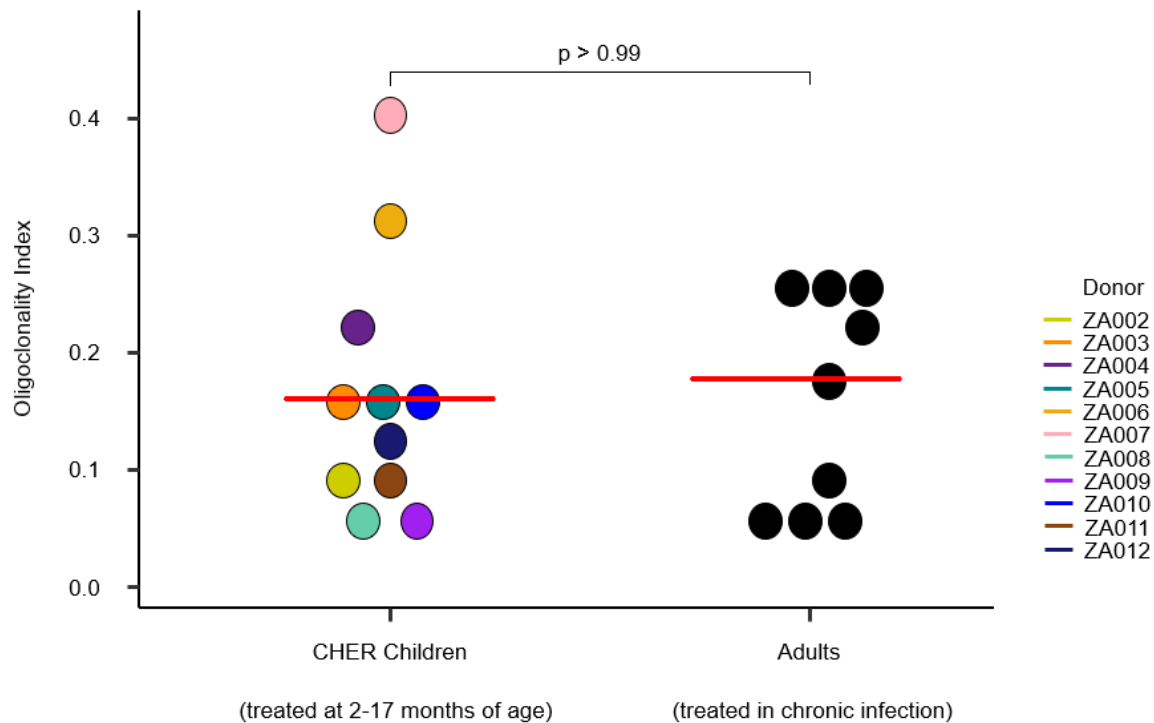
728 Oligoclonality indexes (OCI) were calculated from the pre-ART and on-ART libraries plotted  
729 against Donor Age in months. Pre-ART OCIs were evaluated via linear regression and F-test  
730 against donor age (A) while change in OCI as a function of ART status was evaluated by Wilcoxon  
731 Signed-Rank test (B).

732

733

734

735 **Figure 2. Oligoclonality indexes are comparable between ART-suppressed adults and**  
736 **children.**



737

738 Integration site data from donors whose viremia was suppressed on ART were downloaded from  
739 the Retrovirus Integration Database ([rid.ncifcrf.gov](http://rid.ncifcrf.gov)) (43) from two studies totaling 9 individuals  
740 (6, 32) and the OCIs were calculated. OCIs were compared using Mann-Whitney test. Median  
741 values for each patient group are marked by red lines.

