# Reference Plasmid pHXB2_D is an HIV-1 Molecular Clone that Exhibits Identical LTRs and a Single Integration Site Indicative of an HIV Provirus

Alejandro R. Gener ( ✉ itspronouncedhenner@gmail.com )

Baylor College of Medicine    https://orcid.org/0000-0001-7875-2358

**Wei Zou**

First Affiliated Hospital of Nanchang University

**Brian T. Foley**

Los Alamos National Laboratory

**Deborah P. Hyink**

Baylor College of Medicine

**Paul E. Klotman**

Baylor College of Medicine

**Research**

# Reference plasmid pHXB2_D is an HIV-1 molecular clone that exhibits identical LTRs and a single integration site indicative of an HIV provirus

Alejandro R. Gener[1,2,3,4][§], Wei Zou[5], Brian T. Foley[6], Deborah P. Hyink*[2], Paul E. Klotman*[1,2]

[1]Integrative Molecular and Biomedical Sciences Program, Baylor College of Medicine, Houston, Texas, USA
[2]Margaret M. and Albert B. Alkek Department of Medicine, Nephrology, Baylor College of Medicine, Houston, Texas, USA
[3]Department of Genetics, MD Anderson Cancer Center, Houston, Texas, USA
[4]School of Medicine, Universidad Central del Caribe, Bayamón, Puerto Rico, USA
[5]Division of Infectious Diseases, the 1st Affiliated Hospital of Nanchang University, Nanchang, Jiangxi, China
[6]Theoretical Biology and Biophysics Group T-6, Los Alamos National Laboratory, Los Alamos, New Mexico, USA

*Equal contributions.

[§]Corresponding author: Alejandro R. Gener
One Baylor Plaza
Mail Stop 710
Houston, Texas, 77030, USA
9045715562
gener@bcm.edu ; itspronouncedhenner@gmail.com

**Keywords:** HIV-1, reagent verification, nanopore DNA sequencing, provirus, plasmid, sequence variability, resequencing, LTR phasing

1  **Abstract**

2  **Objective:** To compare long-read nanopore DNA sequencing (DNA-seq) with short-read

3  sequencing-by-synthesis for sequencing a full-length (e.g., non-deletion, nor reporter) HIV-1

4  model provirus in plasmid pHXB2_D.

5  **Design:** We sequenced pHXB2_D and a control plasmid pNL4-3_gag-pol(Δ1443-4553)_EGFP

6  with long- and short-read DNA-seq, evaluating sample variability with resequencing (sequencing

7  and mapping to reference HXB2) and *de novo* viral genome assembly.

8  **Methods:** We prepared pHXB2_D and pNL4-3_gag-pol(Δ1443-4553)_EGFP for long-read

9  nanopore DNA-seq, varying DNA polymerases Taq (Sigma-Aldrich) and Long Amplicon (LA)

10  Taq (Takara). Nanopore basecallers were compared. After aligning reads to the reference HXB2

11  to evaluate sample coverage, we looked for variants. We next assembled reads into contigs,

12  followed by finishing and polishing. We hired an external core to sequence-verify pHXB2_D

13  and pNL4-3_gag-pol(Δ1443-4553)_EGFP with single-end 150 base-long Illumina reads, after

14  masking sample identity.

15  **Results:** We achieved full-coverage (100%) of HXB2 HIV-1 from 5' to 3' long terminal repeats

16  (LTRs), with median per-base coverage of over 9000x in one experiment on a single MinION

17  flow cell. The longest HIV-spanning read to-date was generated, at a length of 11,487 bases,

18  which included full-length HIV-1 and plasmid backbone with flanking host sequences

19  supporting a single HXB2 integration event. We discovered 20 single nucleotide variants in

20  pHXB2_D compared to reference, verified by short-read DNA sequencing. There were no

21  variants detected in the HIV-1 segments of pNL4-3_gag-pol(Δ1443-4553)_EGFP.

22  **Conclusions:** Nanopore sequencing performed as-expected, phasing LTRs, and even covering

23  full-length HIV. The discovery of variants in a reference plasmid demonstrates the need for

24     sequence verification moving forward, in line with calls from funding agencies for reagent

25     verification. These results illustrate the utility of long-read DNA-seq to advance the study of

26     HIV at single integration site resolution.


27     **Introduction**

28         Much of what we know about human acquired immunodeficiency syndrome (AIDS)

29     came after isolating the causative agent – the  human immunodeficiency virus type 1 (HIV-1) –

30     and describing the viral genome information content. The HIV-1 isolate HXB2 (also known as

31     HTLV-III and HIV-1$_{LAI}$ or LAV/BRU [1], [2]) was the first full-length replication-competent

32     HIV genome sequenced [3]. Derivative clones commonly called "HXB2" are still used for *in*

33     *vitro* infection assays, including RNA (almost always cDNA [4]) sequencing (**Figure 1A** and

34     **Supplemental Table 1**). Despite the availability of the HXB2 HIV-1 reference sequence [3], no

35     sequence is available for any complete and readily available HXB2 clone.

36         HIV clones were originally made by choosing non-cutter restriction enzymes to digest

37     intact proviral sequences upstream and downstream of unknown integration sites from infected

38     host cells while sparing HIV-1 sequence, followed by ligation into an *E. coli* cloning vector

39     (plasmid) (**Figure 1B**), allowing for low-error (but not error-free) propagation [5]. These clones

40     became available before tractable sequencing methods permitted routine sequence verification.

41     As such, it was uncommon to sequence them. While funding agencies now require investigators

42     to include in their proposals plans to validate their key reagents, these funders tend to leave the

43     process up to investigators and may not always follow up on whether a given reagent is ever

44     actually validated (or revalidated between changes of hand). Investigators do not regularly

45     validate their clones, in part because there is no universally accepted standard. Instead, a

46	common practice is to assume a given clone, often kindly gifted from a colleague, is as reported.

47	As such, we often do not truly know what we have been working with for 35+ years.

48	       Making sense of the information from HIV sequencing experiments is complicated by

49	many factors, including the cycling that all orterviruses [6] undergo between two major states (as

50	infectious virion RNA and integrated proviral DNA **Figure 1B**), repetitive viral sequences like

51	long terminal repeats (LTRs), non-integrated forms [7], rarity of integration events *in vivo*

52	(reviewed in [8]), and alternative splicing of viral mRNAs [9]. Short-read DNA sequencing

53	(<150 bases (bp) long in most reported experiments, but up to 500 bp for either Illumina (ILMN)

54	sequencing-by-synthesis or <1,000 bp for chain termination sequencing) provides some

55	information, but analyses require high coverage and/or extensive effort (non-exhaustive

56	examples [10], [11]). These factors limit the ability to assign variants to specific loci within each

57	provirus, as well as at the proviral integration site(s) (reviewed in [12]). Despite progress (HIV

58	DNA) [13], (HIV RNA) [14], [15], [16], researchers have yet to observe the genome of HIV-1 as

59	complete provirus (integrated DNA) in a single read, hindering locus-specific studies. To this

60	end, current long-read DNA sequencing clearly surpasses the limitations of read length of

61	leading next-generation/short-read sequencing platforms. Here we used the MinION sequencer

62	to sequence HIV-1 plasmid pHXB2_D in a pilot study focusing on coverage acquisition (as

63	opposed to full-length sequencing), with the goal of evaluating the technology for future

64	applications.

## Methods

**HIV-1 plasmids**

A plasmid, "pHXB2_D" (alternate names pHXB2, pHXB-2D), believed to contain the HIV-1 reference strain HXB2 [17] was acquired from the NIH AIDS Reagent and Reference Program (ARP) via BioServe. pHXB2_D was believed to be a molecular clone (likely a restriction product of HXB2 proviral DNA inserted into an unknown cloning plasmid backbone) from one of the earliest clinical "HXB2" HIV-1 isolates. At the time of this work, it was unknown whether this plasmid was ever sequence-verified before or after the reference sequence for HXB2 was deposited.

The provenance of pNL4-3_gag-pol($\Delta$1443-4553)_EGFP, a reporter construct of pNL4-3 with a gag-pol deletion between base 1443 and 4553 is known. HIV-1 NL4-3 (pNL4-3) was a fusion of NY5 and LAV/HXB2 plasmids [18] that to our knowledge are not readily available. pEVd1443 [19] was a deletion construct made from pNL4-3 used to make several HIV-1 transgenic animals, including the FVB/N-Tg(HIV)26Aln/PkltJ (The Jackson Laboratory stock No: 022354) "Tg26" mouse. The deletion in pEVd1443 was made by SphI cutting between d1443 and 1444 with binding site 1443-1448, and cutting at a BalI site at 4551-4556 with blunt cutting between 4553 and 4554. The EGFP cassette includes additional sequence upstream and downstream of EGFP coding sequence. SphI and BalI may still be used to excise EGFP cassette. A reporter construct was designed mimicking the pEVd1443 deletion: pNL4-3: $\Delta$G/P-EGFP [20]. Dr. Wei Zou rederived pNL4-3: $\Delta$G/P-EGFP at BCM [21]. Both constructs (plasmid and mouse) retained parts of gag and pol, with limited effects on protein-coding capacity, such as expression of p17 [22]. Based on Addgene naming conventions, we suggest pNL4-3_gag-pol($\Delta$1443-4553)_EGFP to replace the previous name pNL4-3: $\Delta$G/P-EGFP for clarity.

**88**     **HIV-1 reference sequences**

**89**         The reference sequence of HXB2 is from the National Center for Biotechnology

**90**     Information (NCBI), Genbank accession number K03455.1. It runs from the beginning of the 5'

**91**     LTR to the end of the 3' LTR, and is 9,719 bp. This is similar to another HIV-1 reference that

**92**     NCBI uses, AF033819.3. This is a 9,181 base HXB2-like sequence that starts at the 97 bp repeat

**93**     in the 5'LTR, continues with the 5'UTR (U5), extends past the 3'UTR (U3) to the end of the 97

**94**     bp repeat in 3'LTR, with one SNV at the *vpu* start codon aTg to aCg at position AF033819.3:560

**95**     or K03455.1:6063. The reference sequence of NL4-3 is as a plasmid with accession number

**96**     AF324493.1. It runs from the beginning of the 5' LTR to the end of the 3' LTR, spanning 9,709

**97**     bp, and includes plasmid backbone with total length 14,825 bp.

**98**     **Long-read DNA sequencing**

**99**         A plasmid containing HXB2 was sequence-verified with long-read nanopore sequencing

**100**     on a MinION Mk1B (Oxford Nanopore Technologies (ONT), Oxford, United Kingdom). Unless

**101**     otherwise noted, reagents (and software) were purchased (or acquired) from ONT. Briefly, stock

**102**     plasmid was diluted to 5 ng final amount in ultrapure water (as two samples) and processed with

**103**     Rapid PCR Barcoding kit SQK-RPB004 along with 10 other barcoded samples (not discussed

**104**     further in this manuscript) following ONT protocol RPB_9059_V1_REVA_08MAR2018

**105**     (**Figure 1C**), a public description of which is here: https://store.nanoporetech.com/us/sample-

**106**     prep/rapid-pcr-barcoding-kit.html. Two DNA polymerases were evaluated (barcode 10 used

**107**     high-fidelity LA (for "long amplicon") Taq (Takara); barcode 11 Taq (Sigma-Aldrich). Libraries

**108**     were loaded onto a MinION flow cell version R9.4.1 and a 48-hour sequencing run was

**109**     completed with MinKNOW (version 1.10.11). Residual reads from subsequent runs were pooled

110    for final analyses. Long read data for pNL4-3_gag-pol(Δ1443-4553)_EGFP was generated in

111    other barcoded experiments (not shown).

112         Raw data was basecalled (converted from FAST5 to FASTQ format) with Albacore

113    version 2.3.4 (older basecaller), Guppy version 2.3.1 (current official at time of work), and

114    FlipFlop (Guppy development config). Mapping to reference was done with Minimap2 [23] and

115    BWA-MEM [24], implemented in Galaxy (usegalaxy.org) [25]. Alignments (.bam and .bai files)

116    were visualized in the Integrative Genomics Viewer [26] unless otherwise noted. For *de novo*

117    assembly, demultiplexed basecalled reads were fed into Canu version 1.8 [27]. Genome size was

118    estimated to be 16 Kb from agarose gel of undigested, but naturally degraded linearized

119    pHXB2_D (data not shown). SnapGene version 4.3.4 was used to manually annotate contigs

120    from Canu. Blastn (NCBI) was used to identify unknown regions of pHXB2_D. Polishing was

121    performed on ONT-only assemblies with Medaka (https://github.com/nanoporetech/medaka), in

122    Galaxy. Medaka models: r941_min_fast_g303, r941_min_high_g303, r941_min_high_g330.

123    Inference batch size (-b) = 100. The final pHXB2_D assembly and other full-length HIV clones

124    from the ARP were aligned to the most recent human reference genome (hg38) with Minimap2

125    in Galaxy with the following parameters: Long assembly to reference mapping (-k19 -w19 -A1 -

126    B19 -O39,81 -E3,1 -s200 -z200 --min-occ-floor=100).

127    **Statistics**

128         Two-tailed Mann-Whitney U tests were used to compare distributions in long-read data.

129    P-values are reported over brackets delineating relevant comparisons. Calculations and graphing

130    were done with GraphPad Prism for macOS version 8.0.2.

**Short-read DNA sequencing**

131     pHXB2_D and control pNL4-3_gag-pol(Δ1443-4553)_EGFP were provided as 35 ul at

133     ~63 ng/ul to the Center for Computational & Integrative Biology DNA Core at Massachusetts

134     General Hospital, an external DNA sequencing core specializing in high-throughput next

135     generation (short-read) plasmid sequencing and assembly. Neither HXB2/pNL4-3 reference

136     sequences nor pHXB2_D/pNL4-3_gag-pol(Δ1443-4553)_EGFP draft assemblies (from this

137     work) were provided to core staff at the time of sequencing so that testing would remain masked.

138     While the core's exact library prep is proprietary, multiplexed library prep and 150 single-end

139     ILMN sequencing were most likely performed on a MiSeq with platform-specific reagents (V2

140     chemistry, per their website) and barcoding. Data was returned as FASTQ. FASTQC [28] was

141     used in Galaxy for in-house data quality control, and read lengths were all 142 bp per this tool.

142     Mapping as above.

**Sequence comparisons**

144     We used MAFFT v7.475 [29], [30] to compare the LTR sequences of pHXB2_D and

145     HXB2, and pNL4-3 and pNL4-3_gag-pol(Δ1443-4553)_EGFP. For cladistics, we used BLAST

146     at HIV-DB (https://www.hiv.lanl.gov/content/sequence/BASIC_BLAST/basic_blast.html) to

147     find other HXB2-like genomes. The top 50 BLAST hits included many sequences pNL43 clones.

148     pNL4-3 is an artificial recombinant of the NY5 clone with LAV and/or the HXB2 clone [18].

149     The recombination point is marked by an EcoRI restriction site. We then made a multi-sequence

150     alignment with the final pHXB2_D assembly, the top BLAST hits, and the HIV-1 M group

151     subtype reference set using GeneCutter

152     (https://www.hiv.lanl.gov/content/sequence/GENE_CUTTER/cutter.html), and built the

153     maximum likelihood tree using IQ-tree

154 (https://www.hiv.lanl.gov/content/sequence/IQTREE/iqtree.html). pNL4-3_gag-pol(Δ1443-

155 4553)_EGFP was not included in the above trees because of absence of divergence from pNL4-3

156 sequences outside of the EGFP cassette.


## Results

157 **Results**

158  Viewing mapped data in IGV, the long reads (median read length >2000 bp, **Figure 1E**)

159 from both pHXB2_D ONT experiments clearly covered each LTR (**Figure 1F, Supplemental**

160 **Figures 1**, **3**), while shorter reads collapsed into one of either LTR (**Figure 1F, Supplemental**

161 **Figures 3D,3E**). This was also seen when long reads were shorter than LTRs (<600 bp).

162 Mappers BWA-MEM and Minimap2 were chosen based on their ability to handle long and short

163 reads. Other mappers were not evaluated. BWA-MEM mapped more ambiguously, piling

164 partially mapped reads between each LTR; Minimap2 mapped with higher fidelity to reference

165 without splitting reads. Coverage as sequencing depth was higher and more even from the

166 higher-fidelity LA Taq library (**Supplemental Figure 1**). pNL4-3 was known to have distinct

167 LTRs because it was a synthetic recombinant. The higher variant density in NL4-3 LTRs enabled

168 mapping and phasing from short-read data only (**Supplemental Figure 2**).

169  We counted 20 single nucleotide variants (SNVs) in this reference clone of HXB2 (**Table**

170 **1, Supplemental Table 3, Supplemental Figure 3E**). These mismatches were seen in all Canu

171 assemblies (**Supplemental Figures 4A,4B**), verified in IGV and/or SnapGene, and were

172 orthogonally verified by short-read sequencing performed by the external core given masked

173 samples (**Supplemental Figure 3E**). These mismatches represent a ~0.21% divergence from

174 reference HXB2 K03455.1 (20/9719), which was assumed to have perfect identity (0%

175 divergence). Transitions were more common (14/20) (**Table 1**), coinciding with a previous

176 report of increased transitions over transversions in infection models, because transversions are

177 more likely to be deleterious to viral replication (i.e.: to cause protein-coding changes) [31].

178 Indeed, almost half (9/20) of the observed SNVs occurred in protein-coding regions, even though

179 92% of HXB2 is coding (791/9719). Of those 9 SNVs in protein-coding regions, 4 caused non-

180 synonymous mutations. One of those occurs in a region overlapping both gag and pol regions,

181 however only pol exhibited a non-synonymous change from valine to isoleucine in p6, at

182 position 2259 relative to HXB2. Other non-synonymous variants occurred at 4609 (in p31

183 integrase, arginine to lysine), 7823 (in ASP antisense protein, glycine to arginine), and 9253 (in

184 nef, isoleucine to valine). 11/20 SNVs were in LTRs (see **Supplemental Figure 3** for counting

185 based on mapping); 8/20 of these would have been missed with mapping-only variant calling or

186 consensus. The longest HIV-mapping read (**Figure 2**) phased 16/20 SNVs (failed at sites

187 2,8,10,12, **Table 1**). pNL4-3_gag-pol($\Delta$1443-4553)_EGFP did not have HIV-1 or plasmid

188 backbone variants supported by long and short reads outside of the EGFP cassette.

