

Serveur Académique Lausannois SERVAL serval.unil.ch

Author Manuscript

Faculty of Biology and Medicine Publication

This paper has been peer-reviewed but does not include the final publisher proof-corrections or journal pagination.

Published in final edited form as:

Title: Analysis of HIV-1 expression level and sense of transcription by high-throughput sequencing of the infected cell.

Authors: Lefebvre G, Desfarges S, Uyttebroeck F, Muñoz M, Beerenwinkel N, Rougemont J, Telenti A, Ciuffi A

Journal: Journal of virology

Year: 2011 Jul

Volume: 85

Issue: 13

Pages: 6205-11

DOI: 10.1128/JVI.00252-11

Analysis of HIV-1 expression level and sense of transcription by high-throughput sequencing of the infected cell

Running title: HIV-1 expression by SAGE-Seq

Gregory Lefebvre¹, Sébastien Desfarges², Frédéric Uyttebroeck², Miguel Muñoz², Niko Beerenwinkel^{3,4}, Jacques Rougemont¹, Amalio Telenti^{2*§} and Angela Ciuffi^{2*§}.

Institute of Microbiology, University Hospital Center and University of Lausanne, Lausanne, Switzerland.

¹Bioinformatics and Biostatistics Core Facility, School of Life Sciences, EPFL, Lausanne, Switzerland.

²Institute of Microbiology, University Hospital Center and University of Lausanne, Lausanne, Switzerland.

³Department of Biosystems Science and Engineering, ETH Zurich, Basel, Switzerland.

⁴SIB – Swiss Institute of Bioinformatics, Basel, Switzerland

§ Last co-authors

*Correspondent footnote:

Angela Ciuffi

Institute of Microbiology

University Hospital Center and University of Lausanne

Bugnon 48

CH-1011 Lausanne

Switzerland

Phone: +4121314.40.99

Fax:+4121314.40.60

e-mail: Angela.Ciuffi@chuv.ch

Amalio Telenti

Institute of Microbiology

University Hospital Center and University of Lausanne

Bugnon 48

CH-1011 Lausanne

Switzerland

Phone: +4121314.40.97

Fax:+4121314.40.60

e-mail: Amalio.Telenti@chuv.ch

1 **Abstract**

2 Next-generation sequencing offers an unprecedented opportunity to jointly analyze cellular and
3 viral transcriptional activity, without prerequisite knowledge on the nature of the transcripts.
4 SupT1 cells were infected with a VSV-G pseudotyped HIV vector. At 24 hours post-infection,
5 both cellular and viral transcriptomes were analyzed by Serial Analysis of Gene Expression
6 followed by high-throughput sequencing (SAGE-Seq). Read mapping resulted in 33 to 44 million
7 tags aligning to the human transcriptome and 0.23 to 0.25 million tags aligning to the genome of
8 the HIV-1 vector. Thus, at peak infection, one transcript in 143 is of viral origin (0.7%),
9 including a small component of antisense viral transcription. Out of the detected cellular
10 transcripts, 826 (2.3%) were differentially expressed between mock and HIV-infected samples.
11 The approach also assessed whether HIV-1 infection modulates expression of repetitive elements
12 or endogenous retroviruses. We observed very active transcription of these elements, with one
13 transcript in 237 being of such origin, corresponding in average to 123,123 reads in mock
14 (0.40%) and 129,149 reads in HIV-1 (0.45%) mapping to the genomic Repbase repository. This
15 analysis highlights key details in the generation and interpretation of high-throughput data in the
16 setting of HIV-1 cellular infection.

17

18 **Introduction**

19 There is significant interest in determining the level of HIV-1 transcription in the context of the
20 infected cell. The impact of HIV-1 infection on cellular gene expression has been investigated in
21 the past by gene expression arrays (3, 7, 9, 10, 15-18, 20, 21, 29, 31, 36, 37, 42, 47, 48, 51).

22 However, use of these arrays is limited by the set of probes that are included in the chip.

23 Recently, it has even been proposed that unbiased analysis of transcription activity by deep
24 sequencing will soon replace gene expression profiling by microarrays (6, 45, 49). The first
25 application of the novel technologies in the field of HIV-1 has been in the assessment of viral
26 sequence variation, in particular mutations present at low frequency in complex (quasispecies)
27 populations (2, 8, 12, 19, 22, 23, 43, 46, 50, 52, 55). Pyrosequencing approaches have also been
28 instrumental in the assessment of small non-coding RNAs in HIV-1 infected cells (54). The next
29 goal is the joint analysis of the viral and host transcriptome of the infected cell.

30

31 Two general methods of deep sequencing are defined on the basis of the mode of sample
32 preparation (28). One approach uses fractionation of polyadenylated RNAs (valid also for non
33 poly(A)+ RNAs) in short fragments, followed by reverse transcription using hexamers. Adapters
34 are ligated to both ends, and used for sequencing. This method allows the detailed
35 characterization of the transcript structure. A second approach uses Serial Analysis of Gene
36 Expression (SAGE) to allow precise quantization of poly(A)+ RNA and to facilitate information
37 on strandedness, a key for the understanding of antisense transcription of the HIV-1 genome. An
38 important consideration regarding these technologies is the growing consensus that they will
39 represent the new gold-standard in the analysis of transcription – deep sequencing-based

40 expression analysis has shown to represent a major advance in robustness, resolution and inter-
41 lab portability over multiple microarray platforms (6, 45, 49).

42
43 In the present study, we used Super-SAGE, followed by high-throughput sequencing (SOLiD,
44 Life Technologies), also known as Digital Gene Expression (DGE) tag profiling, 3' tag DGE, tag
45 sequencing (Tag-Seq), or SAGE-Seq (32, 33, 38, 53), to analyze the transcriptome of a T cell line
46 24 hours post-infection.

47

48 **Materials and Methods**

49 **Cells.** HEK 293T cells were cultured in Dulbecco's Modified Eagle Medium (DMEM,
50 Invitrogen), supplemented with 10% heat-inactivated fetal calf serum (FCS) and 50 µg/ml
51 gentamycin (D-10 culture medium). SupT1 cells (a T-cell line) were cultured in RPMI-1640
52 (Invitrogen), supplemented with 10% heat-inactivated FCS and 50 µg/ml gentamycin (R-10
53 culture medium).