742

743

744

745

746

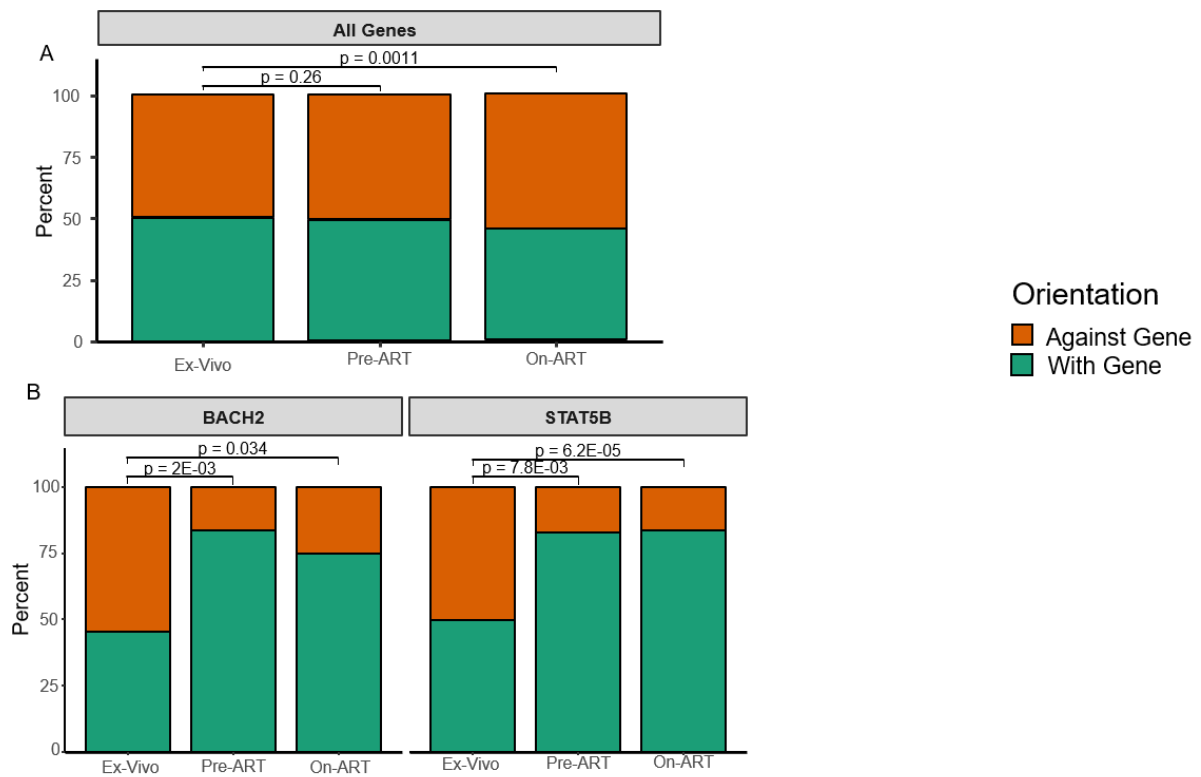
747

748

749

750

751 **Figure 3. Global selection against proviruses oriented with-the-gene and selection in two**  
752 **genes for with-the-gene proviruses.**



753

754 For each of the integration site libraries, unique integration sites for all genes (**A**) and for proviruses  
755 integrated in *BACH2* and *STAT5B* (**B**) were plotted as the percentage of integrations against-the-  
756 gene (orange) and with-the-gene (green). Significance was assessed via Fisher's Exact test  
757 between the *ex vivo* infected PBMC library and the pre-ART and on-ART integration site libraries  
758 from children. p-Values for pre-ART comparisons were post-hoc adjusted. The on-ART  
759 comparisons were not adjusted because of the differences in the number of independent statistical  
760 tests against each library.