189       We assembled the previously undefined plasmid pHXB2_D (**Supplemental Figures**

190 **4A,4B**). Canu's final output was a set of contiguous DNA sequences (contigs) as FASTA files. A

191 consequence of assembling plasmid sequences with this tool was partial redundancy at contig

192 ends (**Supplemental Figure 4C**). Manual end-trimming of contigs was performed in SnapGene

193 based on an estimated length of 16 kb. Top blastn hits from barcode 10/LA Taq pHXB2

194 basecalled with FlipFlop were as follows: for the main backbone (with origin of replication and

195 antibiotic selection cassette for cloning), shuttle vector pTB101-CM DNA, complete sequence

196 (based on pBR322), from 4352-8340; for the upstream element (relative to 5' LTR), Homo

197 sapiens chromosome 3 clone RP11-83E7 map 3p, complete sequence from 58,052 to 59,165; for

198 the downstream element, cloning vector pNHG-CapNM from 10,204 to 11,666. Other identified

199  elements included Enterobacteria phage SP6 (the SP6 promoter, per SnapGene's "Detect

200  common features"), complete sequence from 39,683 to 39,966. Identities of query to HXB2 and

201  hits were all approximately 99%. The MGH CCIB DNA Core's proprietary *de novo* UltraCycler

202  v1.0 assembler (Brian Seed and Huajun Wang, unpublished) was able to assemble both 5' and 3'

203  LTRs with short-read data only but may have collapsed SNVs into an artificial single consensus.

204  Long-read mapping and assembly (and polished assemblies) orthogonally validated LTRs and

205  supported a single HIV-1 HXB2_D haplotype (**Supplemental Figure 4,6**). A final LTR-phased

206  and annotated assembly leveraging short and long reads is provided as pHXB2_D

207  Genbank:MW079479 (embargoed until publication). Importantly, for pHXB2_D, each LTR was

208  identical, which is distinct from the current HXB2 (K03455.1) (**Figure 3A**). Compared to pNL4-

209  3_gag-pol(Δ1443-4553)_EGFP (ACCESSION_TBD) , each LTR was distinct, but identical to

210  pNL4-3's distinct 5' and 3' LTRs (AF324493.1) (**Figures 3B,6**).

211      To determine whether pHXB2_D was an isolated provirus (as opposed to a cDNA clone),

212  the pHXB2_D assembly was aligned to the current human reference hg38, returning a single

213  complete insertion site on 3p24.3 (**Figure 4A, Supplemental Table 2**). As expected, our pNL4-

214  3_gag-pol(Δ1443-4553)_EGFP had homology arms from two chromosomes (**Figures 4B,6,**

215  **Supplemental Table 2**). We sought to put our pHXB2_D assembly into context of other HXB2-

216  like references available (**Figure 5**). pHXB2_D (red) clusters closely with HXB2 reference

217  (K03455) and related clone sequences (green). pNL4-3 clones in blue. The LTR-masked HIV-

218  spanning segment of pHXB2_D is most homologous to B.FR.1983.DM461230 and

219  B.FR.1983.CS793683, which are identical except for areas in nef and a GFP insertion (verified

220  by blastn). This finding suggests they were from the same stock. HIV-1 M group subtype

221  reference set (HIV Sequence Database) was added to put HXB2s and pNL4-3 clones into

222    perspective. HXB2 (believed to be a complete isolate) and NL4-3 (synthetic clone based on two

223    early isolates [18]) are examples of HIV type 1 (HIV-1), group M, subgroup B.

224         As previously reported [32], per-read variability in ONT data was higher near

225    homopolymers (runs of the same base) (**Supplemental Figure 5A**). For the datasets generated in

226    the present study, homopolymers were counted and classified as continuous (unbroken run of a

227    given nucleobase) vs. discontinuous (broken run of a given nucleobase) (**Supplemental Figures**

228    **5B,5D,5F,5H**). A/T (2 hydrogen bonds; 2H) and G/C (3 hydrogen bonds; 3H) were evaluated.

229    Because runs longer than 4 or 5 were rare in these datasets, it was impossible to evaluate longer

230    homopolymers. A simple calculation $Abs(\Delta)=Abs(\#homopolymers_{reference}$ -

231    $\#homopolymers_{assembly})$ helped to evaluate the performance of basecallers, such that better

232    basecallers had smaller $Abs(\Delta)$ (**Supplemental Figures 5C,5E,5G,5I,5K**). At the level of

233    consensus (made from sequences mapped to reference HXB2), homopolymers contributed few,

234    if any, obvious errors. A special case of homopolymer, dimer runs, was noted to cause persistent

235    errors regardless of ONT basecaller (**Supplemental Figures 5J,5K**). While dips occurred at

236    certain points near homopolymers, the consensus did not change much at the sequencing depth

237    used in this study for either barcoded pHXB2_D samples (**Supplemental Figures 1,3,4**).

238    Another interpretation is that homopolymers tend to seem truncated with ONT, with more reads

239    in support of shorter homopolymers. Canu assemblies showed basecaller-dependent variability

240    (**Supplemental Table 3**). That said, newer basecallers tended to produce fewer and smaller per-

241    read truncations. Assemblies without polishing did not correct all homopolymer truncations

242    (**Supplemental Figure 4A**). Polishing assemblies tended to correct these toward the final

243    pHXB2_D assembly (**Supplemental Figures 4B,6**). Data from polished ONT-only assemblies

244    and short-read sequencing do not support the truncations (gaps relative to reference) suggested

245     by unpolished ONT-only assemblies, representing a known current limitation of ONT. These are

246     not the same as the 20 SNVs supported by BOTH long- and short-read sequencing performed in

247     this study. The ratio of per-read deletions to per-read insertions (DEL/INS) was much higher for

248     SNVs occurring at homopolymers and near the same base, and this difference was maintained

249     between all basecallers used (**Supplemental Figure 5L**). These changes created more

250     problematic (longer) homopolymers.


251     **Discussion**

252        This work represents the first instance of complete and unambiguous sequencing of HIV-

253     1 provirus as plasmid and contributed to the identification of single nucleotide variants which

254     may not have been easily determined using other sequencing modalities, illustrating the

255     importance of validating molecular reagents in their entirety, and with complementary

256     approaches. Nanopore sequencing surpassed the read length limitations of traditional sequencing

257     modalities used for HIV such as Sanger sequencing and sequencing-by-synthesis by at least two

258     orders of magnitude. Other long-read DNA sequencing technologies such as PacBio's zero-mode

259     waveguide DNA sequencing were not evaluated in this work, but in principle would be

260     interchangeable for nanopore sequencing. Paired-end sequencing (as either DNA-seq or RNA-

261     seq) was not evaluated in this work, but has shown promise phasing LTRs in our hands [33]–

262     [35].


263     **First complete pass over all HIV information in reference plasmid pHXB2_D**

264        HIV provirus is believed to occur naturally as one or a few copies of reverse-transcribed

265     DNA forms integrated into the host nuclear genome. Depending on where integration occurs,

266     local GC or AT content might cause problems for detecting integrants with PCR. HIV also has

267    conserved transitions from areas of higher GC content (~60%) to content approximating average

268    human GC content (~40%). To limit PCR sequencing bias and to accommodate for the potential

269    heterogeneity of HIV sequences, we fractionated whole sample directly (as opposed to PCR-

270    barcoding select amplicons) with tagmentation provided in the Rapid PCR-Barcoding kit (ONT).

271    Tagmentation in this workd used transposon-mediated cleavage and ligation of barcode adapters

272    for later PCR amplification. A consequence of this fractionation was a distribution of reads

273    (**Figure 1E**) shorter than longer reads reported elsewhere for ONT experiments [36]. Based on

274    this distribution and the level of coverage, it was expected that HIV might be covered from end

275    to end, but this would have been exceptional. That said, an example is presented here (**Figure 2**).

276    The provirus status of pHXB2_D is supported by recovery of both upstream and downstream

277    homology arms which map to a single human integration site.

278    **Long reads enable LTR phasing and HIV haplotype definition**

279    We created 6 assemblies for pHXB2_D from ONT-only data (**Supplemental Figure 4**),

280    each with a common set of 20 SNVs (11 in LTRs), and final assemblies (a single HIV-1

281    HXB2_D haplotype; a single HIV-1 NL4-3_gag-pol(Δ1443-4553)_EGFP haplotype) leveraging

282    long- and short-read data. The external core's *de novo* assembly pipeline identified the same 20

283    SNVs, and variants in the LTRs were supported by ONT unambiguously. That the core's

284    assembler was able to phase LTR variants in these samples may have been because the samples

285    had high amounts of the same upstream and downstream sequences because of coming from one

286    plasmid. The core's assembler thus may have had additional sequencing information at the edges

287    of HXB2, helping it to map deeper into each LTR. This approach would likely fail in samples

288    with multiple integrations (as in various animal models of HIV disease [37]), which have

289    unknown upstream and downstream sequences, or in samples from natural human infection,

290   which is well known to exhibit multiple pseudo-random integration sites between cells [38],

291   [39], but with mostly single integration events per cell [8]. Inverse PCR (iPCR) is an alternative

292   method [40] with its own issues (e.g., PCR biases, HIV concatemers, host repeats). While current

293   PCR reagents have extended the range of what can be seen with iPCR, current approaches are

294   likewise limited by long DNA extraction methods, sample amount, and remain to be optimized.

295   If coverage is sufficient (≥10 reads in non-homopolymers and non-dimer runs), long-read

296   sequencing can provide linked variant information to individual integration sites. Identical 5' and

297   '3 LTRs (**Figure 3**) in the context of a single integration event (**Figure 4A**) support this integrant

298   being a *bona fide* provirus [41]. Other proviruses also had identical LTR pairs (**Supplemental**

299   **Table 2**). Technical limitations such as PCR errors before earlier sequencing may explain the

300   variability in the HXB2 reference LTRs. These were sequenced at a time before paired-end 150

301   or long-read DNA-seq were available to phase LTRs, raising the possibility that these LTRs

302   were incorrectly annotated by depositors assuming identity and copy-and-pasting the sequence of

303   one LTR for both without being able to unambiguously resolve each LTR.

304   **Mutations in a reference HIV-1 plasmid illustrate the need for reagent verification**

305          Up until 2020, HIV had been the most studied human pathogen, but HIV reagents are not

306   routinely re(verified). The pHXB2_D sequenced was allegedly a reference plasmid, with

307   unknown divergence between the published reference HXB2. Three independent experiments

308   (two long-read with PCR-barcoded libraries made with regular and long-amplicon Taq master

309   mixes, one short-read) yielded at least 20 single nucleotide variants in pHXB2_D which differed

310   from the HXB2 reference sequence (**Table 1, Supplemental Figure 3**), which were also

311   concordant across the three basecallers used (**Supplemental Table 3**) and are therefore not PCR

312   errors. By leveraging long reads with the MinION, we were able to find mutations in highly

313   repetitive LTRs relative to HXB2 Genbank:K03455.1 which are often assumed (but until now

314   never proven) to be identical (**Table 1**, **Figure 1**, **Supplemental Figures 1, 3E**), as well as

315   mutations in protein-coding regions (**Table 1**). We were also able to confirm that the backbone

316   of this plasmid is from pSP62 [17], a pBR322 derivative with the SP6 promoter [42], aiding in

317   the continued use of this important reagent, and illustrating the need of full-length reagent

318   validation moving forward. We suggest that all clinical reagents (e.g., vectors) be sequence-

319   verified at the level of single-molecule sequencing as standard quality control to protect against

320   sample heterogeneity.

321   **Improvement in ONT basecallers over time**

322        Albacore, Guppy, and FlipFlop basecallers were compared. Each produced reads of

323   similar length distributions (relative to polymerase used), while Guppy and FlipFlop produced

324   improved and best performance relative to quality score distributions (**Figure 1D**). Interestingly,

325   while read length distributions were affected by fidelity of polymerases evaluated in this work,

326   mean quality distributions were not. This is important because of the differences in cost between

327   higher fidelity Taq and classic Taq enzymes. That said, higher fidelity LA Taq produced much

328   higher coverage compared to Taq (**Supplemental Figure 1**). In consideration of library prep,

329   choice of enzyme used should be based on the desired read-length distribution and coverage.

330   Regarding read mapping, the increase in mean quality score between these basecallers improved

331   overall mapping, in part by facilitating demultiplexing, resulting in approximately ~10%

332   increases number of reads in barcoded libraries before mapping (shift in reads from unclassified

333   to a given barcode). FlipFlop tended to handle homopolymers better than previous basecallers

334   (**Supplemental Figures 5,6**). Homopolymers in HXB2 tended to exhibit apparent deletions near

335   5' ends of homopolymers (upstream due to technical artifact from mapping), but because

336    consensus is conserved (example, at least 80% of base in called read set is identical to reference),

337    and because short-read data lacks INDELS at these sites, it is unlikely that any of these

338    homopolymer deletions are real in these experiments. Dimer runs – stretches of repeating 2-mers

339    (pronounced "two-mers") – proved challenging regardless of basecaller. Mapping as above may

340    be used to aid in manually calling these when they occur. Albacore is currently deprecated, and

341    current versions of Guppy now incorporate a version of FlipFlop called Guppy High-ACcuracy

342    (HAC). Guppy HAC and subsequence versions were not evaluated in this work. Polishing is

343    becoming standard practice for processing assemblies from ONT data because it redresses most

344    homopolymer errors propagated into long-read-only assemblies. The best manually finished and

345    polished contig had 1 error out of 16,722 bases, illustrating the utility of ONT hardware when

346    paired with burgeoning software.


347    **Conclusions**

348         HIV informatics, the study of HIV sequence information, has been limited by the

349    common assumption that sequence fidelity exists between reference genomes available in

350    sequence databases and similarly named HIV clones. Modern DNA sequencing methods, such as

351    long- and short-read sequencing, are available to redress this issue. Long-read sequencing fills in

352    gaps left behind by short-read interrogation of HIV-1. Current limitations of the approaches used

353    in the present work to study HIV are 1.) the cost of long-read sequencing, regardless of platform,

354    compared to the cheaper short reads from sequencing-by-synthesis, 2.) long DNA extraction

355    methods in diseased tissue (Gener, unpublished), and 3.) the lower per-base accuracy (low-mid

356    90's with ONT vs. 98-99% with ILMN or newer PacBio HiFi), including difficulty near

357    homopolymers and dimer runs (**Supplemental Figure 5**). A nontrivial but redressable limitation

358    is availability of personnel trained to prepare sequencing libraries, to run sequencing, and to

359    analyze results. As the price of long-read sequencing decreases, hardware and software used in

360    basecalling and library protocols improve, and with the advent of more user-friendly tools, the

361    cost of obtaining usable data from long reads will become negligible compared to the ability to

362    answer historically intractable questions. This work raises the possibility of being able to detect

363    at least some recombination events, in a reference-free manner requiring only the comparison of

364    LTRs from the same integrants (**Figure 6**). We suggest that pHXB2_D and pNL4-3 constructs

365    may be used as negative and positive controls for the development of such screens. While other

366    HIV reference proviral clones were reported to have identical LTR pairs, this remains to be

367    tested in other clones, since other clones were generated with shorter sequencing methods. For

368    example, pNL4-3_gag-pol($\Delta$1443-4553)_EGFP had distinct LTRs as a plasmid. However, if an

369    NL4-3 virus is made from pNL4-3, the LTR sequences would homogenize to pNL4-3's 3' LTR

370    sequence. Future work will include optimizing DNA extraction protocols with the goal of

371    capturing higher-coverage fuller glimpses of each HIV proviral integration site in *in vivo* HIV

372    models and patient samples. This work has broad implications for all cells infected by both

373    integrating and non-integrating viruses, and for the characterization of targeted regions in the

374    genome which may be recalcitrant to previous sequencing methods. Long-read sequencing is an

375    important emerging tool defining the post-scaffold genomic era, allowing for the characterization

376    of anatomical landmarks of hosts and pathogens at the genomic scale.

## Declarations

**List of Abbreviations**

HIV-1, human immunodeficiency virus.

PCR, polymerase chain reaction.

ONT, Oxford Nanopore Technologies; nanopore sequencing.

RNA-seq, usually refers to cDNA sequencing, unless otherwise specified (e.g., native

RNA-seq is not cDNA sequencing).

cDNA, reverse-transcribed "copy" DNA from single-stranded RNA templates. cDNA is

usually double-stranded DNA.

DNA-seq, DNA sequencing.

EGFP, enhanced green fluorescent protein.

LTR, long terminal repeat.

R, R repeat (a 97bp region in the LTR).

STEM, science, technology, engineering, math.

DNA, deoxyribonucleic acid.

RNA, ribonucleic acid.

LA, Long Amplicon.

SNV, single nucleotide variant.

AIDS, human acquired immunodeficiency syndrome.

mRNA, messenger RNA. These are usually 5' capped and polyadenylated in human

tissues.

BCM, Baylor College of Medicine.

IGV, Integrative Genomics Viewer.

400       bp, bases long.

401       Kb, kilobases long.

402       ASP, antisense protein.

403       MGH CCIB, Center for Computational & Integrative Biology DNA Core at

404       Massachusetts General Hospital.

405       2H, two hydrogen bonds, as in A-T, A-U.

406       3H, three hydrogen bonds, as in G-C.

407       iPCR, inverse PCR.

408       INDELS, insertions and/or deletions.

409       ILMN, Illumina short-read DNA sequencing-by-synthesis.

410       HiFi, High-fidelity.

411       Abs(), absolute value.

412       RRE, rev-response element.

413       LANL, Los Alamos National Laboratory.

414  **Ethics Approval**

415       This work did not include human or animal subjects. Nanopore libraries for this

416  work were prepared in their entirety by ARG in a Biosafety Level 2 laboratory on main campus

417  at Baylor College of Medicine (BCM). Nanopore sequencing was completed between April and

418  May of 2018 as two of several control experiments included in the Student Genomics pilot run

419  (**Supplemental Information**). Short-read sequencing was completed in April 2019.

420  **Consent to publish**

421       All authors give their consent to publish the current work, pending peer review and

422  acceptance.

**Availability of data and materials.**

The sequence for pHXB2_D will be made available after publication with Genbank

accession number MW079479 (embargoed until publication). We suggest this as a replacement

for K03455.1 since MW079479 will have phased LTRs. pNL4-3_gag-pol(Δ1443-4553)_EGFP

will be deposited in Genbank as well. Both plasmids will be made available to the community

through the AIDS Reagent Program. Sequencing data will be deposited into GEO. **Available**

**additional files:** Albacore basecalled barcode 10, Guppy basecalled barcode 10, FlipFlop

basecalled barcode 10, Albacore basecalled barcode 11, Guppy basecalled barcode 11, FlipFlop

basecalled barcode 11, Minimap2 and BWA-MEM alignments (.bam and .bai), Clipboards from

points of interest (verified SNVs; n=20), .dna files of contigs (n=6), MGH data (raw + contig),

Supplemental Tables, Supplemental Figures.

**Competing interests**

ARG received travel bursaries from Oxford Nanopore Technologies (ONT). The present

work was completed independently of ONT. Other authors declare no conflicts of interest.

**Funding**

This work was funded in part by institutional support from Baylor College of Medicine;

the Human Genome Sequencing Center at Baylor College of Medicine; private funding by Bob

Ostendorf, CEO of East Coast Oils, Inc., Jacksonville, Florida; ARG's own private funding,

including Student Genomics (manuscripts in prep). Compute resources from the Computational

and Integrative Biomedical Research Center at BCM ("sphere" cluster managed by Dr. Steven

Ludtke) and the Department of Molecular and Human Genetics at BCM ("taco" cluster managed

by Mr. Tanner Beck and Dr. Charles Lin) greatly facilitated the completion of this work. ARG

has also received the PFLAG of Jacksonville scholarship for multiple years.

**Authors' contributions**

ARG conceived of this project, performed experiments, analyzed results, and drafted the manuscript. WZ rederived pNL4-3_gag-pol(Δ1443-4553)_EGFP. All authors discussed data and edited the manuscript. ARG and PK provided funding.

**Acknowledgements**

As part of a summer bioinformatics internship in the Paul E. Klotman Laboratory at Baylor College of Medicine, Akash Naik supervised by ARG performed *in silico* mapping analyses/experiments, generated and/or aided in the synthesis of **Supplemental Figure 4**, and assisted in writing relevant portions, discussing, and editing this manuscript. During a second summer internship with American Physician Scientists Association Virtual Summer Research Program, the following students were supervised by ARG helped to create **Figure 1A** and **Supplemental Table 1**: Yini Liang, Kirk Niekamp, Maliha Jeba, Delmarie M. Rivera Rodríguez. Orthogonal sequence verification was performed as a service by staff at the Center for Computational & Integrative Biology DNA Core at Massachusetts General Hospital, Boston, MA, USA.