54

55 **HIV-vector production.** To produce HIV-based vector particles, 293T cells were co-transfected
56 with two plasmids (20 µg total) using the calcium phosphate method (Invitrogen) and according
57 to the manufacturer's instructions. One plasmid, pNL4-3Δenv-eGFP, codes for all viral proteins
58 except the envelope, which was disrupted and replaced by GFP (56). The second plasmid,
59 pMD.G, encodes the vesicular stomatitis virus G envelope protein (VSV-G) (39). Forty-eight
60 hours after transfection, culture supernatant containing viral particles was collected, centrifuged
61 to pellet cell debris, filtered through 0.22 µm filters, concentrated using Centricon-Plus 70-100K
62 (Millipore), treated with 100 U/ml DNaseI (Roche) and stored frozen at -80°C. Virion

63 concentration was assessed by measuring the CA (p24) antigen by ELISA (Murex HIV Ag MAB;
64 Abbott), according to instructions.

65

66 **HIV infection.** SupT1 cells (5×10^6) were infected with 15 μg p24 equivalent of HIV NL4-
67 $3\Delta\text{env-eGFP/VSV-G}$ in presence of 5 $\mu\text{g/ml}$ polybrene (Sigma) by centrifugation for 30 min at
68 1500 g. Mock cells were treated similarly but without HIV-based particles. Cells were washed
69 with culture medium, resuspended at 10^6 cells/ml in R-10 and further incubated.

70

71 **Transcriptional profiling by SAGE-Seq and bioinformatic analysis.** At 24h post-exposure,
72 cells were washed once with phosphate-buffered saline (PBS) and resuspended in 0.5 ml
73 RNALater (Ambion). Total RNA was purified using miRvana isolation kit (Ambion), or TRIzol
74 (Invitrogen) for the second experiment, according to recommendations. RNA quality was
75 assessed by electrophoresis on the Bioanalyzer 2100 (Total RNA Nano; Agilent). SAGE libraries
76 were prepared with SOLiD SAGE kit (Applied Biosystems) at the Functional Genomics Center
77 Zurich. High-throughput sequencing was performed using SOLiD3 technology (Applied
78 Biosystems). Sequence mapping was done with the bowtie software (27), using defined analysis
79 parameters (read length (rl), number of mismatches (n) and alignments authorized (m)). Tag hits
80 (*i.e.* successfully aligned reads) were normalized according to the number of locations they
81 mapped to. Alignment used the human genome GRCh37, the human transcriptome RefSeq
82 database, and HIV-1 genome *HIV NL4-3 $\Delta\text{env-eGFP}$* . Differential expression analysis was
83 performed using R programming language (14) and based on negative-binomial modeling of
84 count data as described in the Bioconductor package DESeq (1). Genes with significantly

85 modulated expression in HIV-1 samples compared to mock samples were selected with a false
86 discovery rate (FDR) of 0.05 and annotation was retrieved using the biomaRt package (11).
87

88 **Northern hybridization.** Total RNA (10 µg) from mock-treated or HIV vector-infected cells
89 was supplemented with RNA loading buffer (Bioline), heated 5 min at 65°C and separated by
90 electrophoresis on 1.2% agarose-MOPS-formaldehyde gel (44). The gel was washed 3x10 min
91 with 10xSSC. RNA transfer from gel to nitrocellulose membrane (Hybond N+; GE Healthcare)
92 was performed over night by capillarity using 20x SSC. The membrane was rinsed with 6xSSC,
93 crosslinked with UV (GeneLinker; BioRad), rinsed briefly with water and added to 10 ml
94 preheated QuickHyb Hybridization solution (Agilent) for 30 min at 65°C. 200 pmol
95 oligonucleotide probes were labeled with 20 µCi [γ -³²P]dATP with 10 U T4 polynucleotidyl
96 kinase (T4 PNK; New England Biolabs) for 1h at 37°C. After heat inactivation of the enzyme for
97 10 min at 95°C, probes were purified on Qiaquick Nucleotide removal kit (Qiagen), mixed with
98 100 µl sonicated salmon sperm DNA [10 mg/ml] (Eppendorf), and added to the membrane in
99 QuickHyb Hybridization solution for 1h at 65°C. Following hybridization, membrane was
100 washed 2x15 min with 2xSSC/0.1% SDS at room temperature, 1x30 min in 2xSSC/0.1% SDS at
101 52°C, exposed to an intensifying screen over night at -80°C and revealed using ImageQuantTL
102 v2005 software on Typhoon scanner (GE Healthcare). For re-use, membrane was stripped with
103 boiling 0.1xSSC/0.1%SDS twice for 15 min. Oligonucleotide probes were 5'-
104 TCTCTCTCAGGGTCATCCATTCCA-3' for peak 1 sense transcripts, 5'-
105 GCTCGTCCTTGTACAGCTCGTCCA-3' for peak 3 sense transcripts, 5'-
106 ATGGTGTTTTACTAATCTTTTCCATGTGTT-3' for peak 4 sense transcripts, 5'-
107 CTTGTTACACCCTGTGAGCCTGCA-3' for peak 1 antisense transcripts, 5'-

108 CCGCCGCCGGGATCACTCTCGGCA-3' for peak 3 antisense transcripts and 5'-
109 GTAGACAGGATGAGGATTAACACATGGAAA-3' for peak 4 antisense transcripts. As
110 positive control for hybridization, 2 µg pNL4-3Δenv-eGFP digested with BglII and AflIII was
111 used.

112

113 **Results and Discussion**

114 SupT1 cells (5×10^6 cells) were mock-treated or infected with a VSV-G pseudotyped HIV-based
115 vector expressing GFP (56) in duplicate (**Figure 1A**). At 36h post-infection, FACS analysis
116 revealed that 93% of cells were successfully transduced, expressing GFP; other cell types were
117 tested and found less efficient in transduction (**Supplemental File 1A**). The use of 24h time point
118 in the experimental model takes into consideration cell loss and viability, initiation of translation,
119 as measured by GFP expression (as soon as 16 hours), and the completion of the viral cycle, as
120 measured by the production and release of p24 (as soon as 20 hours), reflecting viral particle
121 release in SupT1 cells (**Supplemental File 1B**). Cells were collected for total RNA extraction
122 and poly(A)+ RNAs were reverse transcribed and processed for SAGE library preparation.
123 cDNAs were digested with NlaIII, a frequent 4-bp cutter, and ligated to an adaptor sequence,
124 containing the recognition site of EcoP15I; an enzyme that cuts DNA asymmetrically, 25/27bp
125 away. A second adapter was ligated leading to a 27bp transcript-specific tag, surrounded by two
126 distinct adapter sequences. High-throughput sequencing was performed using the SOLiD 3
127 system with a universal primer annealing to the adaptor sequence. Reads were mapped either to
128 the human genome (GRCh37), the human RefSeq database (40), the human repetitive-element
129 database (24), or the HIV-1 vector genome (*HIV NL4-3Δenv-eGFP/VSV-G*) (**Figure 1B**).