761

762

763

764

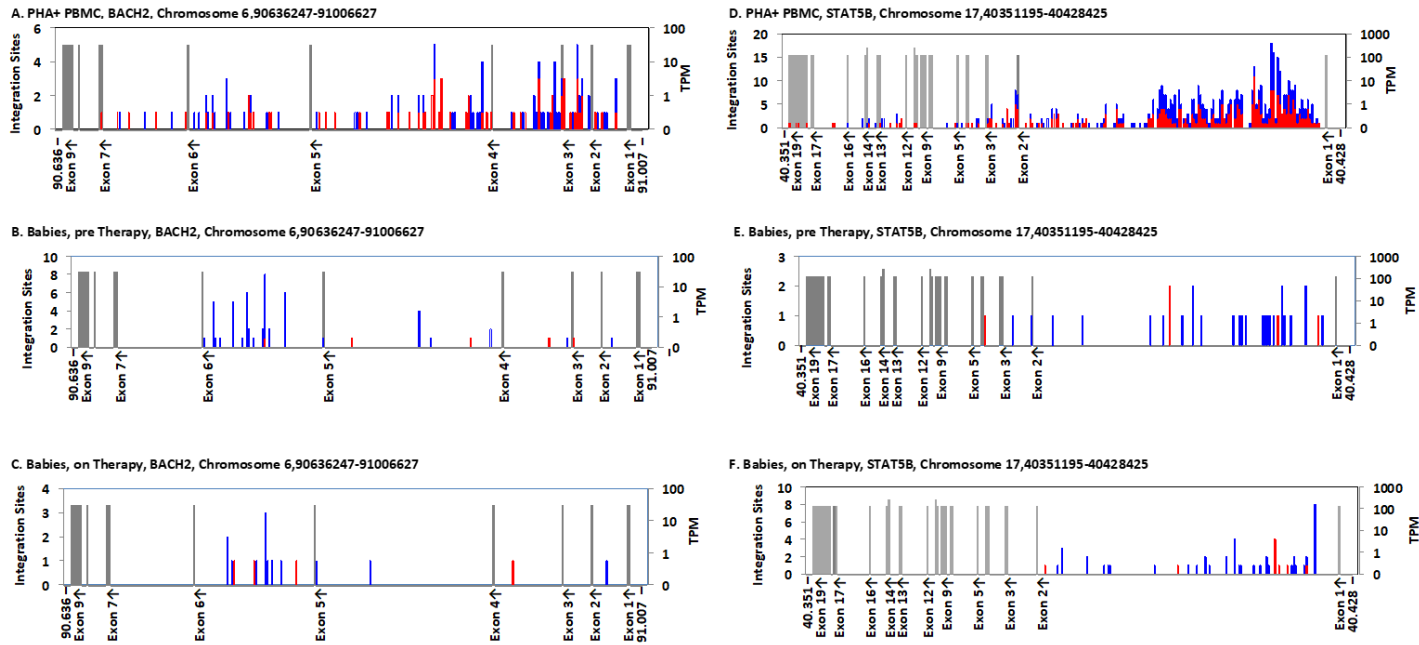
765

766

767



768 **Figure 4. Distribution of Integration sites in *BACH2* and *STAT5B*.**



769

770 The maps show the integration sites of the proviruses in *ex vivo* infected PBMC (A, D) (36), infants  
771 sampled pre-ART (B, E), and children sampled on ART (C, F). Each graph shows the entire gene,  
772 divided into 250 bins. For *BACH2* (A-C), each bin corresponds to ca 1500 NT; and for *STAT5B*  
773 (D-F), ca 300 NT. Exons (labeled on the X axis, with orientation of transcription shown) are shown  
774 as grey bars, whose height indicates the level of expression, in transcripts per million (TPM), as  
775 shown on the scale on the right. Note that the resolution of the text sometimes leads to loss of  
776 labels of closely spaced exons. The numbers of integration sites in each bin are indicated by the  
777 stacked bars, according to the scale on the left, with red indicating the same transcriptional  
778 orientation as the chromosome numbering and blue indicating the opposite orientation. In these  
779 two genes, blue indicates the number of proviruses in each bin integrated in the same orientation  
780 as the gene.