We would like to thank the staff at the DNA Core for their exceptional services, including expert analyses and rapid turnaround time. We would like to thank Drs. Steven Richards, Qingchang Meng and the staff of the Human Genome Sequencing Center Research (HGSC) and Development (R&D) team for their earlier support in nanopore adoption. We would like to thank the team at Oxford Nanopore Technologies for their timely improvements and continued R&D. I would also like to thank Ms. Taneasha Monique Washington (current) and former members of the Paul E. Klotman lab, Dr. Gokul C. Das and Alexander Batista. I would

470 **References**

471 [1]  F. Barré-Sinoussi *et al.*, "Isolation of a T-lymphotropic retrovirus from a patient at risk for

472   acquired immune  deficiency syndrome (AIDS).," *Science*, vol. 220, no. 4599, pp. 868–

473   871, May 1983, doi: 10.1126/science.6189183.

474 [2]  S. Wain-Hobson *et al.*, "LAV revisited: origins of the early HIV-1 isolates from Institut

475   Pasteur.," *Science*, vol. 252, no. 5008, pp. 961–965, May 1991, doi:

476   10.1126/science.2035026.

477 [3]  L. Ratner *et al.*, "Complete nucleotide sequence of the AIDS virus, HTLV-III.," *Nature*,

478   vol. 313, no. 6000, pp. 277–284, Jan. 1985, doi: 10.1038/313277a0.

479 [4]  A. R. Gener and J. T. Kimata, "Full-coverage native RNA sequencing of HIV-1 viruses,"

480   *bioRxiv*, p. 845610, Jan. 2019, doi: 10.1101/845610.

481 [5]  G. M. Shaw, B. H. Hahn, S. K. Arya, J. E. Groopman, R. C. Gallo, and F. Wong-Staal,

482   "Molecular characterization of human T-cell leukemia (lymphotropic) virus type III in the

483   acquired immune deficiency syndrome.," *Science*, vol. 226, no. 4679, pp. 1165–1171,

484   Dec. 1984, doi: 10.1126/science.6095449.

485 [6]  M. Krupovic *et al.*, "Ortervirales: New Virus Order Unifying Five Families of Reverse-

486   Transcribing Viruses," *J. Virol.*, vol. 92, no. 12, pp. e00515-18, May 2018, doi:

487   10.1128/JVI.00515-18.

488 [7]  E. H. Graf *et al.*, "Elite suppressors harbor low levels of integrated HIV DNA and high

489   levels of 2-LTR circular HIV DNA compared to HIV+ patients on and off HAART,"

490   *PLoS Pathog.*, vol. 7, no. 2, 2011, doi: 10.1371/journal.ppat.1001300.

491 [8]  Y. Ito *et al.*, "Number of infection events per cell during HIV-1 cell-free infection," *Sci.*

492   *Rep.*, vol. 7, no. 1, p. 6559, 2017, doi: 10.1038/s41598-017-03954-9.

493    [9]    I. Cuesta, A. Mari, A. Ocampo, C. Miralles, S. Pérez-castro, and M. M. Thomson,

494          "Sequence Analysis of In Vivo -Expressed HIV-1 Spliced RNAs Reveals the Usage of

495          New and Unusual Splice Sites by Viruses of Different Subtypes," pp. 1–24, 2016, doi:

496          10.1371/journal.pone.0158525.

497    [10]   C. Wymant *et al.*, "Easy and accurate reconstruction of whole HIV genomes from short-

498          read sequence data with shiver," *Virus Evol.*, vol. 4, no. 1, pp. 1–13, 2018, doi:

499          10.1093/ve/vey007.

500    [11]   K. M. Bruner *et al.*, "A quantitative approach for measuring the reservoir of latent HIV-1

501          proviruses," *Nature*, vol. 566, no. 7742, pp. 120–125, 2019, doi: 10.1038/s41586-019-

502          0898-8.

503    [12]   M. R. Pinzone and U. O'Doherty, "Measuring integrated HIV DNA ex vivo and in vitro

504          provides insights about how reservoirs are formed and maintained," *Retrovirology*, vol.

505          15, no. 1, pp. 1–12, 2018, doi: 10.1186/s12977-018-0396-3.

506    [13]   K. B. Einkauf *et al.*, "Intact HIV-1 proviruses accumulate at distinct chromosomal

507          positions during prolonged antiretroviral therapy Find the latest version : Intact HIV-1

508          proviruses accumulate at distinct chromosomal positions during prolonged antiretroviral

509          therapy," vol. 129, no. 3, pp. 988–998, 2019.

510    [14]   D. Bonsall *et al.*, "THAA0101 - HIV genotyping and phylogenetics in the HPTN 071

511          (PopART) study: Validation of a high-throughput sequencing assay for viral load

512          quantification, genotyping, resistance testing and high-resolution transmission

513          networking," in *22nd International AIDS Conference (AIDS2018)*, 2018, p. Oral Abstract.

514    [15]   A. N. Banin *et al.*, "Development of a Versatile, Near Full Genome Amplification and

515          Sequencing Approach for a Broad Variety of HIV-1 Group M Variants," *Viruses*, vol. 11,

516          no. 4, p. 317, Apr. 2019, doi: 10.3390/v11040317.

517    [16]   N. Nguyen Quang *et al.*, "Dynamic nanopore long-read sequencing analysis of HIV-1

518          splicing events during the early steps of infection," *Retrovirology*, vol. 17, no. 1, p. 25,

519          2020, doi: 10.1186/s12977-020-00533-1.

520    [17]   A. G. Fisher, E. Collalti, L. Ratner, R. C. Gallo, and F. Wong-Staal, "A molecular clone of

521          HTLV-III with biological activity," *Nature*, vol. 316, no. 6025, pp. 262–265, 1985, doi:

522          10.1038/316262a0.

523    [18]   A. Adachi *et al.*, "Production of acquired immunodeficiency syndrome-associated

524          retrovirus in human and nonhuman cells transfected with an infectious molecular clone.,"

525          *J. Virol.*, vol. 59, no. 2, pp. 284–91, 1986.

526    [19]   P. Dickie *et al.*, "HIV-associated nephropathy in transgenic mice expressing HIV-1

527          genes," *Virology*, 1991. [Online]. Available: http://ac.els-cdn.com/0042682291907595/1-

528          s2.0-0042682291907595-main.pdf?_tid=8f811f10-d10c-11e5-82e8-

529          00000aacb35e&acdnat=1455228938_33d4226549c6410971ced1c4c3573a44. [Accessed:

530          11-Feb-2016].

531    [20]   M. Husain, "HIV-1 Nef Induces Proliferation and Anchorage-Independent Growth in

532          Podocytes," *J. Am. Soc. Nephrol.*, vol. 13, no. 7, pp. 1806–1815, 2002, doi:

533          10.1097/01.ASN.0000019642.55998.69.

534    [21]   H. Li *et al.*, "Epigenetic regulation of RCAN1 expression in kidney disease and its role in

535          podocyte injury," *Kidney Int.*, vol. 94, no. 6, pp. 1160–1176, 2018, doi:

536          10.1016/j.kint.2018.07.023.

537    [22]   S. Curreli *et al.*, "B cell lymphoma in HIV transgenic mice.," *Retrovirology*, vol. 10, p.

538          92, Jan. 2013, doi: 10.1186/1742-4690-10-92.

539     [23]    H. Li, "Minimap2: pairwise alignment for nucleotide sequences," *Bioinformatics*, vol. 34,

540            no. 18, pp. 3094–3100, May 2018, doi: 10.1093/bioinformatics/bty191.

541     [24]    H. Li and R. Durbin, "Fast and accurate long-read alignment with Burrows – Wheeler

542            transform," vol. 26, no. 5, pp. 589–595, 2010, doi: 10.1093/bioinformatics/btp698.

543     [25]    E. Afgan *et al.*, "The Galaxy platform for accessible, reproducible and collaborative

544            biomedical analyses: 2016 update," *Nucleic Acids Res.*, vol. 44, no. W1, pp. W3–W10,

545            2016, doi: 10.1093/nar/gkw343.

546     [26]    J. T. Robinson *et al.*, "Integrative genomics viewer," *Nat Biotechnol*, vol. 29, no. 1, pp.

547            24–26, 2011, doi: 10.1038/nbt0111-24.

548     [27]    B. P. Walenz, S. Koren, N. H. Bergman, A. M. Phillippy, J. R. Miller, and K. Berlin,

549            "Canu: scalable and accurate long-read assembly via adaptive k -mer weighting and repeat

550            separation," *Genome Res.*, vol. 27, no. 5, pp. 722–736, 2017, doi: 10.1101/gr.215087.116.

551     [28]    S. Andrews, "FastQC A Quality Control tool for High Throughput Sequence Data."

552     [29]    K. Katoh and D. M. Standley, "MAFFT multiple sequence alignment software version 7:

553            Improvements in performance and usability," *Mol. Biol. Evol.*, vol. 30, no. 4, pp. 772–780,

554            2013, doi: 10.1093/molbev/mst010.

555     [30]    K. Katoh, J. Rozewicki, and K. D. Yamada, "MAFFT online service: multiple sequence

556            alignment, interactive sequence choice and visualization," *Brief. Bioinform.*, vol. 20, no. 4,

557            pp. 1160–1166, Sep. 2017, doi: 10.1093/bib/bbx108.

558     [31]    D. M. Lyons and A. S. Lauring, "Evidence for the Selective Basis of Transition-to-

559            Transversion Substitution Bias in Two RNA Viruses," *Mol. Biol. Evol.*, vol. 34, no. 12,

560            pp. 3205–3215, 2017, doi: 10.1093/molbev/msx251.

561     [32]    N. J. Loman, J. Quick, and J. T. Simpson, "A complete bacterial genome assembled de

562   novo using only nanopore sequencing data," *Nat. Methods*, vol. 12, no. 8, pp. 733–735,

563   2015, doi: 10.1038/nmeth.3444.

564 [33] A. Gener *et al.*, "PEA0011 - Insights from HIV-1 transgene insertions in the murine

565   model of HIV-associated nephropathy," in *23rd International AIDS Conference*

566   *(AIDS2020)*, 2020, vol. ePoster.

567 [34] A. R. Gener *et al.*, "P39 - Insights from comprehensive transcript models of HIV-1," in

568   *Genome Informatics 2020*, 2020, p. ePoster.

569 [35] A. R. Gener, T. Washington, D. Hyink, and P. Klotman, "3264 - The Multiple HIV-1

570   Transgenes in the Murine Model of HIV-Associated Nephropathy Fail to Segregate as

571   Expected," in *American Society of Human Genetics Annual Meeting*, 2020, p. ePoster.

572 [36] A. Payne, N. Holmes, V. Rakyan, and M. Loose, "Whale watching with BulkVis: A

573   graphical viewer for Oxford Nanopore bulk fast5 files," *bioRxiv*, p. 312256, Jan. 2018,

574   doi: 10.1101/312256.

575 [37] P. Rosenstiel, A. Gharavi, V. D'Agati, and P. Klotman, "Transgenic and infectious animal

576   models of HIV-associated nephropathy.," *J. Am. Soc. Nephrol.*, vol. 20, no. 11, pp. 2296–

577   304, 2009, doi: 10.1681/ASN.2008121230.

578 [38] M. Kvaratskhelia, A. Sharma, R. C. Larue, E. Serrao, and A. Engelman, "Molecular

579   mechanisms of retroviral integration site selection.," *Nucleic Acids Res.*, vol. 42, no. 16,

580   pp. gku769-, 2014, doi: 10.1093/nar/gku769.

581 [39] B. Marini *et al.*, "Nuclear architecture dictates HIV-1 integration site selection," *Nature*,

582   vol. 521, pp. 227–233, 2015, doi: 10.1038/nature14226.

583 [40] H. Ochman, A. S. Gerber, and D. L. Hartl, "Genetic applications of an inverse polymerase

584   chain reaction," *Genetics*, vol. 120, no. 3, pp. 621–623, 1988.

585 [41] W.-S. Hu and S. H. Hughes, "HIV-1 Reverse Transcription," *Cold Spring Harb. Perspect.*

586    *Med.* , vol. 2, no. 10, Oct. 2012, doi: 10.1101/cshperspect.a006882.

587 [42] M. R. Green, T. Maniatis, and D. A. Melton, "Human beta-globin pre-mRNA synthesized

588    in vitro is accurately spliced in Xenopus  oocyte nuclei.," *Cell*, vol. 32, no. 3, pp. 681–

589    694, Mar. 1983, doi: 10.1016/0092-8674(83)90054-5.

590 [43] A. R. Gener, "Full-coverage sequencing of HIV-1 provirus from a reference plasmid,"

591    *bioRxiv*, p. 611848, Jan. 2019, doi: 10.1101/611848.

592 [44] B. Lucic *et al.*, "Spatially clustered loci with multiple enhancers are frequent targets of

593    HIV-1," *bioRxiv*, 2018.

594 [45] W. J. Kent *et al.*, "The Human Genome Browser at UCSC," *Genome Res.* , vol. 12, no. 6,

595    pp. 996–1006, Jun. 2002, doi: 10.1101/gr.229102.

596 [46] Y. Peng, H. C. M. Leung, S. M. Yiu, and F. Y. L. Chin, "IDBA – A Practical Iterative de

597    Bruijn Graph De Novo Assembler BT  - Research in Computational Molecular Biology,"

598    2010, pp. 426–440.

599 [47] P. J. A. Cock, B. A. Grüning, K. Paszkiewicz, and L. Pritchard, "Galaxy tools and

600    workflows for sequence analysis with applications in molecular plant pathology," *PeerJ*,

601    vol. 1, p. e167, 2013, doi: 10.7717/peerj.167.

602 [48] C. B, W. T, and S. S, "Genome sequence assembly using trace signals and additional

603    sequence information.," *Comput. Sci. Biol. Proc. Ger. Conf. Bioinforma.*, vol. 99, pp. 45–

604    56.

605 [49] A. Bankevich *et al.*, "SPAdes: A New Genome Assembly Algorithm and Its Applications

606    to Single-Cell Sequencing," *J. Comput. Biol.*, vol. 19, no. 5, pp. 455–477, Apr. 2012, doi:

607    10.1089/cmb.2012.0021.

608    [50]    G. Cuccuru *et al.*, "Orione, a web-based framework for NGS analysis in microbiology,"

609         *Bioinformatics*, vol. 30, no. 13, pp. 1928–1929, Jul. 2014, doi:

610         10.1093/bioinformatics/btu135.

611    [51]    R. L. Warren, G. G. Sutton, S. J. M. Jones, and R. A. Holt, "Assembling millions of short

612         DNA sequences using SSAKE," *Bioinformatics*, vol. 23, no. 4, pp. 500–501, Feb. 2007,

613         doi: 10.1093/bioinformatics/btl629.

614

615

616 **Tables**

617

618 **Table 1: Summary of pHXB2 sample divergence from reference HXB2.**

| Site | Position | Change | Substitution Class | Change | Mutation Class (Syn/Non/Stop) | Homopolymer-adjacent? | Same as neighbor? | LANL Feature | Subfeature | Frame |
|------|----------|--------|-------------------|--------|------------------------------|----------------------|-------------------|--------------|------------|-------|
| 1 | 24 | C>A | transversion | NA | NA | yes | yes | 5'LTR | U3 | NA |
| 2 | 108 | A>G | transition | NA | NA | yes | yes | 5'LTR | U3 | NA |
| 3 | 164 | G>T | transversion | NA | NA | yes | no | 5'LTR | U3 | NA |
| 4 | 168 | T>G | transversion | NA | NA | yes | yes | 5'LTR | U3 | NA |
| 5 | 176 | A>G | transition | NA | NA | yes | yes | 5'LTR | U3 | NA |
| 6 | 182 | C>T | transition | NA | NA | yes | no | 5'LTR | U3 | NA |
| 7 | 227 | A>G | transition | NA | NA | yes | yes | 5'LTR | U3 | NA |
| 8 | 291 | A>G | transition | NA | NA | no | no | 5'LTR | U3 | NA |
| 9 | 333 | C>T | transition | NA | NA | no | no | 5'LTR | U3 | NA |
| 10 | 654 | C>T | transition | NA | NA | no | no | None | None | NA |
| 11 | 1659 | aaG>aaA | transition | None | Syn | yes | yes | gag | p24, p55 | gag frame 1 |
| 12 | 2259 | gag:agG>agA pol:Gtc>Atc | transition | gag:Arg>Arg pol:Val>Ile | Syn/Non | no | no | gagpol | p6 | gag frame 1 pol frame 3 |
| 13 | 2927 | aaG>aaA | transition | None | Syn | yes | yes | pol | p51 RT | pol frame 3 |
| 14 | 3812 | ccC>ccT | transition | None | Syn | yes | yes | pol | p51 RT | pol frame 3 |
| 15 | 4574 | acT>acA | transversion | None | Syn | no | no | pol | p31 IN | pol frame 3 |
| 16 | 4596 | Ggt>Agt | transition | None | Syn | yes | no | pol | p31 IN | pol frame 3 |
| 17 | 4609 | aGg>aAg | transition | Arg>Lys | Non | yes | yes | pol | p31 IN | pol frame 3 |
| 18 | 7823 | gcC>gcG Ggc>Cgc | transversion | ASP:Gly>Arg | Syn/Non | no | no | gp41 | RRE, also ASP | gp41 frame 3, ASP -2 |
| 19 | 9253 | Ata>Gta | transition | Ile>val | Non | no | yes | nef/3'LTR | also U3 | nef frame 1 |
| 20 | 9418 | C>T | transition | NA | NA | no | no | 3'LTR | U3 | NA |

619 Coverage numbers vary by input (albacore, guppy, FlipFlop basecalled FASTQ) and mapping

620 method (Minimap2 vs. BWA-MEM). This information is provided as Supplemental Digital
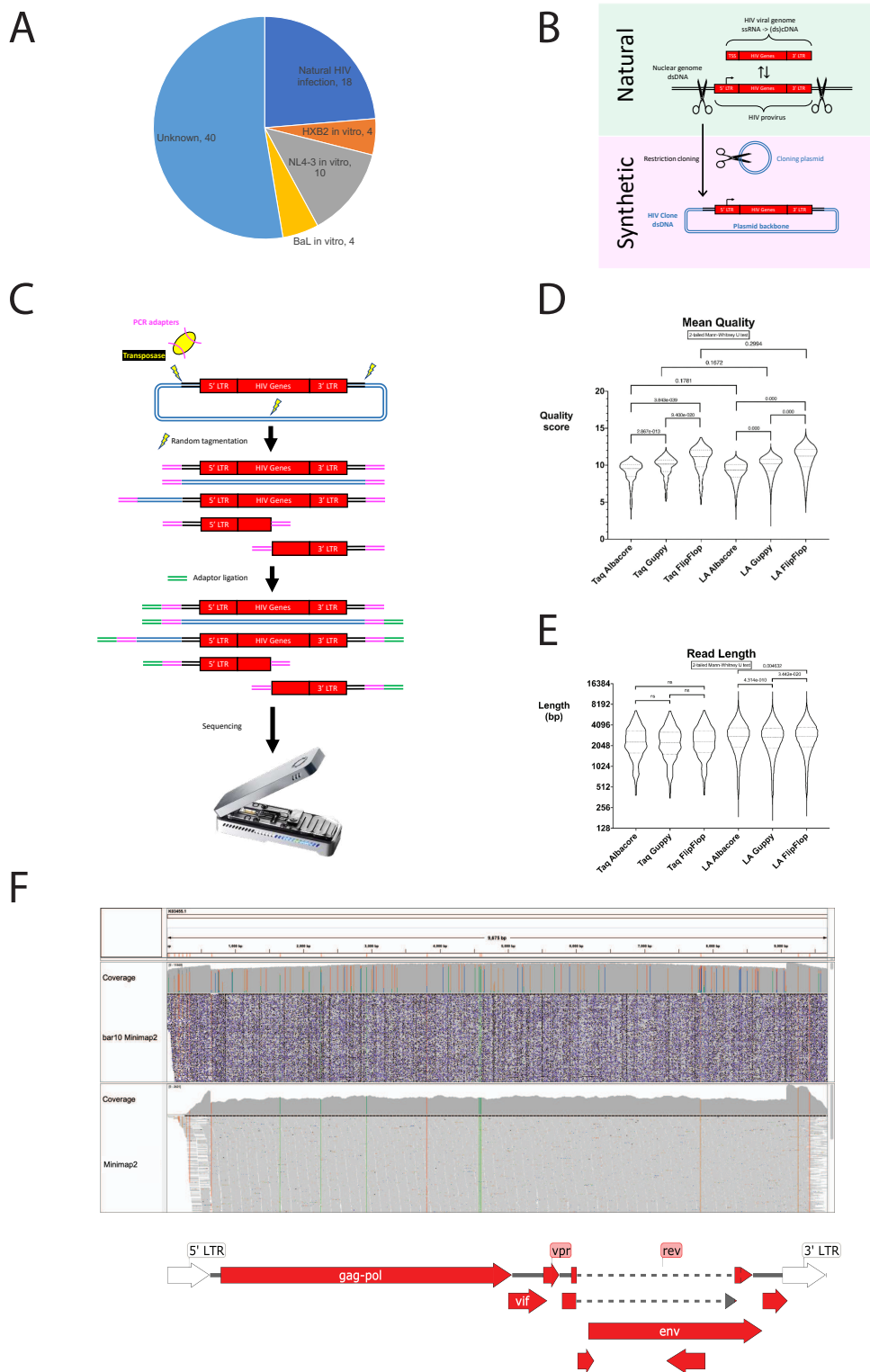
621 Content. Base-1 (first base is numbered 1, 2nd 2, etc.), relative to HXB2, Genbank:K03455.1.