130

131 **Assessment of sequence alignment parameters.** Variation of mapping parameters for HIV-1
132 tags was tested to maximize the number of tags mapping to the viral genome while minimizing
133 the number of tags spuriously mapping to the host genome (**Supplemental File 2**). After
134 shortening the reads by right-trimming, we used the read lengths of 21bp, 24bp, and 27bp. In
135 addition, we tested different number of mismatches (n=0, 1, or 2) and multiple hits (m=1, 2, or
136 10). We retained condition rl=24, n=1 and m=1 as optimal for read mapping to the human
137 genome, rl=24, n=1 and m=10 for mapping to the human transcriptome, and rl=24, n=2 and m=2
138 for mapping to the HIV-1 genome.

139
140 **HIV-1 RNA tags identify sense and antisense transcription.** Read mapping resulted in 33 to
141 44 million tags aligning to the human transcriptome and 0.23 to 0.25 million tags aligning to the
142 genome of the HIV-1 vector (**Table 1**). These data suggest that, at peak infection, one transcript
143 in 143 is of viral origin (0.7%). Overall, 0.33% were HIV-1-specific (1 in 309) and 0.37% were
144 vector GFP-specific (1 in 267).

145
146 The distribution of tags on the HIV-1 vector genome emphasized several relevant aspects of
147 HIV-1 transcription (**Figure 2A**). Five viral genomic regions carried 87.8% of total tags detected.
148 A high proportion of tags (24.2%) mapped to the 3' end of the HIV-1 vector genome (**Figure 2A,**
149 **Peak 1**), corresponding to sense transcription with a known functional polyadenylation site of
150 HIV-1 (genome positions 9455-9460), and the first upstream NlaIII restriction site (pos. 9150).
151 The second signal (**Figure 2A, Peak 2**), carrying 5.1% of total tags, corresponds to sense
152 transcription and the same known functional polyadenylation site of HIV-1 or an alternative
153 putative signal (at pos. 9108-9113), up to the second upstream NlaIII restriction site (pos. 8757).
154 The third signal (**Figure 2A, Peak 3**), with 51.6% total tags, mapped at the end of the *gfp* orf,

155 with sense transcription, ending at a putative polyadenylation signal at pos. 7133-7138, reaching
156 an upstream NlaIII restriction site at pos. 7072. The massive presence of these gfp tags suggests
157 that viral transcription is skewed in favor of stable gfp transcripts, rendering this vector
158 particularly suited as a reporter of HIV infection. Nevertheless, it seems unlikely that this
159 putative poly(A) signal (located in the *env* orf) is used during transcription of wild type HIV-1.
160 Further investigation of viral transcripts using wild-type HIV-1 should confirm this. The fourth
161 signal (**Figure 2A, Peak 4**), representing 4.8% of total tags, is consistent with antisense viral
162 transcription using a polyadenylation signal at pos. 4908-4903, reported previously by Landry *et*
163 *al.* (26) , up to a first downstream NlaIII restriction site at pos. 5099. The fifth signal (**Figure 2A,**
164 **Peak 5**), with 2.1% of total tags, corresponds to antisense transcription ending at an unidentified
165 polyadenylation site, and up to the first downstream NlaIII restriction site at pos. 1226. The
166 remaining 12.2% of tags mapping elsewhere to the genome correspond mostly to the use of
167 further NlaIII restriction sites or alternative, putative poly(A) sites.

168
169 The technical procedure, from SAGE library preparation to high-throughput sequencing, keeps
170 track of strand specificity. Sequences are read from a universal primer in the first adapter (on the
171 poly(A) side) towards the NlaIII site. Sense tags will map 3' from the NlaIII site, while antisense
172 tags will map 5' of the NlaIII site. Globally, 20.5% of tags mapped to the HIV-1 vector genome
173 in the antisense orientation (**Figure 2A**). Peaks 4 and 5 are comprised mostly of antisense tags
174 (88% and 92%, respectively); however, antisense tags were also present at the 3' end of the viral
175 genome (40.5% and 49.0%, for peaks 1 and 2, respectively). In contrast, a negligible amount of
176 tags in the antisense orientation mapped to the third peak, corresponding to GPF (only 1 out of
177 125,970 tags). With the exception of peak four, we did not identified canonical poly(A) signals
178 (AATAAA, ATTAAA, AGTAAA, AAGGAA) that would associate with the observed

179 accumulation of antisense tags. Confirming and understanding HIV-1 antisense transcription is
180 important because of the a possible roles for such transcripts in the regulation of viral expression
181 (54), the generation of antisense proteins (26, 30, 34, 35), or the production of cryptic epitopes
182 (4). However, the assessment of antisense transcription and cognate poly(A) motifs is not
183 straightforward. Specifically, Landry *et al.* indicated that RT-PCR cannot reliably separate sense
184 vs anti-sense transcripts due to endogenous RT priming, and would therefore require the use of 5'
185 LTR-deleted pNL4-3 constructs , and cloning and sequencing of the 3'RACE amplified products
186 (26). First, we repeated a deep sequencing experiment, using barcoded adapters (initially
187 designed for sample multiplexing) during the library preparation of SAGE-Seq, as it was
188 suggested that barcoding samples might generate less artifacts (53). Although the overall
189 distribution profile of the tags was comparable, only about 0.9% antisense tags were detected. To
190 further estimate the proportion of viral antisense transcription, we performed Northern blots
191 using specific sense and antisense probes (**Figure 2B**). 2-kb, 4-kb and 9-kb classes of viral sense
192 transcripts were detected expectedly according to the specific probe used. However, none of the
193 probes designed to anneal to antisense transcripts revealed any detectable signal. These data
194 suggest that viral antisense transcription is more consistent with 0.9% abundance than with
195 20.5% abundance, calling for caution when analyzing antisense transcription by novel high-
196 throughput technologies, and highlighting the necessity to validate data by alternative methods.
197

198 Three additional technical aspects deserve comment. The use of the HIV-1 vector containing a
199 GFP ORF in the place of *env* could disturb transcription and splicing. The Northern blot in **Figure 2B**
200 excludes significant defects in viral transcription and splicing. The integrity of transcription was
201 further assessed by RT-qPCR using primers flanking the major introns of HIV to confirm the detection of
202 multiply, singly and unspliced transcripts (data not shown). We examined the degree of conservation

203 of the tag sequence to assess whether the tags displayed a perfectly matched sequence
204 corresponding to the HIV-1 vector DNA. We tested 0, 1 and 2 mismatches to ensure that we
205 would capture the putative error rate of the viral reverse transcriptase, as well as base miscalling
206 of the SOLiD sequencing software. Upon HIV vector-specific tag analysis, 96.97% of the tags
207 displayed a perfect match with the HIV vector genome, 2.92% presented one mismatch and only
208 0.11% presented two mismatches. We also identified a limited number of reads (<0.01%) that
209 mapped to the HIV-1 vector genome in the mock infected cells (**Table 1**). Upon inspection, all
210 tags were HIV-specific, indicating a small level of contamination occurring through the
211 experimental procedure (from culture to library preparation).