781

782

783

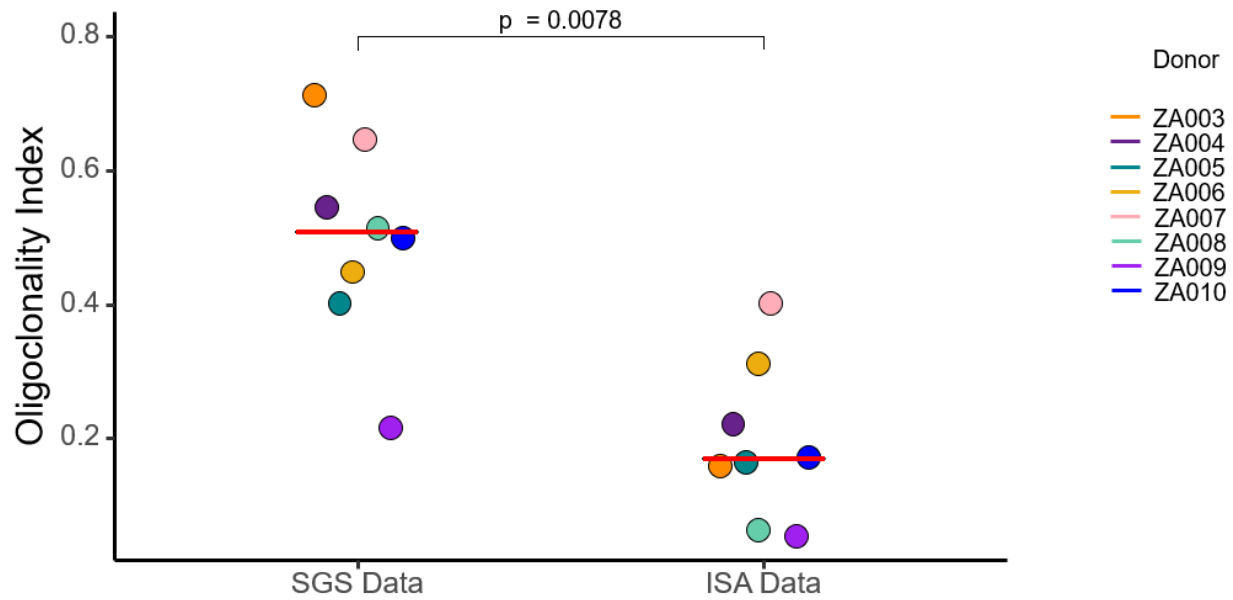
784

785

786

787

788 **Figure 5. OCIs for single-genome sequencing datasets are significantly higher than OCIs**  
789 **derived from integration sites analyses.**



790

791 OCIs were calculated from single-genome sequencing and integration sites data obtained from  
792 PBMC of children suppressed for 6-9 years on ART. Significance was assessed by Wilcoxon  
793 Signed-Rank Test. Median values are noted by red dash for each group.

794

795 **Table 1. Number of Integration Sites and Infected Cell Clones Detected in Children Prior to and On ART.**

Donor ID	Age at ART initiation (months)	No. of integration sites obtained pre-ART <sup>a</sup>	No. of integration sites detected with >2 breakpoints pre-ART (>1 breakpoint)	Oligoclonality index <sup>b</sup> (pre-ART)	Years suppressed on-ART	HIV DNA cps/10 <sup>6</sup> PBMC on-ART <sup>c</sup>	Number of integration sites obtained on-ART <sup>a</sup>	No. of integration sites detected with >1 breakpoint on-ART	Oligoclonality index <sup>b</sup> (on ART)	No. of Integration site matches between pre- and on- ART <sup>d</sup>
ZA002	17.4	1064	27 (50)	0.085	6.87	33	113	9	0.079	7
ZA003	1.8	1386	7 (30)	0.04	8.06	2	148	16	0.161	0
ZA004	2.7	655	5 (13)	0.024	7.92/8.76	24/--	255	25	0.223	3
ZA005	6.0	583	2 (14)	0.027	8.04	9	77	4	0.166	1
ZA006	9.0	486	1 (6)	0.012	7.45	47	137	20	0.313	1
ZA007	9.9	197	4 (5)	0.07	8.24	21	85	8	0.403	2
ZA008	2.2	1293	1 (8)	0.006	6.77	42	225	11	0.065	1
ZA009	2.0	809	0 (16)	0.021	9.13	186	125	5	0.055	0
ZA010	1.8	514	1 (12)	0.027	8.35	5	115	3	0.173	0
ZA011	9.3	432	5 (9)	0.037	7.35	182	149	8	0.092	3
ZA012	5.1	1243	3 (12)	0.01	8.41	12	432	32	0.126	2