622 Changed base represented as upper-case. Annotated as codon if in protein-coding region. No

623 deletions or insertions were predicted from manual inspection or supported by short-read

624 sequencing. Abbreviations, ASP: antisense protein, RRE: rev-response element, NA: not

625 applicable. Syn: synonymous mutation. Non: non-synonymous mutation. Stop: stop codon/non-

626 sense mutation. LTR: long terminal repeat. RT: reverse transcriptase. IN: integrase. LANL: Los

627 Alamos National Laboratory HIV Sequence Database. Data from three separate sequencing

628 experiments on the same plasmid sample support these 20 sites. Note site 1-8 variants in 5'LTR

629 have been previously reported (LANL), albethey ambiguously. These may also be incorrectly

630 annotated as variants in nef.

631 **Figures**

632

633

634    **Figure 1: HIV information in pHXB2_D is recovered by long-read sequencing and**

635    **mapping.**



636

637    Figure 1A: HXB2 is still a commonly used resource. It is the reference HIV-1 genome, derived

638    from one of the earliest clinical isolates. While older HIV samples are occasionally rediscovered,

639    they are not made routinely available to researchers. All public HIV-1 RNA-seq datasets were

640    obtained from the NCBI SRA with the following search phrase: "HIV-1" AND "RNA-seq".

641    Metadata from these 2527 runs (number current as of 7/21/2020) were used to make a pie chart

642    summary.

643    Figure 1B: HIV information comes from three main sources: proviruses (HIV sandwiched

644    between two assumedly identical full-length long terminal repeats (LTRs)), unspliced HIV

645    mRNAs (also known as viral genomes) starting from the transcription start site and ending in the

646    3' LTR [4], and engineered proviruses recovered in their entirety or stitched together from

647    multiple isolates like NL4-3 [18].

648    Figure 1C: ONT library prep pipeline. Tagmentation cleaves double-stranded DNA, ligating

649    barcoded PCR adapters (magenta). PCR-adapted DNA may be amplified. After amplification

650    and cleanup, ONT sequencing adapters (green) are ligated. Barcoded samples may be pooled and

651    sequenced.

652    Figure 1D: Newer basecallers increase read mean quality. Median (big dash) and quartiles (little

653    dash). Effect of enzyme version was not statistically significant.

654    Figure 1E: Read stats with different callers/aligners. Median (big dash) and quartiles (little dash).

655    Read lengths increase with higher fidelity Taq.

656    Figure 1F: Sequencing coverage with long- vs. short-read single-end 150 bp (trimmed to 142 bp)

657    DNA sequencing. Long-read sequencing covers ambiguously mappable areas missed by short-

658    read in HXB2 reference Genbank:K03455.1 (**Supplemental Figures 3D,3E**), but at the expense

659    of accuracy near homopolymers longer than about 4 nucleobases (**Supplemental Figure 5**).

660    Short-read mapping fails at repetitive elements longer than their read lengths (**Supplemental**

661    **Figures 3D,3E**). Long read Minimap2 settings: map-ont -k15. Short read Minimap2 settings:

662    Short reads without splicing (-k21 -w11 --sr -F800 -A2 -B8 -O12,32 -E2,1 -r50 -p.5 -N20 -

663    f1000,5000 -n2 -m20 -s40 -g200 -2K50m --heap-sort=yes --secondary=no) (sr).

664

665

666

667 **Figure 2: Longest read containing complete full-length HIV-1 reference HXB2**

```
@6fbf0205-5195-460e-8e28-930db50e5d79 runid=0b284792282af9a6d7275cfca845556bb4b8ac7f sampleid=Student_Genomics_Run_1 read=87223 ch=460 start_time=2018-05-09T20:57:49Z barcode=barcode10
```



668

669  The 5<sup>th</sup> longest read in the barcode 10 set (read ID 6fbf0205-5195-460e-8e28-930db50e5d79)

670  contained full-length HIV-1. Query (full read) blastn against HIV (taxid:11676) returned 92.95%

671  identity to HIV-1, complete genome (Genbank:AF033819.3). Limiting query to HXB2 (red)

672  blastn against Nucleotide collection nr/nt returned 100% coverage and 93.02% identity to HIV-1

673  HXB2. This read was 11,487 bases long, with mean quality score 11.984396. Basecalled using

674  Guppy 2.3.1 with FlipFlop config.

675

676 **Figure 3A: pHXB2_D has identical LTRs, resolving likely errors in HXB2 (K03455.1)**

677 ```
CLUSTAL format alignment by MAFFT (v7.475)
```
678
679
680 ```
K03455.1_5'LTR   tggaagggctaattcactcccaacgaagacaagatatccttgatctgtggatctaccaca
```
681 ```
pHXB2_D_5'LTR    tggaagggctaattcactcccaaagaagacaagatatccttgatctgtggatctaccaca
```
682 ```
pHXB2_D_3'LTR    tggaagggctaattcactcccaaagaagacaagatatccttgatctgtggatctaccaca
```
683 ```
K03455.1_3'LTR   tggaagggctaattcactcccaaagaagacaagatatccttgatctgtggatctaccaca
```
684 ```
                 ********************** **************************************
```
685
686 ```
K03455.1_5'LTR   cacaaggctacttccctgattagcagaactacacaccagggccagggatcagatatccac
```
687 ```
pHXB2_D_5'LTR    cacaaggctacttccctgattagcagaactacacaccagggccaggggtcagatatccac
```
688 ```
pHXB2_D_3'LTR    cacaaggctacttccctgattagcagaactacacaccagggccaggggtcagatatccac
```
689 ```
K03455.1_3'LTR   cacaaggctacttccctgattagcagaactacacaccagggccaggggtcagatatccac
```
690 ```
                 ********************************************** .************
```
691
692 ```
K03455.1_5'LTR   tgacctttggatggtgctacaagctagtaccagttgagccagagaagttagaagaagcca
```
693 ```
pHXB2_D_5'LTR    tgacctttggatggtgctacaagctagtaccagttgagccagataaggtagaagaggcca
```
694 ```
pHXB2_D_3'LTR    tgacctttggatggtgctacaagctagtaccagttgagccagataaggtagaagaggcca
```
695 ```
K03455.1_3'LTR   tgacctttggatggtgctacaagctagtaccagttgagccagataagatagaagaggcca
```
696 ```
                 ****************************************** *** ******* ****
```
697
698 ```
K03455.1_5'LTR   acaaaggagagaacaccagcttgttacaccctgtgagcctgcatggaatggatgacccgg
```
699 ```
pHXB2_D_5'LTR    ataaaggagagaacaccagcttgttacaccctgtgagcctgcatgggatggatgacccgg
```
700 ```
pHXB2_D_3'LTR    ataaaggagagaacaccagcttgttacaccctgtgagcctgcatgggatggatgacccgg
```
701 ```
K03455.1_3'LTR   ataaaggagagaacaccagcttgttacaccctgtgagcctgcatgggatggatgacccgg
```
702 ```
                 *.*********************************************.************
```
703
704 ```
K03455.1_5'LTR   agagagaagtgttagagtggaggtttgacagccgcctagcatttcatcacatggcccgag
```
705 ```
pHXB2_D_5'LTR    agagagaagtgttagagtggaggtttgacagccgcctagcatttcatcacgtggcccgag
```
706 ```
pHXB2_D_3'LTR    agagagaagtgttagagtggaggtttgacagccgcctagcatttcatcacgtggcccgag
```
707 ```
K03455.1_3'LTR   agagagaagtgttagagtggaggtttgacagccgcctagcatttcatcacgtggcccgag
```
708 ```
                 *************************************************.*********
```
709
710 ```
K03455.1_5'LTR   agctgcatccggagtacttcaagaactgctgacatcgagcttgctacaagggactttccg
```
711 ```
pHXB2_D_5'LTR    agctgcatccggagtacttcaagaactgctgatatcgagcttgctacaagggactttccg
```
712 ```
pHXB2_D_3'LTR    agctgcatccggagtacttcaagaactgctgatatcgagcttgctacaagggactttccg
```
713 ```
K03455.1_3'LTR   agctgcatccggagtacttcaagaactgctgacatcgagcttgctacaagggactttccg
```
714 ```
                 ********************************.***************************
```
715
716 ```
K03455.1_5'LTR   ctggggactttccagggaggcgtggcctgggcgggactggggagtggcgagccctcagat
```

```
717  pHXB2_D_5'LTR    ctggggactttccagggaggcgtggcctgggcgggactggggagtggcgagccctcagat
718  pHXB2_D_3'LTR    ctggggactttccagggaggcgtggcctgggcgggactggggagtggcgagccctcagat
719  K03455.1_3'LTR   ctggggactttccagggaggcgtggcctgggcgggactggggagtggcgagccctcagat
720                   ************************************************************

722  K03455.1_5'LTR   cctgcatataagcagctgctttttgcctgtactgggtctctctggttagaccagatctga
723  pHXB2_D_5'LTR    cctgcatataagcagctgctttttgcctgtactgggtctctctggttagaccagatctga
724  pHXB2_D_3'LTR    cctgcatataagcagctgctttttgcctgtactgggtctctctggttagaccagatctga
725  K03455.1_3'LTR   cctgcatataagcagctgctttttgcctgtactgggtctctctggttagaccagatctga
726                   ************************************************************

728  K03455.1_5'LTR   gcctgggagctctctggctaactagggaacccactgcttaagcctcaataaagcttgcct
729  pHXB2_D_5'LTR    gcctgggagctctctggctaactagggaacccactgcttaagcctcaataaagcttgcct
730  pHXB2_D_3'LTR    gcctgggagctctctggctaactagggaacccactgcttaagcctcaataaagcttgcct
731  K03455.1_3'LTR   gcctgggagctctctggctaactagggaacccactgcttaagcctcaataaagcttgcct
732                   ************************************************************

734  K03455.1_5'LTR   tgagtgcttcaagtagtgtgtgcccgtctgttgtgtgactctggtaactagagatccctc
735  pHXB2_D_5'LTR    tgagtgcttcaagtagtgtgtgcccgtctgttgtgtgactctggtaactagagatccctc
736  pHXB2_D_3'LTR    tgagtgcttcaagtagtgtgtgcccgtctgttgtgtgactctggtaactagagatccctc
737  K03455.1_3'LTR   tgagtgcttcaagtagtgtgtgcccgtctgttgtgtgactctggtaactagagatccctc
738                   ************************************************************

740  K03455.1_5'LTR   agacccttttagtcagtgtggaaaatctctagca
741  pHXB2_D_5'LTR    agacccttttagtcagtgtggaaaatctctagca
742  pHXB2_D_3'LTR    agacccttttagtcagtgtggaaaatctctagca
743  K03455.1_3'LTR   agacccttttagtcagtgtggaaaatctctagca
744                   **********************************
```

747 **Figure 3B: pNL4-3_gag-pol(Δ1443-4553)_EGFP (ACCESSION_TBD) has distinct LTRs,**

748 **consistent with pNL4-3 (AF324493.1)**

```
749   CLUSTAL format alignment by MAFFT (v7.475)
750
751
752   AF324493.1_5LTR  tggaagggctaatttggtcccaaaaaagacaagagatccttgatctgtggatctaccaca
753   ACCESSION_TBD_5  tggaagggctaatttggtcccaaaaaagacaagagatccttgatctgtggatctaccaca
754   AF324493.1_3LTR  tggaagggctaattcactcccaaagaagacaagatatccttgatctgtggatctaccaca
755   ACCESSION_TBD_3  tggaagggctaattcactcccaaagaagacaagatatccttgatctgtggatctaccaca
756                    **************..  *******.********* ************************
757
758   AF324493.1_5LTR  cacaaggctacttccctgattggcagaactacacaccagggccagggatcagatatccac
759   ACCESSION_TBD_5  cacaaggctacttccctgattggcagaactacacaccagggccagggatcagatatccac
760   AF324493.1_3LTR  cacaaggctacttccctgattggcagaactacacaccagggccaggggtcagatatccac
761   ACCESSION_TBD_3  cacaaggctacttccctgattggcagaactacacaccagggccaggggtcagatatccac
762                    ********************************************.************
763
764   AF324493.1_5LTR  tgacctttggatggtgcttcaagttagtaccagttgaaccagagcaagtagaagaggcca
765   ACCESSION_TBD_5  tgacctttggatggtgcttcaagttagtaccagttgaaccagagcaagtagaagaggcca
766   AF324493.1_3LTR  tgacctttggatggtgctacaagctagtaccagttgagccagataaggtagaagaggcca
767   ACCESSION_TBD_3  tgacctttggatggtgctacaagctagtaccagttgagccagataaggtagaagaggcca
768                    ***************** ****.*************.*****  *.*************
769
770   AF324493.1_5LTR  atgaaggagagaacaacagcttgttacaccctatgagccagcatgggatggaggacccgg
771   ACCESSION_TBD_5  atgaaggagagaacaacagcttgttacaccctatgagccagcatgggatggaggacccgg
772   AF324493.1_3LTR  ataaaggagagaacaccagcttgttacaccctgtgagcctgcatggaatggatgaccctg
773   ACCESSION_TBD_3  ataaaggagagaacaccagcttgttacaccctgtgagcctgcatggaatggatgaccctg
774                    **.*********** *****************.****** ******.***** ***** *
775
776   AF324493.1_5LTR  agggagaagtattagtgtggaagtttgacagcctcctagcatttcgtcacatggcccgag
777   ACCESSION_TBD_5  agggagaagtattagtgtggaagtttgacagcctcctagcatttcgtcacatggcccgag
778   AF324493.1_3LTR  agagagaagtgttagagtggaggtttgacagccgcctagcatttcatcacgtggcccgag
779   ACCESSION_TBD_3  agagagaagtgttagagtggaggtttgacagccgcctagcatttcatcacgtggcccgag
780                    **.*******.**** *****.*********** **********.****.*********
781
782   AF324493.1_5LTR  agctgcatccggagtactacaaagactgctgacatcgagctttctacaagggactttccg
783   ACCESSION_TBD_5  agctgcatccggagtactacaaagactgctgacatcgagctttctacaagggactttccg
784   AF324493.1_3LTR  agctgcatccggagtacttcaagaactgctgacatcgagcttgctacaagggactttccg
785   ACCESSION_TBD_3  agctgcatccggagtacttcaagaactgctgacatcgagcttgctacaagggactttccg
786                    ***************** ***..****************** ****************
```
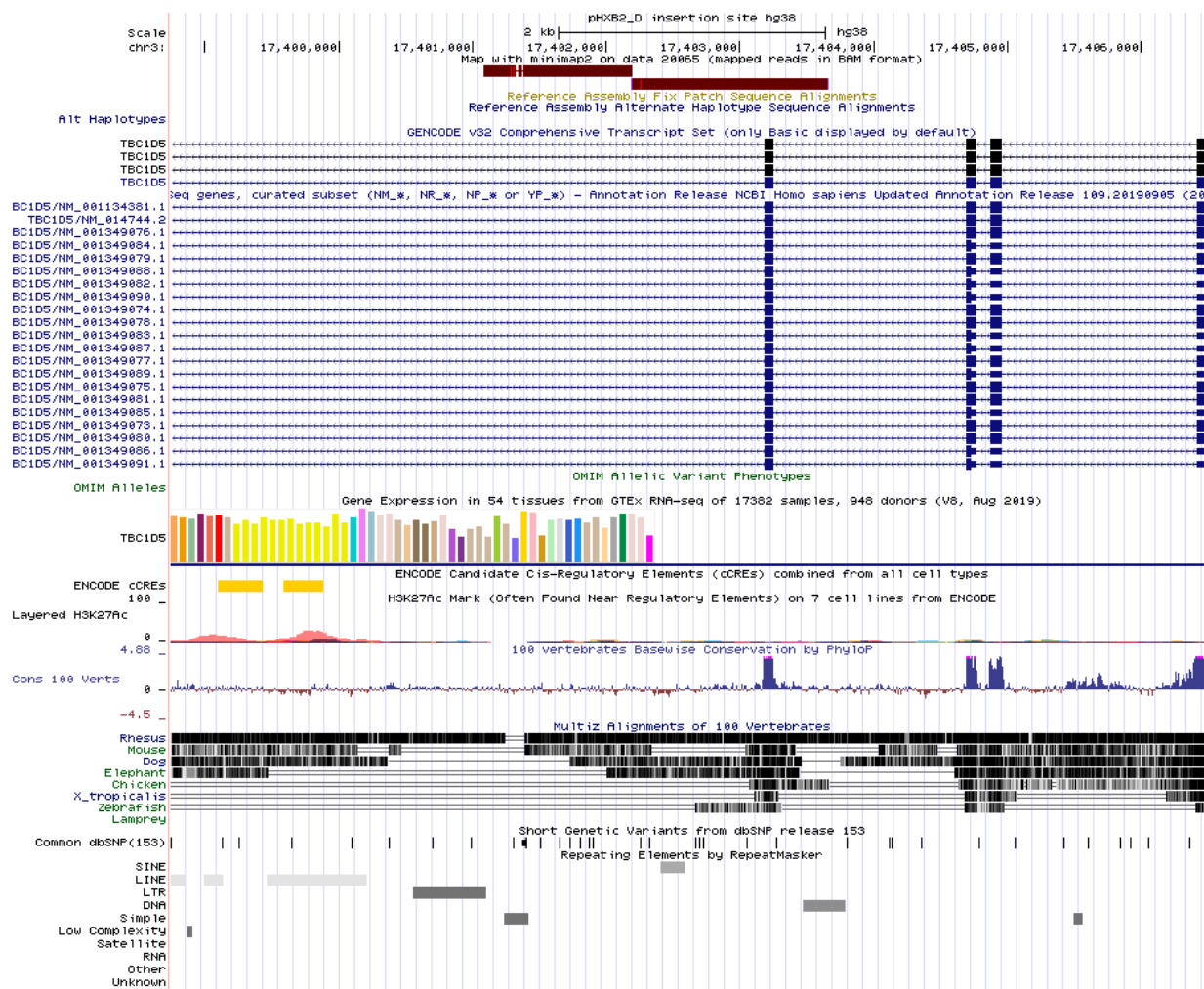
```
AF324493.1_5LTR   ctggggactttccagggaggtgtggcctgggcgggactggggagtggcgagccctcagat
ACCESSION_TBD_5   ctggggactttccagggaggtgtggcctgggcgggactggggagtggcgagccctcagat
AF324493.1_3LTR   ctggggactttccagggaggcgtggcctgggcgggactggggagtggcgagccctcagat
ACCESSION_TBD_3   ctggggactttccagggaggcgtggcctgggcgggactggggagtggcgagccctcagat
                  ********************.****************************************

AF324493.1_5LTR   gctacatataagcagctgcttttttgcctgtactgggtctctctggttagaccagatctga
ACCESSION_TBD_5   gctacatataagcagctgcttttttgcctgtactgggtctctctggttagaccagatctga
AF324493.1_3LTR   gctgcatataagcagctgcttttttgcctgtactgggtctctctggttagaccagatctga
ACCESSION_TBD_3   gctgcatataagcagctgcttttttgcctgtactgggtctctctggttagaccagatctga
                  ***.********************************************************

AF324493.1_5LTR   gcctgggagctctctggctaactagggaacccactgcttaagcctcaataaagcttgcct
ACCESSION_TBD_5   gcctgggagctctctggctaactagggaacccactgcttaagcctcaataaagcttgcct
AF324493.1_3LTR   gcctgggagctctctggctaactagggaacccactgcttaagcctcaataaagcttgcct
ACCESSION_TBD_3   gcctgggagctctctggctaactagggaacccactgcttaagcctcaataaagcttgcct
                  ************************************************************

AF324493.1_5LTR   tgagtgctcaaagtagtgtgtgcccgtctgttgtgtgactctggtaactagagatccctc
ACCESSION_TBD_5   tgagtgctcaaagtagtgtgtgcccgtctgttgtgtgactctggtaactagagatccctc
AF324493.1_3LTR   tgagtgcttcaagtagtgtgtgcccgtctgttgtgtgactctggtaactagagatccctc
ACCESSION_TBD_3   tgagtgcttcaagtagtgtgtgcccgtctgttgtgtgactctggtaactagagatccctc
                  ********. **************************************************

AF324493.1_5LTR   agacccttttagtcagtgtggaaaatctctagca
ACCESSION_TBD_5   agacccttttagtcagtgtggaaaatctctagca
AF324493.1_3LTR   agacccttttagtcagtgtggaaaatctctagca
ACCESSION_TBD_3   agacccttttagtcagtgtggaaaatctctagca
                  **********************************
```