212

213 **The cellular transcriptome during infection.** The high-throughput sequencing analysis of
214 SAGE-RNA allowed the detection of 36,271 human expressed genes (79.8%), according to the
215 RefSeq database (**Table 1**). On average, each transcript was represented by 815 tags (ranging
216 from 0 to 3.8 millions, and with a median of 14.7). While the overall distribution in expression
217 levels was not different for mock and HIV-1 infected cells (**Figure 3**), the identity of transcripts
218 contained in mock or HIV-infected samples was significantly modified. Out of the detected
219 transcripts, 826 (2.3%) were significantly differentially modulated between mock and HIV-
220 infected samples at a stringent adjusted P value of 10^{-4} at 24h post-infection. Ingenuity pathway
221 analysis (www.ingenuity.com) identified “cellular growth and proliferation” and “RNA post-
222 transcriptional modification” as prominent modulated functions, including as relevant networks
223 those related to cell cycle, DNA replication and repair, gene expression, and cell death. There
224 was down regulation in HIV-infected cells of many of the genes associated to those networks,
225 consistent with marked cellular compromise and stress. The complete set of differentially
226 expressed genes is included in **Supplemental File 3**. Of notice, some of the differentially

227 expressed genes include non-protein coding transcripts such as spliceosomal and small nucleolar
228 RNAs that are not generally thought of as polyadenylated, and should not have been captured by
229 SAGE. However, there is increasing evidence that poly(A) tails can be added to such RNAs,
230 possibly marking them for degradation (5, 25).

231
232 We assessed the overlap between the set of genes that were identified as differentially expressed
233 in the current study with those from various microarray studies in the literature that investigated
234 differential expression in CD4+ T cells (42), CD8+ T cells (41), monocytes (17), and lymph
235 nodes (29) during HIV-1 infection *in vivo*, as well as with a set of validated genes compiled by
236 Giri *et al.* (16) from microarray studies published between 2000 and 2006. The overlap varied
237 from 4% in the study of circulating monocytes in HIV-infected individuals (17) to 8% in the
238 analysis of lymphatic tissue in the setting of various stages of HIV-1 infection *in vivo* (29). There
239 was also a 6% overlap with a set of validated genes curated by Giri *et al.* (16). The final shortlist
240 of 52 genes common between the present and one or more of the previous studies (**Supplemental**
241 **File 4**) was also enriched for genes involved in cell cycle, DNA replication and repair, and cell
242 proliferation. However, it should be underscored that differences in the nature of the techniques
243 (deep-sequencing vs. microarray), a lack of recent studies using new generation microarrays for
244 the transcriptome analysis of *in vitro* infection of T cell lines - the comparable experiment as the
245 one completed herein, and general differences in study design makes comparison between studies
246 of unclear significance.

247 The approach also assessed whether HIV-1 infection modulates expression of repetitive elements
248 or endogenous retroviruses. We observed very active transcription of these elements, with one
249 transcript in 237 being of such origin, corresponding in average to 123,123 reads in mock

250 (0.40%) and 129,149 reads in HIV-1 (0.45%) mapping to the genomic Repbase repository
251 (**Figure 3**). Of these, 4021 and 4451 average reads mapped to human endogenous retroviruses
252 (HERV) in mock and in HIV-infected cells, respectively (**Figure 3**). The modest increase in
253 HERV in the HIV-infected cell results mainly from the contribution of HERVK which increases
254 from 3,331 to 3,814 average tags. The biological role of the prominent transcriptional activity of
255 endogenous retroelements, *i.e.* regulatory activity or generation of translated products, needs
256 further analysis (13).

257
258 **Conclusions.** Analysis of deep sequence data using SAGE allowed the precise measurement of
259 expression level of the proviral genome in HIV-1-infected cells and identified a small component
260 of antisense viral transcription. Improvements in RNA-Seq will increasingly deliver information
261 on both strand specificity and nature of the transcript, including splice forms, and sequence
262 variation, that will facilitate the study of the dynamics of viral-host interactions (28). Extending
263 this approach to replication-competent HIV-1 isolates, and to different cellular backgrounds, may
264 reveal differences in viral-host interactions due to specific strain or cellular factors.

265

266 **Authors' contributions**

267 GL and JR performed the bioinformatics analyses, FU and MM carried out HIV-1 production,
268 cell infections and RNA extraction, SD performed Northern blot analyses, AT and AC lead the
269 project and wrote the manuscript, NB edited and proofread the manuscript.

270

271 **Acknowledgements**

272 We thank Dr Marzanna Künzli and Dr Sirisha Aluri from the Functional Genomics Center Zurich
273 (FGCZ) facility for SAGE library preparation and SOLiD sequencing, with the great help and
274 support from Dr Gerrit Kuhn (Life Technologies). The following reagent was obtained through
275 the NIH AIDS Research and Reference Reagent Program, Division of AIDS, NIAID, NIH:
276 pNL4-3-deltaE-EGFP (Cat# 11100) from Drs. Haili Zhang, Yan Zhou, and Robert Siliciano (56).
277 This work was supported by the Swiss National Science Foundation, grant 310030-130699.

278

279 **Competing interests**

280 The authors declare that they have no competing interests.

281

282 **References**

- 283 1. **Anders, S., and W. Huber.** 2010. Differential expression analysis for sequence count
284 data. *Nature Precedings*:1-5.
- 285 2. **Archer, J., M. S. Braverman, B. E. Taillon, B. Desany, I. James, P. R. Harrigan, M.**
286 **Lewis, and D. L. Robertson.** 2009. Detection of low-frequency pretherapy chemokine
287 (CXC motif) receptor 4 (CXCR4)-using HIV-1 with ultra-deep pyrosequencing. *AIDS*
288 (London, England) **23**:1209-1218.
- 289 3. **Arendt, C. W., and D. R. Littman.** 2001. HIV: master of the host cell. *Genome Biol*
290 **2**:REVIEWS1030.
- 291 4. **Bansal, A., J. Carlson, J. Yan, O. T. Akinsiku, M. Schaefer, S. Sabbaj, A. Bet, D. N.**
292 **Levy, S. Heath, J. Tang, R. A. Kaslow, B. D. Walker, T. Ndung'u, P. J. Goulder, D.**
293 **Heckerman, E. Hunter, and P. A. Goepfert.** 2010. CD8 T cell response and
294 evolutionary pressure to HIV-1 cryptic epitopes derived from antisense transcription. *J*
295 *Exp Med* **207**:51-59.
- 296 5. **Bayne, E. H., S. A. White, and R. C. Allshire.** 2007. DegrAAAAded into silence. *Cell*
297 **129**:651-653.
- 298 6. **Blow, N.** 2009. Transcriptomics: The digital generation. *Nature* **458**:239-242.
- 299 7. **Borjabad, A., A. I. Brooks, and D. J. Volsky.** 2010. Gene expression profiles of HIV-1-
300 infected glia and brain: toward better understanding of the role of astrocytes in HIV-1-
301 associated neurocognitive disorders. *J Neuroimmune Pharmacol* **5**:44-62.