Median Values	5.1	655	3 (12)	0.027	8.04	24	137	9	0.161	1
---------------	-----	-----	--------	-------	------	----	-----	---	-------	---

796

797 <sup>a</sup>Value obtained by counting an integration from both the 5' and 3' LTR as a single integration site

798 <sup>b</sup>As described in Gillet, *et al.* and Bangham (41). This value ranges in the interval [0, 1] dependent on the relative size and contribution  
799 of integration site clones to the dataset where 0 signifies a completely uniform distribution while 1 signifies a single integration site.

800 <sup>c</sup>Integrase cell-associated DNA (iCAD) protocol (51)

801 <sup>d</sup>Matches between pre-ART and on-ART are counted as clones in columns 4 and 9.

802 **Table 2. Analysis of Enrichment of Integration into Specific Genes in vivo<sup>a</sup>**

Chromosome	Gene name <sup>b</sup>	Independent integrations <i>ex vivo</i> <sup>c</sup>	Independent integrations in CHER cohort pre-ART	Adjusted p-value <sup>d</sup>	Independent integrations in CHER cohort on-ART	Adjusted p-value <sup>d</sup>
17	<i>STAT5B</i>	562	29	0.14	37	4.0E-29
6	<i>BACH2</i>	132	31	8.9E-17	16	2.7E-15
All	All other genes	334,920	6,979	>0.05	1,149	>0.05

803 <sup>a</sup>Data shown only for integrations into genes and for which at least 1 integration was detected in  
804 both libraries

805 <sup>b</sup>Genic coordinates mapped to hg19

806 <sup>c</sup>*Ex-vivo* dataset contains integration sites from CD8-depleted PBMCs from two healthy donor  
807 patients infected and PHA-stimulated *ex vivo*

808 <sup>d</sup>Adjusted p-value determined by Fisher's Exact Test with post-hoc Benjamini-Hochberg  
809 Correction

810

811

812

813

814

815

816

817

818

819

820

821

822

823 **Supplementary Materials**

824 **Table S1. Donor Characteristics.**

PID	Sex	Time to viral load suppression in years	Pre-ART plasma HIV RNA <sup>a</sup>	ART regimen <sup>b</sup>	Years suppressed on ART	CD4% at on-ART time point
ZA002	Male	1.37	654000	AZT/3TC/LPV/r	6.87	33
ZA003	Male	0.46	>750000	ABC/3TC/LPV/r	8.06	21
ZA004	Male	1.38	>750000	AZT/3TC/LPV/r	7.92/8.76	46
ZA005	Male	0.47	>750000	AZT/3TC/LPV/r	8.04	41
ZA006	Female	0.44	635000	AZT/3TC/LPV/r	7.45	50
ZA007	Male	0.92	>750000	AZT/3TC/EFV	8.24	35
ZA008	Female	0.44	>750000	AZT/3TC/LPV/r	6.77	36
ZA009	Female	3.76	>750000	AZT/3TC/EFV	9.13	29
ZA010	Female	0.46	510000	AZT/3TC/LPV/r	8.35	39
ZA011	Female	2.29	>750000	AZT/3TC/LPV/r	7.35	54
ZA012	Male	0.93	277,000	AZT/3TC/LPV/r	8.41	31