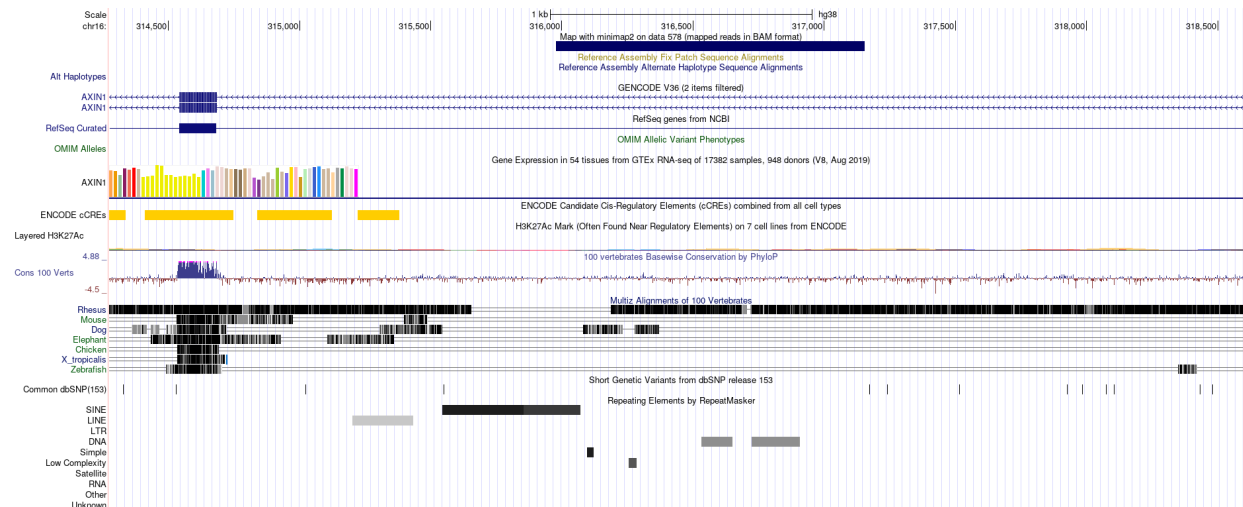818    **Figure 4A: HXB2 integration site**

819



820

821

822
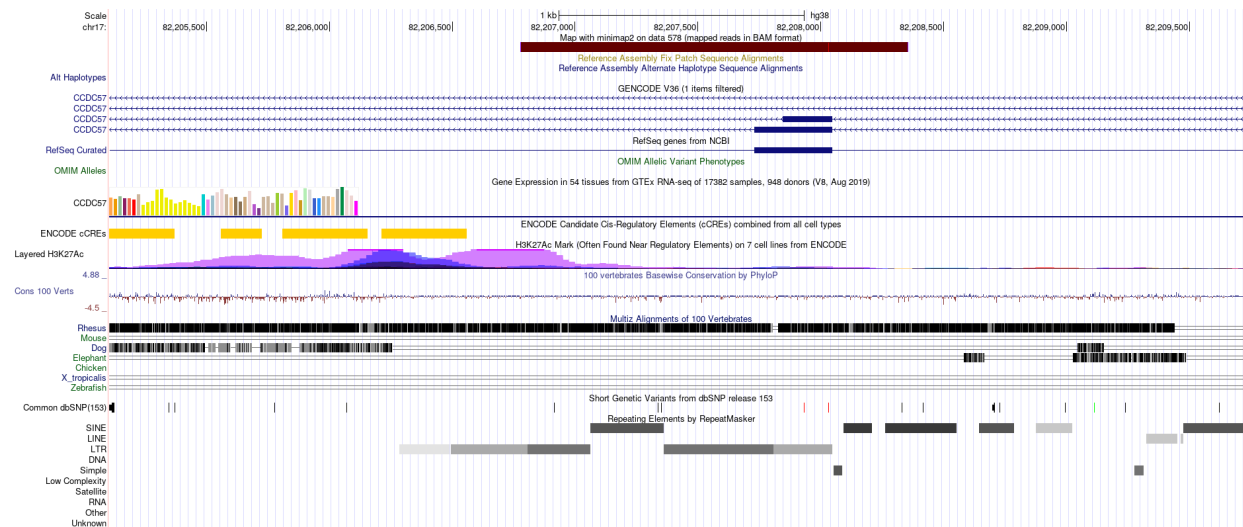
823 **Figure 4B: NL4-3 integration sites**

824

825



826
827

828



829
830

831 Figure 4A: pHXB2_D's, and therefore HXB2's, integration site is unambiguously singular (falls

832 outside of annotated repeat), and in the same orientation (minus strand relative to hg38) as target

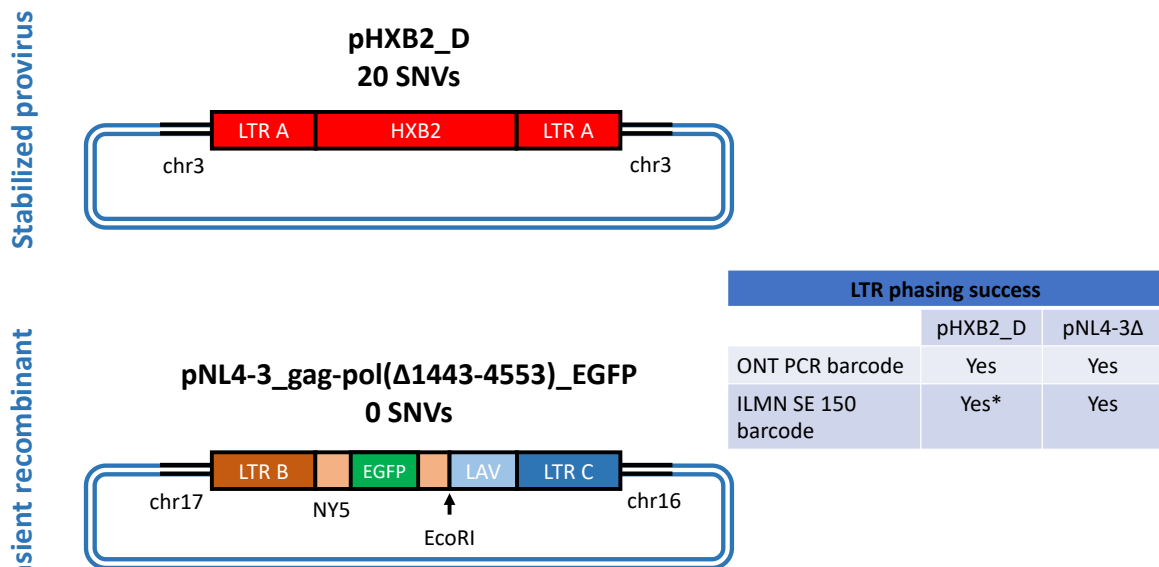833 gene TBC1D5. Alignment quality is 60 for both homology arms (**Supplemental Table 2**).

834 Features captured by homology arms in pHXB2_D and other clones verified as proviruses in the

835 present study are consistent with HIV-1 integration behavior [44]. Visualized in UCSC Genome

836    Browser [45]. Figure 4B: pNL4-3_gag-pol(Δ1443-4553)_EGFP's, and therefore NL4-3's,

837    integration sites fall on annotated repeats, the longer reads help to locate both sites. Alignment

838    quality is 60 for both homology arms (**Supplemental Table 2**). These integration sites would

839    likely be missed by any method leveraging reads shorter than the homology arms.

840

841     **Figure 5: pHXB2_D provenance and top 50 neighbors**



842

843    Figure 6: Summary of long- vs. short-read mapping by ability to phase LTRs



**Stabilized provirus**

**pHXB2_D**
**20 SNVs**

chr3 — LTR A — HXB2 — LTR A — chr3

**Transient recombinant**

**pNL4-3_gag-pol(Δ1443-4553)_EGFP**
**0 SNVs**

chr17 — LTR B — NY5 — EGFP — LAV — LTR C — chr16

EcoRI

LTR C takes over in subsequent viruses.

| LTR phasing success | | |
|---|---|---|
| | pHXB2_D | pNL4-3Δ |
| ONT PCR barcode | Yes | Yes |
| ILMN SE 150 barcode | Yes* | Yes |

*intermittent success with ambiguous mapping at LTR.

844

845 **Supplemental Information**

846 **Data exploration with long- and short-read mapping**

847      To assemble pHXB2_D, we tried the following short read assemblers on short-read data

848 from the external core: IDBA [46], MIRA [47], [48], SPAdes [49], and SSAKE [50], [51]. These

849 were chosen as a convenience because they were already stably implemented in Galaxy

850 (specifically usegalaxy.eu). Of these, SSAKE produced discontinuous assemblies with default

851 parameters. The discontinuous contigs did however map to the core's assembly (not shown).

852 **Enabling STEM outreach**

853      This work was performed as two control experiments with identically prepared libraries

854 for a STEM outreach initiative, Student Genomics (Gener, et al., manuscript in prep). Given the

855 constraints of the Student Genomics pilot, a rapid sequencing kit with tagmentation (explained

856 below) with PCR barcoding was used to pool samples for ONT sequencing, with the

857 consequence of fragmenting plasmid DNA more than what would have been ideal for capturing

858 full-length HIV. That said, these controls could have been just as easily replaced by any

859 samples/experiments benefiting from long-read sequencing at moderate-to-high coverage.

860

861    **Supplemental Tables**

862

863

864　**Supplemental Table 1: HXB2 is still a common HIV clone.**

865　See **Supplemental Digital Content.**

866　See also **Figure 1A**.

867

868

869

870 **Supplemental Table 2: HIV provirus clones**

871 See **Supplemental Digital Content.**

872 Of the HIV clones available through ARP, the table represents the only validated proviruses with

873 both upstream and downstream homology arms mapping to the same integration sites. pNL4-3 is

874 included as a known chimera with two integration half-sites. Other clones were made with

875 cDNA cloning, usually TA cloning (per ARP entries). Note: Reference hg38. Aligner: minimap2

876 with "Long Assembly" mapping settings. All homology arms had Alignment quality = 60.

877 Upstream = host plus strand; independent of integration orientation. Coordinates reported from

878 UCSC. ARP = NIH AIDS Reagent and Reference Program. IS = integration site.

879 **Supplemental Table 3: Variation in assemblies at the feature level.**

| | Mismatches | | | | | | Gaps (INDEL) | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Taq | | | LA Taq | | | Taq | | | LA Taq | | |
| | Albacore | Guppy | FlipFlop | Albacore | Guppy | FlipFlop | Albacore | Guppy | FlipFlop | Albacore | Guppy | FlipFlop |
| 5' LTR | NA | 9 | 9 | 9 | 9 | 9 | NA | NA | 0 | 2 | 0 | 0 |
| gag | 2 | 2 | 2 | 2 | 2 | 2 | 12 | 10 | 9 | 9 | 8 | 8 |
| 5' LTR+ψ | 10 | 10 | NA | NA | NA | NA | 5 | 2 | NA | NA | NA | NA |
| pol | 7 | 6 | 6 | 6 | 6 | 6 | 26 | 22 | 9 | 18 | 11 | 10 |
| vif | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 3 | 2 | 3 | 1 | 1 |
| vpr | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 |
| tat | 2 | 1 | 1 | 1 | 1 | 1 | 10 | 6 | 3 | 7 | 4 | 5 |
| rev | 2 | 1 | 1 | 1 | 1 | 1 | 10 | 7 | 4 | 7 | 4 | 5 |
| vpu | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| gp160 | 2 | 1 | 1 | 1 | 1 | 1 | 11 | 7 | 4 | 8 | 4 | 5 |
| nef | 1 | 1 | 1 | 1 | 1 | 1 | 3 | 2 | 2 | 3 | 2 | 1 |
| 3' LTR | 2 | 2 | NA | NA | NA | NA | 2 | 0 | NA | NA | NA | NA |
| nef+3' LTR | NA | 2 | 2 | 2 | 2 | 2 | NA | 2 | 2 | 5 | 2 | 1 |
| HXB2 | 22 | 20 | 20 | 20 | 20 | 20 | 61 | 46 | 28 | 47 | 27 | 25 |
| Downstream bridge | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 5 | 1 | 2 | 1 | 1 |
| pBR322-related | 0 | 0 | 0 | 0 | 0 | 0 | 19 | 19 | 13 | 18 | 19 | 15 |
| Upstream bridge | 2 | 2 | 3 | 2 | 2 | 2 | 8 | 6 | 5 | 7 | 7 | 3 |

880

881 Assembled with Canu. NA denotes features which may not have matched exactly, but which

882 were collapsed with adjacent features to facilitate counting. Variants called manually by

883 mapping assemblies over HXB2 features with SnapGene.
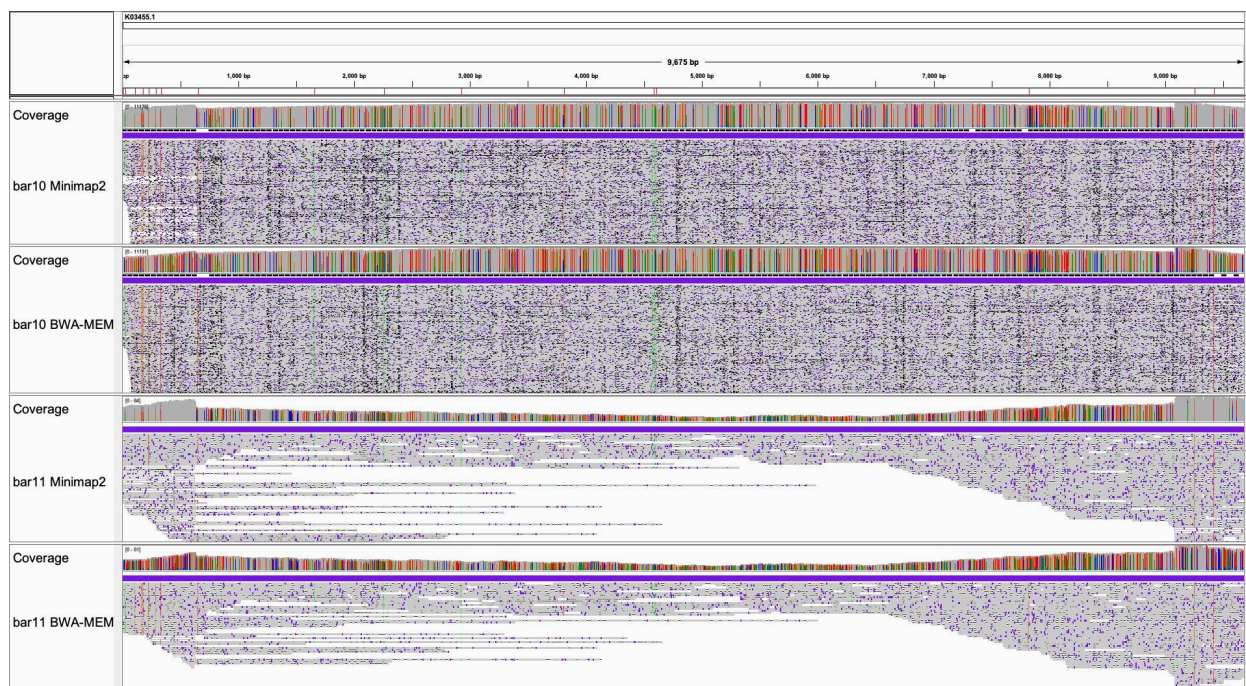
884
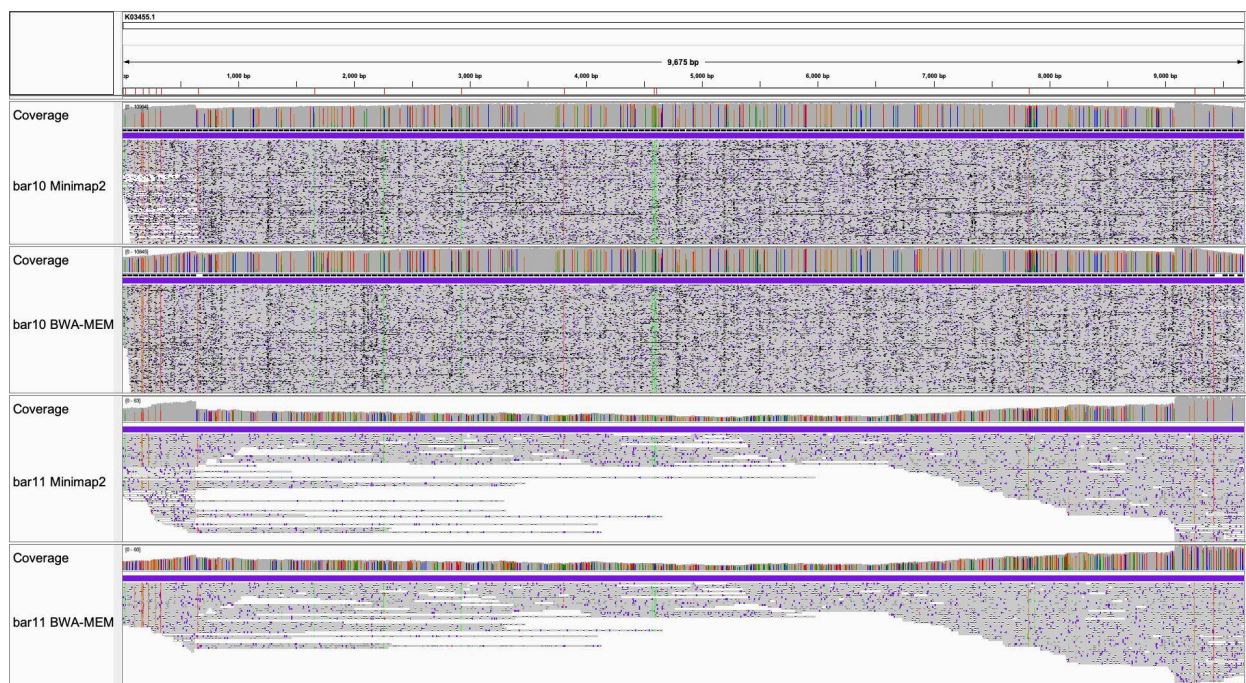
885