- 302 8. **Bushman, F. D., C. Hoffmann, K. Ronen, N. Malani, N. Minkah, H. M. Rose, P.**
303 **Tebas, and G. P. Wang.** 2008. Massively parallel pyrosequencing in HIV research.
304 *AIDS (London, England)* **22**:1411-1415.
- 305 9. **Carroll, E. E., R. Hammamieh, N. Chakraborty, A. T. Phillips, S. A. Miller, and M.**
306 **Jett.** 2006. Altered gene expression in asymptomatic SHIV-infected rhesus macaques
307 (Macacca mulatta). *Virology journal* **3**:74.
- 308 10. **Chung, H. K., C. A. Pise-Masison, M. F. Radonovich, J. Brady, J. K. Lee, S. Y.**
309 **Cheon, P. Markham, A. Cristillo, and R. Pal.** 2008. Cellular gene expression profiles in
310 rhesus macaques challenged mucosally with a pathogenic R5 tropic simian human
311 immunodeficiency virus isolate. *Viral Immunol* **21**:411-423.
- 312 11. **Durinck, S., Y. Moreau, A. Kasprzyk, S. Davis, B. De Moor, A. Brazma, and W.**
313 **Huber.** 2005. BioMart and Bioconductor: a powerful link between biological databases
314 and microarray data analysis. *Bioinformatics* **21**:3439-3440.
- 315 12. **Eriksson, N., L. Pachter, Y. Mitsuya, S. Y. Rhee, C. Wang, B. Gharizadeh, M.**
316 **Ronaghi, R. W. Shafer, and N. Beerenwinkel.** 2008. Viral population estimation using
317 pyrosequencing. *PLoS computational biology* **4**:e1000074.
- 318 13. **Garrison, K. E., R. B. Jones, D. A. Meiklejohn, N. Anwar, L. C. Ndhlovu, J. M.**
319 **Chapman, A. L. Erickson, A. Agrawal, G. Spotts, F. M. Hecht, S. Rakoff-Nahoum, J.**
320 **Lenz, M. A. Ostrowski, and D. F. Nixon.** 2007. T cell responses to human endogenous
321 retroviruses in HIV-1 infection. *PLoS pathogens* **3**:e165.
- 322 14. **Gentleman, R. C., V. J. Carey, D. M. Bates, B. Bolstad, M. Dettling, S. Dudoit, B.**
323 **Ellis, L. Gautier, Y. Ge, J. Gentry, K. Hornik, T. Hothorn, W. Huber, S. Iacus, R.**
324 **Irizarry, F. Leisch, C. Li, M. Maechler, A. J. Rossini, G. Sawitzki, C. Smith, G.**

- 325 **Smyth, L. Tierney, J. Y. Yang, and J. Zhang.** 2004. Bioconductor: open software
326 development for computational biology and bioinformatics. *Genome Biol* **5**:R80.
- 327 15. **George, M. D., D. Verhoeven, Z. McBride, and S. Dandekar.** 2006. Gene expression
328 profiling of gut mucosa and mesenteric lymph nodes in simian immunodeficiency virus-
329 infected macaques with divergent disease course. *J Med Primatol* **35**:261-269.
- 330 16. **Giri, M. S., M. Nebozhyn, L. Showe, and L. J. Montaner.** 2006. Microarray data on
331 gene modulation by HIV-1 in immune cells: 2000-2006. *Journal of leukocyte biology*
332 **80**:1031-1043.
- 333 17. **Giri, M. S., M. Nebozyhn, A. Raymond, B. Gekonge, A. Hancock, S. Creer, C.**
334 **Nicols, M. Yousef, A. S. Foulkes, K. Mounzer, J. Shull, G. Silvestri, J. Kostman, R.**
335 **G. Collman, L. Showe, and L. J. Montaner.** 2009. Circulating monocytes in HIV-1-
336 infected viremic subjects exhibit an antiapoptosis gene signature and virus- and host-
337 mediated apoptosis resistance. *J Immunol* **182**:4459-4470.
- 338 18. **Harman, A. N., M. Kraus, C. R. Bye, K. Byth, S. G. Turville, O. Tang, S. K. Mercier,**
339 **N. Nasr, J. L. Stern, B. Slobedman, C. Driessen, and A. L. Cunningham.** 2009. HIV-
340 1-infected dendritic cells show 2 phases of gene expression changes, with lysosomal
341 enzyme activity decreased during the second phase. *Blood* **114**:85-94.
- 342 19. **Hoffmann, C., N. Minkah, J. Leipzig, G. Wang, M. Q. Arens, P. Tebas, and F. D.**
343 **Bushman.** 2007. DNA bar coding and pyrosequencing to identify rare HIV drug
344 resistance mutations. *Nucleic acids research* **35**:e91.
- 345 20. **Hyrza, M. D., C. Kovacs, M. Loutfy, R. Halpenny, L. Heisler, S. Yang, O. Wilkins,**
346 **M. Ostrowski, and S. D. Der.** 2007. Distinct transcriptional profiles in ex vivo CD4+
347 and CD8+ T cells are established early in human immunodeficiency virus type 1 infection