825 <sup>a</sup>Determined by Roche Amplicor HIV Monitor assay v1.0

826 <sup>b</sup>ART abbreviations: Zidovudine (AZT), Lamivudine (3TC), Ritonavir-boosted Lopinavir  
827 (LPV/r), Efavirenz (EFV)

828

829

830

831

832

833

834

835

836

837

838 **Table S2. Orientation Bias for Genic Integrations Pre-ART<sup>a</sup>.**

Chromosome	Gene name <sup>b</sup>	Unique integrations with the gene (CHER) <sup>c</sup>	Unique integrations against the gene (CHER) <sup>c</sup>	Unique integrations with the gene ( <i>ex vivo</i> ) <sup>c</sup>	Unique integrations against the gene ( <i>ex vivo</i> ) <sup>c</sup>	Adjusted p-value <sup>d</sup>
<i>chr6</i>	<i>BACH2</i>	26	5	60	72	0.0019
<i>chr17</i>	<i>STAT5B</i>	24	5	284	278	0.0078
<i>chr16</i>	<i>ANKRD11</i>	3	16	380	403	0.0281
<i>chr22</i>	<i>HORMAD2</i>	11	4	192	241	0.15
<i>chrX</i>	<i>MECP2</i>	4	12	228	267	0.45
<i>chr17</i>	<i>VMP1</i>	6	12	376	374	0.52
<i>chr17</i>	<i>GRB2</i>	11	5	339	336	0.52
<i>chr17</i>	<i>NPLOC4</i>	5	10	418	408	0.52
<i>chr17</i>	<i>POLR2A</i>	6	10	161	177	0.82
<i>chr11</i>	<i>MALAT1</i>	6	9	113	96	0.82
<i>chr11</i>	<i>PACSI1</i>	18	21	852	821	0.84
<i>chr11</i>	<i>KDM2A</i>	9	10	806	724	0.84
<i>chr19</i>	<i>CARD8</i>	10	7	331	315	0.84
<i>chr19</i>	<i>VAV1</i>	7	8	255	222	0.84
<i>chr17</i>	<i>RPTOR</i>	14	12	820	807	0.89
<i>chr17</i>	<i>CYTH1</i>	10	9	352	363	0.89
<i>chr22</i>	<i>TNRC6B</i>	9	10	444	435	0.89
<i>chr1</i>	<i>ASH1L</i>	8	9	378	385	>0.99

839 <sup>a</sup>Data shown only for integrations into genes for which at least 15 unique integrations were  
840 detected

841 <sup>b</sup>Genic coordinates mapped to hg19

842 <sup>c</sup>“With” gene and “Against” gene defined as orientation of integrated provirus compared with the  
843 sense of the host gene it’s integrated into

844 <sup>d</sup>Adjusted p-value determined by Fisher Test with post-hoc Benjamini-Hochberg Correction

845

846

847 **Table S3. Orientation Bias for Genic Integrations On-ART<sup>a</sup>.**

Chromosome	Gene name <sup>b</sup>	Unique integrations with the gene (CHER) <sup>c</sup>	Unique integrations against the gene (CHER) <sup>c</sup>	Unique integrations with the gene ( <i>ex vivo</i> ) <sup>c</sup>	Unique integrations against the gene ( <i>ex vivo</i> ) <sup>c</sup>	p-value <sup>d</sup>
<i>chr17</i>	<i>STAT5B</i>	31	6	284	278	6.2E-05
<i>chr6</i>	<i>BACH2</i>	12	4	60	72	0.034

848 <sup>a</sup>Data shown only for integrations into genes for which at least 15 unique integrations were  
849 detected in vivo and at least 1 unique integration ex vivo

850 <sup>b</sup>Genic coordinates mapped to hg19

851 <sup>c</sup>“With” gene and “Against” gene defined as orientation of integrated provirus compared with the  
852 sense of the host gene it’s integrated into

853 <sup>d</sup>p-Value determined by Fisher Test – no post-hoc adjustments performed