886 **Supplemental Figures**

887

888 **Supplemental Figure 1A: Unbiased nanopore DNA sequencing coverage over HXB2**

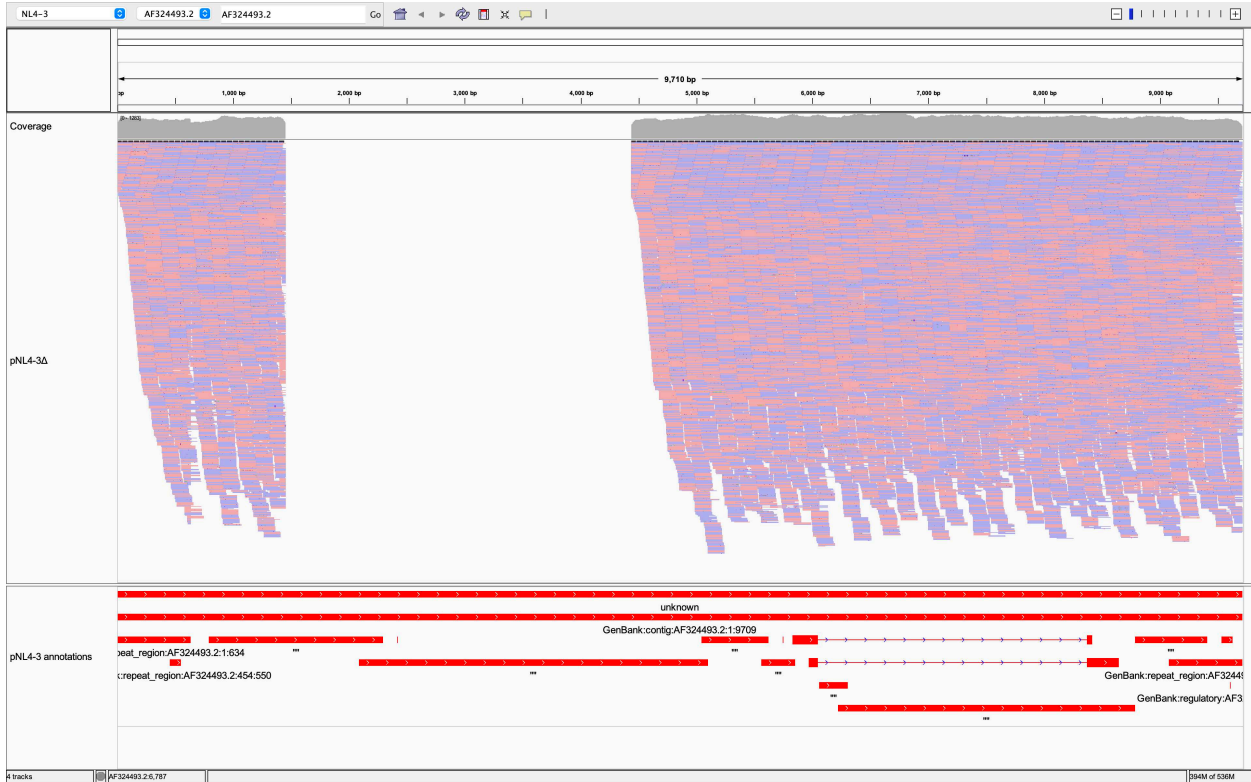889 **depends on DNA polymerase and mapper. ONT basecaller = Albacore (worst).**



890
891

892 **Supplemental Figure 1B: Unbiased nanopore DNA sequencing coverage over HXB2**

893 **depends on DNA polymerase and mapper. ONT basecaller = Guppy.**



894
895

896 **Supplemental Figure 1C: Unbiased nanopore DNA sequencing coverage over HXB2**

897 **depends on DNA polymerase and mapper. ONT basecaller = FlipFlop (best).**



898
899

900  Top two Coverage and Alignment panels from barcoded library 10 (bar10 = LA Taq). Bottom

901  two from Barcode 11 (bar11 = Taq). Minimap2 and BWA-MEM were used to map reads

902  basecalled with Albacore (worst), Guppy, or FlipFlop (best) to HXB2. Color-coding: Red below

903  genome scale marks 20 SNVs across the HIV segment of pHXB2_D. Purple is an insertion in a

904  given read relative to reference. White is either a deletion in a given read or space between two

905  aligned reads. Gray in alignment field means base same as reference, and in coverage field

906  means major allele is at least 95% the same as reference. Per-read "insertions" and "deletions"

907  do not necessarily represent true insertions or deletions actually present in the sample, because

908  each read is likely an imperfect independent observation. Automated assembly followed by

909  manual consensus building converts these overlapping reads into approximations of the ground

910  truth. "Unbiased" refers to not amplifying a given region (e.g., pol) before ligating ONT

911  sequencing adapters. In the present approach, the tagmentation process randomly cuts DNA,

912  creating ~2000 bp pieces. Tagmented DNA is then amplified based on tagmentation adapters.
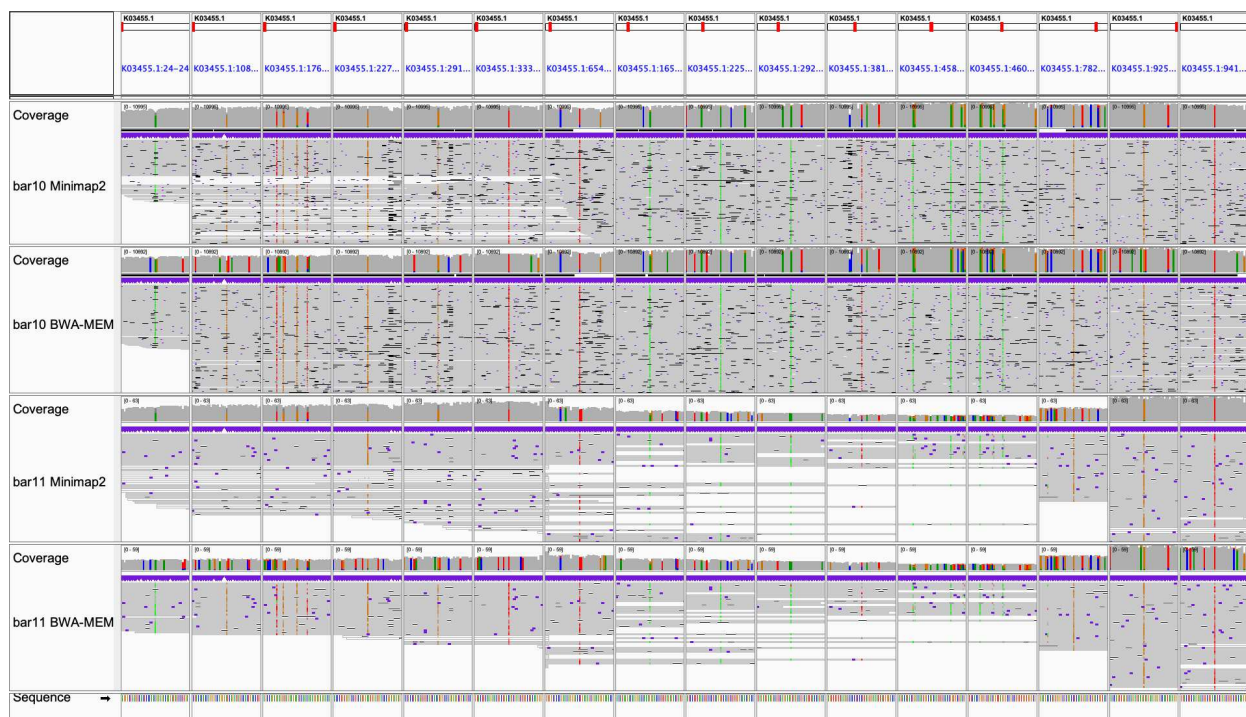
913

914 **Supplemental Figure 2: Reads map well to HIV-1 NL4-3 segment of pNL4-3 assembly**

915 **because NL4-3 LTRs are distinct.**

916



917
918

919 **Supplemental Figure 3A: HIV single nucleotide variants (SNVs) in pHXB2_D. ONT**

920 **basecaller = Albacore (worst).**



921
922

923  **Supplemental Figure 3B: HIV single nucleotide variants (SNVs) in pHXB2_D. ONT**
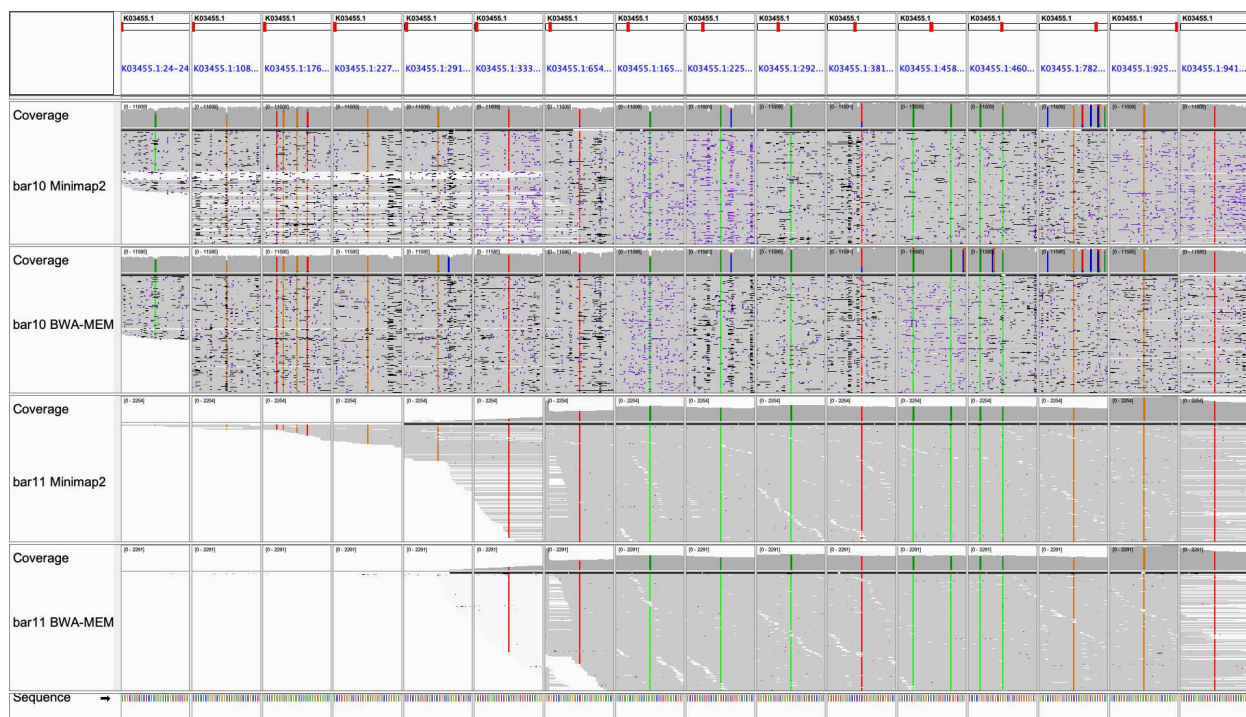
924  **basecaller = Guppy.**



925
926

927    **Supplemental Figure 3C: HIV single nucleotide variants (SNVs) in pHXB2_D. ONT**

928    **basecaller = FlipFlop (best).**



929
930

931 **Supplemental Figure 3D: HIV single nucleotide variants (SNVs) in pHXB2_D, long vs.**

932 **short reads (HIV genome).**



933
934

935    **Supplemental Figure 3E: HIV single nucleotide variants (SNVs) in pHXB2_D, long vs.**

936    **short reads (20 SNV-focused).**



937
938

939    Supplemental Figure 3A: HIV single nucleotide variants (SNVs) in pHXB2_D. ONT basecaller

940    = Albacore (worst). Gray indicates per-base consensus accuracy ≥ 80%. These alignments are

941    the noisiest (less gray and most divergent from reference) between Supplemental Figures 3A,

942    3B, and 3C.

943    Supplemental Figure 3B: HIV single nucleotide variants (SNVs) in pHXB2_D. ONT basecaller

944    = Guppy.

945    Supplemental Figure 3C: HIV single nucleotide variants (SNVs) in pHXB2_D. ONT basecaller

946    = FlipFlop (best). These alignments are the least noisy (most gray and like reference) between

947    Supplemental Figures 3A, 3B, and 3C.

948    Supplemental Figure 3D: HIV single nucleotide variants (SNVs) in pHXB2_D, long vs. short

949    reads (HIV genome). Long reads outperform short reads at HIV-1 LTRs. ONT basecaller =

950    FlipFlop. Short read as single-end 150, clipped to 142, provided by external core. Mappers =

951    Minimap2 (better), BWA-MEM.

952    Supplemental Figure 3E: HIV single nucleotide variants (SNVs) in pHXB2_D, long vs. short
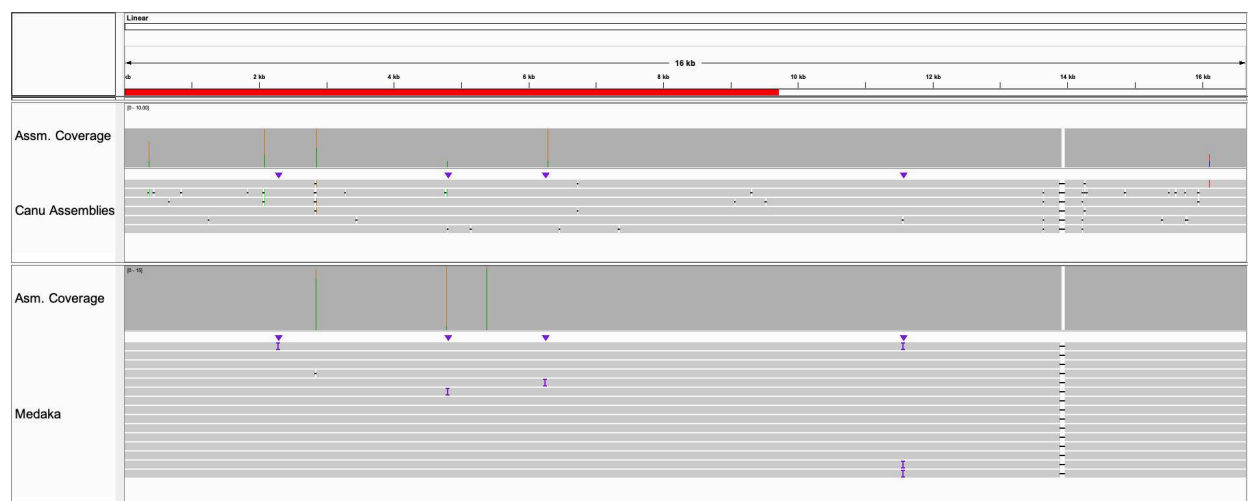
953    reads (SNV-focused).

954

955 **Supplemental Figure 4A: Assembling pHXB2_D from long reads only, varying basecaller**
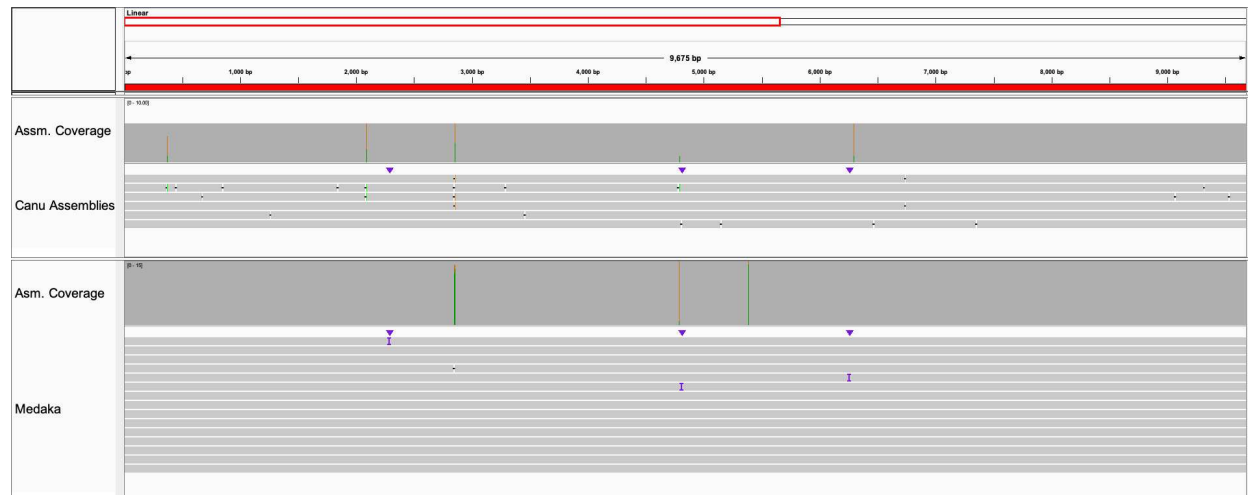
956 **and polymerase.**



957
958 Each pane (n=6) summarizes the results of contig curation. Divergence from reference decreases

959 with newer basecallers, and with long amplicon DNA polymerase (Sigma-Aldrich Taq vs. LA

960 Taq by Takara). Errors in assembly occurred at homopolymers (most often deletions not visible

961 at this resolution; see **Supplemental Figure 6**), dimer or trimer runs. bar10 = LA Taq library.