- 348 and are characterized by a chronic interferon response as well as extensive transcriptional
349 changes in CD8+ T cells. *Journal of virology* **81**:3477-3486.
- 350 21. **Imbeault, M., M. Ouellet, and M. J. Tremblay.** 2009. Microarray study reveals that
351 HIV-1 induces rapid type-I interferon-dependent p53 mRNA up-regulation in human
352 primary CD4+ T cells. *Retrovirology* **6**:5.
- 353 22. **Ji, H., N. Masse, S. Tyler, B. Liang, Y. Li, H. Merks, M. Graham, P. Sandstrom, and**
354 **J. Brooks.** 2010. HIV drug resistance surveillance using pooled pyrosequencing. *PLoS*
355 *ONE* **5**:e9263.
- 356 23. **Jojic, V., T. Hertz, and N. Jojic.** 2008. Population sequencing using short reads: HIV as
357 a case study. *Pac Symp Biocomput*:114-125.
- 358 24. **Jurka, J., V. V. Kapitonov, A. Pavlicek, P. Klonowski, O. Kohany, and J.**
359 **Walichiewicz.** 2005. Repbase Update, a database of eukaryotic repetitive elements.
360 *Cytogenet Genome Res* **110**:462-467.
- 361 25. **LaCava, J., J. Houseley, C. Saveanu, E. Petfalski, E. Thompson, A. Jacquier, and D.**
362 **Tollervey.** 2005. RNA degradation by the exosome is promoted by a nuclear
363 polyadenylation complex. *Cell* **121**:713-724.
- 364 26. **Landry, S., M. Halin, S. Lefort, B. Audet, C. Vaquero, J. M. Mesnard, and B.**
365 **Barbeau.** 2007. Detection, characterization and regulation of antisense transcripts in
366 HIV-1. *Retrovirology* **4**:71.
- 367 27. **Langmead, B., C. Trapnell, M. Pop, and S. L. Salzberg.** 2009. Ultrafast and memory-
368 efficient alignment of short DNA sequences to the human genome. *Genome Biol* **10**:R25.
- 369 28. **Levin, J. Z., M. Yassour, X. Adiconis, C. Nusbaum, D. A. Thompson, N. Friedman,**
370 **A. Gnirke, and A. Regev.** 2010. Comprehensive comparative analysis of strand-specific
371 RNA sequencing methods. *Nature methods* **7**:709-715.

- 372 29. **Li, Q., A. J. Smith, T. W. Schacker, J. V. Carlis, L. Duan, C. S. Reilly, and A. T.**
373 **Haase.** 2009. Microarray analysis of lymphatic tissue reveals stage-specific, gene
374 expression signatures in HIV-1 infection. *J Immunol* **183**:1975-1982.
- 375 30. **Ludwig, L. B., J. L. Ambrus, Jr., K. A. Krawczyk, S. Sharma, S. Brooks, C. B.**
376 **Hsiao, and S. A. Schwartz.** 2006. Human Immunodeficiency Virus-Type 1 LTR DNA
377 contains an intrinsic gene producing antisense RNA and protein products. *Retrovirology*
378 **3**:80.
- 379 31. **Martinez-Marino, B., H. Foster, Y. Hao, and J. A. Levy.** 2007. Differential gene
380 expression in CD8(+) cells from HIV-1-infected subjects showing suppression of HIV
381 replication. *Virology* **362**:217-225.
- 382 32. **Matsumura, H., D. H. Kruger, G. Kahl, and R. Terauchi.** 2008. SuperSAGE: a
383 modern platform for genome-wide quantitative transcript profiling. *Curr Pharm*
384 *Biotechnol* **9**:368-374.
- 385 33. **Matsumura, H., K. Yoshida, S. Luo, E. Kimura, T. Fujibe, Z. Albertyn, R. A.**
386 **Barrero, D. H. Kruger, G. Kahl, G. P. Schroth, and R. Terauchi.** 2010. High-
387 throughput SuperSAGE for digital gene expression analysis of multiple samples using
388 next generation sequencing. *PLoS ONE* **5**.
- 389 34. **Michael, N. L., M. T. Vahey, L. d'Arcy, P. K. Ehrenberg, J. D. Mosca, J. Rappaport,**
390 **and R. R. Redfield.** 1994. Negative-strand RNA transcripts are produced in human
391 immunodeficiency virus type 1-infected cells and patients by a novel promoter
392 downregulated by Tat. *Journal of virology* **68**:979-987.
- 393 35. **Miller, R. H.** 1988. Human immunodeficiency virus may encode a novel protein on the
394 genomic DNA plus strand. *Science (New York, N.Y)* **239**:1420-1422.

- 395 36. **Mitchell, R., C. Y. Chiang, C. Berry, and F. Bushman.** 2003. Global analysis of
396 cellular transcription following infection with an HIV-based vector. *Mol Ther* **8**:674-687.
- 397 37. **Montano, M., M. Rarick, P. Sebastiani, P. Brinkmann, M. Russell, A. Navis, C.**
398 **Wester, I. Thior, and M. Essex.** 2006. Gene-expression profiling of HIV-1 infection and
399 perinatal transmission in Botswana. *Genes Immun* **7**:298-309.
- 400 38. **Morrissy, A. S., R. D. Morin, A. Delaney, T. Zeng, H. McDonald, S. Jones, Y. Zhao,**
401 **M. Hirst, and M. A. Marra.** 2009. Next-generation tag sequencing for cancer gene
402 expression profiling. *Genome Res* **19**:1825-1835.
- 403 39. **Naldini, L., U. Blomer, P. Gally, D. Ory, R. Mulligan, F. H. Gage, I. M. Verma, and**
404 **D. Trono.** 1996. In vivo gene delivery and stable transduction of nondividing cells by a
405 lentiviral vector. *Science (New York, N.Y)* **272**:263-267.
- 406 40. **Pruitt, K. D., T. Tatusova, and D. R. Maglott.** 2007. NCBI reference sequences
407 (RefSeq): a curated non-redundant sequence database of genomes, transcripts and
408 proteins. *Nucleic acids research* **35**:D61-65.
- 409 41. **Rotger, M., J. Dalmau, A. Rauch, P. McLaren, S. Bosinger, R. Martinez, N. G.**
410 **Sandler, A. Roque, J. Liebner, M. Battegay, E. Bernasconi, P. Descombes, I. Erkizia,**
411 **J. Fellay, B. Hirschel, J. M. Miró, E. Palou, M. Hoffmann, M. Massanella, J. Blanco,**
412 **M. Woods, H. F. Günthard, P. de Bakker, D. C. Douek, G. Silvestri, J. Martinez-**
413 **Picado, and A. Telenti.** 2011. Comparative transcriptomics of extreme phenotypes of
414 human HIV-1 infection and SIV infection in sooty mangabey and rhesus macaque. *The*
415 *Journal of clinical investigation* **in press**.
- 416 42. **Rotger, M., K. K. Dang, J. Fellay, E. L. Heizen, S. Feng, P. Descombes, K. V.**
417 **Shianna, D. Ge, H. F. Gunthard, D. B. Goldstein, and A. Telenti.** 2010. Genome-wide