854

855

856

857

858

859

860

861

862

863

864

865

866

867

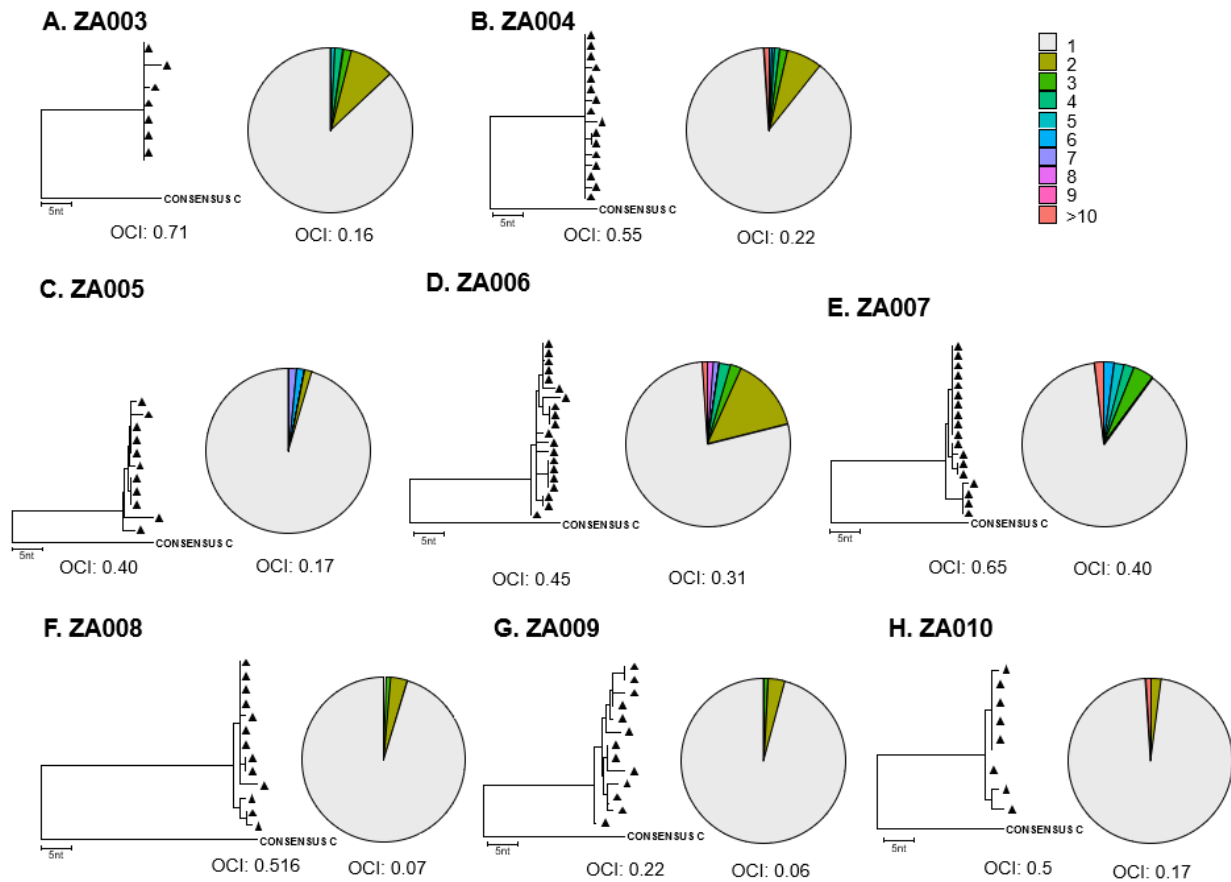
868

869



870 **Figure S1. Number of Detections of Integration Sites.**

871 FF



872

873 For each study participant, a neighbor-joining phylogenetic tree representing gag-pol single  
874 genome sequences with its respective OCI value is shown on the left; on the right, a pie chart  
875 representing the number of detections of integrations sights by ISA and the respective OCI value.