962 bar11 = Taq library. pHXB2_D Genbank:MW079479. Best contigs presented, manually curated

963 to match pHXB2_D coordinates. Note LTRs (beginning and terminal 634 bp of red bar) are

964 resolved in almost all assemblies. See **Supplemental Table 3** for differences between assemblies

965 and the reference (left red). Plasmid backbone (right) differences are not reported.

966

967 **Supplemental Figure 4B: ONT errors corrected by polishing ONT-only assemblies.**



968



969

970 Assemblies polished with Medaka (ONT). Top: pHXB2_D genome. Bottom: HIV-only segment.

971 The best polished assembly had one error in the entire plasmid (1 error out of 16,722 bases), with

972 a corresponding consensus accuracy of 99.99402%. This happened to be in HIV segment (HIV-1

973 between position 1 and 9719; 1 error out of 9719 bases ), with corresponding accuracy of
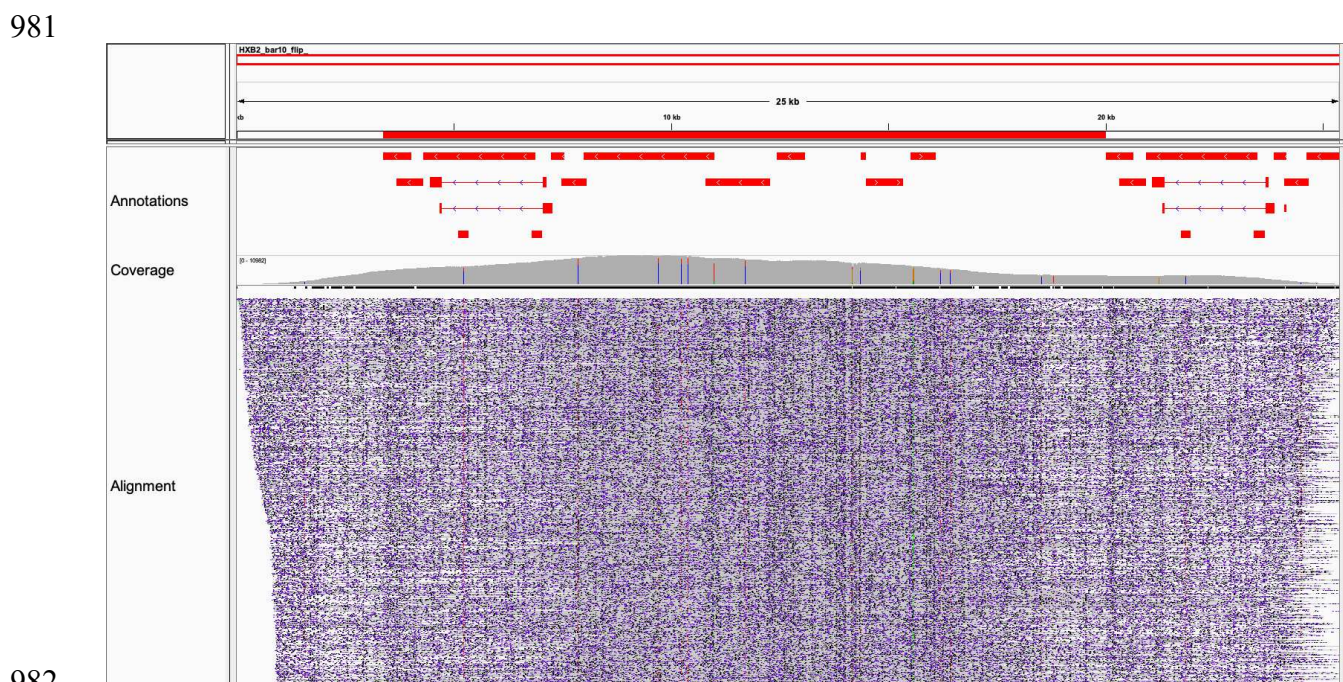
974 99.989711%. Note the conserved 52 bp gap in the backbone of pHXB2_D was redundant

975 sequence included in the short-read assembly from the core. It was not supported by long-read

976 data, and therefor was validated as a technical artifact from the core's pipeline. Reference: short-

977    read assembly. LTRs (beginning and terminal 634 bp of red bar) are resolved in polished
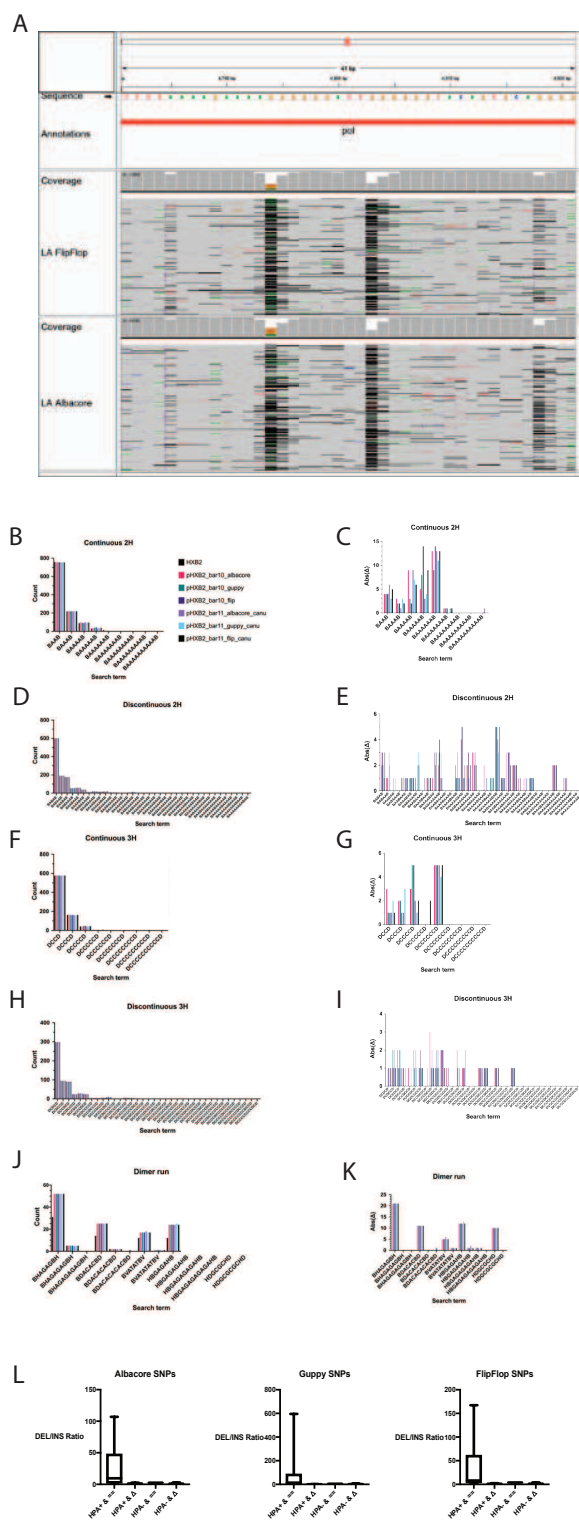
978    assemblies.

979 **Supplemental Figure 4C: Mappability of long reads over contigs during assembly quality**

980 **control.**

981



982

983 Coverage depends on context. Abrupt changes in coverage from terminal regions of HXB2

984 **(Figure 1F, Supplemental Figures 1 and 3)** were artifacts from supplying mappers with an

985 HIV reference without a plasmid backbone. Long reads from barcode 10 (LA Taq) mapped with

986 minimap2 [23] and "reference" contig from assembly (Canu v1.8) with basecalled data

987 (FlipFlop). Stripes in this figure are not SNVs. They represent technical variability at

988 homopolymers. Assemblies were manually curated to start with 5' LTR in the sense orientation,

989 leaving the plasmid backbone on the left. Because there were not real insertions in the HIV

990 segment, the HIV coordinates are the same as HXB2 (both 9719 bp long). Compare with 52 bp

991 technical artifact from the core's short read assembly in **Supplemental Figure 4A, 4B**.

992

993 **Supplemental Figure 5: Homopolymers and dimer runs are ONT artifacts in unpolished**

994 **assemblies.**



995

996     Supplemental Figure 5A: A set of homopolymer tracks from HXB2 plasmid. Alignments with

997     BWA-MEM shown from FlipFlop (top) and Albacore (bottom) basecalled reads. Mapping is

998     pre-assembly.

999     Supplemental Figure 5B: Continuous 2H counts in unpolished assemblies. 2H = A or T

1000    homodimers.

1001    Supplemental Figure 5C: Continuous 2H Absolute Difference.

1002    Supplemental Figure 5D: Discontinuous 2H counts in unpolished assemblies.

1003    Supplemental Figure 5E: Discontinuous 2H Absolute Difference.

1004    Supplemental Figure 5F: Continuous 3H counts in unpolished assemblies. 3H = C or G

1005    homodimers.

1006    Supplemental Figure 5G: Continuous 3H Absolute Difference.

1007    Supplemental Figure 5H: Discontinuous 3H counts in unpolished assemblies.

1008    Supplemental Figure 5I: Discontinuous 3H Absolute Difference.

1009    Supplemental Figure 5J: Dimer run counts in unpolished assemblies.

1010    Supplemental Figure 5K: Dimer run Absolute Difference. Dimer runs as pairs are the most

1011    problematic, with runs as triplets being resolvable by ONT.
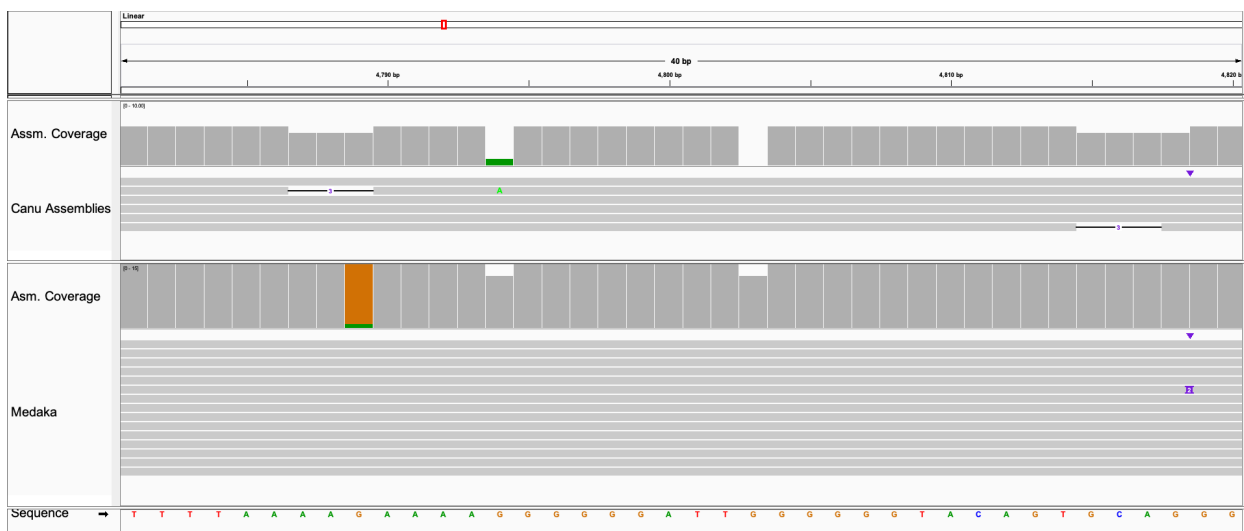
1012    Supplemental Figure 5L: The ratio of deletions to insertions is higher at mismatches both

1013    adjacent to homopolymers and similar to neighbor bases. Box plot shows median ("x" is mean)

1014    and quartile ranges. Y-axis is ratio. HPA: homopolymer-adjacent. ==: same as neighbor base. Δ:

1015    different than neighbor base. Higher coverage (above ~10) usually makes up for current error

1016    profile. Above true for Albacore, Guppy, and FlipFlop.

1017 **Supplemental Figure 6: Assembly partially resolved homopolymers, which are improved**

1018 **by polishing**



1019



1020

1021 Top: Six ONT-only assemblies. Bottom: polished ONT-only assemblies, varying Medaka

1022 models. Deletions at 5' of G homopolymers were not corrected, regardless of basecaller or Taq

1023 isoform. Note that polishing was not performed. IGV window is Linear:4,781-4,820. Bottom:
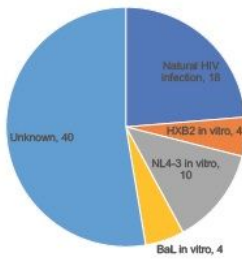
1024    polishing canu assemblies with medaka abrogated most ONT artifacts.  Best medaka setting

1025    tested: r941_min_high_g330.

1026

# Figures



## Figure 1

HIV information in pHXB2_D is recovered by long-read sequencing and mapping. Figure 1A: HXB2 is still a commonly used resource. It is the reference HIV-1 genome, derived from one of the earliest clinical isolates. While older HIV samples are occasionally rediscovered, they are not made routinely available to

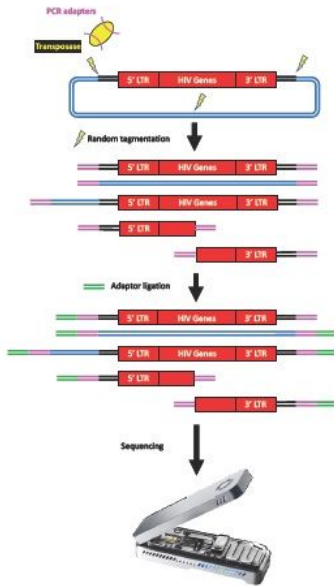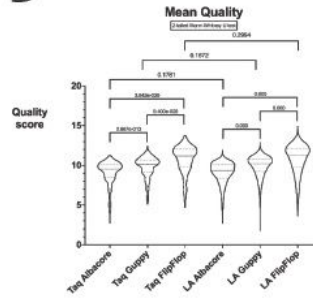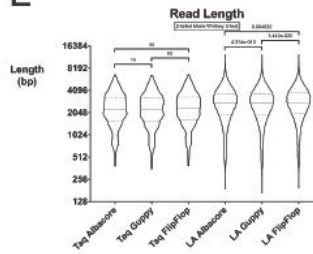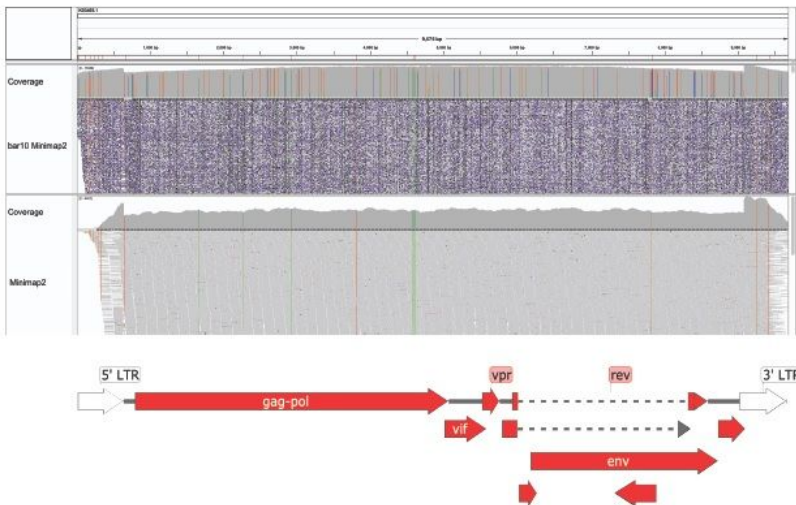researchers. All public HIV-1 RNA-seq datasets were obtained from the NCBI SRA with the following search phrase: "HIV-1" AND "RNA-seq". Metadata from these 2527 runs (number current as of 7/21/2020) were used to make a pie chart summary. Figure 1B: HIV information comes from three main sources: proviruses (HIV sandwiched between two assumedly identical full-length long terminal repeats (LTRs)), unspliced HIV mRNAs (also known as viral genomes) starting from the transcription start site and ending in the 3' LTR [4], and engineered proviruses recovered in their entirety or stitched together from multiple isolates like NL4-3 [18]. Figure 1C: ONT library prep pipeline. Tagmentation cleaves double-stranded DNA, ligating barcoded PCR adapters (magenta). PCR-adapted DNA may be amplified. After amplification and cleanup, ONT sequencing adapters (green) are ligated. Barcoded samples may be pooled and sequenced. Figure 1D: Newer basecallers increase read mean quality. Median (big dash) and quartiles (little dash). Effect of enzyme version was not statistically significant. Figure 1E: Read stats with different callers/aligners. Median (big dash) and quartiles (little dash). Read lengths increase with higher fidelity Taq. Figure 1F: Sequencing coverage with long- vs. short-read single-end 150 bp (trimmed to 142 bp) DNA sequencing. Long-read sequencing covers ambiguously mappable areas missed by short read in HXB2 reference Genbank:K03455.1 (Supplemental Figures 3D,3E), but at the expense of accuracy near homopolymers longer than about 4 nucleobases (Supplemental Figure 5). Short-read mapping fails at repetitive elements longer than their read lengths (Supplemental Figures 3D,3E). Long read Minimap2 settings: map-ont -k15. Short read Minimap2 settings: Short reads without splicing (-k21 -w11 --sr -F800 -A2 -B8 -O12,32 -E2,1 -r50 -p.5 -N20 - f1000,5000 -n2 -m20 -s40 -g200 -2K50m --heap-sort=yes --secondary=no) (sr).

@6fbf0205-5195-460e-8e28-930db50e5d79 runid=0b284792282af9a6d7275cfca845556bb4b8ac7f sampleid=Student_Genomics_Run_1 read=87223 ch=460 start_time=2018-05-09T20:57:49Z barcode=barcode10

## Figure 2

Longest read containing complete full-length HIV-1 reference HXB2. The 5th longest read in the barcode 10 set (read ID 6fbf0205-5195-460e-8e28-930db50e5d79) contained full-length HIV-1. Query (full read) blastn against HIV (taxid:11676) returned 92.95% identity to HIV-1, complete genome (Genbank:AF033819.3). Limiting query to HXB2 (red) blastn against Nucleotide collection nr/nt returned

100% coverage and 93.02% identity to HIV-1 HXB2. This read was 11,487 bases long, with mean quality score 11.984396. Basecalled using Guppy 2.3.1 with FlipFlop config.