- 418 mRNA expression correlates of viral control in CD4+ T-cells from HIV-1-infected
419 individuals. *PLoS pathogens* **6**:e1000781.
- 420 43. **Rozera, G., I. Abbate, A. Bruselles, C. Vlassi, G. D'Offizi, P. Narciso, G. Chillemi,**
421 **M. Prosperi, G. Ippolito, and M. R. Capobianchi.** 2009. Massively parallel
422 pyrosequencing highlights minority variants in the HIV-1 env quasispecies deriving from
423 lymphomonocyte sub-populations. *Retrovirology* **6**:15.
- 424 44. **Sambrook, J. F., E.F.; Maniatis, T.** 1989. Northern Hybridization, p. 7.39-37.52. *In* N.
425 N. Ford, C.; Ferguson, M. (ed.), *Molecular Cloning: A Laboratory Manual*, 2nd ed, vol. 1.
426 Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York.
- 427 45. **Shendure, J.** 2008. The beginning of the end for microarrays? *Nature methods* **5**:585-
428 587.
- 429 46. **Simen, B. B., J. F. Simons, K. H. Hullsiek, R. M. Novak, R. D. Macarthur, J. D.**
430 **Baxter, C. Huang, C. Lubeski, G. S. Turenchalk, M. S. Braverman, B. Desany, J. M.**
431 **Rothberg, M. Egholm, and M. J. Kozal.** 2009. Low-abundance drug-resistant viral
432 variants in chronically HIV-infected, antiretroviral treatment-naive patients significantly
433 impact treatment outcomes. *The Journal of infectious diseases* **199**:693-701.
- 434 47. **Sirois, M., L. Robitaille, R. Sasik, J. Estaquier, J. Fortin, and J. Corbeil.** 2008. R5
435 and X4 HIV viruses differentially modulate host gene expression in resting CD4+ T cells.
436 *AIDS research and human retroviruses* **24**:485-493.
- 437 48. **Solis, M., P. Wilkinson, R. Romieu, E. Hernandez, M. A. Wainberg, and J. Hiscott.**
438 2006. Gene expression profiling of the host response to HIV-1 B, C, or A/E infection in
439 monocyte-derived dendritic cells. *Virology* **352**:86-99.
- 440 49. **t Hoen, P. A., Y. Ariyurek, H. H. Thygesen, E. Vreugdenhil, R. H. Vossen, R. X. de**
441 **Menezes, J. M. Boer, G. J. van Ommen, and J. T. den Dunnen.** 2008. Deep

442 sequencing-based expression analysis shows major advances in robustness, resolution and
443 inter-lab portability over five microarray platforms. *Nucleic acids research* **36**:e141.

444 50. **Tsibris, A. M., B. Korber, R. Arnaout, C. Russ, C. C. Lo, T. Leitner, B. Gaschen, J.**
445 **Theiler, R. Paredes, Z. Su, M. D. Hughes, R. M. Gulick, W. Greaves, E. Coakley, C.**
446 **Flexner, C. Nusbaum, and D. R. Kuritzkes.** 2009. Quantitative deep sequencing reveals
447 dynamic HIV-1 escape and large population shifts during CCR5 antagonist therapy in
448 vivo. *PLoS ONE* **4**:e5683.

449 51. **Vahey, M. T., Z. Wang, Z. Su, M. E. Nau, A. Krambrink, D. J. Skiest, and D. M.**
450 **Margolis.** 2008. CD4+ T-cell decline after the interruption of antiretroviral therapy in
451 ACTG A5170 is predicted by differential expression of genes in the ras signaling
452 pathway. *AIDS research and human retroviruses* **24**:1047-1066.

453 52. **Wang, G. P., A. Ciuffi, J. Leipzig, C. C. Berry, and F. D. Bushman.** 2007. HIV
454 integration site selection: Analysis by massively parallel pyrosequencing reveals
455 association with epigenetic modifications. *Genome Res* **17**:1186-1194.

456 53. **Wu, Z. J., C. A. Meyer, S. Choudhury, M. Shipitsin, R. Maruyama, M. Bessarabova,**
457 **T. Nikolskaya, S. Sukumar, A. Schwartzman, J. S. Liu, K. Polyak, and X. S. Liu.**
458 2010. Gene expression profiling of human breast tissue samples using SAGE-Seq.
459 *Genome Res* **20**:1730-1739.

460 54. **Yeung, M. L., Y. Bennasser, K. Watashi, S. Y. Le, L. Houzet, and K. T. Jeang.** 2009.
461 Pyrosequencing of small non-coding RNAs in HIV-1 infected cells: evidence for the
462 processing of a viral-cellular double-stranded RNA hybrid. *Nucleic acids research*
463 **37**:6575-6586.

- 464 55. **Zagordi, O., R. Klein, M. Daumer, and N. Beerenwinkel.** 2010. Error correction of
465 next-generation sequencing data and reliable estimation of HIV quasispecies. *Nucleic*
466 *acids research*.
- 467 56. **Zhang, H., Y. Zhou, C. Alcock, T. Kiefer, D. Monie, J. Siliciano, Q. Li, P. Pham, J.**
468 **Cofrancesco, D. Persaud, and R. F. Siliciano.** 2004. Novel single-cell-level phenotypic
469 assay for residual drug susceptibility and reduced replication capacity of drug-resistant
470 human immunodeficiency virus type 1. *Journal of virology* **78**:1718-1729.
- 471

472

473 **Figures**

474

475 **Figure 1 - Overview of experimental and analytical procedures.**

476 Panel A, Experimental pipeline. Panel B, Bioinformatics pipeline.

477

478 **Figure 2 - Distribution of HIV-specific tags along the viral genome.**

479 (A) The HIV vector genome is depicted on the top panel with nucleotide positioning, NlaIII
480 restriction sites (vertical black bars), the common sense transcription start site (TSS; red arrow)
481 and the polyadenylation signal (poly(A); red vertical bar) are shown below. Multiple antisense
482 transcription start sites (distributed between positions 9175 and 8714; only two are drawn here
483 for figure simplicity; green arrows), and the poly(A) signal (vertical green bar), as described in
484 Landry *et al.* (26), are indicated. HIV-1 open reading frames are also indicated, including the
485 putative antisense protein (ASP). The bottom panel indicates the density of tags (vertical axis)
486 distributed along the HIV-1 vector genome (x-axis), for both HIV-1 samples (HIV rep1 and HIV
487 rep2) and mock samples (Mock rep1 and Mock rep2). The five major peaks are numbered at the
488 bottom of the figure, displaying the number of tags at each peak, as well as the proportion of
489 these tags representing sense or antisense transcription. (B) Total RNA from mock or HIV-based
490 vector infected cells was extracted using miRvana or Trizol at 22h and 24h post-exposure and
491 subjected to Northern blot hybridization. Left panel, total RNA gel electrophoresis. Right panels,
492 Northern blots using strand-specific oligonucleotide probes aligning to the corresponding peak of
493 tags detected by SAGE-Seq. Probes were designed to anneal with sense transcripts (red probes,
494 top panels) or antisense transcripts (green probes, bottom panels), with nucleic acid sequence

495 complementary to the identified peak tag sequence. The plasmid encoding HIV vector genome
496 digested with BglIII and AflIII was used as positive control for hybridization (+).