Figure 3A:

```
CLUSTAL format alignment by MAFFT (v7.475)

K03455.1_5'LTR   tggaagggctaattcactcccaacgaagacaagatatccttgatctgtggatctaccaca
pHXB2_D_5'LTR    tggaagggctaattcactcccaaagaagacaagatatccttgatctgtggatctaccaca
pHXB2_D_3'LTR    tggaagggctaattcactcccaaagaagacaagatatccttgatctgtggatctaccaca
K03455.1_3'LTR   tggaagggctaattcactcccaaagaagacaagatatccttgatctgtggatctaccaca
                 ***************  ********* ********* ************************

K03455.1_5'LTR   cacaaggctacttccctgattagcagaactacacaccagggccagggatcagatatccac
pHXB2_D_5'LTR    cacaaggctacttccctgattagcagaactacacaccagggccaggggtcagatatccac
pHXB2_D_3'LTR    cacaaggctacttccctgattagcagaactacacaccagggccaggggtcagatatccac
K03455.1_3'LTR   cacaaggctacttccctgattagcagaactacacaccagggccaggggtcagatatccac
                 ********************************************** *.************

K03455.1_5'LTR   tgacctttggatggtgctacaagctagtaccagttgagccagagaagttagaagaagcca
pHXB2_D_5'LTR    tgacctttggatggtgctacaagctagtaccagttgagccagataaggtagaagaggcca
pHXB2_D_3'LTR    tgacctttggatggtgctacaagctagtaccagttgagccagataaggtagaagaggcca
K03455.1_3'LTR   tgacctttggatggtgctacaagctagtaccagttgagccagataagatagaagaggcca
                 *************************************** *** ******.****

K03455.1_5'LTR   acaaaggagagaacaccagcttgttacaccctgtgagcctgcatggaatggatgacccgg
pHXB2_D_5'LTR    ataaaggagagaacaccagcttgttacaccctgtgagcctgcatgggatggatgacccgg
pHXB2_D_3'LTR    ataaaggagagaacaccagcttgttacaccctgtgagcctgcatgggatggatgacccgg
K03455.1_3'LTR   ataaaggagagaacaccagcttgttacaccctgtgagcctgcatgggatggatgacccgg
                 *.*********** ***************** *************** *************

K03455.1_5'LTR   agagagaagtgttagagtggaggtttgacagccgcctagcatttcatcacatggcccgag
pHXB2_D_5'LTR    agagagaagtgttagagtggaggtttgacagccgcctagcatttcatcacgtggcccgag
pHXB2_D_3'LTR    agagagaagtgttagagtggaggtttgacagccgcctagcatttcatcacgtggcccgag
K03455.1_3'LTR   agagagaagtgttagagtggaggtttgacagccgcctagcatttcatcacgtggcccgag
                 ************************************************** *********

K03455.1_5'LTR   agctgcatccggagtacttcaagaactgctgacatcgagcttgctacaagggactttccg
pHXB2_D_5'LTR    agctgcatccggagtacttcaagaactgctgatatcgagcttgctacaagggactttccg
pHXB2_D_3'LTR    agctgcatccggagtacttcaagaactgctgatatcgagcttgctacaagggactttccg
K03455.1_3'LTR   agctgcatccggagtacttcaagaactgctgacatcgagcttgctacaagggactttccg
                 ******************************** .***************************

K03455.1_5'LTR   ctggggactttccaggggaggcgtggcctgggcgggactggggagtggcgagccctcagat
pHXB2_D_5'LTR    ctggggactttccaggggaggcgtggcctgggcgggactggggagtggcgagccctcagat
pHXB2_D_3'LTR    ctggggactttccaggggaggcgtggcctgggcgggactggggagtggcgagccctcagat
K03455.1_3'LTR   ctggggactttccaggggaggcgtggcctgggcgggactggggagtggcgagccctcagat
                 ***********************************************************

K03455.1_5'LTR   cctgcatataagcagctgctttttgcctgtactgggtctctctggttagaccagatctga
pHXB2_D_5'LTR    cctgcatataagcagctgctttttgcctgtactgggtctctctggttagaccagatctga
pHXB2_D_3'LTR    cctgcatataagcagctgctttttgcctgtactgggtctctctggttagaccagatctga
K03455.1_3'LTR   cctgcatataagcagctgctttttgcctgtactgggtctctctggttagaccagatctga
                 ***********************************************************

K03455.1_5'LTR   gcctgggagctctctggctaactagggaacccactgcttaagcctcaataaagcttgcct
pHXB2_D_5'LTR    gcctgggagctctctggctaactagggaacccactgcttaagcctcaataaagcttgcct
pHXB2_D_3'LTR    gcctgggagctctctggctaactagggaacccactgcttaagcctcaataaagcttgcct
K03455.1_3'LTR   gcctgggagctctctggctaactagggaacccactgcttaagcctcaataaagcttgcct
                 ***********************************************************

K03455.1_5'LTR   tgagtgcttcaagtagtgtgtgcccgtctgttgtgtgactctggtaactagagatccctc
pHXB2_D_5'LTR    tgagtgcttcaagtagtgtgtgcccgtctgttgtgtgactctggtaactagagatccctc
pHXB2_D_3'LTR    tgagtgcttcaagtagtgtgtgcccgtctgttgtgtgactctggtaactagagatccctc
K03455.1_3'LTR   tgagtgcttcaagtagtgtgtgcccgtctgttgtgtgactctggtaactagagatccctc
                 ***********************************************************

K03455.1_5'LTR   agacccttttagtcagtgtggaaaatctctagca
pHXB2_D_5'LTR    agacccttttagtcagtgtggaaaatctctagca
pHXB2_D_3'LTR    agacccttttagtcagtgtggaaaatctctagca
K03455.1_3'LTR   agacccttttagtcagtgtggaaaatctctagca
                 **********************************
```

Figure 3B:

```
CLUSTAL format alignment by MAFFT (v7.475)

AF324493.1_5LTR   tggaagggctaatttggtcccaaaaaagacaagagatccttgatctgtggatctaccaca
ACCESSION_TBD_5   tggaagggctaatttggtcccaaaaaagacaagagatccttgatctgtggatctaccaca
AF324493.1_3LTR   tggaagggctaattcactcccaaagaagacaagatatccttgatctgtggatctaccaca
ACCESSION_TBD_3   tggaagggctaattcactcccaaagaagacaagatatccttgatctgtggatctaccaca
                  **************.. *******.********* ************************

AF324493.1_5LTR   cacaaggctacttccctgattggcagaactacacaccagggccagggatcagatatccac
ACCESSION_TBD_5   cacaaggctacttccctgattggcagaactacacaccagggccagggatcagatatccac
AF324493.1_3LTR   cacaaggctacttccctgattggcagaactacacaccagggccaggggtcagatatccac
ACCESSION_TBD_3   cacaaggctacttccctgattggcagaactacacaccagggccaggggtcagatatccac
                  ********************************************** *.************

AF324493.1_5LTR   tgacctttggatggtgcttcaagttagtaccagttgaaccagagcaagtagaagaggcca
ACCESSION_TBD_5   tgacctttggatggtgcttcaagttagtaccagttgaaccagagcaagtagaagaggcca
AF324493.1_3LTR   tgacctttggatggtgctacaagctagtaccagttgagccagataaggtagaagaggcca
ACCESSION_TBD_3   tgacctttggatggtgctacaagctagtaccagttgagccagataaggtagaagaggcca
                  ***************** ****.************** ***** .*************

AF324493.1_5LTR   atgaaggagagaacaacagcttgttacaccctatgagccagcatgggatggaggacccgg
ACCESSION_TBD_5   atgaaggagagaacaacagcttgttacaccctatgagccagcatgggatggaggacccgg
AF324493.1_3LTR   ataaaggagagaacaccagcttgttacaccctgtgagcctgcatggaatggatgaccctg
ACCESSION_TBD_3   ataaaggagagaacaccagcttgttacaccctgtgagcctgcatggaatggatgaccctg
                  **.*********** *****************.****** ****** ***** ***** *

AF324493.1_5LTR   agggagaagtattagtgtggaagtttgacagcctcctagcatttcgtcacatggcccgag
ACCESSION_TBD_5   agggagaagtattagtgtggaagtttgacagcctcctagcatttcgtcacatggcccgag
AF324493.1_3LTR   agagagaagtgttagagtggaggtttgacagccgcctagcatttcatcacgtggcccgag
ACCESSION_TBD_3   agagagaagtgttagagtggaggtttgacagccgcctagcatttcatcacgtggcccgag
                  **.******* **** *****.*********** ********** **** *********

AF324493.1_5LTR   agctgcatccggagtactacaaagactgctgacatcgagctttctacaagggactttccg
ACCESSION_TBD_5   agctgcatccggagtactacaaagactgctgacatcgagctttctacaagggactttccg
AF324493.1_3LTR   agctgcatccggagtacttcaagaactgctgacatcgagcttgctacaagggactttccg
ACCESSION_TBD_3   agctgcatccggagtacttcaagaactgctgacatcgagcttgctacaagggactttccg
                  ***************** ***. ***************** **.****************

AF324493.1_5LTR   ctggggactttccagggaggtgtggcctgggcgggactggggagtggcgagccctcagat
ACCESSION_TBD_5   ctggggactttccagggaggtgtggcctgggcgggactggggagtggcgagccctcagat
AF324493.1_3LTR   ctggggactttccaggggaggcgtggcctgggcgggactggggagtggcgagccctcagat
ACCESSION_TBD_3   ctggggactttccaggggaggcgtggcctgggcgggactggggagtggcgagccctcagat
                  *****************.****.***********************************

AF324493.1_5LTR   gctacatataagcagctgctttttgcctgtactgggtctctctggttagaccagatctga
ACCESSION_TBD_5   gctacatataagcagctgctttttgcctgtactgggtctctctggttagaccagatctga
AF324493.1_3LTR   gctgcatataagcagctgctttttgcctgtactgggtctctctggttagaccagatctga
ACCESSION_TBD_3   cctgcatataagcagctgctttttgcctgtactgggtctctctggttagaccagatctga
                  ***.*****************************************************

AF324493.1_5LTR   gcctgggagctctctggctaactagggaacccactgcttaagcctcaataaagcttgcct
ACCESSION_TBD_5   gcctgggagctctctggctaactagggaacccactgcttaagcctcaataaagcttgcct
AF324493.1_3LTR   gcctgggagctctctggctaactagggaacccactgcttaagcctcaataaagcttgcct
ACCESSION_TBD_3   gcctgggagctctctggctaactagggaacccactgcttaagcctcaataaagcttgcct
                  ***********************************************************

AF324493.1_5LTR   tgagtgcttcaaagtagtgtgtgcccgtctgttgtgtgactctggtaactagagatccctc
ACCESSION_TBD_5   tgagtgcttcaaagtagtgtgtgcccgtctgttgtgtgactctggtaactagagatccctc
AF324493.1_3LTR   tgagtgcttcaagtagtgtgtgcccgtctgttgtgtgactctggtaactagagatccctc
ACCESSION_TBD_3   tgagtgcttcaagtagtgtgtgcccgtctgttgtgtgactctggtaactagagatccctc
                  ********.***************************************************

AF324493.1_5LTR   agacccttttagtcagtgtggaaaatctctagca
ACCESSION_TBD_5   agacccttttagtcagtgtggaaaatctctagca
AF324493.1_3LTR   agacccttttagtcagtgtggaaaatctctagca
ACCESSION_TBD_3   agacccttttagtcagtgtggaaaatctctagca
                  **********************************
```
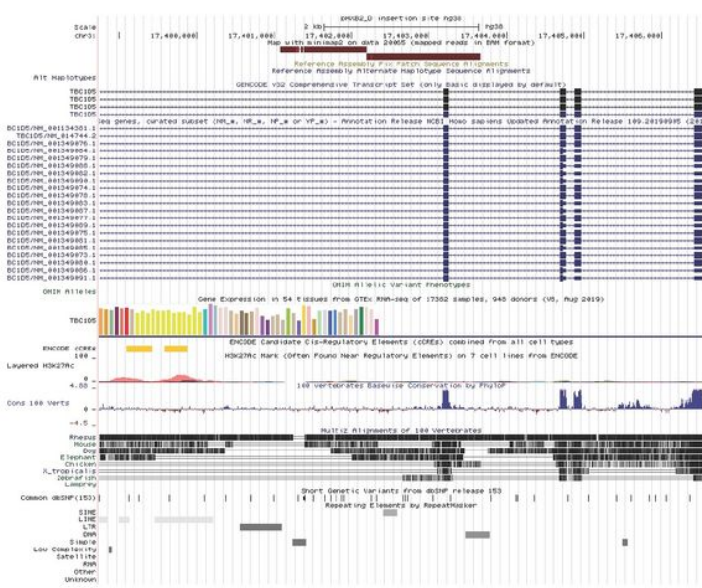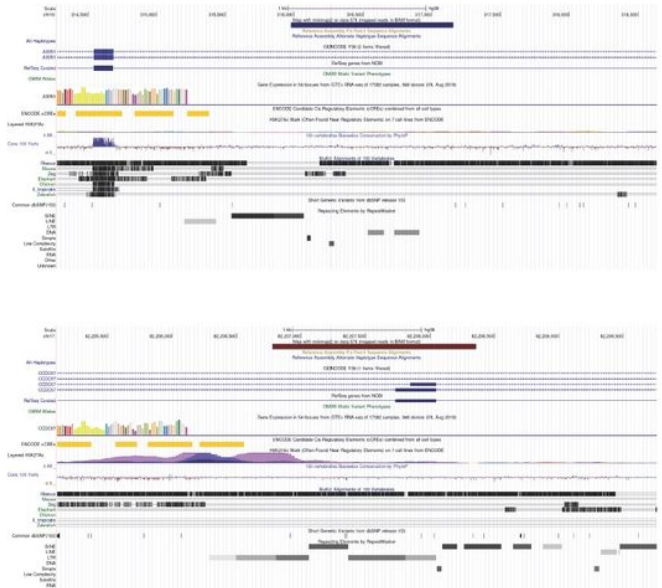
3A

3B

Figure 3

3A: pHXB2_D has identical LTRs, resolving likely errors in HXB2 (K03455.1) 3B: pNL4-3_gag-pol(Δ1443-4553)_EGFP (ACCESSION_TBD) has distinct LTRs, consistent with pNL4-3 (AF324493.1)

4A                4B

**Figure 4**

4A: HXB2 integration site 4B: NL4-3 integration sites Figure 4A: pHXB2_D's, and therefore HXB2's, integration site is unambiguously singular (falls outside of annotated repeat), and in the same orientation (minus strand relative to hg38) as target gene TBC1D5. Alignment quality is 60 for both homology arms (Supplemental Table 2). Features captured by homology arms in pHXB2_D and other clones verified as proviruses in the present study are consistent with HIV-1 integration behavior [44]. Visualized in UCSC Genome Browser [45]. Figure 4B: pNL4-3_gag-pol(Δ1443-4553)_EGFP's, and therefore NL4-3's, integration sites fall on annotated repeats, the longer reads help to locate both sites. Alignment quality is 60 for both homology arms (Supplemental Table 2). These integration sites would likely be missed by any method leveraging reads shorter than the homology arms.
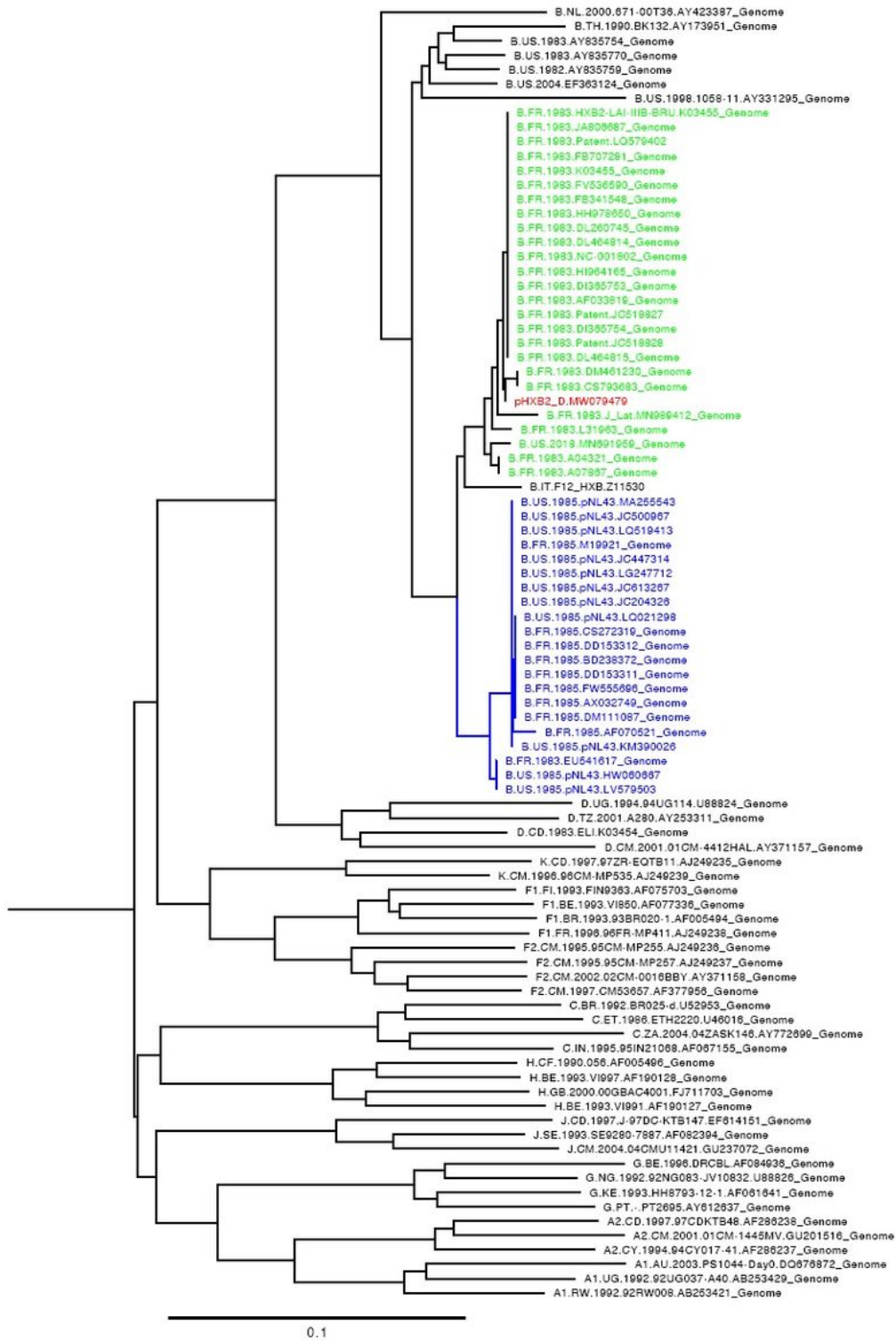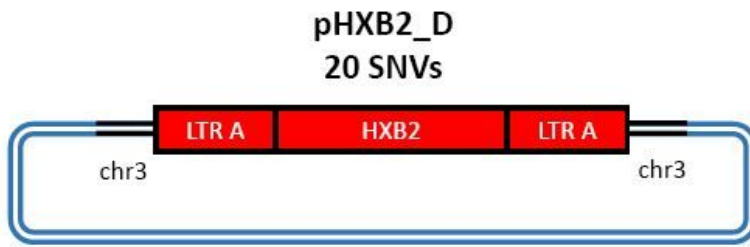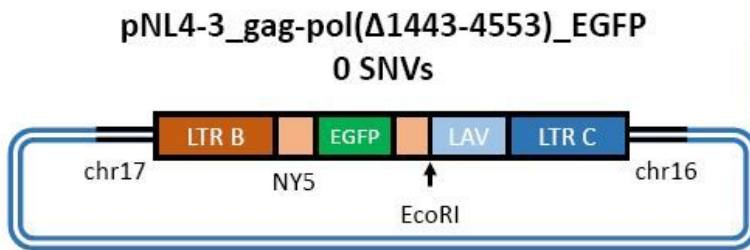
**Figure 5**

pHXB2_D provenance and top 50 neighbors

**Figure 6**

Summary of long- vs. short-read mapping by ability to phase LTRs

# Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- SupplementalInformation.pdf
- Generetal.pHXB2DTablesandSupplementalTables.xlsx