497

498 **Figure 3 - Analysis of cellular transcripts in mock and HIV-infected cells.**

499 Average tag density analysis in mock (black lines) and HIV-1 (red lines) samples for total
500 cellular transcripts (lines), repetitive elements (dashed lines) and HERV elements (dotted lines).
501 The figure indicates that most cellular transcripts were detected with 1 to 1000 tags (log 0 to 3,
502 average 2.9 log), across the complete transcriptome. Only a few transcripts were identified in
503 larger numbers (*e.g.* above 10000, log 4). Lesser number of tags was observed for the expressed
504 repetitive elements (average 2.3 log) including HERVs (average 1.9 log). Statistical differences
505 between mock and HIV-1 distributions were assessed by Wilcoxon test.

506

507 Tables

508 **Table 1 - Mapping to the human genome, human transcriptome, and HIV-1 vector genome.**

509

Sample	Human		HIV-1 vector	
	Genome tags ^a	Transcriptome tags ^b	HIV-specific tags ^c	GFP-specific tags ^d
Mock rep1	30,668,095	38,730,420	45	69
Mock rep2	31,611,573	43,691,287	2,334	2,862
HIV rep1	27,894,076	33,337,832	11,7391	13,5487
HIV rep2	29,484,387	36,666,667	10,9122	12,6344

510

511 ^aAlignment/mapping to GRCh37 allowed 1 mismatch (n=1), 1 hit (m=1), read length of 24 (rl=24).

512 ^bAlignment/mapping to RefSeq allowed 1 mismatch (n=1), 10 multiple hits (m=10), read length of 24 (rl=24).

513 ^cAlignment/mapping to *HIV NLA-3Δenv/eGFP* was performed allowing 2 mismatches (n=2) and 2 multiple hits
514 (m=2), read length of 24 (rl=24).

515

516

517

518 **Supplemental Material**

519 **Supplemental File 1 – Completion of viral cycle in SupT1 cells.**

520 (A) SupT1 cells were highly susceptible to HIV-based vector infection (compared to CEM and
521 Jurkat T cell lines). At 24h and 36h post-infection, 54% and 93% of cells were successfully
522 transduced. (B) To establish the optimal time point to capture the completion of the viral
523 replication cycle, as reflected by particle release (measured by p24 ELISA in the supernatant), we
524 performed a time course collecting materials every 2 hours. The plots represent here non-
525 cumulative estimates, *i.e. de novo* production, calculated by subtracting the measurement at a
526 given time point by the prior measurement. Shown are the percentage of GFP+ cells and mean
527 fluorescence intensity (MFI), representing the success of viral translation (green), which peaks at
528 22 hours, and followed by the peak of extracellular p24 (red) at 24 hours expected 2-4 after the
529 peak of translation.

530 File format: pdf

531 This file can be viewed with: Adobe Acrobat Reader

532

533 **Supplemental File 2 – Optimal mapping parameters for read analysis.**

534 Read length (rl), number of mismatches (n) and multiple hits allowed (m) were tested to
535 maximize the number of tags mapping to HIV-1 while minimizing the number of tags spuriously
536 mapping to both viral and host genome. Y-axis: reads mapping to HIV-1 vector genome; x-axis:

537 number of reads mapping to HIV-1 but not to human genome. Parameters rl, n, and m are
538 indicated only for the best conditions.

539 File format: pdf

540 This file can be viewed with: Adobe Acrobat Reader

541

542 **Supplemental File 3 – Transcripts differentially expressed between mock and HIV-1**
543 **infected cells.**

544 File format: xls

545 This file can be viewed with: Microsoft Excel Viewer

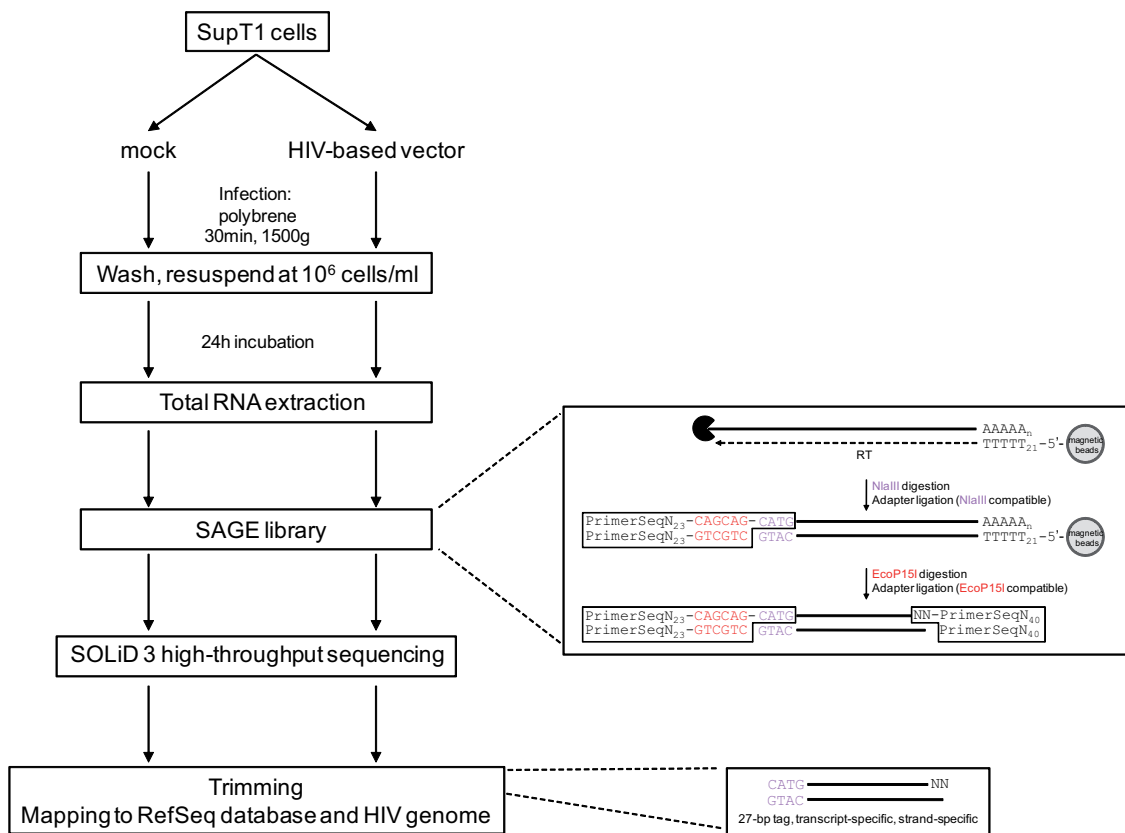
546

547 **Supplemental File 4 – Transcripts identified by SAGE-Seq common to previous studies**
548 **using microarray technology.**

549 File format: xls

550 This file can be viewed with: Microsoft Excel Viewer

551

A**B**